



Deep Triplet Hashing Network for Case-based Medical Image Retrieval

Jiansheng Fang^{a,b,d}, Huazhu Fu^c, Jiang Liu^{b,*}

^aSchool of Computer Science and Technology, Harbin Institute of Technology, Harbin 150001, China

^bDepartment of Computer Science and Engineering, Southern University of Science and Technology, Shenzhen 518055, China

^cInception Institute of Artificial Intelligence, Abu Dhabi, United Arab Emirates

^dCVTE Research, Guangzhou 510530, China

ARTICLE INFO

Article history:

Received 18 May 2020

Received in final form 10 Aug 2020

Accepted 12 Oct 2020

Available online 28 Jan 2021

Communicated by J.Liu

2000 MSC: 41A05, 41A10, 65D05, 65D17

Keywords: Medical Image Retrieval, Deep Hashing Methods, Spatial Attention, Region of Interest, Triplet Labels

ABSTRACT

Deep hashing methods have been shown to be the most efficient approximate nearest neighbor search techniques for large-scale image retrieval. However, existing deep hashing methods have a poor small-sample ranking performance for case-based medical image retrieval. The top-ranked images in the returned query results may be as a different class than the query image. This ranking problem is caused by classification, regions of interest (ROI), and small-sample information loss in the hashing space. To address the ranking problem, we propose an end-to-end framework, called Attention-based Triplet Hashing (ATH) network, to learn low-dimensional hash codes that preserve the classification, ROI, and small-sample information. We embed a spatial-attention module into the network structure of our ATH to focus on ROI information. The spatial-attention module aggregates the spatial information of feature maps by utilizing max-pooling, element-wise maximum, and element-wise mean operations jointly along the channel axis. To highlight the essential role of classification in differentiating case-based medical images, we propose a novel triplet cross-entropy loss to achieve maximal class-separability and maximal hash code-discriminability simultaneously during model training. The triplet cross-entropy loss can help to map the classification information of images and similarity between images into the hash codes. Moreover, by adopting triplet labels during model training, we can utilize the small-sample information fully to alleviate the imbalanced-sample problem. Extensive experiments on two case-based medical datasets demonstrate that our proposed ATH can further improve the retrieval performance compared to the state-of-the-art deep hashing methods and boost the ranking performance for small samples. Compared to the other loss methods, the triplet cross-entropy loss can enhance the classification performance and hash code-discriminability.

© 2021 Elsevier B. V. All rights reserved.

1. Introduction

With the rapid growth of medical images produced by various radiological imaging techniques, significant attention has

been devoted to the application of medical image processing technology. In recent decades, in particular, motivated by pattern recognition and computer vision techniques, such as deep learning methods, medical image processing has played an increasingly important role in assisting the diagnosis and assessment of diseases (Litjens et al., 2017; Fu et al., 2020; Orlando et al., 2020). However, the objective interpretation of medi-

*Corresponding author

e-mail: liuj@sustech.edu.cn (Jiang Liu)

cal images is fraught with high inter-observer variability and limited reproducibility. Further, although the outputs given by medical image classification tasks are only meant to be complementary to clinical decision-making (Doi, 2007), they still inevitably affect the expert decisions when discrepancies occur. To circumvent any discrepancies between expert interpretations, Content-Based Image Retrieval (CBIR) can present prior cases with similar disease manifestations to provide a reference-based assessment. CBIR aims to produce a fine-grained ranking of a large number of candidates according to their relevance to the query by indexing and mining large image databases (Zhou *et al.*, 2017). In effect, this helps create a context similar to the query, thus assisting in evidence-based clinical decision-making.

For better assistance in assessment, CBIR should have access to plenty of cases, which requires the retrieval algorithm to be both scalable and accurate. Hashing methods for CBIR have arisen as a promising solution for this, mapping high-dimensional feature descriptors to compact hash codes (Conjeti *et al.*, 2017a; Wu *et al.*, 2019). The low-dimensional hash codes can preserve the semantic structure of the high-dimensional feature descriptors and are suitable for efficient data storage and fast searching (Zhuang *et al.*, 2016). Hashing methods can be roughly categorized into data-dependent and data-independent methods (Wang *et al.*, 2017). Data-independent methods focus on using random projections to construct random hash functions. Compared with the data-dependent methods, data-independent methods need longer codes to achieve satisfactory performance (Gong *et al.*, 2012). Recent research focus has shifted to data-dependent methods, which learn hash functions in either a two-stage or end-to-end manner. The two-stage manner generates a vector of hand-crafted descriptors followed by learning the hashing functions. The similarity between two independent stages might not be optimally preserved by the hand-crafted features, and thus the learned hash codes are sub-optimal (Lai *et al.*, 2015). To address the drawbacks of the two-stage manner, deep hashing methods (Li *et al.*, 2015) with end-to-end training have been proposed to simultaneously learn image features and hash codes with deep neural networks, and have demonstrated superior performance over traditional hashing methods (Zheng *et al.*, 2017).

Currently, deep hashing methods are widely applied for application-specific medical image retrieval, such as deep multi-instance hashing for tumor assessment (Conjeti *et al.*, 2017a), deep residual hashing for chest X-ray images (Conjeti *et al.*, 2017b), order-sensitive deep hashing for multi-morbidity medical image retrieval (Chen *et al.*, 2018), etc. However, the hash codes learned by existing hashing methods usually lose information related to classification, regions of interest (ROI), and small samples. These three pieces of information play an essential role in enhancing the ranking quality of small samples in case-based medical image retrieval. As shown in Fig. 1, the negative image is more similar to the query image than the positive image in the hashing space. Such a ranking problem originates from the fact that information relating to classification, ROIs, and small samples is not fully mapped into the hash codes. For example, a very prominent wedge-shaped airspace

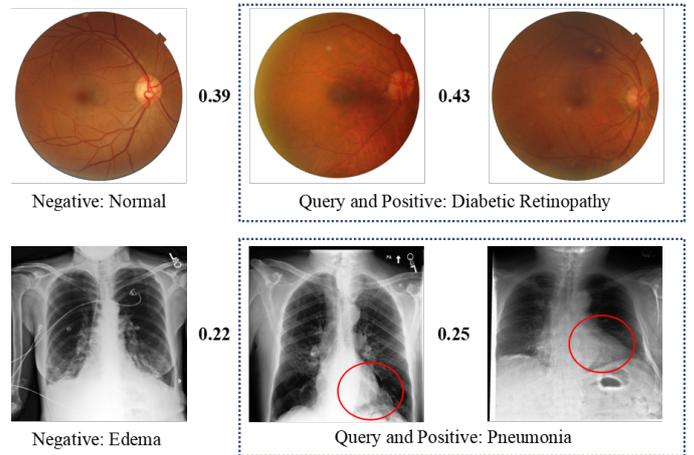


Fig. 1. Schematic of ranking problem. Each row represents a triplet sample, where the query image and positive image enclosed by a rectangle are from the same class. Each image is represented by a 36-bit hash code generated by existing deep hashing methods. Here we can see that the distance of hash codes between the query image and positive image is greater than the query image and negative image.

consolidation in the left lung (red circle) is so small that its information embedded into the compact hash codes is drowned out in the whole X-ray image. Clinically, the manifestation of the wedge-shaped region is characteristic of bacterial pneumonia. If the information of this region could be fully learned and mapped in the hash codes, this region could be used to discriminate between pneumonia and edema. Moreover, the information regarding small samples and their classes should be fully utilized during model training, so that the corresponding information can be mapped into the hash codes to play essential roles in differentiating case-based medical images.

To address the ranking problem, in this work, we present an end-to-end deep triplet hashing framework to learn hash codes with maximal discriminative capability, which we call the Attention-based Triplet Hashing (ATH) network. Inspired by the attention mechanism (Vaswani *et al.*, 2017; Woo *et al.*, 2018; Li *et al.*, 2019), we embed a spatial-attention module into the convolutional neural network (CNN) structure to promote discriminative capability by capturing the ROI information and mapping it into the hashing space. To alleviate the imbalanced-sample problem, we utilize the information of small samples fully, with the help of triplet labels during model training. A novel triplet cross-entropy loss is proposed to preserve the classification and similarity information in the hashing space, simultaneously. Given triplet labels, with the help of the spatial-attention module and triplet cross-entropy loss, the information loss related to classification, ROIs, and small samples can be alleviated to enhance the ranking quality for case-based medical image retrieval. The main contributions of this work are summarized as follows:

- 1) An end-to-end framework, named the Attention-based Triplet Hashing (ATH) network, is presented to address the ranking problem for case-based medical image retrieval. Our ATH aims to promote the discriminative capability of hash codes by preserving the information related to classi-

fication, ROIs, and small samples in the hashing space.

- 2) A spatial-attention module is embedded into the ATH network to boost the ROI representation by focusing on ROI information in the whole medical image. With the help of a novel triplet cross-entropy loss, maximal class-separability and maximal hash code-discriminability are simultaneously achieved during model training. To alleviate the imbalanced-sample problem to some extent, the hash codes are learned with the help of triplet labels in order to fully utilize small samples.
- 3) Extensive experiments on two case-based medical datasets demonstrate that our proposed ATH can further improve the retrieval performance compared to the state-of-the-art deep hashing methods and boost the ranking performance for small samples. Compared to the other loss methods, the triplet cross-entropy loss can enhance the classification performance and hash code-discriminability. Our code and model have been released in <https://github.com/fjssharpword/ATH>.

The rest of this paper is organized as follows: Section 2 discusses related works. Section 3 describes our methodology in detail. Section 4 extensively evaluates the proposed method on two medical images datasets. Section 5 gives concluding remarks.

2. Related Works

Existing hashing methods can be categorized into data-independent methods and data-dependent methods. The representative data-independent methods include Locality Sensitive Hashing (LSH) (Slaney and Casey, 2008) and Shift-Invariant Kernels Hashing (SIKH) (Raginsky and Lazebnik, 2009). The data-dependent methods, also called learning-based hashing methods, can be further categorized into (Chen *et al.*, 2018): **(1)** shallow learning-based hashing methods, like Metric Hashing Forests (MHF) (Conjeti *et al.*, 2016), Kernel Sensitive Hashing (KSH) (Liu *et al.*, 2016), and Spectral Hashing (SH) (Weiss *et al.*, 2009); **(2)** deep learning-based hashing methods, like Convolutional Neural Network Hashing (CNNH) (Xia *et al.*, 2014), Deep Hashing (DH) (Erin Liong *et al.*, 2015), Deep Hashing Network (DHN) (Zhu *et al.*, 2016), Simultaneous Feature Learning and Hashing (SFLH) (Lai *et al.*, 2015), Deep Semantic Ranking based Hashing (DSRH) (Zhao *et al.*, 2015), and Deep Similarity Comparison Hashing (DSCH) (Zhang *et al.*, 2015). The former learn hashing functions in a two-stage manner from hand-crafted features such as SIFT (Lowe, 2004), GIST (Oliva and Torralba, 2001), and the hash codes learning procedure is independent of the image features, which may lead to sub-optimal performance. In contrast to the former, the latter directly tailor features for hashing in an end-to-end manner with a powerful CNN, and have shown great promise recently.

Deep hashing methods leverage ground-truth labels to preserve similarity in the hash codes. Typically, the labeled data for ranking tasks are provided in one of two forms: pairwise labels or triplet labels (Wang *et al.*, 2016). The canonical examples of pairwise labels are Deep Residual Hashing (DRH)

(Conjeti *et al.*, 2017b) and Deep Pairwise-Supervised Hashing (DPSH) (Li *et al.*, 2015). DPSH was the first deep hashing method to simultaneously perform image feature learning and hash code learning with pairwise labels, and achieves the highest performance compared to other deep hashing methods. DRH was designed to preserve similarity and generate compact hashing code by defining a similarity matrix with pairwise labels for medical image retrieval. Representative methods of triplet labels include Deep Supervised Hashing (DSH) (Wang *et al.*, 2016) and Deep Binary Embedding Networks (DBEN) (Zhuang *et al.*, 2016). Because the triplet labels inherently contain richer information than pairwise labels, the DSH outperforms DPSH and other deep hashing methods.

While the aforementioned deep hashing methods have certainly achieved some degree of success, there still exists a ranking problem in the field of case-based medical image retrieval. One of the reasons for this is the classification information loss during model training. Inspired by the circle loss, which combines the triplet loss and cross-entropy loss (Sun *et al.*, 2020), we propose the triplet cross-entropy loss to preserve the classification information. In the circle loss, each similarity score is given different penalties according to its distance to the optimal effect. We argue that the optimal effect is still learned from samples during model training, so the triplet cross-entropy keeps the original form of both the triplet loss and cross-entropy loss. The other reason is the small-sample information loss during model training. We utilize triplet labels for model training to overcome the imbalanced-sample problem. Each triplet label can be naturally decomposed into two pairwise labels. In the hashing space, the query image is simultaneously close to the positive image and far from the negative image. Triplet labels explicitly provide a notion of relative similarities between images, while pairwise labels can only encode this implicitly (Wang *et al.*, 2016). Small samples are not only used as positive images themselves but also as negative images of large samples. Thus we argue that triplet labels, which can better exploit small-sample information during model training, can help to alleviate the imbalanced-sample problem.

Similar to DSH, we propose ATH to perform image feature learning and hash code learning simultaneously by maximizing the likelihood of the given triplet labels. With the given triplet labels, our ATH enhances the ranking quality for small samples by fully utilizing the small-sample information. Different from DSH, with the help of a novel triplet cross-entropy loss, maximal class-separability and maximal hash code-discriminability are simultaneously achieved during model training. We argue that class-separability information may be lost when punishing the similarity loss of global features. By conducting validations on chest X-ray images, DRH was demonstrated its better performance by preserving the class-separability. In DRH, a retrieval loss inspired by neighborhood component analysis is used for learning discriminative hash codes. However, the retrieval loss in DRH is only suitable for the co-occurring manifestation of multiple diseases. In this work, the policy on preserving the class-separability is to directly punish the classification loss and similarity loss simultaneously.

In addition to the triplet cross-entropy loss, a spatial-attention

mechanism (Zhu et al., 2019) is also introduced for case-based medical retrieval. Recently, attention mechanisms have been successfully applied in CNNs, significantly boosting the performance of many medical image tasks (Oktay et al., 2018; Nie et al., 2018), including segmentation, recognition, and classification. For instance, an attention-based CNN (Li et al., 2019) is proposed for glaucoma detection, including an attention prediction subnet, a pathological area localization subnet, and a glaucoma classification subnet. A novel Attention Gate (AG) (Schlemper et al., 2019) can also be easily integrated into standard CNN models to leverage salient regions in medical images for various medical image analysis tasks, including fetal ultrasound classification, and 3D CT abdominal segmentation. Attention mechanisms improve the performance by guiding the model activations to be focused around salient regions. Based on prior research, we argue that the spatial-attention mechanism can also be beneficial to the performance of medical image retrieval by capturing the ROI information. Different from the average-pooling and max-pooling applied in CBAM (Woo et al., 2018), we utilize max-pooling, element-wise maximum, and element-wise mean operations jointly along the channel axis to generate an efficient feature descriptor.

Based on the above discussion related to the novelty of this work, the proposed ATH has two key components: (1) a medical image feature learning component with a spatial-attention module; (2) a hash code learning component for image features, with the triplet cross-entropy loss. Extensive experiments on two medical image datasets demonstrate the effectiveness of our ATH.

3. Proposed Methodology

Our ATH aims to learn compact hash codes from original medical images with the given triplet labels. The hash codes should meet three requirements. (a) The query image should be encoded close to positive images and far from negative images in the hashing space. (b) The ROI information should be effectively encoded in discriminative hash codes. (c) The information related to small samples and their class should be fully mapped into the hash codes. Based on the spatial-attention mechanism and triplet cross-entropy loss, ATH is trained in an end-to-end manner, in which image feature learning and hash code learning from triplet labels are performed simultaneously.

3.1. Attention-based Network Architecture

Analytically, for the task of case-based medical image retrieval, the small-sample ranking problem originates from classification, ROI, and small-sample information loss in the hash codes. The fundamental goal of medical image retrieval tasks is to present prior cases with similar disease manifestations to assist evidence-based clinical decision-making. However, if the prior cases with similar disease manifestations (same class) rank lower than expected in the returned query list, the ranking quality not only impacts the effectiveness of clinical decision-making but also likely leads to error-prone diagnosis. To improve the ranking quality, we propose corresponding solutions for the classification, ROI, and small-sample information loss,

including a novel triplet cross-entropy loss, a spatial-attention mechanism, and triplet labels.

As shown in Fig. 2, we present an attention-based triplet hashing network to jointly learn visual feature extraction and the subsequent mapping to compact hash codes. The architecture of our ATH consists of a net1 module, a spatial-attention module, and a net2 module, and terminates in a dense layer for hash code-generation and classification outputs. The net1 module contains a residual block (He et al., 2016) and a max-pooling layer followed by a spatial-attention module. The spatial-attention module generates an attention map that is multiplied with the input feature maps. After the net2 module, which contains a residual block and an average-pooling layer, the dense layer is designed as a convolutional layer with $32 \times 32 \times 1$ nodes according to class activation mapping (CAM) (Zhou et al., 2016). The dense layer is separately mapped into the hash code-generation layer with k nodes and the classification output layer with c nodes. All the convolutional and pooling layers use 3×3 filters with stride 2 and are followed by batch normalization (Ioffe and Szegedy, 2015). All the convolutional layers and the fully connected layer are equipped with the ReLU (Nair and Hinton, 2010) activation function. The triplet images are input into the ATH at the same time to generate triplet hash codes and share network weights during training.

In ATH, the spatial-attention module inputs the feature maps F with $128 \times 128 \times 16$ dimensions and outputs an attention map $M(F)$ with $128 \times 128 \times 1$ dimensions by utilizing the inter-spatial relationship of features. We argue that each ROI region mainly consists of informative parts and salient points, both of which should respond to the gradient back-propagation. To focus on salient points, we use element-wise maximum (MaxPoint) and element-wise mean (AvgPoint) operations along the channel axis to generate two different spatial context descriptors: F_{avg} and F_{max} . The element-wise maximum and element-wise mean operations compute the maximum and average of each element along the channel axis, respectively. F_{avg} and F_{max} are calculated as:

$$\begin{aligned} F_{max} &= [f_1 \dots f_i \dots f_{128 \times 128}] & f_i &= \max_{1 \leq c \leq 16} X_i(c) \\ F_{avg} &= [f_1 \dots f_i \dots f_{128 \times 128}] & f_i &= \overline{X_i(c)} \end{aligned}, \quad (1)$$

where $X_i(c)$ is the response of the i^{th} point in c^{th} channel. To focus on informative parts, we use a max-pooling (MaxPool) operation to generate a max-pooled feature descriptor: F_{maxp} . All three descriptors are independently fed forward to a shared multi-layer perceptron (MLP) for denoising. Inspired by maximum activation of convolutions (Tolias et al., 2015), the MaxPoint operation encodes the maximal point response across feature maps of the convolutional layers. It is different from the MaxPool operation which encodes the maximal local response of each of the convolutional layers. The shared MLP contains a hidden layer, and the hidden activation size is set to the size of the dense layer. The MLP weights are shared by the three input descriptors and followed by the ReLU activation function. After the MLP, we concatenate and convolve the three descriptors into an attention map that encodes which areas to emphasize or suppress. The convolutional layer has a filter size of 3×3 and is followed by the tangent activation function. In short, the

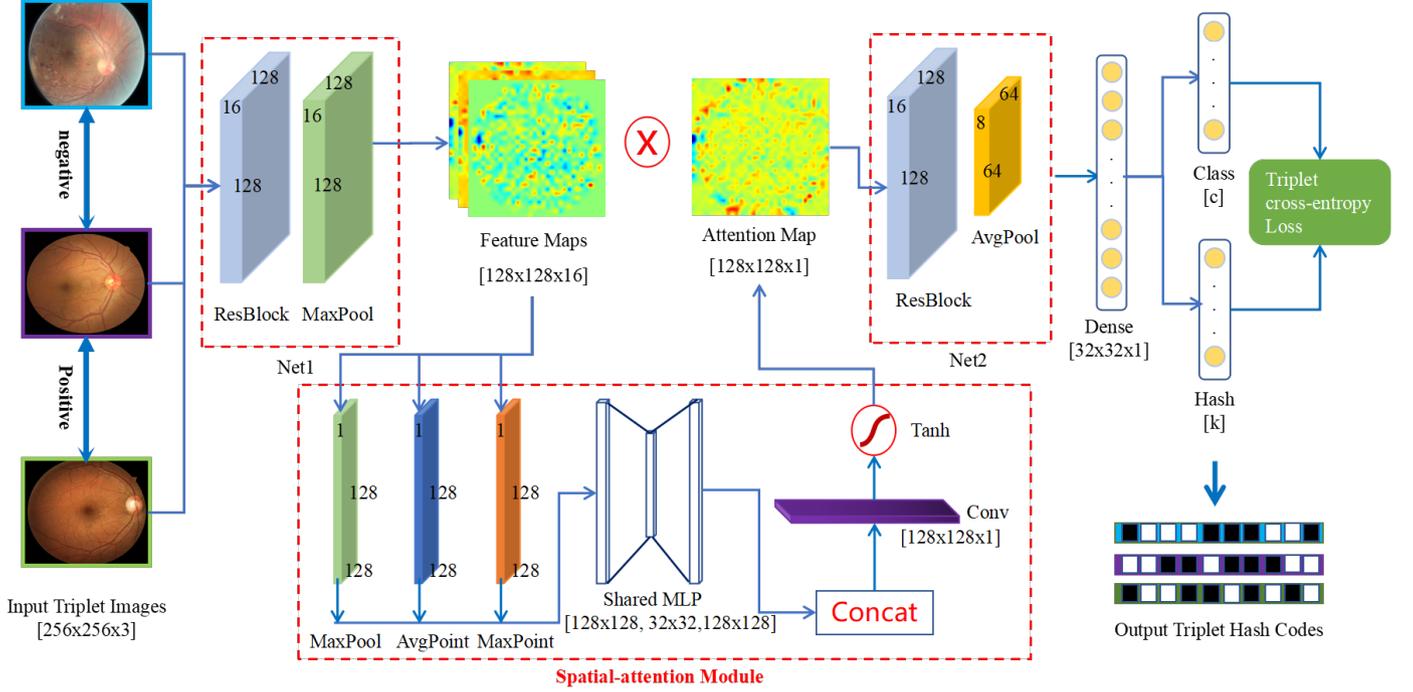


Fig. 2. Illustration of our ATH network structure. A spatial-attention module aggregates spatial information of feature maps by utilizing max-pooling (MaxPool), element-wise maximum (MaxPoint), and element-wise mean (AvgPoint) operations jointly along the channel axis. A dense layer is simultaneously mapped into the hash code-generation layer and classification output layer for training the triplet cross-entropy loss. The input triplet images share the network weights and are mapped into the corresponding triplet hash codes in which the classification, ROI, small-sample information is fully embedded.

spatial-attention is computed as:

$$\mathbf{M}(\mathbf{F}) = \mathcal{T}(\mathcal{F}([\text{MLP}(\mathbf{F}_{avg}); \text{MLP}(\mathbf{F}_{max}); \text{MLP}(\mathbf{F}_{maxp})])), \quad (2)$$

where \mathcal{T} denotes the tangent function, \mathcal{F} represents a convolution operation, and MLP is the operator of the shared MLP.

Based on the feature maps generated by the net1 module, we apply max-pooling, element-wise maximum, and element-wise mean operations jointly along the channel axis to generate a three-channel feature map. Each operation gathers an important clue related to the ROI information to infer a finer spatial-attention by computing spatial statistics. For each pixel, the maximum and average along the channel axis are computed for describing the channel-wise context features. The max-pooling operation prevents the marginal value of the ROI region from weakening. By exploiting both element-wise maximum and element-wise mean outputs independently, the features with distinct levels are input into the shared MLP module for further denoising. Such a spatial-attention module is beneficial for capturing the ROI information in high-resolution medical image. Using an attention-based CNN enables the learning procedure of our ATH to focus on ROI feature extraction on the raw pixels of input images. The hierarchical non-linear function exhibits a powerful learning capacity and encourages the learned feature to capture the ROI information by using the spatial-attention module.

3.2. Triplet Cross-entropy Loss

To overcome the ranking problem shown in Fig. 1, we propose a novel triplet cross-entropy loss to achieve maximal-class

separability and maximal hash code-discriminability by punishing similarity and classification losses simultaneously. The main idea behind the triplet cross-entropy loss is that the similarity and classification information should be simultaneously preserved during model training from triplet input images.

Mathematically, given m training images $\mathbf{I} = \{I_1, \dots, I_m\}$ and class labels $\mathbf{L} = \{1, \dots, c\}$, the triplet labels $\mathbf{T} = \{(Q_1, P_1, N_1), \dots, (Q_i, P_i, N_i), \dots, (Q_m, P_m, N_m)\}$ are generated by randomly selecting two images as a query image and a positive image from the same class (Q_i and P_i) and randomly selecting a negative image from different classes (Q_i and N_i). In the triplet labels, the query image of index Q_i is similar to the positive image P_i and dissimilar to the negative image N_i , where the index $i \in \{1, \dots, m\}$ is randomly selected from the m training images. When sampling triplet labels, small samples are selected as the negative image of large samples. In other words, small samples are reused multiple times during model training to preserve the information in the hash codes.

Generally speaking, ATH aims at learning a mapping from input triplet images to triplet hash codes $\mathbf{H}_Q = \{\mathbf{h}_{1,k}, \dots, \mathbf{h}_{m,k}\}$, $\mathbf{H}_P = \{\mathbf{h}_{1,k}, \dots, \mathbf{h}_{m,k}\}$, and $\mathbf{H}_N = \{\mathbf{h}_{1,k}, \dots, \mathbf{h}_{m,k}\}$. For scalable retrieval, the length of hash code k is much smaller than the dimension of medical image. The distance between \mathbf{H}_Q and \mathbf{H}_P should be smaller than the distance between \mathbf{H}_Q and \mathbf{H}_N . More specifically, the ROI information should be effectively encoded in the hash codes, such that the distances of similar images or dissimilar images can be affected by the ROI information. Based on the design of the triplet cross-entropy loss, ATH is trained with triplet la-

bels and ground-truth labels to perform hash code learning and classification likelihood learning simultaneously. Corresponding to the triplet labels \mathbf{T} , the ground-truth labels are $\mathbf{Y} = \{(y_{Q_1}, y_{P_1}, y_{N_1}), \dots, (y_{Q_i}, y_{P_i}, y_{N_i}), \dots, (y_{Q_m}, y_{P_m}, y_{N_m})\}$, where $y_{Q_i}, y_{P_i}, y_{N_i} \in \mathbf{L}$.

To punish the similarity loss given the triplet labels, a simple distance (e.g. Euclidean distance) is used to compare the similarity in the target space by mapping the input to the target space. In the training phase, the codes of query images and positive images should be as close as possible, while the codes of query images and negative images should be far away. Based on this objective, a hinge ranking loss form is naturally designed to minimize the distance between similar image pairs and maximize the distance between dissimilar image pairs. The loss with respect to triplet labels is defined as:

$$\mathcal{L}(\mathbf{T}) = \max\{r \cdot k - \text{Dist}(\mathbf{H}_Q, \mathbf{H}_N) + \text{Dist}(\mathbf{H}_Q, \mathbf{H}_P), 0\}, \quad (3)$$

where $\text{Dist}(\cdot, \cdot)$ denotes the L2-norm to measure the distance between hash outputs, k is the length of hash codes, and $r \in [0, 1]$ is a weighting parameter that controls the punishment strength of differentiating degrees between dissimilar images. The triplet loss is applied in such a way that only dissimilar pairs with the distance between them being within a specific radius are eligible to contribute to the loss. When $r = 0$, there is no punishment for dissimilar images mapped to close hash codes. When $r = 1$, the hash codes of dissimilar images must be completely different. If $r = 0.5$, this implies that half of the hash code lengths between dissimilar images should be different. The triplet loss is designed to measure how well the given triplet labels are satisfied by the learned hash codes by computing the likelihood of the given triplet labels.

To punish the classification loss given the ground-truth labels, we define the cross-entropy loss by jointly considering the input triplet images, as follows:

$$\mathcal{L}(\mathbf{T}, \mathbf{Y}) = \sum_{i=1}^m \{CE(\hat{y}_{Q_i}, y_{Q_i}) + CE(\hat{y}_{P_i}, y_{P_i}) + CE(\hat{y}_{N_i}, y_{N_i})\}, \quad (4)$$

where $CE(\cdot, \cdot)$ denotes the common cross-entropy loss form, and \hat{y} denotes the predicted class. With the similarity loss $\mathcal{L}(\mathbf{T})$ and classification loss $\mathcal{L}(\mathbf{T}, \mathbf{Y})$, we reverse the sum of both to update the weights of the model. Theoretically, the cross-entropy loss is beneficial for preserving the classification information in the hash codes, and the triplet loss can also help to improve the classification performance by encouraging hash codes to minimize intra-class similarity and maximize inter-class similarity.

4. Experiments

4.1. Datasets

- 1) **Fundus-iSee**: The private ophthalmic Fundus-iSee dataset with four disease classes consists of 10,000 high-resolution images labeled by professional doctors with rich clinical experience specifically. There are 720 images of age-related macular degeneration (AMD), 270 images of diabetic retinopathy (DR), 450 images of Glaucoma, 790

images of Myopia, and 7770 images of Normal. We randomly extract ten percent of each class for the query test, for a total of 1,000 images.

- 2) **MIMIC-CXR**: The public MIMIC-CXR (Johnson et al., 2019a) dataset is a large publicly available dataset of chest radiographs that are used to identify acute and chronic cardiopulmonary conditions and to assist in related medical workups. We randomly select two groups of images with frontal view from four classes, including Normal, Edema, Pneumonia, Fracture. The training set contains 14,555 images of Normal, 1,837 images of Edema, 3,220 images of Pneumonia, 388 images of Fracture. The test set contains 14,854 images of Normal, 944 images of Edema, 3,788 images of Pneumonia, 414 images of Fracture. The ratio of the training set and test set is 1 to 1.

4.2. Evaluation Settings

Three metrics are adopted to measure the precision and retrieval quality in our experiments. (1) **Hit Ratio (HR)**. HR is designed to measure how many images in the returned list are similar to the query image. (2) **Average Precision (AP)**. In the returned list, AP averages the rank positions of images similar to the query image to measure the rank quality. (3) **Reciprocal Rank (RR)**. RR refers to the reciprocal of the ranking of the first similar image in the returned list. The two datasets selected have a serious imbalanced-sample problem from the perspective of classification tasks. To validate the effectiveness of the triplet cross-entropy loss in preserving classification information in the hash codes, we apply **Specificity** and **Sensitivity** to measure the accuracy of classification.

In our comparative study, we use four deep hashing methods to evaluate the retrieval accuracy with mean HR (mHR), mean AP (mAP), mean RR (mRR), including DPSH-pairwise (Li et al., 2015), DRH-pairwise (Conjeti et al., 2017b), DSH-triplet (Wang et al., 2016), and DBEN-triplet (Zhuang et al., 2016). Two of the comparative methods use AlexNet (Krizhevsky et al., 2012) as the backbone, including DPSH-pairwise and DSH-triplet. Recently, the residual block (He et al., 2016) has been used popularly as the backbone in deep hashing methods, such as DRH-pairwise and DBEN-triplet, and shows the advantage of feature extraction. In our ATH, both the net1 and net2 modules also use the residual block as the backbone. Hence, compared to DRH-pairwise and DBEN-triplet, the advantage of our ATH framework can be demonstrated. To compare the effectiveness of the triplet cross-entropy loss, we introduce different loss methods into our ATH network, including cross entropy loss (ATH-CE), focal loss (ATH-focal) (Lin et al., 2017), circle loss (ATH-circle) (Sun et al., 2020), pairwise loss (ATH-pairwise), and triplet loss (ATH-triplet).

Our ATH is implemented with pyTorch, and the network structure of ATH is illustrated in Fig. 2. All deep hashing methods are trained from scratch, setting the batch size as 10 and the iteration number as 50. The parameters of comparative methods are set according to their implementation details in the corresponding papers, and the best performance is reported. The input triplet images are randomly sampled in every iteration. For the triplet cross-entropy loss in our ATH, the margin threshold

$r \cdot k$ should be considerably set to suit the clinical datasets better. Without loss of generality, the weighting parameter r and the length of hash codes k are synchronous to each other, so we empirically set $r = 0.5$ in the experiments. Without a special description, we report all the performances over 36-bit hash codes and weighting parameter r of 0.5, and the top-10 similar images are returned and ranked in every retrieval. Our ATH is implemented under Pytorch framework and experiments are run on Geforce RTX 2080 Ti. In our work, the indexing and similarity calculation for evaluation uses Faiss (Johnson et al., 2019b) which is a library for efficient similarity search and clustering of dense vectors.

4.3. Results and Analysis

To enhance the ranking quality of case-based medical image retrieval, we argue that the classification, ROI, and small-sample information loss in the hashing space need to be overcome. For ROI feature learning, we embed the spatial-attention module into the network to capture the ROI information. To preserve classification information, we propose the triplet cross-entropy loss to punish the similarity and classification losses simultaneously during model training. To overcome the imbalanced-sample problem, we sample the triplet labels from classification datasets in order to fully utilize the small-sample information.

4.3.1. Observation of ranking quality

As shown in Table 1, due to the effectiveness of the triplet cross-entropy loss and the spatial-attention mechanism, our ATH can further improve the performance compared to the second-highest (underlined) deep hashing methods given the triplet labels. The higher mAP means more similar images (same class) in the returned list are ranked ahead, while the higher mHR implies that more similar images are retrieved. The better performance in terms of the mAP and mHR demonstrates that our ATH not only achieves better accuracy of retrieval but helps users to find the required disease case quickly. In terms of mAP, which is used to evaluate the ranking quality, our ATH consistently outperforms the second-highest methods (underlined) by around 8%. From the perspective of mRR measuring ranking quality, our ATH consistently obtains gains above 0.90, which means that the first images in the returned list are nearly all from the same class as the query image. Following the convention in the literature of CBIR (Zhan and Zhao, 2018; Xiao and Zhao, 2020), we further provide the performance of mAP over varying top-k, where k varies from 5, 10 to 20. The results of Table 2 also confirm the benefits of our ATH. When the returned list top-k lengthens, the performance of all methods declines to some extent. We can observe that our ATH all achieves the best performance over the returned list of 5, 10, 20. According to the overall performance, we argue that our ATH can alleviate the information loss to some extent and fully map useful information into the hash codes.

We can conclude two points from Table 1 and Table 2. (1) The two best results of both datasets are related to the network structure of ATH, in which the spatial-attention module plays a significant role. (2) The order of performance (ATH>ATH-circle>ATH-triplet) demonstrates the merit of combining the

Table 1. The performances (mHR, mAP, mRR) on the Fundus-iSee and MIMIC-CXR datasets

Dataset	Methods	mHR	mAP	mRR
Fundus-iSee	DPSH-pairwise	0.6063	0.5182	0.8077
	DRH-pairwise	0.6761	0.5889	0.8555
	DSH-triplet	0.6146	0.5291	0.8187
	DBEN-triplet	0.7074	0.6378	0.8918
	ATH-CE	0.6558	0.5696	0.8441
	ATH-focal	0.6642	0.5807	0.8549
	ATH-circle	<u>0.7180</u>	<u>0.6621</u>	<u>0.9033</u>
	ATH-pairwise	0.6804	0.5907	0.8688
	ATH-Triplet	0.6981	0.6359	0.8845
	ATH(ours)	0.7682	0.7220	0.9256
MIMIC-CXR	DPSH-pairwise	0.7213	0.6810	0.8645
	DRH-pairwise	0.7434	0.7218	0.9028
	DSH-triplet	0.7699	0.7274	0.8853
	DBEN-triplet	0.7535	0.7260	0.9156
	ATH-CE	0.7592	0.7187	0.8943
	ATH-focal	0.7701	0.7333	0.9049
	ATH-circle	<u>0.7861</u>	<u>0.7629</u>	<u>0.9217</u>
	ATH-pairwise	0.7601	0.7386	0.9078
	ATH-triplet	0.7751	0.7467	0.9131
	ATH(ours)	0.8543	0.8260	0.9668

triplet loss and the cross-entropy loss. Due to the triplet labels, which enable the small-sample information to be fully used during training, our triplet cross-entropy loss outperforms the circle loss. Next, we further observe the ranking quality in Fig. 3. The distance between the hash codes of the same class is closer than that between different classes. By querying an AMD image, two AMD images are hit and ranked ahead in the Fundus-iSee dataset. By querying a pneumonia image, two pneumonia images are hit and ranked ahead in the MIMIC-CXR dataset. Clinically, the symptoms and manifestations of edema are similar to pneumonia. On the whole, our ATH can effectively solve the ranking problem in Fig. 1.

Based on the above experimental results and analysis, we have demonstrated the effectiveness of our ATH in addressing the ranking problem by fully mapping useful information into the hash codes. To illustrate the function of the spatial-attention mechanism in capturing ROI information, we provide heat maps of the dense layer outputs in Fig. 4. The heat maps of ROI regions between positive and negative results vary in degree and location. Such differences in ROI heat maps have an effect on class-separability. Each class has its distinctive ROI heat maps. Based on this observation, we argue that the features of ROI regions captured by the spatial-attention module are mapped into the hash codes. To further evaluate the contribution of the proposed spatial-attention module, we substitute

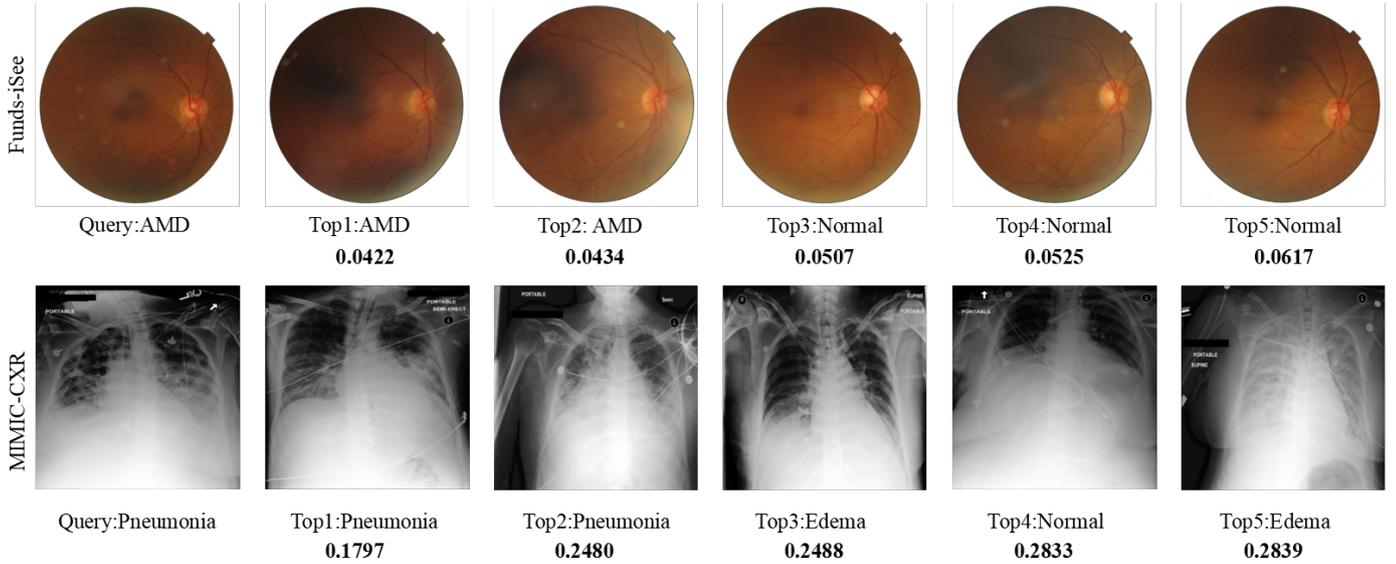


Fig. 3. Qualitative results of ranking quality. Upper row: querying an age-related macular degeneration (AMD) image, top-5 similar images with hash code distance are returned in the Fundus-iSee dataset. Lower row: querying a pneumonia image, top-5 similar images with hash code distance are returned in the MIMIC-CXR dataset.

Table 2. mAP over the varying number of the returned list on Fundus-iSee and MIMIC-CXR datasets

Dataset	Methods	top-5	top-10	top-20
Fundus-iSee	DPSH-pairwise	0.5400	0.5182	0.4946
	DRH-pairwise	0.5970	0.5889	0.5440
	DSH-triplet	0.5416	0.5291	0.5014
	DBEN-triplet	0.6459	0.6378	0.6091
	ATH-CE	0.5752	0.5696	0.4992
	ATH-focal	0.5918	0.5807	0.5174
	ATH-circle	<u>0.6813</u>	<u>0.6621</u>	<u>0.6260</u>
	ATH-pairwise	0.6286	0.5907	0.5262
	ATH-Triplet	0.6696	0.6359	0.5589
	ATH(ours)	0.7421	0.7220	0.6802
MIMIC-CXR	DPSH-pairwise	0.7099	0.6810	0.5901
	DRH-pairwise	0.7436	0.7218	0.6758
	DSH-triplet	0.7595	0.7274	0.6327
	DBEN-triplet	0.7551	0.7260	0.6960
	ATH-CE	0.7432	0.7187	0.6397
	ATH-focal	0.7595	0.7333	0.6862
	ATH-circle	<u>0.7954</u>	<u>0.7629</u>	<u>0.7371</u>
	ATH-pairwise	0.7672	0.7386	0.6990
	ATH-triplet	0.7616	0.7467	0.7137
	ATH(ours)	0.8412	0.8260	0.7628

the spatial-attention module with AG and CBAM in our ATH framework to compare the performance. Both AG and CBAM have been widely applied in networks and achieve competi-

tive performance in many tasks of computer vision. For mAP of the Fundus-iSee dataset and the MIMIC-CXR dataset, our spatial-attention module averagely outperforms AG by 6.9% and CBAM by 7.3%. This experiment shows that the proposed spatial-attention module can achieve better performance than AG and CBAM by integrating the three spatial descriptors: F_{avg} , F_{max} , and F_{maxp} . Based on the above observation, the capability of capturing the salient value of the spatial-attention module has been attested. The spatial-attention mechanism can contribute to the hash code-discriminability by capturing the ROI information.

4.3.2. Observation of the triplet cross-entropy loss

The triplet cross-entropy loss for model training achieves maximal class-separability and maximal hash code-discriminability simultaneously. On the one hand, as shown in Table 3, our ATH achieves the best performance of the mAP over each class on the two datasets. In the ablation study of the loss function, including ATH-focal, ATH-circle, ATH-pairwise, and ATH-triplet, all of them can achieve second-highest performance on some class (underlined). Our ATH can improve the ranking quality of small samples by mapping each class's semantic information into the hash codes without sacrificing the performance of large samples. On the other hand, with the help of the triplet cross-entropy loss, we can get the classification results for the task of medical image retrieval. As shown in Table 4, the sensitivity performance of small samples shows that the triplet cross-entropy loss is superior to the focal loss and circle loss (underlined) without affecting the performance on the large samples. Recently, the focal loss and circle loss have demonstrated their superiority in alleviating the imbalanced-sample problem. Compared to the focal loss and circle loss, the triplet cross-entropy loss achieves maximal class-separability and maximal hash code-discriminability simultaneously during model training.

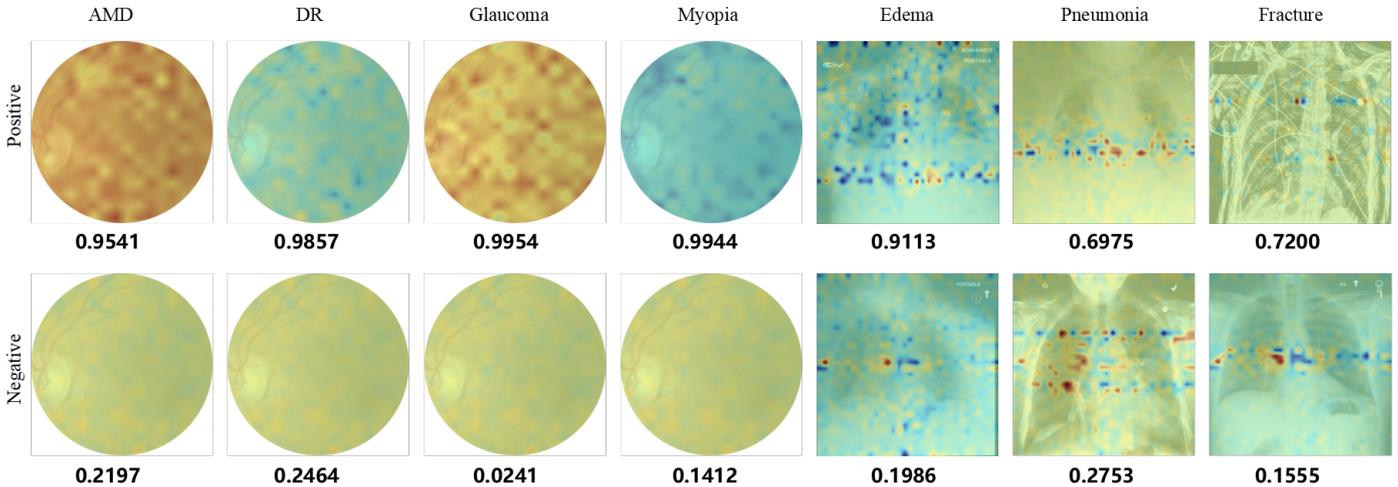


Fig. 4. Qualitative results of heat maps. Heat maps of the dense layer outputs are generated according to the class activation mapping. The upper row and lower row are the positive and negative results with predicting probability, respectively. The first four classes from left to right are Age-related Macular Degeneration (AMD), Diabetic Retinopathy (DR), Glaucoma, and Myopia in the Fundus-iSee dataset, and the last three classes are Edema, Pneumonia, Fracture in the MIMIC-CXR dataset. For example, by query an AMD image, the result predicted as AMD is positive, and predicted as non-AMD is negative. The negative result belongs to AMD with a probability of 0.2197, and the positive result is 0.9541.

Table 3. mAP of each class on the Fundus-iSee and MIMIC-CXR datasets.

Methods	Fundus-iSee Dataset					MIMIC-CXR Dataset			
	Normal	AMD	DR	Glaucoma	Myopia	Normal	Edema	Pneumonia	Fracture
DPSH-pairwise	0.6541	0.0299	0.0018	0.0126	0.0272	0.5940	0.0126	0.1916	0.0052
DRH-pairwise	0.6533	0.0264	0.0094	0.0142	0.1422	0.6732	0.2119	0.7831	0.0099
DSH-triplet	0.6623	0.0236	0.0080	0.0075	0.0216	0.5906	0.0321	0.1637	0.0096
DBEN-triplet	0.6711	0.0182	0.0123	0.0023	0.1735	0.7644	0.2158	0.7989	0.0114
ATH-CE	0.6629	0.0214	0.0091	0.0181	0.0550	0.7340	0.1026	0.5240	0.0147
ATH-focal	0.6865	0.0244	0.0073	<u>0.0272</u>	0.0494	0.7284	0.0547	0.4373	0.0095
ATH-circle	0.7279	<u>0.0360</u>	<u>0.0142</u>	0.0136	0.0491	<u>0.7719</u>	0.1035	<u>0.8368</u>	0.0087
ATH-pairwise	0.6935	0.0328	0.0089	0.0122	0.0375	0.7388	0.0732	0.8123	<u>0.0182</u>
ATH-triplet	<u>0.7734</u>	0.0168	0.0028	0.0062	<u>0.3384</u>	0.7520	<u>0.2771</u>	0.7656	0.0038
ATH(ours)	0.8564	0.0425	0.0185	0.0294	0.5995	0.8220	0.3068	0.8571	0.0332

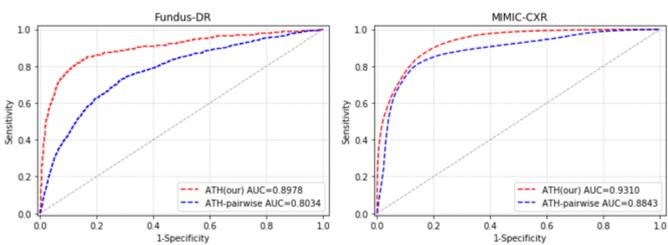


Fig. 5. The ROC curves with AUC scores of our ATH (red line) and ATH-pairwise (blue line) on the Fundus-iSee and MIMIC-CXR datasets.

For case-based medical image retrieval, a common issue is the imbalanced-sample due to the scarce disease cases. Fewer cases of certain types of diseases lead to their low sensitivity and high missed diagnosis rate. To alleviating the imbalanced-sample problem, triplet labels are input into the model to in-

crease the use of the small-sample information. Compared to the pairwise labels, the sampling mechanism of triplet labels demonstrates advantages in fully using small samples, as shown by observing the Receiver Operating Characteristic (ROC) curves with Area-Under-the-Curve (AUC) in Fig. 5. We argue that the triplet labels can play a role in maximizing inter-class distance and minimizing intra-class distance by utilizing small-sample information fully. The highest accuracy of classification means that the feature in the ATH network fully contains the information related to small samples and their classes. The classification output layer and the hash code-generation layer share the feature layers in the ATH network and are generated with the dense layer. Thus, given triplet labels, we can argue that the classification information can be mapped into the hash codes by using the triplet cross-entropy loss and can help to improve the hash code-discriminability.

By preserving the information of small samples and their

Table 4. Sensitivity of each class on the Fundus-iSee and MIMIC-CXR datasets.

Methods	Fundus-iSee Dataset					MIMIC-CXR Dataset			
	Normal	AMD	DR	Glaucoma	Myopia	Normal	Edema	Pneumonia	Fracture
ATH-CE	0.7745	0.2669	0.2107	0.2198	0.2739	0.7868	0.3220	0.6637	0.2188
ATH-focal	0.7747	0.2755	0.2055	0.2201	0.2843	0.8149	<u>0.3305</u>	<u>0.6998</u>	0.2562
ATH-circle	0.7902	<u>0.2855</u>	<u>0.2140</u>	<u>0.2288</u>	<u>0.2906</u>	0.7783	0.3059	0.6858	<u>0.2683</u>
ATH-pairwise	<u>0.8218</u>	0.2594	0.1852	0.1881	0.2394	<u>0.8303</u>	0.2806	0.6375	0.2186
ATH-triplet	0.7918	0.2794	0.2052	0.2081	0.2894	0.8018	0.3194	0.6852	0.2581
ATH(ours)	0.8375	0.4027	0.3132	0.3216	0.3955	0.8511	0.4552	0.7864	0.3728

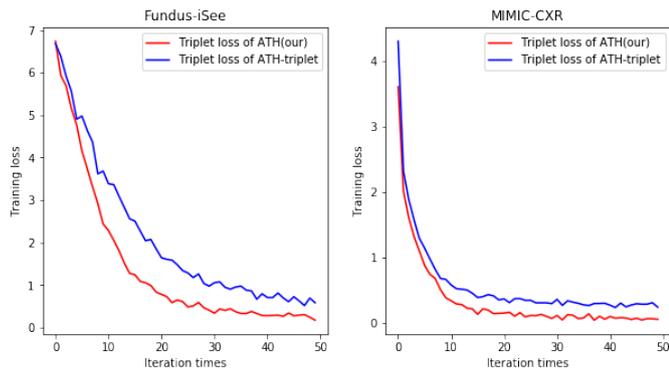


Fig. 6. Qualitative results of training loss. The triplet loss value between our ATH (red line) and ATH-triplet (blue line) is compared on the Fundus-iSee (left plot) and MIMIC-CXR (right plot) datasets.

class, we would like to demonstrate the cross-entropy loss can help to minimize intra-class similarity and maximize inter-class similarity. Although the circle loss and our triplet cross-entropy loss all combine the triplet loss and the cross-entropy loss, our triplet cross-entropy loss keeps both original forms to train on the triplet labels. The original form of the triplet loss punishes the similarity distance, the effect of which can thus be compared by observing the training loss value. We can easily extract the triplet loss value from the sum loss in our ATH and compare it to the triplet loss in the ATH-triplet. As shown in Fig. 6, with the help of the cross-entropy loss, the triplet loss value of our ATH is lower than the ATH-triplet. According to Equation 3, a lower triplet loss value indicates the better hash code-discriminability.

We further investigate the effectiveness of the weighting parameter r and the length of hash codes k of the triplet cross-entropy loss. As shown in Table 5, we observe the performances over hash codes with lengths of 12, 24, 36, and 48 by setting the weighting parameter r as 0.3, 0.5, and 0.7, respectively. The best performance is achieved by setting $r = 0.5$ and $k = 24$ on the Fundus-iSee dataset, and $r = 0.5$ and $k = 36$ on the MIMIC-CXR dataset. Reasonably, $r = 0.5$ refers to that half of the hash code lengths between dissimilar images should be different. With the hash codes lengthen, r can be correspondingly set higher than the short hash codes, then the performance can correspondingly improve at the cost of storage and search efficiency. As a trade-off between performance and search cost,

we set $r = 0.5$ and $k = 36$ in our experiments.

At the last analysis of experiments, the efficiency of the proposed method will be discussed on four-folds by putting the MIMIC-CXR dataset as an example.

- 1) **Feature computation time.** Based on the pre-trained ATH model with 36-bit hash codes, the feature extraction of the training set of 20,000 images can be completed in 28 seconds by using GPU.
- 2) **Retrieval time.** After feature extraction mapping into 36-bit hash codes, the index is built in 1 second by using Faiss. Then the retrieval of the test set of 20,000 images can be done in 1,346 ms by returning top-10 most similar images.
- 3) **Training time.** Training our ATH with an end-to-end manner takes 2,325 seconds by setting the iteration number as 50.
- 4) **Memory cost.** During model training with a batch size of 10, the memory-consuming is about 2,000 Mbps. The on-line search for the index also consumes about 2,000 Mbps.

Compared to the state-of-the-art methods, including DPSH-pairwise, DRH-pairwise, DSH-triplet, and DBEN-triplet, the complexity of our ATH is slightly higher in training time and memory cost due to the added attention module and is fair in feature computation time and retrieval time. According to the above analysis of efficiency, our ATH can provide fair real-time responses with significantly improving the performance, compared to the state-of-the-art deep hashing methods.

5. Conclusions

To enhance the ranking quality of case-based medical image retrieval, the proposed Attention-based Triplet Hashing network (ATH) is able to preserve classification, regions of interest (ROI), and small-sample information in the hashing space. We embed a spatial-attention module into the network to capture the ROI information. A novel triplet cross-entropy loss is proposed to preserve the classification information in the hash codes by punishing the similarity and classification losses simultaneously. Further, the triplet labels can fully utilize small samples to alleviate the imbalanced-sample problem to some

Table 5. mAP over the varying r and k of our method on Fundus-iSee and MIMIC-CXR datasets.

Parameter	Fundus-iSee Dataset				MIMIC-CXR Dataset			
	$k = 12$	$k = 24$	$k = 36$	$k = 48$	$k = 12$	$k = 24$	$k = 36$	$k = 48$
$r = 0.3$	0.6593	0.7051	0.6714	0.6431	0.7562	0.7164	0.7363	0.7433
$r = 0.5$	0.6470	0.7322	0.7220	0.6504	0.7354	0.7780	0.8260	0.7572
$r = 0.7$	0.6345	0.6534	0.6214	0.6659	0.7296	0.7434	0.7645	0.7718

extent. Experiments on two case-based medical datasets, including fundus images and chest X-rays, demonstrate that our ATH can obtain state-of-the-art performances in medical image retrieval. Further analysis confirms that the triplet cross-entropy loss can enhance classification performance and hash code-discriminability. Although our ATH achieves competitive performance, two promising directions for future research can be devoted to case-based medical image retrieval. First, we can integrate the triplet cross-entropy loss with regularizer to further differentiate images better for small samples, referring to existing hashing methods. Second, to focus on ROI information, we can design a new region-wise attention network that weights an attentive score of a region considering attentiveness.

Acknowledgments

The authors would like to thank many members of the Intelligent Medical Imaging (iMED) group for the inspiring knowledge sharing, technical discussions, clinical background infusion. Their helping hands make this paper a reality. We also would like to thank our partner Xiaoi Clinic for the data support. This work was supported in part by Guangdong Provincial Key Laboratory of Brain-inspired Intelligent Computation (Grant No. 2020B121201001).

References

- Chen, Z., Cai, R., Lu, J., Feng, J., Zhou, J., 2018. Order-sensitive deep hashing for multimorbidity medical image retrieval, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 620–628.
- Conjeti, S., Katouzian, A., Kazi, A., Mesbah, S., Beymer, D., Syeda-Mahmood, T.F., Navab, N., 2016. Metric hashing forests. *Medical image analysis* 34, 13–29.
- Conjeti, S., Paschali, M., Katouzian, A., Navab, N., 2017a. Deep multiple instance hashing for scalable medical image retrieval, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 550–558.
- Conjeti, S., Roy, A.G., Katouzian, A., Navab, N., 2017b. Hashing with residual networks for image retrieval, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 541–549.
- Doi, K., 2007. Computer-aided diagnosis in medical imaging: historical review, current status and future potential. *Computerized medical imaging and graphics* 31, 198–211.
- Erin Liang, V., Lu, J., Wang, G., Moulin, P., Zhou, J., 2015. Deep hashing for compact binary codes learning, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2475–2483.
- Fu, H., Li, F., Sun, X., et al., 2020. AGE Challenge: Angle Closure Glaucoma Evaluation in Anterior Segment Optical Coherence Tomography. *Medical Image Analysis*.
- Gong, Y., Lazebnik, S., Gordo, A., Perronnin, F., 2012. Iterative quantization: A procrustean approach to learning binary codes for large-scale image retrieval. *IEEE transactions on pattern analysis and machine intelligence* 35, 2916–2929.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770–778.
- Ioffe, S., Szegedy, C., 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*.
- Johnson, A., Pollard, T., Mark, R., Berkowitz, S., Horng, S., 2019a. MIMIC-CXR database.
- Johnson, J., Douze, M., Jégou, H., 2019b. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*.
- Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. Imagenet classification with deep convolutional neural networks, in: Advances in neural information processing systems, pp. 1097–1105.
- Lai, H., Pan, Y., Liu, Y., Yan, S., 2015. Simultaneous feature learning and hash coding with deep neural networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 3270–3278.
- Li, L., Xu, M., Wang, X., Jiang, L., Liu, H., 2019. Attention based glaucoma detection: A large-scale database and cnn model, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 10571–10580.
- Li, W.J., Wang, S., Kang, W.C., 2015. Feature learning based deep supervised hashing with pairwise labels. *arXiv preprint arXiv:1511.03855*.
- Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P., 2017. Focal loss for dense object detection, in: Proceedings of the IEEE international conference on computer vision, pp. 2980–2988.
- Litjens, G., Kooi, T., Bejnordi, B.E., Setio, A.A.A., Ciompi, F., Ghafoorian, M., Van Der Laak, J.A., Van Ginneken, B., Sánchez, C.I., 2017. A survey on deep learning in medical image analysis. *Medical image analysis* 42, 60–88.
- Liu, H., Wang, R., Shan, S., Chen, X., 2016. Deep supervised hashing for fast image retrieval, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2064–2072.
- Lowe, D.G., 2004. Distinctive image features from scale-invariant keypoints. *International journal of computer vision* 60, 91–110.
- Nair, V., Hinton, G.E., 2010. Rectified linear units improve restricted boltzmann machines, in: Proceedings of the 27th international conference on machine learning (ICML-10), pp. 807–814.
- Nie, D., Gao, Y., Wang, L., Shen, D., 2018. Asdnet: Attention based semi-supervised deep networks for medical image segmentation, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 370–378.
- Oktay, O., Schlemper, J., Folgoc, L.L., Lee, M., Heinrich, M., Misawa, K., Mori, K., McDonagh, S., Hammerla, N.Y., Kainz, B., et al., 2018. Attention u-net: Learning where to look for the pancreas. *arXiv preprint arXiv:1804.03999*.
- Oliva, A., Torralba, A., 2001. Modeling the shape of the scene: A holistic representation of the spatial envelope. *International journal of computer vision* 42, 145–175.
- Orlando, J.I., Fu, H., Barbosa Breda, J., et al., 2020. REFUGE Challenge: A unified framework for evaluating automated methods for glaucoma assessment from fundus photographs. *Medical Image Analysis* 59, 101570.
- Raginsky, M., Lazebnik, S., 2009. Locality-sensitive binary codes from shift-invariant kernels, in: Advances in neural information processing systems, pp. 1509–1517.

- Schlemper, J., Oktay, O., Schaap, M., Heinrich, M., Kainz, B., Glocker, B., Rueckert, D., 2019. Attention gated networks: Learning to leverage salient regions in medical images. *Medical image analysis* 53, 197–207.
- Slaney, M., Casey, M., 2008. Locality-sensitive hashing for finding nearest neighbors [lecture notes]. *IEEE Signal processing magazine* 25, 128–131.
- Sun, Y., Cheng, C., Zhang, Y., Zhang, C., Zheng, L., Wang, Z., Wei, Y., 2020. Circle loss: A unified perspective of pair similarity optimization. *arXiv preprint arXiv:2002.10857*.
- Tolias, G., Sicre, R., Jégou, H., 2015. Particular object retrieval with integral max-pooling of cnn activations. *arXiv preprint arXiv:1511.05879*.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I., 2017. Attention is all you need, in: *Advances in neural information processing systems*, pp. 5998–6008.
- Wang, J., Zhang, T., Sebe, N., Shen, H.T., et al., 2017. A survey on learning to hash. *IEEE transactions on pattern analysis and machine intelligence* 40, 769–790.
- Wang, X., Shi, Y., Kitani, K.M., 2016. Deep supervised hashing with triplet labels, in: *Asian conference on computer vision*, Springer. pp. 70–84.
- Weiss, Y., Torralba, A., Fergus, R., 2009. Spectral hashing, in: *Advances in neural information processing systems*, pp. 1753–1760.
- Woo, S., Park, J., Lee, J.Y., So Kweon, I., 2018. Cbam: Convolutional block attention module, in: *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 3–19.
- Wu, D., Dai, Q., Liu, J., Li, B., Wang, W., 2019. Deep incremental hashing network for efficient image retrieval, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9069–9077.
- Xia, R., Pan, Y., Lai, H., Liu, C., Yan, S., 2014. Supervised hashing for image retrieval via image representation learning, in: *Twenty-eighth AAAI conference on artificial intelligence*.
- Xiao, H.C., Zhao, W.L., 2020. Deeply activated salient region for instance search. *arXiv preprint arXiv:2002.00185*.
- Zhan, Y., Zhao, W.L., 2018. Instance search via instance level segmentation and feature representation. *arXiv preprint arXiv:1806.03576*.
- Zhang, R., Lin, L., Zhang, R., Zuo, W., Zhang, L., 2015. Bit-scalable deep hashing with regularized similarity learning for image retrieval and person re-identification. *IEEE Transactions on Image Processing* 24, 4766–4779.
- Zhao, F., Huang, Y., Wang, L., Tan, T., 2015. Deep semantic ranking based hashing for multi-label image retrieval, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1556–1564.
- Zheng, L., Yang, Y., Tian, Q., 2017. Sift meets cnn: A decade survey of instance retrieval. *IEEE transactions on pattern analysis and machine intelligence* 40, 1224–1244.
- Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A., 2016. Learning deep features for discriminative localization, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2921–2929.
- Zhou, W., Li, H., Tian, Q., 2017. Recent advance in content-based image retrieval: A literature survey. *arXiv preprint arXiv:1706.06064*.
- Zhu, H., Long, M., Wang, J., Cao, Y., 2016. Deep hashing network for efficient similarity retrieval, in: *Thirtieth AAAI Conference on Artificial Intelligence*.
- Zhu, X., Cheng, D., Zhang, Z., Lin, S., Dai, J., 2019. An empirical study of spatial attention mechanisms in deep networks, in: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 6688–6697.
- Zhuang, B., Lin, G., Shen, C., Reid, I., 2016. Fast training of triplet-based deep binary embedding networks, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5955–5964.