

# MVCO-DOT: MULTI-VIEW CONTRASTIVE DOMAIN TRANSFER NETWORK FOR MEDICAL REPORT GENERATION

Ruizhi Wang<sup>1</sup>, Xiangtao Wang<sup>1</sup>, Zhenghua Xu<sup>1,\*</sup>, Wenting Xu<sup>1</sup>, Junyang Chen<sup>2</sup>, Thomas Lukasiewicz<sup>3,4</sup>

<sup>1</sup>State Key Laboratory of Reliability and Intelligence of Electrical Equipment, School of Health Sciences and Biomedical Engineering, Hebei University of Technology, Tianjin, China

<sup>2</sup>College of Computer Science and Software Engineering and Guangdong Laboratory of Artificial Intelligence and Digital Economy (SZ), Shenzhen University, Shenzhen, China

<sup>3</sup>Institute of Logic and Computation, TU Wien, Vienna, Austria

<sup>4</sup>Department of Computer Science, University of Oxford, Oxford, United Kingdom

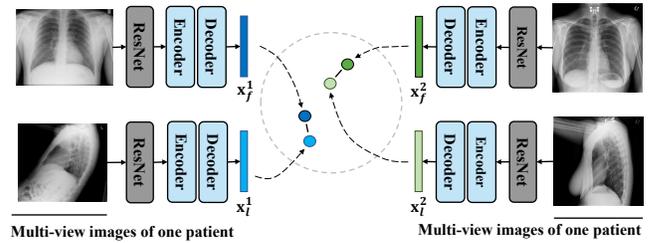
## ABSTRACT

In clinical scenarios, multiple medical images with different views are usually generated at the same time, and they have high semantic consistency. However, the existing medical report generation methods cannot exploit the rich multi-view mutual information of medical images. Therefore, in this work, we propose the first multi-view medical report generation model, called MvCo-DoT. Specifically, MvCo-DoT first propose a multi-view contrastive learning (MvCo) strategy to help the deep reinforcement learning based model utilize the consistency of multi-view inputs for better model learning. Then, to close the performance gaps of using multi-view and single-view inputs, a domain transfer network is further proposed to ensure MvCo-DoT achieve almost the same performance as multi-view inputs using only single-view inputs. Extensive experiments on the IU X-Ray public dataset show that MvCo-DoT outperforms the SOTA medical report generation baselines in all metrics.

**Index Terms**— Multi-view contrastive learning, Domain transfer, Medical report generation, Chest X-Ray

## 1. INTRODUCTION

Medical report generation is a multimodal cross-task in computer vision and natural language processing, which aims to reduce the workload of doctors by automatically generating diagnostic descriptions from medical images. Motivated by the application of deep learning in medical image analysis [1, 2], current medical report generation approaches typically utilize encoder-decoder architecture to learn medical images from different views but independently. And many spatial and channel attention methods [3] have been carefully designed to explore multimodal interactions between image-level and sentence-level semantic features for report generation [4–8], however, this way of understanding images is not ideal for complex cross-modal generation tasks. Then, [9] mitigates textual and visual data bias by exploring prior knowledge and posterior knowledge, [10] models and memorizes reports similar patterns between them, and thus promote Transformer to generate more informative long text interpretation reports. In order to utilize visual information more effectively, inspired by contrastive learning in the domain of



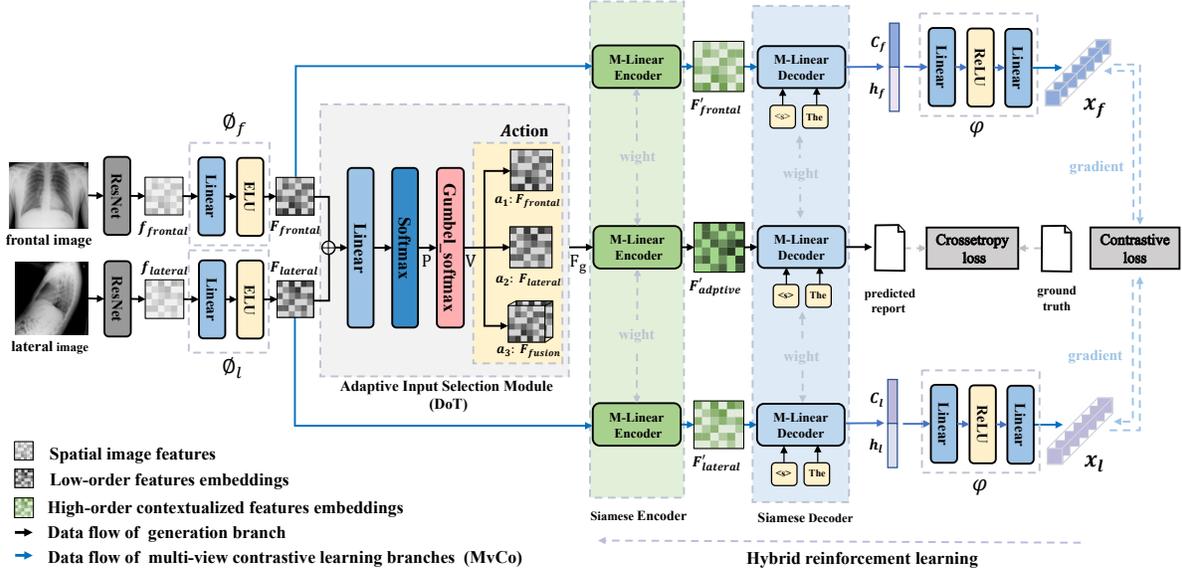
**Fig. 1.** Semantic-based multi-view contrastive learning.

natural images and natural language process [11–14], people started to try this method to improve the visual representation of medical images [15–17]. However, the existing medical report generation methods have a common shortcoming that their inputs are usually single view, which fails to utilize the rich multi-view information in chest X-Ray images.

Consequently, in this work, we propose a novel **Multi-view Contrastive Domain Transfer (MvCo-DoT)** network for better medical report generation. Specifically, we first propose to integrate a Multi-view Contrastive Learning (**MvCo**) strategy into our previous deep reinforcement learning based medical report generation model [7]. Intuitively, as shown in Fig. 1, since the paired multi-view medical images are different imaging results of the same patient, their descriptions of lesions or organs in the patient should have high consistency [18]. Therefore, MvCo is proposed to utilize semantic embeddings of different views of patients’ X-Ray images for contrastive learning. Compared to existing self-supervision based solutions [15, 16] whose contrastive learning modules are applied in encoders, feature representations used in MvCo is located in decoders, which thus have more direct impact on the quality of resulting medical reports.

In addition, our experimental studies show that although using MvCo can greatly improve the performances of medical report generation, it suffers from the problem of domain shift. When we have only single view X-Ray images as inputs, the inference results are greatly degraded because distributions of single-view inputs are very different from that of multi-view inputs. Consequently, we further propose to incorporate a Domain Transfer Network (**DoT**) into our medical report generation model to resolve this problem by closing the performance gaps of multi-view and single-view inputs. Specifically, DoT is achieved by using a sampling-based adaptive input selection module, which enables generation branch to randomly

\*Corresponding author: zhenghua.xu@hebut.edu.cn (Zhenghua Xu).



**Fig. 2.** The architecture of our proposed MvCo-DoT network.

select single or multi-view fused features as final input according to estimated probability. Advantages of DoT are as follows: (i) It ensures that the model learns using a more comprehensive input distribution, which thus closes performance gaps of using multi-view and single-view inputs in inference; (ii) different from random selection, using DoT will not degrade model’s feature learning capability; (iii) it also narrow information gap between contrastive learning branch (single-view input) and generation branch (multi-view input).

The contributions of this work can be summarized as follows. (i) We identify the lack of multi-view input problem of existing medical report generation methods and propose a MvCo-DoT network for better medical report generation. (ii) A multi-view contrastive learning (MvCo) strategy is first proposed to utilize the multi-view information of chest X-Ray images for better model learning, while a domain transfer network is then proposed to ensure model can achieve good performances using only single-view inputs in inference stage. (iii) Extensive experiments on a public dataset (IU X-Ray) show that: first, our proposed MvCo-DoT model greatly outperforms existing medical report generation baselines in all metrics; second, MvCo and DoT are both effective and essential for the model to achieve the superior performances; third, MvCo-DoT can achieve almost the same performance as multi-view inputs using only single-view inputs, which greatly saves the patients’ time and money.

## 2. METHODOLOGY

We propose a multi-view contrastive domain transfer network for medical report generation. As shown in Fig. 2, we will adopt the architecture of generation branch and contrastive learning branch. Contrastive learning branch is used for inter-view mutual information mining, and a generation branch is used for report generation. The two processes are alternately performed during training.

### 2.1. Multi-View Contrastive Learning

In order to mine and utilize the mutual information between medical images of different views, we propose semantic-based multi-view contrastive learning method with [7] as the backbone network. Specifically, we concatenate the

contextual semantic representations  $c_f$ ,  $c_l$  and hidden layer information  $h_f$ ,  $h_l$  decoded from different views and project them onto the same implicit space for comparison.

$$x_f = \psi(\text{Concat}(c_f, h_f)), x_l = \psi(\text{Concat}(c_l, h_l)), \quad (1)$$

where  $c_f$  and  $c_l$  are the context vectors of the last step of the LSTM in the twin M-Linear decoders,  $h_f$  and  $h_l$  are the corresponding hidden layer vectors.  $\psi(\cdot)$  is modeled as two fully connected layers with ReLU activations, according to [19]. Then we maximize the semantic concordance between the frontal and lateral view of the same patient while minimizing the similarity between different patients. The multi-view contrastive loss function is defined as

$$L_{MvCo} = -\log \frac{\exp(\text{sim}(x_l, x_f)/\tau_c)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq l]} \exp(\text{sim}(x_l, x_k)/\tau_c)}, \quad (2)$$

where  $\text{sim}(x_l, x_k)$  represents the cosine similarity,  $\text{sim}(x_l, x_k) = x_l^\top x_k / (\|x_l\| \|x_k\|)$  and  $\tau_c$  is temperature parameter.

### 2.2. Domain Transfer Network

To overcome the domain shift problem caused by the input distribution gap between training and testing phases, We propose a domain transfer network based on adaptive input selection. First, we project the frontal and lateral visual features  $f_{\text{frontal}}$  and  $f_{\text{lateral}}$  extracted by ResNet101 [20] into two different latent spaces, to obtain more discriminative visual embeddings  $F_{\text{frontal}}$  and  $F_{\text{lateral}}$  for different views.

$$F_{\text{frontal}} = \phi_f(f_{\text{frontal}}), F_{\text{lateral}} = \phi_l(f_{\text{lateral}}), \quad (3)$$

where  $\phi_f(\cdot)$  and  $\phi_l(\cdot)$  are modeled as fully connected layers with ELU activations. Afterward, we add the visual embeddings of these different views to obtain  $F_{\text{fusion}}$ , which replaces the operation of direct concatenation of original features commonly used in previous work. Then, in order to make the model get the most useful information input and better balance the use of frontal and lateral view information, we adaptively decide to input a single feature or mixed features through action sampling. The action space  $A \in \mathbb{R}^{1 \times 3}$  is defined as  $F_{\text{frontal}}$  (when  $i = 0$ ),  $F_{\text{lateral}}$  (when  $i = 1$ ), or  $F_{\text{fusion}}$  (when  $i = 2$ ).

Image	Ground-truth	MRMA	X-LAN	HRe-MR	MvCo-DoT(Ours)
	“no acute <b>cardiopulmonary</b> abnormality. the lungs are clear and without <b>focal airspace opacity</b> . the <b>cardiomediastinal silhouette</b> is normal in size and <b>contour</b> and stable. there is no <b>pneumothorax</b> or large <b>pleural effusion</b> .”	“no acute <b>cardiopulmonary</b> abnormality . the lungs are clear . there is no <b>pneumothorax</b> or <b>pleural effusion</b> . the heart and mediastinum are within normal limits . bony structures are intact.”	“no acute <b>cardiopulmonary</b> findings . lungs are clear bilaterally . cardiac and <b>mediastinal silhouettes</b> are normal . pulmonary vasculature is normal . no <b>pneumothorax</b> or <b>pleural effusion</b> . no acute bony abnormality.”	“no acute <b>cardiopulmonary</b> abnormality . the <b>cardiomediastinal silhouette</b> and pulmonary vasculature are normal in size . the lungs are clear . there is no <b>pneumothorax</b> or <b>pleural effusion</b> . no <b>focal airspace consolidation</b> . no acute bony findings.”	“no acute <b>cardiopulmonary</b> abnormality . the <b>cardiomediastinal silhouette</b> is within normal limits and <b>contours</b> are stable. the lungs are clear . there is no <b>focal airspace opacity</b> . no <b>pleural effusion</b> or <b>pneumothorax</b> . there are no acute bony abnormality.”

Fig. 3. Example of reports generated by our MvCo-DoT model and baselines.

Table 1. Results of MvCo-DoT and the SOTA baselines on IU X-Ray, where B, ME and RO stand for BLEU, METEOR and ROUGE-L, and all models are re-implemented by us.

Model	B-1	B-2	B-3	B-4	ME	RO
Top-down [4]	0.2822	0.1866	0.1241	0.0830	0.1455	0.3330
RTMIC [22]	0.3448	0.2188	0.1484	0.1063	0.1509	0.2890
MRMA [5]	0.3820	0.2520	0.1730	0.1200	0.1630	0.3090
X-LAN [6]	0.3826	0.2724	0.1949	0.1405	0.1750	0.3441
R2Gen [10]	0.4349	0.2802	0.1868	0.1510	0.1773	0.3509
HRe-MR [7]	0.4265	0.3025	0.2119	0.1502	0.1871	0.3608
<b>MvCo-DoT</b>	<b>0.4533</b>	<b>0.3180</b>	<b>0.2228</b>	<b>0.1568</b>	<b>0.1958</b>	<b>0.3743</b>

Table 2. Results of ablation experiments on IU X-Ray.

Model	B-1	B-2	B-3	B-4	ME	RO
Base-Cat	0.4175	0.2813	0.1951	0.1400	0.1820	0.3604
MvCo-Cat	0.4373	0.3062	0.2139	0.1482	0.1933	0.3609
MvCo-Fusion	0.4440	0.3130	0.2196	<b>0.1571</b>	0.1953	0.3698
<b>Ours</b>	<b>0.4533</b>	<b>0.3180</b>	<b>0.2228</b>	0.1568	<b>0.1958</b>	<b>0.3743</b>

This non-deterministic approach enables the model to adaptively select optimal input to obtain the maximum amount of visual information for each image. To circumvent the technical problem that binary sampling actions can’t participate in backpropagation, we utilize random sampling based on *Gumbel-Softmax* distribution. This reparameterization trick has been used in reinforcement learning to enable discrete decision [21]. Non-differentiable action values are replaced by differentiable samples from *Gumbel-Softmax* distribution. Specifically, we concatenate the global visual features  $F_{frontal}$  and  $F_{lateral}$ , which are multi-scale fusions of frontal and lateral views in the latent space and taken as a comprehensive information basis for the current action selection. It is sent to a linear layer through the fully connected layer to obtain the action confidence warehouse  $P \in \mathbb{R}^{1 \times 3}$ .

$$P = \text{softmax}(W_c(\text{Concat}(F_{frontal}, F_{lateral}))), \quad (4)$$

where  $W_c$  represents the fully connected layer parameter matrix. Subsequently, the sampling module will generate action values  $V \in \mathbb{R}^{1 \times 3}$ , which defined as

$$V(a) = \frac{\exp(\log(P_i(a)) + g_i(a)/\tau_s)}{\sum_{j=1}^3 \exp(\log(P_j(a)) + g_j(a)/\tau_s)}, \text{ for } i = 1, 2, 3, \quad (5)$$

where  $g$  represents the noise sampled from standard *Gumbel-Softmax* distribution, and  $\tau_s$  is temperature parameter. The final input strategy is gained after  $V$  through *argmax* layer. During inference,  $V$  is generated according to input directly. Sample action whose sample value in  $A$  is calculated to be 1, and reconstruct only the features corresponding to the action into final input feature  $F_g$ .

$$F_g = A(a_i), V(a_i) = 1, \quad (6)$$

### 3. EXPERIMENTS

#### 3.1. Experimental Settings

To evaluate the performance of our proposed MvCo-DoT, extensive experiments are conducted on public chest X-Ray image dataset IU X-Ray [23]. We screen 3,111 groups of cases from the dataset, each containing two X-Ray images of frontal and lateral views and a paired report. The dataset is randomly divided by 7:1:2. In addition, words with frequency of less than 5 are discarded and replaced with "UNK", and reported maximum generated length is set to 114. We reimplement six state-of-the-art image captioning and medical report generation models as baselines, including Top-down [4], RTMIC [22], MRMA [5], X-LAN [6], R2Gen [10] and HRe-MR [7]. We evaluated the models using six common automatic language generation metrics, including BLEU [24], METEOR [25], and ROUGE-L [26], where BLEU includes four n-gram-based metrics (BLEU-1 to BLEU-4).

We utilize ResNet-101 pre-trained on ImageNet [27] to extract 2048 dimensional region-level image features from the last convolutional layer. After being converted to visual embeddings of size 1024, the encoder exploration with four stacks of M-linear attention blocks yields high-order synthetic features. During decoding process, we set size of hidden layer, word embedding dimension, and latent dimension of the projection layer to 1024. During training, we first pre-train model with a batch size of 6 for 60 epochs using NVIDIA RTX 2080Ti GPUs, the model is optimized alternately by generation branch and contrastive learning branch. We set the base learning rate to 0.0001, paired with a Norm decay strategy with 10,000 warm-up steps, and used the ADAM [28] optimizer. We set  $\tau_c$  to 0.1 and  $\tau_s$  to 0.3. Finally, we train model with batch size of 2 for 60 epochs of reinforcement learning using beam search [29] with a beam size of 2 to further improve model performance. We set the indicator-weighted mixed reward as our training reward [7], where weights of BLEU-1, BLEU-2, BLEU-3, BLEU-4, METEOR, and ROUGE-L, are 2, 2, 1, 1, 2, and 2, respectively; and base learning rate is reduced to 0.00001 and decayed by cosine annealing with a period of 15 epochs.

#### 3.2. Main Results

Table 1 shows experimental results of our proposed MvCo-DoT and five baselines on six natural language generation metrics, where all baselines are re-implemented. Furthermore, Fig. 3 presents some examples of generated reports.

In general, MvCo-DoT outperforms all state-of-the-art baselines among all metrics in Table 1, because (i) our multi-view contrastive learning adequately performs multi-view mutual information learning to obtain superior performance, (ii) the same input for multi-view training and single-view

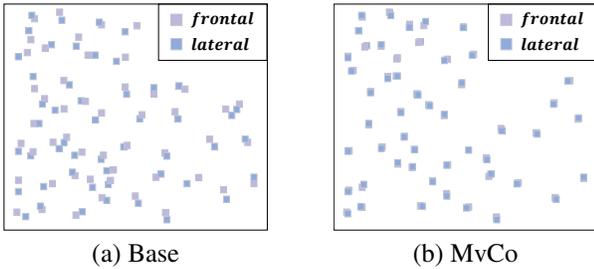


Fig. 4. Semantic feature embeddings in latent space.

testing is maintained, and task gap between contrastive learning branch and generation branch is narrowed, avoiding domain shift problem. Moreover, visualization results shown in Fig. 3 also support our findings, where MvCo-DoT generation obtains a more comprehensive and accurate report description than baselines. Thus, our proposed multi-view contrastive learning and domain transfer network are highly effective in enhancing the quality of report generation.

### 3.3. Ablation Study

We further conduct a series of ablation experiments to demonstrate the effectiveness of each module of our proposed MvCo-DoT. We take generation branch as the base model, which utilizes raw feature concatenation of different views as input, called Base-Cat. Then we further implement two other versions of our MvCo-DoT: (i) Introduce a multi-view contrastive learning branch on base model, called MvCo-Cat, and (ii) using visual embedding fusion instead of original feature direct concatenation, called MvCo-Fusion. In Table 2, we observe that all metrics of MvCo-Cat are superior higher than Base-Cat, as mutual information mined by multi-view contrastive learning helps to focus on salient lesions, explore deep semantic features and enable multimodal inference. The MvCo-Fusion score is further improved, indicating that feature fusion strategy is more suitable for our task. Next, we use *Gumbel-Softmax*-based random sampling to adaptively select input strategy to obtain final model MvCo-DoT. Overcoming input distribution gap enables model to improve cross-domain transferability while reducing the task distance between the contrastive learning branch and generation branch, which makes it the highest score among all versions of the model. In conclusion, our proposed MvCo learning strategy and DoT network are very effective and essential to improve the accuracy of automatic generation of medical image reports.

### 3.4. Additional Results

In this section, we investigate the impact of semantic-based multi-view contrastive learning on model performance and advantages of domain transfer networks based on adaptive inputs. Fig. 4 shows the distance variation of frontal and lateral view semantic embeddings in the implicit space. Compared with generation branch Base in (a), multi-view contrastive learning MvCo in (b) can make the semantic embeddings of different views closer. This is because the model learns mutual information between different views, thereby decoding feature vectors with more semantic consistency. In addition, in Fig. 5, compared with the pure multi-view comparison MvCo in (a), domain transfer network MvCo-DoT with adaptive input selection in (b) can generate high-scoring reports under any single view, which well solves domain shift caused by different input distributions during multi-view training and single-view testing. Moreover, the model performance is further improved compared to (a), which means that the same

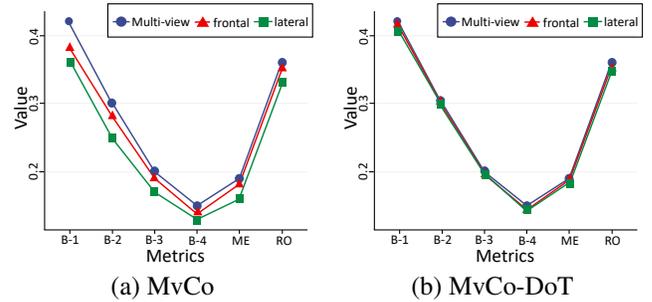


Fig. 5. Comparison of the performance of MvCo and MvCo-DoT for reporting inference using different inputs.

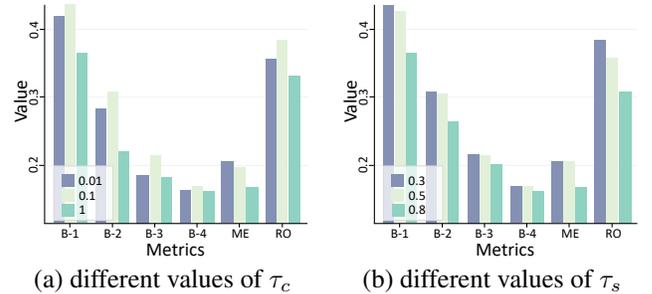


Fig. 6. Results of varying hyperparameters  $\tau_c$  and  $\tau_s$ .

input distribution also reduces the gap. The model obtains the optimal representation of the two, promoting each other.

### 3.5. Effect of Varying Hyper-Parameters $\tau_c$ and $\tau_s$

$\tau_c$  and  $\tau_s$  are the temperature parameter of the contrastive loss and *Gumbel-softmax* distribution, respectively. In Fig. 6, according to (a), in the range of 0.01 to 1, the performance of the model fluctuates with the size of  $\tau_c$ , and according to (b), smaller  $\tau_s$  in the range of 0.3 to 0.8 is better. Therefore, we need certain tuning parameters to achieve the best performance. In this model, we set  $\tau_c$  to 0.1 and  $\tau_s$  to 0.3.

## 4. CONCLUSIONS

In this paper, we proposed a multi-view contrastive domain transfer network (MvCo-DoT) for medical report generation. We mined mutual information between different views of chest X-Ray using multi-view contrastive learning based on semantic information to aid model learning. We also closed the input distribution gap between training and inference stages and contrastive learning branch and generation branch through domain transfer network based on adaptive input selection to address the domain shift problem. We performed extensive experiments on the publicly available dataset IU X-Ray, demonstrating the superiority and effectiveness of our proposed method. Future researches may include the usage of the proposed method in imbalanced learning task [30], and incorporate more deep reinforcement learning techniques [31] to enhance the deep generation model’s learning capabilities.

## 5. ACKNOWLEDGMENTS

This work was supported by the National Natural Science Foundation of China under the grants 62276089, 61906063 and 62102265, by the Natural Science Foundation of Hebei Province, China, under the grant F2021202064, by the “100 Talents Plan” of Hebei Province, China, under the grant E2019050017, by the Open Research Fund from Guangdong Laboratory of Artificial Intelligence and Digital Economy

(SZ) under the grant GML-KF-22-29, and by the Natural Science Foundation of Guangdong Province of China under the grant 2022A1515011474.

## 6. REFERENCES

- [1] Zhenghua Xu, Chang Qi, and Guizhi Xu, "Semi-supervised attention-guided CycleGAN for data augmentation on medical images," in *Proceedings of IEEE BIBM*, 2019, pp. 563–568.
- [2] Zhenghua Xu, Tianrun Li, Yunxin Liu, Yuefu Zhan, Junyang Chen, and Thomas Lukasiewicz, "PAC-Net: Multi-pathway FPN with position attention guided connections and vertex distance IoU for 3D medical image detection," *Frontiers in Bioengineering and Biotechnology*, vol. 11, pp. 1049555, 2023.
- [3] Zhenghua Xu, Shijie Liu, Di Yuan, Lei Wang, Junyang Chen, Thomas Lukasiewicz, Zhigang Fu, and Rui Zhang, " $\omega$ -net: Dual supervised medical image segmentation with multi-dimensional self-attention and diversely-connected multi-scale convolution," *Neurocomputing*, vol. 500, pp. 177–190, 2022.
- [4] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang, "Bottom-up and top-down attention for image captioning and visual question answering," in *Proceedings of CVPR*, 2018, pp. 6077–6086.
- [5] Yuan Xue, Tao Xu, L Rodney Long, Zhiyun Xue, Sameer Antani, George R Thoma, and Xiaolei Huang, "Multimodal recurrent model with attention for automated radiology report generation," in *Proceedings of MICCAI*, 2018, pp. 457–466.
- [6] Yingwei Pan, Ting Yao, Yehao Li, and Tao Mei, "X-linear attention networks for image captioning," in *Proceedings of CVPR*, 2020, pp. 10971–10980.
- [7] Wenting Xu, Zhenghua Xu, Junyang Chen, Chang Qi, and Thomas Lukasiewicz, "Hybrid reinforced medical report generation with m-linear attention and repetition penalty," *arXiv preprint arXiv:2210.13729*, 2022.
- [8] Benjamin Hou, Georgios Kaissis, Ronald M Summers, and Bernhard Kainz, "Ratchet: Medical transformer for chest x-ray diagnosis and reporting," in *Proceedings of MICCAI*, 2021, pp. 293–303.
- [9] Fenglin Liu, Xian Wu, Shen Ge, Wei Fan, and Yuexian Zou, "Exploring and distilling posterior and prior knowledge for radiology report generation," in *Proceedings of CVPR*, 2021, pp. 13753–13762.
- [10] Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan, "Generating radiology reports via memory-driven transformer," *arXiv preprint arXiv:2010.16056*, 2020.
- [11] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton, "A simple framework for contrastive learning of visual representations," in *Proceedings of ICML*, 2020, pp. 1597–1607.
- [12] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of CVPR*, 2020, pp. 9729–9738.
- [13] Yonglong Tian, Dilip Krishnan, and Phillip Isola, "Contrastive multiview coding," in *Proceedings of ECCV*, 2020, pp. 776–794.
- [14] Tianyu Gao, Xingcheng Yao, and Danqi Chen, "Simcse: Simple contrastive learning of sentence embeddings," *arXiv preprint arXiv:2104.08821*, 2021.
- [15] An Yan, Zexue He, Xing Lu, Jiang Du, Eric Chang, Amilcare Gentili, Julian McAuley, and Chun-Nan Hsu, "Weakly supervised contrastive learning for chest x-ray report generation," *arXiv preprint arXiv:2109.12242*, 2021.
- [16] Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz, "Contrastive learning of medical visual representations from paired images and text," in *Proceedings of MLHC*, 2022, pp. 2–25.
- [17] Shuo Zhang, Jiaojiao Zhang, Biao Tian, Thomas Lukasiewicz, and Zhenghua Xu, "Multi-modal contrastive mutual learning and pseudo-label re-learning for semi-supervised medical image segmentation," *Medical Image Analysis*, vol. 83, pp. 102656, 2023.
- [18] Yen Nhi Truong Vu, Richard Wang, Niranjana Balachandrar, Can Liu, Andrew Y Ng, and Pranav Rajpurkar, "Medaug: Contrastive learning leveraging patient metadata improves representations for chest x-ray interpretation," in *Proceedings of MLHC*, 2021, pp. 755–769.
- [19] Vinod Nair and Geoffrey E Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proceedings of ICML*, 2010, pp. 807–814.
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of CVPR*, 2016, pp. 770–778.
- [21] Eric Jang, Shixiang Gu, and Ben Poole, "Categorical reparameterization with gumbel-softmax," *arXiv preprint arXiv:1611.01144*, 2016.
- [22] Yuxuan Xiong, Bo Du, and Pingkun Yan, "Reinforced transformer for medical image captioning," in *Proceedings of MICCAI*, 2019, pp. 673–680.
- [23] Dina Demner-Fushman, Marc D Kohli, Marc B Rosenman, Sonya E Shooshan, Laritza Rodriguez, Sameer Antani, George R Thoma, and Clement J McDonald, "Preparing a collection of radiology examinations for distribution and retrieval," *Journal of the American Medical Informatics Association*, vol. 23, no. 2, pp. 304–310, 2016.
- [24] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of ACL*, 2002, pp. 311–318.
- [25] Satyanjeev Banerjee and Alon Lavie, "Meteor: An automatic metric for mt evaluation with improved correlation with human judgments," in *Proceedings of ACL*, 2005, pp. 65–72.
- [26] Chin-Yew Lin, "Rouge: A package for automatic evaluation of summaries," in *Proceedings of ACL*, 2004, pp. 74–81.
- [27] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proceedings of CVPR*, 2009, pp. 248–255.
- [28] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [29] Ashwin K Vijayakumar, Michael Cogswell, Ramprasad R Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra, "Diverse beam search: Decoding diverse solutions from neural sequence models," *arXiv preprint arXiv:1610.02424*, 2016.
- [30] Jianfeng Wang, Thomas Lukasiewicz, Xiaolin Hu, Jianfei Cai, and Zhenghua Xu, "RSG: A simple but effective module for learning imbalanced datasets," in *Proceedings of CVPR*, 2021, pp. 3784–3793.
- [31] Di Yuan, Yunxin Liu, Zhenghua Xu, Yuefu Zhan, Junyang Chen, and Thomas Lukasiewicz, "Painless and accurate medical image analysis using deep reinforcement learning with task-oriented homogenized automatic pre-processing," *Computers in Biology and Medicine*, vol. 153, pp. 106487, 2023.