

ICON: Improving Inter-Report Consistency in Radiology Report Generation via Lesion-aware Mixup Augmentation

Anonymous ACL submission

Abstract

Previous research on radiology report generation has made significant progress in terms of increasing the clinical accuracy of generated reports. In this paper, we emphasize another crucial quality that it should possess, i.e., *inter-report consistency*, which refers to the capability of generating consistent reports for semantically equivalent radiographs. This quality is even of greater significance than the overall report accuracy in terms of ensuring the system’s credibility, as a system prone to providing conflicting results would severely erode users’ trust. Regrettably, existing approaches struggle to maintain inter-report consistency, exhibiting biases towards common patterns and susceptibility to lesion variants. To address this issue, we propose ICON, which Improves the inter-report CONSistency of radiology report generation. Aiming at enhancing the system’s ability to capture the similarities in semantically equivalent lesions, our approach involves first extracting lesions from input images and examining their characteristics. Then, we introduce a lesion-aware mixup technique to ensure that the representations of the semantically equivalent lesions align with the same attributes, by linearly interpolating them during the training phase. Extensive experiments on three publicly available chest X-ray datasets verify the effectiveness of our approach, both in terms of improving the consistency and accuracy of the generated reports¹.

1 Introduction

Being part of the diagnostic process, radiology report generation (Shin et al., 2016; Zhang et al., 2017; Jing et al., 2018) has garnered significant attention within the research community, due to its large potential to alleviate the heavy strain of radiologists. Recent research (Nishino et al., 2022;

¹We will release our codes and model checkpoints after the review process.

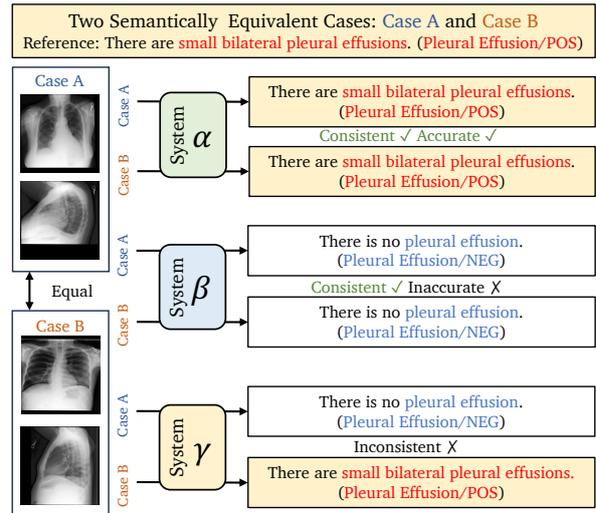


Figure 1: Given two semantically equivalent cases (i.e., Case A and Case B), an example to illustrate the difference between three radiology report generation systems: a consistent and accurate system (i.e., System α) and a consistently inaccurate system (i.e., System β), and an inconsistent system (i.e., System γ).

Tanida et al., 2023; Hou et al., 2023b) has made noteworthy progress in enhancing the clinical accuracy of the generated reports.

However, constructing a credible report generation system goes beyond the overall accuracy. There is another crucial quality for report generation systems that has been largely overlooked in the existing literature of medical report generation, which is, *inter-report consistency* (Elazar et al., 2021). To illustrate the disparity between accuracy and inter-report consistency, we exemplify two semantically equivalent cases as shown in Figure 1. Specifically, System α demonstrates the ability to maintain both inter-report consistency and factual accuracy for two similar cases (i.e., "small bilateral pleural effusions" for positive *Pleural Effusion*), whereas other systems (i.e., β and γ) fail to meet these criteria. These systems might have overfitted to ordinary cases and could be vulnera-

ble to noise or attack. In terms of enhancing the system’s credibility, inter-report consistency might even hold greater significance than the overall accuracy, since a system prone to providing conflicting results would severely undermine users’ trust (Qayyum et al., 2020; Asan et al., 2020). Regrettably, existing report generation systems struggle to maintain this important quality. They tend to exhibit biases towards common patterns, primarily describing normal observations and are extremely susceptible to lesion variants and context noise (Chen et al., 2020; Qin and Song, 2022; Ma et al., 2021; Kaviani et al., 2022). We argue that this is largely due to their limited capability of capturing shared attributes of similar patterns, which arises from the data scarcity of distributed lesions and their semantically equivalent variants, rendering it challenging for neural models to accurately locate and describe abnormalities.

In this paper, we propose ICON, which aims to Improve inter-report CONSistency of radiology report generation. Our proposed method involves first extracting lesions from given input images, followed by examining the attributes of these lesions. Subsequently, both the radiographs and their associated attributes are utilized as inputs for report generation. To further enhance the inter-report consistency, we introduce a lesion-aware mixup technique by learning from linearly interpolated lesions and attributes that belong to the same observation. In summary, the contributions of this paper are as follows:

- To the best of our knowledge, we are the first to introduce *inter-report consistency* in radiology report generation. To this end, we have devised two metrics (CON and R-CON) to measure such consistency.
- We propose ICON, which improves both the *consistency* and *accuracy* in radiology report generation by capturing abnormalities at the region level. ICON only requires coarse-grained labels (i.e., image labels) for training to extract lesions², in contrast to previous methods that require fine-grained labels (i.e., bounding boxes).
- Extensive experiments are conducted on three

²In this context, the term "lesion" generally refers to a specific abnormality. It encompasses most observation categories, excluding *Support Devices*, *Cardiomegaly*, and *Enlarged Cardiomediastinum*. For simplicity, we consider all corresponding regions as lesions.

publicly available datasets, and the results demonstrate the effectiveness of ICON in terms of improving both the consistency and accuracy of the generated reports.

2 Preliminaries

2.1 Problem Formulation

Given a set of radiographs $\mathcal{X} = \{X_1, \dots, X_L\}$ in one study, along with its historical records $\mathcal{X}^p = \{X_1^p, \dots, X_{|p|}^p\}$ or $\mathcal{X}^p = \emptyset$, and its report $\mathcal{Y} = \{y_1, \dots, y_T\}$, the task of radiology report generation (RRG) is formulated as $p(\mathcal{Y}|\mathcal{X}, \mathcal{X}^p)$. We elaborate on the justification of using the historical records as context in Appendix A.8. Our proposed method, denoted as ICON, decomposes the RRG task into two stages: Lesion Extraction (Stage 1) and Report Generation (Stage 2). Specifically, given the input images \mathcal{X} , ICON first extracts lesions $\mathcal{Z} = \{Z_1, \dots, Z_{|O|}\}$ from \mathcal{X} , where the probability of a region $R_{i,j}$ from image X_i being identified as a lesion Z_k is estimated as $p(Z_k|X_i)$. Subsequently, in Stage 2, ICON generates a report based on both the input images and the extracted lesions, modeled as $P(\mathcal{Y}|\mathcal{X}, \mathcal{X}^p, \mathcal{Z})$. Finally, our framework aims to maximize the following probability:

$$P(\mathcal{Y}|\mathcal{X}, \mathcal{X}^p) \propto \underbrace{p(\mathcal{Z}|\mathcal{X})}_{\text{Stage 1}} \cdot \underbrace{P(\mathcal{Y}|\mathcal{X}, \mathcal{X}^p, \mathcal{Z})}_{\text{Stage 2}}.$$

2.2 Observation and Attribute Annotation

Observations for Lesion Extraction. Lesion extraction requires report-level labels, and we adopt CheXbert (Smit et al., 2020) for this purpose. Specifically, CheXbert annotates a report with 14 observation categories $O = \{o_1, \dots, o_{14}\}$ (refer to Appendix A.1 for data statistics). Each observation is assigned one of four statuses: *Present*, *Absent*, *Uncertain*, and *Blank*. During training and evaluation, *Present* and *Uncertain* are merged into the *Positive* category, which represents abnormal observations. Note that for the observation category *No Finding*, only two statuses, *Present* or *Absent*, are applicable. Finally, observation information is utilized for lesion extraction as described in §3.2.

Attributes for Lesion-Attribute Alignment. After extracting observations, we further extract entities that represent their characteristics. Specifically, we adopt the attributes released by Hou et al. (2023a)³, which are entities (with a relation *mod-*

³The attributes are available at <https://github.com/wjhou/Recap>.

ify or *located_at*) extracted from RadGraph (Jain et al., 2021) using PMI (Church and Hanks, 1990). We select the top 30 attributes for each observation and list some of them in Appendix A.2 for a better understanding. These attributes are then utilized for lesion-attribute alignment as described in §3.3.

2.3 Inter-Report Consistency Metrics

To assess the inter-report consistency of a model, we introduce two metrics, CON and R-CON, inspired by Elazar et al. (2021). Semantically equivalent samples should have high observation and entity similarity, which we calculate using the Overlap Coefficient (Simpson, 1943): $\text{Overlap}(A, B) = \frac{|A \cap B|}{\min(|A|, |B|)}$. For a report Q_i and its semantically equivalent samples $\mathcal{K}_i = \{K_{i,1}, \dots, K_{i,N}\}$, the observation similarity should meet $\text{Overlap}(O_{Q_i}, O_{K_{i,j}}) \geq 0.75$ and the entity similarity should meet $\text{Overlap}(Q_i, K_{i,j}) \geq 0.5$. We collect the corresponding outputs of \mathcal{K}_i from a model, denoted as $\hat{\mathcal{K}}_i = \{\hat{K}_{i,1}, \dots, \hat{K}_{i,N}\}$. The similarity between two outputs \hat{Q}_i and $\hat{K}_{i,j}$ is:

$$\text{Overlap}(\hat{Q}_i, \hat{K}_{i,j}) = \frac{|\hat{e}_i \cap \hat{e}_j|}{\min(|\hat{e}_i|, |\hat{e}_j|)},$$

where \hat{e}_i and \hat{e}_j are entities and attributes in \hat{Q}_i and $\hat{K}_{i,j}$ (mentioned in §2.2), respectively. The inter-report consistency is then defined as:

$$\text{CON}(\hat{Q}_i, \hat{\mathcal{K}}_i) = \frac{1}{N} \sum_{j=1}^N \text{Overlap}(\hat{Q}_i, \hat{K}_{i,j}).$$

Since CON only considers inter-report consistency without accounting for the reference quality, we introduce R-CON to consider both consistency and accuracy:

$$\text{R-CON}(\hat{Q}_i, \hat{\mathcal{K}}_i) = \tau_i \cdot \text{CON}(\hat{Q}_i, \hat{\mathcal{K}}_i),$$

where $\tau_i = \text{Overlap}(\hat{Q}_i, Q_i)$ is the similarity between the hypothesis and its reference.

3 Methodology

3.1 Visual Encoding

Given an image X_l , an image processor is first utilized to split X_l into N patches. Then, a visual encoder f_θ (e.g., Swin Transformer (Liu et al., 2021d)) is employed to extract visual representations \mathbf{X}_l and the pooler output $\mathbf{P}_l \in \mathbb{R}^h$:

$$[\mathbf{P}_l, \mathbf{X}_l] = f_\theta(X_l),$$

where $\mathbf{X}_l = \{\mathbf{x}_{l,i}, \dots, \mathbf{x}_{l,N}\}$ and $\mathbf{x}_{l,i} \in \mathbb{R}^h$ is the i -th visual representation.

3.2 Stage 1: Extracting Lesions via Observation Classification (ZOOMER)

Observation Classification. A ZOOMER is a visual encoder parameterized by θ_Z and trained to classify a given input \mathcal{X} into abnormal observations as mentioned in §2.2:

$$p(o_i) = \text{ZOOMER}(\mathcal{X}).$$

Specifically, ZOOMER first encodes images $\mathcal{X} = \{X_1, \dots, X_L\}$ as outlined in §3.1, and then takes the averaged pooler output for classification, following these steps:

$$[\mathbf{P}_l, \mathbf{X}_l] = f_{\theta_Z}(X_l),$$

$$\mathbf{P} = \frac{1}{L} \sum \mathbf{P}_l,$$

$$p(o_i) = \sigma(\mathbf{W}_i \mathbf{P} + b_i),$$

where $\mathbf{W}_i \in \mathbb{R}^h$ is the weight for the i -th observation, $b_i \in \mathbb{R}$ is its bias, and σ is the Sigmoid function.

Zooming In for Lesion Extraction. Upon completing training ZOOMER, we can use it to extract lesions without the need for object detectors (Ren et al., 2015). It is worth noting that our method does not require fine-grained labels, such as bounding boxes (Tanida et al., 2023), making it easily adaptable to other modalities, e.g., FFA images (Li et al., 2021).

For an image X_l , a sliding window with a 0.375 ratio of X_l is applied to extract M region candidates $\mathcal{R}_l = \{R_{l,1}, \dots, R_{l,M}\}$ from X_l , as shown in the left side of Figure 2. These regions are then sequentially fed into ZOOMER for classification. Further details on the extraction of these regions can be found in Appendix A.6. The probability of a region $R_{l,j}$ being classified as an abnormal observation o_i is:

$$p_{l,j}(o_i) = \text{ZOOMER}(R_{l,j}).$$

For each study, all images in \mathcal{X} are iterated, and only the region with the highest $p_{l,j}(o_i)$ is chosen as a lesion Z_i corresponding to the observation o_i . Finally, the set of lesions is denoted as $\mathcal{Z} = \{Z_1, \dots, Z_{|\mathcal{O}|}\}$.

Training ZOOMER. ZOOMER is optimized using the binary cross-entropy (BCE) loss. To handle the class-imbalanced issue (refer to Appendix A.1 for details), a weight factor α_j is applied for each abnormal observation, and the loss function \mathcal{L}_{S1} is:

$$\begin{aligned} \text{BCE}(p(o_j), o_j) = & -\frac{1}{|\mathcal{O}|} \sum_j [\alpha_j \cdot o_j \cdot \log p(o_j) \\ & + (1 - o_j) \cdot \log(1 - p(o_j))], \end{aligned}$$

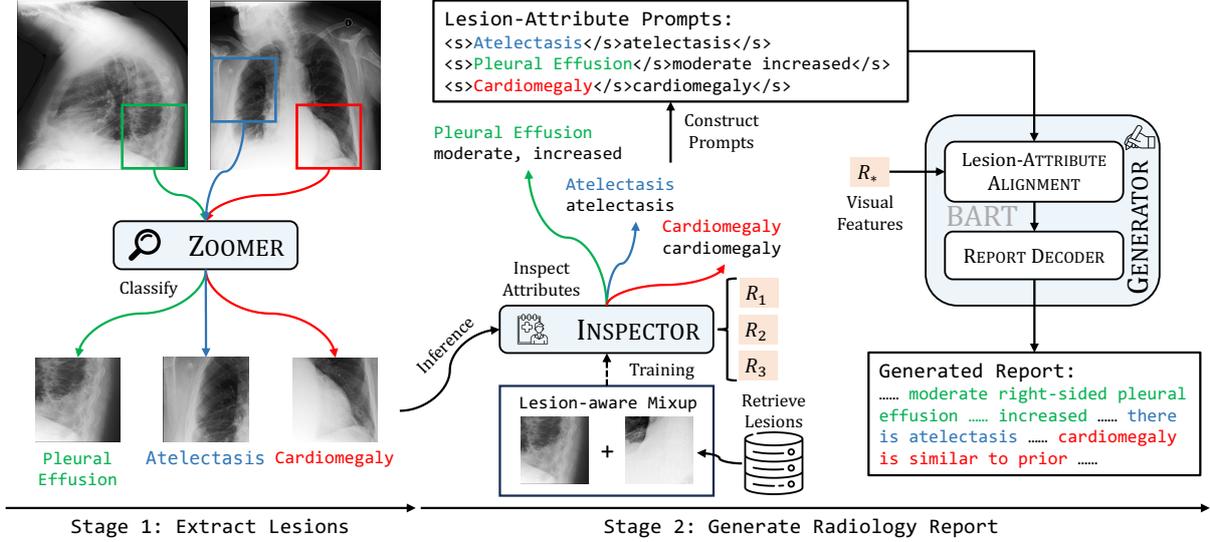


Figure 2: Overview of the ICON framework, which first extracts lesions and then generates reports. Attributes are extracted from RadGraph (Jain et al., 2021).

where $o_j \in \{0, 1\}$ is the label, $\alpha_j = 1 + \log\left(\frac{|\mathcal{D}_{\text{train}}| - w_j}{w_j}\right)$, and $|\mathcal{D}_{\text{train}}|$ and w_j are the number of samples and the number of j -th observations in the training set, respectively.

3.3 Stage 2: Inspecting Lesions (INSPECTOR)

Inspecting Lesions with Attributes. Given that lesions of the same observation can exhibit different characteristics, it is crucial to inspect each lesion and match it with corresponding attributes (§2.2) to differentiate it from other variations. Specifically, an INSPECTOR is a visual encoder parameterized by θ_I , similar to §3.2. $\text{INSPECTOR}(\mathbf{P}^p, \mathbf{P}, Z_j)$ takes prior and current visit chest X-rays as context, along with a lesion region as input:

$$[\mathbf{P}_{Z_j}, \mathbf{Z}_j] = f_{\theta_I}(Z_j),$$

$$p_j(a_k) = \sigma(\text{MLP}(\mathbf{P}^p, \mathbf{P}, \mathbf{P}_{Z_j})),$$

where MLP is a two-layer perceptron with non-linear activation, and $\mathbf{P}^p, \mathbf{P}, \mathbf{P}_{Z_j} \in \mathbb{R}^h$ are pooler outputs of prior images, current images, and the lesion, respectively. The lesion features $\mathbf{Z} = \{Z_1, \dots, Z_{|O|}\}$ are then collected for report generation. For image encoding, we use another visual encoder f_{θ_V} to encode \mathcal{X} into \mathcal{X} and \mathcal{X}^p into \mathcal{X}^p . By inspecting lesion-level features, ICON can capture fine-grained details which are beneficial for generating consistent outputs.

Lesion-aware Mixup. To further improve the consistency of the generated outputs, we adopt the mixup augmentation method (Zhang et al., 2018) and devise a Lesion-aware mixup during the training phase. Specifically, for a lesion-attribute pair

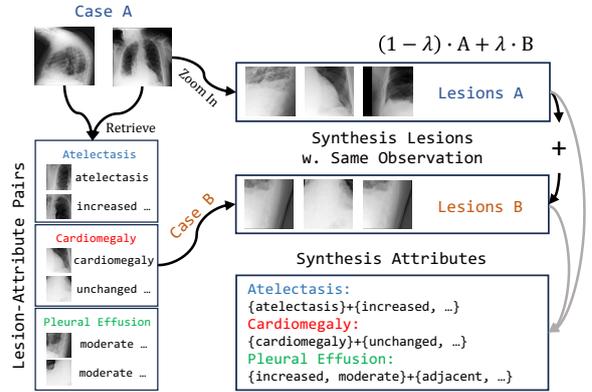


Figure 3: Overview of our proposed lesion-aware mixup augmentation.

(Z_j, A_j) , we retrieve a similar pair (Z_k, A_k) with the same observation from the training data based on report similarity. These lesions are synthesized by a linear combination, as illustrated in Figure 3:

$$Z_j^* = \lambda Z_j + (1 - \lambda) Z_k,$$

where λ is set to 0.75. Note that during training, Z_j^* is used for both INSPECTOR and GENERATOR.

Training INSPECTOR. Similar to §3.2, we adopt a linearly interpolated BCE loss to optimize INSPECTOR:

$$\mathcal{L}_1 = \lambda \text{BCE}_j + (1 - \lambda) \text{BCE}_k,$$

where BCE_j and BCE_k take A_j and A_k as their respective labels. Notably, only the attributes that are shared between Z_j and Z_k are fully optimized. Consequently, our lesion-aware mixup technique facilitates the improvement of output consistency for two semantically equivalent lesions.

Dataset	Model	NLG Metrics						CE Metrics		
		B-1	B-2	B-3	B-4	MTR	R-L	P	R	F ₁
MIMIC-ABN	R2GEN	0.290	0.157	0.093	0.061	0.105	0.208	0.266	0.320	0.272
	R2GENCMN	0.264	0.140	0.085	0.056	0.098	0.212	0.290	0.319	0.280
	ORGAN	0.314	0.180	0.114	0.078	<u>0.120</u>	<u>0.234</u>	0.271	0.342	0.293
	RECAP	<u>0.321</u>	<u>0.182</u>	<u>0.116</u>	<u>0.080</u>	<u>0.120</u>	0.223	<u>0.300</u>	<u>0.363</u>	<u>0.305</u>
	ICON (Ours)	0.337	0.195	0.126	0.086	0.129	0.236	0.332	0.430	0.360
MIMIC-CXR	R2GEN	0.353	0.218	0.145	0.103	0.142	0.270	0.333	0.273	0.276
	R2GENCMN	0.353	0.218	0.148	0.106	0.142	0.278	0.344	0.275	0.278
	M ² TR	0.378	0.232	0.154	0.107	0.145	0.272	0.240	0.428	0.308
	KNOWMAT	0.363	0.228	0.156	0.115	–	0.284	0.458	0.348	0.371
	CMM-RL	0.381	0.232	0.155	0.109	0.151	0.287	0.342	0.294	0.292
	CMCA	0.360	0.227	0.156	0.117	0.148	0.287	0.444	0.297	0.356
	KiUT	0.393	0.243	0.159	0.113	0.160	0.285	0.371	0.318	0.321
	DCL	–	–	–	0.109	0.150	0.284	0.471	0.352	0.373
	METrans	0.386	0.250	0.169	0.124	0.152	<u>0.291</u>	0.364	0.309	0.311
	RGRG	0.373	0.249	0.175	0.126	0.168	0.264	0.380	0.319	0.305
	ORGAN	0.407	0.256	0.172	0.123	0.162	0.293	0.416	0.418	0.385
	RECAP	0.429	0.267	<u>0.177</u>	0.125	<u>0.168</u>	0.288	0.389	<u>0.443</u>	<u>0.393</u>
	ICON (Ours)	0.429	<u>0.266</u>	0.178	0.126	0.170	0.287	<u>0.445</u>	0.505	0.464

Table 1: Experimental results of our model and baselines on the MIMIC-ABN and MIMIC-CXR datasets. The best results are in **boldface**, and the underlined are the second-best results. The listed CE results are macro-weighted, and example-based CE results are provided in Table 9.

3.4 Generating Consistent Radiology Reports (GENERATOR)

Lesion-Attribute Alignment. To bridge the modality gap between lesion representations and text-based attributes, we leverage a BART (Lewis et al., 2020) encoder to extract attribute representations. The attributes associated with each lesion are formulated as a prompt: $\langle s \rangle o_j \langle /s \rangle A_j \langle /s \rangle$, as depicted in the upper part of Figure 2. Then, a cross-attention module (Vaswani et al., 2017) is inserted after every self-attention module. This module aligns the lesion representations with the attribute representations by querying visual representations using attribute representations, similar to Q-Former (Li et al., 2023a):

$$H_j^a = \text{CrossAttention}(H_j^s, Z_j, Z_j),$$

where $H_j^a, H_j^s \in \mathbb{R}^h$ are the aligned attribute representation and the self-attended representation of A_j , respectively. All prompts are encoded, and the attribute representations of Z are denoted as \mathcal{H}^a .

Report Generation. Given the input images \mathcal{X} , images of prior visits \mathcal{X}^p , the lesions Z , and attribute \mathcal{H}^a , we utilize a BART decoder in conjunction with the Fusion-in-Decoder (FiD; (Izacard and Grave, 2021)) that simply concatenates multiple context sequences for report generation. Then, the probability of the t -th step is expressed as:

$$h_t = \text{FiD}([\mathcal{X}; \mathcal{X}^p; Z; \mathcal{H}^a], h_{\langle t \rangle}),$$

$$p(y_t | \mathcal{X}, \mathcal{X}^p, Z, \mathcal{Y}_{\langle t \rangle}) = \text{Softmax}(W_g h_t + b_g),$$

where $h_t \in \mathbb{R}^h$ is the t -th hidden representation, $W_g \in \mathbb{R}^{|\mathcal{V}| \times h}$ is the weight matrix, $b_g \in \mathbb{R}^{|\mathcal{V}|}$ is the bias vector, and \mathcal{V} is the vocabulary.

Training GENERATOR. The generation process is optimized using the negative log-likelihood loss, given each token’s probability $p(y_t | \mathcal{X}, \mathcal{X}^p, Z, \mathcal{Y}_{\langle t \rangle})$:

$$\mathcal{L}_G = - \sum_{t=1}^T \log p(y_t | \mathcal{X}, \mathcal{X}^p, Z, \mathcal{Y}_{\langle t \rangle}).$$

The loss function of Stage 2 is: $\mathcal{L}_{S2} = \mathcal{L}_I + \mathcal{L}_G$.

4 Experiments

4.1 Datasets

Three public datasets are used to evaluate our models, i.e., IU X-RAY⁴ (Demner-Fushman et al., 2016), MIMIC-CXR⁵ (Johnson et al., 2019), and MIMIC-ABN⁶ (Ni et al., 2020). We follow previous research (Chen et al., 2020) to preprocess these datasets, and provide other details in Appendix A.7.

- IU X-RAY consists of 3,955 reports. We follow previous research (Chen et al., 2020) and split the dataset into train/validation/test sets with a ratio of 7:1:2.
- MIMIC-CXR consists of 377,110 chest X-ray images and 227,827 reports.

⁴<https://openi.nlm.nih.gov/>

⁵<https://physionet.org/content/mimic-cxr-jpg/2.0.0/>

⁶<https://github.com/zzxslp/WCL>

Dataset	Model	NLG Metrics		RadGraph		
		B-4	R-L	RG _E	RG _{ER}	RG _{ER}
IU X-RAY	R2GEN	0.120	0.298	—	—	—
	M ² TR	0.121	0.288	—	—	—
	T _{NLL}	0.114	—	0.230	0.202	0.153
	ICON	0.098	0.320	0.342	0.312	0.246
MIMIC -CXR	T _{NLL}	0.105	0.253	0.230	0.202	0.153
	ORGAN	0.123	0.293	0.303	0.275	0.199
	RECAP	0.125	0.288	0.307	0.276	0.205
	ICON	0.126	0.287	0.312	0.278	0.197

Table 2: Radgraph evaluation results on the IU X-RAY and MIMIC-CXR datasets. Results of T_{NLL} are cited from Delbrouck et al. (2022).

- MIMIC-ABN is modified from the MIMIC-CXR dataset and its reports only contain abnormal part. We adopt the data-split as used in Hou et al. (2023a), and the data-split is 71,786/546/806 for train/validation/test sets.

Unlike previous research (Chen et al., 2020) which only used one view for report generation on MIMIC-CXR and MIMIC-ABN, we collect all views for each visit in experiments. The justification is provided in Appendix A.8.

4.2 Evaluation Metrics and Baselines

NLG Metrics. To assess the quality of generated reports, we adopt several natural language generation (NLG) metrics for evaluation. BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), and ROUGE (Lin, 2004) are selected as NLG Metrics, and we use the MS-COCO caption evaluation tool⁷ to compute the results.

CE Metrics. Following previous research (Chen et al., 2020, 2021), we adopt clinical efficacy (CE) metrics to evaluate the observation-level factual accuracy, and CheXbert (Smit et al., 2020) is used in this paper. To measure the entity-level factual accuracy, we leverage the RadGraph (Jain et al., 2021; Delbrouck et al., 2022) and temporal entity matching (TEM) scores for evaluation.

Consistency Metrics. CON and R-CON (§2.3) are utilized to measure the inter-report consistency. Note that entities used in measuring consistency are adopted from RadGraph (Jain et al., 2021). A MAJORITY baseline which outputs the same report for all inputs, is included.

Baselines. We compare our models with the following baselines: R2GEN (Chen et al., 2020), R2GENCMN (Chen et al., 2021), KNOWMAT (Yang et al., 2021), M²TR (Nooralahzadeh et al., 2021), CMM-RL (Qin and Song, 2022), CMCA (Song et al., 2022), CXR-RePaiR-Sel/2 (Endo et al.,

⁷<https://github.com/tylin/coco-caption>

Model	MIMIC-ABN		MIMIC-CXR	
	CON	R-CON	CON	R-CON
MAJORITY	1.000	—	1.000	—
R2GEN	0.280	0.072	0.137	0.042
R2GENCMN	0.302	0.091	0.155	0.049
ORGAN	0.338	0.127	0.345	0.126
RECAP	0.311	0.108	0.345	0.114
ICON (Ours)	0.316	0.140	0.351	0.163
ICON w/o ZOOM	0.183	0.073	0.175	0.066
ICON w/o INSPECT	0.253	0.100	0.245	0.090
ICON w/o MIXUP	0.286	0.119	0.334	0.156

Table 3: The CON score and the R-CON score. MAJORITY: outputs the same report for all inputs.

2021), BioViL-T (Bannur et al., 2023), DCL (Li et al., 2023b), METrans (Wang et al., 2023c), KiUT (Huang et al., 2023), RGRG (Tanida et al., 2023), ORGAN (Hou et al., 2023b), and RECAP (Hou et al., 2023a).

4.3 Implementation Details

The small and tiny versions of Swin Transformer V2 (Liu et al., 2022) are used as the visual backbone for ZOOMER and INSPECTOR, respectively. The GENERATOR is initialized with the base version of BART pretrained on biomedical corpus (Yuan et al., 2022). Other parameters are randomly initialized. For Stage 2 training, the learning rate is $5e-5$ with linear decay, the batch size is 32, and the models are trained for 20 and 5 epochs on MIMIC-ABN and MIMIC-CXR with early stopping, respectively. Since the number of samples in IU X-RAY is too small to train a multimodal model, we only provide results produced by models trained on MIMIC-CXR as a reference, similar to (Delbrouck et al., 2022). For other training details (e.g., training ZOOMER), and the resources used in this paper, we list them in Appendix A.3.

5 Results

5.1 Quantitative Analysis

Inter-Report Consistency Analysis. Table 3 provides CON and R-CON scores of baselines, our model, and its ablated variants. **ICON achieves the highest R-CON on both datasets, indicating the best inter-report consistency.** In terms of the CON score, ICON demonstrates competitive performance when compared with ORGAN. We also notice that introducing mixup augmentation leads to a large improvement on CON, demonstrating the effectiveness of lesion-aware mixup.

NLG and Temporal Modeling Results. The NLG results are presented in Table 1 and the Temporal

Dataset	Model	Components			NLG Metrics						CE Metrics		
		ZOOM	INSPECT	MIXUP	B-1	B-2	B-3	B-4	MTR	R-L	P	R	F ₁
MIMIC-ABN	ICON	✓	✓	✓	0.337	0.195	0.126	0.086	0.129	0.236	0.332	0.430	0.360
	ICON w/o ZOOM	—	—	—	0.310	0.181	0.119	0.084	0.120	0.243	0.306	0.353	0.306
	ICON w/o INSPECT	✓	—	—	0.315	0.182	0.117	0.081	0.121	0.236	0.338	0.401	0.352
	ICON w/o MIXUP	✓	✓	—	0.335	0.192	0.124	0.085	0.129	0.239	0.332	0.413	0.356
MIMIC-CXR	ICON	✓	✓	✓	0.429	0.266	0.178	0.126	0.170	0.287	0.445	0.505	0.464
	ICON w/o ZOOM	—	—	—	0.377	0.237	0.162	0.119	0.149	0.288	0.363	0.280	0.278
	ICON w/o INSPECT	✓	—	—	0.399	0.248	0.168	0.122	0.157	0.287	0.444	0.447	0.423
	ICON w/o MIXUP	✓	✓	—	0.427	0.264	0.176	0.124	0.169	0.285	0.444	0.502	0.462

Table 4: Ablation results of our model and its variants on the MIMIC-ABN and MIMIC-CXR datasets.

Model	B-4	R-L	CE-F ₁	TEM
CXR-RePaiR-2	0.021	0.143	0.281	0.125
BioViL-NN	0.037	0.200	0.283	0.111
BioViL-T-NN	0.045	0.205	0.290	0.130
BioViL-AR	0.075	0.279	0.293	0.138
BioViL-T-AR	0.092	0.296	0.317	0.175
RECAP	0.118	0.279	0.400	0.304
ICON (Ours)	0.120	0.279	0.468	0.335

Table 5: Progression modeling results on the MIMIC-CXR dataset. Results of BioViL-* are cited from Ban-nur et al. (2023).

Modeling results are listed in Table 5. Among all models, **ICON achieves SOTA performance on the NLG and Temporal metrics.** As shown in Table 1, our model demonstrates significant improvements on the MIMIC-ABN dataset and achieves competitive performance on the MIMIC-CXR dataset. Additionally, we provide experimental results on the IU X-RAY dataset as a reference in Table 2. Regarding temporal modeling, ICON exhibits significant improvements over other baselines in terms of BLEU score, clinical accuracy, and TEM score while maintaining competitive performance on ROUGE, indicating its enhanced capacity to effectively utilize historical records.

Clinical Efficacy Results. In the right section of Table 1, we observe that **ICON achieves SOTA clinical accuracy**, increasing CE F₁ from 0.393 to 0.464 on the MIMIC-CXR dataset and rising by 5.5% on the MIMIC-ABN dataset. These results indicate that our model is capable of generating accurate and consistent radiology reports. Furthermore, Table 2 presents the RadGraph F₁ on both the IU X-RAY and MIMIC-CXR datasets. Our model achieves competitive performance compared with the non-RL-optimized baselines.

Ablation Results. The ablation results for MIMIC-ABN and MIMIC-CXR are listed in Table 3 and Table 4. We study three variants: (1) *w/o ZOOM*, where all components are removed, (2) *w/o INSPECT*, where both the INSPECTOR and MIXUP are removed, and (3) *w/o MIXUP*, where only MIXUP

is removed. The performance of the ablated model *w/o ZOOM* drops significantly for both datasets, while the variant *w/o INSPECT* achieves competitive results on clinical accuracy. This suggests that the ZOOMER effectively extracts lesions and provides relevant abnormal information for report generation. In addition, the variant *w/o MIXUP* further improves the performance, demonstrating the effectiveness of INSPECTOR in transforming concise lesion information into precise free-text reports. Moreover, introducing lesion-aware mixup augmentation strengthens the consistency of generated outputs, indicating the effectiveness of ICON.

5.2 Qualitative Analysis

Case Study. Figure 4 showcases two semantically equivalent cases, i.e., Case A and Case B, extracted from the test set of MIMIC-CXR. In both instances, ICON successfully identifies abnormal observations (e.g., *Cardiomegaly*, *Pleural Effusion*, and *Atelectasis*) and generates consistent phrases including "*pulmonary vascular congestion*", "*bilateral pleural effusions*", and "*compressive atelectasis*." Conversely, the variant *w/o ZOOM* fails to produce these descriptions in Case A. This demonstrates that ZOOMER plays a crucial role in identifying lesions and highlights the ability of the mixup augmentation to ensure the alignment of lesions with their corresponding attributes.

Error Analysis. Figure 5 presents an error case produced by ICON. Although ZOOMER successfully identifies *Pneumonia* in the given radiographs, the GENERATOR fails to realize it into descriptions like "*multifocal pneumonia*" (i.e., a false negative observation). We notice that the region of this observation is inaccurately identified. Additionally, ZOOMER outputs a false positive observation *Lung Opacity*, leading to an inaccurate phrase "*increased opacity*". To mitigate this issue, a better ZOOMER trained with larger datasets could be beneficial.

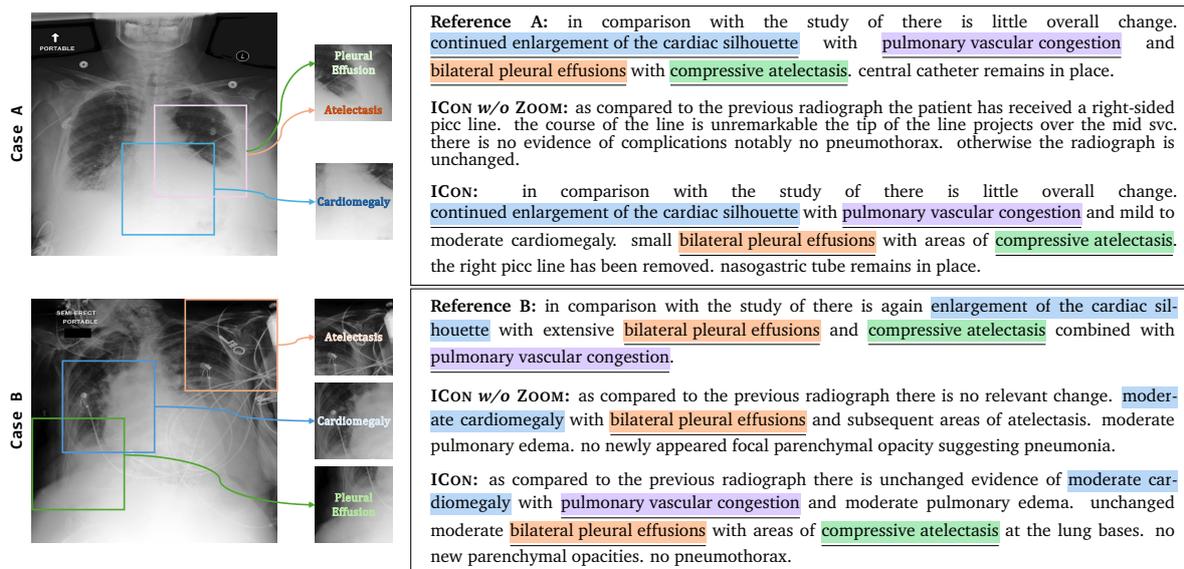


Figure 4: A case study of ICON on two semantically equivalent cases (i.e., Case A and Case B), given their radiographs and lesions. Spans with the same color (*Cardiomegaly*, *Pleural Effusion*, *Atelectasis*, and *Others*) represent the same positive observation. Consistent and accurate outputs are highlighted with underline.

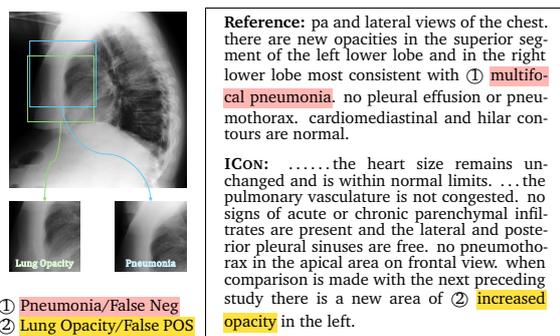


Figure 5: An error case produced by ICON. The span and the span denote false negative observation and false positive observation, respectively.

6 Related Works

Radiology report generation (Jing et al., 2018; Li et al., 2018; Jing et al., 2019) has gained much attention. Prior research has either devised various memory mechanisms to record key information (Chen et al., 2020, 2021; Qin and Song, 2022; Wang et al., 2023c; Zhao et al., 2023) or proposed different learning methods to enhance performance (Liu et al., 2021c,a,b). In addition, Yang et al. (2021); Li et al. (2023b); Huang et al. (2023); Yan et al. (2023) proposed utilizing knowledge graphs for report generation. Liu et al. (2019); Lovelace and Mortazavi (2020); Miura et al. (2021); Nishino et al. (2022); Delbrouck et al. (2022) designed various rewards for reinforcement learning to improve clinical accuracy. Tanida et al. (2023) proposed an explainable framework for report generation. Hou et al. (2023b) introduced observations to improve

factual accuracy. Kale et al. (2023) proposed a template-based approach to improve the quality and accuracy of radiology reports. Additionally, Ramesh et al. (2022); Bannur et al. (2023); Hou et al. (2023a); Dalla Serra et al. (2023) focused on exploring the temporal structure. Wang et al. (2023b,a) utilized CLIP (Radford et al., 2021) to bridge the modality gap. Mixup is also closely related to this research (Zhang et al., 2018), and this method has been adopted in NLP research (Sun et al., 2020; Yoon et al., 2021; Yang et al., 2022).

Although consistency has been studied in many domains (Thimm, 2013; Ribeiro et al., 2019; Camburu et al., 2019; Elazar et al., 2021), it remains unexplored in medical report generation.

7 Conclusion and Future Works

In this paper, we propose ICON, comprising three components to improve both accuracy and inter-report consistency. ICON first extracts lesions and then matches fine-grained attributes for report generation. A lesion-aware mixup method is devised for attribute alignment. Experimental results on three datasets demonstrate the effectiveness of ICON. In the future, we plan to explore incorporating large language models (LLMs) into our framework, given their advanced capabilities in planning and generation, to further enhance the performance of radiology report generation. Leveraging the strengths of LLMs could provide more refined signals to enhance the performance of ICON.

530 Limitations

531 Although ICON can improve the consistency of
532 radiology report generation, it still exhibits some
533 limitations. Since our lesion extraction method is
534 based on coarse-grained labels (i.e., image labels),
535 training such a model requires annotations for im-
536 ages. However, obtaining these annotations can be
537 challenging in some medical settings. Recent ad-
538 vances in foundation vision models (Kirillov et al.,
539 2023) and open-set learning (Zara et al., 2023)
540 could be a potential direction to handle this is-
541 sue. Additionally, since our framework consists of
542 two stages, prediction errors can propagate through
543 the pipeline, making the final performance of our
544 framework largely dependent on Stage 1. Rein-
545 forcement learning (Nishino et al., 2022) that takes
546 factual improvement as a reward could be a solu-
547 tion to optimize the framework in an end-to-end
548 manner.

549 Ethics Statement

550 The IU X-RAY (Demner-Fushman et al., 2016),
551 MIMIC-ABN (Ni et al., 2020), and MIMIC-
552 CXR (Johnson et al., 2019) datasets are publicly
553 available and have been automatically de-identified
554 to protect patient privacy. Our goal is to enhance
555 the inter-report consistency of radiology report gen-
556 eration systems. Despite the substantial improve-
557 ment of our framework over state-of-the-art base-
558 lines, the performance still lags behind the require-
559 ments for real-world deployment and could lead
560 to unexpected failures in untested environments.
561 Thus, we urge readers of this paper and potential
562 users of this system to cautiously check the gen-
563 erated outputs and seek expert advice when using
564 it.

565 References

566 Onur Asan, Alparslan Emrah Bayrak, and Avishek
567 Choudhury. 2020. [Artificial intelligence and human
568 trust in healthcare: Focus on clinicians](#). *J Med Inter-
569 net Res*, 22(6):e15154.

570 Satanjeev Banerjee and Alon Lavie. 2005. [METEOR:
571 An automatic metric for MT evaluation with im-
572 proved correlation with human judgments](#). In *Pro-
573 ceedings of the ACL Workshop on Intrinsic and Ex-
574 trinsic Evaluation Measures for Machine Transla-
575 tion and/or Summarization*, pages 65–72, Ann Arbor,
576 Michigan. Association for Computational Linguis-
577 tics.

Shruthi Bannur, Stephanie Hyland, Qianchu Liu, Fer- 578
nando Pérez-García, Maximilian Ilse, Daniel C. Cas- 579
tro, Benedikt Boecking, Harshita Sharma, Kenza 580
Bouزيد, Anja Thieme, Anton Schwaighofer, Maria 581
Wetscherek, Matthew P. Lungren, Aditya Nori, Javier 582
Alvarez-Valle, and Ozan Oktay. 2023. [Learning
583 to exploit temporal structure for biomedical vision-
584 language processing](#). 585

Oana-Maria Camburu, Brendan Shillingford, Pasquale 586
Minervini, Thomas Lukasiewicz, and Phil Blunsom. 587
2019. Make up your mind! adversarial generation 588
of inconsistent natural language explanations. *arXiv
589 preprint arXiv:1910.03065*. 590

Zhihong Chen, Yaling Shen, Yan Song, and Xiang Wan. 591
2021. [Cross-modal memory networks for radiology
592 report generation](#). In *Proceedings of the 59th An-
593 nual Meeting of the Association for Computational
594 Linguistics and the 11th International Joint Confer-
595 ence on Natural Language Processing, ACL/IJCNLP
596 2021, (Volume 1: Long Papers), Virtual Event, Au-
597 gust 1-6, 2021*, pages 5904–5914. Association for
598 Computational Linguistics. 599

Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xi- 600
ang Wan. 2020. Generating radiology reports via 601
memory-driven transformer. In *Proceedings of the
602 2020 Conference on Empirical Methods in Natural
603 Language Processing*. 604

Kenneth Ward Church and Patrick Hanks. 1990. [Word
605 association norms, mutual information, and lexicog-
606 raphy](#). *Computational Linguistics*, 16(1):22–29. 607

Francesco Dalla Serra, Chaoyang Wang, Fani Deli- 608
gianni, Jeff Dalton, and Alison O’Neil. 2023. [Con-
609 trollable chest X-ray report generation from longi-
610 tudinal representations](#). In *Findings of the Associa-
611 tion for Computational Linguistics: EMNLP 2023*,
612 pages 4891–4904, Singapore. Association for Com-
613 putational Linguistics. 614

Jean-Benoit Delbrouck, Pierre Chambon, Christian 615
Bluethgen, Emily Tsai, Omar Almusa, and Curtis 616
Langlotz. 2022. [Improving the factual correctness of
617 radiology report generation with semantic rewards](#).
618 In *Findings of the Association for Computational
619 Linguistics: EMNLP 2022*, pages 4348–4360, Abu
620 Dhabi, United Arab Emirates. Association for Com-
621 putational Linguistics. 622

Dina Demner-Fushman, Marc D Kohli, Marc B Rosen- 623
man, Sonya E Shooshan, Laritza Rodriguez, Sameer 624
Antani, George R Thoma, and Clement J McDon- 625
ald. 2016. Preparing a collection of radiology ex- 626
aminations for distribution and retrieval. *Journal
627 of the American Medical Informatics Association*,
628 23(2):304–310. 629

Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhi- 630
lasha Ravichander, Eduard Hovy, Hinrich Schütze, 631
and Yoav Goldberg. 2021. [Measuring and improving
632 consistency in pretrained language models](#). *Transac-
633 tions of the Association for Computational Linguis-
634 tics*, 9:1012–1031. 635

636	Mark Endo, Rayan Krishnan, Viswesh Krishna, Andrew Y. Ng, and Pranav Rajpurkar. 2021. Retrieval-based chest x-ray report generation using a pre-trained contrastive language-image model . In <i>Proceedings of Machine Learning for Health</i> , volume 158 of <i>Proceedings of Machine Learning Research</i> , pages 209–219. PMLR.	691
637		692
638		693
639		694
640		695
641		696
642		
643	Wenjun Hou, Yi Cheng, Kaishuai Xu, Wenjie Li, and Jiang Liu. 2023a. Recap: Towards precise radiology report generation via dynamic disease progression reasoning .	697
644		698
645		699
646		700
647	Wenjun Hou, Kaishuai Xu, Yi Cheng, Wenjie Li, and Jiang Liu. 2023b. ORGAN: Observation-guided radiology report generation via tree reasoning . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8108–8122, Toronto, Canada. Association for Computational Linguistics.	701
648		702
649		703
650		704
651		705
652		706
653		707
654	Zhongzhen Huang, Xiaofan Zhang, and Shaoting Zhang. 2023. Kiut: Knowledge-injected u-transformer for radiology report generation . In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 19809–19818.	708
655		709
656		710
657		711
658		712
659	Gautier Izacard and Edouard Grave. 2021. Leveraging passage retrieval with generative models for open domain question answering . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 874–880, Online. Association for Computational Linguistics.	713
660		714
661		715
662		716
663		
664		
665		
666	Saahil Jain, Ashwin Agrawal, Adriel Saporta, Steven Q. H. Truong, Du Nguyen Duong, Tan Bui, Pierre Chambon, Yuhao Zhang, Matthew P. Lungren, Andrew Y. Ng, Curtis P. Langlotz, and Pranav Rajpurkar. 2021. Radgraph: Extracting clinical entities and relations from radiology reports . <i>CoRR</i> , abs/2106.14463.	717
667		718
668		719
669		720
670		
671		
672	Baoyu Jing, Zeya Wang, and Eric Xing. 2019. Show, describe and conclude: On exploiting the structure information of chest X-ray reports . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 6570–6580, Florence, Italy. Association for Computational Linguistics.	721
673		722
674		723
675		724
676		725
677		726
678	Baoyu Jing, Pengtao Xie, and Eric P. Xing. 2018. On the automatic generation of medical imaging reports . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers</i> , pages 2577–2586. Association for Computational Linguistics.	727
679		728
680		729
681		
682		
683		
684		
685	Alistair EW Johnson, Tom J Pollard, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Yifan Peng, Zhiyong Lu, Roger G Mark, Seth J Berkowitz, and Steven Horng. 2019. Mimic-cxr-jpg, a large publicly available database of labeled chest radiographs . <i>arXiv preprint arXiv:1901.07042</i> .	730
686		731
687		732
688		733
689		734
690		735
	Kaveri Kale, Pushpak Bhattacharyya, and Kshitij Jadhav. 2023. Replace and report: NLP assisted radiology report generation . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 10731–10742, Toronto, Canada. Association for Computational Linguistics.	736
		737
		738
		739
		740
		741
		742
	Sara Kaviani, Ki Jin Han, and Insoo Sohn. 2022. Adversarial attacks and defenses on ai in medical imaging informatics: A survey . <i>Expert Systems with Applications</i> , 198:116815.	743
		744
		745
		746
	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar, and Ross Girshick. 2023. Segment anything . In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 4015–4026.	
	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880, Online. Association for Computational Linguistics.	
	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models .	
	Mingjie Li, Wenjia Cai, Rui Liu, Yuetian Weng, Xiaoyun Zhao, Cong Wang, Xin Chen, Zhong Liu, Caineng Pan, Mengke Li, yingfeng zheng, Yizhi Liu, Flora D. Salim, Karin Verspoor, Xiaodan Liang, and Xiaojun Chang. 2021. FFA-IR: Towards an explainable and reliable medical report generation benchmark . In <i>Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)</i> .	
	Mingjie Li, Bingqian Lin, Zicong Chen, Haokun Lin, Xiaodan Liang, and Xiaojun Chang. 2023b. Dynamic graph enhanced contrastive learning for chest x-ray report generation . In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 3334–3343.	
	Yuan Li, Xiaodan Liang, Zhiting Hu, and Eric P. Xing. 2018. Hybrid retrieval-generation reinforced agent for medical image report generation . In <i>Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada</i> , pages 1537–1547.	
	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	

747	Fenglin Liu, Shen Ge, and Xian Wu. 2021a.	<i>Association for Computational Linguistics: Human</i>	804
748	Competence-based multimodal curriculum learning	<i>Language Technologies</i> , pages 5288–5304, Online.	805
749	for medical report generation. In <i>Proceedings of the</i>	Association for Computational Linguistics.	806
750	59th Annual Meeting of the Association for Com-		
751	putational Linguistics and the 11th International		
752	Joint Conference on Natural Language Processing,	Jianmo Ni, Chun-Nan Hsu, Amilcare Gentili, and Julian	807
753	ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual	McAuley. 2020. Learning visual-semantic embed-	808
754	Event, August 1-6, 2021, pages 3001–3012. Associa-	dings for reporting abnormal findings on chest X-rays.	809
755	tion for Computational Linguistics.	In <i>Findings of the Association for Computational Lin-</i>	810
		guistics: EMNLP 2020, pages 1954–1960, Online.	811
756	Fenglin Liu, Xian Wu, Shen Ge, Wei Fan, and Yuexian	Association for Computational Linguistics.	812
757	Zou. 2021b. Exploring and distilling posterior and		
758	prior knowledge for radiology report generation. In	Toru Nishino, Yasuhide Miura, Tomoki Taniguchi,	813
759	<i>IEEE Conference on Computer Vision and Pattern</i>	Tomoko Ohkuma, Yuki Suzuki, Shoji Kido, and	814
760	<i>Recognition, CVPR 2021, virtual, June 19-25, 2021,</i>	Noriyuki Tomiyama. 2022. Factual accuracy is not	815
761	pages 13753–13762. Computer Vision Foundation /	enough: Planning consistent description order for	816
762	IEEE.	radiology report generation. In <i>Proceedings of the</i>	817
		2022 Conference on Empirical Methods in Natural	818
763	Fenglin Liu, Changchang Yin, Xian Wu, Shen Ge, Ping	Language Processing, Online. Association for Com-	819
764	Zhang, and Xu Sun. 2021c. Contrastive attention	putational Linguistics.	820
765	for automatic chest x-ray report generation. In <i>Find-</i>		
766	ings of the Association for Computational Linguis-	Farhad Nooralahzadeh, Nicolas Perez Gonzalez,	821
767	tics: ACL/IJCNLP 2021, Online Event, August 1-6,	Thomas Frauenfelder, Koji Fujimoto, and Michael	822
768	2021, volume ACL/IJCNLP 2021 of <i>Findings of ACL,</i>	Krauthammer. 2021. Progressive transformer-based	823
769	pages 269–280. Association for Computational Lin-	generation of radiology reports. In <i>Findings of the</i>	824
770	guistics.	Association for Computational Linguistics: EMNLP	825
		2021, pages 2824–2832, Punta Cana, Dominican Re-	826
771	Guanxiong Liu, Tzu-Ming Harry Hsu, Matthew B. A.	public. Association for Computational Linguistics.	827
772	McDermott, Willie Boag, Wei-Hung Weng, Peter		
773	Szolovits, and Marzyeh Ghassemi. 2019. Clini-	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	828
774	cally accurate chest x-ray report generation. <i>CoRR,</i>	Jing Zhu. 2002. Bleu: a method for automatic evalu-	829
775	abs/1904.02633.	ation of machine translation. In <i>Proceedings of the</i>	830
		40th Annual Meeting of the Association for Compu-	831
776	Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda	tational Linguistics, pages 311–318, Philadelphia,	832
777	Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang,	Pennsylvania, USA. Association for Computational	833
778	Li Dong, Furu Wei, and Baining Guo. 2022. Swin	Linguistics.	834
779	transformer v2: Scaling up capacity and resolution.		
780	In <i>Proceedings of the IEEE/CVF Conference on Com-</i>	Adnan Qayyum, Junaid Qadir, Muhammad Bilal, and	835
781	puter Vision and Pattern Recognition (CVPR), pages	Ala Al-Fuqaha. 2020. Secure and robust machine	836
782	12009–12019.	learning for healthcare: A survey. <i>IEEE Reviews in</i>	837
		Biomedical Engineering, 14:156–180.	838
783	Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei,		
784	Zheng Zhang, Stephen Lin, and Baining Guo. 2021d.	Han Qin and Yan Song. 2022. Reinforced cross-modal	839
785	Swin transformer: Hierarchical vision transformer	alignment for radiology report generation. In <i>Find-</i>	840
786	using shifted windows. In <i>Proceedings of the</i>	ings of the Association for Computational Linguistics:	841
787	<i>IEEE/CVF International Conference on Computer</i>	ACL 2022, Dublin, Ireland, May 22-27, 2022, pages	842
788	<i>Vision (ICCV)</i> , pages 10012–10022.	448–458. Association for Computational Linguistics.	843
789	Justin Lovelace and Bobak Mortazavi. 2020. Learning	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	844
790	to generate clinically coherent chest X-ray reports.	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	845
791	In <i>Findings of the Association for Computational</i>	try, Amanda Askell, Pamela Mishkin, Jack Clark,	846
792	Linguistics: EMNLP 2020, pages 1235–1243, Online.	Gretchen Krueger, and Ilya Sutskever. 2021. Learn-	847
793	Association for Computational Linguistics.	ing transferable visual models from natural language	848
		supervision.	849
794	Xingjun Ma, Yuhao Niu, Lin Gu, Yisen Wang, Yitian		
795	Zhao, James Bailey, and Feng Lu. 2021. Under-	Vignav Ramesh, Nathan Andrew Chi, and Pranav Ra-	850
796	standing adversarial attacks on deep learning based	jpurkar. 2022. Improving radiology report generation	851
797	medical image analysis systems. <i>Pattern Recognition,</i>	systems by removing hallucinated references to non-	852
798	110:107332.	existent priors.	853
799	Yasuhide Miura, Yuhao Zhang, Emily Tsai, Curtis Lan-	Shaoqing Ren, Kaiming He, Ross Girshick, and Jian	854
800	glotz, and Dan Jurafsky. 2021. Improving factual	Sun. 2015. Faster r-cnn: towards real-time object	855
801	completeness and consistency of image-to-text	detection with region proposal networks. In <i>Proceed-</i>	856
802	radiology report generation. In <i>Proceedings of the 2021</i>	ings of the 28th International Conference on Neural	857
803	Conference of the North American Chapter of the	Information Processing Systems - Volume 1, NIPS’15,	858
		page 91–99, Cambridge, MA, USA. MIT Press.	859

860	Marco Tulio Ribeiro, Carlos Guestrin, and Sameer Singh. 2019. Are red roses red? evaluating consistency of question-answering models . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 6174–6184, Florence, Italy. Association for Computational Linguistics.	918
861		919
862		920
863		921
864		922
865		923
866		924
867	Hoo-Chang Shin, Kirk Roberts, Le Lu, Dina Demner-Fushman, Jianhua Yao, and Ronald M Summers. 2016. Learning to read chest x-rays: Recurrent neural cascade model for automated image annotation . In <i>2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 2497–2506.	925
868		926
869		927
870		928
871		929
872		930
873	George Gaylord Simpson. 1943. Mammals and the nature of continents. <i>American Journal of Science</i> , 241(1):1–31.	931
874		932
875		933
876	Akshay Smit, Saahil Jain, Pranav Rajpurkar, Anuj Pareek, Andrew Ng, and Matthew Lungren. 2020. Combining automatic labelers and expert annotations for accurate radiology report labeling using BERT . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1500–1519, Online. Association for Computational Linguistics.	934
877		935
878		936
879		937
880		938
881		939
882		940
883		941
884		942
885	Xiao Song, Xiaodan Zhang, Junzhong Ji, Ying Liu, and Pengxu Wei. 2022. Cross-modal contrastive attention model for medical report generation . In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 2388–2397, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.	943
886		944
887		945
888		946
889		947
890		948
891	Lichao Sun, Congying Xia, Wenpeng Yin, Tingting Liang, Philip Yu, and Lifang He. 2020. Mixup-transformer: Dynamic data augmentation for NLP tasks . In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 3436–3440, Barcelona, Spain (Online). International Committee on Computational Linguistics.	949
892		950
893		951
894		952
895		953
896		954
897		955
898	Tim Tanida, Philip Müller, Georgios Kaissis, and Daniel Rueckert. 2023. Interactive and explainable region-guided radiology report generation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 7433–7442.	956
899		957
900		958
901		959
902		
903	Matthias Thimm. 2013. Inconsistency measures for probabilistic logics . <i>Artif. Intell.</i> , 197:1–24.	960
904		961
905	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17</i> , page 6000–6010, Red Hook, NY, USA. Curran Associates Inc.	962
906		963
907		964
908		965
909		966
910		967
911		968
912	Siyuan Wang, Zheng Liu, and Bo Peng. 2023a. A self-training framework for automated medical report generation . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 16443–16449, Singapore. Association for Computational Linguistics.	969
913		970
914		971
915		972
916		973
917		974
		975
	Siyuan Wang, Bo Peng, Yichao Liu, and Qi Peng. 2023b. Fine-grained medical vision-language representation learning for radiology report generation . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 15949–15956, Singapore. Association for Computational Linguistics.	
	Zhanyu Wang, Lingqiao Liu, Lei Wang, and Luping Zhou. 2023c. Metransformer: Radiology report generation by transformer with multiple learnable expert tokens. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 11558–11567.	
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	
	Benjamin Yan, Ruochen Liu, David Kuo, Subathra Adithan, Eduardo Reis, Stephen Kwak, Vasantha Venugopal, Chloe O’Connell, Agustina Saenz, Pranav Rajpurkar, and Michael Moor. 2023. Style-aware radiology report generation with RadGraph and few-shot prompting . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 14676–14688, Singapore. Association for Computational Linguistics.	
	Huiyun Yang, Huadong Chen, Hao Zhou, and Lei Li. 2022. Enhancing cross-lingual transfer by manifold mixup . In <i>International Conference on Learning Representations</i> .	
	Shuxin Yang, Xian Wu, Shen Ge, Shaohua Kevin Zhou, and Li Xiao. 2021. Knowledge matters: Radiology report generation with general and specific knowledge . <i>CoRR</i> , abs/2112.15009.	
	Soyoung Yoon, Gyuwan Kim, and Kyumin Park. 2021. SSMix: Saliency-based span mixup for text classification . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 3225–3234, Online. Association for Computational Linguistics.	
	Hongyi Yuan, Zheng Yuan, Ruyi Gan, Jiaying Zhang, Yutao Xie, and Sheng Yu. 2022. BioBART: Pretraining and evaluation of a biomedical generative language model . In <i>Proceedings of the 21st Workshop on Biomedical Language Processing</i> , pages 97–109, Dublin, Ireland. Association for Computational Linguistics.	
	Giacomo Zara, Subhankar Roy, Paolo Rota, and Elisa Ricci. 2023. Autolabel: Clip-based framework for open-set video domain adaptation .	

Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. 2018. *mixup: Beyond empirical risk minimization*. In *International Conference on Learning Representations*.

Zizhao Zhang, Yuanpu Xie, Fuyong Xing, Mason McGough, and Lin Yang. 2017. *Mdnet: A semantically and visually interpretable medical image diagnosis network*. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3549–3557.

Guosheng Zhao, Yan Yan, and Zijian Zhao. 2023. *Normal-abnormal decoupling memory for medical report generation*. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1962–1977, Singapore. Association for Computational Linguistics.

A Appendix

A.1 Abnormal Observation Statistics

The abnormal observation statistics of MIMIC-ABN, MIMIC-CXR, and IU X-RAY are listed in Table 6.

#Observation	MIMIC-ABN	MIMIC-CXR	IU X-RAY
No Finding	5002/32/22	64,677/514/229	744/108/318
Cardiomegaly	16,312/118/244	70,561/514/1,602	244/38/61
Pleural Effusion	10,502/80/186	56,972/477/1,379	60/13/15
Pneumothorax	1,452/24/4	8,707/62/106	9/2/5
Enlarged Card.	5,202/40/90	49,806/413/1,140	159/29/28
Consolidation	4,104/36/96	14,449/119/384	17/1/3
Lung Opacity	22,598/166/356	67,714/497/1,448	295/35/57
Fracture	4,458/32/76	11,070/59/232	84/6/15
Lung Lesion	5,612/54/112	11,717/123/300	85/14/17
Edema	8,704/76/168	33,034/257/899	28/2/7
Atelectasis	19,132/134/220	68,273/515/1,210	143/15/37
Support Devices	9,886/58/196	60,455/450/1,358	89/20/16
Pneumonia	17,826/138/260	23,945/184/503	20/2/1
Pleural Other	2,850/30/62	7,296/70/184	32/4/7

Table 6: Observation distribution in train/valid/test split of three datasets. *Enlarged Card.* refers to *Enlarged Cardiomeastinum*.

A.2 Attributes of Observations

We list top-5 attributes for each observation for a better understanding in Table 7.

Observation	Top-5 Attributes
Cardiomegaly	cardiomegaly, borderline, moderately, severely, mildly
Pleural Effusion	layering, subpulmonic, thoracentesis, trace, small
Pneumothorax	hydropneumothorax, apical, tiny, tension, component
Enlarged Card.	mediastinum, widening, contour, widened, lymphadenopathy
Consolidation	consolidative, collapse, underlying, developing, consolidations
Lung Opacity	opacification, opacifications, patchy, heterogeneous, scarring
Fracture	healed, fractured, healing, nondisplaced, posterolateral
Lung Lesion	nodular, nodule, mass, nodules, mm
Edema	indistinctness, asymmetrical, haziness, asymmetric, interstitial
Atelectasis	atelectatic, atelectasis, collapsed, subsegmental, collapse
Support Devices	sidehole, carina, coiled, tunneled, duodenum
Pneumonia	infectious, infection, atypical, supervening, developing
Pleural Other	fibrosis, thickening, biapical, blunting, scarring

Table 7: Top-5 attributes for each observation.

A.3 Additional Implementation Details

For Stage 1, all three datasets use the same hyperparameters for training ZOOMER, with a learning rate of $1e - 4$, batch size of 128, and dropout rate of 0.1, and the number of training epochs is adjusted accordingly. We train ZOOMER for 5, 10, and 15 epochs on MIMIC-CXR, MIMIC-ABN, and IU X-RAY, respectively. During training, several data augmentation methods are applied. The input resolution of Swin Transformer is 256×256 , and we first resize an image to 288×288 , and then randomly crop it to 256×256 with random horizontal flip. All experiments are conducted using one NVIDIA-3090 GTX GPU. For Stage 2, no data augmentation is applied, and we conduct experiments on MIMIC-ABN and IU X-RAY using two NVIDIA-3090 GTX GPUs, and on MIMIC-CXR using four NVIDIA-V100 GPUs, both with half precision. Our model has 328.38M trainable parameters, and the implementations are based on the HuggingFace’s Transformers (Wolf et al., 2020). Here are the pretrained models we used:

- Small version of Swin Transformer V2: <https://huggingface.co/microsoft/swinv2-small-patch4-window8-256>
- Tiny version of Swin Transformer V2: <https://huggingface.co/microsoft/swinv2-tiny-patch4-window8-256>
- Base Version of Biomedical BART: <https://huggingface.co/GanjinZero/biobart-v2-base>

A.4 Additional CE Results on the MIMIC-CXR and MIMIC-ABN Datasets

Observation	Image Classification			Report Classification		
	P	R	F ₁	P	R	F ₁
Enlarged Card.	0.426	0.540	0.476	0.442	0.525	0.428
Cardiomegaly	0.635	0.838	0.722	0.630	0.822	0.714
Lung Opacity	0.535	0.725	0.616	0.542	0.563	0.552
Lung Lesion	0.318	0.187	0.235	0.321	0.177	0.228
Edema	0.471	0.851	0.607	0.464	0.784	0.583
Consolidation	0.283	0.227	0.251	0.275	0.162	0.204
Pneumonia	0.367	0.396	0.381	0.341	0.350	0.345
Atelectasis	0.541	0.660	0.595	0.539	0.620	0.577
Pneumothorax	0.392	0.481	0.432	0.400	0.444	0.421
Pleural Effusion	0.719	0.842	0.776	0.721	0.827	0.770
Pleural Other	0.289	0.440	0.349	0.295	0.315	0.304
Fracture	0.266	0.198	0.227	0.225	0.164	0.190
Support Devices	0.747	0.850	0.795	0.785	0.784	0.785
No Finding	0.366	0.459	0.407	0.263	0.535	0.352
Macro Average	0.454	0.550	0.491	0.445	0.505	0.464

Table 8: Experimental results of each observation on the MIMIC-CXR dataset.

Model	MIMIC-ABN			MIMIC-CXR		
	P	R	F ₁	P	R	F ₁
R2GEN	0.340	0.413	0.348	0.390	0.336	0.337
R2GENCMN	0.360	0.363	0.336	0.358	0.276	0.290
RGRG	—	—	—	0.461	0.475	0.447
ORGAN	0.418	0.471	0.412	0.493	0.560	0.493
RECAP	0.366	0.468	0.382	0.447	0.558	0.464
ICON	0.512	0.428	0.436	0.513	0.597	0.522
ICON <i>w/o</i> ZOOM	0.397	0.406	0.372	0.440	0.362	0.373
ICON <i>w/o</i> INSPECT	0.430	0.479	0.424	0.506	0.553	0.500
ICON <i>w/o</i> MIX-UP	0.433	0.509	0.438	0.507	0.590	0.517

Table 9: Example-based CE results on the MIMIC-ABN and MIMIC-CXR datasets.

A.5 Experimental Results of Stage 1

The experimental results are provided in Table 10. Results on the IU X-RAY dataset are only provided for reference.

Dataset	P	R	F ₁
IU X-RAY	0.223	0.243	0.225
MIMIC-ABN	0.379	0.472	0.411
MIMIC-CXR	0.454	0.550	0.491

Table 10: Abnormal observation prediction results of ZOOMER at Stage 1.

A.6 Lesion Extraction

There are two steps in extraction lesions: candidate generation and candidate classification. Given an image with a resolution of 1024×1024 , padding if needed, we apply a sliding window of 384×384 , with a step size of 128 to extract candidates for classification. This operation results in 36 regions. Then, each region is fed into the ZOOMER for classification, and only the top-1 lesion is selected for each observation. Note that before extracting lesions, each input case is first assigned with their observations by ZOOMER, and as a result, the number of lesions corresponds to the number of observations.

The *No Finding* observation is excluded for lesion extraction, as it estimates the overall conditions of a patient, which makes it difficult to locate at specific regions.

A.7 Other Preprocessing Details

We adopt the same preprocessing setup used in Chen et al. (2020), and the minimum count of each token is set to 3/3/10 for IU X-RAY/MIMIC-ABN/MIMIC-CXR, respectively. Other tokens are replaced with a special token <unk>.

A.8 Justifications for Additional Data Processing

Justification for Using Historical Records. As stated in Hou et al. (2023a), without historical information, it is unreasonable to generate reports with comparisons between two consecutive visits and will lead to hallucinations (Ramesh et al., 2022). As a result, we include historical records as context information for report generation.

Justification for Using All Views. Prior research (Chen et al., 2020, 2021; Hou et al., 2023b,a) treated different views of radiographs in one visit as different samples. However, this is unreasonable to generate a report with only one view position, since different diseases could be observed from different view positions. For example, most of the devices can not be observed from a Lateral view. Given a lateral view radiograph, writing a sentence of "A right chest tube is in unchanged position." is not acceptable.

In addition, some reports describe how many views are provided at the beginning, e.g., "PA and lateral views are provided." Above all, we have justified reasons to use all the views in one visit of a patient to generate the target report. Note that previous work treated each image as a sample and their settings have more samples than ours. For a fair comparison, each generated output of a study with L images is duplicated L times so that the number of samples in evaluation is consistent with previous research.