

OSCARS: An Outlier-Sensitive Content-Based Radiography Retrieval System

Xiaoyuan Guo^a, Jiali Duan^b, Saptarshi Purkayastha^c, Hari Trivedi^{c,d}, Judy Wawira Gichoya^{c,d}, Imon Banerjee^{f,g,*}

^aDepartment of Computer Science, Emory University, GA, USA

^bMing Hsieh Department of Electrical and Computer Engineering, University of Southern California, CA, USA

^cSchool of Medicine, Emory University, GA, USA

^dDepartment of Radiology and Imaging Sciences, Emory University, GA, USA

^eIndiana University-Purdue University Indianapolis, School of Informatics and Computing, IN, USA

^fDepartment of Radiology, Mayo clinic, Phoenix, AZ, USA

^gSchool of Computing and Augmented Intelligence, Arizona State University, AZ, USA

Abstract

Improving the retrieval relevance on noisy datasets is an emerging need for the curation of a large-scale clean dataset in the medical domain. While existing methods can be applied for class-wise retrieval (aka. inter-class), they cannot distinguish the granularity of likeness within the same class (aka. intra-class). The problem is exacerbated on medical external datasets, where noisy samples of the same class are treated equally during training. Our goal is to identify both intra/inter-class similarities for fine-grained retrieval. To achieve this, we propose an **Outlier-Sensitive Content-based rAdiology Retrieval System (OSCARS)**, consisting of two steps. First, we train an outlier detector on a clean internal dataset in an unsupervised manner. Then we use the trained detector to generate the anomaly scores on the external dataset, whose distribution will be used to bin intra-class variations. Second, we propose a quadruplet ($a, p, n_{intra}, n_{inter}$) sampling strategy, where intra-class negatives n_{intra} are sampled from bins of the same class other than the bin anchor a belongs to, while n_{inter} are randomly sampled from inter-classes. We suggest a weighted metric learning objective to balance the intra and inter-class feature learning. We experimented on two representative public radiography datasets. Experiments show the effectiveness of our approach. The training and evaluation code can be found in <https://github.com/XiaoyuanGuo/oscars>.

Keywords: Medical image retrieval, Deep metric learning, Outlier detection, Radiography

1. Introduction

With the widespread adoption of radiology in diagnosis and treatment planning, the amount of medical image data is rapidly increasing Hwang et al. (2012). Fast and effective retrieval in large-scale medical image repositories has been demanding to support data management, re-

search and clinical applications Sotomayor et al. (2021). One common way to retrieval images is content-based, which has been widely researched and applied to the medical field Wang et al. (2014); Dubey (2021); Chowdhury et al. (2016); Chen et al. (2022).

For a given query image, a content-based image retrieval (CBIR) system returns a ranked list of images from the database based on a similarity measure between the query and retrieved images Duan and Kuo (2021); Revaud et al. (2019). The core idea behind CBIR is to minimize the distance of an anchor image a to its positive counterparts ps and maximize the distance to the corresponding negative images ns in the feature space. Usually, the

*Corresponding author

Email addresses: xiaoyuan.guo@emory.edu (Xiaoyuan Guo), jialidua@usc.edu (Jiali Duan), saptpurk@iupui.edu (Saptarshi Purkayastha), hari.trivedi@emory.edu (Hari Trivedi), judywawira@emory.edu (Judy Wawira Gichoya), banerjee.imon@mayo.edu (Imon Banerjee)

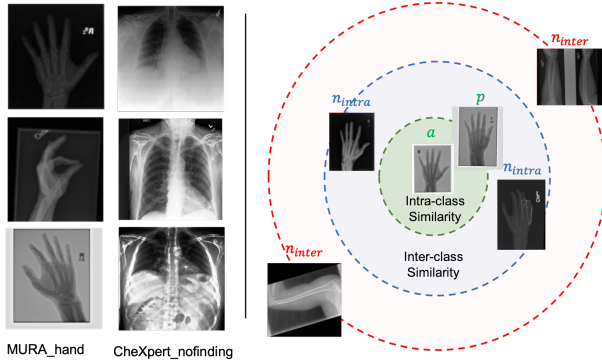


Figure 1: (Left) Examples of intra-class variations. 1st column shows samples from Stanford MURA *HAND* class and 2nd column presents CheXpert *No Finding* class data. (Right) Oscars learns the intra-class and inter-class similarity simultaneously. Images with intra-class similarity p should be closer to the given image a than the samples that show inter-class similarity n_{intra} in the feature space.

positive images are in the same class as the anchor image. However, adopting this strategy can be problematic as it only considers the inter-class variation. The assumption - as long as a and p are from the same class, they show similar visual features - is not realistic as samples from one class often exhibit certain intra-class variations. Noisy, under-represented data can exist, also called *outliers*. This phenomenon is more common in radiology as images are often acquired via different equipment from different sources and varies based on acquisition protocols. These variations, as shown in the left part of Fig. 1, pose specific challenges in the consumer domain and need to be recognized in assessing image similarity Akgül et al. (2011).

Although there have been multiple studies for radiograph retrieval, few of them pay attention to the intra-class similarity problem. Anavi et al. (2015) investigated X-ray image retrieval with both distance-based and probability-based approaches. Chowdhury et al. (2016) proposed a content-based medical image retrieval (CBMIR) system for radiographic images and employed a CNN to obtain high-level image representations. Qayyum et al. (2017) proposed a CBIR framework by training a CNN for the classification task. Layode and Rahman (2020) developed a chest X-ray image retrieval system for COVID-19 detection with deep denoising autoencoders as feature extractors. Zhong et al. (2021) designed an image re-

trieval system for COVID-19 chest radiograph via optimizing a multi-similarity loss. Outside medical domain, methods including FastAP Cakir et al. (2019), MultiSimilarity Wang et al. (2019), CircleLossSun et al. (2020) and SupCon Khosla et al. (2020) try to discover challenging negative data to improve the retrieval accuracy. Nevertheless, these existing efforts all emphasize the inter-class similarity but neglect the intra-class similarity.

In this paper, we focus on relevant radiograph image retrieval in external datasets which can contain lots of noisy data compared to the clean internal dataset. Such a system will help to *collect cleaner external image dataset with minimal human effort and accelerate AI evaluation*. To achieve the goal, we propose an **Outlier-Sensitive Content-based rAdiology Retrieval System (OSCARS)**, which takes both the intra-class and inter-class variations into consideration. To acquire the intra-class variation information, we adopt the unsupervised anomaly detectors trained on the internal dataset and utilize the assigned anomaly scores to the external dataset to split each class into several bins, with each bin in a certain range regarding the anomaly scores. Based on which, we construct the quadruplet data $(a, p, n_{intra}, n_{inter})$ with an anchor image a , a positive image p from the same class and same bin, an intra-class negative image n_{intra} from the same class but different bins, and an inter-class negative image n_{inter} that is from a different class.

With the proposed quadruplet sampling strategy, we incorporate the intra-class discriminative information into the training data and hence improve the retrieval of sensitivity outlier-related queries after model training. All the images in a quadruplet are fed into the feature extractor to learn their latent embeddings $(e_a, e_p, e_{n_{intra}}, e_{n_{inter}})$. As illustrated in the right of Fig. 1, we then learn the intra-class embedding similarity to achieve $(Sim(e_a, e_p) > Sim(e_a, e_{n_{intra}}))$ with an intra-class triplet loss L_{intra} and the inter-class similarity for $(Sim(e_a, e_{n_{intra}}) > Sim(e_a, e_{n_{inter}}))$ with an inter-class triplet loss L_{inter} in a weighted way. Our summarized contributions are:

1. We introduce the task of outlier-sensitive image retrieval for noisy external medical image dataset and propose an effective image retrieval system **OSCARS** to enhance the relevance of outlier-related results.
2. We propose to acquire intra-class information of ex-

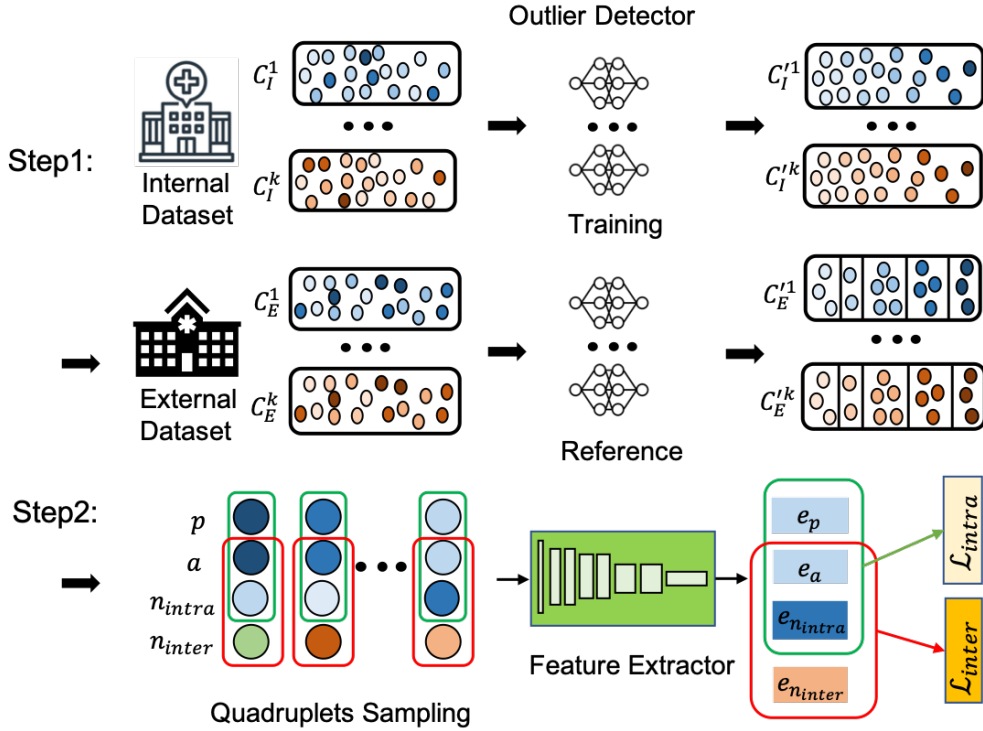


Figure 2: OSCARS architecture involves two main steps. Step1: train anomaly detectors on the internal dataset for each class C_I^i ; learn clean in-distributions with anomaly scores assigned to $C_I^{\prime i}$; apply the trained anomaly detectors on each class C_E^i of the external dataset and split the data into several bins $C_E^{\prime i}$ according to the anomaly scores. (Dark colors mean more distribution shifts.) Step2: generate quadruplets $(a, p, n_{intra}, n_{inter})$ by sampling the intra-class positive, negative and inter-class negative simultaneously; learn the intra-class and inter-class similarity in feature space with the intra-class triplet loss L_{intra} and inter-class triplet loss L_{inter} .

ternal datasets via anomaly detectors trained unsupervised. By training on clean internal datasets, the anomaly detectors assign each sample of the external dataset with a specific anomaly score. Based on which, we split each class into several bins with different intra-class variations.

3. We sample both the intra-class and inter-class negative images to construct quadruplets for intra-class and inter-class similarity learning.
4. We demonstrate the model effectiveness with two public representative radiography datasets - Stanford Musculoskeletal Radiography (MURA) Rajpurkar et al. (2017) and CheXpert Irvin et al. (2019).

2. Methodology

Given a clean internal dataset D_I and a noisy external dataset D_E , the external data of class c can contain outliers visually different from the internal class. Therefore, a conventional image retrieval system for the external dataset will be insufficient as it merely treats all the samples from one class as the same without considering the intra-class variations. Thus, the system will lack sensitivity to the outliers, undermining the retrieval accuracy. Our objective is to train an image retrieval model that will prioritize the images with both intra-class and inter-class dissimilarity during retrieval ranking. Figure 2 summarizes the whole framework of our model. There are mainly two steps involved. First, we design to learn intra-class information in an unsupervised way (introduced in Sec. 2.1). Second, we propose to sample training data that are with intra-class bin information and inter-class information (introduced in Sec. 2.2). With these steps, images with the same labels and similar contents are pulled together by maintaining intra-class similarity.

2.1. Learning intra-class information

Due to the difficulties of collecting annotated data with intra-class information provided in the medical domain, the outlier-sensitivity research on medical images has been delayed. To overcome the problem, we propose to generate intra-class labels automatically inspired by a recent work - MedShift Guo et al. (2021b). Given a clean internal dataset D_I , MedShift has suggested an approach

to identify outliers for noisy external dataset D_E . Following the same steps of MedShift, we first obtain the internal distribution information by training an unsupervised outlier detector named CVAD Guo et al. (2021a) for each class on the same internal datasets used in Guo et al. (2021b). Then, the trained anomaly detectors are evaluated on the external datasets as they have learnt intra-class discriminative features. Thus, each external data has its anomaly score, based on which we split each class into B bins with the K-Means clustering techniques Lloyd (1982); MacQueen et al. (1967). B (5 in our paper) is determined by the Elbow method Thorndike (1953). The resulting bins are in different anomaly score ranges. With the data from different bins, we get the intra-class labels. Given that both the intra-class and inter-class labels are available, for each image a , we randomly sample one intra-class positive image p , one intra-class negative sample n_{intra} and one inter-class negative sample n_{inter} accordingly, thus collecting the quadruplets $(a, p, n_{intra}, n_{inter})$ for training.

2.2. Balancing the inter- and intra-class influence

With the sampled quadruplets data, we feed each of the image to a CNN-based feature extractor to acquire latent embeddings $(e_a, e_p, e_{n_{intra}}, e_{n_{inter}})$. For simplicity, we adopt the ResNet18 He et al. (2016) pre-trained on ImageNet Deng et al. (2009) as the network backbone. OSCARS is designed to consider both the inter-class similarity and the intra-class similarity at the same time, which brings the model advantages of acquiring the sensitivity of intra-class outlier relevance during image retrieval. However, balancing the effect of the two parts is a challenging problem. Too much weight on intra-class information will distract the general retrieval accuracy of inter-class data. Therefore, we design an intra-class triplet margin loss and an inter-class triplet margin loss to optimize the model architecture. To balance the influence of intra-class and intra-class information on final ranking, we adopt a weighted loss formulated as:

$$\begin{aligned}
 \mathcal{L} &= \lambda \mathcal{L}_{intra}(e_a, e_p, e_{n_{intra}}) + (1 - \lambda) \mathcal{L}_{inter}(e_a, e_{n_{intra}}, e_{n_{inter}}) \\
 &= \lambda(\max\{d(e_a, e_p) - d(e_a, e_{n_{intra}}) + \mathcal{M}_{intra}, 0\}) \\
 &\quad + (1 - \lambda)(\max\{d(e_a, e_{n_{intra}}) - d(e_a, e_{n_{inter}}) + \mathcal{M}_{inter}, 0\})
 \end{aligned} \tag{1}$$

where $d(x_i, y_i) = \|x_i - y_i\|_2$. λ , \mathcal{M}_{intra} and \mathcal{M}_{inter} are set as 0.05, 1 and 2 in our experiments respectively.

When we have a query image unseen during training, we first acquire the query representation with the trained image feature backbone and then compute the cosine similarity between the representative features of the query image and dataset images. Images are ranked based on the similarity scores in the descending order.

3. Experiments

We have evaluated our approach on two publicly available large-scale radiograph image datasets. The first is Stanford MURA dataset, a large dataset of bone X-rays, which contains seven classes - *HAND*, *FORARM*, *FINGER*, *SHOULDER*, *ELBOW*, *WRIST*, *HUMERUS*. The second is CheXpert dataset, which in total has 14 classes - *No Finding*, *Enlarged Cardiomediastinum*, *Cardiomegaly*, *Lung Lesion*, *Lung Opacity*, *Edema*, *Consolidation*, *Pneumonia*, *Atelectasis*, *Pneumothorax*, *Pleural Effusion*, *Pleural Other*, *Fracture*, *Support Devices*. As the chest x-ray images are with two views - frontal and lateral. We here only use frontal view and leave the lateral view for future studies. See more details in the supplementary materials.

3.1. Evaluation Metrics

For the retrieval task, we report the retrieval *recall* at rank K ($R@K$, $K \in \{1, 5, 10, 50, 100\}$), *precision* at rank K ($P@K$, $K \in \{1, 5, 10, 50, 100\}$), outlier sensitivity ($S@K$, $K \in \{1, 5, 10, 50, 100\}$). The metric *recall* is the percentage of relevant images retrieved over the total number of retrieved images, defined as $recall = \frac{N_R}{K}$ where R represents the relevant images retrieved. The metric *precision* is assigned based on the existence of the same labels between the query image and the retrieved images. If $\delta(\cdot) \in \{0, 1\}$ is an indicator function, the *precision* is defined as $precision = \frac{\sum_{i=1}^K \delta(R^i > 0)}{K}$. Additionally, we evaluate the outlier sensitivity by calculating the anomaly score difference with $sensitivity = \sum_{i=1}^{N_R} \frac{|\mathcal{A}_R^i - \mathcal{A}_Q|}{N_R}$, where \mathcal{A} means anomaly score. We scale the anomaly scores of MURA dataset into $[0, 1]$ with the sigmoid function due to the large variations of its anomaly scores.

3.2. Implementation Details

The pipelines are developed using Pytorch 1.9.0, Python 3.7.3 and Cuda compilation tools V11.4 on a machine with 4 NVIDIA Quadro RTX A6000 GPUs with 48GB memory. The training for all the models is run for 50 epochs with a start learning rate 0.001 and a SGD optimizer.

3.3. Search Results

As a representative image retrieval method with triplet data in training, we select DeepRank as our baseline. State-of-the-art CBIR approaches including FastAP, MultiSimilarity, CircleLoss and SupCon are used to compare the model performance. Notably, we keep the feature extractor consistent for all the methods to ensure fair comparisons.

Quantitative Results: Table 1 presents the recall and precision performance for both Stanford MURA and CheXpert datasets respectively. Since the data in CheXpert can have multiple labels, we calculate the correct hit with the strategy - loose match, which means that for a query chest X-ray with multiple labels, a retrieval image is relevant as long as it has one label matched. Compared to the baseline DeepRank, Oscars can enhance the recall and precision performances on both datasets and achieve the best recall at 1 and precision at 1. In general, SupCon has the highest recall for MURA dataset. Nonetheless, Oscars achieves the best precision for MURA dataset and recall for CheXpert. Additionally, we report the sensitivity results in the supplementary materials.

Qualitative Results: Figure D.7 shows an example of a *HAND* query image in MURA dataset. The corresponding retrieval results including ours are present in different rows. As can be seen, although many methods can achieve high recall and precision (see Table 1), they fail to distinguish the intra-class variations. Especially, MultiSimilarity and SupCon exhibit little sensitivity to the noisy query. Comparatively, our method can prioritize intra-class similarity and rank the images with similar anomaly semantics ahead. Please refer to the supplementary materials for more results.

Table 1: Quantitative performance on Stanford MURA and CheXpert datasets . Bold indicates the best.

Method	MURA										CheXpert									
	R@1↑	R@5↑	R@10↑	R@50↑	R@100↑	P@1↑	P@5↑	P@10↑	P@50↑	P@100↑	R@1↑	R@5↑	R@10↑	R@50↑	R@100↑	P@1↑	P@5↑	P@10↑	P@50↑	P@100↑
DeepRank Wang et al. (2014)	0.912	0.914	0.912	0.906	0.903	0.912	0.964	0.972	0.984	0.988	0.734	0.694	0.716	0.721	0.442	0.734	0.911	0.961	1.000	1.000
FastAP Cakir et al. (2019)	0.927	0.930	0.931	0.932	0.933	0.927	0.956	0.961	0.973	0.977	0.734	0.742	0.733	0.716	0.710	0.734	0.943	0.968	1.000	1.000
MultiSimilarity Wang et al. (2019)	0.923	0.921	0.919	0.915	0.913	0.923	0.955	0.968	0.977	0.980	0.695	0.680	0.677	0.676	0.682	0.695	0.957	0.975	1.000	1.000
CircleLossSun et al. (2020)	0.929	0.932	0.933	0.934	0.934	0.929	0.960	0.964	0.979	0.985	0.727	0.703	0.718	0.717	0.726	0.727	0.936	0.968	1.000	1.000
SupCon Khosla et al. (2020)	0.930	0.933	0.936	0.938	0.937	0.930	0.964	0.971	0.981	0.985	0.776	0.730	0.720	0.734	0.726	0.776	0.936	0.950	1.000	1.000
OSCARS (ours)	0.931	0.922	0.920	0.913	0.910	0.931	0.965	0.974	0.986	0.991	0.787	0.763	0.747	0.745	0.743	0.787	0.908	0.950	1.000	1.000

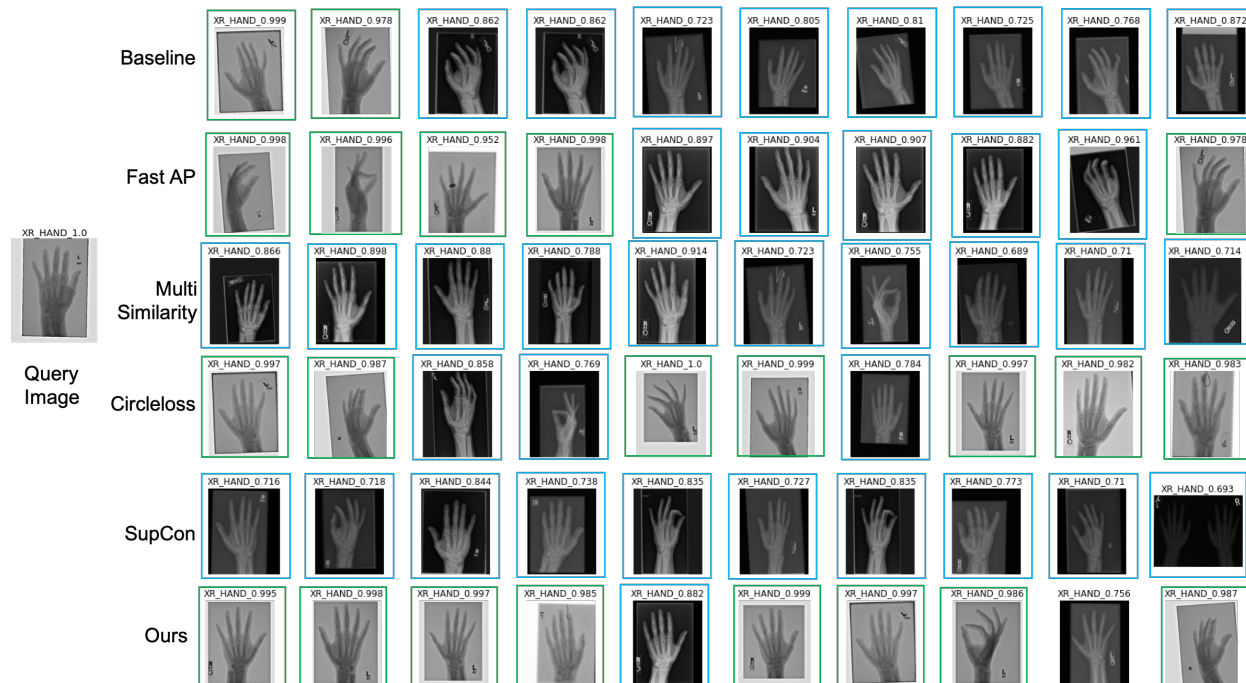


Figure 3: Hand results, left is the query image, right shows retrieval results. Green boxes mean both intra- and inter-class correct; blue boxes are for inter-class correct predictions. Each retrieval image has its label on top of itself. For correct predictions, we also put the anomaly scores on them. Closer anomaly scores mean more similarity.

Impact of Lambda: We also explore the impacts of applying different λ values to the loss function (Eqn. 1). A good balance between the intra-class and inter-class information will enable the retrieval system to acquire both accurate inter-class and outlier-sensitive intra-class results. Figure 4 illustrates the performance variations in different datasets under different settings. λ decides how the model learns to weight the intra-class and inter-class information simultaneously. We observe that too much weight on the intra-class similarity will degrade the inter-class similarity predictions. Experiments suggest 0.05 can work well.

4. Conclusion

In this work, we propose an outlier-sensitive radiography image retrieval system **OSCARS**, which goes beyond retrieving images with the most inter-class similarity but also inspects the intra-class similarity implicitly when query images show certain variations. Utilizing the automatic learning of clean internal distribution, the intra-class variations of external sources are captured and used to generate intra-class labels by splitting the class into several groups. Feeding the sampled quadruplets consisting of both the intra, inter-class positive and negative samples to the image feature learner, a weighted margin loss is adopted to optimize the retrieval network. The resulting

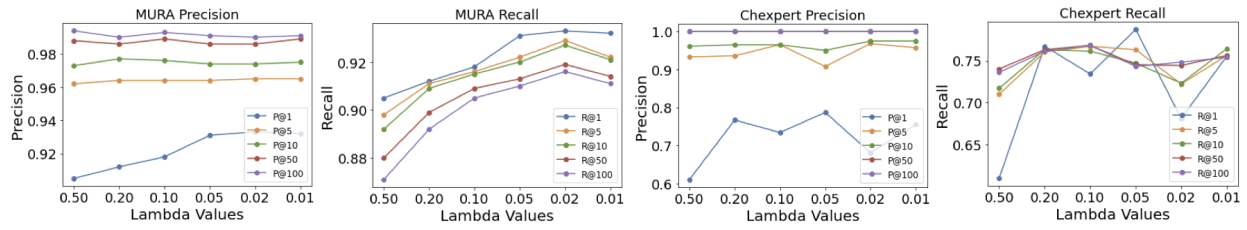


Figure 4: Effects with different lambda values on different datasets.

retrieval system is sensitive to outlier-related queries as it has learnt to rank the retrieved results based on both intra-class and inter-class similarities. This outlier-sensitive image retrieval approach provides clinical users the access to receive more relevant medical images and allow radiologists to process and analyze radiography images more effectively.

Appendix A. Dataset Details

We here introduce the details about the Stanford MURA and CheXpert dataset, and the samples we used in training and evaluation.

Stanford MURA contains 21,471 images. We sample a quadruplet for each image and split the quadruplets into training and validation with the ratio of 8:2. We evaluate the retrieval performance with the left unseen 1873 images.

CheXpert in total has 223,414 training images, of which 138,358 are in frontal view. By filtering out the invalid samples, 118,286 are left. For each image, we sample a quadruplet, resulting 118,286 training quadruplets. We thus split the quadruplets into training(80%) and validation(20%) parts. The left 282 frontal chest X-rays are used for testing.

Since both the two datasets are in varied sizes, we resize all the images into a fixed size of $224 \times 224 \times 3$ to fit the feature extractor network.

Appendix B. Sensitivity Results

We report the sensitivity results on both MURA and CheXpert datasets. We only calculate the sensitivity for the correct hits. Therefore, even in some situations, a model has lower sensitivity values, the general evaluation of the model performance should take the recall and precision into consideration. Because the anomaly score ranges of MURA can vary a lot, we scale its score into the range of [0,1] with a sigmoid function. For CheXpert dataset, we keep the original anomaly scores in use. As CheXpert data is with multi-labels, one sample with more than one label can have multiple anomaly scores considering each class variations. Therefore, when there are multiple hits, we take the minimum difference of the anomaly scores between the query image and the database images. Compared to MURA classes, chest X-rays are often similar with each other and thus difficult to retrieve. We here present the results for CheXpert with higher float precision. Generally, the lower the sensitivity values the better.

Table B.2: Sensitivity Similarity of Stanford MURA and CheXpert datasets. Best values are in bold.

Method	MURA					CheXpert				
	S@1↓	S@5↓	S@10↓	S@50↓	S@100↓	S@1↓	S@5↓	S@10↓	S@50↓	S@100↓
DeepRank	0.034	0.035	0.036	0.039	0.041	0.01157	0.01252	0.01217	0.01221	0.01217
FastAP	0.036	0.037	0.038	0.041	0.043	0.01369	0.01287	0.01281	0.01293	0.01285
MultiSimilarity	0.035	0.037	0.038	0.041	0.043	0.01115	0.01250	0.01265	0.01303	0.01338
CircleLoss	0.036	0.038	0.039	0.041	0.043	0.01600	0.01498	0.01529	0.01469	0.01518
SupCon	0.038	0.039	0.040	0.043	0.045	0.01639	0.01366	0.01309	0.01348	0.01348
OSCARS (<i>ours</i>)	0.030	0.032	0.033	0.036	0.038	0.01090	0.01015	0.00998	0.01039	0.01037

Appendix C. Visualization of MURA Retrieval

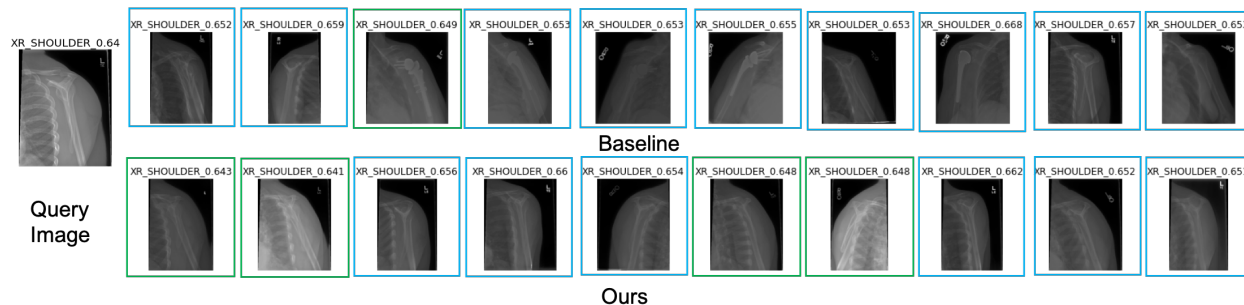


Figure C.5: Query results, left is the query image, retrieval results are shown in the right part. Green boxes mean both intra- and inter-class correct; blue boxes are for inter-class correct predictions. Each retrieval image has its corresponding label on top of itself, and for the correct predictions, we also put the anomaly scores on them. Closer anomaly scores mean more similarity.



Figure C.6: More query results for MURA dataset. Captions follow the same style as Fig. C.5.

Appendix D. Visualization of CheXpert Retrieval

Since a sample of CheXpert can have more than one labels, we encode the labels into a binary code with a length of 14, of which 1 means the data belongs to the class and 0 means irrelevant. The 14-bit label corresponds to the classes *No Finding*, *Enlarged Cardiomeastinum*, *Cardiomegaly*, *Lung Lesion*, *Lung Opacity*, *Edema*, *Consolidation*, *Pneumonia*, *Atelectasis*, *Pneumothorax*, *Pleural Effusion*, *Pleural Other*, *Fracture*, *Support Devices* in order. For simplicity, we only show labels, not with the anomaly scores.

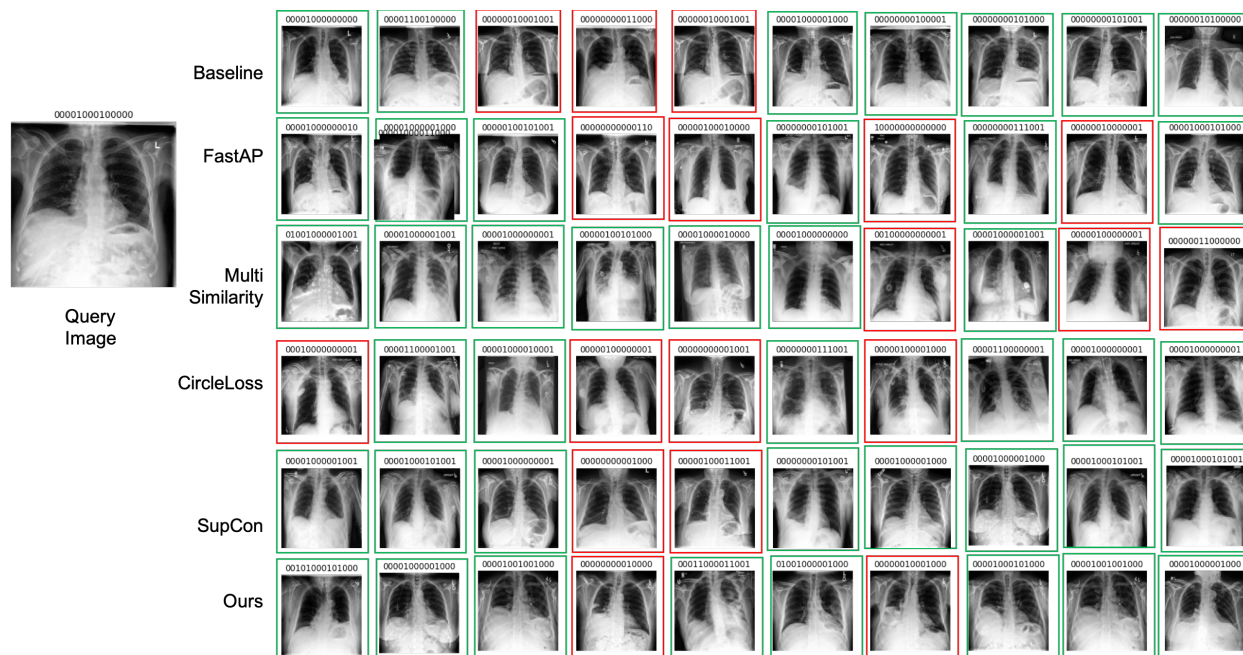


Figure D.7: CheXpert query results, left is the query image, retrieval results are shown in the right part. Green boxes mean both intra-class and inter-class correct; blue boxes are for inter-class correct predictions and red boxes are for wrong predictions. Each retrieval image has its corresponding label on top of itself.

References

- Akgül, C.B., Rubin, D.L., Napel, S., Beaulieu, C.F., Greenspan, H., Acar, B., 2011. Content-based image retrieval in radiology: current status and future directions. *Journal of digital imaging* 24, 208–222.
- Anavi, Y., Kogan, I., Gelbart, E., Geva, O., Greenspan, H., 2015. A comparative study for chest radiograph image retrieval using binary texture and deep learning classification, in: 2015 37th annual international conference of the IEEE engineering in medicine and biology society (EMBC), IEEE. pp. 2940–2943.
- Cakir, F., He, K., Xia, X., Kulis, B., Sclaroff, S., 2019. Deep metric learning to rank, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1861–1870.
- Chen, W., Liu, Y., Wang, W., Bakker, E., Georgiou, T., Fieguth, P., Liu, L., Lew, M.S., 2022. Deep learning for instance retrieval: A survey. *arXiv:2101.11282*.
- Chowdhury, M., Buló, S.R., Moreno, R., Kundu, M.K., Smedby, Ö., 2016. An efficient radiographic image retrieval system using convolutional neural network, in: 2016 23rd International Conference on Pattern Recognition (ICPR), IEEE. pp. 3134–3139.
- Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L., 2009. Imagenet: A large-scale hierarchical image database, in: 2009 IEEE conference on computer vision and pattern recognition, Ieee. pp. 248–255.
- Duan, J., Kuo, C.C.J., 2021. Bridging gap between image pixels and semantics via supervision: A survey. *arXiv preprint arXiv:2107.13757*.
- Dubey, S.R., 2021. A decade survey of content based image retrieval using deep learning. *IEEE Transactions on Circuits and Systems for Video Technology*.
- Guo, X., Gichoya, J.W., Purkayastha, S., Banerjee, I., 2021a. Cvad: A generic medical anomaly detector based on cascade vae. *arXiv preprint arXiv:2110.15811*.
- Guo, X., Gichoya, J.W., Trivedi, H., Purkayastha, S., Banerjee, I., 2021b. Medshift: identifying shift data for medical dataset curation. *arXiv preprint arXiv:2112.13885*.
- 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.
- Hwang, K.H., Lee, H., Choi, D., 2012. Medical image retrieval: past and present. *Healthcare informatics research* 18, 3–9.
- Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Ilcus, S., Chute, C., Marklund, H., Haghgoo, B., Ball, R., Shpanskaya, K., et al., 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison, in: Proceedings of the AAAI conference on artificial intelligence, pp. 590–597.
- Khosla, P., Teterwak, P., Wang, C., Sarna, A., Tian, Y., Isola, P., Maschinot, A., Liu, C., Krishnan, D., 2020. Supervised contrastive learning. *Advances in Neural Information Processing Systems* 33.
- Layode, O., Rahman, M., 2020. A chest x-ray image retrieval system for covid-19 detection using deep transfer learning and denoising auto encoder, in: 2020 International Conference on Computational Science and Computational Intelligence (CSCI), IEEE. pp. 1635–1640.
- Lloyd, S., 1982. Least squares quantization in pcm. *IEEE transactions on information theory* 28, 129–137.
- MacQueen, J., et al., 1967. Some methods for classification and analysis of multivariate observations, in: Proceedings of the fifth Berkeley symposium on mathematical statistics and probability, Oakland, CA, USA. pp. 281–297.
- Qayyum, A., Anwar, S.M., Awais, M., Majid, M., 2017. Medical image retrieval using deep convolutional neural network. *Neurocomputing* 266, 8–20.
- Rajpurkar, P., Irvin, J., Bagul, A., Ding, D., Duan, T., Mehta, H., Yang, B., Zhu, K., Laird, D., Ball, R.L., et al., 2017. Mura: Large dataset for abnormality detection in musculoskeletal radiographs. *arXiv preprint arXiv:1712.06957*.

- Revaud, J., Almazán, J., Rezende, R.S., Souza, C.R.d., 2019. Learning with average precision: Training image retrieval with a listwise loss, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 5107–5116.
- Sotomayor, C.G., Mendoza, M., Castañeda, V., Farías, H., Molina, G., Pereira, G., Härtel, S., Solar, M., Araya, M., 2021. Content-based medical image retrieval and intelligent interactive visual browser for medical education, research and care. *Diagnostics* 11, 1470.
- Sun, Y., Cheng, C., Zhang, Y., Zhang, C., Zheng, L., Wang, Z., Wei, Y., 2020. Circle loss: A unified perspective of pair similarity optimization, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6398–6407.
- Thorndike, R.L., 1953. Who belongs in the family? *Psychometrika* 18, 267–276.
- Wang, J., Song, Y., Leung, T., Rosenberg, C., Wang, J., Philbin, J., Chen, B., Wu, Y., 2014. Learning fine-grained image similarity with deep ranking, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1386–1393.
- Wang, X., Han, X., Huang, W., Dong, D., Scott, M.R., 2019. Multi-similarity loss with general pair weighting for deep metric learning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5022–5030.
- Zhong, A., Li, X., Wu, D., Ren, H., Kim, K., Kim, Y., Buch, V., Neumark, N., Bizzo, B., Tak, W.Y., et al., 2021. Deep metric learning-based image retrieval system for chest radiograph and its clinical applications in covid-19. *Medical Image Analysis* 70, 101993.