Debiased Medical Report Generation with High-Frequency Amplification

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026 027 028

029

Paper under double-blind review

ABSTRACT

In recent years, automated medical report generation (MRG) has gained significant research value for its potential to reduce workload and prevent diagnostic errors. However, generating accurate radiology reports remains challenging due to the prevalence of normal regions in X-ray images and normal descriptions in medical reports. Despite various efforts to address these issues, the definitions of visual bias and textual bias remain unclear and there is still a lack of comprehensive analysis of how these biases affect model behavior. In this work, we rigorously define and conduct an in-depth examination of visual and textual biases inherent in MRG datasets. Our analysis emphasizes that global patterns, such as normal regions and findings, contribute to visual and textual bias. Further, we discuss how these biases make MRG models especially prone to frequency bias, where models tend to prioritize low-frequency signals that capture global patterns, while neglecting high-frequency signals. To debiase the frequency bias, we propose the high-frequency amplification layer (HAL), aimed at enhancing the model's perceptiveness to fine-grained details. Our extensive experiments show that by amplifying high-frequency signals, HAL reduces both visual and textual biases, leading to improved performance in MRG tasks.

1 INTRODUCTION

The automation of diagnosis and treatment using medical images has received growing attention in both academia and industry (Wolleb et al., 2022; Manzari et al., 2023; Jiang et al., 2022). In particular, medical report generation (MRG) is one of the most promising tasks as it can alleviate the heavy burden of radiologists and reduce diagnostic errors. MRG aims to automatically generate a free-text description given a medical image (e.g., chest X-ray), describing the detailed findings on both normal and abnormal regions.

Generating diagnostically accurate and domain-specific radiology reports is challenging due to the
presence of severe visual and textual biases. From the perspective of data, most medical images are
dominated by normal regions, making it difficult to capture distinct features (see Figure 1a). Similarly, medical reports primarily describe normal findings, complicating the explanation of abnormal
findings (see Figure 1b). Recently, several methods have been proposed to address visual and textual
biases (You et al., 2021; Liu et al., 2021a; Tanida et al., 2023; Zhang et al., 2020; Liu et al., 2021a;
Huang et al., 2023; Li et al., 2023). However, the definitions of visual bias and textual bias have not
been clearly established and there remains a lack of comprehensive analytical understanding of how
these biases affect model behavior.

Our work focuses on rigorously defining and identifying the fundamental challenges in MRG, analyzing how visual and textual biases hinder model performance. Further, from the perspective of the model, we relate these biases to frequency bias, where the model tends to capture low-frequency signals, while neglecting high-frequency signals. In this context, we associate normal features with low-frequency signals and abnormal features with high-frequency signals. In MRG, where transformers are widely used, this issue is exacerbated by inherent visual and textual biases. To address this fundamental challenge, we introduce a simple method called high-frequency amplification, which amplifies high-frequency signals to better capture abnormal features. We demonstrate that this simple approach effectively debiases frequency bias through extensive experiments, including pseudo-spectrogram analysis, cross-attention analysis, and representation analysis. We evaluate our model on two benchmarks, MIMIC-CXR (Johnson et al., 2019) and IU X-ray (Demner-Fushman et al., 2016). The contributions of our study can be summarized as follows:

- 056 057
- 058 059 060

061

062

063

064 065

- We precisely define visual bias and textual bias, which are crucial but underexplored challenges in MRG. Through comprehensive analysis, we empirically confirm the presence of each bias and show how they exacerbate frequency bias. We emphasize that debiasing and mitigating frequency bias is a fundamental challenge that must be addressed in MRG tasks.
- We introduce a simple yet effective method named high-frequency amplification, specifically designed to mitigate the dominance of normal features in medical images and reports. By amplifying high-frequency signals, which correspond to abnormal features, our approach enables models to effectively capture both global and local patterns.
- We validate the effectiveness of our approach through extensive experiments, including pseudo-spectrogram analysis, cross-attention analysis, and representation analysis. We demonstrate our simple approach achieves performance superior or comparable to state-of-the-art models across both natural language generation and clinical efficacy metrics.
- 068 069 070

067

2 RELATED WORKS

071 072

Most existing MRG methods follow standard image captioning approaches due to the similarities between the two tasks. Despite remarkable success in image captioning models, MRG still faces significant challenges due to severe visual and textual biases inherent in medical images and reports.

Medical images, often captured from consistent angles (e.g., frontal), tend to have similar appear-076 ances but contain subtle, localized abnormal regions. To better identify these abnormal regions, 077 some studies enhanced the alignment between abnormal regions and corresponding disease tags 078 (You et al., 2021; Liu et al., 2021a), generating disease-grounded visual features. Liu et al. (2021b) 079 introduced a differentiated attention mechanism that subtracts common features from the input image, enabling the model to better focus on abnormal regions. Tanida et al. (2023) utilized a scene 081 graph dataset to detect anatomical regions and describe corresponding abnormal regions, enhancing 082 the explainability of the model. All of these prior studies aimed to overcome the limitations of med-083 ical images that are visually biased due to localized abnormal regions. However, none of them have 084 thoroughly analyzed or empirically shown the existence of visual bias.

Medical reports are relatively lengthy, comprising multiple sentences that describe both normal and abnormal findings. Early approaches attempted to generate long reports by integrating relational memory into transformers or incorporating memory matrices. (Chen et al., 2020; 2022). However, these methods often struggled to accurately describe abnormal findings, as they prioritized generating extended narratives over capturing specific abnormalities. To improve the precision of abnormal findings, more recent works have incorporated prior knowledge into MRG models using medical knowledge graphs (Zhang et al., 2020; Liu et al., 2021a; Huang et al., 2023; Li et al., 2023). All of these studies aimed to address so-called textual bias, which has been inconsistently defined sometimes based on text length and at other times on the articulation of abnormal findings. That is, none of the previous works have provided a clear definition of textual bias.

In this paper, we establish precise definitions for visual bias and textual bias and rigorously confirm
 the presence of each bias. We believe that this attempt will promote more focused and productive
 discussions in future MRG research.

098

102

099 3 PRELIMINARIES

3.1 MEDICAL REPORT GENERATION

Advances in deep learning for computer vision (CV) and natural language processing (NLP) have spurred progress in natural image captioning, which involves generating descriptive text given images (Lin et al., 2014). This success has been extended into the healthcare domain, particularly through medical report generation (MRG). MRG aims to assist radiologists by automatically generating diagnostic reports from medical images. The goal of MRG is not only to ensure accurate disease identification but also to generate context-rich reports.



118 normal region with a small abnor- a medical report containing many negative sample (i.e., a normal case), 119 mal region. The red bounding box normal findings (NF) and few ab- while the lower image displays a 120 is the radiologist annotation.



case with multiple diseases.

Figure 1: Illustration of characteristics in the MRG dataset. (a) presents a typical example of a chest X-ray image, highlighting localized abnormal regions. (b) visualizes the imbalanced ratio of normal to abnormal findings. (c) shows two unusual cases where abnormal findings are absent or abundant.

3.2 TERMS AND NOTATIONS

128 $X \in \mathbb{R}^{W \times H \times C}$ represents a medical image, specifically a chest X-ray as shown in Figure 1a, where W, H, and C denote the width, height, and number of channels, respectively. Each medical image is 129 paired with a corresponding medical report $Y = [y_1, \dots, y_t, \dots, y_T] \in \{0, 1\}^{|v|}$, where $y_t \in \mathbb{N}_0^+$ 130 represents the t-th token and |v| indicates the size of vocabulary. The (X, Y)-pair is provided along 131 with a disease label $Z \in \{0,1\}^K$ in which K-1 classes are disease-related and the rest one is a 132 non-disease class (i.e., normal class). Let $X^{(z)}$ and $Y^{(z)}$ represent the abnormal region and finding 133 in the (X, Y)-pair, with their respective size and amount denoted by $|X^{(z)}|$ and $|Y^{(z)}|$, while $X^{(-z)}$ 134 and $Y^{(-z)}$ indicate the normal regions and findings, with their respective amount given by $|X^{(-z)}|$ 135 and $|Y^{(-z)}|$. For the positive samples, i.e., Z|X = 1 and Z|Y = 1, the image and report are defined 136 137 as $X = X^{(z)} \cup X^{(-z)}$ and $Y = Y^{(z)} \cup Y^{(-z)}$, respectively. For the negative samples, i.e., Z|X = 0and Z|Y = 0, each image and report is defined as $X = X^{(-z)}$ and $Y = Y^{(-z)}$, respectively. 138

139 140 141

142

143

121

122

123

124 125 126

127

4 **PROBLEM STATEMENT**

THREE IMBALANCES AND TWO BIASES IN MRG DATASET 4.1

144 Figure 1 illustrates three key imbalances in MRG datasets. First, X-ray images are mostly composed 145 of normal regions, with only a small portion representing abnormal areas. This visual imbalance makes model performance heavily dependent on the normal regions, resulting in visual bias. Second, 146 medical reports are asymmetrically written, with far more sentences describing normal findings than 147 abnormal ones. This textual imbalance causes model performance to rely on the normal findings, 148 leading to textual bias. Finally, the distribution of disease labels is highly skewed; certain diseases 149 are common (e.g., cardiomegaly), while others are relatively rare (e.g., pneumothorax). Such a label 150 imbalance can further deteriorate model performance, but we do not explicitly address it given that 151 mitigating visual and textual biases will inherently resolve this issue. Formal definitions of visual 152 bias and textual bias are provided below. 153

Definition 4.1 (Visual Bias). Let $f_{Z|X}$ denote an image classifier trained to predict a disease label 154 Z given an X-ray image X. Given that the classification accuracy is highly sensitive to the size of the 155 abnormal region $|X^{(z)}|$, the model exhibits a bias towards classifying images as normal. This bias 156 arises because normal regions typically represent global patterns, while abnormal ones are local.

157 **Definition 4.2** (Textual Bias). Let $f_{Z|\hat{Y}}$ denote a text classifier trained to predict a disease label Z 158 from a generated report $\hat{Y} \sim G_{Y|X}(\hat{y}_t|\hat{Y}_{1:t-1},X)$ where $G_{Y|X}$ is a report generator. Given that 159 the classification accuracy is highly sensitive to the number of abnormal findings $|Y^{(z)}|$, the model 160 exhibits a bias towards classifying the generated report as normal. This bias arises because normal 161

findings typically represent global patterns, while abnormal ones are local.



4.2 EXISTENCE OF VISUAL AND TEXTUAL BIASES

175 176

177

208

209

210

211

212

215

178 To demonstrate the existence of visual and textual biases, we analyzed IoU and classification accu-179 racy (e.g., precision and F1) in relation to the size of abnormal regions and the number of abnormal findings. The image encoder and text decoder, followed by $f_{Z|X}$ and $f_{Z|\hat{Y}}$, were examined inde-181 pendently to assess the impact of visual and textual biases, respectively. The IoU (Intersection-over-Union) score quantifies the overlap between ground truth bounding boxes and predicted attention 182 regions, as identified by the Grad-CAM heatmap (Selvaraju et al., 2017; Li et al., 2021; Xiao et al., 183 2023).¹ This metric allows us to evaluate how well the model captures representations relevant to MRG tasks. Classification accuracy measures the performance of the image and text classifiers, 185 $f_{Z|X}$ and $f_{Z|\hat{Y}}$, with values ranging from 0 to 1. A higher score indicates that the image encoder or text decoder has been effectively aligned with MRG tasks. 187

188 Figure 2 presents evidence of visual bias. Specifically, Figure 2a shows that as the size of the bound-189 ing box (i.e., abnormal region) increases, both IoU and classification accuracy improve. This suggests that as the proportion of normal regions in the X-ray increases, the image classifier $f_{Z|X}$ is 190 more likely to misclassify, indicating that the image encoder is influenced by visual bias. This find-191 ing is further supported by Figure 2b. The left plot shows the number of samples misclassified as 192 negative,² suggesting that smaller bounding boxes (i.e., larger normal regions) tend to trigger mis-193 classification. The right plot shows a positive correlation between IoU and F1 scores, implying that 194 reduced attention to abnormal regions increases the likelihood of misclassification. This highlights 195 that the abnormal region size contributes to visual bias. 196

Figure 3 presents evidence of textual bias. The left plot in Figure 3a compares classification accuracy 197 before and after training $G_{Y|X}$, with and without the inclusion of negative samples. The results indicate that the text classifier, $f_{Z|\hat{Y}}$, is more prone to misclassification when negative samples dominate 199 the training data, where the number of abnormal findings is relatively low. The right plot further re-200 inforces this observation: classification accuracy improves as the number of diseases increases. This 201 demonstrates that the number of abnormal findings significantly affects classification accuracy, em-202 phasizing that the text decoder suffers from textual bias.³ Figure 3b confirms the presence of textual 203 bias. The left plot shows the number of samples misclassified as negative samples, suggesting that 204 fewer abnormal findings are more likely to trigger misclassification. The right plot shows a linear 205 correlation between classification accuracy with and without negative samples, with a slope greater 206 than 0.5⁴⁵ This indicates that excluding negative samples improves classification accuracy, further 207 highlighting that the number of abnormal findings is a significant factor contributing to textual bias.

¹See Figure 8 in Appendix A.1.

²Negative samples indicate the cases with no documented abnormal regions and findings. The upper image of Figure 1c showcases an example of a negative sample.

³More diseases typically correspond to more abnormal findings, making multi-disease cases easier to classify correctly. Refer to the lower image in Figure 1c for an example of a multi-disease case.

 ⁴For visual clarity, we grouped the F1 predictions by interval along the y-axis, where the bubble size represents the number of predictions in each group.

⁵Note that the x-axis denotes the group-wise average F1 predictions, and a slope greater than 0.5 indicates that the negative samples are imposing a text bias on the model.

4.3 FREQUENCY BIAS IN TRANSFORMER ARCHITECTURE

5.1 PRETRAINED ENCODER-DECODER NETWORK

As discussed in previous sections, global patterns, such as normal regions and findings, contribute to visual and textual biases. This bias towards global patterns has been extensively studied from the model's perspective, commonly known as frequency bias or spectral bias. *Frequency bias* refers to the phenomenon where models tend to prioritize low-frequency signals that capture global patterns across multiple samples, while neglecting high-frequency signals that represent local patterns unique to each sample (Schwarz et al., 2021; Tian et al., 2023).

224 Transformers (Vaswani, 2017) are particularly vulnerable to frequency bias, as the self-attention module functions as a low-pass filter, inherently paying more attention to low-frequencies than high-225 frequencies (Wang et al., 2022; Park & Kim, 2022; Piao et al., 2024). This globality-seeking behavior 226 of the self-attention module has also been discussed in relation to Principal Component Analysis 227 (PCA) (Zhou et al., 2023; Teo & Nguyen, 2024). Given this, MRG models, where transformers are 228 dominantly used, are especially susceptible to frequency bias, because as described in $\S4.1$ and $\S4.2$, 229 the training data itself is inherently biased towards global patterns. Therefore, mitigating frequency 230 bias is an important and obvious challenge in MRG tasks. The following sections introduce our 231 simple yet powerful approach to addressing this issue. 232

5 Method

234 235

233

236 237

Vision Transformer for Image Encoder Vision Transformer (ViT) (Dosovitskiy, 2020) was the first to successfully apply the transformer architecture directly to image recognition tasks. ViT processes images as sequences of patches, enabling it particularly effective for medical imaging, where abnormalities may span large or subtle regions, such as in X-rays. Accordingly, we implemented an image encoder using the ViT-B model pre-trained on ImageNet (Russakovsky et al., 2015), a widely used approach for medical image encoders. The key ingredient of the ViT encoder is the attention module, encoding each image by aggregating all patchified views. An image embedding, $U \in \mathbb{R}^{N \times |d|}$, processed by a ViT encoder is computed as:

245 246

247

$$U = \text{Attention}(X_p) = \text{softmax}\left(\frac{EW_Q(EW_K)^{\mathrm{T}}}{\sqrt{d}}\right) EW_V \quad \text{where} \quad E = X_p W_E .$$
(1)

Here, $X_p \in \mathbb{R}^{N \times (P^2 \times C)}$ denotes a patchified image sequence with N, P and C as the number of patches, the patch size, and the number of channels, respectively. $W_E \in \mathbb{R}^{(P^2 \times C) \times |d|}$ represents the weight matrix mapping each image to the embedding vector. $W_Q \in \mathbb{R}^{|d| \times |d_q|}$, $W_K \in \mathbb{R}^{|d| \times |d_k|}$, $W_V \in \mathbb{R}^{|d| \times |d|}$ are the query, key, and value weights, respectively, and \sqrt{d} is a scaling factor.

Biomedical GPT for Text Decoder Pre-training models on domain-specific data, such as biomedical text, has been shown to significantly enhance downstream task performance (Peng et al., 2019; Lee et al., 2020; Beltagy et al., 2019). Following this approach, we used a biomedical GPT model as the text decoder. Specifically, we initialized the weights of the text decoder based on Papanikolaou & Pierleoni (2020), which fine-tuned the GPT model using biomedical relations extracted from the PubMed corpus. By doing so, the text decoder can better capture domain-specific details or knowledge and is expected to improve the quality and fluency of medical reports accordingly.

Cross-Attention Module The cross-attention module aligns the image embedding with the text 261 embedding. Specifically, the query vector is derived from the text decoder, while the key and value 262 vectors are sourced from the image encoder. As shown in Eq. (2), the cross-attention mechanism 263 computes attention weights, obtained via the softmax(\cdot) function, by aligning the text embedding 264 V with the image embedding U. These attention weights are then used to re-weight the image 265 embedding, allowing the model to aggregate visual features based on their relevance to the textual 266 context. As a result, the aligned representation $A \in \mathbb{R}^{T \times |d|}$ is generated, representing the fused 267 information of both image and text embeddings in a unified space: 268

269

$$A = \text{Attention}(U, V) = \text{softmax}\left(\frac{VW_Q(UW_K)^{\mathrm{T}}}{\sqrt{d}}\right)UW_V .$$
(2)

270 5.2 HIGH-FREQUENCY AMPLIFICATION LAYER271

As described in §4.3, the attention module introduces an inductive bias, so-called the frequency bias, having transformer-based models less focused on local patterns. To address this, we introduce a *highfrequency amplification layer* (HAL), wherein Fourier transform, high-pass filtering, and inverse Fourier transform are applied subsequently. This layer enhances the model's ability to capture finegrained details, thereby mitigating its bias towards global patterns.

Fourier Transform The Fourier transform decomposes a function into its constituent frequencies using sinusoids as basis functions (Heckbert, 1995). Since both patches and tokens are discrete data, we applied the discrete Fourier Transform (DFT) which is denoted as an operator \mathcal{T} :

291 292 293

294

295

296

297

298

299 300

301 302

307 308

310

311

312

313

314

315

316

 $\mathcal{T}: A \to F$ where $F_c = \sum_{t=0}^{T-1} A_t e^{-\frac{2\pi i}{T}tc}$, $0 \le c \le T-1$.

Here, F_c is the *c*-th frequency component, x_t is the *t*-th time-domain signal, and *i* is the imaginary unit. Computing the DFT directly has a complexity of $O(T^2)$, which is inefficient for large datasets. To overcome this, the Fast Fourier Transform (FFT) was proposed, reducing the complexity to $O(T \log T)$ (Cooley & Tukey, 1965; Brigham, 1988). We apply the FFT to the aligned representation $A \in \mathbb{R}^{T \times |d|}$ using a two-dimensional DFT: one 1D DFT along the time axis, $\mathcal{T}_{\text{time}}$, and another along the feature axis, $\mathcal{T}_{\text{feature}}$, as in (Lee-Thorp et al., 2021; Lee & Lee, 2024). This yields the frequency of the aligned representation denoted as $F \in \mathbb{C}^{T \times |d|}$:

$$F = \mathcal{T} \circ A = \mathcal{T}_{\text{time}}(\mathcal{T}_{\text{feature}}(A))$$
.

High-Pass Filtering and Inverse FFT The frequency representation, F, consists of low and high frequencies, corresponding to global and local patterns, respectively. In our context, global patterns capture normal regions and findings across (X, Y)-pairs, while local patterns represent abnormal ones unique to each pair. High-pass filtering (HPF) is applied to emphasize these local patterns by removing low-frequency components, thus enabling the model to focus on fine-grained details (Pollack, 1948; Costen et al., 1996; Tamkin et al., 2020). Specifically, HPF eliminates frequency components below a certain threshold α by setting $F_{c,d} \leftarrow 0$ for all $c, d \leq \alpha$.⁶ This operation is implemented using a binary mask $F_{\text{HPF}} = F \odot M$, where $M = \{m_{c,d} \mid m_{c,d} \in \{0,1\}, 0 \leq c \leq T - 1, 1 \leq d \leq |d|\}$, with $m_{c,d} = 1$ for high-frequency components and $m_{c,d} = 0$ otherwise. Finally, the original representation A is reconstructed by transforming F_{HPF} back to the original domain through an inverse FFT (iFFT):

$$A_{\rm HPF} = \mathcal{T}^{-1} \circ F_{\rm HPF} = \mathcal{T}_{\rm feature}^{-1}(\mathcal{T}_{\rm time}^{-1}(F_{\rm HPF}))$$

6 EXPERIMENTAL SETUP

Dataset We evaluate our model on two widely used medical report generation benchmarks, i.e., MIMIC-CXR and IU X-ray. 1) **MIMIC-CXR** is the largest radiography dataset with 377,110 chest X-ray images and 227,827 reports from 65,379 patients. We followed the data split and preprocessing steps from (Chen et al., 2020), and used only frontal view images and reports with more than three tokens, resulting in 153,130 images for the training set, 1,201 for the validation set, and 2,193 for the test set. 2) **IU-Xray** is a relatively small public radiography dataset that comprises 7,470 chest X-ray images and 3,955 reports from a total of 3,955 patients. Following the approach of (Chen et al., 2020); Li et al., 2023), we used the dataset only when both frontal and lateral view images were available for each report, resulting in 2,069 images for the training set, 296 for the validation set, and 590 for the test set.

317 318 319

320

Baselines We compare our model with state-of-the-art models on two benchmark datasets. R2Gen (Chen et al., 2020), R2GenCMN (Chen et al., 2022) have been widely used as baseline MRG models.

⁶ α represents the distance from the origin within the 2D frequency space $F \in \mathbb{C}^{T \times |d|}$. In this domain, components closer to the center represent lower frequencies, while those further from the center represent higher frequencies. Therefore, a lower α removes only a small number of low-frequency components near the origin, whereas a higher α eliminates components up to relatively higher frequencies, farther from the origin.

324 AlignTransformer (You et al., 2021), CA (Liu et al., 2021b), RGRG (Tanida et al., 2023) are pro-325 posed to address visual bias, while PPKED (Liu et al., 2021a), KiUT (Huang et al., 2023), DCL (Li 326 et al., 2023) are designed to address textual bias. Additionally, we include R2GenGPT (Wang et al., 327 2023b), METransformer (Wang et al., 2023a), and PromptMRG (Jin et al., 2024) which are widely 328 regarded as SOTA models. For the IU X-ray dataset, we include two additional baselines, CVT2Dis (Nicolson et al., 2023), and M2KT (Yang et al., 2023) which have been used for comparing clinical 329 efficacy performance. 330

331 332

333

334

335

336 337

341

342

344

345

347 348

349 350

351

Evaluation Metrics To measure the fluency and quality of the generated reports, we evaluate them using natural language generation (NLG) metrics, including BLEU (Papineni et al., 2002), METEOR (Denkowski & Lavie, 2011), and ROUGE-L (Lin, 2004).⁷ For the clinical efficacy (CE), we include metrics such as precision, recall, and F1. The CheXbert labeling tool (Smit et al., 2020) is used to convert each report into 14 disease classification labels.

338 **Implementation Details** For the MIMIC-CXR dataset, we use a single frontal view image, while 339 for the IU X-ray dataset, we utilize a pair of images captured from different views of the patient as input. To ensure compatibility across both datasets, all images are resized to 224 and transformed 340 into visual tokens. For the IU X-ray dataset, an additional step is performed where the paired images are concatenated along the embedding dimension and projected back to the original embedding dimension. The hyperparameter α for the HPF is set to 8, and the sensitivity to α is analyzed in 343 Figure 4. We use AdamW optimizer (Loshchilov, 2017) with a learning rate of 5e-6 and a weight decay of 0.05. The learning rate is scheduled using a cosine annealing scheduler, with warm restarts every 5 iterations. The model is trained on an A100 GPU with a batch size of 64 for 39 epochs. 346

7 RESULTS

The results discussed in this section are primarily based on our model trained on the MIMIC-CXR dataset. In §7.3, we provide an additional evaluation conducted on the IU-Xray dataset. We hypothesize that the results will generalize well to other MRG datasets.

352 353 354

355

7.1 GENERALIZATION ASSESSMENT

356 As discussed in §5.2, HAL reconstructs the 357 original feature representation, A, using a lim-358 ited number of filtered high-frequency compo-359 nents, F_{HPF} . Since high-frequency components 360 capture fine-grained local details of the input signals, the reconstructed representation, $A_{\rm HPF}$, 361 may be more prone to overfitting. To assess 362 this risk, we computed the average accuracies 363 and losses and analyzed their trends across both 364 training and validation sets. Figure 4 illustrates the training and validation performance accord-366 ing to different α over 20 epochs. We calculated 367 hit accuracy and categorical cross-entropy loss 368 between the ground truth and predicted tokens. 369 The results indicate that applying HPF with a 370 higher α does not lead to overfitting but con-371 sistently improves generalization. We attribute this outcome to the balancing effect between 372 the low-frequency bias inherent in the model 373



Figure 4: Training and validation performance according to different α .

and the high-frequency bias introduced by HAL, which allows the model to learn balanced rep-374 resentations that enhance its emergent generalizability. Based on this result, the default setting for 375 the analyses in the following sections is fixed at $\alpha = 8$. 376

³⁷⁷

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





Figure 5: Comparison of neuron activation inten- Figure 6: Comparison of cross-attention map sity before and after HAL ($\alpha = 8$)

without and with HAL ($\alpha = 8$)



Figure 7: Comparison of T-SNE embeddings before and after HAL ($\alpha = 8$)

7.2 ABLATION STUDIES

In this section, we present a comprehensive analysis of the impact of HAL on internal model dynamics, focusing on key aspects such as neuron activation patterns, cross-attention distributions, and representation topology.

407 408

387

388

389 390

391

392

393

394

395

396

397

398

399 400

401 402 403

404 405

406

409 Pseudo-spectrogram Analysis A spectrogram provides a visual representation of how frequency 410 components evolve over time, typically depicting frequency on the x-axis and time on the y-axis, 411 with color indicating the intensity of each frequency component. Inspired by this approach, we con-412 ducted a pseudo-spectrogram analysis of neuron activation. Figure 5 compares the neuron activation intensity before and after HAL, where the x-axis represents the top 30 neurons ranked by activation 413 level and the y-axis denotes the temporal sequence of tokens. The figure shows that in the layer 414 before HAL, only a subset of neurons are strongly activated, with most neurons remaining inactive. 415 Furthermore, those few active neurons exhibit uniform activation across all tokens in the sequence, 416 suggesting an indistinguishable activation pattern. In contrast, there is a rich and non-uniform acti-417 vation after HAL. An activation spectrum indicates that HAL allows the model to effectively cap-418 ture fine-grained details by amplifying high-frequency signals, which might otherwise be ignored. 419 Consequently, HAL produces a richer representation so that neuron activation forms a spectrum, 420 ultimately improving the model's discriminative perceptiveness.

421

422 **Cross-attention Analysis** HAL is placed after the cross-attention layer, making it highly depen-423 dent on the influence of HAL. Therefore, comparing the cross-attention map with and without HAL 424 helps illustrate how it has affected the image-to-text alignment. Figure 6 shows a comparison of the 425 cross-attention distributions across (224×224) images for models trained with and without HAL. 426 The results reveal clear advantages of using HAL. In the model without HAL, cross-attention tends 427 to focus on the periphery, especially the mid-abdominal part, which contains little information about 428 chest disease. This may be due to the "common" appearance of grey areas around the abdomen on 429 most X-ray images (see Figures 1 and 8), representing a global pattern across all samples regardless of disease type. On the other hand, the model trained with HAL shows that the image-to-text 430 cross-attention is concentrated around the center of the image (i.e., the upper-mid-thoracic region), 431 typically containing the "specific" information of chest disease. This may be understood as evidence

435 NLG Metrics CE Metrics 436 Baselines BLEU-1 BLEU-2 BLEU-3 BLEU-4 METEOR ROUGE-L F1 Precision Recall 437 R2Gen 0.353 0.218 0.145 0.103 0.142 0.277 0.333 0.273 0.276 R2GenCMN 0.353 0.142 0.218 0 1 4 8 0 106 0 278 0.334 0.275 0.278 438 0.378 0.235 0.156 0.112 0.158 0.283 AlignTransformer 439 CA 0.350 0.219 0.152 0.109 0.151 0.283 0.352 0.298 0.303 RGRG 0.373 0.249 0.175* 0.126* 0.168 0.264 0.461* 0.475 0.447 440 PPKED 0.360 0.2240.149 0.106 0.149 0.284 441 0.371 0.321 0.243 0.160* 0.318 KiUT 0.393 0.159 0.113 0.285 442 DCL 0.109 0.150 0.284 0.471 0.352 0.373 0.411 0.267 R2GenGPT 0.186 0.134 0.160* 0.297 0.392 0.387 0.389 443 0.386 0.152 0.291* 0.364 METransformer 0.250* 0.169 0.124 0.309 0.311 444 PromptMRG 0 398* 0 1 5 7 0.501 0.509 0.112 0.268 0.476 0.264 0.189 0.299 445 Ours (HAL) 0.399 0.143 0.170 0.434 0.410* 0.392*

432 Table 1: The comparison of model performance on MIMIC-CXR dataset. Note that **bold** numbers 433 highlight the best performance, underlined numbers indicate the second-best performance, and as-434 terisked (*) numbers denote the third-best performance, respectively.

Table 2: The comparison of model performance on IU-Xray dataset. Note that **bold** numbers highlight the best performance, underlined numbers indicate the second-best performance, and asterisked (*) numbers denote the third-best performance, respectively.

Decelines	NLG Metrics					CE Metrics			
Dasennes	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	Precision	Recall	F1
R2Gen	0.470	0.304	0.219	0.165	0.187	0.371	0.141	0.136	0.136
R2GenCMN	0.475	0.309	0.222	0.170	0.191	0.375	-	-	-
AlignTransformer	0.484	0.313	0.225	0.173	0.204	0.379	-	-	-
CA	0.492	0.314	0.222	0.169	0.193	0.381	-	-	-
RGRG	-	-	-	-	-	-	0.183*	0.187*	0.180*
PPKED	0.483	0.315	0.224	0.168	0.190	0.376	-	-	-
KiUT	0.525	0.360	0.251	0.185	0.242	0.409	-	-	-
DCL	-	-	-	0.163	0.193	0.383	0.168	0.167	0.162
R2GenGPT	0.488	0.316	0.228	0.173	0.211*	0.377	-	-	-
METransformer	0.483	0.322*	0.228	0.172	0.192	0.38	-	-	-
PromptMRG	0.401	-	-	0.098	0.160	0.281	0.213	0.229	0.211
CVT2Dis	0.473	0.304	0.224	0.175*	0.200	0.376	0.174	0.172	0.168
M2KT	0.497*	0.319	0.230*	0.174	-	0.399*	0.153	0.145	0.145
Ours (HAL)	0.521	0.425	0.371	0.336	0.263	0.507	0.418	0.415	0.414

that HAL enhances robustness to the frequency bias. In summary, Figure 6 demonstrates that HAL improved the model to attend to diagnostically significant regions by mitigating the frequency bias.

466 **Representation Analysis** Comparing the topology of representation before and after a specific 467 layer provides an intuitive explanation of how it works as an operator, and proves the utility it yields 468 from the perspective of representation quality. In this regard, we performed the T-SNE embedding 469 (Van der Maaten & Hinton, 2008) and visualized representation for both single-disease and multi-470 disease cases, as shown in Figure 7. For single-disease cases (see Figure 7a), which encompass 12 471 distinct diseases, the embedding vector before HAL produces entangled clusters, indicating poor 472 feature discrimination by diseases. In contrast, the embeddings after HAL form well-separated clusters, suggesting a marked improvement in representation quality. This improvement is likely due to 473 HAL, where amplified high-frequency signals highlight the local patterns unique to each sample, 474 but erase the global patterns shared across samples, which contribute to entangled representations. 475

476 For multi-disease cases (see Figure 7b), we cannot conduct cluster analysis as T-SNE embeddings 477 fail to build distinguishable representations due to the high complexity of disease combinations—the complex nature of these combinations results in highly entangled feature representations-making 478 it challenging to achieve well-separated clusters.⁸ Instead, we can do scatter analysis to demonstrate 479 whether HAL makes a dispersed representation—the larger dispersion means that the model treats 480 each point more uniquely. Before HAL, the T-SNE embeddings show compact representation, except 481 for the |Z| = 4 case that exhibits dispersed representation. This dispersion is likely due to the 482

483 ⁸In multi-disease cases, there are many samples that have the same disease in common. In these cases, 484 the overlapping diseases among samples dilute or mix the distinctive local patterns. That is, the local patterns become nothing but noise, and only the global patterns survive. As a result, the frequency bias in models 485 becomes more pronounced compared to single-disease scenarios, making it challenging to do cluster analysis.

9

446

447

448

486 reduced influence of normal findings, as illustrated in Figure 1c, enabling clearer differentiation 487 between samples. After HAL, the embeddings appear to spread more dispersely, implying that each 488 sample is embedded finely enough to be distinguishable. That is, the scatter analysis suggests that 489 HAL increases the model's perceptiveness to the local details and thus mitigates frequency bias.

490 491

492

7.3 PERFORMANCE COMPARISON

493 Through the previous sections has it been shown that HAL is an effective tool for addressing fre-494 quency bias, especially crucial in MRG tasks where visual and textual biases are already prevalent. Nevertheless, one might argue that HAL, due to its simplicity, may not be competitive against exist-495 ing methods for MRG. To address this concern, this section presents a comparative evaluation on two 496 MRG benchmark datasets: MIMIC-CXR and IU-Xray. Tables 1 and 2 summarize the performance 497 of MRG models using both NLG and CE metrics, illustrating that our model performs competi-498 tively against other baselines. Specifically, our model outperforms baseline models on NLG metrics, 499 demonstrating its ability to generate high-quality reports containing featured medical terminology 500 found in real-world clinical texts. Furthermore, when comparing the top-3 ranks for each metric, our 501 model ranks consistently high across almost all metrics. This balanced achievement across diverse 502 metrics suggests that HAL not only enhances the overall quality of generated reports but also pro-503 vides robustness in capturing key clinical concepts, making it a reliable tool for MRG tasks. Note 504 that the performance of baselines was taken directly from the results reported in the original papers.⁹

505 506 507

508

LIMITATIONS AND FUTURE WORK 8

It is important to note that our current results were obtained without extensive hyperparameter op-509 timization. We believe that a systematic exploration of hyperparameters could further enhance the 510 model's performance and stability, providing stronger evidence of HAL's effectiveness. The primary 511 goal of HAL is to reduce the impact of global patterns by amplifying high-frequency signals. How-512 ever, this approach may be less effective if the training data is either insufficient or contains too much 513 randomness, making it difficult for dominant global patterns to emerge. In such cases, HAL might 514 even introduce a bias towards local patterns instead. Accordingly, future research should explore 515 methods to balance global and local patterns, especially when training data is limited or noisy. In 516 addition, while this study empirically links visual and textual biases with frequency bias, additional 517 theoretical grounding is needed to strengthen those empirical findings. We hope that future research 518 will further explore this area. Meanwhile, further improvements could also be explored in the pretraining phase. We anticipate that pre-training the encoder or decoder models on chest X-ray data 519 would yield greater performance. 520

521 522

523

531

9 CONCLUSION

524 In this work, we demonstrated the existence of visual and textual biases in the MRG dataset ($\S4.1$) 525 and discussed how these biases make MRG models especially prone to frequency bias, with a ten-526 dency to prioritize low-frequency components. To counter this vulnerability, we introduced the high-527 frequency amplification layer (HAL) (§5.2), designed to mitigate the model's predisposition towards 528 such biases. Our results showed that HAL significantly enhances various aspects, including neuron 529 activation, cross-attention map, and representation quality, as detailed in ablation studies ($\S7.2$). Despite its simplicity, HAL exhibited outstanding performance in comparative evaluations (§7.3). All 530 these findings strongly support our arguments that: (1) debiasing is a fundamental issue for improving MRG tasks, and (2) mitigating frequency bias is crucial for enabling models to capture both 532 global and local patterns of medical images more effectively. We believe this work will pave the 533 way for more robust MRG models that are better equipped to handle the complexities of real-world 534 medical data, ultimately contributing to advanced medical imaging tasks. 535

⁵³⁶ ⁹It is important to note that only Jin et al. (2024) reported CE performance on the IU-Xray dataset. Therefore, the CE metrics presented in Table 2 are all borrowed from the results reported in Jin et al. (2024). Unlike 538 the NLG metrics, the CE metrics were evaluated in a zero-shot setting, where the models were not trained on the IU-Xray dataset. This context explains why the baseline models exhibit relatively lower CE performance compared to our model.

540	REFERENCES
541	

550

565

- Iz Beltagy, Kyle Lo, and Arman Cohan. Scibert: A pretrained language model for scientific text. 542 arXiv preprint arXiv:1903.10676, 2019. 543
- 544 E Oran Brigham. The fast Fourier transform and its applications. Prentice-Hall, Inc., 1988.
- Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. Generating radiology reports via 546 memory-driven transformer. arXiv preprint arXiv:2010.16056, 2020. 547
- 548 Zhihong Chen, Yaling Shen, Yan Song, and Xiang Wan. Cross-modal memory networks for radiol-549 ogy report generation. arXiv preprint arXiv:2204.13258, 2022.
- James W Cooley and John W Tukey. An algorithm for the machine calculation of complex fourier 551 series. Mathematics of computation, 19(90):297-301, 1965. 552
- 553 Nicholas P Costen, Denis M Parker, and Ian Craw. Effects of high-pass and low-pass spatial filtering 554 on face identification. Perception & psychophysics, 58:602-612, 1996.
- Dina Demner-Fushman, Marc D Kohli, Marc B Rosenman, Sonya E Shooshan, Laritza Rodriguez, 556 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 558 Association, 23(2):304-310, 2016. 559
- Michael Denkowski and Alon Lavie. Meteor 1.3: Automatic metric for reliable optimization and 560 evaluation of machine translation systems. In Proceedings of the sixth workshop on statistical 561 machine translation, pp. 85–91, 2011. 562
- 563 Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. 564 arXiv preprint arXiv:2010.11929, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-566 nition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 567 770-778, 2016. 568
- Paul Heckbert. Fourier transforms and the fast fourier transform (fft) algorithm. *Computer Graphics*, 2(1995):15-463, 1995. 570
- 571 Zhongzhen Huang, Xiaofan Zhang, and Shaoting Zhang. Kiut: Knowledge-injected u-transformer 572 for radiology report generation. In Proceedings of the IEEE/CVF Conference on Computer Vision 573 and Pattern Recognition, pp. 19809–19818, 2023. 574
- Yun Jiang, Yuan Zhang, Xin Lin, Jinkun Dong, Tongtong Cheng, and Jing Liang. Swinbts: A method 575 for 3d multimodal brain tumor segmentation using swin transformer. Brain sciences, 12(6):797, 576 2022. 577
- 578 Haibo Jin, Haoxuan Che, Yi Lin, and Hao Chen. Promptmrg: Diagnosis-driven prompts for medical 579 report generation. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pp. 2607-2615, 2024. 580
- 581 Alistair EW Johnson, Tom J Pollard, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ving Deng, 582 Yifan Peng, Zhiyong Lu, Roger G Mark, Seth J Berkowitz, and Steven Horng. Mimic-cxr-jpg, a 583 large publicly available database of labeled chest radiographs. arXiv preprint arXiv:1901.07042, 584 2019. 585
- Davood Karimi, Serge Didenko Vasylechko, and Ali Gholipour. Convolution-free medical im-586 age segmentation using transformers. In Medical image computing and computer assisted 587 intervention-MICCAI 2021: 24th international conference, Strasbourg, France, September 27-588 October 1, 2021, proceedings, part I 24, pp. 78-88. Springer, 2021. 589
- 590 Changhun Lee and Gyumin Lee. Repurformer: Transformers for repurposing-aware molecule generation. In Carl Edwards, Qingyun Wang, Manling Li, Lawrence Zhao, Tom Hope, and Heng Ji (eds.), Proceedings of the 1st Workshop on Language + Molecules (L+M 2024), pp. 116–127, 592 Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/ 2024.langmol-1.14. URL https://aclanthology.org/2024.langmol-1.14.

594 Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jae-595 woo Kang. Biobert: a pre-trained biomedical language representation model for biomedical text 596 mining. Bioinformatics, 36(4):1234-1240, 2020. 597 James Lee-Thorp, Joshua Ainslie, Ilya Eckstein, and Santiago Ontanon. Fnet: Mixing tokens with 598 fourier transforms. arXiv preprint arXiv:2105.03824, 2021. 600 Mingjie Li, Wenjia Cai, Rui Liu, Yuetian Weng, Xiaoyun Zhao, Cong Wang, Xin Chen, Zhong Liu, 601 Caineng Pan, Mengke Li, et al. Ffa-ir: Towards an explainable and reliable medical report gen-602 eration benchmark. In Thirty-fifth conference on neural information processing systems datasets 603 and benchmarks track (round 2), 2021. 604 Mingjie Li, Binggian Lin, Zicong Chen, Haokun Lin, Xiaodan Liang, and Xiaojun Chang. Dynamic 605 graph enhanced contrastive learning for chest x-ray report generation. In Proceedings of the 606 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3334–3343, 2023. 607 608 Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In Text summarization branches out, pp. 74-81, 2004. 609 610 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr 611 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer 612 Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Pro-613 ceedings, Part V 13, pp. 740-755. Springer, 2014. 614 Fenglin Liu, Xian Wu, Shen Ge, Wei Fan, and Yuexian Zou. Exploring and distilling posterior and 615 prior knowledge for radiology report generation. In Proceedings of the IEEE/CVF conference on 616 computer vision and pattern recognition, pp. 13753–13762, 2021a. 617 618 Fenglin Liu, Changchang Yin, Xian Wu, Shen Ge, Yuexian Zou, Ping Zhang, and Xu Sun. Con-619 trastive attention for automatic chest x-ray report generation. arXiv preprint arXiv:2106.06965, 620 2021b. 621 I Loshchilov. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017. 622 623 Omid Nejati Manzari, Hamid Ahmadabadi, Hossein Kashiani, Shahriar B Shokouhi, and Ahmad 624 Ayatollahi. Medvit: a robust vision transformer for generalized medical image classification. 625 Computers in Biology and Medicine, 157:106791, 2023. 626 Christos Matsoukas, Johan Fredin Haslum, Magnus Söderberg, and Kevin Smith. Is it time to replace 627 cnns with transformers for medical images? arXiv preprint arXiv:2108.09038, 2021. 628 629 Aaron Nicolson, Jason Dowling, and Bevan Koopman. Improving chest x-ray report generation by 630 leveraging warm starting. Artificial intelligence in medicine, 144:102633, 2023. 631 Yannis Papanikolaou and Andrea Pierleoni. Date: Data augmented relation extraction with gpt-2. 632 arXiv preprint arXiv:2004.13845, 2020. 633 634 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic 635 evaluation of machine translation. In Proceedings of the 40th annual meeting of the Association for Computational Linguistics, pp. 311–318, 2002. 636 637 Namuk Park and Songkuk Kim. How do vision transformers work? arXiv preprint 638 arXiv:2202.06709, 2022. 639 Yifan Peng, Shankai Yan, and Zhiyong Lu. Transfer learning in biomedical natural language 640 processing: an evaluation of bert and elmo on ten benchmarking datasets. arXiv preprint 641 arXiv:1906.05474, 2019. 642 643 Xihao Piao, Zheng Chen, Taichi Murayama, Yasuko Matsubara, and Yasushi Sakurai. Fredformer: 644 Frequency debiased transformer for time series forecasting. In Proceedings of the 30th ACM 645 SIGKDD Conference on Knowledge Discovery and Data Mining, pp. 2400–2410, 2024. 646 Irwin Pollack. Effects of high pass and low pass filtering on the intelligibility of speech in noise. 647 The Journal of the Acoustical Society of America, 20(3):259–266, 1948.

648 649 650	Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? <i>Advances in neural information processing systems</i> , 34:12116–12128, 2021.
652 653 654	Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. <i>International journal of computer vision</i> , 115:211–252, 2015.
655 656 657	Katja Schwarz, Yiyi Liao, and Andreas Geiger. On the frequency bias of generative models. <i>Advances in Neural Information Processing Systems</i> , 34:18126–18136, 2021.
658 659 660 661	Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based local- ization. In <i>Proceedings of the IEEE international conference on computer vision</i> , pp. 618–626, 2017.
662 663 664	Akshay Smit, Saahil Jain, Pranav Rajpurkar, Anuj Pareek, Andrew Y Ng, and Matthew P Lungren. Chexbert: combining automatic labelers and expert annotations for accurate radiology report la- beling using bert. <i>arXiv preprint arXiv:2004.09167</i> , 2020.
666 667 668	Alex Tamkin, Dan Jurafsky, and Noah Goodman. Language through a prism: A spectral approach for multiscale language representations. <i>Advances in Neural Information Processing Systems</i> , 33: 5492–5504, 2020.
669 670 671 672	Tim Tanida, Philip Müller, Georgios Kaissis, and Daniel Rueckert. Interactive and explainable region-guided radiology report generation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 7433–7442, 2023.
673 674	Rachel SY Teo and Tan M Nguyen. Unveiling the hidden structure of self-attention via kernel principal component analysis. <i>arXiv preprint arXiv:2406.13762</i> , 2024.
675 676 677 678	Yuandong Tian, Yiping Wang, Beidi Chen, and Simon S Du. Scan and snap: Understanding train- ing dynamics and token composition in 1-layer transformer. <i>Advances in Neural Information</i> <i>Processing Systems</i> , 36:71911–71947, 2023.
679 680	Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. Journal of machine learning research, 9(11), 2008.
681 682	A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017.
683 684 685	Peihao Wang, Wenqing Zheng, Tianlong Chen, and Zhangyang Wang. Anti-oversmoothing in deep vision transformers via the fourier domain analysis: From theory to practice. <i>arXiv preprint arXiv:2203.05962</i> , 2022.
687 688 689	Zhanyu Wang, Lingqiao Liu, Lei Wang, and Luping Zhou. Metransformer: Radiology report gen- eration by transformer with multiple learnable expert tokens. In <i>Proceedings of the IEEE/CVF</i> <i>Conference on Computer Vision and Pattern Recognition</i> , pp. 11558–11567, 2023a.
690 691 692	Zhanyu Wang, Lingqiao Liu, Lei Wang, and Luping Zhou. R2gengpt: Radiology report generation with frozen llms. <i>Meta-Radiology</i> , 1(3):100033, 2023b.
693 694 695	Julia Wolleb, Florentin Bieder, Robin Sandkühler, and Philippe C Cattin. Diffusion models for med- ical anomaly detection. In <i>International Conference on Medical image computing and computer-</i> <i>assisted intervention</i> , pp. 35–45. Springer, 2022.
696 697 698 699	Junfei Xiao, Yutong Bai, Alan Yuille, and Zongwei Zhou. Delving into masked autoencoders for multi-label thorax disease classification. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pp. 3588–3600, 2023.
700 701	Shuxin Yang, Xian Wu, Shen Ge, Zhuozhao Zheng, S Kevin Zhou, and Li Xiao. Radiology report generation with a learned knowledge base and multi-modal alignment. <i>Medical Image Analysis</i> , 86:102798, 2023.

702 703 704 705 706	Di You, Fenglin Liu, Shen Ge, Xiaoxia Xie, Jing Zhang, and Xian Wu. Aligntransformer: Hier- archical alignment of visual regions and disease tags for medical report generation. In <i>Medical</i> <i>Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Con-</i> <i>ference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part III 24</i> , pp. 72–82. Springer, 2021.
707 708 709 710	Yixiao Zhang, Xiaosong Wang, Ziyue Xu, Qihang Yu, Alan Yuille, and Daguang Xu. When ra- diology report generation meets knowledge graph. In <i>Proceedings of the AAAI conference on</i> <i>artificial intelligence</i> , volume 34, pp. 12910–12917, 2020.
711 712 713 714	Zhicheng Zhang, Lequan Yu, Xiaokun Liang, Wei Zhao, and Lei Xing. Transct: dual-path trans- former for low dose computed tomography. In <i>Medical Image Computing and Computer Assisted</i> <i>Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–</i> <i>October 1, 2021, Proceedings, Part VI 24</i> , pp. 55–64. Springer, 2021.
715 716 717	Daquan Zhou, Zhiding Yu, Enze Xie, Chaowei Xiao, Animashree Anandkumar, Jiashi Feng, and Jose M Alvarez. Understanding the robustness in vision transformers. In <i>International Conference on Machine Learning</i> , pp. 27378–27394. PMLR, 2022.
718 719 720 721	Tian Zhou, Peisong Niu, Liang Sun, Rong Jin, et al. One fits all: Power general time series analysis by pretrained lm. <i>Advances in neural information processing systems</i> , 36:43322–43355, 2023.
722	
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
740	
747	
748	
749	
750	
751	
752	
753	
754	
755	

756 APPENDIX А

758 VISUAL BIAS EXPERIMENTS A.1 759

760 According to existing studies (Liu et al., 2021b; You et al., 2021; Tanida et al., 2023), localized abnormal regions are difficult to capture. In this experiment, we investigated several potential factors 761 that may affect visual bias, including the bounding box size. We used VinDr-CXR dataset which 762 provides 18,000 chest X-ray images annotated with bounding box information for disease regions and 23 disease labels.¹⁰ We utilized ViT-S pre-trained on 510K X-ray images for this section (Xiao 764 et al., 2023). To evaluate the model's ability to capture disease-relevant regions, we employed IoU 765 (Intersection-over-Union) metric. Specifically, we utilized the IoU Multiple metric, which compares 766 all actual bounding boxes to all predicted bounding boxes, and the IoU Largest metric, which com-767 pares the largest actual bounding box and the largest predicted bounding box. Classification perfor-768 mance was assessed using precision and F1.11 769

770 **Grad-CAM in ViTs** In this paragraph, we briefly explain how we calculated the IoU score based 771 on Grad-CAM. Grad-CAM (Selvaraju et al., 2017) is a method for generating visual explanations 772 to identify which parts of the input image have the most influence on a given class prediction. 773 It was originally proposed for CNN architecture and is now applicable to ViTs as well. Similar 774 to CNN-based models, which extract feature maps from the last convolutional layer, Grad-CAM for ViTs extracts feature maps from the norm1 layer of the final block. These feature maps have 775 the shape of $f \in \mathbb{R}^{N \times (P^2 \times C)}$, where N, P, and C represent the number of patches, the patch 776 size, and the number of channels, respectively. The feature maps are then projected back into the 777 original image space by reversing the patchification process, allowing for spatial interpretation. In 778 this experiment, we extracted activation maps based on the ground truth disease label. To generate 779 predicted bounding boxes, we retained only the regions of the activation maps that exceed 75% of the maximum activation value. We then utilized OpenCV's findContours function to detect the contours 781 of these regions, followed by the minAreaRect function to generate the minimum area bounding 782 rectangles that enclose each contour (see Figure 8). 783



Figure 8: Bounding box generation from Grad-CAM visualizations

Results In $\S4.2$, we demonstrated that as the (1) bounding box size (i.e., abnormal regions) de-794 creases, both the ability to capture abnormal regions and accuracy of disease classification deterio-795 rate. Additionally, the number of positive samples misclassified as negative increases, highlighting 796 the presence of visual bias. We then analyzed the effect of (2) biased distribution of diseases on 797 the model performance. In Figure 9a, the IoU and classification scores of the vanilla classification 798 model for different diseases are sorted in ascending order based on the number of training samples. 799 Notably, disease ID 22 represents the negative samples. While classification performance tends to improve with an increasing number of training samples, the IoU score does not consistently reflect 800 the classification accuracy. To examine potential correlations, we plotted a regression using bootstrapped samples, ensuring a minimum of 100 samples per class and down-weighting outliers. As 802 shown in Figure 9b, the regression plot reveals no significant correlation between IoU and classifi-803 cation performance. The above analysis is consistently observed when comparing performance by the (3) number of diseases per image, as shown in Figure 10. For cases where medical images 805 contain multiple diseases (excluding classes with fewer than 30 samples), the results show that as 806 the number of diseases per sample increases, the classification performance improves, while the IoU

801

804

790

791 792 793

¹¹In this analysis, precision is calculated by considering only the top-k predicted labels as the predicted disease, where k corresponds to the number of ground truth disease labels for each sample.

⁸⁰⁷ 808

¹⁰https://physionet.org/content/vindr-cxr/1.0.0/



Figure 11: Location of bounding boxes



scores remain unaffected. We further analyzed whether the (4) location of bounding box might influence the model performance. The image is divided into nine areas $(3 \times 3 \text{ grid})$ where the y-axis represents upper, middle, and lower regions, and the x-axis represents left, center, and right regions. In Figure 11b, a weak positive correlation is observed due to outliers, but in Figure 11a, the upper center (UC), middle center (MC), and lower right (LR) areas show significant differences in IoU performance, despite having nearly identical classification performance. This implies there is little correlation between these two metrics. Additionally, no significant patterns were identified when analyzing the (5) scatterness of bounding box based on the number of differently located bounding boxes, as shown in Figure 12. This experiment confirmed that the primary factor influencing visual bias is the size of the bounding box.

848 A.2 **TEXTUAL BIAS EXPERIMENTS**

836 837

838

839

840

841

842

843

844

845

846 847

849 In this experiment, we investigated several potential factors that may contribute to textual bias in 850 MRG, including the number of abnormal findings. We used the MIMIC-CXR dataset and the base-851 line MRG model without any HPF. We prepared ground truth labels using CheXbert labeling tool 852 (Smit et al., 2020) for the analysis, but the "support devices" label was excluded as it does not repre-853 sent an actual disease. The results were evaluated using precision and F1 score for CE metrics, and 854 BLEU-4, METEOR, and ROUGE-L for NLG metrics. 855

856 **Results** In $\S4.2$, we demonstrate that as the (1) number of abnormal findings decreases, the 857 model's diagnostic performance degrades. In other words, the model exhibits a textual bias towards 858 dominant normal findings. This is further supported by the increasing misclassification tendency 859 as negative samples dominate the training data. Additionally, we analyzed the effect of (2) biased 860 distribution of diseases. As shown in Figure 13a, the NLG and CE scores of MRG model for 861 different diseases are sorted in ascending order based on the number of training samples, with disease label 13 representing negative samples. In contrast to the visual bias experiment, both NLG and 862 CE metrics showed no significant differences for abnormal diseases, except in the case of negative 863 samples. Additionally, Figure 13b does not appear to have a strong correlation. Next, we analyzed





Figure 14: Length of reports

the impact of the (3) report length, categorizing reports as short (20 words or less), medium (up to 50 words), long (up to 80 words), and extra long (more than 80 words) given the max length is 100. As shown in Figure 14a, short to medium-length generated reports tend to perform better on both NLG and CE metrics. However, the scores do not show any clear upward or downward trend based on the generated length, suggesting that the results may not be statistically significant. Figure 14b shows a positive correlation between the NLG score and CE score. This experiment confirmed that the primary factor influencing textual bias is the number of abnormal findings.

A.3 VIT OUTPERFORMS RESNET IN MEDICAL IMAGING TASKS

888 Convolutional Neural Networks (CNNs) have been widely used in various computer vision tasks. 889 After the advent of Vision Transformers (ViTs), ViTs have shown its potential as a competitive 890 alternative to CNNs. Since CNNs and ViTs each exhibit distinct advantages and limitations, it is 891 general to choose the appropriate backbone model based on the downstream tasks. For instance, 892 CNNs possess a high inductive bias, while ViTs effectively capture long-range dependencies. Al-893 though CNN-based models might seem suitable for the localized nature of abnormalities in medical images, numerous studies have demonstrated the effectiveness of ViTs in automated medical image 894 diagnosis, ranging from medical image segmentation (Karimi et al., 2021), medical image classifi-895 cation (Matsoukas et al., 2021), and medical image reconstruction (Zhang et al., 2021), especially 896 when pre-trained on ImageNet. Since diagnosis often requires consideration of distant organs and 897 tissues, the long-range dependencies captured by ViTs are particularly advantageous in the medical 898 domain. In medical images, abnormal regions may span large or subtle regions. Unlike CNNs, which 899 focus on local features in the lower layers and global features in the higher layers, ViTs capture 900 both local and global features at every layer, preserving fine-grained details and contextual relation-901 ships (Raghu et al., 2021). Additionally, sparse attention mechanisms in ViTs are known to enhance 902 robustness against noise (Zhou et al., 2022). ViTs perform particularly well when pre-trained on 903 large-scale datasets or via self-supervision.

This finding is further validated through our simple experiment. Specifically, we compared the classification performance of ResNet-50 (He et al., 2016) and ViT-S (Dosovitskiy, 2020), both pretrained on ImageNet and fine-tuned on the Vindr-CXR dataset for 50 epochs. The results demonstrated that ViT-S outperformed ResNet-50, despite the latter having a slightly higher number of parameters (See Table 3).

Model	Params	F1	Hit@k
ResNet-50	23.9M	0.682	0.546
ViT-S	22.2M	0.699	0.624

915

916 917

877 878 879

880

882

883

884

885 886