

JUNO: Optimizing High-Dimensional Approximate Nearest Neighbour Search with Sparsity-Aware Algorithm and Ray-Tracing Core Mapping

Zihan Liu^{1,2}, Wentao Ni¹, Jingwen Leng^{*1,2}, Yu Feng³

Cong Guo^{1,2}, Quan Chen^{1,2}, Chao Li^{1,2}, Minyi Guo^{*1,2} and Yuhao Zhu³

¹ Shanghai Jiao Tong University, ² Shanghai Qi Zhi Institute, ³ University of Rochester
{altair.liu, wennitao, guocong}@sjtu.edu.cn, { leng-jw, chen-quan, lichao, guo-my}@cs.sjtu.edu.cn
yfeng28@ur.rochester.edu, yzhu@rochester.edu

Abstract

Approximate nearest neighbor (ANN) search is a widely applied technique in modern intelligent applications, such as recommendation systems and vector databases. Therefore, efficient and high-throughput execution of ANN search has become increasingly important. In this paper, we first characterize the state-of-the-art product quantization-based method of ANN search and identify a significant source of inefficiency in the form of unnecessary pairwise distance calculations and accumulations. To improve efficiency, we propose JUNO, an end-to-end ANN search system that adopts a carefully designed sparsity- and locality-aware search algorithm. We also present an efficient hardware mapping that utilizes ray tracing cores in modern GPUs with pipelined execution on tensor cores to execute our sparsity-aware ANN search algorithm. Our evaluations on four datasets from 1 to 100 million search points demonstrate $2.2\times$ - $8.5\times$ improvements in search throughput. Moreover, our algorithmic enhancements alone achieve a maximal $2.6\times$ improvement on the hardware without the acceleration of the RT core.

1. Introduction

Computer applications are becoming more intelligent with the breakthrough from deep learning [47, 70]. One of the key data structures in these applications is high-dimensional embedding vector [2, 15, 59, 65, 71], which are usually generated from some kind of transformation or learning to the raw data like text, images, audio, video, and others. Using embeddings allows for fast and accurate nearest neighbor (NN) search and retrieval of data based on their distance. For example, we can perform the NN search to find images similar to a given image based on their content and style.

Embedding vectors are often high-dimensional, ranging from tens to thousands. The curse of dimension makes the exact NN search computationally expensive so that approximate nearest neighbor (ANN) search becomes increasingly popular in industrial practice [40, 7, 69, 22, 11]. ANN search trades search quality (measured in recall) for search performance (measured in throughput) in different scenarios, with GPUs being widely used to achieve higher throughput.

The IVFPQ technique, which combines the inverted file index (i.e., IVF) with product quantization (i.e., PQ), is the most commonly used ANN method [40, 36]. This method is often combined with other techniques to improve performance [43]. The PQ-based approach encodes search point projections using a codebook in every subspace offline and then searches for the nearest neighbors by accumulating the distance information distributed in these subspaces online. However, this process requires a significant number of pairwise distance calculations between entry and query projections in low-dimensional subspaces, as well as frequent look-ups to calculate the total distance for every query.

In this work, we use the state-of-the-art ANN framework FAISS [36] to study the efficiency of the IVFPQ method. Although the framework employs hundreds of codebook entries to encode search points in each subspace, only a fraction of these entries is used by the top-100 results returned by a query. Notably, in some subspaces, all top-100 search points are encoded with just one entry. This entry-level sparsity provides opportunities to avoid calculating the pairwise distances of unused entries and to skip the distance look-up and accumulation of search point projections encoded by these unused entries.

In addition, our findings suggest that entries with high usage frequency demonstrate a significant degree of spatial locality, rendering the exploitation of sparsity more advantageous. Despite being sparsely distributed in memory, these entries are closely positioned in Euclidean space. In some situations, by selecting the 25% closest entries, we can obtain all the used entries in response to a query within a subspace. Therefore, only entries that are close to the query projections are essential.

We propose a selective codebook construction algorithm to exploit the above sparsity and spatial locality to accelerate high-dimensional ANN search. Our approach involves using an adaptive distance threshold in each subspace to only select necessary entries. We have identified a strong correlation between the distance threshold and search point density, and we exploit it by training a simple regression model offline, with the density serving as the input. During runtime, we utilize the threshold value determined by the regression model to choose a small fraction of the search points.

Our sparsity-aware ANN method heavily uses the distance comparison operation, which is well suited for ray tracing (RT)

* Jingwen Leng and Minyi Guo are corresponding authors of this paper.

cores in modern GPUs [49, 50, 54]. The RT cores implement the tree-based intersection check [66, 17], which has a reduced logarithmic time complexity. As a result, RT cores can find objects intersected with a given ray with high efficiency [48]. Obviously, this is quite similar to our algorithmic intuition that finds close entries for a query projection in every subspace. Intuitively, we can organize entries as objects and query projection as rays in every subspace and let the RT core efficiently find the intersected entries of a query ray.

We present JUNO, a fast and high-throughput ANN search system that employs algorithmic enhancement and RT core mapping to leverage sparsity and spatial locality. However, utilizing the RT core in an enhanced ANN algorithm still poses unique challenges. Firstly, calculating distances for selected codebook entries after filtering is still necessary. Secondly, naively implementing adaptive dynamic radius in RT cores would cause unacceptable scene preparation overhead during runtime. To overcome these issues, we exploit the concept of "time" in the ray tracing scenario. Specifically, we utilize the hit time results from RT cores to efficiently calculate distances, thereby avoiding costly global memory access. Additionally, we convert the dynamic distance threshold to a dynamic maximum travel time for rays, thereby avoiding online scene preparation. Finally, we optimize JUNO to support inner product similarity and efficient RT-Tensor core pipelining [55].

We evaluate JUNO on multiple datasets sizing from 1M to 100M with both L2 distance and inner product similarity. Together they deliver an average (a maximum) of $4.4\times$ ($8.5\times$) and $2.1\times$ ($3.2\times$) improvement in search throughput in low and high search quality against the baseline. Moreover, the improvement is bound by the performance of RT cores.

We make the following main contributions in this work:

- We study the inefficiency of the typical $IVFPQ$ pipeline and identify sparsity and spatial similarity in codebook usage.
- We design a threshold-based selective algorithm to rapidly filter out the unnecessary search points leveraging the sparsity and spatial locality and propose a mapping for our algorithm to run on the RT core.
- To the best of our knowledge, we are the first to study how to generalize the existing kNN-RT core mapping to ANN search with arbitrary dimensions, in aspects of approximation method, metrics and system design. Based on our experimental analysis and insights, we propose JUNO, an end-to-end high-dimensional ANN search engine with both algorithmic enhancement and optimized hardware mapping.
- We quantify the effectiveness of JUNO over existing ANN framework FAISS [36], with detailed breakdown analysis.

2. Background

This section introduces the typical process of ANN search, ray tracing pipeline, and its application in 2D/3D ANN search.

Table 1: Notations used in this work.

N	# search points	D	search points dimension
C	# clusters	E	# codebook entries
x	search point vector	M	subspaces dimension
q	query vector	r	bounding radius
s	subspace id	e	codebook entry id
$nprobs$		# chosen clusters in <i>filtering</i>	

2.1. Approximate Nearest Neighbor (ANN) Search

The objective of the nearest neighbor (NN) search is to identify the top-k most similar points to a query within a given set. Typically, the similarity between two points is determined using L2 distance or inner (or dot) product, as illustrated in Equ. 2.1. L2 distance is lower-is-better and widely used in measuring image similarity [68]. Inner product is higher-is-better and frequently used in large language models (LLMs) and Transformer architectures [56, 60].

$$L2(\mathbf{q}, \mathbf{x}) = \sum_{i=0}^{D-1} (x_i - q_i)^2, \quad IP(\mathbf{q}, \mathbf{x}) = \sum_{i=0}^{D-1} x_i \cdot q_i$$

The exact NN search is often expensive and thus slow. Many practical cases can tolerate certain search inaccuracies, which can be exploited to improve the search throughput. This is referred to as approximate NN (ANN) search [58]. One of the most popular ANN techniques is inverted file index with product quantization ($IVFPQ$), which is used by top-performing ANN frameworks such as FAISS [36] and ScANN [28]. Notices that there are other indexing and encoding techniques, and we will discuss their details in Sec. 7.

Fig. 1 shows an example of the $IVFPQ$ technique, with notations defined in Tbl. 1, which has an offline (top) and an online component (bottom). The offline component relies on the inverted file index (IVF) and product quantization (PQ), and the online component consists of *filtering*, *L2-LUT construction* and *distance calculation*, as described as follows.

Inverted File Index (IVF) ① Given a query point q , the $IVFPQ$ technique performs the first coarse-grained filtering step to identify a set of candidates from all search points. This step is commonly implemented through k -means, which generates C clusters on N search points of full dimension D . The inverted file index (IVF) stores the associated search points for each cluster centroid. This step calculates the distance between query point q and each centroid and identifies the closest centroid(s). The IVF data structure lets us quickly locate the associated search points for selected centroid(s).

Product Quantization (PQ) The PQ method is widely used for vector compression. ② Initially, the D -dimensional space is divided into D/M M -dimensional subspaces. ③ Next, in each subspace of total D/M subspaces, E clusters are generated with projections of residuals (in this specific subspace) between search points and first cluster centroids. We refer to these E clusters as the "second" clusters. The centroids from the second clusters are combined to form the codebook,

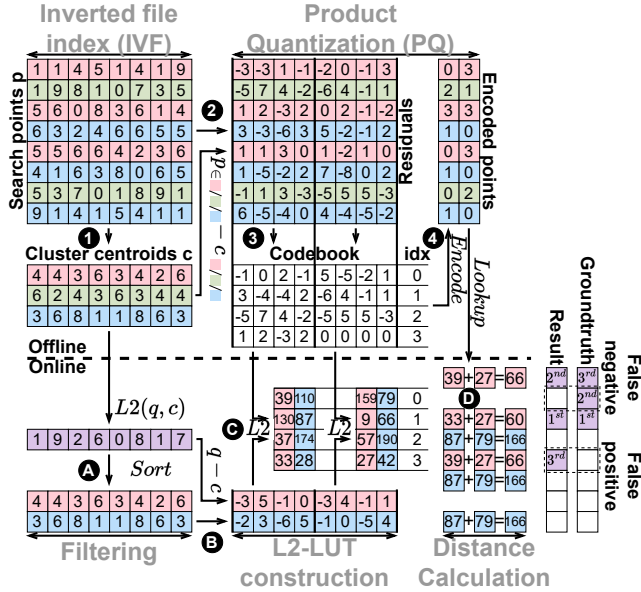


Figure 1: The example of offline training (top) and online searching (bottom) component of IVFPQ-based ANN search.

with each centroid serving as a codebook entry. ④ Finally, search points are encoded using this codebook by replacing the projection on a subspace with the cluster ID to which the residual projection belongs. PQ reduces the storage space required for each search point originally in float format from $D \times \text{sizeof(float)} \times 8$ bits to $(D/M) \times \log_2 E$ bits.

Online Search. After the IVF and the PQ are trained, the online search process begins. ① When a query arrives, it first calculates the pairwise distances with C cluster centroids in IVF. The n probs closest centroids and their corresponding clusters are chosen. Subsequent processes only occur on the search points that belong to selected clusters. We refer to this initial stage as the *filtering stage* following previous literature [40], as illustrated in the bottom left of Fig. 1.

② Next, the query calculates residuals with the n probs selected cluster centroids. ③ For each residual, E pairwise distances are computed between every codebook entry within each subspace. These n probs $\times E \times (D/M)$ pairwise distances are then organized into a look-up table (LUT). We refer to this second stage as *L2-LUT construction stage* [40], as shown in the bottom middle of Fig. 1. ④ Finally, the query iterates over all the encoded search points that belong to the n probs chosen clusters to calculate the overall distance. For a given encoded search point with s^{th} subspace encoded with codebook entry e , assuming it belongs to the np^{th} cluster, the total distance is calculated by summing up all the $L2\text{-LUT}[np][s][e]$ values. We refer to this third stage as *distance calculation stage* [40], as shown in the bottom right of Fig. 1.

After calculating the total distance between all encoded search points, the query then sorts the results and selects the top- k closest neighbors. It’s important to note that an ANN search may yield false positives and false negatives. For ex-

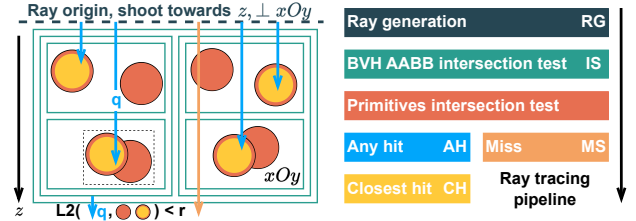


Figure 2: The ray tracing pipeline for RT core.

ample, the 2nd point could be the second closest to the query, but is ignored in ANN search, as depicted in Fig. 1.

2.2. Accelerating NN Search with Ray Tracing Core

Ray tracing (RT) is a distinct rendering pipeline from traditional rasterization [29, 44]. In the RT pipeline, rays are cast from a camera through each pixel into a virtual scene, simulating their interactions with objects to render pixels with accurate colors. However, the RT pipeline can be computationally expensive and time-consuming due to the need to track interactions between rays and objects, resulting in a large-scale operation. To accelerate this process, NVIDIA introduced GPUs with dedicated hardware RT accelerators, known as RT cores [49]. The RT core utilizes a BVH tree-based RT algorithm in hardware, which efficiently finds intersections between lines and surfaces in 3D Euclidean space. Fig. 2 illustrates the typical BVH tree-based RT pipeline.

The NVIDIA Ray Tracing (RT) cores consist of two core hardware components for accelerating the typical ray tracing pipeline: the Axis Aligned Bounding Box (AABB) intersection test and the Bounding Volume Hierarchy (BVH) traversal. The AABB-based intersection test first creates bounding boxes (with edges parallel to x -, y -, or z -axes) to bound objects that need to test the intersection with the ray. Then, the ray will conduct a cheap interval-based calculation to test the intersection status with bounding boxes. If a box is intersected with the ray, the ray will further test the intersection status with objects bounded by this tested box. Otherwise, all objects bounded by this box will be ignored. Notice that one AABB can recursively contain smaller AABBs, and finally form a tree-like structure with log-scale depth to total objects number, where a node represents an AABB and its child nodes represent smaller AABBs contained by it. This structure is called the BVH tree and makes the process of finding intersected bounding box recursive. Obviously, there can be a huge amount of conditions and divergences in the tree traversal process. The NVIDIA RT core also implements corresponding hardware to accelerate the traversal process of the BVH tree.

Researchers have utilized the capability of identifying intersections in three-dimensional Euclidean space to apply the RT core in two-dimensional/three-dimensional nearest neighbor search [75]. The basic idea is to first organize N search points into N circles located in the xOy plane, each with a bounding radius r . Subsequently, the queries can be converted into

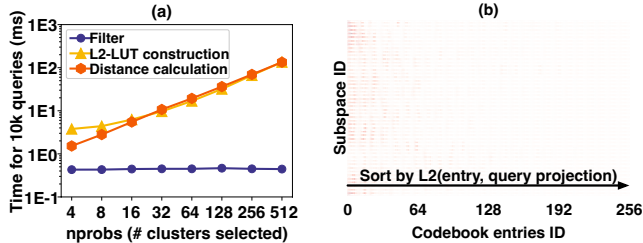


Figure 3: (a) Execution time breakdown of searching queries in DEEP1M using FAISS. (b) Codebook entries usage of a single query: higher usage frequency leads to darker color.

rays that originate from the query coordinates and are directed towards the z-axis, as illustrated in Fig. 2. Any circles intersected by a ray indicate that the distance between the query (represented by the ray) and the search points (represented by the circles) is lower than r , and thus has the potential to be the nearest neighbors of the query. For example, in Fig. 2, the ray q intersects with the two circles in the bottom left quadrant, signifying their potential as the nearest neighbors of the query.

The aforementioned approach, which involves utilizing RT core hardware for accelerating NN search, is limited to low-dimensional (2 and 3) spaces, thereby restricting its practicality. Instead, our work endeavors to explore the effective utilization of RT core for more general NN search tasks, with a specific focus on ANN search in high-dimensional spaces.

3. Motivation

This section provides an analysis of the state-of-the-art library of high-dimensional approximate nearest neighbor (ANN) search, FAISS [36]. We begin by measuring the breakdown of execution time for FAISS queries, followed by an analysis of the identified inefficiencies. Finally, we propose optimization takeaways based on the findings of our analysis.

3.1. Execution Time Breakdown

We use the latest version of FAISS [36] and the DEEP1M dataset [4] in our study. Specifically, we configure FAISS with $\langle \text{IVF}4096, \text{PQ}48 \rangle$, where 1,000,000 search points are grouped into 4096 clusters, and the 96-dimensional space is divided into 48 2-dimensional subspaces. We measure the execution time of three stages as mentioned in Sec. 2.1 (*filtering*, *L2-LUT construction*, *distance calculation*) on an NVIDIA Geforce RTX 4090 GPU [52]. For each query, we evaluate its execution time with different $nprobs$, which is a hyper-parameter and means the number of selected clusters in *filtering*.

The experimental results are presented in Figure 3(a), which illustrates the breakdown of execution time. The majority of the execution time is consumed by the *L2-LUT construction* and *distance calculation* stages, accounting for approximately 90% to 99.9% of the total time. Additionally, the time taken by these stages increases linearly with the hyper-parameter $nprobs$, which is set smaller for higher performance and larger for better search quality [41]. On the other hand, the *filtering*

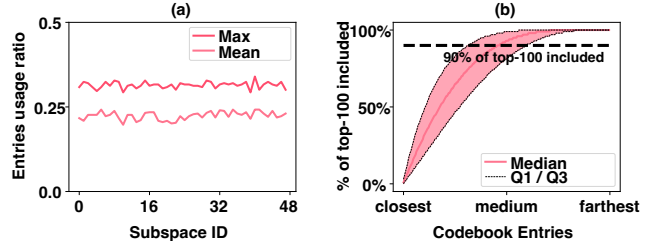


Figure 4: (a) Max used ratio of codebook entries on every subspace of 100 queries. (b) CDF of entries to contain top-100 from closest to farthest. Figures are plotted using DEEP1M dataset.

* $Q1, Q3$ are 1st (25%), 3rd (75%) quantile (percentile).

stage remains relatively stable, as its computation depends on $Q \times D \times C$, where C is independent of $nprobs$. This observation motivates us to focus on optimizing the *L2-LUT construction* and *distance calculation* stages, as we will analyze the inefficiencies in the current approach in below.

3.2. Sparsity of Codebook Entries Used by Top-k Neighbors

The current ANN search implementation, like FAISS, calculates the pairwise distance between the query projection and codebook entries in all subspaces during the *L2-LUT construction* stage. However, our findings indicate that for a single query, only limited codebook entries are used to identify the top-100 search points with the closest proximity.

To illustrate the above point, we calculate the usage frequency of each codebook entry by the top-100 search points. The results in the form of heatmap are presented in Fig. 3(b), where the statistics for all entries from different subspaces are depicted. The shading of the cells corresponds to the frequency of usage, ranging from 0 to 100, with darker colors indicating higher frequency of usage. For example, a $cell[31][114] = 75$ means that in the 31st subspace, 75 among top-100 search points are encoded using the 114th codebook entry. A value of 0 signifies that none of the top-100 search points are encoded with the respective codebook entry.

The codebook entries utilization in each subspace for the DEEP1M dataset [4] with the PQ48 configuration is shown in Fig. 4(a), indicating that, on average, only 25% of the entries (at most 30%) are utilized. To investigate the impact of data distribution on this sparsity, we analyze the SIFT1M (PQ64) [35] and TTI1M [35] (PQ40) datasets, as presented in Fig. 5(a). It can be observed that, on average, these datasets also exhibit less than 30% codebook entries utilization. Exploiting this sparsity can potentially result in a significant reduction of up to 70% floating point operations in pairwise distance calculation for each query, thereby substantially reducing the time required for the *L2-LUT construction* stage.

With above analyses, we derive the first key takeaway: **Codebook entries used by top-100 neighbours are sparse.**

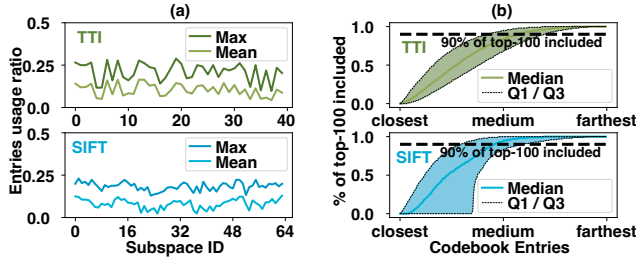


Figure 5: (a) Codebook entries average and maximal usage ratio and (b) CDFs of entries to contain top-100 from closest to farthest.

3.3. Spatial Locality of Codebook Entries Used by Top-k Neighbors

Although we have shown that there exists a high sparsity degree in the codebook entries used by top-k points, it maybe still challenging to convert the sparsity to actual performance speedup as sparsity often results in irregular access patterns. However, as depicted in Fig. 3(b), the codebook entries that are actually used are densely concentrated in the front half of the frequency heatmap. This suggests that the used codebook entries are closer to the query projection compared to others, as the heatmap is sorted based on the distance between the entry and the query projection in each subspace.

To validate this claim, we calculate and plot the cumulative distribution function (CDF) of actual top-100 search points from the closest to farthest entries in every subspace. As depicted in Fig. 4(b), we observe that considering using half of the codebook entries enables us to obtain over 90% of the top-100 search points. The CDFs for the SIFT1M and TTI1M datasets are shown in Fig. 5(b), respectively. Despite having different patterns, both datasets exhibit a similar property, with approximately 50% of the closest entries containing over 90% of the top-100 ground truth.

With above analyses, we derive another key takeaway: **Codebook entries used by top-100 neighbours are closely distributed in the space.**

4. Selective L2-LUT Construction on RT Core

In this section, we present our algorithm enhancement and RT core mapping for an efficient alternative to the original *L2-LUT construction* stage. We begin by introducing our algorithmic intuition and design details, then we introduce how we map the enhanced algorithm to the RT core.

4.1. Threshold Based Selective L2-LUT Construction

In this subsection, we describe and explain our algorithm design to leverage the features above of the current ANN search. The new algorithm aims to replace current *L2-LUT construction* in our searching pipeline and provide significant operation saving with acceptable approximation.

Intuition of Algorithm Enhancement. Our algorithm aims to exploit the sparsity and spatial locality in subspaces, result-

ing in improved search performance through additional approximation techniques alongside product quantization. Specifically, we propose a distance checking-based approach to efficiently prune unnecessary entries in each subspace. This approach involves setting distance thresholds in all subspaces, defining the interested region of a query projection as the union of all points with distances from the query projection being less than the distance threshold, and subsequently discarding codebook entries that fall outside the interested region of a query projection. This pruning strategy is motivated by the findings presented in Sec. 3.

Pruning unnecessary codebook entries can significantly reduce the number of distance calculations. As in Fig. 6, the remaining search point projections that require accessing the L2-LUT and distance calculations decrease linearly with the threshold. This finding suggests that the projections of the top-100 search points are closely distributed around the query projection in a subspace. Consequently, a substantial amount of L2-LUT lookups and distance calculations can be saved by filtering out codebook entries far away from the query projection in each subspace.

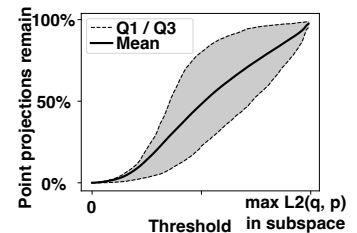


Figure 6: Remained search points that need to accumulate the distance.

Choose Necessary Entries with Efficiency. As discussed in Sec. 4.1, our algorithm leverages the sparsity and spatial locality by selectively choosing the codebook entries that fall within the region of interest in each subspace. The distances between these selected entries and the query projection are then calculated for constructing a *L2-LUT* in a selective manner. To achieve this, we bound smaller entries/boxes into a larger box and hierarchically organize these bounding boxes into a tree structure. As such, we can avoid pairwise distance calculations between all E entries and the query projection in each subspace, and instead perform inside/outside checks with logarithmic complexity of $\log E$.

Determining Proper Threshold. The selection of an appropriate threshold is crucial in our approach, as it is used to distinguish whether an entry falls within the interested region of a query projection. A tight threshold may filter out too many entries, resulting in missed true neighbors, while a relaxed threshold may include unnecessary entries, leading to a waste of time and resources. To determine the optimal threshold, we conduct a thorough study of the relationship between query features and the threshold that can contain the top-100 search points in each subspace.

As shown in Fig. 7(a), we observe a negative correlation between the threshold for containing the top-100 search points

and the region density of the query projection in a subspace. To calculate the region density, we divide the entire subspace into a 100×100 grid and define the density of each cell as the quotient obtained by dividing the number of search point projections falling into that cell by the area of the cell. This finding is reasonable, as a higher threshold is likely to contain many search point projections in high-density regions, resulting in a higher probability of containing the top-100 search points. Conversely, in low-density regions, even the same threshold may not be able to contain 100 search point projections, let alone the top-100 points.

Based on our findings, we propose incorporating dynamic radius mechanics to determine an appropriate threshold during runtime. Initially, we generate a density map and a simple regression model offline. The density map consists of a 100×100 grid for each subspace, where each cell records the density computed as previously described. Subsequently, we randomly select several search point projections to train the regression model, with the region density as input and the threshold to contain the top-100 search points as output. During runtime, we look up the density of the query projection in each subspace and use the regression model to infer a query-specific threshold. Notably, we find that a simple polynomial regression model accurately captures the relationship, resulting in minimal runtime overhead for threshold determination. Once the threshold is determined, we disregard all codebook entries beyond the threshold and compute the distance between the query projection and the remaining codebook entries to create an L2 lookup table (L2-LUT).

Furthermore, we have observed a power-law pattern in the variation of the threshold value, as illustrated in Fig. 7(b). When the distance threshold is scaled down to half of its original value to accommodate the top-100 search points, approximately 90% of these points are retained. This indicates the potential for using a smaller distance threshold to further prune codebook entries. In other words, we can trade-off slightly lower search quality for significantly enhanced search throughput, such as queries per second (QPS). In our work, we provide users with the flexibility to set the threshold value through a dedicated interface, enabling them to make this trade-off according to their specific needs.

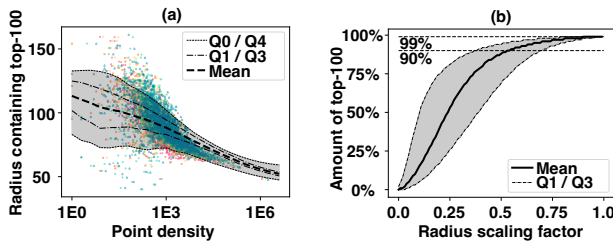


Figure 7: (a) The relation between threshold to contain top-100 search points and region density a query falls into. (b) The amount of top-100 search points contained when threshold scales smaller.

$$*Q0=Q1-1.5 \times IQR, Q4=Q3+1.5 \times IQR, IQR=Q3-Q1.$$

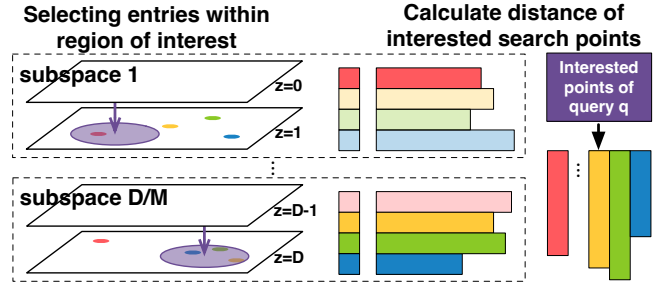


Figure 8: Intuition of hardware mapping of JUNO.

4.2. Mapping the Selective Algorithm to the RT Core

In this subsection, we explain how to map the computation of selective L2-LUT construction described in Sec. 4.1 to the ray-tracing (RT) core in modern GPUs. The key approach is to utilize the RT core to accurately, efficiently determine the intersections of rays with the selected search points within each subspace. Additionally, we exploit the capabilities of the RT core for efficient hit distance calculation, which can be extended to support other distance metrics like inner product.

Algorithm Mapping. The RT core implements bounding box intersection check (AABB) and hierarchical data structure traversing (BVH) in the hardware, which aligns with the computational characteristics of our proposed threshold-based LUT construction. Similar to the original ANN search algorithm discussed in Sec. 2.1, our RT-core-based ANN algorithm comprises both an offline and online component.

Previous research [75] has demonstrated the efficient inside/outside check capability of the RT core through its BVH tree-based algorithm. With this capability, our method has the potential to save computational resources in determining whether an entry falls inside the region of interest within a query projection in a subspace. To implement this, we align our distance-checking approach with the RT core by placing a sphere at the coordinates of the query projection, with a radius set to the distance threshold. This sphere naturally defines the region of interest. Subsequently, rays are cast from codebook entries towards the region of interest in the RT core, enabling rapid determination of whether an entry falls within the region of interest, as the left part of Fig. 8 shows.

In our approach, we utilize a sphere to delineate the region of interest in a query projection. However, this approach necessitates adding spheres into the runtime scene, resulting in excessive overhead. To mitigate this issue, we take advantage of the commutativity of the L2 distance metric. Specifically, we pre-generate spheres at the center of each codebook entry during the offline processing stage, and cast query projection towards these spheres at runtime. This yields identical inside/outside results as the online sphere construction.

Next, we construct the L2-LUT specifically for the entries located within the region of interest. In this process, we also utilize RT cores for the fast and efficient distance calculation of these selected entries, employing the concept of hit distance

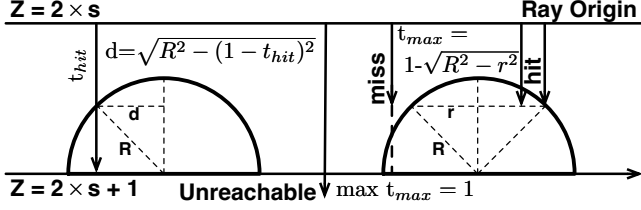


Figure 9: (left) We use the t_{hit} to calculate the hit distance between query projection and sphere centroid. (right) We dynamically adjust the radius of the sphere that a ray can hit by adjusting t_{max} .

which will be detailed in the next paragraph. Once the L2-LUT is constructed, it can be employed for performing distance lookup and accumulation for the points of interest in the search, as demonstrated in the right part of Fig. 8.

Calculating Hit Distance. The construction of L2-LUT requires calculating the distance between each hit entry and the sphere center. The naive approach is to directly use the coordinates of the hit entries stored in global memory, which may have irregular memory accesses due to the sparsity of hit entries. To address this issue, we propose utilizing the results from the RT core to efficiently calculate the distance, thus avoiding global memory accesses.

The RT core incorporates the notion of time in ray tracing, signifying the duration of a ray’s travel. By default, a ray traverses one unit of space within a one unit time interval. There are two defined crucial time points: t_{hit} and t_{max} . The former denotes the time interval from the initiation of ray travel to the point when it intersects an object, while the latter represents the maximum duration that a ray can travel. We can obtain t_{hit} in hit shader without global memory access.

To accurately calculate the distance, we employ the t_{hit} of a ray and rapidly determine the distance between the hit point and the center of the sphere using the radius of the hit sphere. This allows us to obtain the precise distance between the query projection and the codebook entry, as depicted in the left portion of Fig. 9. Notice that R in the figure means the radius of the spheres and are now set as identical constants.

Dynamic Threshold. As mentioned in the previous subsection, our approach also incorporates a dynamic distance threshold to obtain top-100, which can also be adjusted by the user to balance search quality and throughput. A naive implementation would involve modifying the radius of spheres at runtime, which may result in unacceptable overhead.

To eliminate the need for online sphere creation, we convert the dynamic distance threshold into a dynamic maximum travel time for rays. For smaller regions, we set a correspondingly smaller value for t_{max} , as shown on the right side of Fig. 9. For example, if the original threshold is 0.6, a value of $t_{max} = 0.64$ corresponds to a scaling factor of 0.8. Thus, we support both dynamic distance threshold and user-defined scaling factor by adjusting t_{max} only, and we can set the radius of all spheres to be identical, which will significantly simplify the

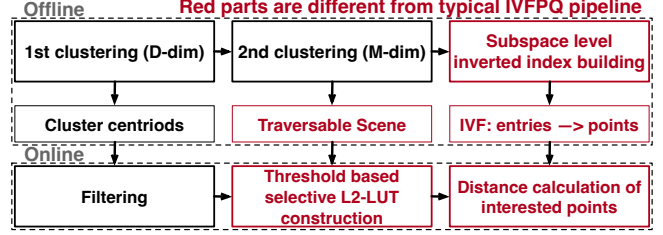


Figure 10: Overview design of JUNO.

forementioned distance calculation and further processing.

Inner Product Similarity Support. Our approach efficiently supports maximal inner product similarity (MIPS) at minimal computational cost through a well-designed transformation mechanism. Previous methods rely on introducing extra dimensions so that they can maximize inner product by minimizing L2 distance between transformed query and search points [5, 62]. The extra dimensions affect training, aligning, search performance [36]. Instead, we propose a method that is free of extra dimensions and RT-core friendly.

Notice that we use t_{hit} to calculate $L2(e, q)$ in a 2D subspace where e is the codebook entry, that is:

$$L2^2(\mathbf{e}, \mathbf{q}) = R^2 - (1 - t_{hit})^2 = (x_e - x_q)^2 + (y_e - y_q)^2$$

$$IP(\mathbf{e}, \mathbf{q}) = x_e x_q + y_e y_q = (\underline{x}_e^2 + \underline{y}_e^2 + x_q^2 + y_q^2 - L2^2(\mathbf{e}, \mathbf{q}))/2$$

$$= (\underline{x}_e^2 + \underline{y}_e^2 + x_q^2 + y_q^2 - R^2 + (1 - t_{hit})^2)/2$$

The codebook entry related part $\underline{x}_e^2, \underline{y}_e^2$ requires accessing global memory lookup for their coordinates at runtime to calculate inner product $IP(e, q)$. While with the RT core, we can eliminate $\underline{x}_e^2, \underline{y}_e^2$ by replacing R with $R' = \sqrt{R^2 + \underline{x}_e^2 + \underline{y}_e^2}$ without any extra dimension(s) as follows:

$$t_{hit}^{new} = 1 - \sqrt{R'^2 - \underline{x}_e^2 - \underline{y}_e^2 - x_q^2 - y_q^2 + 2 \times IP(\mathbf{e}, \mathbf{q})}$$

$$= 1 - \sqrt{R'^2 - x_q^2 - y_q^2 + 2 \times IP(\mathbf{e}, \mathbf{q})}$$

$$\Rightarrow IP(\mathbf{e}, \mathbf{q}) = (x_q^2 + y_q^2 - R'^2 + (1 - t_{hit}^{new})^2)/2$$

where x_q, y_q, R are all constant values for a single query. Thus, the inner product $IP(e, q)$ can be directly calculated using t_{hit}^{new} , without any sphere coordinate accesses. Furthermore, the term $x_q^2 + y_q^2$ can be disregarded as it remains constant for all codebook entries. Consequently, only the radiuses of spheres need to be adjusted from R for the L2 distance metric to $\sqrt{R^2 + \underline{x}_e^2 + \underline{y}_e^2}$ for the inner product metric offline. Notice that we also need to change metric of the cluster in *filtering* from L2 distance to the inner product.

5. JUNO System Design

Based on aforementioned insights, we propose JUNO, an end-to-end search system for efficient ANN search in high-dimensional space. JUNO consists of offline and online phases leveraging the algorithmic enhancement and hardware mapping mentioned in Sec. 4, which are described in Sec. 5.2 and

Sec. 5.3. Besides, we propose pipelining and aggressive approximation leveraging the hardware features of the RT core, which are discussed in Sec. 5.3 and Sec. 5.4.

5.1. System Design Overview

We present the overview of JUNO in Fig. 10, which is also formally described in supplemental materials by algorithm ‘JUNO end to end’. In the offline phase, we prepare the traversable scene and inverted indices from codebook entries to search points in every subspace. Once a batch of queries arrive, we first conduct the filtering that is identical to the original IVFPQ approach. Then we use the RT core to do threshold-based selective L2-LUT construction to obtain only necessary entries falling inside the interested region of query projection in every subspace. Finally, we conduct the distance calculation where only interested search points are considered for the final accumulated distance, by using the inverted indices prepared offline and L2-LUT constructed online.

5.2. Offline Preparation Phase

In the offline phase, we need to prepare a traversable scene and subspace-level inverted indices for later online searching. Alg. 1 shows its details. We first use the typical IVFPQ offline training process. Specifically, we first obtain the C cluster centroids and labels of search points needed by *filtering* (line 3-4), and then generate the codebook trained by the residual between search points and their centroids (line 5-9).

The conventional IVFPQ approach encodes and stores all search points with the codebook entries, which is not efficient for our selective L2-LUT construction. To address this issue, we maintain an inverted index from codebook entries to search points (lines 14-16) in each subspace. For example, $Map[114][19][24]$ contains all search points that satisfy the following conditions: i) the search point belongs to the 114th cluster, and ii) its projection is encoded with the 24th codebook entry in the 19th subspace. This inverted index enables us to only iterate through necessary search points from entries that are close to the query projection in a subspace.

In the d^{th} subspace, the spheres representing the entries of this subspace are positioned at the corresponding x and y coordinates. For the z coordinate, we place the entries from different subspaces at different depths, specifically $z = 2s + 1$ for the s^{th} subspace. This approach prevents interference from rays originating from other subspaces during the ray tracing process. Lines 10-13 show the codebook entries placement in D/M subspaces, which is also illustrated in the left of Fig. 8.

Besides the position of spheres, we also need to determine their radius (i.e., distance threshold). Recall that the distance threshold depends on the density of the grid a query projection falls into (Sec. 4.1), and we use the travel time-based method to enable dynamic threshold. So, we set the same radius for all spheres for the convenience of further computation.

5.3. Online Searching

We describe the case of a single query search in the online phase, which is detailed in Alg. 2. It straightforward to generalize the single-query case to the multi-query case.

Given a query, we first perform the filtering and select n_{probs} clusters (line 2). We then calculate the residuals between the query and centroids of the selected clusters (line 4-5). For dynamic distance threshold, and we calculate the threshold for every query projection with density map, polynomial regressor trained offline and user-defined scaling factor, and transform the threshold to ray’s maximal travel time t_{max} (line 6-7). For each selected cluster, a ray is created from the coordinate of the residual between the query projection and the cluster centroid projection in each subspace (line 8-9). As a reminder, the spheres representing the codebook entries in the s^{th} subspace are placed at $z = 2s + 1$. By placing the ray origin at $z = 2s$ and restricting the t_{max} of the rays to be 1.0, these rays can accurately interact with the spheres in the same subspace without interfering with entries in other subspaces. This is illustrated in the left part of Fig. 8.

Line 10 registers a callback called `RT_HitShader`, which would be invoked when the ray hits the sphere, i.e., an entry falls within the interested region of a query projection in a subspace. Inside the hit callback, we calculate and record the actual distance between the entry and the query projection using the variable t_{hit} (line 14-16). To process multiple queries, we create and shoot rays of all projections in parallel, maintaining for each query a list that records the hit entry IDs and their corresponding distances in every subspace (line 17-18). These lists collectively form the L2 lookup table.

After obtaining the L2-LUT, we perform distance calculations for the search points (formally described in supplemental materials). For each subspace, we access the inverted index

Algorithm 1 Build a traversable scene, prepare cluster centroids of filter and entry-search points mapping offline.

Input: $points[N][D], M = 2, E, C, metric$
Output: $Map[C][\frac{D}{M}]\{entry_id : points_id[\]\}$
Output: $centroids, labels, Scene$

- 1: **function** BUILDRTSCENE($points[N][D], M, E, C$)
- 2: $Scene, filter \leftarrow \emptyset, kmeans(points, n_cluster = C)$
- 3: $centroids, labels \leftarrow filter.centroids, filter.labels$
- 4: $residual \leftarrow [x - centroids[x.label] \text{ for } x \text{ in } points]$
- 5: **for** $s \in [0, \frac{D}{M})$ **do**
- 6: $res \leftarrow residual[:, 2s : 2s + 2]$
- 7: $codebook[s] \leftarrow kmeans(res, n_cluster = E)$
- 8: $entries \leftarrow codebook[s].centroids$
- 9: **for** $e \in [0, E)$ **do**
- 10: $x, y, z \leftarrow entries[e].x, entries[e].y, 2s + 1$
- 11: $Scene.add(sphere(pos = (x, y, z), r = Const))$
- 12: **for** $e \in [0, E), c \in [0, C)$ **do**
- 13: $res_c \leftarrow [labels[p] = c \text{ for } p \text{ in } res]$
- 14: $Map[c][e] \leftarrow [p \text{ encoded by } e \text{ for } p \text{ in } res_c]$
- 15: **return** $Scene, Map, centroids, labels$

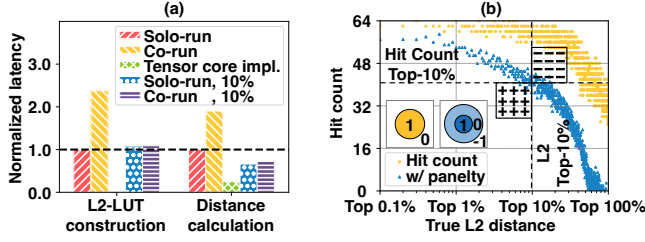


Figure 11: (a) Latency breakdown of three stages. (b) Relationship between hit count and exact distance of query and search points.

to retrieve the search points whose entry is matched by the query projection, and then accumulate their distances with the results in the L2-LUT. The remaining search points would be directly assigned with a large constant without performing any L2-LUT lookup. Finally, a list containing the search points and their distances are returned for selecting the top-k results.

Pipelining on Heterogeneous Cores. In our JUNO framework, the *L2-LUT construction* utilizes the RT core, while the *distance calculation* is performed on the CUDA core. This configuration has the potential for pipelined execution, which can improve the search throughput. This capability is supported by NVIDIA GPUs starting from the Ampere architecture, which allows for co-running of the Tensor core, RT core, and CUDA core [50]. However, naive co-running without proper optimization can result in severe interference and slowdown, as demonstrated in Fig. 11(a). The reason is that long latency of distance calculation on the CUDA core leads to severe resource contention.

JUNO achieves efficient hardware pipelining by mapping the accumulation in the distance calculation stage to Tensor cores. The distances of the selected search points in each subspace are organized into rows to form matrix A with shape $M, K =$

Algorithm 2 Construct L2-LUT with the RT core and conduct distance calculation for interested search points.

Input: $queries[Q][D]$, $index$, $query_select_clusters$, $nprobs$
Input: $density_map$, $poly_regressor$, $thres_scale$ (user defined)
Output: $L2_LUT[Q][nprobs][\frac{D}{M}]\{entry_id : distance\}$

```

1: function L2_LUT( $query$ ,  $index$ ,  $query\_select\_clusters$ )
2:   for  $q \in [0, Q]$ ,  $s \in [0, \frac{D}{M}]$  do
3:      $x, y, z \leftarrow q[q][2s : 2s + 1], 2s$ 
4:     for  $c$  in  $query\_select\_clusters[q]$  do
5:        $thres \leftarrow poly\_regressor(density\_map(x, y))$ 
6:        $t \leftarrow 1 - \sqrt{1.0^2 - (thres \times thres\_scale)^2}$ 
7:        $x, y \leftarrow (x, y) - index.centroids[c][2s : 2s + 1]$ 
8:        $rays.add(x, y, z, t_{max} = t, dir = (0, 0, 1))$ 
9:    $index.scene.set\_hit\_callback(RT\_HitShader)$ 
10:  return RayTracing( $rays, scene$ )
11: function RT_HITSHADER( $index, query\_select\_clusters$ )
12:   $ray, sphere, t_{hit} \leftarrow GetRay(), GetHitSphere(), GetTime()$ 
13:   $q, s, e \leftarrow ray.query\_id, ray.subspace\_id, sphere.entry\_id$ 
14:   $distance \leftarrow \sqrt{R^2 - (1 - t_{hit})^2} // R = Const$ 
15:  for  $c$  in  $query\_select\_clusters$  do
16:     $L2\_LUT[q][c][s].add(\{e : distance\})$ 

```

$Q \times \text{sizeof}(\text{selected points}) \times nprobs, D/M$ (with padding for simplicity). Then, matrix B is created with shape $K, N = D/M, 1$, and all elements are set to 1.0. The accumulation is then performed by calculating the $\text{matmul } A \times B$, utilizing the `cublas` library [53] on Tensor cores. Thus, the reduced latency mitigates the resource contention.

We utilize CUDA MPS [51] to partition the SM resources in a 9:1 ratio, allocating 90% of the resources to *L2-LUT construction* (using RT cores) and 10% to *distance calculation* (using Tensor cores). This partitioning results in similar latencies for the two stages, maximizing the overlap of the two stages. We also apply proper data padding and transformation to enable the pipeline, with an overhead of less than 5% of the latency, compared to the solo-run of Fig. 11(a).

For *filtering*, L2 distance can be calculated with $\|x - q\|_2^2 = x^2 - 2xq^T + q^2$, where $x^2 \leftarrow \sum_{i=0}^{D-1} x_i^2$ is calculated ahead of time. We just calculate the $q^2 \leftarrow \sum_{i=0}^{D-1} q_i^2$ for every query. To calculate $2xq^T$, we use `cublas` and Tensor core too with $\alpha = -2, \beta = 1, A = x, B = q^T, C = x^2(q^2)^T$. Inner product is even simpler by calling a `matmul` with $A = x, B = q^T$.

5.4. Aggressive Approximation: hit count-based Method

The above *L2-LUT construction* still needs to conduct several floating point operations to calculate the radius of a hit sphere and hit distance of a ray. We propose a more aggressive approximated ANN search that uses the hit/miss result from the RT core, inspired by previous work [39].

We first study the relationship between hit count and exact distance, where all spheres are set with a radius threshold that encompasses the top-100 search points. As Fig. 11(b) shows, there is strong correlation between these two factors. The reason is that higher hit count implies being close to the query projection in more subspaces. This finding inspires us to implement a hit count-based ANN search method.

Specifically, we employ a reward/penalty-based model, which involves the extra sphere with half the radius for each original sphere, as in Fig. 11(b). The hit count is incremented by one only when the ray successfully intersects the inner sphere, and decremented by one as a penalty when the ray misses both spheres. As Fig. 11(b) shows, the hit count (blue ▲) calculated using this approach exhibits a stronger correlation compared to the original hit count (yellow ●). While this approximation may result in false positives/negatives (marked as region of ‘-’ and ‘+’ respectively in Fig. 11(b)), it offers users a new parameter to trade-off between reasonable search quality degradation and improved throughput.

6. Evaluation

We demonstrate the effectiveness of the proposed algorithm and hardware mapping of JUNO throughout experiments.

6.1. Experimental Setup

Setup. The ray tracing (RT) part of JUNO is implemented with NVIDIA OptiX 7.6 [48]. We evaluate JUNO on

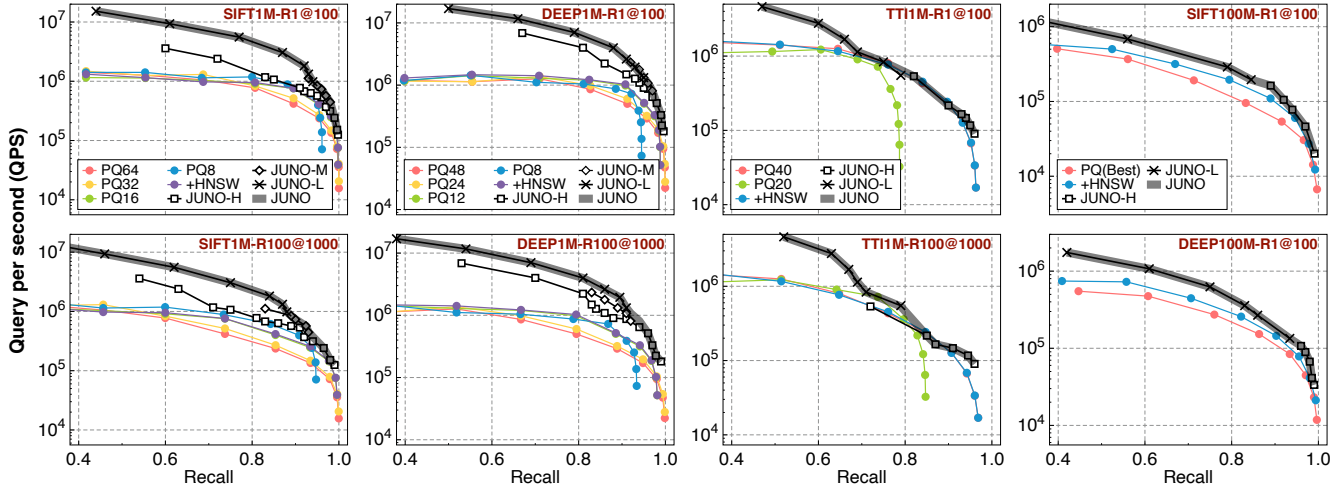


Figure 12: Result of QPS (query per second) and search quality of JUNO on various datasets including SIFT1M, DEEP1M, TTI1M, SIFT100M and DEEP100M. The bolded grey line labeled JUNO is the Pareto frontier of our search engine under different configurations (i.e., the configuration of JUNO-L, JUNO-M, and JUNO-H), standing for the optimal performance of JUNO at a given search quality requirement.

different NVIDIA GPUs (with CUDA/RT cores): RTX 4090 (16384/128), A100 (6912/0), and NVIDIA Tesla A40 (10752/84). Notice that OptiX offloads the RT computation to CUDA cores if the GPU does not have RT cores. This makes JUNO also compatible with the GPU without the RT core. As such, we can conduct the sensitivity analysis to study i) the effectiveness of algorithm enhancement only and ii) the impact of ray tracing hardware performance.

Dataset. Our work focuses on improving the performance of high-dimensional ANN search on a single GPU. We use popular datasets, including SIFT1M/100M [35], DEEP1M/100M [4] and TTI1M[61]. The suffix 1M stands for 1 million search points and the suffix 100M stands for 100 million points. Their embedding sizes are 128/128, 96/96 and 200, respectively. TTI1M uses the inner product metric (MIPS), and the rest uses the L2 distance metric. Notice that larger datasets (>1B) cannot be held in the memory of a single GPU, and they need chunking, partitioning, or other storage-related techniques proposed in other orthogonal works [10, 63, 57]. For example, GGNN utilizes 8 GPUs, and ANNA utilizes 12 accelerators to process 1B datasets [24, 40].

Baseline and Configurations. We mainly compare JUNO against FAISS, the state-of-the-art GPU-accelerated ANN search library [36]. To comprehensively evaluate the advantages of our work, we use identical IVF cluster numbers as those used in FAISS within JUNO. Next, we conduct evaluations on various PQ configurations. Finally, we augment our approach by incorporating the widely applied HNSW (Hierarchical Navigable Small World graphs) optimization [43] on top of the IVF and PQ methods. It is important to note that HNSW represents an orthogonal index optimization technique, compatible with IVF. Since in the search process of HNSW, distance calculation and sorting are still necessary among neighbors of vertices in the navigable small world graphs,

and can still benefit from optimized PQ process proposed in JUNO. While our primary focus in JUNO is on optimizing the PQ component, we defer the integration of HNSW optimization into the current IVF indexing for future research. Nevertheless, we make a thorough comparison of our work against baselines that include HNSW optimizations. In FAISS, this is implemented by calling the `index_factory` with the parameter `IVFx_HNSWy_PQz`. There are other works, such as CPU-centric [28], RAM-centric [10], and disk-centric [10, 63] ANN search optimization, which are orthogonal to JUNO.

Metric. In our evaluation, we assess the search quality using two metrics: Recall-1@100 (R1@100) and Recall-100@1000 (R100@1000). The R1@100 metric is defined as follows. For a set of Q queries, each with 100 retrieved neighbors, R1@100 represents the count of queries, among the set of Q , where their 100 retrieved neighbors include the true nearest neighbor. It is important to note that R1@100 does not consider the specific order of the 100 retrieved neighbors. Considering a scenario with ten queries. If the retrieved neighbors of eight queries include the true nearest neighbor, then the R1@100 score would be calculated as 8/10. The R100@1000 metric measures the averaged number of retrieved neighbors, among a total of 1000 retrieved neighbors, that belong to the 100 true nearest neighbors for each query.

Evaluation Plan. In this study, we conduct an evaluation of Query Per Second (QPS) and search quality for JUNO using various configurations. For each configuration, we apply a scaling factor, as referenced in Section 4, to achieve an optimal balance between search quality and system performance.

- **JUNO-H:** We employ hit time based exact hit distance calculation for high quality requirement.
- **JUNO-M:** We employ finer-grained hit count-based selection with multiple spheres for medium quality requirement.
- **JUNO-L:** We employ hit count-based selection only for low

quality requirement.

It is worth noting that there is some overlap in the search quality among these configurations. And we have empirically divided the search quality requirements into three intervals: $[0.0, 0.95]$, $[0.95, 0.97]$, and $[0.97, 1.0]$. Accordingly, we designate JUNO-L, JUNO-M, and JUNO-H for configurations that correspond to these intervals, respectively. In instances where JUNO-L or JUNO-M fails to meet the requirements of the default intervals, we default to using JUNO-H. We then evaluate the improvement breakdown of two optimization techniques. Finally, we conduct a sensitivity analysis to evaluate the effectiveness of design decisions in JUNO.

6.2. Search Quality and Throughput

Fig. 12 shows the overall search quality and throughput on different datasets. We aggregate the various configurations of JUNO-L/M/H into the grey bold line in the plot as JUNO allows users to make trade-off between search quality and throughput. The lines labelled with $PQ \times$ are from the FAISS baseline stand for dividing the entire space into \times subspaces. The lines labelled with $+HNSW$ stand for the performance of adding HNSW optimization to the best performed PQ configuration. Notice that the HNSW is configured as the best performed parameter.

Justifications of Baseline Configurations. We conduct a detailed analysis on the performance of various baseline methods and provide a rationale for their chosen configurations. Upon investigating the application of HNSW optimization, we observe that it brings limited improvement in smaller datasets (1M) but exhibits more significant enhancements in larger datasets (100M). Moreover, these observations align well with the benchmark results reported in FAISS [45]. Regarding the PQ configuration, we find that incorporating more smaller subspaces results in improved search quality; however, it also leads to lower throughput. We optimize the baseline on different datasets so that the baseline performance approaches closer to the Pareto optimal. Despite the varying performances of the baseline methods, in our subsequent analysis, we will compare JUNO against the best-performing one to ensure a fair and rigorous evaluation.

SIFT1M and DEEP1M We achieve $7.8\times$ higher QPS than the baseline in the low search quality scenario, i.e., JUNO-L with $R1@100 \leq 0.95$. The improvements are from the exploiting of sparsity and aggressive approximation. Firstly, by setting a tighter threshold, we can discard more search points, thereby improving the filtering effectiveness. Secondly, we perform hit count-based selection without the actual distance calculation. Additionally, using a tighter threshold leads to fewer hit events due to the reduced size of the spheres. As a result, JUNO is able to fully exploit the sparsity and spatial locality of the data, which is particularly advantageous when dealing with low search quality requirements.

The high requirement for search quality, such as $R1@100=0.99$, results in increased selected search points.

Additionally, it is necessary to compute the actual distance to meet this stringent accuracy requirement. These two factors limit the advantages of leveraging sparsity. However, despite these limitations, the JUNO approach still achieves a throughput improvement of $2.4\times$ compared to the baseline, owing to the tree-based search methodology employed in the RT core.

It is worth noting that JUNO-L only achieves 0.95 recall for these two datasets as it employs a pure hit count-based approximation approach. JUNO-M is able to improve the search quality to 0.97 by employing the reward/penalty-based approximation approach with extra inner spheres (Sec. 5.4). The throughput improvement over the baseline is $2.9\times$.

TTI1M This dataset uses the inner product metric (i.e., MIPS). Since JUNO-H calculates the exact distance in every subspace, its performance improvement is $2.04\times$, similar to the two previous datasets with L2 metric. Notice that the FAISS baseline can also only reach the 0.96 recall. While the hit count-based method abandons the t_{hit} information, and intersecting only implies being close in aspect of L2 distance, rather than being similar in aspect of inner product. So, the search quality rapidly drops in dataset with inner product similarity. Thus, the line representing JUNO-L moves to left.

SIFT100M and DEEP100M The JUNO-H and JUNO-L configurations achieve an averaged improvement of $1.5\times$ and $2.1\times$ over the baseline, respectively. The performance gain of JUNO-H diminishes as the *distance calculation* becomes the bottleneck. Notice that these improvements are calculated by comparing JUNO **without** HNSW against FAISS **with** HNSW optimization. We do not implement HNSW in JUNO owing to its high code complexity in current FAISS framework. Meanwhile, implementing it provides no extra insight to guide the optimization of PQ. Notably, JUNO-H outperforms the baseline without HNSW optimization by a factor of $3.0\times$.

Results of Different Metrics. The results for $R100@1000$ of SIFT1M, DEEP1M, and TTI1M are presented in Figure 12. It is evident that JUNO exhibits comparable improvements over the baseline, demonstrating the effectiveness of the approximation techniques proposed by JUNO even under more challenging metrics. Specifically, on average, the top 100 retrieved neighbors out of 1000 contain 65% of the true top 100 nearest neighbors. This performance is also influenced by the clustering quality of IVF and the PQ quality, which remain identical in JUNO when compared to the baseline.

6.3. Effectiveness of Different Optimizations

We now evaluate the effectiveness of optimizations used in JUNO, including the pipelining among CUDA-tensor-RT cores, hit count-based L2-LUT selective construction, and dynamic radius (i.e., distance threshold). Fig. 13(a) shows the both the overall improvement and effects without applying the first two optimizations, while (b) shows the effects of dynamic radius.

Overall Improvement. JUNO achieves averaged $2.1\times$ - $4.4\times$ QPS improvement on five datasets from high to low search

quality requirement over the baseline. Although we do not tune JUNO for a particular dataset, JUNO still achieves $8.5 \times - 3.2 \times$ maximum improvement on these datasets.

Effectiveness of Pipelining. The third bar in Fig. 13(a) presents the performance improvement achieved without pipelining. Specifically, in scenarios where high search quality is required, the critical path is the *L2-LUT construction*, which has a longer latency compared to the distance calculation. Consequently, without pipelining, the improvement decreases by 44%. On the other hand, in situations where lower search quality is acceptable, the latency between the *L2-LUT construction* and *distance calculation* is similar, resulting in a 50% decrease in the improvement without pipelining.

Effectiveness of Hit Count-based Selection. The final bar in Fig. 13(a) displays the improvement achieved when hit count-based selection is not utilized. It is evident that for extremely high search quality requirements, hit count-based selection has no influence as it is unable to achieve such high quality. However, as the search quality requirement decreases, the influence of hit count-based selection increases, as lower recall rates necessitate fewer exact distance calculations. In conclusion, a combination of hit count-based selection and exact distance calculation is necessary to deliver high search throughput across different search quality requirements.

Effectiveness of Dynamic Threshold Strategy. We evaluate the search quality and search throughput (QPS) using a small and large static threshold. We conduct this evaluation on SIFT1M dataset with JUNO-H, where small and large static thresholds are determined with the minimum and maximal threshold values of dynamic threshold. As Fig. 13(b) shows, utilizing a large static threshold results in a decreased search throughput but achieves higher search quality. This can be attributed to the fact that a larger threshold leads to a larger sphere, causing a ray to intersect with more spheres and trigger more hit shader functions. Conversely, using a small static threshold improves the search throughput, but at the cost of degraded search quality. Even for scenarios with low search quality requirements, selecting more clusters to address the recall issue becomes necessary, thereby negating the potential performance gain from reducing hit shader invocations. For cases with high search quality requirements, this small static

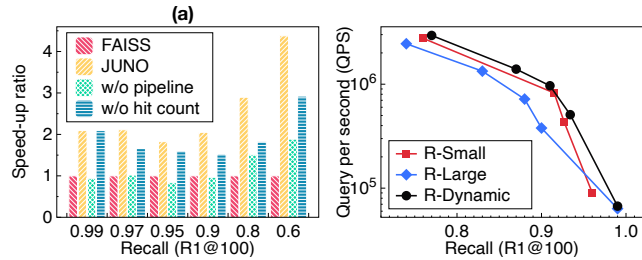


Figure 13: (a) Improvement breakdown of JUNO against FAISS. (b) Performance of different threshold strategy, evaluated on A40.

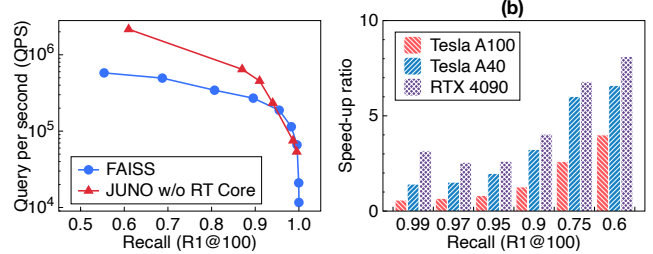


Figure 14: (a) QPS and recall of JUNO and FAISS on A100. (b) Average advantage against FAISS on different GPUs.

threshold proves to be inadequate, as it results in missing too many true top-k neighbors, thereby jeopardizing recall. In contrast, our dynamic threshold strategy outperforms both the small and large static threshold-based approaches in terms of both search quality and throughput.

6.4. Sensitivity to RT Core Performance

Finally, we evaluate the performance of JUNO with and without the acceleration of the RT core using different GPUs. Fig. 14(a) shows the detailed performance of JUNO and the baseline on Tesla A100, a GPU without the RT core. We conduct this evaluation on SIFT1M dataset. The baseline is configured with $PQ_{16}+HNSW$ (the best performed configuration among others). We find that JUNO still gains significant improvement with low search quality requirement on Tesla A100, which means the advantage is purely contributed by the threshold-based selective algorithm. This result also verifies that it is reasonable to leverage the sparsity and spatial similarity in the typical *IVFPQ* process. For high search quality requirements, JUNO gradually performs worse than the baseline since the overhead to simulate ray tracing with the CUDA core suppresses the tiny positive effect brought by the sparsity.

Results in Fig. 11(a) imply that the performance of JUNO is bound by the performance of RT cores. So, we expect a performance improvement with a faster RT core. According to the white paper of NVIDIA Ada architecture [54], the Gen.3 RT core of Ada GPUs has $2 \times$ throughput compared to a Gen.2 RT core of Ampere GPUs. As shown in Fig. 14(b), on average of three 1M datasets, RTX4090 has $1.5 \times$ higher improvement over the baseline than Tesla A40. Notice that the throughput of CUDA cores and Tensor cores of RTX4090 is $1.4 \times$ of A40 GPU per SM [50, 54].

6.5. Robustness Discussions

Since JUNO adopts an approximate search methodology, we delve deeper into its robustness aspect, focusing on concerns that researchers may have pertaining to accuracy guarantees and dataset characteristics.

Accuracy Guarantee. Similar to our baseline ANN search framework, JUNO does not offer a theoretical accuracy guarantee and relies on the quality of offline clustering. Nonetheless, it has the capability to support lossless searching through a

series of minor adjustments: i) exhaustively searching all `IVF` clusters, equivalent to a `Flat` index, ensures that all search points are thoroughly explored, leaving no potential neighbors missed. ii) projecting all original search points (rather than PQ codebook entries) into a two-dimensional space, iii) and employing ray tracing to calculate precise distances between search points (rather than PQ codebook entries) and queries. As a result, we can guarantee the search accuracy and use JUNO in scenarios with stringent search quality requirements.

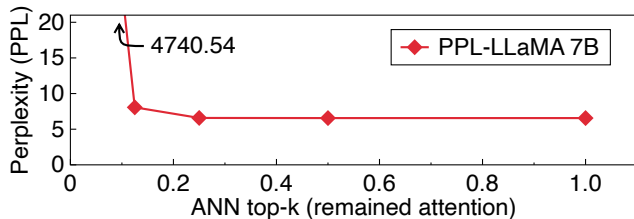


Figure 15: Word perplexity of Llama-7B with different amount of attention remained.

Dataset Characteristics. JUNO leverages sparsity and locality inherent in high-dimensional ANN search to improve efficiency. Critically, JUNO does NO dataset-specific tuning and optimization. As we have shown in the above results, JUNO consistently delivers significant performance improvements across a number of popular datasets, despite their different levels of sparsity and redundancy characteristics. Furthermore, in the era of training expansive foundational models using vast datasets, the data exhibit an escalating trend of becoming increasingly sparse and redundant [19, 18, 12, 31, 25]. Those foundational models adopt the decoder-based Transformer architecture that adopts the multi-head attention mechanism. The memory and computation complexity of attention mechanism scales quadratically with the sequence length, making it the main bottleneck when dealing with long input sequences. On the other hand, the attention mechanism essentially calculates the inner product between query vectors and key vectors, which correspond to the query points and search points in vector search. Prior works have shown that keeping the most significant attention values can still preserve the model accuracy [37, 8], making LLM an ideal candidate to be accelerated by ANN search. This phenomenon underscores the promising potential of JUNO in future.

To verify, we conduct an experiment on Llama-7B [64]. Recall that 30%-50% nearest entries should be preserved to maintain search qualities in ANN search according to Fig. 4 and Fig. 5, while as shown in Fig. 15, a commendable quality can be sustained with less than 20% nearest tokens attended.

7. Related Work

Since there are two aspects, i.e., indexing and encoding, in an ANN algorithm, we compare JUNO with existing works in these two algorithmic aspects and hardware acceleration.

Indexing. Indexing techniques aim to compress the search space by eliminating unnecessary search points. Except for the trivial flat index that store the complete database, one example of such technique is the inverted file index (`IVF`) [42], which organizes search points into clusters and selects several closest clusters when searching (search points in other unselected clusters are ignored). Another type of technique involves graph-based methods that construct a nearest neighbor graph to quickly prune the search space [16, 30, 67]. To accelerate these methods, heuristic-based approaches [6, 20, 32] have been proposed. Currently, among these works, hierarchical navigable small world (`HNSW`) [43] and navigating spread-out graph (`NSG`) [21] are the most representative. The `HNSW` construct the neighbor graph hierarchically, with the search going deeper, the graph have higher degree and shorter edges. So that the search can quickly choose a good search direction and get good enough results in log-scale. The `NSG` further reduce the size and degree of graphs, thus reduce the length of search path via setting navigation points from which the search begins. Additionally, tree-based techniques including kd-tree and octree [46, 9], and locality-sensitive hashing (`LSH`) [14, 13] are also commonly used for indexing. Noted that JUNO is not limited to specific indexing methods, and is compatible with `Flat`, `IVF`, `HNSW`, etc.

Encoding. Encoding techniques aim to reduce the memory consumption of search points. The most commonly used technique is product quantization (`PQ`) [34], which splits the space into several subspaces and encodes the search point projections. Several techniques have been proposed to optimize the codebook quality thanks to its strong relation between search quality, such as `DPQ` [38] and `OPQ` [23]. Additionally, scalar quantization (`SQ`) [74] maps vector components separately and linearly, similar to traditional quantization in DNNs [27, 26]. Additive quantization (`AQ`) [3] encodes search points as a sum of codebook entries. JUNO currently supports product quantization (`PQ`) only.

Hardware Acceleration. There are several specialized architecture designs including hardware support for hierarchical product quantization [1] and fused high-performance k-selection [73]. ANNA proposed an end-to-end hardware solution for PQ-based ANN search [40]. Several architectural designs based on tree-like data structures [72, 9] are proposed for low-dimensional ANN search. In addition to computation, large-scale ANN search presents severe challenges to the memory and storage subsystems. DiskANN presented a graph-based indexing that can search with limited RAM and cheap solid-state drives [33]. SPANN presented a memory-disk hybrid indexing following the inverted file index [10]. Those solutions employ partitioning techniques to divide large datasets into smaller chunks for index creation and processing, in order to overcome resource limitations. Particularly in the case of GPU-accelerated ANNS, the GPU’s memory capacity is often inadequate to load the entire dataset. As such, we

believe that evaluating a 100M dataset would be adequate to showcase the superior performance of JUNO.

Researchers have also proposed various techniques to map the ANN search to existing hardware. ScANN optimizes the ANN search with the AVX ISA on CPUs [28]. RTNN maps the low-dimensional ANN search to the RT core [75]. Our system seeks to harness the power of the RT core to accelerate more general ANN search in high-dimensional spaces.

8. Conclusion

In this work, we have presented JUNO, an end-to-end approximate nearest neighbor (ANN) search system that incorporates a sparsity-aware codebook entry selection algorithm and a highly efficient RT core mapping. The key in our algorithm is to exploit the opportunities of sparsity and spatial locality that we have identified through detailed profiling. Specifically, we employ a distance threshold filtering that can be efficiently mapped to RT cores. Additionally, we optimize the system with time-based hit distance calculation, hit count-based aggressive approximation, and Tensor-RT core pipelining. Evaluation of JUNO on multiples datasets demonstrates a $2.1 \times - 8.5 \times$ improvement over existing product quantization based ANN search in search throughput across various scenarios.

Acknowledgement

This work was supported by the National Key R&D Program of China under Grant 2022YFB4501401, the National Natural Science Foundation of China (NSFC) grant (62222210, and 62072297, and U21B2017). We would like to thank the anonymous reviewers for their constructive feedback and comments to improve this work. We also thank our shepherd, Professor Sasa Misailovic, for his support and guidance in preparing this paper for publication. We also thank Yue Guan, Yangjie Zhou, Weihao Cui, Sean Yao, Tinghan Qian, Yingzhe Lyu, and many other colleagues for beneficial discussion and useful comments. Special thanks to Vega Jiang for continuous help and support. Any opinions, findings, and conclusions in this paper are those of the authors only and do not necessarily reflect the views of our sponsors.

References

- [1] Ameer MS Abdelhadi, Christos-Savvas Bouganis, and George A Constantinides. Accelerated approximate nearest neighbors search through hierarchical product quantization. In *2019 International Conference on Field-Programmable Technology (ICFPT)*, pages 90–98. IEEE, 2019.
- [2] Oron Ashual, Shelly Sheynin, Adam Polyak, Uriel Singer, Oran Gafni, Eliya Nachmani, and Yaniv Taigman. Knn-diffusion: Image generation via large-scale retrieval. *CoRR*, abs/2204.02849, 2022.
- [3] Artem Babenko and Victor S. Lempitsky. Additive quantization for extreme vector compression. In *2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*, pages 931–938. IEEE Computer Society, 2014.
- [4] Artem Babenko and Victor S. Lempitsky. Efficient indexing of billion-scale datasets of deep descriptors. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2055–2063. IEEE Computer Society, 2016.
- [5] Yoram Bachrach, Yehuda Finkelstein, Ran Gilad-Bachrach, Liran Katzir, Noam Koenigstein, Nir Nice, and Ulrich Paquet. Speeding up the xbox recommender system using a euclidean transformation for inner-product spaces. In Alfred Kobsa, Michelle X. Zhou, Martin Ester, and Yehuda Koren, editors, *Eighth ACM Conference on Recommender Systems, RecSys '14, Foster City, Silicon Valley, CA, USA - October 06 - 10, 2014*, pages 257–264. ACM, 2014.
- [6] Dmitry Baranchuk, Dmitry Persiyonov, Anton Sinitin, and Artem Babenko. Learning to route in similarity graphs. In *International Conference on Machine Learning*, pages 475–484. PMLR, 2019.
- [7] Jeffrey S. Beis and David G. Lowe. Shape indexing using approximate nearest-neighbour search in high-dimensional spaces. In *1997 Conference on Computer Vision and Pattern Recognition (CVPR '97), June 17-19, 1997, San Juan, Puerto Rico*, pages 1000–1006. IEEE Computer Society, 1997.
- [8] Amanda Bertsch, Uri Alon, Graham Neubig, and Matthew R. Gormley. Unlimiformer: Long-range transformers with unlimited length input. *CoRR*, abs/2305.01625, 2023.
- [9] Faquan Chen, Rendong Ying, Jianwei Xue, Fei Wen, and Peilin Liu. Parallelnn: A parallel octree-based nearest neighbor search accelerator for 3d point clouds. In *IEEE International Symposium on High-Performance Computer Architecture, HPCA 2023, Montreal, QC, Canada, February 25 - March 1, 2023*, pages 403–414. IEEE, 2023.
- [10] Qi Chen, Bing Zhao, Haidong Wang, Mingqin Li, Chuanjie Liu, Zengzhong Li, Mao Yang, and Jingdong Wang. Spann: Highly-efficient billion-scale approximate nearest neighborhood search. *Advances in Neural Information Processing Systems*, 34:5199–5212, 2021.
- [11] Rihan Chen, Bin Liu, Han Zhu, Yaouxuan Wang, Qi Li, Buting Ma, Qingbo Hua, Jun Jiang, Yunlong Xu, Hongbo Deng, and Bo Zheng. Approximate nearest neighbor search under neural similarity metric for large-scale recommendation. In Mohammad Al Hasan and Li Xiong, editors, *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17-21, 2022*, pages 3013–3022. ACM, 2022.
- [12] Hongrong Cheng, Miao Zhang, and Javen Qinfeng Shi. A survey on deep neural network pruning-taxonomy, comparison, analysis, and recommendations. *CoRR*, abs/2308.06767, 2023.
- [13] Anirban Dasgupta, Ravi Kumar, and Tamás Sarlós. Fast locality-sensitive hashing. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1073–1081, 2011.
- [14] Mayur Datar, Nicole Immorlica, Piotr Indyk, and Vahab S. Mirrokni. Locality-sensitive hashing scheme based on p-stable distributions. In *Proceedings of the twentieth annual symposium on Computational geometry*, pages 253–262, 2004.
- [15] Min Dong, Zhe Wang, Chenghui Dong, Xiaomin Mu, and Yide Ma. Classification of region of interest in mammograms using dual contourlet transform and improved KNN. *J. Sensors*, 2017:3213680:1–3213680:15, 2017.
- [16] 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*, pages 577–586, 2011.
- [17] J. Elseberg, S. Magnenat, R. Siegart, and A. Nüchter. Comparison of nearest-neighbor-search strategies and implementations for efficient shape registration. *Journal of Software Engineering for Robotics (JOSER)*, 3(1):2–12, 2012.
- [18] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.

- [19] Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 10323–10337. PMLR, 2023.
- [20] Cong Fu, Chao Xiang, Changxu Wang, and Deng Cai. Fast approximate nearest neighbor search with the navigating spreading-out graph. *arXiv preprint arXiv:1707.00143*, 2017.
- [21] Cong Fu, Chao Xiang, Changxu Wang, and Deng Cai. Fast approximate nearest neighbor search with the navigating spreading-out graph. *Proc. VLDB Endow.*, 12(5):461–474, jan 2019.
- [22] Jianyang Gao and Cheng Long. High-dimensional approximate nearest neighbor search: with reliable and efficient distance comparison operations. *CoRR*, abs/2303.09855, 2023.
- [23] Tiezheng Ge, Kaiming He, Qifa Ke, and Jian Sun. Optimized product quantization. *IEEE transactions on pattern analysis and machine intelligence*, 36(4):744–755, 2013.
- [24] Fabian Groh, Lukas Ruppert, Patrick Wieschollek, and Hendrik P. A. Lensch. GGNN: graph-based GPU nearest neighbor search. *IEEE Trans. Big Data*, 9(1):267–279, 2023.
- [25] Cong Guo, Bo Yang Hsueh, Jingwen Leng, Yuxian Qiu, Yue Guan, Zehuan Wang, Xiaoying Jia, Xipeng Li, Minyi Guo, and Yuhao Zhu. Accelerating sparse DNN models without hardware-support via tile-wise sparsity. In Christine Cui, Irene Qualters, and William T. Kramer, editors, *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC 2020, Virtual Event / Atlanta, Georgia, USA, November 9-19, 2020*, page 16. IEEE/ACM, 2020.
- [26] Cong Guo, Jiaming Tang, Weiming Hu, Jingwen Leng, Chen Zhang, Fan Yang, Yunxin Liu, Minyi Guo, and Yuhao Zhu. Olive: Accelerating large language models via hardware-friendly outlier-victim pair quantization. In Yan Solihin and Mark A. Heinrich, editors, *Proceedings of the 50th Annual International Symposium on Computer Architecture, ISCA 2023, Orlando, FL, USA, June 17-21, 2023*, pages 3:1–3:15. ACM, 2023.
- [27] Cong Guo, Chen Zhang, Jingwen Leng, Zihan Liu, Fan Yang, Yunxin Liu, Minyi Guo, and Yuhao Zhu. ANT: exploiting adaptive numerical data type for low-bit deep neural network quantization. In *55th IEEE/ACM International Symposium on Microarchitecture, MICRO 2022, Chicago, IL, USA, October 1-5, 2022*, pages 1414–1433. IEEE, 2022.
- [28] Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Kumar. Accelerating large-scale inference with anisotropic vector quantization. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 3887–3896. PMLR, 2020.
- [29] Eric Haines and Tomas Akenine-Möller, editors. *Ray Tracing Gems*. Apress, 2019. <http://raytracinggems.com>.
- [30] Kiana Hajebi, Yasin Abbasi-Yadkori, Hossein Shahbazi, and Hong Zhang. Fast approximate nearest-neighbor search with k-nearest neighbor graph. In *Twenty-Second International Joint Conference on Artificial Intelligence*, 2011.
- [31] Torsten Hoefler, Dan Alistarh, Tal Ben-Nun, Nikoli Dryden, and Alexandra Peste. Sparsity in deep learning: Pruning and growth for efficient inference and training in neural networks. *J. Mach. Learn. Res.*, 22:241:1–241:124, 2021.
- [32] 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.
- [33] 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.
- [34] 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, 2010.
- [35] Hervé Jégou, Romain Tavenard, Matthijs Douze, and Laurent Amsaleg. Searching in one billion vectors: Re-rank with source coding. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP 2011, May 22-27, 2011, Prague Congress Center, Prague, Czech Republic*, pages 861–864. IEEE, 2011.
- [36] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547, 2019.
- [37] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- [38] Benjamin Klein and Lior Wolf. End-to-end supervised product quantization for image search and retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5041–5050, 2019.
- [39] Jon M. Kleinberg. Two algorithms for nearest-neighbor search in high dimensions. In Frank Thomson Leighton and Peter W. Shor, editors, *Proceedings of the Twenty-Ninth Annual ACM Symposium on the Theory of Computing, El Paso, Texas, USA, May 4-6, 1997*, pages 599–608. ACM, 1997.
- [40] Yejin Lee, Hyunji Choi, Sunhong Min, Hyunseung Lee, Sangwon Beak, Dawoon Jeong, Jae W. Lee, and Tae Jun Ham. ANNA: specialized architecture for approximate nearest neighbor search. In *IEEE International Symposium on High-Performance Computer Architecture, HPCA 2022, Seoul, South Korea, April 2-6, 2022*, pages 169–183. IEEE, 2022.
- [41] Conglong Li, Minjia Zhang, David G. Andersen, and Yuxiong He. Improving approximate nearest neighbor search through learned adaptive early termination. In David Maier, Rachel Pottinger, AnHai Doan, Wang-Chiew Tan, Abdussalam Alawini, and Hung Q. Ngo, editors, *Proceedings of the 2020 International Conference on Management of Data, SIGMOD Conference 2020, online conference [Portland, OR, USA], June 14-19, 2020*, pages 2539–2554. ACM, 2020.
- [42] Yuchen Liu, Zhibin Pan, Liangzhuang Wang, and Yang Wang. A new fast inverted file-based algorithm for approximate nearest neighbor search without accuracy reduction. *Inf. Sci.*, 608:613–629, 2022.
- [43] 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.
- [44] Adam Marrs, Peter Shirley, , and Ingo Wald, editors. *Ray Tracing Gems II*. Apress, 2021. <http://raytracinggems.com/rtg2>.
- [45] Meta. Indexing 1g vectors. <https://github.com/facebookresearch/faiss/wiki/Indexing-1G-vectors>, 2023.
- [46] Marius Muja and David G Lowe. Scalable nearest neighbor algorithms for high dimensional data. *IEEE transactions on pattern analysis and machine intelligence*, 36(11):2227–2240, 2014.
- [47] Maxim Naumov, Dheevatsa Mudigere, Hao-Jun Michael Shi, Jianyu Huang, Narayanan Sundaraman, Jongsoo Park, Xiaodong Wang, Udit Gupta, Carole-Jean Wu, Alisson G. Azzolini, Dmytro Dzhulgakov, Andrey Malleevich, Iliia Cherniavskii, Yinghai Lu, Raghuraman Krishnamoorthi, Ansha Yu, Volodymyr Kondratenko, Stephanie Pereira, Xianjie Chen, Wenlin Chen, Vijay Rao, Bill Jia, Liang Xiong, and Misha Smelyanskiy. Deep learning recommendation model for personalization and recommendation systems. *CoRR*, abs/1906.00091, 2019.
- [48] NVIDIA. Nvidia optix™ ray tracing engine. <https://developer.nvidia.com/rtx/ray-tracing/optix>.
- [49] NVIDIA. Nvidia turing gpu architecture. <https://images.nvidia.com/aem-dam/en-zz/Solutions/design-visualization/technologies/turing-architecture/NVIDIA-Turing-Architecture-Whitepaper.pdf>, 2018.
- [50] NVIDIA. Nvidia ampere ga102 gpu architecture. <https://www.nvidia.com/content/PDF/nvidia-ampere-ga-102-gpu-architecture-whitepaper-v2.pdf>, 2021.
- [51] NVIDIA. Multi-process service. https://docs.nvidia.com/deploy/pdf/CUDA_Multi_Process_Service_Overview.pdf, 2022.
- [52] NVIDIA. Nvidia ada craft the engineering marvel of the rtx 4090. <https://images.nvidia.com/aem-dam/Solutions/geforce/ada/ada-lovelace-architecture/nvidia-ada-gpu-craft.pdf>, 2022.
- [53] NVIDIA. Basic linear algebra on nvidia gpus. <https://developer.nvidia.com/cublas>, 2023.
- [54] NVIDIA. Nvidia ada gpu architecture. <https://images.nvidia.com/aem-dam/Solutions/Data-Center/14/nvidia-ada-gpu-architecture-whitepaper-v2.0.pdf>, 2023.
- [55] NVIDIA. Nvidia tensor cores unprecedented acceleration for hpc and ai. <https://www.nvidia.com/en-us/data-center/tensor-cores/>, 2023.
- [56] OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023.

- [57] Zhen Peng, Minjia Zhang, Kai Li, Ruoming Jin, and Bin Ren. iqa: Fast and accurate vector search with efficient intra-query parallelism on multi-core architectures. In Maryam Mehri Dehnavi, Milind Kulkarni, and Sriram Krishnamoorthy, editors, *Proceedings of the 28th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming, PPOPP 2023, Montreal, QC, Canada, 25 February 2023 - 1 March 2023*, pages 313–328. ACM, 2023.
- [58] Sakti Pramanik and Jinhua Li. Fast approximate search algorithm for nearest neighbor queries in high dimensions. In Masaru Kitsuregawa, Michael P. Papazoglou, and Calton Pu, editors, *Proceedings of the 15th International Conference on Data Engineering, Sydney, Australia, March 23-26, 1999*, page 251. IEEE Computer Society, 1999.
- [59] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 77–85. IEEE Computer Society, 2017.
- [60] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- [61] Yandex Research. Benchmarks for billion-scale similarity search. <https://research.yandex.com/blog/benchmarks-for-billion-scale-similarity-search>, 2021.
- [62] Anshumali Shrivastava and Ping Li. Asymmetric LSH (ALSH) for sublinear time maximum inner product search (MIPS). In Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2321–2329, 2014.
- [63] Harsha Vardhan Simhadri, Ravishankar Krishnaswamy, Gopal Srinivasa, Suhas Jayaram Subramanya, Andrija Antonijevic, Dax Pryce, David Kaczynski, Shane Williams, Siddarth Gollapudi, Varun Sivashankar, Neel Karia, Aditi Singh, Shikhar Jaiswal, Neelam Mahapatro, Philip Adams, Bryan Tower, and Yash Patel.
- [64] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971, 2023.
- [65] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008, 2017.
- [66] Ingo Wald and Vlastimil Havran. On building fast kd-trees for ray tracing, and on doing that in $o(n \log n)$. In *2006 IEEE Symposium on Interactive Ray Tracing*, pages 61–69, 2006.
- [67] Jing Wang, Jingdong Wang, Gang Zeng, Zhuowen Tu, Rui Gan, and Shipeng Li. Scalable k-nn graph construction for visual descriptors. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1106–1113. IEEE, 2012.
- [68] Liwei Wang, Yan Zhang, and Jufu Feng. On the euclidean distance of images. *IEEE Trans. Pattern Anal. Mach. Intell.*, 27(8):1334–1339, 2005.
- [69] Mengzhao Wang, Xiaoliang Xu, Qiang Yue, and Yuxiang Wang. A comprehensive survey and experimental comparison of graph-based approximate nearest neighbor search. *Proc. VLDB Endow.*, 14(11):1964–1978, 2021.
- [70] Ruoxi Wang, Rakesh Shivanna, Derek Zhiyuan Cheng, Sagar Jain, Dong Lin, Lichan Hong, and Ed H. Chi. DCN V2: improved deep & cross network and practical lessons for web-scale learning to rank systems. In Jure Leskovec, Marko Grobelnik, Marc Najork, Jie Tang, and Leila Zia, editors, *WWW '21: The Web Conference 2021, Virtual Event / Ljubljana, Slovenia, April 19-23, 2021*, pages 1785–1797. ACM / IW3C2, 2021.
- [71] Wenxuan Wu, Zhongang Qi, and Fuxin Li. Pointconv: Deep convolutional networks on 3d point clouds. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 9621–9630. Computer Vision Foundation / IEEE, 2019.
- [72] Tiancheng Xu, Boyuan Tian, and Yuhao Zhu. Tigris: Architecture and algorithms for 3d perception in point clouds. In *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture, MICRO 2019, Columbus, OH, USA, October 12-16, 2019*, pages 629–642. ACM, 2019.
- [73] Jialiang Zhang, Soroosh Khoram, and Jing Li. Efficient large-scale approximate nearest neighbor search on opencl fpga. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4924–4932, 2018.
- [74] Wengang Zhou, Yijuan Lu, Houqiang Li, and Qi Tian. Scalar quantization for large scale image search. In Noboru Babaguchi, Kiyoharu Aizawa, John R. Smith, Shin'ichi Satoh, Thomas Plagemann, Xian-Sheng Hua, and Rong Yan, editors, *Proceedings of the 20th ACM Multimedia Conference, MM '12, Nara, Japan, October 29 - November 02, 2012*, pages 169–178. ACM, 2012.
- [75] Yuhao Zhu. RTNN: accelerating neighbor search using hardware ray tracing. In Jaejin Lee, Kunal Agrawal, and Michael F. Spear, editors, *PPOPP '22: 27th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, Seoul, Republic of Korea, April 2 - 6, 2022*, pages 76–89. ACM, 2022.