
MedTrinity-25M: A Large-scale Multimodal Dataset with Multigranular Annotations for Medicine

Yunfei Xie^{1,*}, Ce Zhou^{1,*}, Lang Gao^{1,*}, Juncheng Wu^{2,*}, Xianhang Li³,
Hong-Yu Zhou⁴, Sheng Liu⁵, Lei Xing⁵, James Zou⁵, Cihang Xie³, Yuyin Zhou³
*equal technical contribution

¹Huazhong University of Science and Technology, ²Tongji University,
³UC Santa Cruz, ⁴Harvard University, ⁵Stanford University

Abstract

1 This paper introduces MedTrinity-25M, a comprehensive, large-scale multimodal
2 dataset for medicine, covering over 25 million images across 10 modalities, with
3 multigranular annotations for more than 65 diseases. These enriched annotations
4 encompass both global textual information, such as disease/lesion type, modality,
5 region-specific descriptions, and inter-regional relationships, as well as detailed
6 local annotations for regions of interest (ROIs), including bounding boxes, seg-
7 mentation masks. Unlike existing approach which is limited by the availability
8 of image-text pairs, we have developed the first automated pipeline that scales
9 up multimodal data by generating multigranular visual and textual annotations (in
10 the form of image-ROI-description triplets) without the need for any paired text
11 descriptions. Specifically, data from over 90 different sources have been collected,
12 preprocessed, and grounded using domain-specific expert models to identify ROIs
13 related to abnormal regions. We then build a comprehensive knowledge base
14 and prompt multimodal large language models to perform retrieval-augmented
15 generation with the identified ROIs as guidance, resulting in multigranular tex-
16 tual descriptions. Compared to existing datasets, MedTrinity-25M provides the
17 most enriched annotations, supporting a comprehensive range of multimodal tasks
18 such as captioning and report generation, as well as vision-centric tasks like clas-
19 sification and segmentation. This dataset can be utilized to support large-scale
20 pre-training of multimodal medical AI models, contributing to the development of
21 future foundation models in the medical domain. The dataset is publicly available
22 at <https://yunfeixie233.github.io/MedTrinity-25M/>.

23 1 Introduction

24 Large-scale multimodal foundation models [1, 2, 3, 4, 5] have demonstrated remarkable success
25 across various domains due to their ability to understand complex visual patterns in conjunction with
26 natural language. This success has sparked significant interest in applying such models to medical
27 vision-language tasks. Much progress has been made to improve the medical capacity of general
28 domain multimodal foundation models by constructing medical datasets with image-text pairs and
29 fine-tuning general domain models on these datasets [6, 7, 8, 9, 10].

30 However, current medical datasets have several limitations. Firstly, these datasets lack **multigranular**
31 annotations that reveal the correlation between local and global information within medical images.

32 Medical images often contain detailed cues, such as regional abnormal textures or structures, which
33 may indicate specific types of lesions. Therefore, multimodal models need the ability to infer global
34 information, such as disease or lesion type, from local details. The absence of such data limits
35 the models’ capacity to comprehensively understand medical images. Moreover, current dataset
36 construction methods heavily rely on medical images paired with reports or captions, which restricts
37 their scalability.

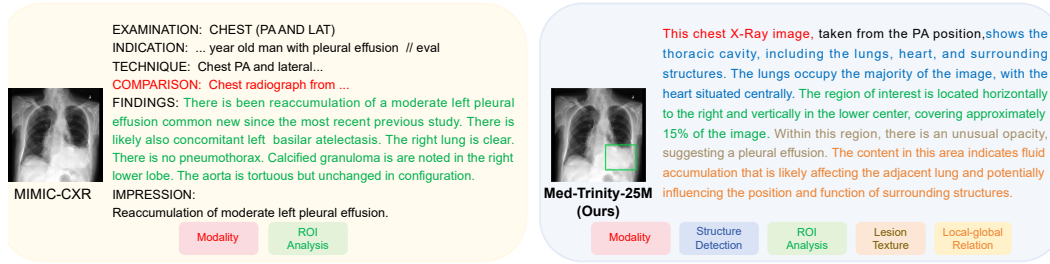
38 In this paper, we address the above challenges by proposing an automated data construction pipeline
39 using multimodal large language models (MLLMs) without relying on paired text descriptions. To
40 address the lack of comprehensive medical knowledge in general-purpose MLLMs, we leverage
41 domain-specific expert grounding models and retrieval-augmented generation (RAG) to extract
42 relevant medical knowledge. We then prompt MLLMs to generate multigranular visual and textual
43 annotations enriched with this knowledge based on identified regions of interest (ROIs). We utilize
44 this pipeline to transform the collected data, including large-scale unpaired images, into image-
45 ROI-description triplets. These triplets provide multigranular annotations that encompass both
46 global textual information, such as disease/lesion type, modality, and inter-regional relationships,
47 as well as detailed local annotations for ROIs, including bounding boxes, segmentation masks, and
48 region-specific textual descriptions. Using the proposed pipeline, we create a large-scale multimodal
49 multigranular medical dataset containing over 25 million triplets, named **MedTrinity-25M**. To our
50 best knowledge, this is the largest multimodal dataset in medicine to date.

51 Initially, we assemble a large amount of medical data from over 90 online resources such as TCIA,
52 Kaggle, Zenodo, Synapse, etc. In addition to images with a small amount of high-quality paired
53 manual reports, this assembled data also includes two types of coarse medical data: 1) Image
54 data with segmentation masks, lesion bounding boxes, or only disease types but lacking detailed
55 textual descriptions, and 2) Images paired with coarse captions that describe only global modality
56 or disease information, but lack detailed descriptions of local regions. To generate multigranular
57 annotations from the massive coarse medical data, we first identify ROIs that contain disease or lesion
58 patterns by applying expert grounding models. We then build a comprehensive knowledge base from
59 online corpora (e.g., PubMed) and retrieve image-related medical knowledge. Finally, we prompt
60 MLLMs to integrate medical knowledge with guidance of identified ROIs to generate multigranular
61 textual descriptions.

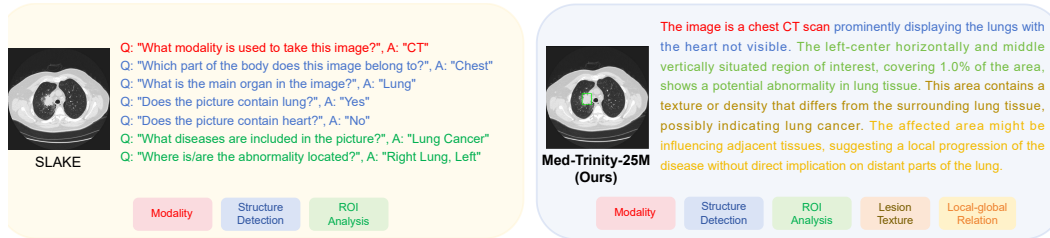
62 2 Related Work

63 **Medical Multimodal Foundation Models.** Due to the effectiveness of multimodal foundation
64 models in understanding visual features, adapting these models to perform medical vision-language
65 tasks has garnered increasing attention in recent years [11, 12, 9, 5]. Several papers attempt to
66 adapt general domain multimodal foundation models with varying architecture to medical domain
67 through end-to-end training on medical datasets. For example, Med-Flamingo [11] enhances the
68 medical capacity of OpenFlamingo-9B [13] by fine-tuning it with 0.8M interleaved and 1.6M
69 paired medical image-text data. While Med-PaLM [12] adapts PaLM-E [14] to medical domain
70 using approximately 1M medical data points, demonstrating competitive or surpassing performance
71 compared to state-of-the-art models. Additionally, LLaVA-Med [9] employs end-to-end visual
72 instruction tuning [1] with two stages, achieving remarkable results in medical Visual Question
73 Answering (VQA) tasks. Similarly, Med-Gemini [15] employs a long-form question answering
74 dataset to enhance the multimodal and long-context capabilities of baseline Gemini [16]. Although
75 these models have achieved remarkable performance, they are still limited by the scale of training
76 data. Prior research [17] has shown that scaling up the training data improves the performance of
77 large multimodal foundation models. In this paper, we aim to build a large-scale medical dataset to
78 facilitate the development of more powerful medical multimodal foundation models.

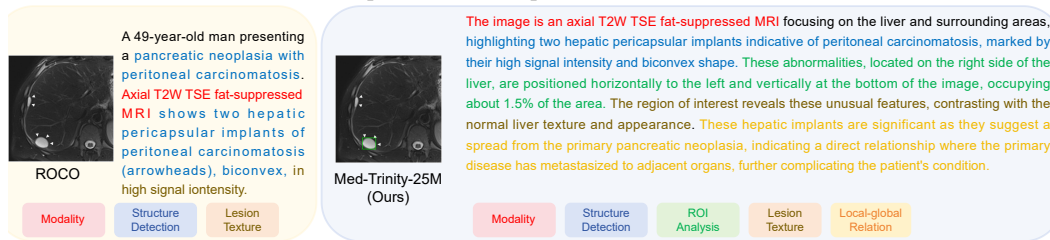
79 **Multimodal Datasets for medicine.** The significance of constructing comprehensive medical
80 multimodal datasets has garnered considerable attention [9, 18, 19, 7]. Several works attempt to
81 collect images and paired clinical reports prepared by pathology specialist [19, 7, 8], which provide



(a) Qualitative Comparison with sample in radiology report of chest x-rays dataset MIMIC-CXR [21].



(b) Qualitative Comparison with sample in visual QA dataset SLAKE [22].



(c) Qualitative Comparison with sample in radiology objects caption dataset ROCO [18].

Figure 1: Qualitative comparison with different types of dataset.

82 comprehensive descriptions of images, including disease types and corresponding reasoning. For
 83 example, MIMIC-CXR[8] comprises 227,835 images for 65,379 patients, containing pathological
 84 findings and impressions in reports paired with each images. However, manually constructing such
 85 reports is both time-consuming and expensive, thereby limiting the scale of these datasets. PMC-
 86 OA [20] aims to expand the dataset scale by extracting a large number of image-caption pairs from
 87 medical papers, increasing the number of data samples to 1.65 million. However, the extracted
 88 captions are less detailed compared to manual clinical reports, resulting in a lack of multigranular
 89 annotations. RadGenome-Chest CT [19] includes more detailed annotations, such as segmentation
 90 masks and medical reports generated by MLLMs. Nonetheless, its construction method still relies
 91 on paired image-text data, which limits its scalability. Unlike these existing methods, we devise the
 92 first automated data construction pipeline to generate multigranular annotations for unpaired images,
 93 achieving a comprehensive multigranular dataset with 25 million data samples.

94 3 MedTrinity-25M Dataset

95 3.1 Data Triplet

96 Our dataset comprises triplets of {image, ROI, description}. Each ROI is associated with an
 97 abnormality and is represented by a bounding box or a segmentation mask, specifying the relevant
 98 region within the image. For each image, we provide a multigranular textual description, which
 99 includes the disease/lesion type, modality, region-specific description, and inter-regional relationships
 100 as illustrated in Figure 2.

101 **Images.** We use the original medical image in the source dataset, we extensively collected medical
102 datasets from the following sources: (1) online resources such as TCIA, Kaggle, Zenodo, Synapse,
103 Hugging Face, Grand Challenge, GitHub, etc. (2) relevant medical dataset research, such as CheX-
104 pert [7] and DeepLesion [23]. These datasets were first categorized into two types: (1) datasets
105 containing local annotations, such as MIMIC-CXR [8] with corresponding radiology reports, and
106 PMC-OA [24] with corresponding captions, where the reports or captions provide analysis of specific
107 local conditions in the images; another example is the 3D image segmentation dataset BraTS2024 [25],
108 which marks the tumor regions in CT scans with masks. (2) datasets containing global annotations:
109 such as image classification datasets ISIC2019 [26] and ISIC2020 [27], whose classification labels
110 reflect the overall pathological condition of tissue sections; another example is the CheXpert [7]
111 dataset, which provides detailed classification of disease types for each chest X-ray. We collect
112 25,001,668 samples spanning 10 modalities and over 65 diseases. For 3D volumetric images stored
113 in DICOM or NIfTI formats, we converted each 2D slice to PNG format. Additional caption and
114 annotations like masks and bounding boxes from these datasets were utilized to construct ROIs and
115 corresponding textual descriptions as below.

116 **ROIs.** For each image, ROIs are highlighted using segmentation masks or bounding boxes. These
117 ROIs mostly contain pathological findings such as lesions, inflammation, neoplasms, infections, or
118 other potential abnormalities. In the few cases without abnormalities, the ROIs generally indicate the
119 primary object or organ in the image, as shown in examples in the supplementary material.

120 **Textual Descriptions.** The textual descriptions for each image are provided with detailed infor-
121 mation across various aspects. Unlike the unstructured free-text descriptions found in previous
122 medical report datasets [7, 8, 6] or simple short sentences in visual QA dataset [28, 22] and caption
123 dataset [18, 24], our textual descriptions are multigranular and structured. General attributes related to
124 the image are described first, including the image modality, the specific organ depicted, and the type
125 of disease presented. Subsequently, ROI-related information is provided, including their locations
126 and the abnormal characteristics within them that indicate underlying pathology, such as distinctive
127 color and texture. Additionally, comparisons between the ROIs and surrounding regions are presented
128 to highlight differences in features and the extent of disease progression.

129 We also demonstrate the multigranular textual descriptions in our dataset with those in other common
130 forms. As illustrated in Figure 1, our textual description is multigranular with more attributes
131 than radiology report of chest x-rays dataset MIMIC-CXR [21], visual QA dataset SLAKE [22] and
132 radiology objects caption dataset ROCO [18].

133 3.2 Data Construction Pipeline

134 Given a medical image, we aim to generate corresponding multigranular visual and textual annotations
135 by leveraging MLLMs. Specifically, as shown in Figure 2, our pipeline can be decomposed into two
136 stages - **Data Processing** and **Generation of Multigranular Text Description**. In the **Data Pro-**
137 **cessing** stage (Section 3.2.1), we address the lack of domain-specific knowledge in general-purpose
138 MLLMs by leveraging expert grounding models and retrieval-augmented generation (RAG). This
139 stage includes three key steps: 1) **Metadata Integration** to produce coarse captions encapsulating
140 fundamental image information such as modality and disease types; 2) **ROI Locating** to identify
141 regions of abnormalities; and 3) **Medical Knowledge Retrieval** to extract relevant fine-grained
142 medical details. Based on the processed data, we then prompt MLLMs to generate multigranular text
143 descriptions, resulting in the creation of fine-grained captions, as detailed in Section 3.2.2.

144 3.2.1 Data Processing

145 **Coarse Caption Generation via Metadata Integration.** We aim to generate coarse captions that
146 provide fundamental information for a given image, including modality, organ labels, disease types,
147 and optionally, camera views and equipment information. Instead of extracting features directly from
148 the images, we generate these captions by integrating dataset metadata. We first extract metadata from

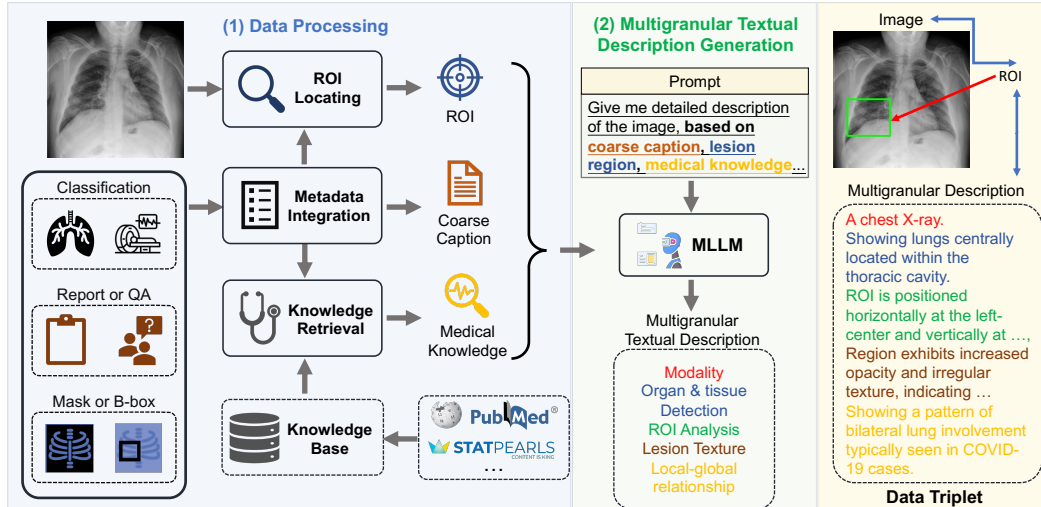


Figure 2: **Data construction pipeline.** 1) Data processing: extracting essential information from collected data, including **metadata integration** to generate coarse caption, **ROI locating**, and **medical knowledge collection**. 2) Multigranular textual description generation: using this information to prompt MLLMs to generate fine-grained captions.

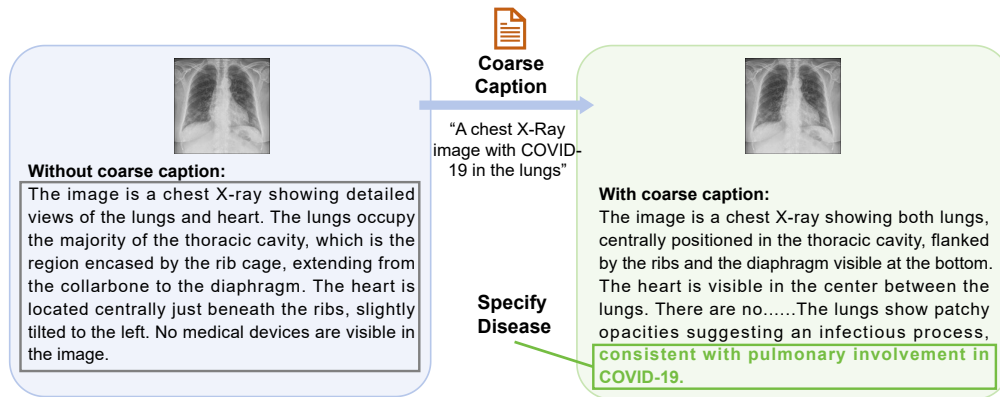


Figure 3: **A qualitative comparison example of generated textual description with and without coarse caption.** Without a coarse caption, MLLMs fails to detect diseases. On the contrary, providing a caption mentioning “COVID-19” allows MLLMs to identify and categorize the disease, facilitating further analysis.

149 the datasets and then apply a fixed rule to integrate this information into coarse captions. For example,
 150 for an image from the QaTa-COV19 dataset¹, we derive metadata from the dataset’s accompanying
 151 paper or documentation, indicating that it consists of COVID-19 chest X-ray images. Next, we
 152 construct coarse captions like “A chest X-ray image with COVID-19 in the lungs” highlighting the
 153 modality, organ types, and disease labels. If the image contains additional textual information like
 154 radiological findings, this is also integrated to enhance the richness of the caption. The effectiveness
 155 of adding coarse captions when generating fine-grained captions is illustrated in Figure 3. In contrast
 156 to the scenario without a coarse caption where MLLMs fails to recognize the disease, providing
 157 MLLMs with a coarse caption that includes the disease type “COVID-19” enables it to identify and
 158 categorize the disease, thereby laying the foundation for further analysis.

159 **ROI Locating.** We employ various strategies to locate Regions of Interest (ROIs) in images. For
 160 datasets that already include localization annotations, such as segmentation masks or bounding boxes,
 161 we derive the ROIs from these existing annotations. Specifically, bounding boxes are directly used

¹<https://www.kaggle.com/aysendegerli/qatacov19-dataset>.

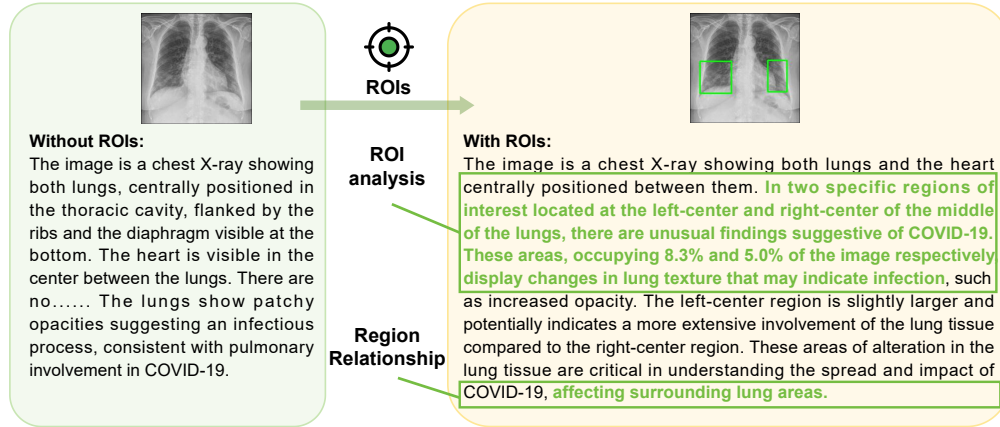


Figure 4: A qualitative comparison example of generated textual description with and without locating ROIs. Without ROIs, the caption offers only a brief global analysis; with ROIs, MLLMs conducts detailed local analysis and assesses the impact of lesion ROIs on adjacent normal regions.

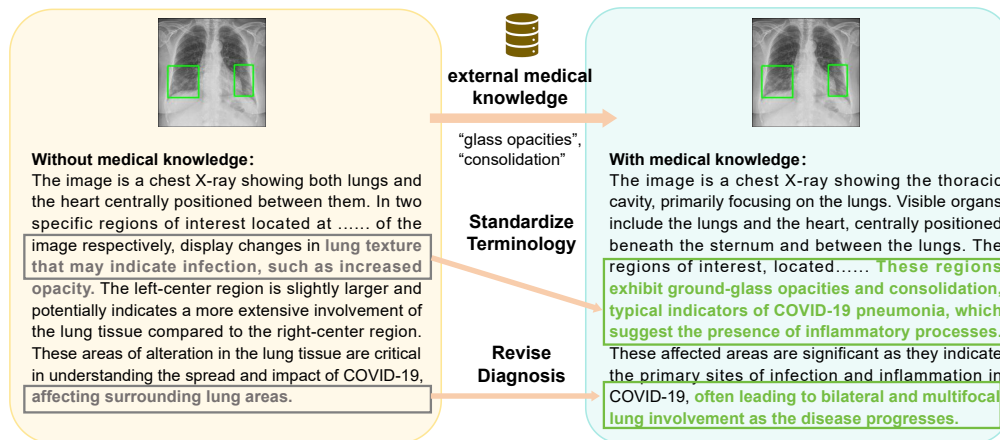


Figure 5: A qualitative comparison example of generated textual description with and without external medical knowledge. MLLMs can standardize medical terminology in its expressions and refine its diagnosis based on disease progressions detailed in medical literature.

162 as the ROIs, while segmentation masks are converted to ROIs by creating the smallest bounding
 163 box that covers the mask. When such localization annotations are not available, we apply different
 164 pretrained expert models listed in the Appendix to generate ROIs. For text-prompt driven grounding
 165 model[29], we use disease and organ information in coarse captions as text prompts to guide the
 166 model in segmenting specific parts. Examples of generated ROIs from various modalities with
 167 different models are demonstrated in Figure 6.

168 Without ROIs, the original description is limited to a brief global analysis of the image. However,
 169 with ROIs, MLLMs can perform a more detailed local analysis of the ROIs and assess the impact of
 170 lesion ROIs on the surrounding normal regions, as demonstrated in Figure 4.

171 **Medical Knowledge Retrieval.** General-purpose MLLMs often produce content that lacks spe-
 172 cialized medical terminology and professional expression. To address this issue, we build a medical
 173 knowledge database following the approach in MedRAG [32]. We collect three main corpora:
 174 PubMed² for biomedical knowledge, StatPearls³ for clinical decision support, and medical text-
 175 books [33] for domain-specific knowledge. We segment these corpora into short snippets and encode

²<https://pubmed.ncbi.nlm.nih.gov/>

³<https://www.statpearls.com/>

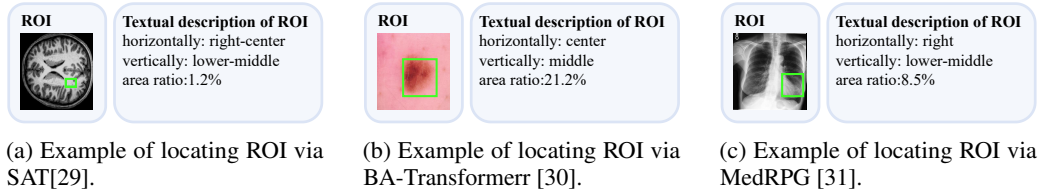


Figure 6: Example of ROIs and their corresponding textual descriptions.

Knowledge 1:
Title: Mobile chest X-ray manifestations of 54 deceased patients with coronavirus disease 2019: Retrospective study.
 Content: We found that 50 (93%) patients with **lesions occurred in the bilateral lung**, 4 (7%) patients occurred in the right lung, 54 (100%) patients were **multifocal involvement**. The number of lung fields involved was 42 (78%) patients in 6 fields, 3 (6%) patients in 5 lung fields, 4 (7%) patients in 4 lung fields, and 5 (9%) patients in 3 lung fields. Fifty-three (98%) patients had **patchy opacities**, 3 (6%) patients had round or **oval solid nodules**, 9 (17%) patients had fibrous stripes, 13 (24%) patients had **pleural effusion**, 8 (15%) patients had **pleural thickening**, 6 (11%) patients had **pneumothorax**, 3 (6%) patients had **subcutaneous emphysema**. Among the 24 patients who had serial mobile chest X-rays, 16 (67%) patients had the progression of the lesions, 8 (33%) patients had no significant change of the lesions, and there was no case of reduction of the lesions. The mobile chest X-ray manifestations of deceased patients with COVID-19 were **mostly bilateral lung, multifocal involvement, and extensive lung field, and pleural effusion, pleural thickening, and pneumothorax probably could be observed**. The serial mobile chest X-ray showed that the chest lesions were progressive with a high probability.

Figure 7: **An example of the Top-8 retrieval results.** By leveraging COVID-19-related medical knowledge, MLLMs can standardize medical terminology and enhance diagnoses according to the disease progressions described in medical literature.

176 them into high-dimensional vectors using the text encoder from Med-CPT [34]. These vectors
 177 are then indexed into a specialized vector knowledge base using Faiss[35], optimized for efficient
 178 retrieval.

179 For a given image, we retrieve relevant medical knowledge by using its coarse caption, which is
 180 generated through metadata integration. Specifically, we encode the coarse captions, including disease
 181 and organ classifications, into vectors using the Med-CPT text encoder. We then perform a vector
 182 similarity search in the medical vector database, retrieving the top eight medical knowledge snippets
 183 that semantically match the query. These snippets provide the external medical knowledge paired
 184 with the image. A qualitative example demonstrating the effectiveness of incorporating external
 185 medical knowledge is shown in Figure 7. With access to COVID-19-related medical knowledge,
 186 MLLMs can standardize medical terminology and refine diagnoses based on the disease progressions
 187 outlined in medical literature.

188 3.2.2 Generation of Multigranular Text Description

189 After data processing, a comprehensive prompt is utilized to guide the MLLMs in generating multi-
 190 granular descriptions. The prompt template consists of a three-level hierarchical framework with
 191 questions to instruct MLLMs: (1) a global description that captures all details of the image; (2) a
 192 local-focused analysis of specific ROIs that potentially are unusual; and (3) a local-global examination
 193 of the interaction between local and global attributes to understand the impact of local abnormalities
 194 on the entire organ. Detailed prompt template is presented in supplementary materials.

195 To ensure that the MLLMs are guided by relevant medical information not inherently present
 196 in their training data, we incorporate the processed data (coarse captions, ROIs, and retrieved
 197 medical knowledge) into the prompts. Specifically, for global information, coarse captions are
 198 directly integrated into the prompt. For local information, ROIs on images are converted into textual
 199 descriptions based on their coordinates and area ratio within the images. Examples of these textual
 200 descriptions are shown in Figure 6, using terms such as "left-center" and "area ratio: 1.2%."

201 To refine terminology and diagnosis within ROIs, relevant medical knowledge about specific diseases
 202 is incorporated into the prompt. Instead of merely inserting this knowledge, we instruct MLLMs to
 203 identify and align the relevant knowledge to ROIs that require analysis.

204 **Choice of MLLMs** We first prompt GPT-4V with the provided medical coarse captions, ROIs,
 205 and medical knowledge to generate a subset of 200,000 samples, maintaining a similar modality
 206 and organ distribution to our full 25 million dataset. The goal of curating this subset is to calibrate
 207 a medical knowledge-guided MLLM to adhere to the formatting instructions specified for our text.

Dataset	Modality	Lesion Type	Lesion BBox/Mask	Color Texture Description	Region Relationship
MedMNIST [39]	✗	✓	✗	✗	✗
DeepLesion [40]	✓	✗	✓	✗	✗
BraTS 2024 [41]	✓	✗	✓	✗	✗
MIMIC-CXR [21]	✓	✓	✓	✓	✗
Quilt-1M [10]	✓	✓	✗	✓	✓
VQA-RAD [42]	✓	✓	✗	✓	✗
CRC100K [43]	✓	✓	✗	✗	✗
SA-Med2D-20M [44]	✓	✓	✓	✗	✗
MedTrinity-25M(Ours)	✓	✓	✓	✓	✓

Table 1: Comparison of dataset types based on provided attributes of annotations.

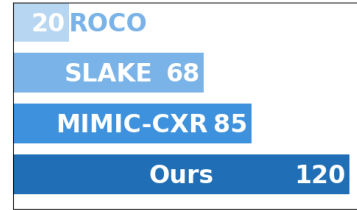


Figure 10: Comparison of the average word count of text descriptions.

219 4 Dataset Analysis

220 **Diversity** Our dataset encompasses a wide range of 10 imaging modalities, with more than 65
221 diseases across various anatomical structures in human. The distribution of Anatomical and biological
222 structures in MedTrinity-25M is shown in Figure 9b. Meanwhile, the number of samples in the dataset
223 for each modality are shown in Figure 9a, spanning from common ones with over 1 million samples
224 each (CT, MRI, X-ray) to rare modalities(ultrasound, dermoscopy) with at least more than 100,000
225 samples, demonstrating a much more balanced distribution compared to other large-scale dataset like
226 SA-Med2D-20M[38], which only contain thousands of ultrasound and dermoscopy samples.

227 **Scale** Figure 9c shows the amount of our dataset, which is significantly larger than previous datasets.
228 To the best of our knowledge, this is the largest open-source, multi-modal multigranular medical
229 dataset to date.

230 **Diseases** The datasets involved in constructing MedTrinity-25M primarily focus on disease diagno-
231 sis and medical discovery. In MedTrinity-25M, diseases are given in the free-form text. The same
232 disease may be referred to using different terms, allowing for elaborate identification and analysis.
233 Figure 9d illustrates the frequently used words related to diseases in our dataset.

234 **Richness** We provide both quantitative analysis and qualitative examples to show the richness
235 of our generated multigranular compare to other medical dataset. Qualitative examples are shown
236 in Figure 1, our textual description is multigranular with more attributes than radiology report of
237 chest x-rays dataset MIMIC-CXR [21], visual QA dataset SLAKE[22] and radiology objects caption
238 dataset ROCO[18]. To demonstrate the multi-granularity of our data, we compared the average word
239 count of text descriptions in our dataset, MedTrinity-25M, with those in other medical datasets, as
240 illustrated in Figure 10. The word count in our dataset is significantly higher, indicating greater
241 richness.

242 **Alignment with human** We leverage GPT-4 to quantify the alignment of generated text descriptions
243 compared to clinical reports from pathologist, which is set as the ground-truth. Specifically, we
244 utilize GPT-4 to score the helpfulness, relevance, accuracy, and level of details of the our generated
245 text descriptions based on clinical reports, and give an overall score on a scale of 1 to 10, where
246 a higher score indicates better overall performance. Additionally, GPT-4 is required to provide a
247 comprehensive explanation for the evaluation score. Detailed experiment results are presented in
248 supplementary materials.

249 5 Conclusion

250 This paper introduces MedTrinity-25M, a large-scale multimodal medical dataset comprising over
251 25 million image-ROI-description triplets sourced from more than 90 online resources, spanning
252 10 modalities and covering over 65 diseases. Unlike existing dataset construction methods that rely
253 on image-text pairs, we have developed the first automated pipeline to scale up multimodal data by
254 generating multigranular visual and textual annotations from unpaired image inputs, leveraging expert
255 grounding models, retrieval-augmented generation techniques, and advanced MLLMs. MedTrinity-
256 25M’s enriched annotations have the potential to support a wide range of multimodal tasks, such as
257 captioning, report generation, classification, and segmentation, as well as facilitate the large-scale
258 pre-training of multimodal medical AI models.

259 References

- 260 [1] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural*
261 *information processing systems*, 36, 2024.
- 262 [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
263 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv*
264 *preprint arXiv:2303.08774*, 2023.
- 265 [3] Tao Tu, Shekoofeh Azizi, Danny Driess, Mike Schaekermann, Mohamed Amin, Pi-Chuan Chang, Andrew
266 Carroll, Charles Lau, Ryutaro Tanno, Ira Ktena, et al. Towards generalist biomedical ai. *NEJM AI*,
267 1(3):AIoa2300138, 2024.
- 268 [4] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu
269 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable
270 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 271 [5] Hong-Yu Zhou, Subathra Adithan, Julián Nicolás Acosta, Eric J Topol, and Pranav Rajpurkar. A generalist
272 learner for multifaceted medical image interpretation. *arXiv preprint arXiv:2405.07988*, 2024.
- 273 [6] Aurelia Bustos, Antonio Pertusa, Jose-Maria Salinas, and Maria De La Iglesia-Vaya. Padchest: A large
274 chest x-ray image dataset with multi-label annotated reports. *Medical image analysis*, 66:101797, 2020.
- 275 [7] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik
276 Marklund, Behzad Haghighi, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest radiograph
277 dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial*
278 *intelligence*, volume 33, pages 590–597, 2019.
- 279 [8] Alistair EW Johnson, Tom J Pollard, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Yifan
280 Peng, Zhiyong Lu, Roger G Mark, Seth J Berkowitz, and Steven Horng. Mimic-cxr-jpg, a large publicly
281 available database of labeled chest radiographs. *arXiv preprint arXiv:1901.07042*, 2019.
- 282 [9] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Nau-
283 mann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision assistant for
284 biomedicine in one day. *Advances in Neural Information Processing Systems*, 36, 2024.
- 285 [10] Wisdom Ikezogwo, Saygin Seyfioglu, Fatemeh Ghezloo, Dylan Geva, Fatwir Sheikh Mohammed, Pa-
286 van Kumar Anand, Ranjay Krishna, and Linda Shapiro. Quilt-1m: One million image-text pairs for
287 histopathology. *Advances in Neural Information Processing Systems*, 36, 2024.
- 288 [11] Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Yash Dalmia, Jure Leskovec, Cyril Zakka,
289 Eduardo Pontes Reis, and Pranav Rajpurkar. Med-flamingo: a multimodal medical few-shot learner. In
290 *Machine Learning for Health (ML4H)*, pages 353–367. PMLR, 2023.
- 291 [12] Tao Tu, Shekoofeh Azizi, Danny Driess, Mike Schaekermann, Mohamed Amin, Pi-Chuan Chang, Andrew
292 Carroll, Charles Lau, Ryutaro Tanno, Ira Ktena, et al. Towards generalist biomedical ai. *NEJM AI*,
293 1(3):AIoa2300138, 2024.
- 294 [13] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe,
295 Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An open-source framework for training
296 large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023.
- 297 [14] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan
298 Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language
299 model. In *International Conference on Machine Learning*, pages 8469–8488. PMLR, 2023.
- 300 [15] Khaled Saab, Tao Tu, Wei-Hung Weng, Ryutaro Tanno, David Stutz, Ellery Wulczyn, Fan Zhang, Tim
301 Strother, Chunjong Park, Elahe Vedadi, et al. Capabilities of gemini models in medicine. *arXiv preprint*
302 *arXiv:2404.18416*, 2024.
- 303 [16] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu
304 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable
305 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 306 [17] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray,
307 Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint*
308 *arXiv:2001.08361*, 2020.

- 309 [18] Obioma Pelka, Sven Koitka, Johannes Rückert, Felix Nensa, and Christoph M Friedrich. Radiology objects
310 in context (roco): a multimodal image dataset. In *Intravascular Imaging and Computer Assisted Stenting
311 and Large-Scale Annotation of Biomedical Data and Expert Label Synthesis: 7th Joint International
312 Workshop, CVII-STENT 2018 and Third International Workshop, LABELS 2018, Held in Conjunction with
313 MICCAI 2018, Granada, Spain, September 16, 2018, Proceedings 3*, pages 180–189. Springer, 2018.
- 314 [19] Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Jiayu Lei, Ya Zhang, Yanfeng Wang, and Weidi Xie.
315 Radgenome-chest ct: A grounded vision-language dataset for chest ct analysis. *arXiv preprint
316 arXiv:2404.16754*, 2024.
- 317 [20] Weixiong Lin, Ziheng Zhao, Xiaoman Zhang, Chaoyi Wu, Ya Zhang, Yanfeng Wang, and Weidi Xie. Pmc-
318 clip: Contrastive language-image pre-training using biomedical documents. In *International Conference
319 on Medical Image Computing and Computer-Assisted Intervention*, pages 525–536. Springer, 2023.
- 320 [21] AlistairEW Johnson, TomJ Pollard, SethJ Berkowitz, NathanielR Greenbaum, MatthewP Lungren, Chih-
321 ying Deng, RogerG Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database of
322 chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019.
- 323 [22] Bo Liu, Li-Ming Zhan, Li Xu, Lin Ma, Yan Yang, and Xiao-Ming Wu. Slake: A semantically-labeled
324 knowledge-enhanced dataset for medical visual question answering. In *2021 IEEE 18th International
325 Symposium on Biomedical Imaging (ISBI)*, pages 1650–1654. IEEE, 2021.
- 326 [23] Ke Yan, Xiaosong Wang, Le Lu, and Ronald M Summers. Deeplesion: Automated deep mining, categor-
327 ization and detection of significant radiology image findings using large-scale clinical lesion annotations.
328 *arXiv preprint arXiv:1710.01766*, 2017.
- 329 [24] axiong/pmc_oa datasets at hugging face. https://huggingface.co/datasets/axiong/pmc_oa.
- 330 [25] Alexandros Karargyris, Renato Umerton, Micah J Sheller, Alejandro Aristizabal, Johnu George, Anna
331 Wuest, Sarthak Pati, Hasan Kassem, Maximilian Zenk, Ujjwal Baid, et al. Federated benchmarking of
332 medical artificial intelligence with medperf. *Nature Machine Intelligence*, 5(7):799–810, 2023.
- 333 [26] Noel CF Codella, David Gutman, M Emre Celebi, Brian Helba, Michael A Marchetti, Stephen W Dusza,
334 Aadi Kallou, Konstantinos Liopyris, Nabin Mishra, Harald Kittler, et al. Skin lesion analysis toward
335 melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi),
336 hosted by the international skin imaging collaboration (isic). In *2018 IEEE 15th international symposium
337 on biomedical imaging (ISBI 2018)*, pages 168–172. IEEE, 2018.
- 338 [27] Veronica Rotemberg, Nicholas Kurtansky, Brigid Betz-Stablein, Liam Caffery, Emmanouil Chousakos,
339 Noel Codella, Marc Combalia, Stephen Dusza, Pascale Guitera, David Gutman, et al. A patient-centric
340 dataset of images and metadata for identifying melanomas using clinical context. *Scientific data*, 8(1):34,
341 2021.
- 342 [28] Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. Pmc-
343 vqa: Visual instruction tuning for medical visual question answering. *arXiv preprint arXiv:2305.10415*,
344 2023.
- 345 [29] Ziheng Zhao, Yao Zhang, Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. One
346 model to rule them all: Towards universal segmentation for medical images with text prompts. *arXiv
347 preprint arXiv:2312.17183*, 2023.
- 348 [30] Jiacheng Wang, Lan Wei, Liansheng Wang, Qichao Zhou, Lei Zhu, and Jing Qin. Boundary-aware trans-
349 formers for skin lesion segmentation. In *Medical Image Computing and Computer Assisted Intervention-
350 MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021,
351 Proceedings, Part I 24*, pages 206–216. Springer, 2021.
- 352 [31] Zhihao Chen, Yang Zhou, Anh Tran, Junting Zhao, Liang Wan, Gideon Su Kai Ooi, Lionel Tim-Ee Cheng,
353 Choon Hua Thng, Xinxing Xu, Yong Liu, et al. Medical phrase grounding with region-phrase context
354 contrastive alignment. In *International Conference on Medical Image Computing and Computer-Assisted
355 Intervention*, pages 371–381. Springer, 2023.
- 356 [32] Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. Benchmarking retrieval-augmented generation
357 for medicine. *arXiv preprint arXiv:2402.13178*, 2024.
- 358 [33] Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease
359 does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied
360 Sciences*, 11(14):6421, 2021.

- 361 [34] Qiao Jin, Won Kim, Qingyu Chen, Donald C Comeau, Lana Yeganova, W John Wilbur, and Zhiyong Lu.
362 Medcpt: Contrastive pre-trained transformers with large-scale pubmed search logs for zero-shot biomedical
363 information retrieval. *Bioinformatics*, 39(11):btad651, 2023.
- 364 [35] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. *IEEE Transac-*
365 *tions on Big Data*, 7(3):535–547, 2019.
- 366 [36] Meta LLaMA Team. Introducing meta llama 3: The most capable openly available llm to date. [https:](https://ai.meta.com/blog/meta-llama-3/)
367 [//ai.meta.com/blog/meta-llama-3/](https://ai.meta.com/blog/meta-llama-3/), 2024.
- 368 [37] Baifeng Shi, Ziyang Wu, Maolin Mao, Xin Wang, and Trevor Darrell. When do we not need larger vision
369 models? *arXiv preprint arXiv:2403.13043*, 2024.
- 370 [38] Jin Ye, Junlong Cheng, Jianpin Chen, Zhongying Deng, Tianbin Li, Haoyu Wang, Yanzhou Su, Ziyang
371 Huang, Jilong Chen, Lei Jiang, et al. Sa-med2d-20m dataset: Segment anything in 2d medical imaging
372 with 20 million masks. *arXiv preprint arXiv:2311.11969*, 2023.
- 373 [39] Jiancheng Yang, Rui Shi, Donglai Wei, Zequan Liu, Lin Zhao, Bilian Ke, Hanspeter Pfister, and Bingbing
374 Ni. Medmnist v2-a large-scale lightweight benchmark for 2d and 3d biomedical image classification.
375 *Scientific Data*, 10(1):41, 2023.
- 376 [40] Ke Yan, Xiaosong Wang, Le Lu, and Ronald M Summers. Deeplesion: Automated deep mining, categor-
377 ization and detection of significant radiology image findings using large-scale clinical lesion annotations.
378 *arXiv preprint arXiv:1710.01766*, 2017.
- 379 [41] Maria Correia de Verdier, Rachit Saluja, Louis Gagnon, Dominic LaBella, Ujjwall Baid, Nourel Hoda
380 Tahon, Martha Foltyn-Dumitru, Jikai Zhang, Maram Alafif, Saif Baig, et al. The 2024 brain tumor segmen-
381 tation (brats) challenge: Glioma segmentation on post-treatment mri. *arXiv preprint arXiv:2405.18368*,
382 2024.
- 383 [42] Jason J Lau, Soumya Gayen, Asma Ben Abacha, and Dina Demner-Fushman. A dataset of clinically
384 generated visual questions and answers about radiology images. *Scientific data*, 5(1):1–10, 2018.
- 385 [43] Jakob Nikolas Kather, Niels Halama, and Alexander Marx. 100,000 histological images of human colorectal
386 cancer and healthy tissue. <https://doi.org/10.5281/zenodo.1214456>.
- 387 [44] Jin Ye, Junlong Cheng, Jianpin Chen, Zhongying Deng, Tianbin Li, Haoyu Wang, Yanzhou Su, Ziyang
388 Huang, Jilong Chen, Lei Jiang, et al. Sa-med2d-20m dataset: Segment anything in 2d medical imaging
389 with 20 million masks. *arXiv preprint arXiv:2311.11969*, 2023.

390 Checklist

391 The checklist follows the references. Please read the checklist guidelines carefully for information on
392 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
393 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
394 the appropriate section of your paper or providing a brief inline description. For example:

- 395 • Did you include the license to the code and datasets? **[Yes]** See Section xxx.
- 396 • Did you include the license to the code and datasets? **[No]** The code and the data are
397 proprietary.
- 398 • Did you include the license to the code and datasets? **[N/A]**

399 Please do not modify the questions and only use the provided macros for your answers. Note that the
400 Checklist section does not count towards the page limit. In your paper, please delete this instructions
401 block and only keep the Checklist section heading above along with the questions/answers below.

402 1. For all authors...

- 403 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
404 contributions and scope? **[Yes]**
- 405 (b) Did you describe the limitations of your work? **[Yes]** See Supplementary materials.
- 406 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]** This
407 research is foundational works, do not include potential negative impacts.
- 408 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
409 them? **[Yes]**

410 2. If you are including theoretical results...

- 411 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]** This paper do
412 not include theoretical results.
- 413 (b) Did you include complete proofs of all theoretical results? **[N/A]** This paper do not
414 include theoretical results.

415 3. If you ran experiments (e.g. for benchmarks)...

- 416 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
417 mental results (either in the supplemental material or as a URL)? **[Yes]** Refer to project
418 page in abstract.
- 419 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
420 were chosen)? **[Yes]**
- 421 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
422 ments multiple times)? **[No]** This paper does not report error bars
- 423 (d) Did you include the total amount of compute and the type of resources used (e.g., type
424 of GPUs, internal cluster, or cloud provider)? **[Yes]** See Supplementary materials.

425 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 426 (a) If your work uses existing assets, did you cite the creators? **[Yes]** We cite all utilized
427 assets in reference.
- 428 (b) Did you mention the license of the assets? **[Yes]**
- 429 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
430 We propose a new dataset, which can be access in our project page.
- 431 (d) Did you discuss whether and how consent was obtained from people whose data you're
432 using/curating? **[Yes]** We follow corresponding licences.
- 433 (e) Did you discuss whether the data you are using/curating contains personally identifiable
434 information or offensive content? **[N/A]** We collect only medical data.

435 5. If you used crowdsourcing or conducted research with human subjects...

- 436 (a) Did you include the full text of instructions given to participants and screenshots, if
437 applicable? [N/A] This paper did not use crowdsourcing.
- 438 (b) Did you describe any potential participant risks, with links to Institutional Review
439 Board (IRB) approvals, if applicable? [N/A]
- 440 (c) Did you include the estimated hourly wage paid to participants and the total amount
441 spent on participant compensation? [N/A]