SUPPLEMENTARY MATERIAL

Anonymous authors

Paper under double-blind review

As described in our main paper, we will provide more details about our dataset CoVT-CRX, annotation tool Labelme-CoVT, prompt generation process, and implementation details in Sec. A, Sec. B, Sec. C, and Sec. D, respectively. Results on task **T1**, **T2**, and **T3** are provided in Sec. E, proving the multi-task ability of our CoVT method on various clinical-related tasks. Finally, we refer the readers to Sec. F for visual examples of our CoVT results on **T5**. Our code, data, demo, and annotation tool are partially available in this anonymous link and will be fully accessible upon acceptance.

A DATASET DETAILS

Besides the visualization in our main paper, we provide more quantitative measures on CoVT-CXR, including word frequency of description, the pre-defined semantic label space C as well as the statistics at each step. Fig. 1 offers a striking word cloud visualization. In part (a), we see findings like cardiac silhouette and hilum contours, depicted through points, lines, and masks. In part (b), common pathologies stand out, with terms such as 'enlargement,' 'low,' 'mild,' and 'crowding' frequently appearing. These subtle gradations are particularly elusive for current vision-language models, making them all the more fascinating and challenging to interpret.



Figure 1: The word cloud illustrating descriptions of pulmonary structures and pathologies.

Category of Annotated Visual Cues. The visual cues in CoVT-CXR annotations encompass 112 categories, covering five annotation types: bounding box, mask, point, line, and linestrip. The details of these categories can be found at here. We further demonstrate the number of mask annotations in Fig. 2 where the original 112 categories are grouped into 21 classes for visualization purposes.



Figure 2: Number of mask annotations in \mathcal{T} based on 21 high-level categories.

072 **Rationale Step Statistics.** As described in our main paper, each diagnostic report \mathcal{R} can be divided 073 into multiple sentence segments S. By definition, each S corresponds to one field, referred to as the target field. Our CoVT-CXR dataset explicitly analyzes a single chain-of-thought reasoning process 074 in one target field. As can be found in Fig. 3, we have at most 12 steps of reasoning for each segment 075 S, indicating the complexity of the reasoning process. The projection tag, which describes the type of 076 imaging views, requires only one step of reasoning. This is reasonable because frontal imaging with 077 AP and PA views, as well as lateral imaging, can be directly defined by the given CXR images. The 078 bone field descriptions exhibit a relatively uniform distribution of step lengths ranging from 1 to 11, 079 likely due to the variability in the types and quantities of bones being examined. Reports related to the lung field show either short or very long step lengths, possibly due to the presence of simple (single 081 nodule) or complex (multiple abnormalities) diagnoses. Cardiac-related reports tend to cluster around 082 intermediate step lengths of 7-8, reflecting the predetermined heart measurement protocols, which typically require 6 steps to complete the cardiothoracic ratio calculation. This analysis highlights 084 the nuanced relationship between diagnostic complexity and reasoning processes, underscoring the 085 importance of tailored approaches in medical image interpretation. We hope that this observation can guide the development of more efficient and accurate diagnostic tools in the future.





105

087

090

092

093

094

095

096

097 098

099 100

070

071



Ribs Distribution Notation. Some may notice a chunk of ribs or vertebrae presented in the sunburst figure in the main paper. However, reports of bone lesions are relatively scarce in the MIMIC-CXR (Johnson et al., 2019). Although rib lesions account for only a small portion of the

108 MIMIC-CXR dataset, ribs play a crucial role in intermediate reasoning steps, thereby are significant in our metadata. For instance, ribs help localize lesions and indicate signs of pulmonary collapse, 110 hyperinflation, or atelectasis, as they are consistently present and orderly arranged within the lungs. 111 These abnormalities, frequently reported by MIMIC-CXR, require precise identification of ribs. 112 Statistically, we analyzed the frequency and percentage of occurrences of these pathologies in 120k findings extracted from the official dataset. Results in Tab. 1 indicate that rib involvement in related 113 lesion descriptions spans a wide range, with both direct and indirect participation in diagnoses 114 exceeding 40%. 115

Class	Ribs	Collapse	Hyperinflation	Atelectasis	Involvement of ribs
Count	10133	15456	4611	29485	48117
Percentage	8.44%	12.88%	3.84%	24.57%	40.10%

Table 1: The count and proportion of diagnoses involving the ribs in the raw dataset.

В ANNOTATION PHASE

Findings

"Portable semi-upright radiograph of the chest demonstrate low lung volumes results in bronchovascular crowding. The cardiomediastinal contours are unchanged. There is no pneumothorax, pleural effusion, or consolidation. A right subclavian central venous line terminates in the cavoatrial junction. Nasogastric tube courses into the stomach and out of the field of view."



142

121

122 123

124 125

126 127

128

143 144

145

146

147

148

149

150

151

152

156

157

159

161

Figure 4: Illustration of five tasks derived from CoVT-CXR, which align our five steps well.

Obtaining Metadata. In Fig. 4, we provide an example that indicate how medical trainees collaborate to obtain metadata \mathcal{T} from $\{I, \mathcal{R}\}$. Our \mathcal{R} is represented as "Findings," with CXR images shown in the bottom-right of Fig. 4. Specifically, \mathcal{R} is broken down into multiple S segments. We further emphasize our semantic disentanglement by bolding sentences where a single sentence in \mathcal{R} can be divided into two S segments. For each sentence S, medical trainees are required to provide four different \mathcal{T} corresponding to various stages derived from the ultimate task T5. For example, step 1 is the segmentation task where annotators provide visual cues v, or masks in this case, for the class c, such as 4-1 and 4-2 for the left and right lungs, along with the description d, "Both lungs are detected".

153 **Annotation Tool Features.** To collect \mathcal{T} from I, \mathcal{R} , we adapted Labelme (Wada) to develop our 154 annotation tool, Labelme-CoVT, incorporating the following key features: 155

- Textual Component Integration: This establishes connections between the CXR image I and the report \mathcal{R} .
- Flag Module Modification: This serves as a reference for reported sentences S, allowing for the disentanglement or merging of sentences.
- SAM Integration for CXR: The SAM-CXR model, trained for T1, assists in expediting the labeling process. This feature allows users to interactively create masks through simple

clicks or points. Additionally, SAM for CXR within the annotation tool can accept labelers' modifications as human feedback for further incremental learning.

163 164

162

165 166

More annotation features can be seen in Labelme-CoVT demo.

Annotation Simplify and Efficiency Improvement. Given the intricate rationale workflow, the Labelme-CoVT tool plays a pivotal role in curating the CoVT-CXR dataset. It not only facilitates the visualization of sequential interleaved cross-modal data for simplification but also accelerates the labor-intensive manual annotation process by streamlining verification and improving efficiency. We randomly selected four groups from a pool of 32 medical trainees, who were required to complete a questionnaire assessing time consumption. Additionally, we measured inter-annotator agreement among the trainees to validate the completed cases. Tab.2 summarizes the efficiency and accuracy outcomes of the annotation phase.

175			
	1	7	5

177

179

181 182 183

185 186

187

188

189

196 197

207 208

Group	#1	#2	#3	#4
Avg. Time (w/o tool)	43	40	40	45
Avg. Time	14	8	9	12
Speed up	3.1x	5x	4.4 x	3.8x
Cross Person Agreement (w/o tool)	62.12%	45.12%	50.44%	75.41%
Cross Person Agreement	85.12%	91.23%	87.33%	90.53%
Agreement Improvement	+23.00%	+46.11%	+36.89%	+15.12%

Table 2: Evaluation of Labelme-CoVT in annotation simplification and efficiency improvement.

The time consumption, measured in minutes, shows that the tool significantly increases annotation efficiency, reducing the time required by three to five times. Furthermore, it enhances annotation accuracy by approximately 1.5 times, with data usability surpassing 80%.

190 Data Assignment for Multi-stage Training.

We follow a progressive training scheme to gradually improve the model's multi-step reasoning capabilities. For instance, we incorporate single-step reasoning tasks like segmentation or image description during pretraining, and further fine-tune the model with chain-of-visual-thought reasoning using long contextual data. As a result, our multi-stage training setup and data volumes are arranged as shown in Tab.3.

Stage	Task	Format	Volume
	T1	$\langle I, p, v' \rangle \to v$	0.4M
Pretrain	T2	$\langle I, p, v \rangle \to t$	2.4M
	T3	$\langle I,q\rangle \to \langle v_1,\ldots,v_i\rangle$	0.2M
Finatuna	T4	$\langle I, q' \rangle \rightarrow \langle (v_1, t_1), \dots, (v_i, t_i), S \rangle$	30k
Filletulle	Т5	$I \rightarrow \{ \langle \{v^m_{i_j}, t^m_{i_j}\}_{i_j=1}^{i_m}, S^m \rangle \}_{m=1}^M$	6k

Table 3: The count and proportion of diagnoses involving the ribs in the raw dataset.

Notably, there are more entries for T2 than the others. This is due to our design, where versatile pand v of specific image I are randomly combined to enhance the robustness of text generation.

Separating the SAM-CXR and CoVT Pretraining Data. The T1 task data is used both to assist
 the build SAM-CXR and CoVT pretrain after the annotation process. Initially, one-third of the data is
 annotated manually, and this data is then used to train the SAM-CXR for CoVT-CXR data curation.
 The trained SAM-CXR is subsequently employed to assist in annotating the remaining two-thirds of
 the data. Once the annotation process is complete, the fully-annotated dataset is then used to update
 all relevant models, including SAM-CXR and CoVT.

C GPT PROMPT FOR CONTEXTUAL REFINEMENT & GENERATION

2	21	7
2	21	8
2	21	9

220 221 222

247

216

	intermediate_step_tmp[f*rationale of step {step_id + 1}"] = \langle
	{"type": "string", "description": f"the rationale description of step {step_id + 1}"}
inter	mediate_step_des = str()
for s	intermediate step des \pm f'step (step id \pm 1): (des) \n"
# De	fine the tools
Fun	ction Name: Intermediate diagnostic rationale generation
Fun	ction Description:
	"There is one or several highlighted image-text pair showing intermediate steps with a final diagnostic report for interpreting a chest X-ray. Based on this report, please provide the finding clues and rationale of diagnosing this report for every mentioned step in this case. Do not mention the given specific step number and provided raw description, and retain the numerical content appearing in each step. The current rationale could be able to reference the previous images and generated rationales, and should gradually lead towards the final report. The given step input and the generated rationale of step must be one-to-one "
Par	ameters:{
	"type": "object",
	"properties": intermediate_step_tmp,
	"required": list(intermediate_step_tmp.keys())
}	
	fine the messages
# D0	
# Do mes	sages =
# Do mes	sages = [{"role": "system", "content": f"This is auxiliary information: Request ID: {request_id}"},
# Do mes	<pre>sages = [{"role": "system", "content": f"This is auxiliary information: Request ID: {request_id}"}, {"role": "user", "content":</pre>
# Do mes	<pre>sages = [{"role": "system", "content": f"This is auxiliary information: Request ID: {request_id}"}, {"role": "user", "content": The final report description states that these findings suggest that "{final_reported_sentence}".</pre>

Table 4: Rationale refinement for d.

248 **Rationale Refinement.** Unlike machines, human descriptions \mathcal{R} can vary in levels of abstraction and 249 may omit background information, considering the conversational context during communication. 250 While this approach is efficient among doctors, it is machine-unfriendly, particularly when deep reasoning is required based on the provided descriptions. To this end, we propose incorporating 251 external knowledge, such as LLMs, to augment chain-of-thought reasoning descriptions. In practice, we use the GPT-4 (OpenAI et al., 2024) API to refine descriptive content, as illustrated in Tab. 4. 253 Specifically, we first convert the class c to a textual description with the phrase, "the highlighted area 254 showing the <c>". We then input this converted description, the raw textual descriptions provided by annotators at the current reasoning step, and the corresponding sentence segment S into the 256 LLMs. The LLMs can understand the causal relationship between the annotated content at the current 257 diagnostic step and the formation of the given final diagnostic sentence. For example, when we 258 provide the annotation of five ribs along with a statement indicating reduced lung capacity, GPT-4 259 may generate a rationale like, "Only five ribs are visible; normally, 9-10 ribs can be seen with fully 260 expanded lungs." This sentence becomes the final textual description d at the current reasoning step. 261 This design helps in forming a more reliable chain-of-thought process, leading to a more confident diagnostic reasoning process by leveraging the strengths of LLMs to support medical professionals in 262 making informed decisions. 263

Instruction Generation. As stated in our main paper, **T4** requires instructions q', or equivalently questions, as input, and GPT-4 is utilized to generate these instructions. However, simply applying GPT-4 on whatever sentence S would lead to overly specific or vague questions as sentences from pathological or normal reports vary a lot. This discrepancy can mislead the CoVT model by causing it to focus on the questions rather than considering the CXR images. For instance, sentences from a pathological report typically offer detailed descriptions of a lesion, leading to highly specific questions about that lesion. Conversely, sentences from a normal report are often oversimplified, causing GPT-4



Table 5: Generated question instruction q' for visual question answering.

to generate vague questions more frequently. To avoid this, GPT is asked to identify the relevant field for each report sentence first. Then it formulates questions based on this field, as demonstrated in Tab. 5. During inference, our CoVT takes $\{I, q'\}$ as input, generates a chain-of-thought process $\langle (v_1, t_1), \ldots, (v_i, t_i) \rangle$, and finally provides sentence S.

D IMPLEMENTATION DETAILS

301 302

303

304

305

306 307 308

309

310 Implementation of SAM for CXR (T1). We apply the SAM (Kirillov et al., 2023) model for 311 CXR segmentation, choosing the smallest configuration with 12 vision encoder layers and 768 312 visual embedding channels, along with the default mask decoder setup. To adapt SAM for CXR 313 segmentation, we modify the prompt encoder by removing the bounding box prompt due to the 314 irregular pulmonary structures and using a point sample strategy within 5 iterations. Additionally, 315 we introduce 76 category prompts to compensate for the absence of a text prompt encoder, as SAM reportedly uses CLIP (Radford et al., 2021). The model is trained in mixed precision with the AdamW 316 optimizer (Loshchilov & Hutter, 2019), using a learning rate of 1e-3 and a weight decay of 0.1. 317 Gradient accumulation is employed to increase the batch size to 128. The training process of SAM 318 for CXR takes more than 48 hours with 12 RTX 4090 GPUs. 319

Implementation of VQGAN for Tokenization. We set the autoencoder (Esser et al., 2021) to have a discrete codebook size of 8192 with a latent dimension of 256. Given a resized image of resolution 336, the image encoder reduces the resolution to 21, resulting in an image or mask taking a contextual length of $21^2 = 441$ in the subsequent foundation model. The model employs L1 reconstruction loss and LPIPS perceptual loss (Zhang et al., 2018). We train the tokenizer on 1.4M samples of

various masks and image combinations from the CoVT-CXR dataset for 250k iterations. Generative adversarial loss is enabled at the 10,000*th* step when the reconstructed image achieves a certain level of quality. We set the batch size to 48 and it takes us 30 hours on 12 RTX 4090 GPUs to complete the tokenizer training process.

Implementation of Foundation Model for CoVT. We adopt Llama3-Instruct-7B (Meta, 2024) as our base model and redefine the model's embedding layer to map visual tokens from VQGAN and textual tokens into a unified embedding space, encompassing representations for both modalities. Subsequently, we perform full-parameter pretraining on a dataset of 1.3 million extracted samples. During the pretraining phase, we use 8 H100 GPUs with a batch size of 128 and a learning rate of 1e-4, taking 24 hours to complete one epoch. In the SFT phase, we fine-tune the model on the CoVT-CXR dataset with a batch size of 64 and a learning rate of 4e-5, training for 4 epochs in 3 hours.

E MORE DETAILS ABOUT EXPERIMENTAL RESULTS

T1: SAM-CXR. As a reminder, **T1** is a **image-to-image** task aiming to mimic the recognition process of doctors, which is fulfilled with SAM framework. In experiments, we take Meta SAM (Kirillov et al., 2023) and MedSAM (Ma et al., 2024) as our baselines; the former performs strong zero-shot capabilities on various natural images, while the latter has undergone extensive fine-tuning for 2D medical images. Given the significant variability in image distributions of different structures in CXR, we choose five of them, namely the trachea, lungs, heart, bones, and mediastinum, to represent the majority of visible structures in CXR. We exploit point-wise prompts for SAM and SAM-CXR for its simplicity and effectiveness.

Model	Trachea	Heart	Lungs	Ribs	Clavicles	Mediastinum	Avg.
SAM(Kirillov et al., 2023)	3.89	23.66	78.79	5.23	9.08	24.91	24.26
MedSAM(Ma et al., 2024)	49.42	75.01	89.70	15.82	10.87	70.23	51.85
CXR-SAM	74.11	86.88	98.41	84.70	91.70	88.16	90.76

Table 6: The comparison of Meta SAM, medical image fine-tuned MedSAM, and SAM-CXR in CXR applications. The evaluation metric used is Intersection over Union (IoU), with 'Avg.' referring to the average IoU score when considering all reported structures

In contrast, MedSAM is trained exclusively with a bounding box (bbox) prompt encoder thus bboxwise prompts are used in our experiment. Results in Tab. 6 show that the general SAM model performs reasonably well in lung segmentation, likely due to the distinct segmentation boundaries of the lungs. However, it fails in segmenting the heart, mediastinum, and bones. We attribute this to the overlapping structures in X-ray imaging, which obscure the target image features. In addition, it is worth noting that the MedSAM fine-tuned on medical images seems to perform better than SAM. This might be because the bbox-wise prompts, which are used by MedSAM, provide more precise location information compared to point-wise ones.

Again, the low accuracy in segmenting the most challenging bones suggests that MedSAM also fails on CXR segmentation. Our approach successfully addresses the shortcomings of the SAM model in CXR segmentation, providing a robust visual foundation model for downstream CXR tasks.

 Metric BLEU-1	BLEU-2	BLEU-3	BLEU-4	ROUGE-1	ROUGE-L	METEOR	CIDEr
Value 0.3877	0.2404	0.1563	0.1020	0.5154	0.3328	0.3221	0.3529

Table 7: Benchmark performance for the T2 task

T2: Imaging Caption. The design of T2 is inspired by the process where doctors proceed to draw conclusions or summarize results in T1, thereby is a image-to-text task.

Similar to (Huang et al., 2023), our CoVT focuses on summarizing a prompted region rather than
 describing the entire CXR image, which turns out to be beneficial for understanding finer details in CXR images.



Table 8: Exsamples for visual question answering. The model highlights specific regions based on the question and further generates region-specific relevant descriptions based on the content of the highlighted areas.

Fig. 8 and Tab. 7 show the qualitative and quantitative results of **T2** on CoVT-CXR respectively. Furthermore, thanks to SAM-CXR, we can interactively generate captions for specific regions in CXR images using points, masks, or text prompts. Please check the demo for the interactive caption generation process.

>Locate the carina of the bronchus	>Detect the bottom contours of the lateral lung cavity
>Locate the contours of the left and right heart	>Detect the blunted costophrenic angle in the lateral view
➤Setting board line of cardiac contours	➢Frontal view right rib blunted contour
Measuring the width of heart/cardia	≻Frontal view left rib blunted contour
>Detect abnormal heart contours	≻Find the tip of the ET tube
>Locate the board line of the outer contours of both lungs	\succ Measure the distance between the tip of ET tube
Measure the width of both lungs	and the carina
>Detect the contour of the aortic arch	≻Check if the NG tube is in the stomach
>Detect mediastinal contours	Detect the pacemaker with electrode
>Find abnormal protruding mediastinal contours	≻Find the tip of the pacemaker electrode
≻Trace the vertebral line of the thoracic spine	> Find the end of the venous tube
≻Detect the hilum contours	≻Locate the tip of the chest tube or drainage tube

Table 9: Predefined instruction for building visual reasoning chain.

T3: Instruction-following Visual Thought. Apart from the cross-modal chain of visual thought, we recognize that single visual reasoning is also beneficial for accomplishing clinical-related tasks, such as coarse-to-fine detection, marking subtle points, and performing complex medical measurements. Tab. 9 displays our predefined instructions used to guide single-step visual reasoning. The specific meanings of the instructions can be found in here, where class id mentioned are following categories of A. It is worth noting that predefined instructions can be combined to form a series of visual reasoning steps to complete more complex clinical-related tasks. More details for visual reasoning guided by instructions can be seen in visual reasoning pipeline.

Why is the multi-step visual thought design helpful? To explicitly showcase the impact of multi-step reasoning, we provide case-wise explorations on the average sampling probability on the ground truth (GT) for final findings across varying numbers of visual cues. As evidenced in Fig. 5, incorporating additional intermediate visual cues and translating them into meaningful sentences significantly increases the average sampling probability of GT tokens, thereby improving the likelihood of generating accurate GT reports.



Figure 5: Ground truth sampling gains in token-level as intermediate steps increases.

F VISUALIZATION OF COVT RESULTS

 Generalization for Unseen Data. CXR images, unlike CT or MRI, exhibit a more consistent distribution, distinguishing them from other medical imaging. Consequently, SAM-CXR and CoVT have shown strong generalization capabilities on unseen data, as illustrated in Fig.6. SAM-CXR excels in segmenting bones and airways via prompt-based segmentation, which are typically challenging to discern by Meta SAM (Kirillov et al., 2023) and medical SAM (Ma et al., 2024) due to overlapping shadows in radiographs. CoVT also performs well on datasets like Vindr (Nguyen et al., 2022), being less affected by distribution shifts. This is due to the precise visual prompts in Chain-of-Visual-Thought, which reduce noise by focusing on highlighted bone pathology regions, avoiding interference from surrounding tissues. Additionally, CoVT aids in identifying CXR views and lesions, enabling interpretable reasoning process through single or multi-step diagnostic tasks.



Figure 6: Generalization test visualization from several unseen data.

486

487

488

489

491 492 493

494

495

496

497

498

499

500

501

502

504 505

506

507 508

509

510 511

512

513

514 515

516

517 518

519

520

521

522 523

524 525

534 535

538 539

Chain of Visual Thought Visualization. We provide some examples of our CoVT results on CoVT-CXR on task T4 and T5 in Fig. 7 and Fig. 8 respectively. Clearly, our CoVT generates more accurate answers and reports compared to existing methods, demonstrating its superiority over the baselines. More importantly, CoVT provides interpretable cross-modal intermediate reasoning steps, 490 enhancing the reliability and interpretability of state-of-the-art methods.

Q: What is the projection view of the chest? GT: AP view of the chest Gemini-1.5-Pro: The projection of this chest X-ray is anteroposterior (AP). This means the X-ray beam passed from the front (anterior) of the patient to the back (posterior). You can tell this is an AP view because Scapulae are not cleared: The scapulae (shoulder blades) are superimposed over the lung fields, which is typical of an AP projection where the patient can't always pull their shoulders forward effectively. GPT-4V: The projection view of the chest X-ray provided is an anteroposterior (AP) view, typically taken in the upright position, as indicated by the labels "PORTABLE" and "UPRIGHT ". CoVT: The bounding box of the word "portble" in image indicates an AP view of the chest. <findings>Single AP view of the chest</findings>. Q: What are the findings of the mediastinal silhouettes? GT: Calcification of the aortic knob is seen. Gemini-1.5-Pro: The mediastinal contours are within normal limits. There is no mediastinal widening. GPT-4V: The mediastinal contours appear normal with no signs of widening. CoVT: There is a The mediastinal The aortic arch structures appear calcification is not prominent. to be normal. of the aortic <findings>calcification of the aortic knob is seen.</findings> O: What is the condition of the lungs? Gemini-1.5-Pro: The lungs are clear bilaterally. There is no evidence of consolidation, pleural effusion, or pneumothorax. GPT-4V: Based on the provided chest X-ray image, the lungs appear clear without signs of consolidation, effusion, or pneumothorax. There are no obvious abnormalities such as masses or lesions observed within the lung fields. The costophrenic angles are sharp, indicating no pleural effusion. Overall, the lungs seem to be in a normal condition. the right and left lung areas are examined on the frontal view. both lungs exhibit CoVT: clear fields without any signs of focal consolidation. additionally, there are no visible signs of pleural effusion or pneumothorax, which supports the final report <findings>There is no focal consolidation, pleural effusion, or pneumothorax identified.</findings>

Figure 7: One example of our CoVT for VQA task based on T4



Figure 8: One example of our CoVT for report generation task based on T5

594 REFERENCES

601

- Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image
 synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pp. 12873–12883. Computer Vision Foundation / IEEE, 2021.
- Xiaoke Huang, Jianfeng Wang, Yansong Tang, Zheng Zhang, Han Hu, Jiwen Lu, Lijuan Wang, and
 Zicheng Liu. Segment and caption anything. *CoRR*, abs/2312.00869, 2023.
- Alistair E. W. Johnson, Tom J. Pollard, Seth J. Berkowitz, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih-ying Deng, Roger G. Mark, and Steven Horng. MIMIC-CXR: A large publicly available database of labeled chest radiographs. *CoRR*, abs/1901.07042, 2019.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloé Rolland, Laura Gustafson, Tete
 Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross B. Girshick.
 Segment anything. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pp. 3992–4003. IEEE, 2023.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019.
 OpenReview.net, 2019.
- Jun Ma, Yuting He, Feifei Li, Lin Han, Chenyu You, and Bo Wang. Segment anything in medical images. *Nature Communications*, 15(1):654, 2024.
- 615 AI Meta. Meta llama 3, 2024.
- Ngoc Huy Nguyen, Hieu H. Pham, Thanh T. Tran, Tuan N. M. Nguyen, and Ha Q. Nguyen. Vindr pcxr: An open, large-scale chest radiograph dataset for interpretation of common thoracic diseases
 in children. *CoRR*, abs/2203.10612, 2022.
- 620 OpenAI, Josh Achiam, Steven Adler, et al. Gpt-4 technical report, 2024.
- 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.
 Learning transferable visual models from natural language supervision. In Marina Meila and Tong
 Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp.
 8748–8763. PMLR, 2021.
- Kentaro Wada. Labelme: Image Polygonal Annotation with Python. URL https://github.
 com/wkentaro/labelme.
- Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018, pp. 586–595. Computer Vision Foundation / IEEE Computer Society, 2018.
- 635

630

- 636 637
- 638
- 639

- 642
- 643
- 644
- 645
- 646
- 647