# Medical diffusion on a budget: textual inversion for medical image generation

Bram de Wilde*          Anindo Saha          Richard P. G. ten Broek

Henkjan Huisman

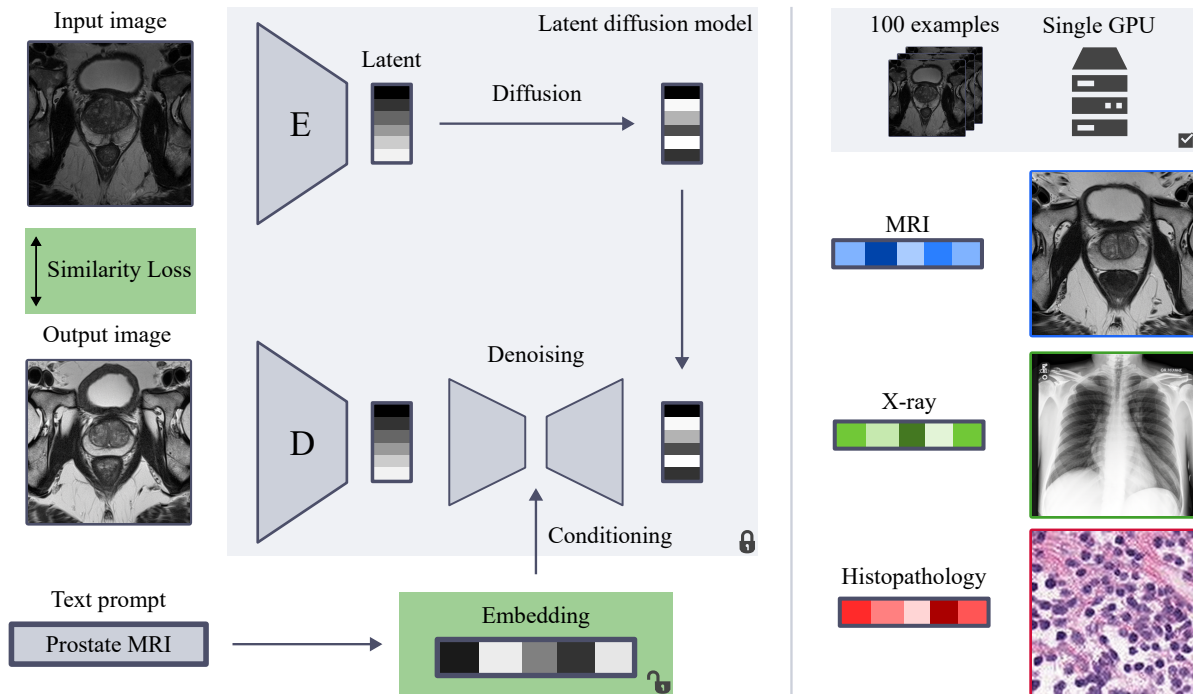Radboud University Medical Center

*contact@bramdewilde.com

Figure 1: The textual inversion fine-tuning process for diffusion models trains a text conditioning embedding for a new token (e.g. Prostate MRI) using a small set of example images, while keeping the rest of the architecture frozen. In this work we show that this allows adaption of latent diffusion models to a variety of medical imaging modalities, using only 100 examples and a single consumer-grade GPU.

## Abstract

*Diffusion-based models for text-to-image generation have gained immense popularity due to recent advancements in efficiency, accessibility, and quality. Although it is becoming increasingly feasible to perform inference with these systems using consumer-grade GPUs, training them from scratch still requires access to large datasets containing captioned data and significant computational resources. In the case of medical image generation, the availability of large, publicly accessible datasets that include text reports is limited due to legal and ethical concerns. While training a diffusion model on a private dataset may address this is-sue, it is not always feasible for institutions lacking the necessary computational resources. This work demonstrates that pre-trained Stable Diffusion models, originally trained on natural images, can be adapted to various medical imaging modalities by training text embeddings with textual inversion. In this study, we conducted experiments using small medical datasets comprising only 100 samples from three medical modalities. The embeddings were trained in a matter of hours, while still retaining diagnostic relevance in image generation. Our experiments were designed to achieve several objectives. Firstly, we aimed to fine-tune the training and inference processes of textual inversion, which revealed that larger embeddings and more examples*

*are required for the medical domain. Secondly, we validated our approach by demonstrating a 2% increase in the diagnostic accuracy (AUC) for detecting prostate cancer on MRI, which is a challenging multi-modal imaging modality, from 0.78 to 0.80. Thirdly, we performed simulations by interpolating between healthy and diseased states, combining multiple pathologies, and inpainting to show the flexibility of the trained embeddings and to demonstrate potential for fine control of disease appearance. Finally, the embeddings trained in this study are small (less than 1 MB), which facilitates easy sharing of medical data with reduced privacy concerns. The code and all embeddings trained in this work will be made available online for future research.*

## 1. Introduction

Image generation has increasingly captured the attention of many researchers, spurring an impressive progression in text-to-image generation. In particular, diffusion models have gained enormous popularity through their ability to generate high-quality and diverse images, conditioned on a text prompt.[12, 8, 25, 24, 29] Of all text-to-image model implementations, Stable Diffusion has by far generated the biggest impact in terms of users, owing to the fact that it is both released under a permissive license and operable using a single GPU.[26] The unprecedented availability and performance of Stable Diffusion is perhaps best demonstrated by artists ringing the alarm bell, for fear of AI systems replicating their work or style without their consent, seemingly even leading to a lawsuit against the Stable Diffusion team.[35] Whatever your stance in this debate, it shows that humans are starting to have difficulty distinguishing AI-generated art from human art.

Generating art or even photorealistic images, however, has some room for error. Even if generated images are not completely physically realistic, they may still be appealing or impressive. The medical imaging field, on the other hand, places a higher bar on generation quality.[37, 32] Images need not only be anatomically correct, but diagnostically as well. In theory, there is no reason why a system like Stable Diffusion for medical imaging is impossible. You simply need a large, varied and ideally public dataset of images with captions.[31] In practice, however, ethical and legal issues often get in the way of sharing medical data.[30, 4] This issue is especially difficult to tackle for radiology reports due to their unstructured nature, whereas they are the natural source for high quality captions. For one of the few public datasets of this caliber that exists, MIMIC-CXR, Chambon et al. have indeed demonstrated that it is possible to train a latent diffusion model capable of generating chest x-ray images with high fidelity and diversity through free text prompts.[16, 5] They trained the system using up to 170,000 images on 64 A100 GPUs.

On top of data sharing issues, some modalities and pathologies are inherently scarce: certain types of scans can be expensive or experimental and some diseases are rare or tied to specific demographics. For these reasons, especially in the medical domain, it is essential to have computationally feasible methods which can fine-tune existing models towards a smaller set of a specific modality or disease. In this paper, we pick one such method, Textual Inversion, and rigorously explore its capacities for adapting Stable Diffusion to medical imaging.[10] All experiments are done using a single RTX2070 GPU and all training scripts and embeddings are shared online. We summarize our contributions as follows:

1. We show, through careful tuning and experiments, that a diffusion model trained on natural images can be adapted to produce a wide variety of realistic medical images.

2. We demonstrate the practical value of our approach, by improving cancer classification models in the low-data regime, using synthetic data.

3. We demonstrate that the trained embeddings are highly flexible, by showing (1) interpolation between healthy and diseased state, (2) inpainting for fine control of disease appearance and (3) that multiple embeddings can be combined to generate images with multiple pathologies.

## 2. Related work

Since our study explores a specific fine-tuning method for diffusion models applied to medical image generation, we briefly review popular fine-tuning methods for diffusion in general and cover studies applying diffusion models in the medical domain.

### 2.1. Fine-tuning diffusion models

Various methods have different computational requirements and output files of different sizes. We review three popular methods in decreasing order of compute and output size.

In [27] Ruiz et al. fine-tune the denoising U-net component of a diffusion model, using a handful of images for introducing a new concept. This method is typically employed on 24 GB GPUs and results in checkpoints of several GB when fine-tuning Stable Diffusion models.

LoRA was introduced as a method to fine-tune large language models, which freezes original model weights and introduces rank decomposition matrices into the Transformer architecture.[14] Recently, this fine-tuning method was incorporated into the popular *diffusers* library. It can be deployed on a 11 GB GPU and results in shareable

files of under 5 MB for a fully fine-tuned latent diffusion model.[36, 1]

In [10], Gal et al. introduce textual inversion, which fine-tunes a diffusion model by finding a new word embedding for newly introduced concepts, resulting in extremely small files of under 1 MB. Like LoRA, this method is deployable on a 11 GB GPU. In this study, we adopt textual inversion as fine-tuning method, because it has low computational requirements and the smallest file output, but we have no reason to expect our results to be limited to this fine-tuning method.

## 2.2. Medical image generation

Several papers have applied diffusion to medical imaging, with a wide range of applications including anomaly detection, segmentation, registration and modality transfer with image-to-image translation.[18] Specifically for medical image generation, several recent works have trained diffusion models for image generation. Pre-trained models are often trained on 2D RGB datasets, but many medical imaging modalities are 3D. Recently, studies such as [19] and [23] have trained diffusion models from scratch on 3D data, or even on 4D data.[20] Several other works studied text-to-image latent diffusion models for medical imaging.[5, 2]

Closest to our work is [6], where Chambon et al. explore various methods to adapt a pre-trained Stable Diffusion model to chest X-ray generation. They performed experiments with both textual inversion and fine-tuning the U-net component of Stable Diffusion, similar to [27]. They find that textual inversion works, but that fine-tuning the U-net is more effective, especially with more complex prompts. They fine-tune using 5 examples per class.

Our work extends on this by exploring textual inversion more deeply, by training with more examples and bigger embeddings. Additionally, we demonstrate the flexibility of the approach through example applications and by adapting to multiple and more complex modalities beyond chest X-ray. In contrast to most other studies, we intentionally do not train from scratch and use small datasets to explore the feasibility of diffusion in low-data and low-compute environments.

## 3. Methods

### 3.1. Image generation

All images are generated with Stable Diffusion v2.0, using an interactive open-source web interface.[26, 3] Images are sampled using the ancestral Euler scheduler.[17] The main inference parameters which influence image generation quality are (1) the number of steps for the sampling scheduler and (2) the classifier-free guidance (CFG) scale.[13] Using more steps for sampling typically leads to better image quality, but increases the inference time. The CFG scale can be used to set the trade-off between sample quality and sample diversity. Practically, a high CFG scale means images will follow the text prompt more closely at the expense of diversity. Conversely, a low CFG scale results in images that deviate more from the prompt and consequently have lower fidelity, but higher diversity.

To introduce a medical modality as a new concept to a pre-trained diffusion model, we use the textual inversion process.[10] Put simply, this process finds a vector in the text embedding space that optimally represents the concept. Practically, this is done by freezing the entire architecture apart from this embedding vector and doing backpropagation with a similarity loss, as illustrated in Figure 1. We train embeddings with a constant learning rate of 0.005 for 50,000 steps with a batch size of 1. In the work of Gal et al., during training prompts are generated from a list of templates, for instance: "*a photo of a <embedding>*" or "*a rendering of a <embedding>*". Since these templates do not necessarily make sense in a medical imaging context, we simplify this by always prompting the model only with the embedding name during training. We experiment with the amount of images used to train an embedding and the vector size of the embedding. To evaluate the impact of inference and training parameters on generation quality, we compute the Fréchet Inception Distance using 100 generated samples compared to 100 real examples for each parameter setting.[33]

To investigate the usability of the trained embeddings, we also experiment with combining multiple trained embeddings using composable diffusion.[22] This method allows prompting with a combination of embeddings using an AND operator in the prompt, e.g. "*<cardiomegaly> AND <pleural effusion>*" to generate an image with both cardiomegaly and pleural effusion present. Additionally, this methods allows a weight to be given to each embedding, to tune the strength of each embedding separately. In this study, we use this to experiment with interpolating between healthy and diseased states and to generate images with multiple diseases present.

### 3.2. Classification

For classification experiments, we train ResNet-18 models, pre-trained on ImageNet.[11, 7] Models are trained with a fixed learning rate of $10^{-4}$ with the Adam optimizer for 6250 batches of 32 images.[21] This corresponds to 100 epochs for the biggest synthesized dataset (2000 synthesized cases). AUC is evaluated on a separate validation set during training and performance of the best validation checkpoint on a separate test set is reported. We apply random horizontal flipping, gaussian noise, intensity transformations, channel dropout, translation, scaling and rotation as data augmentation. Detailed training scripts will be shared online.

### 3.3. Datasets

To showcase the wide applicability of textual inversion for adaptation to medical imaging, we demonstrate results on three different types of data: multi-modal MRI, chest X-ray and histopathology.

#### 3.3.1 Multi-modal MRI - PI-CAI

The main dataset used in this work is a recently released public dataset of 1500 prostate MRI cases. This dataset was released as part of the PI-CAI (Prostate Imaging: Cancer AI) challenge, where the task is to detect clinically significant prostate cancer.[28] Each case is a 3D MRI scan featuring three modalities: T2-weighted imaging (T2W), apparent diffusion coefficient maps (ADC) and diffusion-weighted imaging (DWI). Since this work adapts a pre-trained 2D diffusion model, we extract one 2D axial slice per case. Each case is first resampled to a resolution of $3 \times 0.5 \times 0.5$ mm and then center-cropped to a $90 \times 150 \times 150$ mm ($30 \times 300 \times 300$ px) region. For negative cases, we select the median prostate slice using provided full prostate segmentations. For positive cases, we select the slice with maximum tumor area, according to the provided tumor segmentation maps. Each slice is finally upsampled to $512 \times 512$ px. Each modality is encoded as one of the RGB channels when training multi-modal embeddings. The training, validation and test set or the classification experiments each consist of 100 randomly sampled negative slices and 100 randomly sampled positive slices. The embeddings are trained on the training set.

Prostate MRI is special in that it provides multiple images of the pelvic region that each depict a unique pathophysiologically relevant aspect for the specific disease purpose. This multi-modal dataset deviates from the natural image distribution seen during pre-training, and thus makes it highly challenging for image generation. To be relevant, the generated medical images need to be diagnostically consistent across the modalities. For example, prostate cancer should appear dark on ADC, bright on DWI, and show a blurry structure in T2W.

#### 3.3.2 Chest X-ray - CheXpert

CheXpert is a large public dataset of 224,316 chest radiographs, with corresponding labels for 14 different observations.[15] In principle, this is a multi-label classification task, but since we explicitly investigate compositional prompting with the learned embeddings, we only sample images with a single class. Specifically, we sample 100 AP-view radiographs to train embeddings for the following four observations: No Finding (healthy), Cardiomegaly, Pleural Effusion and Pneumonia. Each radiograph is first cropped to non-zero borders. Then the longest edge is resized to 512 px, while keeping the aspect ratio fixed. Finally, the image is zero-padded to a square resolution of $512 \times 512$ px.

#### 3.3.3 Histopathology - Patch Camelyon

Patch Camelyon is a public dataset of 327,680 $96 \times 96$ px patches extracted from histopathology whole-slide images of lymph node sections, originally released as part of the Camelyon16 challenge.[34, 9] Each patch has a corresponding binary label indicating the presence of metastatic tissue. We randomly selected 100 negative and 100 positive patches, upsampled them to $512 \times 512$ px and used them to train embeddings.

## 4. Experiments

### 4.1. Adapting TI settings to medical imaging

We start by tuning the settings of the Textual Inversion training process, using the PI-CAI prostate cancer dataset. All embeddings in this section were trained using 2D T2-weighted healthy prostate slices, meaning they feature only one modality of the three modalities present for each 2D slice. The T2-weighted image clearly shows the anatomy and is therefore easiest to judge qualitatively.

Firstly, we determine the optimal inference parameters for an embedding with vector size 64, trained with 100 cases. In particular, we tune the number of sampling steps and the classifier-free guidance (CFG) scale. The full results are shown in Table 1. Following this table, in the remainder of the paper we do inference with 100 steps and a CFG scale of 2, unless specified otherwise. Even though a CFG scale of 1 gives slightly lower FID, images are much less accurate, as will be discussed later in this section.

Secondly, we study the impact of the vector size of the trained embeddings, varying it from 8 to 64, again using 100 cases during training. The limit of the CLIP encoder used for Stable Difffusion is 75, so a bigger embedding size would need adaptation of the framework. The results are shown in Table 2 and show that using a larger embedding size is better.

Thirdly, we study the impact of the number of cases used during training, using an embedding size of 64. We vary the number of cases from 5, as proposed in [10], to 100. The results are shown in Table 3 and show that using more cases leads to better embeddings.

Figure 2 shows the effect of the parameters studied in this section visually on a single random seed. A high number of sampling steps improves generation quality, with the generations for 25 and 50 steps showing incorrect anatomy for the bladder. Although a CFG scale of 1 results in the lowest FID score in Table 1, visually the results are much worse, featuring inaccurate general anatomy. A high CFG scale (e.g.
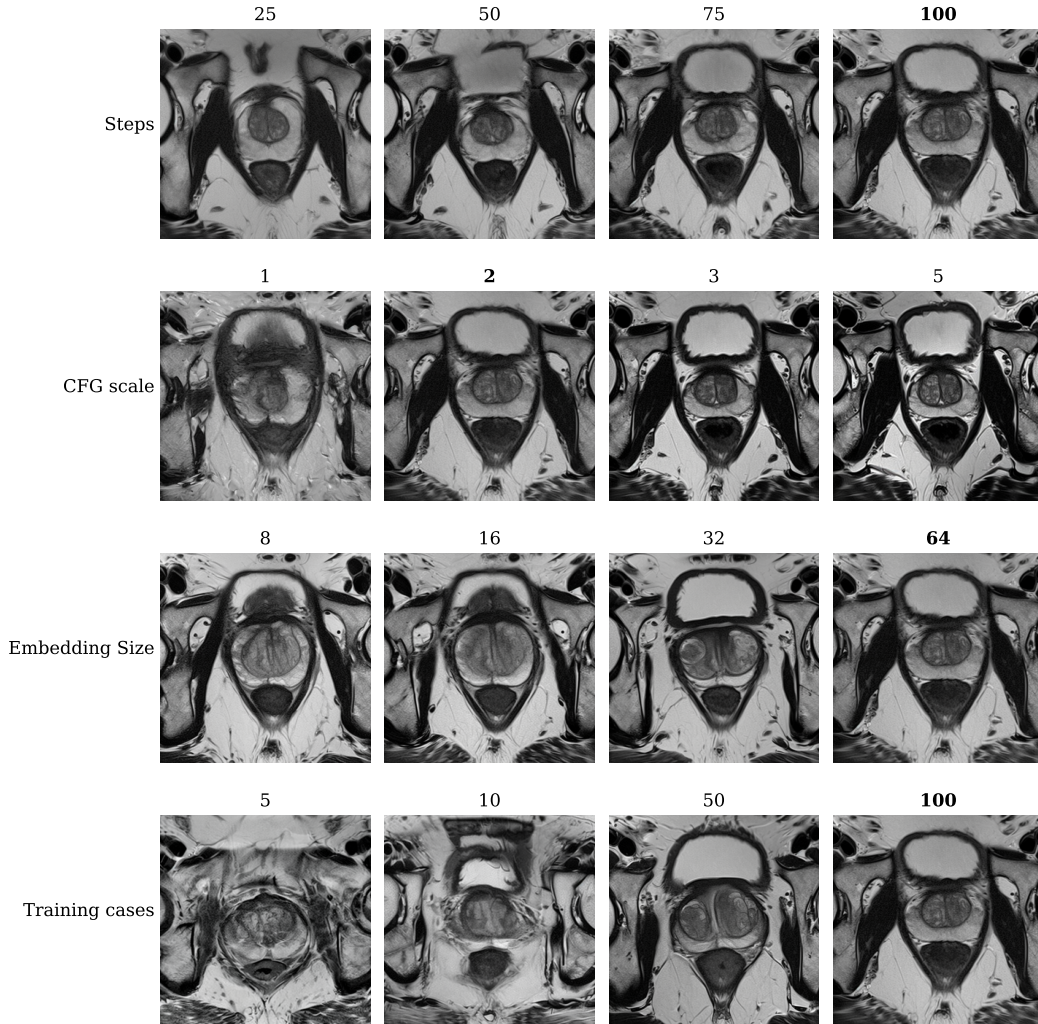
Figure 2: Visual examples illustrating the effect of varying inference and training settings for T2-weighted prostate MRI, all generated using the same random seed. Columns with a bold title indicate optimal values. Row labels indicate the parameter that changes along the column, with bold values set for the other parameters. For example: in the top row the number of steps changes, but CFG scale, embedding size and training cases are 2, 64 and 100, respectively.

5 in the Figure) also leads to bad results, showcased here by simplified structure inside the prostate and a curious fractured pelvic bone. The difference between CFG scale 2 and 3 is not that large, but upon manual inspection we find that a CFG scale of 2 gives better generations overall, as seems to be confirmed by the lower FID score in Table 1. The embedding size is clearly optimally chosen to be large, with sizes 8 and 16 showing inaccurate generation, particularly of the bladder. Although size 32 looks better, the structure of the prostate itself is not nearly as good as generated by the size 64 embedding. Finally, the impact of the amount of training cases seems to trump all other settings, where 5 and 10 cases produce very unrealistic images. The embedding trained with 100 cases generates images with the most realistic prostate structure. In general we found that the FID score, although it identifies general trends, is not the most suitable metric to judge generation quality. We chose optimal parameters by inspecting generation results visually as well.

Following the results presented in this section, all embeddings in the remainder of this paper were trained with a size of 64 vectors per token, using 100 cases per class. For inference, we sample with 100 steps and a CFG scale of 2, unless specified otherwise.

| Steps | CFG scale | FID ↓ |
|---|---|---|
| 25 | 2 | 194 |
| 50 | 2 | 184 |
| 75 | 2 | 176 |
| 100 | 2 | **171** |
| 100 | 1 | **168** |
| 100 | 2 | 171 |
| 100 | 3 | 202 |
| 100 | 4 | 211 |
| 100 | 5 | 222 |

Table 1: FID score for embeddings generated with varying number of sampling steps and CFG scale.

| Embedding size | FID ↓ |
|---|---|
| 8 | 177 |
| 16 | 184 |
| 32 | 218 |
| 64 | **171** |

Table 2: FID score for embeddings trained with varying size (vectors per token).

| Training cases | FID ↓ |
|---|---|
| 5 | 203 |
| 10 | 181 |
| 50 | 181 |
| 100 | **171** |

Table 3: FID score for embeddings trained with varying number of training cases.

| Real cases | Synthetic cases | AUC ↑ |
|---|---|---|
| 200 | 0 | $0.780 \pm 0.017$ |
| 0 | 200 | $0.737 \pm 0.019$ |
| 0 | 2000 | $0.766 \pm 0.020$ |
| 200 | 200 | $0.773 \pm 0.015$ |
| 200 | 2000 | $\mathbf{0.803 \pm 0.009}$ |
| 0 | 2000* | $0.562 \pm 0.036$ |
| 200 | 2000* | $0.745 \pm 0.012$ |

Table 4: AUC for binary prostate cancer classifiers, trained with real or synthetic data, or a combination of both. Synthetic cases marked with an asterisk (*) were generated with an embedding trained on only 10 cases, instead of 100. The mean test AUC over 10 training runs is shown, together with the standard deviation.

## 4.2. Improving classification with synthetic data

In this section, we experiment with using synthetic data to train prostate cancer classification models on multi-modal prostate MRI. Embeddings are trained on two sets of 100 cases, with only negative or only positive cases. With these embeddings, we generate up to 1000 cases for each class and use various combinations of real and synthetic data to train classification models. Results are shown in Table 4, with the primary result being that augmenting the 200-case training set with 2000 synthesized cases leads to a 2% improvement in AUC, from 0.78 to 0.80. Note that these 2000 synthesized cases are based on embeddings trained with the same 200-case set used to train the classification models. This shows that the generated cases actually add non-trivial variation to the data distribution and that the embedding does not simply reproduce training cases. Furthermore, models trained with only synthetic cases do not see a large drop in performance, indicating that the synthetic cases are diagnostically accurate. Finally, to confirm visual results from section 4.1, classification models trained with synthetic cases generated with embeddings trained on 10 cases instead of 100 show a dramatic drop in performance. This confirms our finding in section 4.1 that more cases are needed for textual inversion for medical data.

## 4.3. Composability of embeddings

One particularly attractive property of textual inversion is that you can use your newly introduced concept in text prompts, even by combining multiple trained embeddings. In this section we give visual evidence that this works, to an extent, for medical data as well. Firstly, in Figure 3, disease state is gradually increased from healthy to diseased, using composable diffusion. For instance, the cardiomegaly radiograph in the second column (25% diseased) is generated with a prompt like "0.25*<healthy> AND 0.75*<cardiomegaly>". This seems to work well across the modalities studied in this paper: the tumors in the prostate example become gradually more prominent (darker on ADC, brighter on DWI); the heart in the cardiomegaly example appears to grow from left to right; the tissue in the lymph node metastasis example becomes gradually more abnormal.

Even combining multiple conditions in a single image appears to work, as illustrated by Figure 4. Here we start with a healthy image and consequently add pleural effusion, pneumonia and cardiomegaly to the prompt for a single random seed. For the image with all three diseases, we gave each embedding a strength of 0.5 and found that it works better to increase the CFG scale to 3.

These examples show that the trained embeddings are flexible and go beyond fine-tuning towards a single new concept. The fact that rendering of disease progression or
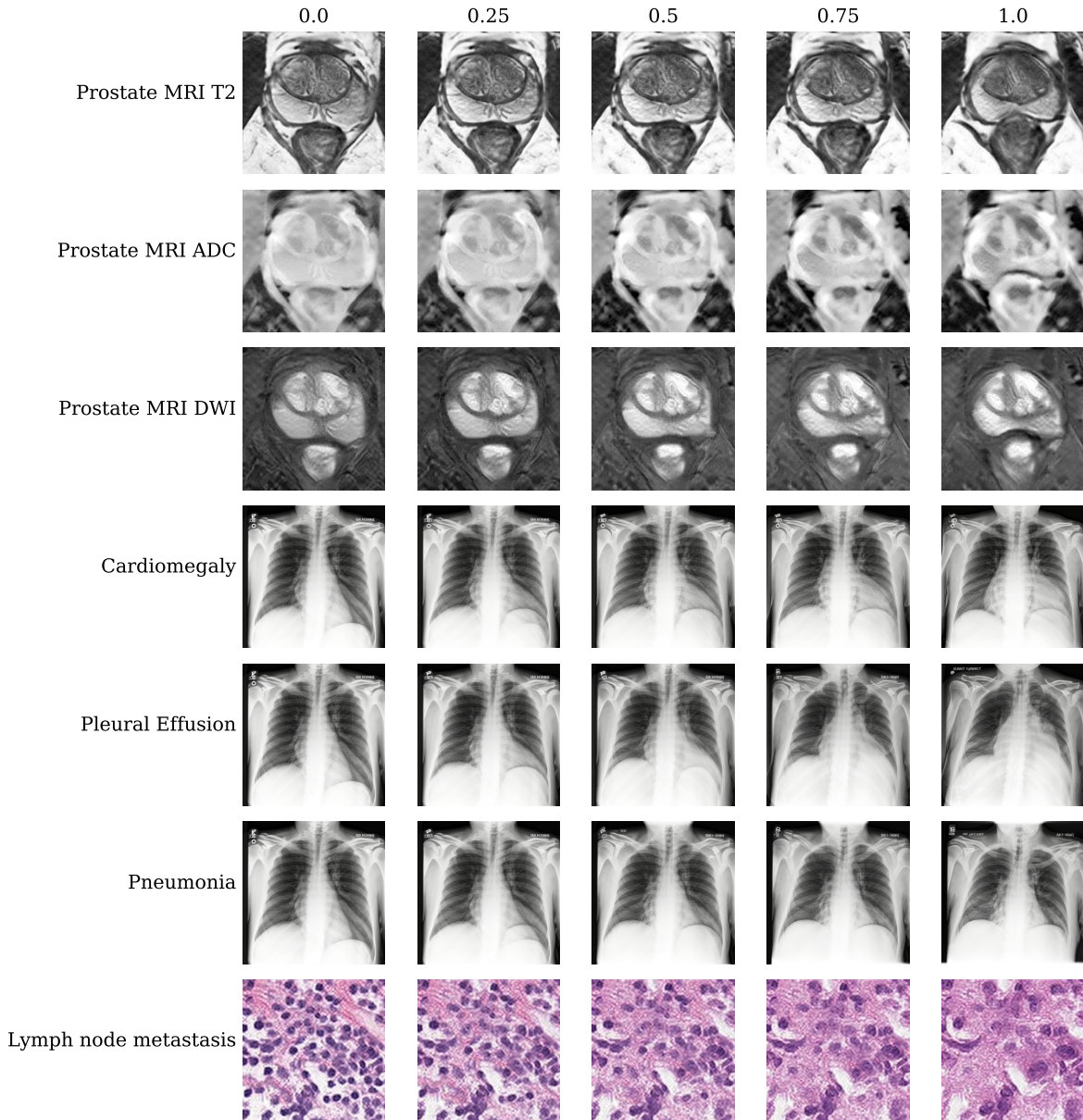
Figure 3: Visual examples illustrating interpolation between healthy and diseased states for multi-modal Prostate MRI, various pathologies on Chest X-Ray and lymph node metastasis in histopathology. The column titles show the trade-off between healthy and diseased. The Chest X-Ray examples are all generated using the same random seed. The prostate images are cropped to the prostate region for visibility.

accurate depiction of multiple conditions is possible, while the embeddings have only been trained on cases with a single condition present is a promising result. Practically, this could be useful to generate cases with rare combinations of conditions in a single image, or to simulate disease progression for medical surveillance settings.

## 4.4. Controlling disease appearance with inpainting

In this section, we demonstrate the potential of inpainting to precisely control where disease shows in an image. Starting with a generated healthy example, a portion of the image is masked. The diffusion model then denoises the masked part of the image, while conditioned on a specific disease embedding. For example, for the top row in Fig-
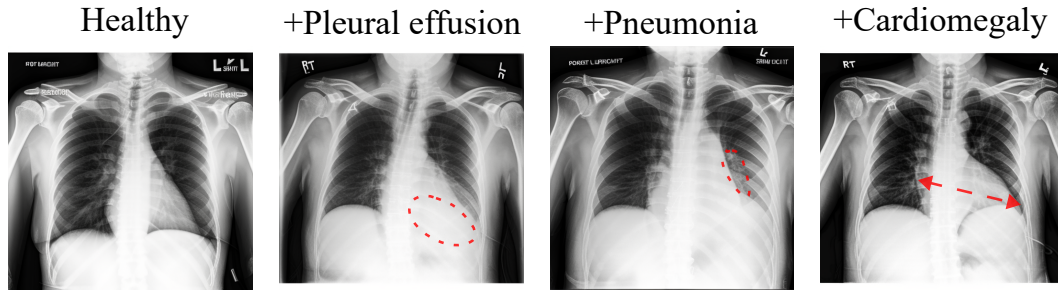
**Figure 4:** Visual example illustrating that multiple embeddings can be composed, to show multiple pathologies in a single image. From left to right, pleural effusion, pneumonia and cardiomegaly are progressively added to a healthy generated example.
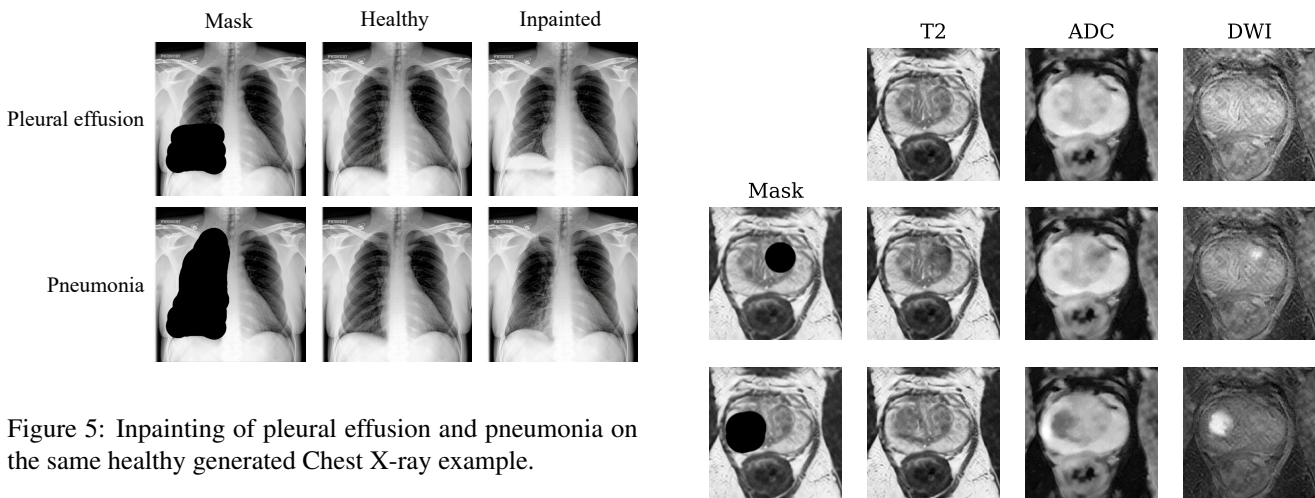


**Figure 5:** Inpainting of pleural effusion and pneumonia on the same healthy generated Chest X-ray example.



**Figure 6:** Inpainting of prostate cancer in different locations on the same healthy generated Prostate MRI example. The top row shows the original healthy case, with the bottom rows showing inpainting in different locations with varying mask size.

ure 5, the bottom of the left lung of the healthy image is masked. The diffusion model inpaints it, while it is conditioned with the pleural effusion embedding. This results in pleural effusion appearing in the masked region. In Figure 5, it is used to force pleural effusion or pneumonia appears at the left lung. Similarly, in Figure 6, the same healthy prostate example is masked in two different locations, with a different mask size. When inpainting conditioned on the positive embedding, this generates tumors at those locations of corresponding sizes. Similar to the examples in section 4.3, this allows engineering of examples with specific disease appearance and could for instance be useful to generate more cases with rare tumor locations.

## 5. Conclusion

In this paper, we show that pre-trained latent diffusion models can be adapted to a variety of modalities in the medical domain, using textual inversion. High quality images can be generated with embeddings trained on 100 examples on a single consumer-grade GPU. We showcased various possible applications: improvement of diagnostic mod-els in the low-data regime by adding synthetic cases during training, simulation of disease progression and generation of images with specific disease appearance. Although a dedicated diffusion model trained on a large captioned medical dataset would likely generate better images, our results are promising for institutions with limited computational resources. Especially for situations where collecting a large dataset is not feasible, such as rare diseases, this approach is suitable and would also be compatible with a medically pre-trained model. Finally, since the trained embeddings are extremely small files, they may facilitate sharing of medical information with reduced privacy concerns.

# References

[1] Using LoRA for Efficient Stable Diffusion Fine-Tuning. https://huggingface.co/blog/lora. 3

[2] Mohamed Akrout, Bálint Gyepesi, Péter Holló, Adrienn Poór, Blága Kincső, Stephen Solis, Katrina Cirone, Jeremy Kawahara, Dekker Slade, Latif Abid, Máté Kovács, and István Fazekas. Diffusion-based Data Augmentation for Skin Disease Classification: Impact Across Original Medical Datasets to Fully Synthetic Images, Jan. 2023. 3

[3] AUTOMATIC1111. Stable Diffusion web UI, Mar. 2023. 3

[4] Jasper Bovenberg, David Peloquin, Barbara Bierer, Mark Barnes, and Bartha Maria Knoppers. How to fix the GDPR's frustration of global biomedical research. *Science*, 370(6512):40–42, Oct. 2020. 2

[5] Pierre Chambon, Christian Bluethgen, Jean-Benoit Delbrouck, Rogier Van der Sluijs, Małgorzata Połacin, Juan Manuel Zambrano Chaves, Tanishq Mathew Abraham, Shivanshu Purohit, Curtis P. Langlotz, and Akshay Chaudhari. RoentGen: Vision-Language Foundation Model for Chest X-ray Generation, Nov. 2022. 2, 3

[6] Pierre Chambon, Christian Bluethgen, Curtis P. Langlotz, and Akshay Chaudhari. Adapting Pretrained Vision-Language Foundational Models to Medical Imaging Domains, Oct. 2022. 3

[7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, June 2009. 3

[8] Prafulla Dhariwal and Alex Nichol. Diffusion Models Beat GANs on Image Synthesis, June 2021. 2

[9] Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes van Diest, Bram van Ginneken, Nico Karssemeijer, Geert Litjens, Jeroen A. W. M. van der Laak, and and the CAMELYON16 Consortium. Diagnostic Assessment of Deep Learning Algorithms for Detection of Lymph Node Metastases in Women With Breast Cancer. *JAMA*, 318(22):2199–2210, Dec. 2017. 4

[10] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion, Aug. 2022. 2, 3, 4

[11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, Las Vegas, NV, USA, June 2016. IEEE. 3

[12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. In *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020. 2

[13] Jonathan Ho and Tim Salimans. Classifier-Free Diffusion Guidance, July 2022. 3

[14] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-Rank Adaptation of Large Language Models, Oct. 2021. 2

[15] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, Jayne Seekins, David A. Mong, Safwan S. Halabi, Jesse K. Sandberg, Ricky Jones, David B. Larson, Curtis P. Langlotz, Bhavik N. Patel, Matthew P. Lungren, and Andrew Y. Ng. CheXpert: A Large Chest Radiograph Dataset with Uncertainty Labels and Expert Comparison, Jan. 2019. 4

[16] 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 de-identified publicly available database of chest radiographs with free-text reports. *Scientific Data*, 6(1):317, Dec. 2019. 2

[17] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the Design Space of Diffusion-Based Generative Models, Oct. 2022. 3

[18] Amirhossein Kazerouni, Ehsan Khodapanah Aghdam, Moein Heidari, Reza Azad, Mohsen Fayyaz, Ilker Hacihaliloglu, and Dorit Merhof. Diffusion Models for Medical Image Analysis: A Comprehensive Survey, Nov. 2022. 3

[19] Firas Khader, Gustav Mueller-Franzes, Soroosh Tayebi Arasteh, Tianyu Han, Christoph Haarburger, Maximilian Schulze-Hagen, Philipp Schad, Sandy Engelhardt, Bettina Baessler, Sebastian Foersch, Johannes Stegmaier, Christiane Kuhl, Sven Nebelung, Jakob Nikolas Kather, and Daniel Truhn. Medical Diffusion: Denoising Diffusion Probabilistic Models for 3D Medical Image Generation, Jan. 2023. 3

[20] Boah Kim and Jong Chul Ye. Diffusion Deformable Model for 4D Temporal Medical Image Generation, June 2022. 3

[21] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization, Jan. 2017. 3

[22] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B. Tenenbaum. Compositional Visual Generation with Composable Diffusion Models, Jan. 2023. 3

[23] Walter H. L. Pinaya, Petru-Daniel Tudosiu, Jessica Dafflon, Pedro F. da Costa, Virginia Fernandez, Parashkev Nachev, Sebastien Ourselin, and M. Jorge Cardoso. Brain Imaging Generation with Latent Diffusion Models, Sept. 2022. 3

[24] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical Text-Conditional Image Generation with CLIP Latents, Apr. 2022. 2

[25] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-Shot Text-to-Image Generation, Feb. 2021. 2

[26] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-Resolution Image Synthesis with Latent Diffusion Models, Apr. 2022. 2, 3

[27] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation, Aug. 2022. 2, 3

[28] Anindo Saha, Jasper Jonathan Twilt, Joeran Sander Bosma, Bram van Ginneken, Derya Yakar, Mattijs Elschot, Jeroen Veltman, Jurgen Fütterer, Maarten de Rooij, and Henkjan Huisman. Artificial Intelligence and Radiologists at Prostate Cancer Detection in MRI: The PI-CAI Challenge (Study Protocol). Technical report, Zenodo, June 2022. 4

[29] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding, May 2022. 2

[30] James Scheibner, Jean Louis Raisaro, Juan Ramón Troncoso-Pastoriza, Marcello Ienca, Jacques Fellay, Effy Vayena, and Jean-Pierre Hubaux. Revolutionizing Medical Data Sharing Using Advanced Privacy-Enhancing Technologies: Technical, Legal, and Ethical Synthesis. *Journal of Medical Internet Research*, 23(2):e25120, Feb. 2021. 2

[31] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5B: An open large-scale dataset for training next generation image-text models, Oct. 2022. 2

[32] Youssef Skandarani, Pierre-Marc Jodoin, and Alain Lalande. GANs for Medical Image Synthesis: An Empirical Study, July 2021. 2

[33] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the Inception Architecture for Computer Vision, Dec. 2015. 3

[34] Bastiaan S. Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. Rotation Equivariant CNNs for Digital Pathology, June 2018. 4

[35] James Vincent. Getty Images is suing the creators of AI art tool Stable Diffusion for scraping its content. https://www.theverge.com/2023/1/17/23558516/ai-art-copyright-stable-diffusion-getty-images-lawsuit, Jan. 2023. 2

[36] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, and Thomas Wolf. Diffusers: State-of-the-art diffusion models, Mar. 2023. 3

[37] Xin Yi, Ekta Walia, and Paul Babyn. Generative adversarial network in medical imaging: A review. *Medical Image Analysis*, 58:101552, Dec. 2019. 2