Report Generation

Anonymous ACL submission

Abstract

Automated radiology report generation (RRG) offers the potential to reduce clinical workload and enhance diagnostic consistency. However, existing models struggle with degraded visual 006 representations caused by long-tailed lesion distributions, and suffer from limited alignment between image features and diagnostic semantics. We propose **VDGen**, a unified framework for calibrating visual and diagnostic representations to improve disease-aware report generation. VDGen integrates two complementary modules: Vision Self-Equilibration (VSE), a self-supervised contrastive module that mitigates visual feature degradation by promoting structured representation learning; and Disease Information Distillation (DID), a cross-modal distillation mechanism that uses diagnostic reports as teacher signals to guide the extraction 020 of disease-sensitive semantics from visual features. Our end-to-end architecture incorporates a LoRA-adapted large language model (LLM) decoder to generate clinically accurate reports. Experiments on the IU-Xray and MIMIC-CXR datasets show that VDGen achieves state-ofthe-art performance on MIMIC-CXR and maintains competitive results on IU-Xray. Code and models will be released upon acceptance.

1 Introduction

001

002

017

022

024

035

040

042

043

Automated radiology report generation (RRG) holds significant promise for alleviating radiologists' workload while improving the interpretability and consistency of diagnostic reports (Wang and Summers, 2012). With rapid advancements in artificial intelligence, automated RRG has emerged as a critical research frontier in multimodal AI. Existing studies have attempted to adapt image caption methods to RRG by training models to generate medical reports from chest X-ray images through supervised learning (Liu et al., 2021a). However, such direct migration faces inherent limitations due to the unique characteristics of radiological data (Gu et al., 2024).

As illustrated in Fig. 1(a), lesion-related pixels in chest X-ray images typically occupy only small regions, resulting in highly imbalanced intra-image distributions. At the dataset level (Fig. 1(b)), normal samples overwhelmingly dominate, limiting the model's exposure to abnormal findings. These dual imbalances-spatial sparsity within images and class imbalance across the dataset-collectively contribute to a longtailed data distribution that challenges conventional vision-language models (Bu et al., 2024). This distributional bias leads to degraded visual representations, especially for abnormal regions, where traditional visual encoders struggle to extract discriminative features (Liu et al., 2021b). As shown in Fig. 1(c, left), embeddings of lesion areas are often disorganized in the latent space, which in turn confuses the report decoder and impairs the generation of clinically meaningful text. We argue that effective radiology report generation requires addressing both: (i) intra-modal degradation of visual features caused by distributional imbalance, and (ii) inter-modal misalignment between image features and diagnostic language semantics leads to difficulties in disease information extraction.

044

045

046

047

051

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

078

081

In this paper, we propose **VDGen**, a unified framework for Visual and Diagnostic representation calibration in radiology report generation. VD-Gen integrates two complementary modules: Vision Self-Equilibration (VSE) introduces a selfsupervised contrastive mechanism that learns to pull together visual embeddings from similar pathological patterns while pushing apart those of dissimilar ones. By performing instance-level regularization within each batch, VSE mitigates the effects of long-tailed and sparse lesion distributions, enhancing intra-modal discrimination without relying on manual labels. Disease Information Distillation (DID) introduces a cross-modal distillation mechanism, where diagnostic reports serve as teacher signals to supervise the learning of disease-



- We propose **VDGen**, a unified framework for calibrating visual and diagnostic representations to improve the quality and interpretability of radiology report generation.

out relying on manual labels.

2 Related Work

2.1 Image Caption

Image caption aims to automatically generate text descriptions of input natural images. The com-mon approach in this field is to use an end-to-end framework for caption generation. After the Trans-former architecture became popular, researchers often used Swin-Transformer for image encoding, and then employed multi-layer Text Transformer to decode the hidden states and finally generate the

corresponding caption text. With the increasing de-131 mand for automated report generation, AI scholars 132 have attempted to apply image captioning methods 133 to the task of radiology report generation. However, 134 due to the fundamental differences between natural images and radiological images, general caption 136 models are difficult to be directly transferred to 137 radiology report generation. In this paper, we pro-138 pose VSE and DID to address the impact of the gap between natural images and radiological images on 140 the performance of the generation model. 141

2.2 Radiology Report Generation

142

143

144 145

146

147

148

149

150

151

152

153

155

156

157

158

159

160

161

162

163

165

168

169

170

171

172

173

174

Radiology report generation is dedicated to helping εstable stable sta AI models are used to automatically analyze X-ray images and generate corresponding reports. With the rapid development of the AI field, the accuracy of radiology report generation has been continuously improved. Initially, the R2Gen model pro-based Transformers for report analysis. This technology aimed to solve the consistency problem in report generation through the memory module and achieved good results, laying the foundation for R2Gen to be a baseline in the RRG field. Later, AI researchers discovered the difference in pixel distribution between X-ray image datasets and natural image datasets, and people began to focus on solving this problem. Some works using attention mechanisms to address this difference were proposed one after another, such as Contrastive (Liu et al., 2021b), COMG (Gu et al., 2024), etc. Subsequently, people tried to introduce some prior medical knowledge to improve the interpretability and reliability of model generation, such as PPKED, KiUT, etc. In this paper, we innovatively propose VSE and DID to bridge the gap between radiological images and natural images, thereby enhancing the interpretability and reliability of the model.

Methodology 3

In this section, we introduce the proposed VDGen ciple of VDGen from two aspects: VSE and DID.

Overview of Radiology Report Generation 3.1

175 The radiology report generation task can be formalized as follows: Given a radiological image 176 $I \in \mathbb{R}^{C \times H \times W}$, where C denotes the number of 177 channels, and H, W represent the height and width 178 of the X-ray image, respectively, the image I is first 179

encoded by the visual encoder M_v of the radiology report generation model M to obtain its feature representation F_e . Typically, M follows an encoderdecoder architecture, where M_v is implemented using ResNet or Swin Transformer. Subsequently, the textual decoder M_t decodes F_e to generate the diagnostic report $\mathcal{R} \in \mathbb{R}^{L \times D_t} = \{T_1, T_2, T_3, ..., T_L\}$ where T_i denotes a token, L is the report length (number of tokens), and D_t represents the embedding dimension of each token. M_t is commonly a large-scale Transformer model. The full report generation process can be recursively formulated as:

180

181

182

183

184

185

186

187

188

189

190

191

192

195

196

197

198

199

200

201

202

204

205

207

209

221

$$p(\mathcal{R} \mid I) = \prod_{i=1}^{L} p(T_{i+1} \mid T_1, \dots, T_i, I) \quad (1)$$

For the training of model M, given a batch of images $X \in \mathbb{R}^{B \times C \times H \times W} = \{I_1, I_2, \dots, I_B\}$ and corresponding reports $Y \in \mathbb{R}^{\tilde{B} \times L \times \tilde{D}_t} =$ $\{Y_1, Y_2, \ldots, Y_B\},$, where b = B denotes the batch size, the predominant training objective in medical report generation is to minimize the cross-entropy loss of model M. Thus, the optimization goal for M can be expressed as:

$$\mathcal{L}_{\text{report}} = \mathcal{L}_{\text{CE}}(\theta) = -\sum_{i=1}^{\mathcal{N}} \log(p_{\theta}(\mathcal{T}_{n}^{*} \mid \mathcal{T}_{1:n-1}^{*})) \quad (2)$$

where $i \in (1, \mathcal{N})$, \mathcal{N} is the vocabulary size. When i = 1, n - 1 = 0, and the input is $\mathcal{T}_{1:n-1}^{\star} = X$; otherwise, the input is $\mathcal{T}_{1:n-1}^{\star} =$ $\{X, T_1, T_2, \ldots, T_{n-1}\}.$

3.2 Vision Self-Equilibration (VSE)

To address the representation challenges caused 208 by imbalanced pixel distributions in radiology report generation, we propose the Vision Self-210 Equilibration (VSE) module. This module auto-211 matically calibrates pixel-level representations of 212 the visual encoder while mitigating distribution 213 bias. During training, given a batch of X-ray im-214 ages X, VDGen first splits X into patches, ap-215 pends a [CLS] token, and feeds them into a Swin 216 Transformer (Liu et al., 2021c) for visual encod-217 ing. The encoded semantic features are defined 218 as $X^e \in \mathbb{R}^{B \times S \times D} = \{I_1^e, I_2^e, I_3^e, \dots, I_i^e\}$, where 219 S denotes the sequence length, D_v is the embed-220 ding dimension, and B represents the batch size. To ensure alignment in a unified embedding space, 222 we apply Layer Normalization to X^e before VSE 223 processing. For each X-ray image, we compute the 224



mean embedding along the sequence dimension S, yielding $h_i^e = \frac{1}{S} \sum_{s=1}^{S} I_i^{e(s)}$, where $h_i^e \in \mathbb{R}^D$ represents the global visual semantics of the *i*th image. To enhance discriminative power, h_i^e is mapped to a higher-dimensional space. The mapped image representation can be expressed as $z_i^v = \mathcal{W}_1^v \cdot h_i^e + \mathcal{B}$, where \mathcal{W}_1^v is learnable matrix, \mathcal{B} is the bias. We expect that X-ray images of similar diseases have similar visual semantic representations, while X-ray images of different diseases have significant differences in semantic representations. In this way, we can alleviate the problem of uneven pixel distribution and achieve a better soft-alignment of visual semantic representations, ultimately achieving a self-equilibrating effect. Our specific approach is as follows: First, we calculate the similarity between different samples within a batch, $S_{ij} = \frac{z_i^{vT} z_j^v}{\tau}$, where τ is the temperature hyperparameter. We set the optimization reward target as itself z_i^v , and other samples z_i^v as penalty terms. Thus, $S_{ii} = -\infty, \forall i \in \{1, 2, 3, \dots, B\}.$ Finally, the objective function that VSE needs to optimize can be expressed as:

227

231

241

242

243

244

247

248

249

251

252

$$\mathcal{L}_{vse} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp(S_{ii})}{\sum_{j=1}^{B} \exp(S_{ij})} \quad (3)$$

 253

254

256

257

258

259

260

261

262

263

264

265

267

268

269

270

271

272

273

274

275

276

277

278

3.3 Disease Information Distillation (DID)

While VSE enables discriminative representations of disease symptoms by addressing pixel distribution imbalances in X-ray images, it remains insufficient for characterizing disease presence and extent (e.g., lesion size or scope). Precise identification of abnormalities, such as determining pathological regions in cardiac imaging, critically impacts model performance. To enhance this capability, we propose Disease Information Distillation (DID), which leverages textual embeddings from diagnostic reports (e.g., "minor pleural effusion") as teacher signals to guide VDGen in extracting diseasespecific features from visual semantic representations, thereby refining the model's ability to localize and quantify pathologies. For a given target report, after being encoded by a language model, we can obtain $Y = \{Y_1, Y_2, Y_3, \cdots, Y_i\} \in \mathbb{R}^{B \times L \times D_t}$. Before performing the DID operation, we first perform some simple normalization and alignment operations on the visual semantic representations. The operation of visual semantic representations can be expressed as:

$$z_{i}^{img} = \frac{\frac{1}{S} \sum_{s=1}^{S} I_{i}^{e(s)}}{\left\| \frac{1}{S} \sum_{s=1}^{S} I_{i}^{e(s)} \right\|_{2}} \in \mathbb{R}^{D}$$
(4)

$$z_{i}^{text} = \frac{W_{2}\hat{z}_{i}^{text} + B_{2}}{\|W_{2}\hat{z}_{i}^{text} + B_{2}\|_{2}} \in \mathbb{R}^{D}$$
(5)

After obtaining the aligned visual semantics \hat{z}_i^{img} and the text representation \hat{z}_i^{text} of the target report, 295 through the knowledge distillation technique, we 296 can use the text representation \hat{z}_i^{text} of the target 297 report as the output of the teacher model to guide the student model VDGen in the disease information distillation of visual semantics \hat{z}_i^{img} . First, we calculate the similarity between the visual representation and the target report representation as probabilities for each X-ray image and the target report features respectively, with the goal of optimizing and aligning the diagonal elements of the similarity matrix. Thus, we can obtain:

$$\mathcal{L}_{v2t} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp(S_{ii}/\tau)}{\sum_{j=1}^{B} \exp(S_{ij}/\tau)}$$

$$\mathcal{L}_{t2v} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp(S_{ii}/\tau)}{\sum_{j=1}^{B} \exp(S_{ji}/\tau)}$$
(6)

Finally, the objective function that DID needs to
optimize can be expressed as:
$$\mathcal{L}_{DID} = \frac{1}{2}(\mathcal{L}_{v2t} + \mathcal{L}_{t2v})$$
. Under the constraint of \mathcal{L}_{DID} , the represen-
tation I_i^e of the X-ray image, guided by the target
report representation Y_i , is further enhanced in rep-
resenting various disease symptoms, thus achieving
the extraction of disease information. Specifically
it further enhances the representation of whether
a disease occurs and the size of the disease scope
This strengthens the disease perception ability of
the downstream LLM during decoding, thereby
improving the accuracy of report generation.

3.4 Large Language Model Decoder

After going through VSE and DID, the visual semantic representations are greatly enhanced. We only need to decode them through a Large Language Model (LLM) to generate the final report. Before decoding, we need to construct a prompt based on the visual semantic representations to enhance the generation ability of the large-scale model (Jia et al., 2022). This kind of prompt is more conducive to the reading and understanding of the large-scale model. The prompt constructed in this paper is "Human: $\langle img \rangle I_i^e \langle /img \rangle$ Generate a comprehensive and detailed diagnosis report for this chest X-ray image.\nAssistant:". Here, I_i^e is the visual representation corrected by VSE and DID. Then this prompt is input into Llama for decoding. The entire decoding process can be represented recursively as:

322

323

324

325

326

327

329

330

331

332

334

335

336

338

340

341

342

343

344

345

346

347

349

350

$$p(\mathcal{R} \mid I_i^e) = \prod_{i=1}^{L} p(T_{i+1} \mid T_1, \dots, T_i, I_i^e) \quad (7)$$

$$\mathcal{L} = \mathcal{L}_{report} + \mathcal{L}_{VSE} + \mathcal{L}_{DID} \tag{8}$$

4 Experiment

4.1 Experiment Setting

4.1.1 Dataset

The MIMIC-CXR dataset (Johnson et al., 2019), 351 jointly released by the Massachusetts Institute of 352 Technology and Beth Israel Deaconess Medical 353 Center (BIDMC), is a large-scale resource con-354 taining approximately 370,000 chest X-ray images 355 from 227,000 patient samples, covering 14 com-356 mon thoracic diseases with over 200,000 associ-357 ated reports. Notably, 23% of the samples exhibit 358 abnormalities, reflecting a pronounced class imbal-359 ance. In contrast, the IU-XRay dataset (Demner-360 Fushman et al., 2016), curated by Indiana Uni-361 versity Hospital, comprises 7,470 chest X-ray im-362 ages from 3,851 patients paired with 3,955 re-363 ports, of which 22.2% (877 samples) are anno-364 tated as abnormal. While MIMIC-CXR typically 365 uses single-image-to-single-report pairs for radiology report generation (RRG) training, IU-XRay 367

Dataset		Year	BLEU1	BLEU2	BLEU3	BLEU4	Meteor	Rouge-L	CIDEr
		2020	0.441	0.291	0.203	0.147	_	0.367	0.304
IU-Xray		2020	0.470	0.304	0.219	0.165	0.187	0.371	-
		2021	0.402	0.284	0.168	0.143	0.170	0.328	-
	Contrastive (Liu et al., 2021b)	2021	0.492	0.314	0.222	0.169	0.193	0.381	-
		2022	0.473	0.305	0.217	0.162	0.186	0.378	-
		2022	0.475	0.305	0.221	0.171	0.188	0.375	-
	<pre>DCL (Li et al., 2023)</pre>	2023	-	-	-	0.163	0.193	0.383	<u>0.586</u>
		2023	0.483	0.322	0.228	0.172	0.192	0.380	0.435
		2023	<u>0.515</u>	0.347	0.241	0.178	0.237	0.406	0.592
		2024	0.487	0.325	0.234	0.178	0.204	0.401	0.573
		2024	0.482	0.316	0.233	0.184	0.198	0.382	0.529
		2024	0.488	0.316	0.228	0.173	0.211	0.377	0.438
		2024	0.517	0.351	0.258	0.191	0.211	0.409	-
		2025	0.491	0.325	0.230	0.169	0.195	0.362	0.386
		2020	0.353	0.218	0.145	0.103	0.142	0.277	_
		2021	0.332	0.210	0.142	0.101	0.134	0.264	-
MIMIC-CXR		2021	0.350	0.219	0.152	0.109	0.151	0.283	-
		2022	0.334	0.217	0.140	0.097	0.133	0.281	-
		2022	0.348	0.206	0.135	0.094	0.136	0.266	-
		2022	0.344	0.215	0.146	0.105	0.138	0.279	-
		2023	-	-	-	0.109	0.150	0.284	0.281
		2023	0.371	0.233	0.152	0.107	0.146	0.286	0.368
		2024	0.382	0.227	0.148	0.105	0.157	0.284	0.375
		2024	0.346	0.216	0.145	0.104	0.137	0.279	0.352
		2024	0.411	0.267	0.186	0.134	0.160	0.297	0.269
	<mbddddddddddddddddddddddddddddddddddddd< td=""><td>2024</td><td>0.415</td><td>0.254</td><td>0.166</td><td>0.117</td><td>0.154</td><td>0.285</td><td>-</td></mbddddddddddddddddddddddddddddddddddddd<>	2024	0.415	0.254	0.166	0.117	0.154	0.285	-
		2025	0.418	0.274	0.190	0.136	0.164	0.301	0.162

4.1.2 Implementation Details

368

371

374

We trained VDGen on a single Nvidia H100 with 375 80GB of memory. The Swin Transformer used is the base version proposed by Microsoft Corpora-377 tion. The input image size is set to 224x224. The 378 LLM decoder is based on Llama2-7B (Touvron 379 et al., 2023) that has been aligned through RLHF (Reinforcement Learning from Human Feedback). During the training phase, we used LoRA (Low-382 Rank Adaptation) (Hu et al., 2022) to fine-tune Llama. We configured the LoRA attention dimension to 16. The alpha hyperparameter for LoRA scaling was also set to 16. Both the training and validation batch sizes were set to 16. The learning rate was set to 1e-4. When VDGen generates reports, the maximum length of the text is controlled to be 100 tokens, and the number of newly added 390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

4.1.3 Evaluation Metrics

In order to compare the performance of VDGen fairly with other methods in the same field, we adopted the mainstream metrics in the RRG field to evaluate the performance of the model. These include commonly used metrics such as BLEU n (Papineni et al., 2002), Meteor (Banerjee and Lavie, 2005), Rouge-L (Chin-Yew, 2004), and CIDEr. BLEU is usually used to measure the similarity of n-grams between the predicted report and the ground truth report. The higher the BLEU score, the closer the generated report is to the ground truth report. Other metrics like Rouge-L are usually used in caption translation tasks. Here, they are used to measure the overall similarity between the report generated by the model and the ground truth report, without calculating from a fine-grained perspective.

416

417

444

445

446

447

448

449

450

451

452

453

454

455

456

457 458

459

460

461

462

By conducting a fair comparison of RRG through these popular metrics, we can observe the performance of VDGen.

4.2 Main Result

To demonstrate the effectiveness of VDGen, we 418 conducted a large number of experiments on the IU-419 Xray and MIMIC-CXR datasets. The experimental 420 results are shown in Table 1. The best results are 421 presented in bold font, and the second-best results 422 are underlined. The VDGen we proposed did not 423 perform outstandingly on the IU-XRay dataset, dif-424 fering from the best result by nearly 1.5%. A pos-425 sible reason is that the IU-Xray dataset is a small-426 sample dataset with only about 7,000 data points, 427 which is difficult to support the amount of data re-428 quired for the training of the LLM decoder, result-429 ing in a slight decrease in performance. However, 430 compared with other results on the IU-Xray dataset, 431 VDGen is still comparable. VDGen performed re-432 markably on the MIMIC-CXR dataset and achieved 433 a new state-of-the-art (SOTA). Specifically, VD-434 Gen outperformed the second-best result by 1% 435 on the MIMIC-CXR dataset. This fully demon-436 strates that VDGen has good performance when 437 there is an adequate amount of training data. The 438 439 SOTA result on the MIMIC-CXR dataset also establishes the novelty and feasibility of our method. 440 In general, we achieved comparable performance 441 on the IU-XRay dataset and a new SOTA on the 442 MIMIC-CXR dataset. 443

4.3 Ablation Study

rriser varier vari Gen, we designed ablation experiments for the VSE, DID, and LoRA modules of VDGen. Table 2 shows that removing VSE led to a 1.26% average performance drop, proving it can address radiological image pixel distribution issues and boost model performance. Removing DID caused a 1.34% average decline, indicating using the ground truth report as a teacher signal for VDGen to extract disease info is effective. VSE is better at capturing fine-grained info, while DID handles both fine and coarse-grained data. The LoRA ablation experiment showed its fine-tuning benefits the downstream LLM decoder. Overall, VSE and DID are essential for VDGen, and their absence will reduce performance. At the same time, using the LoRA method for training can bring further performance improvements.

Clinical Efficacy Analysis 4.4

To demonstrate the potential of VDGen in clinical applications, we evaluated its performance on the MIMIC-CXR dataset using clinical efficacy 466 metrics (Precision, Recall, and F1-score). The evaluation results are shown in Table 1. Experiments show that VDGen achieved the best Recall and F1score in the MIMIC-CXR benchmark test, and ob-470 tained comparable Precision. Although VDGen's 471 Precision is slightly lower than that of DCL, its 472 Recall and F1-score are 4.5% and 2.9% higher than 473 those of DCL respectively. Therefore, consider-474 ing comprehensively, VDGen outperforms DCL in clinical applications. Meanwhile, VDGen leads the 476 second-ranked model by 1% in Recall and 1.3% 477 in F1-score, which also highlights the excellent 478 clinical potential of VDGen. Overall, VDGen has shown significant improvements in clinical efficacy metrics, demonstrating good clinical effects and highlighting its potential in clinical applications.

463

464

465

467

468

469

475

479

480

481

482

483

4.5 **Case Studies**

To further demonstrate the effectiveness and clin-484 ical significance of VDGen, we conducted quali-485 tative experiments on the MIMIC - CXR dataset. 486 We compared VDGen with the most popular LLM -487 based methods on the RRG dataset. The experimen-488 tal results are shown in Figure 3. We selected two 489 groups of classic cases to analyze the differences 490 between the two methods. The first group mainly 491 consists of normal samples, that is, the images do 492 not contain any diseases. The second group is ab-493 normal samples, that is, the images contain one 494 or more diseases that are difficult to observe with 495 the naked eye. For the first-group samples, the 496 reports generated by VDGen are basically indis-497 tinguishable from those described by professional 498 physicians. However, the current mainstream meth-499 ods fail to diagnose the previous labels, that is, they 500 do not recognize that the patient has had a previ-501 ous surgery. For the second - group samples, both 502 VDGen and R2GenGPT identified the degenera-503 tive changes in the thoracic vertebrae. However, 504 VDGen additionally reported that the patient had a 505 previous surgery, which is crucial information for 506 formulating a treatment plan. R2GenGPT missed 507 this key information. Overall, the reports generated 508 by VDGen are more comprehensive and of greater 509 reference value. This also indirectly reflects the 510 clinical significance of VDGen. 511

Model	BLEU 1 \downarrow	BLEU2↓	BLEU3↓	BLEU4↓	Meteor↓	CIDEr↓
Ours	0.418	0.274	0.190	0.136	0.164	0.162
w/o VSE w/o DID	0.395	0.257	0.176	0.128	0.158	0.158 0.154
w/o LoRA	0.406	0.263	0.182	0.132	0.161	0.158

	 Target: The lungs are clear without focal consolidation, effusion, or edema. The cardiomediastinal consolidation, effusion, or edema. The clear the cardiomediastinal consolidation, effusion, or edema. The cardiomediastinal consolidation, effusion, or edema and consolidation, or edema and consolidat	> R2GenGPT : heart size is normal. R2GenGPT : heart size is normal. R2GenGPT : heart size is normal. The mediatismal and hidar and heart size is normal. The mediatismal and heart size is normal. The mediate size is normalities.	> WDGen :pa and lateral views of the chest provided. there is no focal consolided in the chest provided the the chest provided. The chest provided the the chest provided the chest provided the chest provided the the chest provided the chest provided the chest provided the the chest provided the chest provided the chest provided the the chest provided the ches
ධ		> R2GenGPT :heart size is normal. the mediastinal and man and mediatinal and normal. the mediastinal and normal. the mediastinal and normal the mediastical mediast	 VDGen:pa and lateral views of the chest provided. there is no focal consolidation effuers of the chest provided. there is no focal consolidation effuers of the chest provided. There is no focal consolidation effects on the chest provided is no for distribution of the chest provided is not provided. VDGUE views and vi

	Precision↑	Recall↑	F1-score↑	
	0.334	0.275	0.278	
́SEBTSAT+KG	0.356	0.297	0.304	
METransformer	0.364	0.309	0.311	
	0.471	0.352	0.373	
	0.371	0.318	0.321	
	0.392	0.387	0.389	
	0.424	0.326	0.345	
VDGen	0.426	0.397	0.402	

ፗដ

5 Conclusion

512

513

514

515

 Gen incorporates two complementary modules. Vision Self-Equilibration(VSE) mitigates visual feature degradation through self-supervised contrastive learning, enhancing intra-modal discriminability. Disease Information Distillation(DID) leverages diagnostic reports as teacher signals to guide cross-modal representation learning, improving the model's ability to capture disease-specific semantics. These modules are integrated into an end-to-end generation pipeline with a LoRAadapted large language model decoder. Experiments on MIMIC-CXR and IU-Xray show that VDGen improves the state-of-the-art by 1.0% on MIMIC-CXR. Ablation studies demonstrate that removing VSE and DID causes performance drops of 1.26% and 1.34%, respectively, highlighting the importance of both components. Overall, VDGen narrows the performance gap between natural and radiological domains and offers a scalable, interpretable solution for clinical report generation.

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

Limitations

Our study introduces a novel multimodal model, VDGen. This model aims to automatically adjust the model's attention regions via the Visual Self Equilibration(VSE) and extract disease informa-tion through the Disease Information Distillation (DID), thereby enhancing the performance of radi-ology report generation. However, this approach has certain limitations. Firstly, the model is tai-lored for the highly specialized task of radiology report generation, which may limit its adaptability in broader natural language generation tasks. Additionally, current evaluations are based on specific datasets and may not fully reflect the framework's applicability in diverse scenarios.

References

- ऱiin in itees in itees