# **Can Medical Vision-Language Pre-training Succeed with Purely Synthetic** Data?

**Anonymous ACL submission** 

#### Abstract

Medical Vision-Language Pre-training (Med-001 VLP) has made significant progress in enabling 002 zero-shot tasks for medical image understand-003 004 ing. However, training MedVLP models typically requires large-scale datasets with paired, 005 high-quality image-text data, which are scarce 006 in the medical domain. Recent advancements 007 in Large Language Models (LLMs) and dif-008 009 fusion models have made it possible to gener-010 ate large-scale synthetic image-text pairs. This raises the question: Can MedVLP succeed us-011 ing purely synthetic data? To address this, we 012 use off-the-shelf generative models to create 013 014 synthetic radiology reports and paired Chest X-015 ray (CXR) images, and propose an automated pipeline to build a diverse, high-quality syn-016 thetic dataset, enabling a rigorous study that 017 018 isolates model and training settings, focusing entirely from the data perspective. Our results 019 020 show that MedVLP models trained *exclusively* on synthetic data outperform those trained on 021 022 real data by 3.8% in averaged AUC on zeroshot classification. Moreover, using a com-023 024 bination of synthetic and real data leads to a 025 further improvement of 9.07%. Additionally, 026 MedVLP models trained on synthetic or mixed data consistently outperform those trained on 027 028 real data in zero-shot grounding, as well as in fine-tuned classification and segmentation 029 tasks. Our analysis suggests MedVLP trained 030 on well-designed synthetic data can outperform 031 models trained on real datasets, which may be 032 limited by low-quality samples and long-tailed 033 distributions<sup>1</sup>. 034

#### 1 Introduction 035

In medical image analysis, learning representa-036 tive features typically requires labor-intensive and 037 costly image annotations (Ronneberger et al., 2015; 038 Liu et al., 2023b). Medical Vision-Language Pre-039 training (MedVLP) addresses this challenge by 040

aligning vision and language content using paired 041 datasets of images and clinical reports, reducing the need for manual annotations (Radford et al., 043 2021; Zhang et al., 2020; Wu et al., 2023; Liu 044 et al., 2023a). However, existing MedVLP mod-045 els rely heavily on large-scale, high-quality paired 046 data (Liu et al., 2023e), which is scarce in prac-047 tice. Real-world datasets often contain noisy data, 048 such as low-quality images and unpaired image-049 text samples, degrading model performance (Xie 050 et al., 2024; Bannur et al., 2023). Recent advance-051 ments in Large Language Models (LLMs) and dif-052 fusion models enable the generation of large-scale 053 synthetic image-text datasets, offering an alterna-054 tive to traditional data collection. Although these 055 techniques have shown promise in medical tasks, 056 they are primarily used as auxiliary support for 057 real data via augmentation (Chen et al., 2024a; Yao 058 et al., 2021; Chen et al., 2022; Qin et al., 2023), 059 and are often limited to single-modality settings. To the best of our knowledge, no studies have 061 fully explored the potential of using synthetic mul-062 timodal data for MedVLP or considered training 063 exclusively on synthetic data (Liu et al., 2023e). 064

042

060

065

066

067

068

069

070

071

072

To bridge this gap and showcase synthetic data's potential for MedVLP, our contributions are:

- We propose an automated pipeline to create the SynCXR dataset, which contains 200,000 synthetic images and text generated with quality and distribution control using off-the-shelf models, without relying on real data or manual curation.
- We successfully demonstrate that MedVLP 073 models trained on our SynCXR dataset, con-074 taining only synthetic data, outperform those 075 trained on real data. Moreover, combining 076 synthetic and real data further improves per-077 formance, showcasing the effectiveness of our 078 synthetic data generation pipeline. 079

<sup>&</sup>lt;sup>1</sup>All data and code will be released upon acceptance.

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

130

131

132

• We identify several issues in the most com-080 monly used real dataset for MedVLP, MIMIC-081 CXR (Johnson et al., 2019b), that degrade 082 MedVLP performance, including low-quality 083 images and unpaired image-text samples. Fur-084 thermore, we identify the long-tailed distribu-085 tion problem in multimodal datasets, as shown 086 in Fig 1, 2. 087

We conduct an extensive analysis of the key factors contributing to MedVLP's success using purely synthetic data. Our method is evaluated on seven downstream tasks using zeroshot learning and linear probing, demonstrating that MedVLP can effectively perform with synthetic data alone.

### 2 Methods

095

096

#### 2.1 Exploring Imperfections in Real Data

For MedVLP, the most commonly used dataset is 097 MIMIC-CXR (Johnson et al., 2019a,b), a collection 098 of chest x-ray (CXR) images paired with their cor-099 responding textual reports. after following the pre-100 processing steps outlined in previous works (Zhang 101 et al., 2023; Wang et al., 2022; Huang et al., 2021), 102 this dataset provides a total of 213,384 image-text 103 pairs for pre-training. And all images must be 104 frontal views according to the preprocessing steps 105 outlined in (Huang et al., 2021). 106

Previous work on VLP with natural images (Xu 107 et al., 2023a) has shown that data quality, including 108 image fidelity and long-tailed distribution, signifi-109 cantly impacts model performance. However, the 110 quality of MedVLP datasets remains underexplored 111 due to ambiguity in defining medical image quality, 112 stemming from diverse imaging protocols. Addi-113 tionally, quantifying data distribution is complex, 114 as radiology reports often describe patterns across 115 multiple anatomical regions rather than distinct cat-116 egories. To address these challenges, we develop a 117 systematic pipeline to thoroughly analyze the data 118 issues in the MIMIC-CXR (Johnson et al., 2019b) 119 dataset, as detailed in the following sections. 120

Low-Quality and Mismatched Image-Text Pairs. 121 Our aim is to explore and identify issues related to 122 image quality in the MIMIC-CXR dataset (John-123 son et al., 2019a), rather than to completely clean 124 125 the dataset, as creating a perfect dataset and filtering out all low-quality samples is infeasible for 126 large-scale multimodal datasets (Xu et al., 2023b). 127 Inspired by (Bannur et al., 2023), which high-128 lights various issues with poor-quality images, we 129

design six queries for a Multimodal Large Language Model (MLLM), utilizing the InternVL2-26B model<sup>2</sup> (Chen et al., 2023, 2024b). Each CXR image from the MIMIC-CXR dataset is paired with these six queries, and the MLLM process each query independently. The process is depicted in Fig 2 (b). We described the queries for each function in detail in Sec. B.

After this process, we filter out the CXR images where the answers are all 'NO' across the six queries. Fig 2 (a) shows examples of images where the answer was 'NO'. We identified and removed 1,448 such images and their corresponding reports from the preprocessed MIMIC-CXR dataset, leaving us with 211,936 image-text pairs.

To further refine the dataset, we use the CXRspecific vision encoder, RAD-DINO (Pérez-García et al., 2024), to extract image features from the remaining 211,936 CXR images and from the 1,448 samples identified as bad by MLLM filtering. We then compute the similarity between each image in the cleaned dataset and each of the bad samples. Since each image comes from a different clinical case, we only compare image quality rather than the clinical content (e.g., diagnoses or abnormalities). To do this, we set a similarity threshold of 0.5 and remove all images with a similarity score greater than 0.5. This step resulted in the removal of an additional 5,512 images and their paired reports, reducing the dataset to 206,424 image-text pairs. Fig 2 (a) also shows the samples removed based on their similarity to bad images using visual features from RAD-DINO<sup>3</sup> (Pérez-García et al., 2024).

In our exploration of the MIMIC-CXR dataset, we utilized a rough approach to identify problematic images, such as non-chest images, wrong views, overprocessing, and low-fidelity scans. Our results confirm that many images in the dataset exhibit these issues. While our approach identifies numerous problematic images, fully curating and removing all low-quality cases is unfeasible due to the substantial human effort required and the absence of well-defined criteria for an automated cleaning pipeline. Furthermore, addressing all instances of low-quality images remains highly challenging through automated processes alone. **Uncovering Long-tailed Distribution in MIMIC-CXR.** As demonstrated in previous work on natural

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/OpenGVLab/InternVL2-26B

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/microsoft/rad-dino



Figure 1: Comparison of real image-text datasets and synthetic datasets. (a): The real image-text dataset, MIMIC-CXR (Johnson et al., 2019b), while authentic, often contains imperfections such as long-tailed data distribution, unpaired images and text, and low-quality CXR images, which limit the performance of MedVLP models pretrained on this dataset. (b): The synthetic dataset generation process uses clinical entities as prompts to an LLM (e.g., Llama3.1 (AI@Meta, 2024)) to generate synthetic reports. These reports are then used to create synthetic images through RoentGen (Bluethgen et al., 2024). We propose an automated pipeline to control the dataset distribution, ensuring it is balanced and includes paired image-text samples.



Figure 2: (a): Examples of invalid or low-quality images filtered out by the proposed image curation method described in Sec 2.1. (b): The image curation pipeline uses InternVL2 (Chen et al., 2023), a Multimodal Large Language Model (MLLM), to assess CXR image quality. Images that meet the criteria are retained; others are discarded. (c): Entity frequency distribution in the MIMIC-CXR dataset. Due to space constraints, only the top 50 frequent entities for four categories (Abnormality, Non-Abnormality, Disease, Non-Disease) are shown. A more detailed distribution is presented in Fig 6,7,10,8,9.

179 image-text data (Xu et al., 2023b; Hammoud et al., 2024), a long-tailed distribution in VLP datasets 180 negatively impacts model performance. There-181 fore, we aim to explore the data distribution of the 182 MIMIC-CXR dataset. However, directly evaluat-183 ing the text distribution at the sample level, as done 184 in (Xu et al., 2023b), is challenging because each 185 radiology report often describes multiple patterns 186 or anatomical regions, unlike natural image cap-187 tions that typically focus on a single object (Zhang 188 et al., 2024). 189

using an off-the-shelf Named Entity Recognition (NER) tool to extract all medical entities, treating them as representatives of the report's concepts and exploring the dataset distribution at the entity level. For this, we use RaTE<sup>4</sup> (Zhao et al., 2024), a model specifically designed for NER tasks on radiology reports. RaTE automatically classifies the extracted entities into five categories: [ABNORMALITY, NON-ABNORMALITY, DISEASE, NON-DISEASE, ANATOMY]. We dis191

192

<sup>190</sup> Instead, we adopt an alternative approach by

<sup>&</sup>lt;sup>4</sup>https://huggingface.co/Angelakeke/RaTE-NER-Deberta

play the top 50 frequent entiites distribution of each 201 entity type in Fig 2 (c). We display the top 50 fre-202 quent entiites distribution of each entity type in 203 Fig 6,7,10,8,9. As shown, all entity types exhibit a 204 severe long-tailed distribution. As shown, all entity 205 types exhibit a severe long-tailed distribution in the 206 MIMIC-CXR (Johnson et al., 2019b), which con-207 tains 154,049 unique entities, with 55,047 Abnor-208 mality, 36,365 Non-Abnormality, 23,017 Disease, 209 22,103 Non-Disease, and 40,517 Anatomy entities. 210

# 211 2.2 Generating Synthetic CXR reports and 212 Paired Images.

213 Since the MIMIC-CXR dataset (Johnson et al., 2019a) contains various data issues, we generate 214 synthetic radiology reports and CXR images, con-215 trolling data quality and distribution during gen-216 eration to alleviate these problems. In this work, 217 we aim to explore the effectiveness of pretraining 218 MedVLP on a purely synthetic dataset, rather than 219 attempting to create a perfect dataset, as noisy data 220 221 is unavoidable in real-world scenarios and an ideal dataset is unrealistic. 222

CXR Report Generation. To generate the synthetic reports, the pipeline is depicted in Fig 5.
We select a general LLM, Llama3.1-70B-Instruct as the report generator, and we extensively ablate the performance of the report generator with other LLMs in Fig 3. We query the LLM using prompts that include the entity list, as shown in Fig 5.

Since we aim to build a synthetic dataset without a long-tailed distribution, we design a balanced sampling strategy to ensure that the appearance frequency of each entity type is approximately equal across the synthetic dataset. Let  $\mathcal{E}$  be the set of entities, categorized into five types: ABNORMALITY, NON-ABNORMALITY, DISEASE, NON-DISEASE, and ANATOMY.

For each generation, we sample:

239 
$$\mathcal{S}_1 = \{e_1^{(i)}, e_2^{(i)}, \dots, e_k^{(i)}\},\$$

230

231

232

233

234

235

236

237

238

240

241 242

243

244 245  $\forall e_j^{(i)} \in \{ \texttt{ABNORMALITY}, \texttt{NON-ABNORMALITY}, \\ \texttt{DISEASE}, \texttt{NON-DISEASE} \}.$ 

where k is the number of entities sampled from the first four categories. Additionally, we sample:

246 
$$\mathcal{S}_2 = \{a_1^{(i)}, a_2^{(i)}, \dots, a_m^{(i)}\}, \quad \forall a_j^{(i)} \in \texttt{ANATOMY}$$

247 where m is the number of entities sampled from 248 the ANATOMY category. Thus, the total sampled entity set for each generation is:

S

$$=\mathcal{S}_1\cup\mathcal{S}_2$$
 250

249

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

We impose a maximum frequency threshold,251 $\tau_{\max}$ , for each entity  $e \in \mathcal{E}$ . If an entity  $e_j^{(i)}$  in252 $\mathcal{S}$  reaches this threshold, we resample  $e_j^{(i)}$  while253keeping the remaining entities in  $\mathcal{S}$  unchanged:254

if 
$$f(e_j^{(i)}) \ge \tau_{\max}$$
, then resample  $e_j^{(i)}$ . 255

Here, f(e) denotes the current frequency of entity ein the dataset. This ensures a balanced distribution of entities across the synthetic dataset. We set k =9, m = 3, and  $\tau_{max} = 15$  in our work during the generation stage.

After sampling, we input the selected entities  $S = S_1 \cup S_2$  into the LLM and indicate their type. Let the output of the LLM be denoted as  $R_{gen}$ , which represents the synthetic report generated by the model based on the sampled entities. To ensure that the LLM-generated report  $R_{gen}$  covers and only includes the entities in S (since the inclusion of non-specified entities would disrupt the frequency balance), we use the RaTE model (Zhao et al., 2024) to extract entities from  $R_{gen}$ , denoted as  $\mathcal{E}_{gen}$ .

We then verify the entity set  $\mathcal{E}_{gen}$  by comparing it with the originally sampled set  $\mathcal{S}$ . If  $\mathcal{E}_{gen} \neq \mathcal{S}$ , we regenerate the report  $R_{gen}$  by repeating the generation process until  $\mathcal{E}_{gen} = \mathcal{S}$ :

if 
$$\mathcal{E}_{gen} \neq \mathcal{S}$$
, regenerate  $R_{gen}$  until  $\mathcal{E}_{gen} = \mathcal{S}$ .

Once the synthetic report is successfully generated, it is used as the 'FINDINGS' section of the CXR report. We then query the LLM to summarize  $R_{gen}$  into the 'IMPRESSION' section, denoted as  $R_{imp}$ . To ensure consistency between the entities in the 'FINDINGS' and 'IMPRESSION' sections, we extract entities from the summary  $R_{imp}$  using RaTE, denoted as  $\mathcal{E}_{imp}$ . We verify that:

$$\mathcal{E}_{imp} = \mathcal{S}.$$
 285

If the entities in  $R_{imp}$  do not match S, we regenerate the "IMPRESSION" section until  $\mathcal{E}_{imp} = S$ : 287

if  $\mathcal{E}_{imp} \neq \mathcal{S}$ , regenerate  $R_{imp}$  until  $\mathcal{E}_{imp} = \mathcal{S}$ . 288

Given that the number of samples in the original289MIMIC-CXR dataset cannot be perfectly divided290by k and m, we generate a total of 200,000 syn-291thetic samples to ensure a balanced distribution292

346 347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

using only off-the-shelf tools, without any specificdesign for CXR data.

While RadGraph (Delbrouck et al., 2024) could 295 be used for entity extraction, it relies on human-296 annotated data from MIMIC-CXR and is limited 297 to 16,117 entities. In contrast, RaTE (Zhao et al., 298 2024) extracts 154,049 entities, making it more 299 suitable for our goal of creating a general and easily 300 transferable pipeline for synthetic data generation. 301 Thus, we chose RaTE for its broader applicability 302 to various radiology reports. 303

CXR Image Generation. After generating the syn-304 thetic radiology reports, we aim to generate paired 305 CXR images conditioned on the synthetic reports. 306 Since general text-to-image (T2I) models (e.g., Sta-307 ble Diffusion) are not designed for CXR image 308 generation and demonstrate poor performance, as 309 shown in (Liu et al., 2023e; Bluethgen et al., 2024), 310 we select RoentGen<sup>5</sup> (Bluethgen et al., 2024), the 311 most recent and validated CXR-specific T2I model, 312 verified by clinicians, as our image generator. We 313 use RoentGen's (Bluethgen et al., 2024) official 314 pretrained weights to generate images. Following 315 their implementation, we use only the 'IMPRES-316 SION' section from the synthetic reports as the text 317 prompt for the T2I model. The generation process 318 is controlled using the official hyperparameters pro-319 vided by RoentGen, where the classifier-free guid-320 321 ance (CFG) is set to 4 and the number of denoising steps is set to 50. 322

To prevent the synthetic images from exhibiting 323 the same issues found in the real dataset (as dis-324 cussed in Sec. 2.1), we apply a similar curation 325 procedure. First, we use the MLLM to filter syn-326 thetic images, and then we compute the similarity 327 of visual features between synthetic images and 328 the problematic samples identified from the real 329 dataset. If the visual similarity exceeds a threshold 330  $\delta = 0.5$ , we regenerate the images by re-querying 331 the T2I model with the same text prompt until they 332 pass the curation procedure. 333

We generate 200,000 synthetic CXR images, each paired with a corresponding synthetic report, using only general-purpose, open-source models (e.g., Llama3.1 (AI@Meta, 2024), InternVL2 (Chen et al., 2023)) and vision models pre-trained with self-supervised learning (e.g., RAD-DINO (Pérez-García et al., 2024)). No annotated CXR images or MedVLP models pre-trained on specific CXR image-text datasets are used in this

334

335

336

337

338

339

340

341

342

process. This ensures our approach is adaptable and can easily incorporate future advancements in general-purpose models. We refer to this dataset as **SynCXR**.

### 2.3 Synthetic Data Training for MedVLP

In this work, we use the synthetic dataset, SynCXR, to train a MedVLP model and investigate how effectively a model can learn from purely synthetic data. Given the abundance of existing MedVLP methods, we focus on simple baseline models for the following reasons:

**ConVIRT** (Zhang et al., 2020) jointly trains vision and text encoders on paired medical images and reports using global contrastive learning.

**GLoRIA** (Huang et al., 2021) extends ConVIRT by adding both global and regional contrastive learning, enabling more effective training of encoders on paired medical images and reports.

These models are open-source, straightforward, and minimize the influence of external factors, which is crucial for evaluating synthetic data in the context of MedVLP. For retraining these models on our synthetic dataset, SynCXR, we adhere strictly to their official codebases<sup>67</sup>. More complex models may introduce unnecessary complications. We are aware that recent methods (Boecking et al., 2022; Bannur et al., 2023; Wu et al., 2023; Zhang et al., 2023; Phan et al., 2024b) either lack publicly available code or rely on additional human annotations, which make direct implementation with synthetic data challenging and introduce unwanted variables.

## 3 Experiments Configurations

For pre-training, we apply the official configurations provided by ConVIRT (Zhang et al., 2020) and GLoRIA (Huang et al., 2021) on the MIMIC-CXR dataset to our synthetic CXR image-text dataset, SynCXR.

# 3.1 Downstream Task Datasets and Configurations

For downstream tasks, we evaluate the effectiveness of synthetic data for MedVLP across four tasks. We strictly follow the downstream setting described in (Phan et al., 2024b) to evaluate our method. Due to the page limit, detailed information on the datasets and implementation is provided in Appendix, Sec. C.

<sup>&</sup>lt;sup>5</sup>https://stanfordmimi.github.io/RoentGen/

<sup>&</sup>lt;sup>6</sup>https://github.com/marshuang80/gloria

<sup>&</sup>lt;sup>7</sup>https://github.com/edreisMD/ConVIRT-pytorch

#### **390 3.2** Experimental Results

Since the MIMIC-CXR dataset already includes 391 several diseases present in downstream tasks, as 392 393 mentioned in (Phan et al., 2024b; Zhang et al., 2023), we split the zero-shot classification task 394 into seen and unseen categories, strictly follow-395 396 ing (Phan et al., 2024b). Note that all experimental results for ConVIRT and GLoRIA pre-trained with 397 real data (MIMIC-CXR) are directly referenced 398 from (Phan et al., 2024b) to ensure a fair compari-399 400 son.

Zero-shot Classification on Seen Diseases. Tab 1 401 shows the zero-shot classification performance on 402 seen diseases. Across all datasets, both MedVLP 403 methods pretrained on SynCXR (our purely syn-404 thetic dataset) consistently outperform or achieve 405 comparable performance to their counterparts pre-406 trained on real datasets, with an average improve-407 ment of 4.7% in AUC and 4.53% in F1 scores. 408 Furthermore, the methods pretrained on the mixed 409 410 dataset, which directly combines real and synthetic data, achieve even greater improvements, 411 with 10.08% AUC and 7.62% F1 scores on average 412 across all datasets and methods. This demonstrates 413 that the SynCXR dataset effectively enables Med-414 VLP models to learn representative cross-modal 415 features, enhancing their zero-shot classification 416 capability. 417

Zero-shot Classification on Unseen Diseases. Tab 418 2a reports the zero-shot classification performance 419 on unseen diseases. Similar to the results for seen 420 diseases, MedVLP models pretrained on the syn-421 thetic dataset consistently outperform those pre-422 trained on real data, with an average improvement 423 424 of 2.96% AUC and 0.51% F1 scores. Additionally, models pretrained on the mixed dataset show sub-425 stantial gains over those trained on real data, with 426 7.39% AUC and 1.52% F1 scores on average. This 427 428 indicates that the SynCXR dataset, generated with meticulous quality control and balanced distribu-429 tion, can increase the generalizability of MedVLP 430 models for unseen diseases prediction. 431

Zero-shot Visual Grounding. We further evalu-432 ate the effectiveness of synthetic data in improving 433 MedVLP models' local visual understanding capa-434 bilities through zero-shot grounding tasks. Tab 2b 435 436 presents the performance of zero-shot grounding on RSNA (Shih et al., 2019), Covid-19 Rural (De-437 sai et al., 2020), and SIIM (Steven G. Langer and 438 George Shih, 2019). Across all datasets, MedVLP 439 models pretrained on the SynCXR dataset achieve 440

superior performance compared to those trained on 441 the real dataset, with an average increase of 1.42% 442 IoU and 0.97% Dice scores. The mixed dataset 443 further enhances performance, with 4.06% IoU and 444 2.92% Dice scores on average. This demonstrates 445 that the SynCXR dataset not only benefits global 446 cross-modal feature learning but also improves lo-447 cal visual understanding for MedVLP models. 448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

**Fine-tuning Tasks.** To evaluate the representation quality learned by MedVLP, we report the finetuned classification and segmentation performance in Tab 3. Similar to the zero-shot task, MedVLP models pre-trained on SynCXR consistently outperform those trained on the real dataset across all data ratios for both classification and segmentation tasks. Furthermore, the combination of real and synthetic datasets (Mix) further boosts performance, demonstrating that SynCXR data not only enhances cross-modal representation learning but also improves performance in single-modal tasks.

#### 4 Analysis

Effect of Balanced Entity Sampling in Generating Synthetic Reports. We evaluate the impact of balanced sampling entities when generating synthetic reports using LLMs. For the synthetic dataset without balanced sampling, we adjust entity frequencies to match their distribution in MIMIC-CXR, leading to a long-tailed distribution. As shown in Tab 4a, for both MedVLP methods, the performance improves significantly when using synthetic datasets generated from balanced sampled entities. This demonstrates that balanced sampling of entities leads to a more representative dataset, benefiting MedVLP performance.

**Evaluating the Contribution of Synthetic Images and Reports.** We aim to assess the individual impact of synthetic images and synthetic reports on MedVLP performance. As shown in Tab 4b, we generate two partially synthetic datasets by replacing either the image or the text with synthetic data, while keeping the other components real, to evaluate their respective contributions.

• Real Image, Synthetic Report: In this setting, we use MedVersa<sup>8</sup> (Zhou et al., 2024), a state-of-the-art radiology report generation model, to generate synthetic reports for each real CXR image. We then train MedVLP mod-

<sup>&</sup>lt;sup>8</sup>https://huggingface.co/hyzhou/ MedVersa

Method	Pre-training	CheXpert		ChestX	ChestXray-14		st-seen	RSI	NA	SIIM		
	Data	AUC ↑	$F1\uparrow$	AUC ↑	$F1\uparrow$	AUC $\uparrow$	$F1\uparrow$	AUC ↑	$F1\uparrow$	AUC ↑	F1 $\uparrow$	
ConVIRT	MIMIC-CXR	52.10	35.61	53.15	12.38	63.72	14.56	79.21	55.67	64.25	42.87	
	SynCXR	59.49	40.51	56.07	15.43	63.43	15.10	82.08	58.38	75.55	57.43	
	Mix	71.54	47.11	61.28	18.52	68.48	16.67	83.86	61.28	78.51	59.10	
GLoRIA	MIMIC-CXR	54.84	37.86	55.92	14.20	64.09	14.83	70.37	48.19	54.71	40.39	
	SynCXR	61.38	41.05	57.47	15.60	64.26	15.02	72.34	49.50	67.32	53.86	
	Mix	72.32	48.54	61.06	17.33	68.35	17.00	74.32	51.10	73.49	56.09	

Table 1: Performance of zero-shot classification on five datasets for diseases present in the MIMIC-CXR dataset, evaluated on two MedVLP models pretrained on MIMIC-CXR (real) and SynCXR (**pure synthetic**). 'Mix' denotes the direct combination of real and synthetic data for MedVLP pretraining. Best results are highlighted in bold.

Method	Pre-training	Covid-19 CXR-2		XR-2 PadChest-unseen		PadChest-rare		-	Method	Pre-training	RS	SNA	Covid-19 Rural		SIIM	
	Data	AUC ↑	F1 ↑	AUC $\uparrow$	F1 ↑	AUC ↑	F1↑	_		Data	loU↑	Dice $\uparrow$	IoU↑	Dice $\uparrow$	IoU ↑	Dice $\uparrow$
	MIMIC-CXR	62.78	71.23	51.17	4.12	50.37	3.31	-		MIMIC-CXR	18.93	28.45	7.42	10.55	3.01	8.74
ConVIRT	SynCXR	64.41	72.03	54.47	4.51	53.70	3.69		ConVIRT	SynCXR	22.98	31.45	8.62	10.83	3.43	9.67
	Mix	69.23	72.85	58.53	5.35	57.68	4.40	_		Mix	25.97	34.25	12.78	14.12	4.58	11.43
	MIMIC-CXR	64.52	70.78	49.96	4.07	48.25	3.41	-		MIMIC-CXR	21.82	34.68	8.18	12.49	3.11	10.23
GLoRIA	SynCXR	66.70	71.90	54.24	4.10	51.26	3.75		GLoRIA	SynCXR	23.00	35.25	9.47	13.00	3.50	10.75
	Mix	68.76	73.22	58.60	5.60	58.58	4.62	_		Mix	26.34	36.52	12.67	14.63	4.51	11.73

(a) Performance of zero-shot classification on three datasets for (b) Performance of zero-shot grounding on RSNA, SIIM, and unseen diseases. Covid-19 Rural.

Table 2: Zero-shot tasks performance of MedVLP models on disease classification (a) and grounding (b) across multiple datasets, using MIMIC-CXR, SynCXR, and Mix datasets for pretraining.

Task	Classification								Segmentation												
Dataset		RSNA			SIIM		Cov	id19 CX	R-2	Ch	nestXray	-14		RSNA		Co	vid-19 R	ural		SIIM	
Data Ratio	1%	10%	100%	1%	10%	100%	1%	10%	100%	1%	10%	100%	1%	10%	100%	1%	10%	100%	1%	10%	100%
ConVIRT-Real	78.86	85.42	87.64	72.39	80.41	91.67	90.30	97.74	99.70	57.23	72.53	79.13	56.48	63.94	71.87	16.97	30.79	42.71	28.75	47.21	65.75
ConVIRT-Syn	79.01	85.58	87.90	73.51	81.10	91.84	91.50	98.80	99.73	57.45	73.60	80.20	58.00	65.10	72.90	17.10	32.00	43.90	29.90	48.50	66.81
ConVIRT-Mix	79.75	86.21	88.45	73.00	82.80	92.31	91.81	99.00	99.81	57.61	74.20	80.51	58.50	65.81	73.30	18.40	32.50	44.21	30.10	48.81	67.11
GLoRIA-Real	79.13	85.59	87.83	75.85	86.20	91.89	92.74	97.18	99.54	58.94	72.87	79.92	58.13	67.71	72.06	16.12	31.20	43.85	31.87	40.61	64.82
GLoRIA-Syn	80.30	86.75	88.00	76.01	87.40	92.11	94.01	98.41	99.75	60.11	74.01	81.11	60.41	70.01	73.51	17.31	32.51	45.01	32.91	41.91	66.01
GLoRIA-Mix	81.01	87.50	88.61	77.51	88.01	92.51	94.51	99.61	99.86	60.31	74.51	81.51	61.01	70.51	74.01	17.51	33.01	45.31	33.51	42.21	67.51

Table 3: Results from two MedVLP methods pre-trained on real, synthetic, and mixed datasets are reported for classification (AUC) and segmentation (Dice) tasks. 'ConVIRT-Real' and 'GLORIA-Real' refer to models pre-trained on MIMIC-CXR using real data, while 'ConVIRT-Syn' and 'GLORIA-Syn' indicate models pre-trained on SynCXR using synthetic data. 'ConVIRT-Mix' and 'GLORIA-Mix' represent models trained on a combination of MIMIC-CXR and SynCXR. Best results are in bold.

Mathad	Entity Sampling	Avg. Zero-shot		Method	Real Image	Syn. Image	Real Report	Syn. Report	Avg. Zero-shot Classification
Method	Strategy	Classification	Classification		1	1	1		59.59 61.04
ConVIRT (Zhang et al., 2020)	w/o balance Sampling	<b>63.65</b>		ConVIRT (Zhang et al., 2020)	1	1	-	1	59.36 63.65
GLoRIA (Huang et al., 2021)	w/balance Sampling w/o balance Sampling	<b>61.87</b> 58.42		GLoRIA (Huang et al., 2021)	'   '	1	1	1	57.83 58.62 57.69 61 87
	1 10	I			1	•		•	01107

(a) Impact of Entity Sampling Strategies

(b) Impact of Different Synthetic Data

499

500

501

502

503

504

505

506

507

508

509

Table 4: Evaluation of entity sampling strategies for synthetic report generation and the impact of synthetic data types on MedVLP.

els using these real image and synthetic report pairs.

488

489

490

491

492

493

494

495

496

497

• Real Report, Synthetic Image: In this setting, we use RoentGen (Bluethgen et al., 2024), a text-to-image model, to generate synthetic CXR images for each real report. The 'IMPRESSION' section of each report serves as the prompt for generating synthetic CXR images. These synthetic image and real report pairs are used to train MedVLP models.

498 According to Tab 4b, for both MedVLP methods,

using real images with synthetic reports results in decreased performance, likely due to the persistent long-tailed distribution, as the synthetic reports are generated based on real images. However, using real reports with synthetic images slightly improves performance, as synthetic images can be curated using our image filtering procedure to ensure high quality, avoiding issues commonly found in real datasets. Using both synthetic images and synthetic reports achieves the highest performance, indicating that a well-curated synthetic dataset can



Figure 3: Effectiveness of various factors on SynCXR dataset. **Top:** Impact of entity usage ratio on MedVLP performance for ConVIRT and GLoRIA methods. **Bottom Left:** Effectiveness of different LLMs for report generation on both MedVLP methods. **Bottom Right:** Effectiveness of different CXR image generation models for both MedVLP methods.

significantly enhance MedVLP performance. Furthermore, We evaluate the MedVLP method on
the cleaned MIMIC-CXR dataset, as detailed in
Section D and Table 5

Impact of Entity Diversity. We evaluate the im-514 pact of entity diversity by varying the number of 515 entities used for generating the SynCXR dataset. 516 We generate synthetic datasets using 25%, 50%, 517 and 75% of these entities, following the same pro-518 cedure each time. The results, shown in Fig 3 519 (Top), indicate that zero-shot classification perfor-520 mance improves as more entities are used for report 521 generation. This suggests that increasing dataset 522 523 diversity positively influences downstream performance. 524

Impact of Different Report Generators. We also 525 examine the impact of using different LLMs for 526 synthetic report generation. As shown in Fig 3 527 (Bottom Left), we compare two general LLMs, 528 Llama 3.1 (8B and 70B), and two medical-specific 529 LLMs, Meditron3 (8B and 70B) (OpenMeditron, 530 2025). Despite Meditron3 being trained specifi-531 cally on medical corpora and inheriting weights 532 from Llama, the dataset generated by Llama 3.1-533 70B-Instruct achieves the best performance. This 534 indicates that a powerful general LLM is effec-535 tive for generating synthetic datasets, and using 536 domain-specific fine-tuned versions may degrade 537 the quality of the synthetic data. 538

Impact of Different Image Generators. We evaluate various text-to-image models for synthetic CXR
image generation, including CXR-IRGen (Shentu and Al Moubayed, 2024), LLM-CXR (Lee et al., 2023), and RoentGen (Bluethgen et al., 2024). As
shown in Fig 3 (Bottom Right), datasets generated

by RoentGen lead to the best performance for both545MedVLP methods. This is likely because Roent-546Gen is the only image generation model verified547by clinicians, suggesting that the quality of image548generation models is crucial for building synthetic549datasets, and models should be validated by clinical550experts.551

552

## 5 Conclusion

In this work, we tackle the question: *Can MedVLP* 553 succeed using purely synthetic data? Our findings 554 demonstrate that the answer is: Yes. To the best 555 of our knowledge, this is the first study to compre-556 hensively explore the potential of synthetic data 557 for MedVLP models. We also identify key limita-558 tions in existing real-world datasets and introduce 559 SynCXR—a synthetic dataset of 200,000 image-560 text pairs generated without any manual quality 561 checks. Our findings show that MedVLP models 562 trained on purely synthetic data outperform those 563 trained on real data. Moreover, combining syn-564 thetic and real data further boosts model perfor-565 mance, demonstrating the potential of synthetic 566 data to overcome limitations in real-world datasets. 567 We systematically analyze key factors in SynCXR 568 and validate its effectiveness through extensive ab-569 lation studies. In summary, we show that Med-570 VLP achieves strong performance using a purely 571 synthetic image-text dataset and benefits signifi-572 cantly from a combination of real and synthetic 573 data. We believe this work will inspire the commu-574 nity to fully leverage synthetic data and mitigate the 575 challenges posed by noisy and limited real-world 576 datasets. 577

## 578 Limitation

While our work demonstrates that MedVLP can 579 successfully operate using purely synthetic data, 580 there are several limitations to consider. Firstly, 581 the success of the synthetic data approach heavily 582 relies on the capabilities of the language and im-583 age generation models. Moreover, the synthetic 584 data generation process requires a filtering stage, 585 which introduces additional computational over-586 head. Although MedVLP is capable of handling 587 noisy data-an issue also present in real-world 588 datasets-the imperfect pairing of synthetic data 589 may still present challenges. Evaluating the real-590 ism of synthetic data through human judgment is 591 valuable but costly. In future, we aim to design 592 593 more efficient data filtering methods and develop metrics that can better simulate human evaluation 594 to enhance data quality. 595

#### 596 References

625

626 627

628

629

630

631

632 633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

- 597 AI@Meta. 2024. Llama 3 model card.
- 598 Shekoofeh Azizi, Simon Kornblith, Chitwan Saharia,
  599 Mohammad Norouzi, and David J. Fleet. 2023. Syn600 thetic data from diffusion models improves imagenet
  601 classification. *TMLR*.
- Shruthi Bannur, Stephanie Hyland, Qianchu Liu, Fer-602 nando Perez-Garcia, Maximilian Ilse, Daniel C Cas-603 tro, Benedikt Boecking, Harshita Sharma, Kenza 604 Bouzid, Anja Thieme, et al. 2023. Learning to 605 606 exploit temporal structure for biomedical vision-607 language processing. In Proceedings of the 608 IEEE/CVF Conference on Computer Vision and Pat-609 tern Recognition, pages 15016–15027.
- 610 Christian Bluethgen, Pierre Chambon, Jean-Benoit Delbrouck, Rogier van der Sluijs, Małgorzata Połacin, 611 Juan Manuel Zambrano Chaves, Tanishq Mathew 612 Abraham, Shivanshu Purohit, Curtis P Langlotz, and 613 Akshay S Chaudhari. 2024. A vision-language foun-614 615 dation model for the generation of realistic chest 616 x-ray images. Nature Biomedical Engineering, pages 1 - 13.617
- 618Benedikt Boecking, Naoto Usuyama, Shruthi Bannur,619Daniel C Castro, Anton Schwaighofer, Stephanie620Hyland, Maria Wetscherek, Tristan Naumann, Aditya621Nori, Javier Alvarez-Valle, et al. 2022. Making the622most of text semantics to improve biomedical vision-623language processing. In European conference on624computer vision, pages 1–21. Springer.
  - Aurelia Bustos, Antonio Pertusa, Jose-Maria Salinas, and Maria de la Iglesia-Vayá. 2020. Padchest: A large chest x-ray image dataset with multi-label annotated reports. *Medical image analysis*, 66:101797.
    - Chen Chen, Chen Qin, Cheng Ouyang, Zeju Li, Shuo Wang, Huaqi Qiu, Liang Chen, Giacomo Tarroni, Wenjia Bai, and Daniel Rueckert. 2022. Enhancing mr image segmentation with realistic adversarial data augmentation. *Medical Image Analysis*, 82:102597.
      - Qi Chen, Xiaoxi Chen, Haorui Song, Zhiwei Xiong, Alan Yuille, Chen Wei, and Zongwei Zhou. 2024a. Towards generalizable tumor synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11147–11158.
  - Yuhua Chen, Wen Li, Xiaoran Chen, and Luc Van Gool. 2019. Learning semantic segmentation from synthetic data: A geometrically guided input-output adaptation approach. In *CVPR*.
  - Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. 2024b. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*.
- 649 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su,650 Guo Chen, Sen Xing, Muyan Zhong, Qinglong

Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. 2023. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*. 651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

- Jean-Benoit Delbrouck, Pierre Chambon, Zhihong Chen, Maya Varma, Andrew Johnston, Louis Blankemeier, Dave Van Veen, Tan Bui, Steven Truong, and Curtis Langlotz. 2024. Radgraph-xl: A large-scale expert-annotated dataset for entity and relation extraction from radiology reports. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 12902–12915.
- Shivang Desai, Ahmad Baghal, Thidathip Wongsurawat, Piroon Jenjaroenpun, Thomas Powell, Shaymaa Al-Shukri, Kim Gates, Phillip Farmer, Michael Rutherford, Geri Blake, et al. 2020. Chest imaging representing a covid-19 positive rural us population. *Scientific data*, 7(1):414.
- Lijie Fan, Kaifeng Chen, Dilip Krishnan, Dina Katabi, Phillip Isola, and Yonglong Tian. 2023. Scaling laws of synthetic images for model training... for now. *arXiv preprint arXiv:2312.04567*.
- Hasan Abed Al Kader Hammoud, Hani Itani, Fabio Pizzati, Philip Torr, Adel Bibi, and Bernard Ghanem. 2024. Synthclip: Are we ready for a fully synthetic clip training? *arXiv preprint arXiv:2402.01832*.
- Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip Torr, Song Bai, and XIAOJUAN QI. 2023. Is synthetic data from generative models ready for image recognition? In *ICLR*.
- Shih-Cheng Huang, Liyue Shen, Matthew P Lungren, and Serena Yeung. 2021. Gloria: A multimodal global-local representation learning framework for label-efficient medical image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3942–3951.
- Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 590–597.
- Ali Jahanian, Xavier Puig, Yonglong Tian, and Phillip Isola. 2022. Generative models as a data source for multiview representation learning. In *ICLR*.
- Qiao Jin, Won Kim, Qingyu Chen, Donald C Comeau, Lana Yeganova, W John Wilbur, and Zhiyong Lu.
  2023. Medcpt: Contrastive pre-trained transformers with large-scale pubmed search logs for zero-shot biomedical information retrieval. *Bioinformatics*, 39(11):btad651.
- Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz,<br/>Nathaniel R Greenbaum, Matthew P Lungren, Chih-<br/>ying Deng, Roger G Mark, and Steven Horng.704704705705706

- 2019a. Mimic-cxr, a de-identified publicly available
  database of chest radiographs with free-text reports. *Scientific data*, 6(1):1–8.
- Alistair EW Johnson, Tom J Pollard, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Yifan Peng, Zhiyong Lu, Roger G Mark, Seth J Berkowitz, and Steven Horng. 2019b. Mimic-cxr-jpg, a large publicly available database of labeled chest radiographs. arXiv preprint arXiv:1901.07042.
- 716 Matthew Johnson-Roberson, Charles Barto, Rounak
  717 Mehta, Sharath Nittur Sridhar, Karl Rosaen, and Ram
  718 Vasudevan. 2017. Driving in the matrix: Can virtual
  719 worlds replace human-generated annotations for real
  720 world tasks? In *ICRA*.
- 721 Bardia Khosravi, Frank Li, Theo Dapamede, Pouria
  722 Rouzrokh, Cooper U Gamble, Hari M Trivedi,
  723 Cody C Wyles, Andrew B Sellergren, Saptarshi
  724 Purkayastha, Bradley J Erickson, et al. 2024. Synthetically enhanced: unveiling synthetic data's potential
  726 in medical imaging research. *EBioMedicine*, 104.

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

- Lennart R Koetzier, Jie Wu, Domenico Mastrodicasa, Aline Lutz, Matthew Chung, W Adam Koszek, Jayanth Pratap, Akshay S Chaudhari, Pranav Rajpurkar, Matthew P Lungren, et al. 2024. Generating synthetic data for medical imaging. *Radiology*, 312(3):e232471.
- Ira Ktena, Olivia Wiles, Isabela Albuquerque, Sylvestre-Alvise Rebuffi, Ryutaro Tanno, Abhijit Guha Roy, Shekoofeh Azizi, Danielle Belgrave, Pushmeet Kohli, Taylan Cemgil, et al. 2024. Generative models improve fairness of medical classifiers under distribution shifts. *Nature Medicine*, pages 1–8.
- Suhyeon Lee, Won Jun Kim, Jinho Chang, and Jong Chul Ye. 2023. Llm-cxr: Instruction-finetuned llm for cxr image understanding and generation. *arXiv preprint arXiv:2305.11490*.
- 743 Guohao Li, Hasan Abed Al Kader Hammoud, Hani
  744 Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023.
  745 CAMEL: Communicative agents for "mind" exploration of large language model society. In *NeurIPS*.
- 747 Zhe Li, Laurence T Yang, Bocheng Ren, Xin Nie,
  748 Zhangyang Gao, Cheng Tan, and Stan Z Li. 2024.
  749 Mlip: Enhancing medical visual representation
  750 with divergence encoder and knowledge-guided con751 trastive learning. *arXiv preprint arXiv:2402.02045*.
- 752 Che Liu, Sibo Cheng, Chen Chen, Mengyun Qiao,
  753 Weitong Zhang, Anand Shah, Wenjia Bai, and
  754 Rossella Arcucci. 2023a. M-flag: Medical vision755 language pre-training with frozen language modr56 els and latent space geometry optimization. arXiv
  757 preprint arXiv:2307.08347.
- 758 Che Liu, Sibo Cheng, Miaojing Shi, Anand Shah, Wenjia Bai, and Rossella Arcucci. 2023b. Imitate: Clinical prior guided hierarchical vision-language pretraining. *arXiv preprint arXiv:2310.07355*.

Che Liu, Cheng Ouyang, Yinda Chen, Cesar César Quilodrán-Casas, Lei Ma, Jie Fu, Yike Guo, Anand Shah, Wenjia Bai, and Rossella Arcucci. 2023c. T3d: Towards 3d medical image understanding through vision-language pre-training. *arXiv preprint arXiv:2312.01529*.

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

- Che Liu, Cheng Ouyang, Sibo Cheng, Anand Shah, Wenjia Bai, and Rossella Arcucci. 2023d. G2d: From global to dense radiography representation learning via vision-language pre-training. *arXiv preprint arXiv:2312.01522*.
- Che Liu, Anand Shah, Wenjia Bai, and Rossella Arcucci. 2023e. Utilizing synthetic data for medical visionlanguage pre-training: Bypassing the need for real images. *arXiv preprint arXiv:2310.07027*.
- OpenMeditron. 2025. OpenMeditron. https:// huggingface.co/OpenMeditron. Accessed: 2025-01-28.
- Maya Pavlova, Naomi Terhljan, Audrey G Chung, Andy Zhao, Siddharth Surana, Hossein Aboutalebi, Hayden Gunraj, Ali Sabri, Amer Alaref, and Alexander Wong. 2022. Covid-net cxr-2: An enhanced deep convolutional neural network design for detection of covid-19 cases from chest x-ray images. *Frontiers in Medicine*, 9:861680.
- Fernando Pérez-García, Harshita Sharma, Sam Bond-Taylor, Kenza Bouzid, Valentina Salvatelli, Maximilian Ilse, Shruthi Bannur, Daniel C Castro, Anton Schwaighofer, Matthew P Lungren, et al. 2024. Rad-dino: Exploring scalable medical image encoders beyond text supervision. *arXiv preprint arXiv:2401.10815*.
- Minh Hieu Phan, Yutong Xie, Yuankai Qi, Lingqiao Liu, Liyang Liu, Bowen Zhang, Zhibin Liao, Qi Wu, Minh-Son To, and Johan W Verjans. 2024a. Decomposing disease descriptions for enhanced pathology detection: A multi-aspect vision-language matching framework. *arXiv preprint arXiv:2403.07636*.
- Vu Minh Hieu Phan, Yutong Xie, Yuankai Qi, Lingqiao Liu, Liyang Liu, Bowen Zhang, Zhibin Liao, Qi Wu, Minh-Son To, and Johan W Verjans. 2024b. Decomposing disease descriptions for enhanced pathology detection: A multi-aspect vision-language pretraining framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11492–11501.
- Chen Qin, Shuo Wang, Chen Chen, Wenjia Bai, and Daniel Rueckert. 2023. Generative myocardial motion tracking via latent space exploration with biomechanics-informed prior. *Medical Image Analysis*, 83:102682.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *ICML*.

Vladlen Koltun. 2016. Playing for data: Ground and Dilip Krishnan. 2023b. Stablerep: Synthetic im-820 truth from computer games. In ECCV. ages from text-to-image models make strong visual 821 representation learners. In NeurIPS. Robin Rombach, Andreas Blattmann, Dominik Lorenz, 822 Patrick Esser, and Björn Ommer. 2022. 823 High-Ekin Tiu, Ellie Talius, Pujan Patel, Curtis P Langlotz, resolution image synthesis with latent diffusion mod-824 Andrew Y Ng, and Pranav Rajpurkar. 2022a. Expert-825 els. In Proceedings of the IEEE/CVF conference level detection of pathologies from unannotated chest 826 on computer vision and pattern recognition, pages x-ray images via self-supervised learning. Nature 10684-10695. 827 Biomedical Engineering, 6(12):1399–1406. Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 828 Ekin Tiu, Ellie Talius, Pujan Patel, Curtis P Langlotz, 2015. U-net: Convolutional networks for biomedical 829 Andrew Y Ng, and Pranav Rajpurkar. 2022b. Expertimage segmentation. In Medical Image Computing 830 level detection of pathologies from unannotated chest and Computer-Assisted Intervention-MICCAI 2015: 831 x-ray images via self-supervised learning. Nature 832 18th International Conference, Munich, Germany, Biomedical Engineering, pages 1–8. 833 October 5-9, 2015, Proceedings, Part III 18, pages 834 234–241. Springer. Gül Varol, Javier Romero, Xavier Martin, Naureen Mahmood, Michael J. Black, Ivan Laptev, and Cordelia 835 German Ros, Laura Sellart, Joanna Materzynska, David Schmid. 2017. Learning from synthetic humans. In 836 Vazquez, and Antonio M. Lopez. 2016. The synthia CVPR. 837 dataset: A large collection of synthetic images for 838 semantic segmentation of urban scenes. In CVPR. Zhongwei Wan, Che Liu, Mi Zhang, Jie Fu, Benyou Wang, Sibo Cheng, Lei Ma, César Quilodrán-Casas, Nick Rossenbach, Albert Zeyer, Ralf Schlüter, and Her-839 and Rossella Arcucci. 2024. Med-unic: Unifying mann Ney. 2020. Generating synthetic audio data 840 cross-lingual medical vision-language pre-training for attention-based speech recognition systems. In 841 by diminishing bias. Advances in Neural Information ICASSP. 842 Processing Systems, 36. 843 Mert Bulent Sariyildiz, Karteek Alahari, Diane Larlus, Fuying Wang, Yuyin Zhou, Shujun Wang, Varut and Yannis Kalantidis. 2023. Fake it till you make it: 844 Vardhanabhuti, and Lequan Yu. 2022. Multi-845 Learning transferable representations from synthetic imagenet clones. In CVPR. granularity cross-modal alignment for generalized 846 medical visual representation learning. arXiv Sahand Sharifzadeh, Christos Kaplanis, Shreya Pathak, 847 preprint arXiv:2210.06044. 848 Dharshan Kumaran, Anastasija Ilic, Jovana Mitrovic, Charles Blundell, and Andrea Banino. 2024. Synth2: 849 Linda Wang, Zhong Qiu Lin, and Alexander Wong. Boosting visual-language models with synthetic 850 2020. Covid-net: A tailored deep convolutional neucaptions and image embeddings. arXiv preprint 851 ral network design for detection of covid-19 cases arXiv:2403.07750. 852 from chest x-ray images. Scientific reports, 10(1):1-12. Junjie Shentu and Noura Al Moubayed. 2024. Cxr-853 irgen: An integrated vision and language model 854 Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mo-855 for the generation of clinically accurate chest x-ray hammadhadi Bagheri, and Ronald M Summers. 2017. 856 image-report pairs. In Proceedings of the IEEE/CVF Chestx-ray8: Hospital-scale chest x-ray database and Winter Conference on Applications of Computer Vi-857 benchmarks on weakly-supervised classification and sion, pages 5212-5221. 858 localization of common thorax diseases. In Proceedings of the IEEE conference on computer vision and George Shih, Carol C Wu, Safwan S Halabi, Marc D 859 pattern recognition, pages 2097-2106. Kohli, Luciano M Prevedello, Tessa S Cook, Arjun 860 Sharma, Judith K Amorosa, Veronica Arteaga, Maya 861 Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, Galperin-Aizenberg, et al. 2019. Augmenting the 862 and Weidi Xie. 2023. Medklip: Medical knowledge national institutes of health chest radiograph dataset 863 enhanced language-image pre-training. medRxiv, with expert annotations of possible pneumonia. Ra-864 pages 2023-01. diology: Artificial Intelligence, 1(1):e180041. 865 Linshan Wu, Jiaxin Zhuang, Xuefeng Ni, and Hao Konstantin Shmelkov, Cordelia Schmid, and Karteek 866 Chen. 2024. Freetumor: Advance tumor segmen-Alahari. 2018. How good is my gan? In ECCV. 867 tation via large-scale tumor synthesis. arXiv preprint arXiv:2406.01264. CIIP Steven G. Langer, PhD and MS George Shih, MD. 868 2019. Siim-acr pneumothorax segmentation. 869 Yutong Xie, Qi Chen, Sinuo Wang, Minh-Son To, Iris Yonglong Tian, Lijie Fan, Kaifeng Chen, Dina Katabi, Lee, Ee Win Khoo, Kerolos Hendy, Daniel Koh, 870 Yong Xia, and Qi Wu. 2024. Pairaug: What can Dilip Krishnan, and Phillip Isola. 2023a. Learning 871 vision from models rivals learning vision from data. augmented image-text pairs do for radiology? arXiv 872 preprint arXiv:2404.04960. arXiv preprint arXiv:2312.17742. 873 12

Yonglong Tian, Lijie Fan, Phillip Isola, Huiwen Chang,

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906 907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

Stephan R Richter, Vibhav Vineet, Stefan Roth, and

- 928 Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang,
  929 Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi
  930 Ghosh, Luke Zettlemoyer, and Christoph Feichten931 hofer. 2023a. Demystifying clip data. *arXiv preprint*932 *arXiv:2309.16671*.
- 933 Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang,
  934 Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi
  935 Ghosh, Luke Zettlemoyer, and Christoph Feichten936 hofer. 2023b. Demystifying clip data. *arXiv preprint*937 *arXiv:2309.16671.*
- 938 Yiben Yang, Chaitanya Malaviya, Jared Fernandez,
  939 Swabha Swayamdipta, Ronan Le Bras, Ji-Ping Wang,
  940 Chandra Bhagavatula, Yejin Choi, and Doug Downey.
  941 2020. Generative data augmentation for common942 sense reasoning. In *EMNLP*.
- 943Qingsong Yao, Li Xiao, Peihang Liu, and S Kevin Zhou.9442021. Label-free segmentation of covid-19 lesions945in lung ct. IEEE transactions on medical imaging,94640(10):2808–2819.

948

949 950

951

952

953

954

955

956

957

958

959

960

961

962

968

969

970

- Zhuoran Yu, Chenchen Zhu, Sean Culatana, Raghuraman Krishnamoorthi, Fanyi Xiao, and Yong Jae Lee. 2023. Diversify, don't fine-tune: Scaling up visual recognition training with synthetic images. arXiv preprint arXiv:2312.02253.
- Jianhao Yuan, Jie Zhang, Shuyang Sun, Philip Torr, and Bo Zhao. 2024. Real-fake: Effective training data synthesis through distribution matching. In *ICLR*.
- Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. 2024. Long-clip: Unlocking the long-text capability of clip. *arXiv preprint arXiv:2403.15378*.
- Xiaoman Zhang, Chaoyi Wu, Ya Zhang, Weidi Xie, and Yanfeng Wang. 2023. Knowledge-enhanced visuallanguage pre-training on chest radiology images. *Nature Communications*, 14(1):4542.
- 963 Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz. 2020.
  965 Contrastive learning of medical visual representations from paired images and text. *arXiv preprint arXiv:2010.00747*.
  - Weike Zhao, Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2024. Ratescore: A metric for radiology report generation. *arXiv preprint arXiv:2406.16845*.
- 972 Hong-Yu Zhou, Subathra Adithan, Julián Nicolás
  973 Acosta, Eric J Topol, and Pranav Rajpurkar. 2024.
  974 A generalist learner for multifaceted medical image
  975 interpretation. *arXiv preprint arXiv:2405.07988*.
- 976 Yongchao Zhou, Hshmat Sahak, and Jimmy Ba.
  977 2023. Training on thin air: Improve image classification with generated data. arXiv preprint 979 arXiv:2305.15316.

#### A Related Work

980

**Representation Learning with Synthetic Data.** 981 Synthetic data has been widely employed across 982 various deep learning fields (Rossenbach et al., 983 2020; Varol et al., 2017; Jahanian et al., 2022; Zhou 984 985 et al., 2023; Yang et al., 2020; Li et al., 2023). In visual representation learning, synthetic data 986 has improved model performance in a range of 987 tasks (Richter et al., 2016; Ros et al., 2016; Chen 988 et al., 2019; Johnson-Roberson et al., 2017; Yuan 989 et al., 2024; Shmelkov et al., 2018). Recent ef-990 forts have also focused on using synthetic data 991 from text-to-image models to augment real-world 992 data during training (Azizi et al., 2023; Sariyildiz 993 et al., 2023; He et al., 2023). For example, (Yu 994 et al., 2023) introduced a framework to generate 995 synthetic images to diversify existing datasets. No-996 tably, methods utilizing text-to-image generative 997 models (Rombach et al., 2022) have demonstrated 998 that synthetic images guided by real captions can 999 effectively train self-supervised models, achiev-1000 1001 ing performance comparable to that of real images (Tian et al., 2023b). 1002

Further advancements like SynCLR (Tian et al., 1003 2023a) have focused on visual representation learn-1004 ing using only synthetic images, generated with 1005 conditioning on various categories. Meanwhile, 1006 1007 other recent works (Fan et al., 2023; Sharifzadeh et al., 2024; Xie et al., 2024) have explored joint 1008 image and text generation for enhanced vision-1009 language pretraining (VLP). However, only one 1010 study, SynthCLIP (Hammoud et al., 2024), inves-1011 tigates VLP exclusively with synthetic data, and 1012 even that work is limited to natural images. To date, 1013 no research has explored the potential of MedVLP 1014 trained solely on synthetic data. 1015

Medical Vision Language Pre-training. Recent 1016 work on MedVLP has focused on integrating visual 1017 and textual modalities, particularly for chest X-ray 1018 (CXR) images. Studies such as (Zhang et al., 2020; 1019 Huang et al., 2021; Wang et al., 2022; Liu et al., 1020 1021 2023b,d,c; Wan et al., 2024) have concentrated on aligning CXR images with paired radiology reports. 1022 Some methods also leverage external datasets to 1023 boost performance, raising concerns about gener-1024 alizability (Wu et al., 2023; Zhang et al., 2023; 1025 1026 Li et al., 2024; Phan et al., 2024a). However, all current MedVLP approaches rely heavily on large-1027 scale, real image-text paired datasets like MIMIC-1028 CXR (Johnson et al., 2019b). Some even require 1029 additional human-annotated datasets or manual in-1030

terventions to improve model performance (Wu1031et al., 2023; Zhang et al., 2023; Phan et al., 2024a),1032which limits their scalability and accessibility.1033

Synthetic Data for Medical Image Tasks. Given 1034 the scarcity of annotated data, high costs, and pri-1035 vacy concerns in medical data collection, synthetic 1036 data has been explored to support various medical 1037 image tasks (Koetzier et al., 2024). However, most 1038 prior work focuses on image modality and super-1039 vised learning (Chen et al., 2024a; Yao et al., 2021; 1040 Chen et al., 2022; Qin et al., 2023), using synthetic 1041 data solely as augmentation for real datasets (Khos-1042 ravi et al., 2024; Ktena et al., 2024). Few studies 1043 have evaluated models trained entirely on synthetic 1044 medical data (Wu et al., 2024). Recent efforts have 1045 generated synthetic text and images for MedVLP 1046 (Xie et al., 2024), but still restrict synthetic data 1047 usage to augmentation. Consequently, the full po-1048 tential of synthetic data in MedVLP remains largely 1049 unexplored. 1050

In this work, we generate both synthetic CXR images and reports, then training a MedVLP model solely on synthetic data. We conduct an extensive evaluation of the impact of large-scale synthetic medical data on MedVLP, exploring its performance across various downstream tasks. 1051

1052

1053

1054

1055

1056

1057

1058

## B Queries for Using MLLM to Assess Issued CXR Images

This section provides detailed queries used to guide1059the MLLM in assessing issued CXR images. Each1060query is designed to evaluate specific aspects of the1061images, ensuring their quality and suitability for1062diagnostic purposes.1063

- Detecting Non-CXR Images: <CXR 1064 Image>, Please check if the 1065 given image is a chest X-ray 1066 scan. If it is a chest X-ray, 1067 return 'YES'. Otherwise, 1068 return 'NO'. 1069
- Detecting Non-Human CXR Images: 1070 <CXR Image>, Please verify if 1071 the given image is a human 1072 chest X-ray scan. If it is 1073 a chest X-ray, return 'YES'. 1074 Otherwise, return 'NO'. 1075
- Detecting Wrong Views: <CXR Image>, 1076 Please check if the given 1077 image is a frontal chest X-ray 1078

- 1079view. If it is a frontal1080view, return 'YES'. If it is1081a lateral view or any other1082view, return 'NO'.
- Assessing Image Quality: <CXR Image>, 1083 Please analyze the provided 1084 chest X-ray (CXR) image and 1085 respond with 'NO' if the image 1086 quality is poor, such as being 1087 blurry, containing artifacts, 1088 or having poor contrast. 1089 Respond with 'YES' if the 1090 1091 image quality is acceptable.
- Detecting Artifacts 1092 and **Overpro**cessing: <CXR Image>, Please 1093 1094 analyze the following chest X-ray image. Respond with 1095 'YES' if the image is clear, 1096 correctly oriented, and free 1097 of artifacts or imperfections 1098 that could affect its 1099 diagnostic quality. 1100 Respond with 'NO' if the image is 1101 blurry, incorrectly oriented, 1102 contains artifacts, or has 1103 imperfections that make 1104 it unsuitable for further 1105 analysis. 1106
- 1107 Checking High-Fidelity: <CXR Image>, 1108 Please check if the given 1109 image is a high-fidelity human 1110 chest X-ray scan. If it is 1111 a high-fidelity chest X-ray, 1112 return 'YES'. Otherwise, 1113 return 'NO'.
- 1114 C Downstream Tasks Configuration
- 1115 C.1 Dataset Details
- 1116In this section, we provide details on all datasets1117used. The dataset splits are publicly accessible at9.1118ChestX-ray14 (Wang et al., 2017) includes1119112,120 frontal X-rays from 30,805 patients, la-1120beled for 14 diseases. We use the official1121split and partition it into 80%/10%/10% for1122train/validation/test.
- 1123PadChest (Bustos et al., 2020) includes 160,868 X-1124rays from 67,000 patients, annotated with over 1501125findings. As in (Phan et al., 2024b), three subsets

are built based on PadChest: 14 common diseases1126as PadChest-seen, rare diseases from the NORD1127database<sup>10</sup> as PadChest-rare, and the remaining1128diseases as PadChest-unseen. We use the official1129split provided by (Phan et al., 2024b).1130RSNA (Shih et al., 2019) contains over 260,0001131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1174

**RSNA** (Shih et al., 2019) contains over 260,000 frontal X-rays annotated with pneumonia masks. We divide it into training (60%), validation (20%), and test (20%) sets for segmentation and classification tasks (Huang et al., 2021; Wu et al., 2023).

**CheXpert** (Irvin et al., 2019) contains 224,316 chest X-rays from 65,240 patients at Stanford Hospital, with an official validation set of 200 studies and a test set of 500 studies, both annotated by board-certified radiologists. Our evaluation on the five observations in the official test set follows protocols from earlier studies (Tiu et al., 2022a; Irvin et al., 2019).

**SIIM** (Steven G. Langer and George Shih, 2019) consists of over 12,000 frontal X-rays annotated with pneumothorax masks, split into training (60%), validation (20%), and test (20%) sets.

**COVID**x **CXR-2** (Wang et al., 2020) includes 29,986 X-rays from 16,648 COVID-19 patients, divided into training (70%), validation (20%), and test (10%) (Pavlova et al., 2022).

**COVID Rural** (Desai et al., 2020) contains over 200 X-rays with segmentation masks, divided into training (60%), validation (20%), and test (20%).

## C.2 Implementation Details

Zero-shot Image Classification. The CXR im-1156 ages undergo a two-step preprocessing: resizing 1157 to  $256 \times 256$ , followed by center cropping to 1158  $224 \times 224$ . As per (Huang et al., 2021), pixel values 1159 are normalized to [0, 1]. The processed image is 1160 passed through a visual encoder and projector to 1161 generate the image embedding  $\hat{\mathbf{v}}_i$ . Simultaneously, 1162 the text prompts are processed through a text en-1163 coder to obtain text embeddings  $l_i$ . Classification is 1164 based on cosine similarity between image and text 1165 embeddings. If the similarity between the image 1166 embedding and the positive prompt (e.g., *disease*) 1167 is higher than that with the negative prompt (e.g., 1168 *No disease*), the classification is positive, and vice 1169 versa. The prompt design follows (Tiu et al., 2022b) 1170 for both ConVIRT and GLoRIA. 1171 Zero-shot Visual Grounding. For this task, we 1172 follow the BioViL pipeline as described in (Phan 1173

et al., 2024b), since ConVIRT (Zhang et al., 2020)

<sup>&</sup>lt;sup>9</sup>https://github.com/HieuPhan33/CVPR2024\_MAVL/tree/main/datahttps://rarediseases.org/rare-diseases/

and GLoRIA (Huang et al., 2021) do not provide 1175 code for visual grounding. This pixel-level classi-1176 fication task relies on the similarity between text 1177 embeddings and the dense visual feature map from 1178 the final convolutional layer. The cosine similarity 1179 generates a similarity map, resized to match the im-1180 age, and used as segmentation results for grounding 1181 evaluation. 1182

#### Medical Image Fine-tuned Classification. 1183

For fine-tuning, we follow the experimental 1184 setup from (Phan et al., 2024b), updating both the 1185 visual encoder and linear layer. Images are resized 1186 to  $256 \times 256$ , and data augmentation is applied as 1187 recommended in (Zhang et al., 2023). We use the 1188 AdamW optimizer with a learning rate of  $1 \times 10^{-4}$ , 1189 batch size of 64, for 50 epochs on a single A100 1190 GPU. Early stopping is applied, with a learning rate 1191 of 5e-4 and batch size of 8. AdamW is configured 1192 with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , and weight decay of 1193 1194 1e-6.

Medical Image Fine-tuned Segmentation. For 1195 segmentation tasks on the RSNA (Shih et al., 2019), 1196 SIIM (Steven G. Langer and George Shih, 2019), 1197 and Covid-19 Rural (Wang et al., 2020) datasets, 1198 we fine-tune both the pre-trained vision encoder 1199 and decoder. Training is performed with early stop-1200 ping at 50 epochs, using a learning rate of 2e-4 and 1201 weight decay of 0.05. AdamW is the optimizer, 1202 with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . Batch sizes are 1203 8 for SIIM and 16 for RSNA. All configurations 1204 follow the protocol from (Huang et al., 2021). 1205

#### **Results on Cleaned MIMIC-CXR** D

1206

1221

We re-pretrained ConVIRT on the cleaned version 1207 of the MIMIC-CXR dataset, and the results are 1208 presented in Table 5. As shown, the VLP model 1209 pre-trained on the cleaned MIMIC-CXR dataset 1210 achieves slightly better performance than the model 1211 trained on the original, uncleaned dataset. How-1212 ever, it still falls short when compared to the model 1213 pre-trained on our synthetic SynCXR dataset. This 1214 performance gap can be attributed to two key fac-1215 tors: 1216

• Despite filtering out a large number of low-1217 1218 quality samples from the real dataset, it remains challenging to completely remove all 1219 poor or unpaired samples. Consequently, 1220 some problematic samples persist in the dataset, negatively affecting the VLP process. 1222

• The cleaning process for MIMIC-CXR pri-1223 marily targeted image quality, but did not ad-1224



Figure 4: Distribution of Synthetic and Real Data. (a): Comparison of the first principal component distribution of features extracted from RAD-DINO for synthetic and real images. (b): Comparison of the first principal component distribution of features extracted from Med-CPT for synthetic and real reports.

dress the long-tailed distribution of entities in 1225 the dataset. Simply downsampling or over-1226 sampling image-text pairs is not a viable so-1227 lution. This issue limits the representational 1228 balance of the cleaned dataset and impacts 1229 overall model performance. 1230

1231

1245

1246

1247

1248

1249

1250

1251

1252

1253

#### **Extra Visualization** Е

Distribution of Synthetic and Real Data. We il-1232 lustrate the distribution of synthetic and real data 1233 in Fig 4. For visualization, we use RAD-DINO 1234 (Pérez-García et al., 2024) to extract image features 1235 and Med-CPT (Jin et al., 2023) to extract report 1236 features. We then apply Principal component anal-1237 ysis (PCA) to reduce the feature dimensions and 1238 visualize the first principal component. As shown 1239 in Fig 4, the synthetic data covers a broader range 1240 than the real data, indicating greater diversity. In 1241 contrast, the real data shows a more concentrated 1242 distribution, which may limit the generalizability 1243 of MedVLP models. 1244

Pipeline of Synthetic Report Generation. The pipeline for generating synthetic reports using LLMs and balanced sampled clinical entities is illustrated in Fig 5.

Entities Distribution. We visualize the distribution of each type of entity in the MIMIC-CXR dataset. Due to space constraints, only the top 200 most frequent entities are shown, revealing a clear long-tailed distribution in Fig 6, 10, 8, 7, and 9.

Method	Pre-training Data	CheX	pert	ChestX	ray-14	PadChe	st-seen	RSI	NA	SII	М
		AUC ↑	F1 ↑	AUC ↑	$F1\uparrow$	AUC ↑	F1 ↑	AUC ↑	F1 ↑	AUC ↑	F1 ↑
ConVIRT	MIMIC-CXR	52.10	35.61	53.15	12.38	63.72	14.56	79.21	55.67	64.25	42.87
ConVIRT	MIMIC-CXR (cleaned)	53.85	36.45	53.80	13.25	63.51	14.73	80.15	56.10	65.70	44.10
ConVIRT	SynCXR	59.49	40.51	56.07	15.43	63.43	15.10	82.08	58.38	75.55	57.43

Table 5: Performance comparison of ConVIRT pre-trained on different datasets.



Figure 5: Pipeline for generating synthetic reports. The process begins by generating the 'FINDINGS' section, followed by summarizing it into the 'IMPRESSION' section. Both sections are checked to ensure they contain the specified entities; if not, the generation process is repeated. The final dataset includes 200,000 synthetic reports, each containing both 'FINDINGS' and 'IM-PRESSION' sections.



Figure 6: Top 200 most frequent abnormality entities.



Figure 7: Top 200 most frequent non-abnormality entities.



Figure 8: Top 200 most frequent disease entities.



Figure 9: Top 200 most frequent non-disease entities.



Figure 10: Top 200 most frequent anatomy entities.