
Adapting Pretrained Vision-Language Foundational Models to Medical Imaging Domains

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Multi-modal foundational models are trained on millions of pairs of natural images
2 and texts, frequently obtained through web-crawling approaches. Although their
3 performance is excellent, these models do not generalize well to other domains,
4 such as medical imaging, especially when these domains do not resemble the
5 centric-like images that can be found on the web. In this study, we assess the ability
6 of the stable diffusion model to generate domain-specific images in the particular
7 case of medical imaging. Based on quantitative and qualitative evaluations of the
8 main components of the stable diffusion pipeline (the variational autoencoder, the
9 U-Net and the text-encoder), we explore several approaches to fine-tune stable
10 diffusion to generate radiological images, which accurately represent the clinical
11 content of conditional text prompts. Our best-performing model improves upon the
12 stable diffusion baseline and can be correctly conditioned to insert an abnormality
13 on a synthetic radiology image.

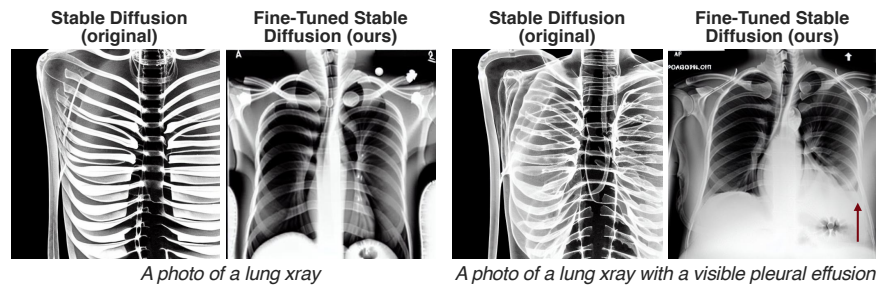


Figure 1: Generated images by both the original stable diffusion model and our fine-tuned model on radiology images. The prompts are designed to compare a standard radiology image with no particular findings, and the insertion of a frequent abnormality "pleural effusion" (red arrows).

14 1 Introduction

15 In recent months, latent diffusion models have gained immense popularity by enabling state-of-the-art
16 image generation amenable to fine-grained control of the image generation process at inference time
17 via conditioning of the denoising process (e.g., using text prompts) (Ramesh et al., 2022; Rombach
18 et al., 2022; Saharia et al., 2022). Such models, termed foundation models (Bommasani, 2021), have
19 been trained with large multi-modal curated datasets such as LAION-5B that consists of natural
20 images and their captions (Schuhmann et al., 2022). The impressive generative capabilities of such
21 models permits creation of high-fidelity synthetic datasets that may be used to augment traditional
22 supervised machine learning pipelines in scenarios that lack training datasets.

23 One particular area that such an advance would be beneficial in is the domain of medical imaging,
24 where there is a paucity of high-quality labeled datasets. Annotating such medical imaging datasets
25 typically requires trained medical experts who are capable in interpreting subtle, but semantically
26 meaningful, image features. Despite the lack of large curated medical imaging datasets, one benefit
27 that such medical imaging examinations have is that there is typically a text-based radiology report that
28 describes pertinent findings from the imaging study. Leveraging the vision-language understanding
29 capabilities of latent diffusion models could potentially provide an intuitive mechanism to create
30 synthetic medical imaging data by prompting with relevant medical keywords or concepts of interest.

31 In this study, we explore the representational bounds of large vision-language foundation models
32 and evaluate how to utilize pretrained foundational models to represent medical imaging studies and
33 concepts, despite models never having been explicitly trained on these concepts. We utilize chest
34 x-rays (CXRs) for this study as they are most common imaging modality globally. CXRs are fast to
35 acquire, inexpensive, can provide important patient health insights, and can identify and monitor a
36 variety of pathologies. We explore and quantify the representational capacity of the stable diffusion
37 model (Rombach et al., 2022) to characterize the efficacy of both its language and vision encoders as
38 applied to CXRs. We further explore different strategies for improving the representational capacity
39 of non-domain-specific foundational models for representing medical concepts specific to CXRs.
40 These experiments help provide novel decision making insights regarding whether such foundational
41 models can accurately represent complex biomedical concepts for clinically-relevant downstream
42 tasks, without explicit training on such concepts. In this study, we specifically show the following:

- 43 1. Training Stable Diffusion on LAION-2B learns a variational autoencoder (VAE) that can
44 reconstruct CXR images out-of-the-box
- 45 2. A frozen CLIP text encoder can generate powerful medical embeddings with enough clinical
46 context to allow accurate generated images, in conjunction with the methods below
- 47 3. Replacing the frozen CLIP encoder with a frozen in-domain text encoder with a projection
48 head trained on LAION to map in-domain embeddings to CLIP embeddings, is not adequate
49 to generate better images
- 50 4. Textual inversion can be used to learn complex medical concepts like pleural effusion in a
51 few-shot manner
- 52 5. Fine-tuning the UNet component enables high-fidelity CXR image generation with the
53 capability to insert custom pathologies (see examples in Figure 1).

54 We verify all our findings using using quantitative metrics of image quality as well as qualitative and
55 domain-specific radiological interpretation from an expert thoracic radiologist.

56 **2 Materials and Methods**

57 **2.1 Datasets**

58 A large, publicly available chest x-ray (CXR) dataset (MIMIC-CXR, version 2.0.0) was used in this
59 work, under institutional review board approval (Johnson et al., 2019). MIMIC-CXR contains a total
60 of 377,110 images from studies performed at the Beth Israel Deaconess Medical Center in Boston,
61 MA, USA, of which 700 frontal (i.e., anterior-posterior or posterior-anterior projection) radiographs
62 were sampled randomly for this study. These images and their associated reports were used for
63 experiments and study of the variational autoencoder and of text encoders.

64 In addition, we manually select 5 images with no findings, as well as 5 images that have visible pleural
65 effusion, discarding any improperly cropped or colorized images (verified by a thoracic radiologist).
66 Along a set of simple prompts generated synthetically, these form pairs of images and texts that are
67 used for fine-tuning the stable diffusion model with various approaches. Finally, a sample of one
68 million text prompts from the LAION-400M dataset (Schuhmann et al., 2021) is used for textual
69 projection training and experiments.

70 **2.2 Stable Diffusion**

71 The stable diffusion model (depicted in Figure 2) is composed of a CLIP text encoder that parses text
72 prompts to create a 768-dimensional latent representation (Radford et al., 2021a). This latent text

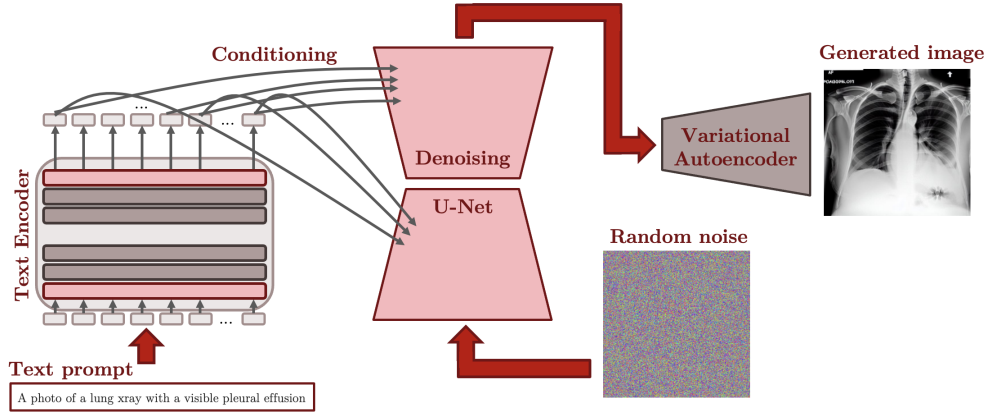


Figure 2: Stable diffusion architecture, run in the radiology setting to generate synthetic radiology images.

73 representation is used to condition a generative U-Net to generate images in the latent image space
 74 using random noise as an additional conditioning. Finally, the decoder component of a variational
 75 autoencoder is used to map this latent image projection to the output image space. While the original
 76 generative model has been trained with image and text captions arising from natural imaging domains,
 77 the extent of its capabilities for representing medical concepts and images remains unclear. To
 78 adapt the stable diffusion model for in-domain image generation, especially for radiology images
 79 and prompts, we can leverage each component and train it, or not, depending on its capabilities to
 80 represent in-domain data. More particularly, we can assess:

- 81 • Whether the variational autoencoder (VAE) alone is capable of reconstructing radiology
 82 images without losing general visual aspect as well as clinically important features.
- 83 • Whether the text encoder alone is capable of projecting clinical prompts to the text latent
 84 space while preserving clinically important features.

85 Section 2.3 presents the methods used to assess the reconstruction quality of the VAE, assessing
 86 whether it requires in-domain fine-tuning; Section 2.4 describes the experiments researching the
 87 quality of the CLIP text encoder and other in-domain text encoders; and Sections 2.5, 2.5, 2.7 present
 88 methods to fine-tune various components of the stable diffusion model for the radiology domain.

89 2.3 Variational Autoencoder

90 As latent diffusion model, stable diffusion translates image inputs into a latent space before performing
 91 the denoising process, using an encoder trained to remove perceptually negligible features (“perceptual
 92 compression”)(Rombach et al., 2022). To analyze how well medical imaging information is preserved
 93 while passing through the VAE, CXR images sampled from MIMIC (“originals”) were encoded to
 94 latent representations and reconstructed into images (“reconstructions”).

95 Reconstruction quality was quantitatively assessed by calculating the root mean square error (RMSE),
 96 the peak signal-to-noise ratio (PSNR) and the structural similarity index measure (SSIM) for each
 97 image-reconstruction pair. Additionally, the Fréchet inception distance (FID, underlying model:
 98 Inception V3, 2048 features) was calculated on minibatches (batch size = 32) to compare the
 99 distribution of reconstructions to the distribution of original images(Szegedy et al., 2015; Heusel
 100 et al., 2017).

101 Qualitatively, the reconstruction quality compared to the original image input was assessed by a
 102 radiologist with 7 years of experience in reading CXR studies, using a scoring system ranging
 103 from 1 to 5 (5: Very good reconstruction with essentially non-inferior diagnostic quality to the
 104 original, 4: Good reconstruction with noticeable errors not negatively influencing diagnostic quality,
 105 3: Moderate reconstruction errors with possible negative effects to diagnostic performance, 2: Severe
 106 reconstruction errors or errors of any level leading to hallucinated lesions, 1: Severe reconstruction
 107 errors yielding the image undiagnostic) on 100 randomly sampled original-reconstruction pairs.

108 The effect of the reconstruction process on classification performance was analyzed using a model pre-
 109 trained to detect 18 different pathologies commonly encountered in CXR (DenseNet-121, torchxrayvi-
 110 sion library, version 0.0.37)(Cohen et al., 2022)(Cohen et al., 2020). Classification accuracy and F1
 111 score were calculated for 12 of the labels included in both MIMIC-CXR and the pretrained model.
 112 For this step, uncertain findings (=‘-1’) were considered positive findings, while missing values were
 113 treated as absence of the corresponding finding. Additionally, latent representations of original and
 114 reconstructed images were compared by calculating their pairwise cosine similarity.

115 2.4 Text Encoder

116 In the domain-specific setting of radiology reports and images, the goal is to be able to condition the
 117 generation of images on associated medical conditions, that can be represented through a text prompt
 118 or report. Therefore, the capability of the text encoder to correctly represent medical features in the
 119 latent space is critical for the rest of the stable diffusion process, in particular the U-Net operating
 120 in the latent space, to be able to generate images that are anatomically correct and representing the
 121 correct set of abnormalities.

122 A set of potential text encoders that could be interesting to accurately represent medical features
 123 was found through study of the previously published pre-trained language models in the field:
 124 PubMedBERT (Gu et al., 2022), ClinicalBERT (Huang et al., 2019), SapBERT-from-PubMedBERT-
 125 full text (Liu et al., 2021), RadBERT (huggingface.co/StanfordAIMI/RadBERT), CXR-BERT-general
 126 (Boecking et al., 2022), CXR-BERT-specialized (Boecking et al., 2022) and finally the Clip text
 127 encoder (Radford et al., 2021b).

128 As described in section 2.1, we can gather radiology report data from CXR, and the corresponding
 129 abnormality labels as output by the CheXpert model (Irvin et al., 2019). Then, for each particular
 130 *text_encoder* model and the corresponding *report_list* of elements *report*, one can run the *report*
 131 through the model and get a representation $text_encoder(report)$. Nevertheless, there exist several
 132 ways to extract embeddings from these text encoders, all based on a transformer architecture:
 133 extracting the last layer hidden state of the associated CLS token, "CLS hidden state"; extracting the
 134 last layer hidden states of each tokens and averaging these representations, "mean hidden states";
 135 using the pooler output, "pooler output"; Using a model specific extraction method, if available,
 136 "model specific".

137 The combination of a *text_encoder* model and the associated extraction method *extraction* gives
 138 a function $extraction \circ text_encoder$ that takes an input report and outputs a document-level
 139 representation. This way, for a defined *text_encoder* model and *extraction* method, one can obtain
 140 the document-level embeddings of radiology reports and then assess the quality of these embeddings
 141 and therefore the capability of a text-encoder to encode medical content.

For the evaluation, we first obtain the document-level embeddings on the impression section of each
 radiology report, obtained through regex parsing. This gives:

$$impression_embeddings = extraction \circ text_encoder(impression_sections)$$

142 For all the text-encoders that we study, the latent representations are of dimension 768. Therefore for
 143 700 impression sections, *impression_embeddings* is a 700×768 matrix.

Then, we can compute the impression-impression similarities in the latent space

$$similarities = impression_embeddings \times impression_embeddings^T$$

We then compute a metric, that we denote the *CheXpert@k* metric, that for each report *i* find the *k*
 most similar reports, and then measure the proportion of reports that share the same CheXpert label.
 If *chexpert_labels* is a list of the chexpert labels corresponding to the reports, we have:

$$CheXpert@k_i = \frac{sum(chexpert_labels[argsort(similarities[i])[-k :]])}{k} == chexpert_labels[i]$$

144 And then over all reports we get $CheXpert@k = \frac{\sum_{i=1}^n CheXpert@k_i}{n}$

145 Notice that in the implementation of this metric, a filter is added to *CheXpert@k_i* so that among
 146 the *k* most similar reports, the report being compared to is not retrieved. In addition, the metric
 147 *CheXpert@k* can be computed over each class instead: so for each abnormality class, we can
 148 average the *CheXpert@k_i* scores, where the similarities are still computed over the reports of all
 149 classes. A macro-averaged score can then be retained for comparison purposes.

150 2.5 Textual Projection

151 Building up upon the stable diffusion work, we propose as a first method to generate domain-specific
152 images to replace the CLIP text encoder, kept as frozen during the stable diffusion original training,
153 with a domain-specific text encoder, typically pre-trained on biomedical or radiology data. The goal
154 behind this architecture change is to hopefully rely on the better understanding the new text encoder
155 has of radiology inputs and therefore provide better latent representations, that the U-Net will then be
156 conditioned upon to generate synthetic images.

157 Nevertheless, simply replacing the CLIP text encoder with a new one should lead to catastrophic
158 performance, given that latent spaces can be structured in a very different manner. There are no
159 guarantee that any latent feature is redundant between the two text encoders. We therefore propose to
160 train a projection capable of translating, in parts, the latent representations of one text encoder to the
161 other. So that running radiology prompt through the in-domain text encoder, and then projecting these
162 latent representations through this trained projection, should allow embeddings to be well-enough
163 aligned for the U-Net conditioning to work, but still provide enhanced clinical representations through
164 the in-domain text encoder added knowledge.

165 To train this projection, we use the LAION-400M dataset and define a *projection* as a MLP model,
166 taking a 768-dimensional input and mapping it to a 768-dimensional translated output. As a first
167 approach, we take $projection = Linear \circ ReLU \circ LayerNorm \circ Linear$ and train it using MLE
168 loss. At inference time, images can be generated by using the in-domain text encoder along the
169 projection, and hopefully having enough clinical features passing through while keeping most of the
170 CLIP latent space structure so that the U-Net conditioning allows for clinically correct generated
171 images.

172 Notice that the prompts the model is trained on can have an impact on the performance, that we try to
173 measure: we explore object-oriented prompts of the form "a photo of a ..." and style-oriented prompts
174 of the form "a photo in the style of a ...", with lexical variants of these two base prompts.

175 2.6 Textual Embeddings Fine-tuning

176 Following the approach of Gal et al. (2022), the stable diffusion model can be further trained to
177 generate better looking images for the radiology setting by focusing on the embeddings of the text
178 encoder. In this case, during training, the variational autoencoder, the U-Net, as well as all the other
179 layers of the text encoder are frozen. In addition, a new token gets introduced, that can either describe:
180 patient-level features, such as gender, age and body weight; procedure-level features, such as body
181 part and modality; abnormality-level features, such as "no findings" or "pleural effusion".

182 As an example, we could introduce the token $\langle lung - xray \rangle$ that is supposed to describe both a
183 body part, lungs, and a modality, X-ray. This learning approach, denoted *Textual Inversion*, zero out
184 all the gradients associated with the embeddings of the already existing tokens, and in the end only
185 learn the embedding of this newly introduced token.

186 Then, during training, input prompts with these new tokens are introduced, along associated radiology
187 images. The rest is very similar to original training of the stable diffusion model, in that the model
188 gets used to generate a synthetic image, and the noise at several timesteps in both the forward and
189 backward process of the U-Net are passed through a MSE loss. Gradients are then used to only
190 update the embeddings of the newly introduced tokens.

191 2.7 U-Net fine-tuning

192 Finally, in a similar approach to Ruiz et al. (2022), one can improve the baseline stable diffusion
193 model to generate better domain-specific images by relying on a U-Net fine-tuning. Instead of
194 switching text encoders and using a projection (see Section 2.5) or training the embeddings of new
195 tokens (see Section 2.6), one could keep all components frozen and the original CLIP text encoder, to
196 only further train the U-Net part. In this sense, the setting is very similar to the approach of Section
197 2.5, except that no new token gets added, and the freezing is over the set of parameters of the U-Net.

198 Then, the training is similar to the training of the original stable diffusion model, relying on MSE
199 loss at several time steps of the denoising process to progressively converge to better generation of
200 in-domain images.

201 **3 Results**

202 **3.1 Training details**

203 Experiments were conducted on several devices, depending on their compute hungriness. VAE and
204 text encoder experiments were run in local, with M1 Pro and M1 Max GPUs. Textual projection
205 relied on 3 NVIDIA Quadro P5000 GPUs, with a single run taking 3 hours for 10k training steps
206 in the case of document-level training, and 8 hours for 10k training steps in the case of token-level
207 training, when using only one of these GPUs. Textual embedding fine-tuning and U-net fine-tuning
208 used a NVIDIA V100 GPU and took respectively 1 hour for 3k training steps and 15 minutes for 400
209 training steps.

210 To conduct our experiments and in particular access model weights, we relied heavily on the *Hugging*
211 *Face* library (Wolf et al., 2019) and the recently released *diffusers* (von Platen et al., 2022). The
212 stable diffusion weights we used come from the *CompVis/stable-diffusion-v1-4* repo. Weights of
213 other in-domain text encoders are the ones associated with each corresponding publication.

214 **3.2 Variational autoencoder**

215 700 CXR images from MIMIC were encoded and decoded using the pretrained VAE from the
216 Stable-Diffusion-v1.4 pipeline. Quantitative assessment showed a low reconstruction error (RMSE
217 41.0 ± 8 ; median, 41.4; range, 20.4 - 76.3; PSNR, 33.6 ± 1.8 ; median 33.3; range, 28.1 - 39.5) and
218 a high structural similarity of original and reconstructed images (SSIM, 0.92 ± 0.02 ; median, 0.93;
219 range, 0.8 - 0.96). See Figure 3 for details. Image quality metrics did not depend on the class labels
220 of the images (data not shown).

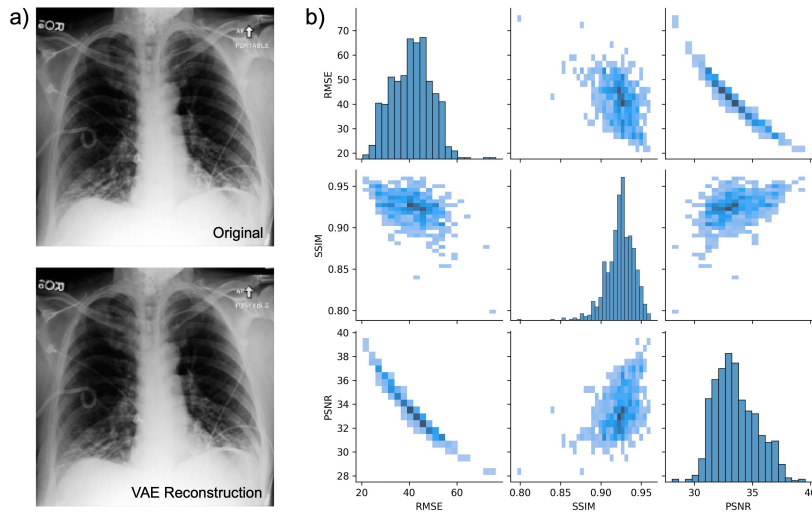


Figure 3: Image reconstruction analysis. a) Original (top) and reconstructed (bottom) image. The small burnt-in annotations in the top right corner get scrambled (seen in almost all samples), while the vast majority of other features (e.g., rib contours, devices) are well-preserved. b) Distribution of image quality metrics assessed for each image-reconstruction pair. RMSE: Root mean square error. SSIM: Structural similarity index measure. PSNR: Peak signal-to-noise ratio.

221 Visual analysis yielded a generally good perceived reconstruction quality (Mean visual score
222 4.51 ± 0.54 ; median score, 5; range, 3 - 5). No reconstruction resulted in a completely non-diagnostic
223 image (score 1) or altered the diagnostic information in a potentially problematic way (score 2).
224 Almost all burnt-in text annotations were scrambled beyond recognition, however, diagnostic features
225 were well preserved in almost all cases. Most of the score deductions were for blurred device
226 components, cerclages and wires that couldn't be traced reliably after reconstruction, or blurred rib
227 contours.

Table 1: Classification results for original and reconstructed CXR images from the MIMIC-CXR dataset

Label	Prevalence	Acc. (%)			F1 (%)		
		original	recon.	%change	original	recon.	%change
Atelectasis	33.3	40.1	40.7	1.4	52.0	52.4	0.7
Cardiomegaly	34.0	45.9	47.4	3.4	54.4	54.8	0.7
Consolidation	13.1	22.9	23.9	4.4	25.0	25.2	1.0
Edema	20.7	37.1	42.7	15.0	39.6	41.3	4.4
En. Mediastinum	13.1	23.4	23.6	0.6	23.4	22.8	-2.7
Fracture	2.4	15.7	19.0	20.9	4.5	4.7	3.9
Lung Lesion	4.6	21.3	25.1	18.1	5.8	5.8	-1.0
Lung Opacity	32.4	39.4	39.6	0.4	50.0	49.6	-0.8
Pleural Effusion	38.1	51.4	52.7	2.5	60.6	61.3	1.1
Pleural Other	2.1	55.4	50.3	-9.3	5.5	4.9	-9.8
Pneumonia	14.6	30.0	33.9	12.9	26.9	25.4	-5.3
Pneumothorax	9.1	24.3	24.7	1.8	18.0	18.0	0.5

228 The reconstruction process negatively impacted the classification performance for the "Pleural
 229 Other" label (accuracy 50.3% for reconstructed vs. 55.4% for reconstructed and original images,
 230 respectively). Interestingly, most other labels were predicted with similar (Atelectases, Cardiomegaly,
 231 Enlarged Cardiomediastinum, Lung Opacity, Pleural Effusion, Pneumothorax) or higher accuracy
 232 (Edema, Fracture, Lung Lesion, Pneumonia) from the reconstructed images. See Table 1 for details.
 233 The latent embeddings generated by the pretrained DenseNet-121 pairs were highly similar for
 234 image-reconstruction pairs (mean cosine similarity, 0.99 ± 0.01 ; median, 0.99; range, 0.94 - 1.00).

235 3.3 Text Encoder

236 Various text encoders and associated embedding methods are assessed on radiology reports in order
 237 to evaluate which method can retain optimum clinical knowledge in the latent representations.

238 Following the definitions in section 2.4 of the *text_encoder* models, the *extraction* methods and
 239 the metric $CheXpert@k$, we compute in Table 2 for each model and each method the macro-average
 240 of the $CheXpert@k$ score aggregated per abnormality class, with $k = 10$. As seen in the table,
 241 the method "CLS hidden state" is in general the one that works best to maximize the quality of
 242 the document-level representations of the impression sections. In addition, the model CXR-BERT-
 243 specialized is the one that reaches highest performance, taking for each model the corresponding
 244 *extraction* method that worked best.

245 Then, using the *extraction* method that works best for each model, we can compute class-wise
 246 $CheXpert@k$ scores as well as the macro-averaged ones. These results are aggregated in Table 3.
 247 As a baseline, we use a bag-of-words approach that outputs a similarity measure between two reports
 248 using an intersection over union measure. This baseline does not create any embeddings, but provide
 249 a token-based similarity measure: we observe that the latent representations of the best models, on
 250 top of contracting the text space, better encode document-level content and result in higher scores.

251 We remark that PubMedBERT, ClinicalBERT and CXR-BERT-general are three models that perform
 252 significantly less well than the other models, and should therefore, if possible, not be preferred for
 253 tasks that involve radiology reports. On the contrary, the two best performing models are CXR-
 254 BERT-specialized and the CLIP text encoder. As CLIP text encoder was not specifically trained
 255 on radiology reports, this underlines the quality of the training and the associated model. Using
 256 CXR-BERT-specialized instead would only improve performance by +15%.

257 For these reasons, we explore the textual projection with CXR-BERT-specialized, but also assess
 258 CLIP performance to be high enough to not justify replacing the text encoder in the various textual
 259 inversion and U-Net fine-tuning experiments.

Table 2: Macro-average of *CheXpert@k* scores computed per abnormality class, over the impression sections of a set of radiology reports. Models that are better at retaining medical features get a higher score.

Model	CLS hidden state	Mean hidden states	Pooler output	Model specific
PubMedBERT	30.8	23.6	20.6	None
ClinicalBERT	26.3	35.1	14.3	None
SapBERT	49.1	48.7	41	None
RadBERT	54.2	32.8	34.7	None
CXRBERTgeneral	32.4	25.4	31.6	None
CXRBERTspecialized	61.1	34.5	None	50.3
ClipTextEncoder	7.0	42.8	52.9	None

Table 3: For each text encoder and the associated best *extraction* method as computed in Table 2, class-wise and macro-averaged *CheXpert@k* scores are computed. Higher scores denote better capability at retaining important clinical features in the structure of the latent space.

Abnormality	Base	Pub.	Clin.	Sap.	Rad.	gen.	spe.	Clip.
Atelectasis	33.4	21.8	19.2	54.2	53	23.4	64.2	52.8
Cardiomegaly	21.6	20.8	10.2	51	53	21.6	67.2	47.6
Consolidation	36	13.4	35.8	39.6	38	35.4	27	38.4
Edema	62.8	54	62.6	64.6	67.2	47.4	85.4	72
Enlarged Cardiomeastinum	38	21.2	30.2	41.8	44.8	35.2	37.6	42.6
Fracture	49	36.2	35.6	73.2	72.6	50.8	83.2	74.2
Lung Lesion	30.2	24	21.2	32	37.8	24.8	56.2	33.8
Lung Opacity	20.4	16.2	20.6	20.4	34.2	20.4	23.2	25.6
No Finding	78.4	82.2	75.4	74.8	79.8	75.4	76.8	80.6
Pleural Effusion	46.4	25	39.4	42.6	65.8	24.2	72.2	68
Pleural Other	21.6	13.6	17.8	36	43.4	16.8	54	34.6
Pneumonia	53.8	33.8	40.4	42.6	44.4	24	45	54
Pneumothorax	56.4	39.6	60.6	65.2	73.6	28.6	92.8	72
Support Devices	32.6	29.2	23	49.6	50.8	25.8	70.4	44.8
Macro	41.5	30.8	35.1	49.1	54.2	32.4	61.1	52.9

260 3.4 Radiology Image Generation

261 Comparing the various methods introduced in Section 2, we use the Fréchet inception distance as
 262 introduced in Section 2.3 to measure the quality of the reconstructed images. The results are compiled
 263 in Table 4, along an empirical sample of images as produced by each method in Figure 4.

264 For the most simple prompt "A photo of a lung xray", progress is done only with the last method
 265 that consists in training the U-Net. In particular, no progress is observed with the token embedding
 266 training (also known as textual inversion). For more complex prompts such as "A photo of a lung xray
 267 with a visible pleural effusion", the stable diffusion baseline shows limitations, being outperformed
 268 by both textual inversion and U-Net fine-tuning.

269 The textual projection does not seem to converge well enough: samples from Figure 4 shows the
 270 generated images to be out-of-domain. Nevertheless, we estimate that a more complex architecture,
 271 instead of our simple 1-hidden-layer projection, could be worth exploring: if projection-based
 272 domain-adaptation turns out to produce interesting examples, this could open the door to very quick
 273 domain-adaptation for the large amount of pre-trained text encoders that are now available.

274 Out of all the methods, the U-Net fine-tuning seems by far the most promising: it gets the lowest
 275 FID-scores and obviously the most realistic outputs. Nevertheless, we notice that this underlines the
 276 limitations of our non-medical-based metric: samples clearly show that U-Net fine-tuning with prior
 277 leads the model to learn the difference between "no findings" and "pleural effusion", something a
 278 model trained without a prior can not do. As seen in Table 4, FID fails to capture this improvement.
 279 We assess that further progress in the domain-specific generation of images for radiology would

Table 4: Evaluation of the quality of generated images with different methods for adapting stable diffusion to the radiology domain. Scores represent the Fréchet inception distance (FID), and lower scores mean better generated images.

Training Strategy	A photo of a lung xray	A photo of a lung xray with a visible pleural effusion	A photo in the style of a lung xray
<i>Original model</i>			
Stable Diffusion	0.097	0.151	
<i>Textual Projection</i>			
<i>CXR-BERT-specialized</i>			
No Projection	0.124	0.144	
Document-level projection	0.266	0.104	
Token-level projection	0.201	0.257	
<i>Token embedding training</i>			
Object, radiology	0.108	0.058	0.092
Object, lung	0.135		0.135
Style, radiology	0.101	0.057	0.084
Style, lung	0.130		0.083
<i>U-Net training</i>			
Trained on no findings	0.057	0.043	
Trained on no findings and abnormality	0.034	0.041	
Trained on no findings and abnormality with prior	0.170	0.086	

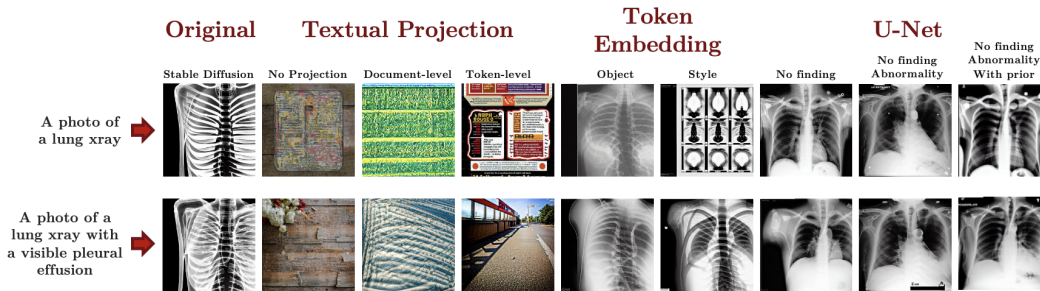


Figure 4: Images generated by various methods conditioned on radiology text prompts.

280 require the design/use of domain-specific metrics, that would be able to capture the ability of the
 281 model to correctly insert abnormalities that are coherent with the conditioning text prompt.

282 4 Conclusion

283 In this paper, we assessed the recently released stable diffusion model, including its variational autoen-
 284 coder, the U-Net and the associated CLIP text encoder, and its capacity to produce clinically relevant
 285 images based on prompts that describe observable abnormalities. We conducted quantitative and
 286 qualitative evaluations, showing that: the variational autoencoder is powerful enough to reconstruct
 287 radiological images, including abnormalities and clinically relevant features; the CLIP text encoder
 288 accurately represents simple radiology-specific text prompts, outperforming 4 out of the 6 reviewed
 289 domain-specific text encoders. We explored textual projection, a domain-adaptation method that we
 290 designed, textual inversion and U-Net fine-tuning, and, with the latter, obtained a model capable of
 291 generating synthetic radiology images that are visually and quantitatively exceeding the baseline, and
 292 that can correctly represent abnormalities.

293 Building upon this work, we would like to further explore the potential of diffusion-based model to
 294 learn a wide-range of abnormalities, being able to combine them, as well as extending the research
 295 to other modalities and body parts. A limitation of our approach is that the employed metrics have
 296 limited capacity to assess the clinical correctness of the generated images. In addition, our fine-tuned
 297 stable diffusion model lacks diversity in the images it generates, probably due to the small range of
 298 samples they were trained on. Finally, the text prompts the models are conditioned on are synthetic
 299 and do not fully correspond to the wording used in the clinical setting, so that models capable of
 300 being conditioned on entire or partial radiology reports are an area of future research.

301 References

- 302 Benedikt Boecking, Naoto Usuyama, Shruthi Bannur, Daniel C. Castro, Anton Schwaighofer, Stephanie
303 Hyland, Maria Wetscherek, Tristan Naumann, Aditya Nori, Javier Alvarez-Valle, Hoifung Poon, and Ozan
304 Oktay. Making the most of text semantics to improve biomedical vision–language processing, 2022. URL
305 <https://arxiv.org/abs/2204.09817>.
- 306 Rishi et al Bommasani. On the opportunities and risks of foundation models, 2021. URL <https://arxiv.org/abs/2108.07258>.
- 308 Joseph Paul Cohen, Mohammad Hashir, Rupert Brooks, and Hadrien Bertrand. On the limits of cross-domain
309 generalization in automated x-ray prediction, 2020. URL <https://arxiv.org/abs/2002.02497>.
- 310 Joseph Paul Cohen, Joseph D. Viviano, Paul Bertin, Paul Morrison, Parsa Torabian, Matteo Guarrera, Matthew P
311 Lungren, Akshay Chaudhari, Rupert Brooks, Mohammad Hashir, and Hadrien Bertrand. TorchXRyVision:
312 A library of chest X-ray datasets and models. In *Medical Imaging with Deep Learning*, 2022. URL
313 <https://github.com/mlmed/torchxrayvision>.
- 314 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or.
315 An image is worth one word: Personalizing text-to-image generation using textual inversion, 2022. URL
316 <https://arxiv.org/abs/2208.01618>.
- 317 Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng
318 Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language
319 processing. *ACM Transactions on Computing for Healthcare*, 3(1):1–23, jan 2022. doi: 10.1145/3458754.
320 URL <https://doi.org/10.1145/2F3458754>.
- 321 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Günter Klambauer, and Sepp Hochre-
322 iter. Gans trained by a two time-scale update rule converge to a nash equilibrium. *CoRR*, abs/1706.08500,
323 2017. URL <http://arxiv.org/abs/1706.08500>.
- 324 Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. Clinicalbert: Modeling clinical notes and predicting hospital
325 readmission, 2019. URL <https://arxiv.org/abs/1904.05342>.
- 326 Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilicus, Chris Chute, Henrik Marklund,
327 Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, Jayne Seekins, David A. Mong, Safwan S. Halabi, Jesse K.
328 Sandberg, Ricky Jones, David B. Larson, Curtis P. Langlotz, Bhavik N. Patel, Matthew P. Lungren, and
329 Andrew Y. Ng. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison,
330 2019. URL <https://arxiv.org/abs/1901.07031>.
- 331 Alistair E. W. Johnson, Tom J. Pollard, Seth J. Berkowitz, Nathaniel R. Greenbaum, Matthew P. Lungren,
332 Chih ying Deng, Roger G. Mark, and Steven Horng. MIMIC-CXR, a de-identified publicly available
333 database of chest radiographs with free-text reports. *Scientific Data*, 6(1), December 2019. doi: 10.1038/
334 s41597-019-0322-0. URL <https://doi.org/10.1038/s41597-019-0322-0>.
- 335 Fanguy Liu, Ehsan Shareghi, Zaiqiao Meng, Marco Basaldella, and Nigel Collier. Self-alignment pretraining for
336 biomedical entity representations. In *Proceedings of the 2021 Conference of the North American Chapter of*
337 *the Association for Computational Linguistics: Human Language Technologies*, pages 4228–4238, Online,
338 June 2021. Association for Computational Linguistics. URL [https://www.aclweb.org/anthology/](https://www.aclweb.org/anthology/2021.naacl-main.334)
339 [2021.naacl-main.334](https://www.aclweb.org/anthology/2021.naacl-main.334).
- 340 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish
341 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning
342 transferable visual models from natural language supervision. *CoRR*, abs/2103.00020, 2021a. URL <https://arxiv.org/abs/2103.00020>.
- 344 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
345 Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable
346 visual models from natural language supervision, 2021b. URL <https://arxiv.org/abs/2103.00020>.
- 347 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional
348 image generation with clip latents, 2022.
- 349 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Bjorn Ommer. High-resolution image
350 synthesis with latent diffusion models. In *2022 IEEE/CVF Conference on Computer Vision and Pattern*
351 *Recognition (CVPR)*. IEEE, June 2022. doi: 10.1109/cvpr52688.2022.01042. URL [https://doi.org/10.](https://doi.org/10.1109/cvpr52688.2022.01042)
352 [1109/cvpr52688.2022.01042](https://doi.org/10.1109/cvpr52688.2022.01042).

- 353 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth:
 354 Fine tuning text-to-image diffusion models for subject-driven generation, 2022. URL [https://arxiv.org/](https://arxiv.org/abs/2208.12242)
 355 [abs/2208.12242](https://arxiv.org/abs/2208.12242).
- 356 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed
 357 Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho,
 358 David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language
 359 understanding, 2022. URL <https://arxiv.org/abs/2205.11487>.
- 360 Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta,
 361 Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million
 362 image-text pairs, 2021. URL <https://arxiv.org/abs/2111.02114>.
- 363 Christoph Schuhmann, Romain Beaumont, Cade W Gordon, Ross Wightman, mehdi cherti, Theo Coombes,
 364 Aarush Katta, Clayton Mullis, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Mitchell
 365 Wortsman, Richard Vencu, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5b: An open
 366 large-scale dataset for training next generation image-text models. In *Thirty-sixth Conference on Neural*
 367 *Information Processing Systems Datasets and Benchmarks Track*, 2022. URL [https://openreview.net/](https://openreview.net/forum?id=M3Y74vmsMcY)
 368 [forum?id=M3Y74vmsMcY](https://openreview.net/forum?id=M3Y74vmsMcY).
- 369 Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the
 370 inception architecture for computer vision. *CoRR*, abs/1512.00567, 2015. URL [http://arxiv.org/abs/](http://arxiv.org/abs/1512.00567)
 371 [1512.00567](http://arxiv.org/abs/1512.00567).
- 372 Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj,
 373 and Thomas Wolf. Diffusers: State-of-the-art diffusion models. [https://github.com/huggingface/](https://github.com/huggingface/diffusers)
 374 [diffusers](https://github.com/huggingface/diffusers), 2022.
- 375 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac,
 376 Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine
 377 Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and
 378 Alexander M. Rush. Huggingface’s transformers: State-of-the-art natural language processing, 2019. URL
 379 <https://arxiv.org/abs/1910.03771>.

380 Checklist

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

- 385 • Did you include the license to the code and datasets? **[Yes]**
- 386 • Did you include the license to the code and datasets? **[No]** The code and the data are
 387 proprietary.
- 388 • Did you include the license to the code and datasets? **[N/A]**

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

- 392 1. For all authors...
 - 393 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 394 contributions and scope? **[Yes]**
 - 395 (b) Did you describe the limitations of your work? **[Yes]**
 - 396 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]**
 - 397 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 398 them? **[Yes]**
- 399 2. If you are including theoretical results...
 - 400 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - 401 (b) Did you include complete proofs of all theoretical results? **[Yes]**

- 402 3. If you ran experiments...
- 403 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 404 mental results (either in the supplemental material or as a URL)? [No] Will be included
- 405 in the non-double-blinded submission.
- 406 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 407 were chosen)? [Yes] We included data collection, splits, compute details, number of
- 408 training steps, code details when they could be explained at a high-level. We did not
- 409 include hyperparameter details.
- 410 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 411 ments multiple times)? [Yes] When relevant, especially for the autoencoder experi-
- 412 ments.
- 413 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 414 of GPUs, internal cluster, or cloud provider)? [Yes]
- 415 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 416 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 417 (b) Did you mention the license of the assets? [Yes] We mentioned the reference, where
- 418 the license can be found.
- 419 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 420 No new assets.
- 421 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 422 using/curating? [Yes] No consent was needed.
- 423 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 424 information or offensive content? [No] Data from only already publicly available
- 425 datasets was used.
- 426 5. If you used crowdsourcing or conducted research with human subjects...
- 427 (a) Did you include the full text of instructions given to participants and screenshots, if
- 428 applicable? [No] Not applicable
- 429 (b) Did you describe any potential participant risks, with links to Institutional Review
- 430 Board (IRB) approvals, if applicable? [Yes]
- 431 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 432 spent on participant compensation? [No] Not applicable