

490 A Related Work

491 A.1 ANNs and MIPS

492 ANNS achieves highly efficient vector search by allowing a small number of errors. Generally,
493 there are two kinds of ANNS algorithms: *non-exhaustive ANNS methods* [37, 34, 36, 28] and *vector*
494 *compression methods* [12, 29, 3, 18]. Specifically, *Non-exhaustive ANNS methods* do not compress
495 the index. They reduce the number of candidates for each query to speed up the retrieval process.
496 Popular algorithms include tree search [37, 11] and graph search [34, 36, 28]. *Vector compression*
497 *methods* mainly aim to compress the index to accelerate retrieval. Popular algorithms include hashing
498 [12, 38, 41] and quantization [29, 3, 18, 17].

499 Under the constraints of storage, compressed methods are widely investigated by researchers. Prod-
500 uct quantization [22, 6] decomposes the vector representation space into the Cartesian product of
501 subspaces. Optimized product quantization (OPQ) [16] jointly learns space decomposition and
502 subspace quantization. Multi-scale Quantization [46] includes a multi-scale framework so that it can
503 learn a separate scalar quantizer. Composite Quantization [50] and Additive Quantization [3] do not
504 decompose space but directly learn multiple codebooks. There are also some algorithms that take
505 query information into account. NEQ [10] decomposes the quantization error into norm error and
506 direction error and improves existing VQ techniques for MIPS. ScaNN [18] computes the weight for
507 each pair of vectors. Different from NEQ and ScaNN, KDindex utilizes query and corresponding
508 top-k candidates. BLISS [19] regards ground truth as labels. However, the ground truth is difficult to
509 obtain in huge quantities of databases. Interested readers could refer to the surveys [43, 31].

510 A.2 Knowledge Distillation

511 *Knowledge Distilling* (KD) was first proposed in [20], in which a complex neural network was firstly
512 trained and then transferred to a small model. Following this, DarkRank [7] proposed a method
513 combining deep metric learning and *Learning to rank* technique with KD to solve image retrieval
514 and image clustering tasks. In addition, a few recent methods [30, 39] have adopted knowledge
515 distillation to RS. RD [42] firstly proposes a KD method that makes the student give high scores on
516 the top-ranked items of the teacher’s recommendation list. Similarly, CD [30] makes the student
517 imitate the teacher’s prediction scores with particular emphasis on the items ranked highly by the
518 teacher. The most recent work, RRD [26], formulates the distillation process as a relaxed ranking
519 matching problem between the ranking list of the teacher and that of the student. However, there are
520 limited works focusing on index building under knowledge distillation.

521 In the context of quantization problems under distillation, the most relevant work is Distill-VQ [47],
522 which uses knowledge distillation for ranking candidates in web search tasks. This method applies
523 the sampling technique to rank a sample of the document from all data each time. But this technique
524 is not applicable to training a ranking model when documents and queries are represented with no
525 content information. In this case, the labeled model training cannot be easily generalized to all
526 documents and queries. In contrast, KDindex relaxes the requirement on labeled data and can be
527 trained purely with unlabeled data.

528 B More Details of Experimental Settings

529 B.1 Dataset Statistics

Table 5: Statistics of the datasets.

| Datasets | #Database | #Train | #Test | Dim |
|------------------|-----------|---------|---------|-----|
| SIFT1M | 1,000,000 | 100,000 | 10,000 | 128 |
| GIST1M | 1,000,000 | 500,000 | 1,000 | 960 |
| MS MARCO Doc | 3,213,833 | 367,013 | 5,193 | 768 |
| MS MARCO Passage | 8,841,823 | 808,731 | 101,093 | 768 |

530 B.2 Baselines

531 The two groups of baseline ANNS models are compared to KDindex.

532 The first group is *Non-quantization-based ANNS methods*, which accelerate the search by reducing the
533 number of candidates. **BLISS** [19] adopts the learning-to-index framework to learn the hashing-based
534 compressed functions. **ScaNN** [18] quantizes with anisotropic quantization loss functions which
535 greatly penalizes the parallel component of a datapoint’s residual relative to its orthogonal component.
536 **HNSW** [35] builds a hierarchical set of proximity graphs.

537 The second is *Quantization-based ANNS methods*, which compress the embeddings by hashing
538 or quantization functions. **PQ** [22] decomposes the vector representation space into the Cartesian
539 product of subspaces. **OPQ** [16] jointly learns space decomposition and subspace quantization.
540 **AQ** [3] represents each vector as a sum of several components each coming from a separate codebook.
541 The baselines are implemented based on the Faiss ANNS library [25]. The parameters B and W
542 are set to be the same as KDindex. **DiffPQ** [5], differentiable product quantization, a generic and
543 end-to-end learnable compression framework. **DeepPQ** [15], deep progressive quantization, end-to-
544 end learns the quantization codes sequentially. **PQ-VAE** [45], an unsupervised model for learning
545 discrete representations by combining product quantization and auto-encoders. The CNN blocks are
546 replaced with MLP because the image datasets have been extracted as 512-dimension features. **GCD**
547 [23] learns rotation matrix via a family of block Givens coordinate descent algorithms. **RepCONC**
548 [49] requires data points to be uniformly clustered around the quantization centroids.

549 B.3 Implemental details

Table 6: Details of teacher model (HNSW).

| Teacher (HNSW) | SIFT1M | GIST1M | MS MARCO Doc | MS MARCO Passage |
|-------------------|--------|---------|--------------|------------------|
| M | 32 | 32 | 32 | 32 |
| efConstruction | 40 | 40 | 100 | 100 |
| efSearch | 100 | 512 | 1024 | 1024 |
| Search Time (s) | 0.5862 | 1.3082 | 1.4805 | 4.7689 |
| Building Time (s) | 20.1s | 2m25.4s | 17m52s | 98m17.2s |
| Recall@10 | 0.9865 | 0.9859 | 0.9292 | 0.9182 |

550 Teacher model is instantiated as HNSW. The details are described as Tab. 6, where M denotes
551 the number of neighbors each node, efConstruction denotes expansion factor at construction time
552 and efSearch denotes expansion factor at search time. To obtain good recall performance, M,
553 efConstruction and efSearch are tuned.

554 B.4 Complexity Analysis

555 For simplicity, we discuss the complexity of each codebook with W centroids. Posting List Balance
556 requires $\mathcal{O}(MWD)$ to calculate the similarities between the M document vectors and the centroids
557 and the space complexity is $\mathcal{O}(D^2)$. Besides, the query encoder brings an extra time cost of $\mathcal{O}(D^2)$
558 and space cost of $\mathcal{O}(WD)$. Overall, the time complexity and space complexity of KDindex is
559 $\mathcal{O}(D^2 + MWD)$ and $\mathcal{O}(D^2 + WD)$, respectively, which is acceptable since W and D are small
560 constants. As for the iterative initialization, the index assignment of documents only needs to be
561 updated after several epochs of centroids optimization. For the differentiable training, both the index
562 assignment and centroids are updated **every mini batch**.

563 C Varying Distillation Loss

564 C.1 Distillation Losses

565 Knowledge distillation was first proposed for classification tasks, where the probabilities of each class
566 attained from the large-scale teacher network are considered as soft labels to supervise the learning
567 of the small-size student network. The cross-entropy loss is commonly used as the distillation loss to
568 minimize the difference between the teacher and student networks. Here, the teacher search model

569 provides the top-k relevant candidates rather than the continuous value of probabilities. Thus, three
 570 ranking-oriented losses are designed to distill knowledge from the more accurate indexes to guide the
 571 student indexes to return the same nearest results.

572 **Lambdarank loss:** The pair-wise ranking-based loss is widely used to learn the ranking list by
 573 leading the high-ranked candidate to have higher similarity scores. Lambdarank [4] further introduces
 574 the change of the indicators, e.g., NDCG, to put more attention on more important candidates that
 575 have not been well ranked. The loss follows as:

$$\mathcal{L}(\mathbf{q}, \mathcal{D}_K^T; \mathbf{C}) = \sum_{i,j \in \mathcal{D}_K^T} \log(1 + \exp(p_i - p_j)) |\Delta NDCG@10_{ij}| \quad (4)$$

576 where \mathcal{D}_K^T denotes the top-k results retrieved from the teacher model and $p_i = S(\mathbf{q}, Q(\mathbf{d}_i))$ is the
 577 similarity score between the query vector and the quantized vector of the candidate i . Q is the
 578 quantizer function related to the codebooks \mathbf{C} . $\Delta NDCG@10_{ij}$ denotes the change with respect to
 579 NDCG@10 if changing the i -th ranked and j -th ranked candidate.

580 **Weighted KL loss:** Similar to the class distribution in classification tasks, the similarity distribution
 581 over the top-k retrieved candidates can also be obtained. One is based on the ground-truth vectors and
 582 the other one is based on the quantized vectors. In order to ensure the ranking orders correspond to the
 583 top-k list, the rank information is also considered where the high-ranked items are more concerned.
 584 Finally, the loss function follows as:

$$\mathcal{L}(\mathbf{q}, \mathcal{D}_K^T; \mathbf{C}) = - \sum_{i \in \mathcal{D}_K^T} \tilde{p}_i^g \log \frac{\tilde{p}_i^g}{\tilde{p}_i} \quad (5)$$

585 where \tilde{p}_i denotes the normalized ranked similarity score with the quantized vector and \tilde{p}_i^g with the
 586 ground-truth vector. Specifically,

$$p_i = w_i \cdot S(\mathbf{q}, Q(\mathbf{d}_i)), \quad p_i^g = w_i \cdot S(\mathbf{q}, \mathbf{d}_i),$$

587 \tilde{p}_i and \tilde{p}_i^g are the normalized values over the top-k retrieved candidates depending on the softmax
 588 function. $w_i = \frac{1}{rank(i)}$ denotes the ranking weight according to the ranking orders among the top-k
 589 results from the teacher model. The weighted KL loss attempts to minimize the distance between the
 590 ground-truth vector and the quantized vector for the top-k relevant candidates to learn better centroids.
 591 The introduced rank-oriented weight further guides the student index to return the same ranking list.

592 **Distributed-based loss:** Instead of being oriented by the score between query and candidates as
 593 above, we attempt to minimize the distance between the queries and top-k neighbors by calculating
 594 the similarity scores with all the centroids. Thus, we could obtain more information from centroids
 595 and focus on the top-K nearest neighbors. The distributed-based loss function follows as:

$$\mathcal{L}(\mathbf{q}, \mathcal{D}_K^T; \mathbf{C}) = - \sum_{i \in \mathcal{D}_K^T} \sum_{b=1}^B \sum_{k=1}^W \tilde{p}_{bk}^q \log(\tilde{p}_{bk}^{d_i} \cdot w_i) \quad (6)$$

596 where B denotes the number of codebooks and W is the number of codewords in each codebook.
 597 $w_i = \frac{1}{rank(i)}$ corresponds to the top-k list given from the teacher model. p_{bk}^q denotes the similarity
 598 score between the query \mathbf{q} and the codeword \mathbf{c}_k^b , i.e., $p_{bk}^q = S(\mathbf{q}, \mathbf{c}_k^b)$, and $p_{bk}^{d_i}$ denotes the similarity
 599 score between the candidate \mathbf{d}_i and the codeword \mathbf{c}_k^b , i.e., $p_{bk}^{d_i} = S(\mathbf{d}_i, \mathbf{c}_k^b)$. The normalized value
 600 \tilde{p}_{bk}^q and $\tilde{p}_{bk}^{d_i}$ are calculated over the W codewords for each codebook through the softmax function.

601 This loss requires the enumeration of all the centroids, while the Weighted KL loss only includes
 602 parts of the centroids corresponding to the quantized function. It also aligns with the goal of nearest
 603 searching for the query with the learnable centroids as the bridge.

604 C.2 Experimental Performances

605 We compare the effectiveness of the three different distillation losses, i.e., Weighted KL loss,
 606 Distributed-based loss, and Lambdarank loss as reported in Table 7.

607 **Findings.** Overall, the Distributed-based loss leads to comparatively better performances than
 608 Weighted KL loss and Lambdarank loss. Compared with Weighted KL Loss and Lambdarank

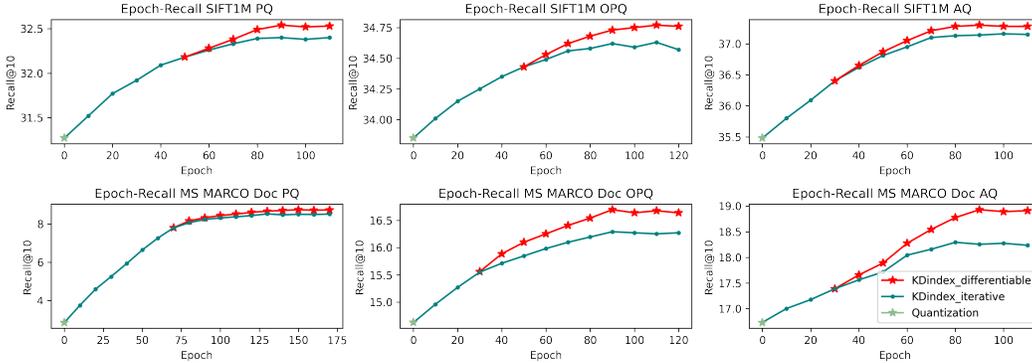


Figure 4: Curves for recall during training warmed up by initialization.

609 Loss, Distributed-based Loss gains the 2.03% and 3.54% improvements of Recall@10, 0.38%, and
 610 0.70% of NDCG@10, respectively. The Lambdarank Loss concerns more about the relationships
 611 between the pair of items, while the other two care about the whole ranking order of the list. The
 612 weighted KL loss actually optimizes parts of the centroids, depending on which query vectors and
 613 candidate vectors are quantized, to match the ranking list, while all of the centroids are updated in the
 614 Distributed-based loss since the probabilities are calculated over all the centroids. Furthermore, the
 615 Distributed-based loss requires the similarity calculation between the original input vectors and the
 616 centroids, which eliminates the error caused by the compressed functions. The last observation is that
 617 Distributed-based Loss works better on inner product metric datasets, since it obtains the average
 618 improvements of 2.24% and 3.33% of Recall@10 on L2 distance and inner product, respectively,
 619 wherein the overall improvements for inner-product similarities.

Table 7: The results of KDIndex under different distillation loss functions.

| Loss Function | SIFT1M | | GIST1M | | MS MARCO Doc | | MS MARCO Passage | |
|------------------------|-----------|---------|-----------|---------|--------------|--------|------------------|--------|
| | Recall@10 | NDCG@10 | Recall@10 | NDCG@10 | Recall@10 | MRR@10 | Recall@10 | MRR@10 |
| Lambdarank Loss | 36.32 | 79.18 | 20.69 | 62.06 | 17.60 | 40.01 | 11.10 | 34.86 |
| Weighted KL Loss | 36.68 | 79.33 | 21.02 | 62.75 | 18.24 | 40.93 | 11.06 | 34.63 |
| Distributed-based Loss | 37.30 | 80.01 | 21.33 | 63.17 | 18.93 | 41.69 | 11.19 | 35.23 |

620 D Performance of Differentiable Training

621 It is extremely difficult for the model to learn the codebooks as well as the index at the same time
 622 during the initialization phase in a differentiable training manner. Thus, we perform experiments by
 623 warming up the codebooks by *Initialization* and we get the following results in terms of Recall@10
 624 on four datasets as Fig. 4. We adopt the early stop strategy to get the best model.

625 **Initialization.** We obtain the pre-trained codebooks by iterative training manners and continue differ-
 626 entiable training when the index assignment is approaching being balanced ($\max|\mathcal{P}_i| - \min|\mathcal{P}_j| <$
 627 $\frac{N}{W}$, $i, j \in W$) where $|\mathcal{P}_i|$ denotes the length of the i -th posting list. To accelerate the iterative training,
 628 codebooks are warmed by original quantization methods such as PQ, OPQ, and AQ.

629 **Findings.** KDIndex converges to a better solution through the differential training manner. Within
 630 the dozens of epochs, the index assignment of iterative training becomes balanced, which warms
 631 up the centroids for later easier learning and thus relatively reduces the learning difficulty for both
 632 codebooks and indexes. Starting from this point, KDIndex with differentiable training consecutively
 633 outperforms that with iterative training, which achieves a relative improvement of 1.63% in terms of
 634 Recall@10 on both datasets, demonstrating the effectiveness of synchronizing updates for codebooks
 635 and indexes. As for the different quantization functions, the improvements of Recall@10 among
 636 different student models (PQ, OPQ, and AQ) are 1.49%, 1.46%, and 1.94%, respectively. The better
 637 performance of KDIndex(AQ) may be attributed to its better expressiveness with more parameters.
 638 Finally, the improvement of Recall@10 on MS MARCO Doc by KDIndex(PQ) is 0.40%, which
 639 is smaller than the other model since the express ability of PQ is limited. The improvement of

640 Recall@10 on SIFT1M by KDindex(AQ) is 0.39% since the express ability of AQ is strong on the
641 L2 distance dataset and no more improvements can be obtained easily.

642 **E Limitations and Future Works**

643 In this paper, we propose a novel knowledge distillation framework for high dimension index, which
644 reduce storage obviously and can learn neighbor information from the teacher model. Especially,
645 KDindex is independent to label (such as interaction information in the recommendation system or
646 ground-truth neighbors in ANNS), which makes it flexible to be applied in more label-free scenarios.
647 In the future, we will try more student models such as lattice quantization, whose codes already
648 imply neighbors relationship. And we will take labels into account to improve retrieval performance
649 progressively. We will further improve our work to benefit the broad community.