
Medical Flow Matching CT Translation

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 Medical imaging still lacks reliable methods to translate contrast-enhanced
2 arterial-phase CT scans into their non-contrast (native) form. We introduce
3 **Medical Flow Matching (MFM)**, which combines an efficient image-translation
4 paradigm with a bottleneck attention mechanism designed for medical images.
5 On the held-out test set, MFM achieves **MAE** = 6.229 HU, **SSIM** = 0.992,
6 and **PSNR** = 38.933, markedly outperforming traditional methods while cutting
7 generation time by 400 times compared with DDPM. With two-way conversion
8 using MFM, we can double the size of available public pathology-image CT
9 datasets. Training nnUNetv2 on MFM-generated (from contrast) images achieves
10 a **Dice score** of 0.926 versus 0.966 when trained on real (non-contrast) images –
11 95 % of baseline performance. By removing the need for an additional native scan,
12 MFM can reduce the radiation dose per examination by about 50 %, lowering pa-
13 tients' cumulative exposure and thereby decreasing the risk of radiation-associated
14 diseases.

15

1 Introduction

16 In recent years, medical imaging has become central to modern healthcare, providing essential
17 information for patient diagnostics, treatment planning, and monitoring disease progression. Machine
18 learning methods have significantly revolutionized this domain by automating routine radiological
19 tasks [1, 2], showing strong predictive abilities for disease diagnosis [3, 4], and helping reduce the
20 clinical workload of medical specialists [5, 6].

21 However, building effective machine-learning models for medical imaging still faces big challenges.
22 First, the lack of annotated medical datasets is a major barrier. Rare diseases are usually under-
23 represented because of their low prevalence, making it harder to gather enough data to train specialized
24 models. Public datasets often include only one modality, for example, contrast-enhanced series,
25 whereas native (non-contrast) images may be scarce or missing. Second, labeling medical images
26 demands specialized expertise, so annotation costs are much higher than in general image-labelling
27 tasks. High-quality, consistent annotations are vital for accurate predictive models, which makes the
28 problem even worse.

29 Computed tomography (CT) is a key tool in clinical diagnosis, giving important insights into anatomy
30 through both native (non-contrast) and contrast-enhanced scans. The latter involves administering a
31 contrast agent to enhance vascular visibility and support precise diagnosis. Currently, due to the trend
32 toward reducing patient radiation exposure during CT examinations, non-contrast (native) series are
33 not always performed. However in certain cases, they are still required specific findings. Therefore,
34 building reliable methods to convert contrast-enhanced CT scans (arterial phase) into native images
35 without contrast media remains important but challenging.

36 Our approach matters not only because it can help researchers and developers expand the modality
37 of existing datasets, but also because it offers a new way to utilize contrast-enhanced images when

38 native scans are unavailable. [7] estimates that CT studies in the United States in 2023 could lead
39 to about 103 000 future skin cancers. Our approach could help reduce this number by enabling the
40 creation of potentially diagnostically useful scan sequences without increasing radiation exposure.
41 Although dual-energy CT scanners are already capable of generating virtual non-contrast images,
42 they remain relatively rare [8]. Our method could help move the field in that direction by making this
43 opportunity more accessible.

44 In this paper, we present a Flow Matching method and TimeResNet model designed to translate
45 arterial-phase contrast-enhanced CT images into corresponding native CT scans. Our proposed
46 method and model address the challenges of limited data high annotation costs by effectively
47 leveraging existing data, which may reduce the need for additional medical image annotations and
48 lower the risks associated with radiation exposure.

49 This paper offers three main contributions:

- 50 1. We propose a new medical image-to-image translation method, that needs far less memory
51 that typical 3D approaches yet keeps consistency across axial slices.
- 52 2. We propose a neural network architecture (TimeResNet) for this method, which achieves
53 strong results, compared with existing architectures in MONAI [9].
- 54 3. The method can translate contrasts to natives and natives to contrast, letting users train a
55 single network instead of two specific tasks; this halves the training time.

56 To simplify further narration, by a Native Image or Native we mean a CT scan that was performed
57 without any contrast enhanced. Under Contrast Image or Contrast, we refer to the Arterial phase of
58 CT examination, usually after 15-35 seconds after administration of iodine-containing contrast.

59 **2 Related work**

60 Given the challenges outlined above, it is therefore unsurprising that the generation of synthetic
61 medical images has become increasingly popular. Initially, classical techniques [10, 11, 12, 13, 14]
62 were employed; these methods proved both effective and applicable across a range of tasks, yet
63 they nonetheless faced limitations in producing truly realistic medical imagery. The advent of
64 deep learning subsequently enabled the synthesis of more complex images. Soon thereafter, deep
65 neural networks began to appear, initially demonstrating marked improvements in image quality and
66 generalization performance relative to classical approaches.

67 As generative adversarial networks (GANs) matured, they found extensive application within medical
68 imaging—for example, in MRI-to-CT translation [15], in MRI-PET translation [16, 17, 18, 19, 20,
69 21, 22], image reconstruction [23, 24], and super-resolution [25, 26]. One of the most widespread
70 uses of GANs has been as a data-augmentation strategy to bolster the training of models for both
71 segmentation and classification tasks.

72 Diffusion models have gained considerable traction in recent years, demonstrating substantial promise
73 in the synthesis of medical images owing to their ability to generate high-quality outputs, their stable
74 training dynamics, and their flexibility during optimization [27]. Studies [28, 29, 30, 31, 32, 33, 34,
75 35, 36, 37] have confirmed the efficacy of these diffusion-based approaches. In two-dimensional
76 medical image generation, state-of-the-art models capture intricate anatomical details with negligible
77 distortion, to the extent that they are occasionally deemed suitable for clinical deployment. In the
78 domain of three-dimensional imaging, NVIDIA’s recent publication [38] on whole-patient generation
79 employed diffusion processes within a latent representation space to achieve remarkably realistic
80 volumetric outputs.

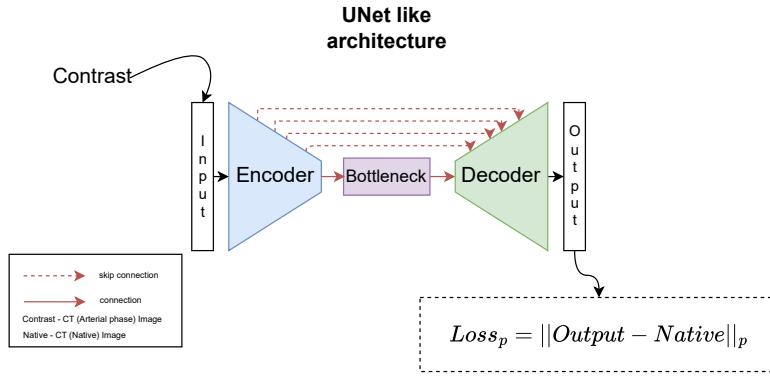
81 In the broader field of computer vision, there has been a marked trend toward generative techniques
82 based on flow matching [39, 40] as well as conditional generation frameworks [41, 42]. Such
83 methodologies are increasingly being adopted within medical imaging, particularly for tasks involving
84 image translation between different modalities [43, 44].

85 In this work, we focus on building a robust image-translation pipeline. We introduce targeted tweaks
86 to a proven architecture so it works well across different translation tasks without heavy tuning. For
87 instance, we have an AbdomenAtlas [45] dataset for 8,448 CT volumes with masks. Most of them
88 are studies with contrasts. Our method can almost double this dataset, which will allow us to use

89 data for specific tasks. This approach not only conserves time and computational resources but also
90 enables researchers to focus on other critical aspects, thereby enhancing overall efficiency in a range
91 of clinical scenarios.

92 3 Methods

93 To show how well our approach works, we compare it with several alternatives. Our simplest baseline
94 is a U-Net-style network: the contrast image goes into the encoder, and the decoder tries to reproduce
95 the native image. We train this network with mean absolute error (MAE). To be sure that any
96 performance gap is not caused by the choice of architecture alone, we test multiple network setups
97 and loss functions, following the procedure in [46]. $\mathcal{L}_{\text{MAE}}(\theta) = \|g_\theta(x_C) - x_N\|_1$, where g_θ – neural
network that we will train for regression task. In our experiments we set $p = 1$.



98 Figure 1: Baseline regression scheme.

99 3.1 Diffusion

100 So if we want to use diffusion for conditional image to image translation, we train a UNet denoiser
101 for the following loss. We will apply DDPM [47] to conditional image-to-image generation. π_0 –
102 distribution of target images, π_1 – distribution of conditional images. $\pi_0 \times \pi_1$ – joint distribution of
103 paired target and conditional images.

$$\mathcal{L}_{\text{DDPM}}(\theta) = \mathbb{E}_{\varepsilon \sim \mathcal{N}(0,1)} \mathbb{E}_{t \sim \mathcal{U}\{0, \dots, T\}} \mathbb{E}_{(\mathbf{x}_0, \mathbf{x}_1) \sim \pi_0 \times \pi_1} \|\varepsilon - \varepsilon_{\theta,t}(\alpha(t) \cdot \mathbf{x}_0 + \sigma(t) \cdot \varepsilon, \mathbf{x}_1)\| \quad (1)$$

104 We have conditions on $\alpha(t)$ and $\sigma(t)$. $\alpha(t), \sigma(t)$ such that $\forall t \in \{0, \dots, T\} \rightarrow \alpha(t)^2 + \sigma(t)^2 = 1$.
105 $\mathcal{U}\{0, \dots, T\}$ – is a uniform integer distribution.

106 If we apply this method in our case, we will get that: π_N – distribution of Native images (target
107 images), π_C – distribution of Contrast images (conditional images), $\pi_N \times \pi_C$ – joint distribution of
108 paired Native and Contrast images. So we choose standard linear scheme for β_t [47], using
109 the notation $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$, so we have $\alpha(t)^2 = \bar{\alpha}_t$, $\sigma(t)^2 = 1 - \bar{\alpha}_t$, $\beta_t =$
110 $\beta_{\text{start}} + (\beta_{\text{end}} - \beta_{\text{start}}) \cdot \frac{t}{T}$, where $T, \beta_{\text{start}}, \beta_{\text{end}}$ – hyperparameters that are set before training.

$$\mathcal{L}_{\text{DDPM}}(\theta) = \mathbb{E}_{\varepsilon \sim \mathcal{N}(0,1)} \mathbb{E}_{t \sim \mathcal{U}\{0, \dots, T\}} \mathbb{E}_{(\mathbf{x}_N, \mathbf{x}_C) \sim \pi_N \times \pi_C} \|\varepsilon - \varepsilon_{\theta,t}(\sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_N + \sqrt{1 - \bar{\alpha}_t} \cdot \varepsilon, \mathbf{x}_C)\| \quad (2)$$

111 For sampling we used DDIM-like scheme, to obtain deterministic images from the initial noise.

$$\mathbf{x}_{t-1} = \sqrt{\frac{1}{1 - \beta_t}} \cdot \mathbf{x}_t + \left[\sqrt{1 - \bar{\alpha}_t} - \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \beta_t}} \right] \cdot \varepsilon_{\theta}([\mathbf{x}_{\text{img}}, \mathbf{x}_t], t)$$

112 It is also important to understand that we must use the same noise from which the image will be
113 generated for all axial images of a single patient, otherwise this will lead to a very large heterogeneity
114 in the relative and axial projections.

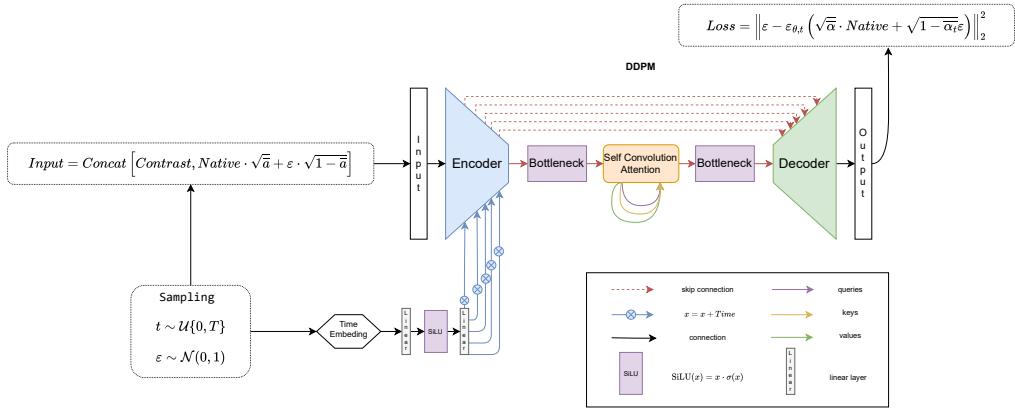


Figure 2: DDPM scheme. Neural network architecture and a visual diagram of training.

115 3.2 Flow Matching

116 Let $\{\pi_t\}_{t \in [0,1]}$ be a *probability path* (*probability path is a continuous trajectory of probability*
 117 *distributions that interpolates between an initial distribution and a target distribution over time.*) that
 118 continuously interpolates between an easy-to-sample reference π_0 (e.g. Gaussian noise) and the data
 119 distribution π_1 . Samples $\mathbf{x}_t \sim \pi_t$ evolve according to an ordinary differential equation (ODE)

$$\frac{d\mathbf{x}_t}{dt} = f_\theta(\mathbf{x}_t, t), \quad \mathbf{x}_0 \sim \pi_0, \quad \mathbf{x}_1 \sim \pi_1. \quad (3)$$

120 In our case we define $\pi_0 \times \pi_1$ - as pairs of images of the same patient but from different series, where
 121 π_0 - distribution of native images, π_1 - distribution of contrast images.

122 Define the convex combination $t \in [0, 1]$, $\mathbf{x}_t = (1 - t)\mathbf{x}_0 + t\mathbf{x}_1$
 123 and the *target velocity*

$$\mathbf{v}_t = \frac{\partial \mathbf{x}_t}{\partial t} = \mathbf{x}_1 - \mathbf{x}_0. \quad (4)$$

124 Flow Matching (FM) trains a neural vector field f_θ by regressing onto \mathbf{v}_t over random $(\mathbf{x}_0, \mathbf{x}_1)$ pairs
 125 and $t \sim \mathcal{U}[0, 1]$:

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t \sim \mathcal{U}[0,1]} \mathbb{E}_{(\mathbf{x}_0, \mathbf{x}_1) \sim \pi_0 \times \pi_1} [\|f_\theta(\mathbf{x}_t, t) - \mathbf{v}_t\|_2^2]. \quad (5)$$

126 Crucially, Eq. (5) is *simulation-free*: no stochastic differential equations need to be solved during
 127 training.

128 Below is a universal diagram of how the model works and learns using the Flow Matching method.

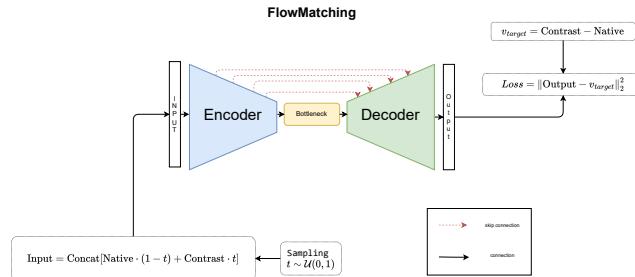


Figure 3: Flow Matching scheme.

129 We also developed and proposed our own model architecture, TimeResNet.

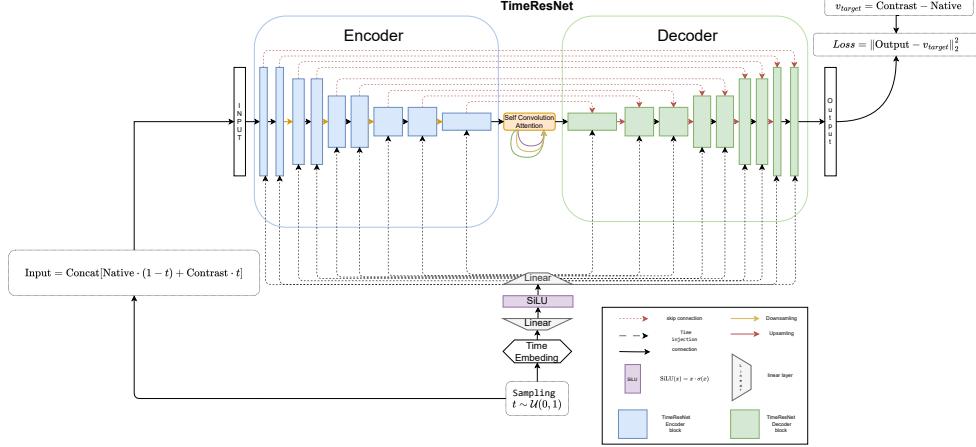


Figure 4: TimeResNet in Flow Matching scheme.

130 A more detailed description of the Self-Convolution Attention block is given below. We determined
 131 that an explicit attention mechanism was necessary to achieve high performance; however, widely
 132 adopted architectures such as SwinUNET [48] – which employ a limited, patch-based attention
 133 scheme—did not yield substantial gains. Consequently, we elected to insert our attention module
 134 between the Bottleneck layers, drawing inspiration from Yang et al. (2019) [49]. To stabilize training,
 135 we further incorporated Group Normalization [50] into the block.

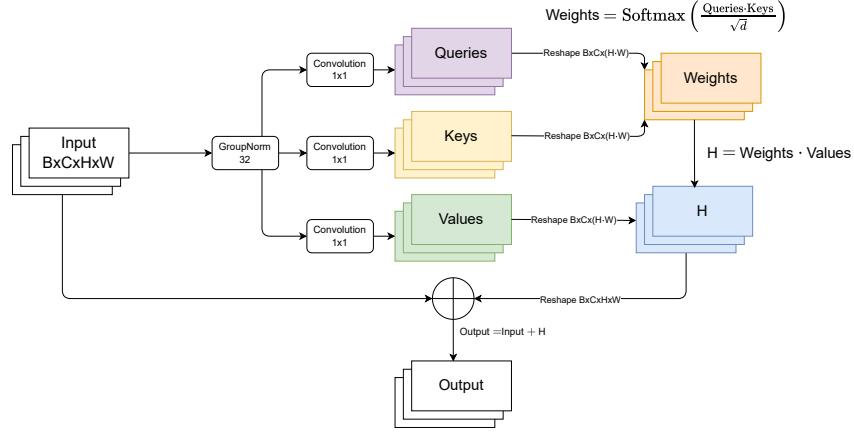


Figure 5: Self convolution attention scheme.

136 After training, sampling reduces to integrating Eq. (3) from $t = 1$ to $t = 0$ with a numerical ODE
 137 solver. We have the opportunity to sample both motifs into contrasts and natives from contrasts.
 138 We make full condition on linear combination (3.2) to improve stability of the work and remove
 139 stochasticity. To correctly add time intervals to our model, we do the following. We use a standard
 140 scheme with sine and cosine embeddings (link to the article), then we feed it into the MLP head,
 141 the results of which are fed into each of the ResNet blocks, which initially also uses a linear layer
 142 to translate the size to the current number of channels, and then the intermediate values are used as
 143 normalization. $h = \mathbf{A} \cdot h + \mathbf{B}$, where \mathbf{A} and \mathbf{B} are linear and trainable parameters in each block. A
 144 more detailed scheme for adding Time embeddings is presented in B

145 **4 Experiments**

146 A total of 120 abdominal CT studies (52561 images) were retrospectively collected, each containing
 147 both native and arterial phases (one study per patient), from 3 different CT stations from 3 different
 148 hospitals. The dataset was split into 80 (34992 images) training, 20 (8197 images) test, and 20
 149 (9372 images) held-out patients. Seventy imaging studies were obtained as follows: first, the
 150 anatomical region of interest was localized in both the native and contrast-enhanced series, next, the
 151 contrast-enhanced series were registered to the native series.

152 All images were taken in the original (512×512) shape and clipped to $[-1000, 1000]$ HU before
 153 scaling to $[-1, 1]$. As augmentations, we used axis-aligned flips and affine transformations. No
 154 intensity-based perturbations were used in order to preserve Hounsfield integrity.

155 All models were trained for 30 epochs with Adam optimizer, batch size was equal 2. We conducted all
 156 the experiments on a node of four NVIDIA RTX 3090 GPU (24 GB). To implement the experiments,
 157 we used the PyTorch [51] and MONAI framework.

158 **4.1 Training Dynamics**

159 Detailed model learning trajectory is presented below. We see that the graphs have reached a plateau,
 160 but we assume that further scaling in terms of the amount of data, training time, and model size may
 161 yield better results, but we tried to use all the resources available to us to present the most valid results
 162 in our opinion.

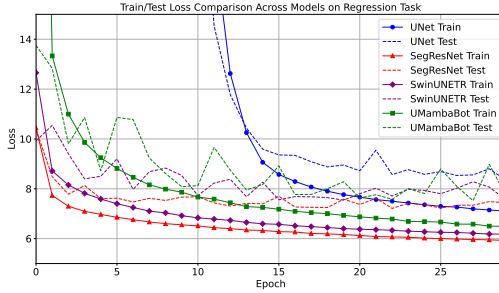


Figure 6: Dynamics of training different models in the regression task.

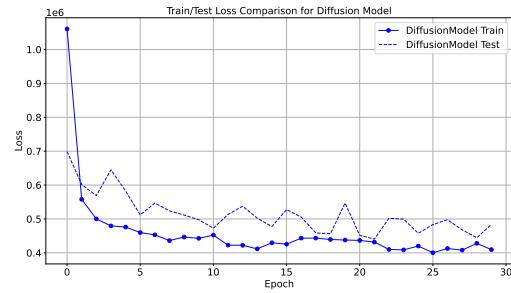


Figure 7: Dynamics of training model in the diffusion task.

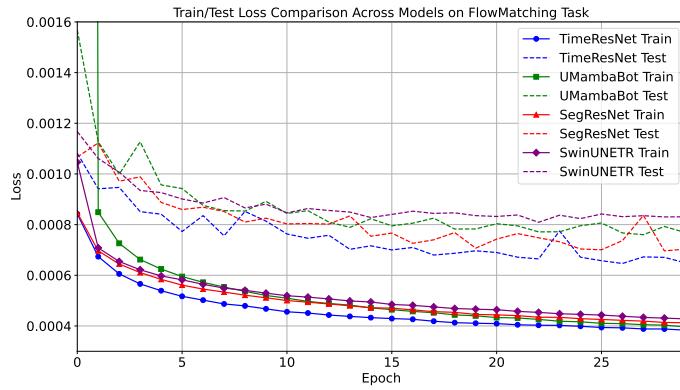


Figure 8: Dynamics of training different models in the flow matching method.

163 We can note that the TimeResNet model we have proposed shows the best convergence compared to
 164 the most popular options among CNN, Transformer, and SSM. Umamba architecture needs a certain
 165 number of epochs to start showing good results.

166 **4.2 Visual results**

167 In order to evaluate the generation quality of our model in an additional way, we conducted the
 168 following experiment. We generated native images from contrast using the SegResNet model for
 169 the regression task, generated images using diffusion, and generated images using the TimeResNet
 170 model. Then we transferred 20 pairs consisting of a real native and a generated sample to qualified
 171 radiologists to blindly determine where the generated image is and where the original one is. Of the
 172 20 examples generated by regression, all 20 examples were correctly classified, and the diffusion
 173 examples had the same result. But using the example of our TimeResNet model, the doctor correctly
 174 selected 17 cases out of 20. Next, we decided to conduct an experiment where we provide only one
 175 image and ask the reader to label it as a generated or real. We took 20 new cases and generated
 176 them in a similar way, so that 10 images were generated and 10 were real. In the same experiment
 177 with images generated by diffusion and regression, similar results were obtained and all 20 images
 178 were selected correctly. But in the pictures generated using TimeResNet, 14 of the 20 pictures were
 179 already correctly selected. It is worth mentioning that our 2 doctors in both experiments used a
 180 professional application for viewing medical images - Slicer3D, which allows you to carefully view
 181 the image in all 3 projections simultaneously, the doctors performed the screening in the abdomen
 182 window. The following conclusions can be drawn from the generated images: in those produced by
 183 the regression-based approach, organs such as the kidneys, liver, and aorta are recognizable, but the
 184 model tends to blur fine details and smooth out artifacts present in the original CT scans. Diffusion
 185 is the easiest to guess due to the fact that the image is unnaturally shifted when viewed in a sagittal
 186 projection. But there are difficulties in the pictures obtained using Flow Matching, so with paired
 187 studies you can find some details in the liver and kidneys that simply cannot be seen in native images
 188 due to the specifics of CT scans. Further improvement of this method can be achieved by improving
 189 the quality of the training dataset. You can see from the example that we presented that the difference
 190 between a native and contrast image is not only in terms of organs that differ in intensity, but also
 191 along the border of the human body. This is due to the fact that the person still moves a little during
 192 the CT scan. An example of generation can be found in 10.

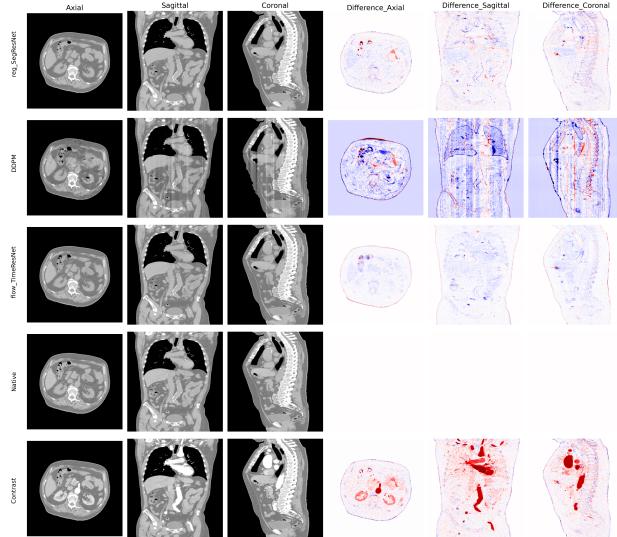


Figure 9: Compare different methods and models for translation from Contrast to Native image. The original Contrast and Native are shown below

193 **4.3 Evaluation Metrics**

194 To evaluate image-to-image translation task we used MAE, although it may not reflect the visual
 195 component, but provided that the picture looks visually correct, it is an excellent metric for comparing
 196 two models.

$$MAE(x, y) = \frac{1}{H \cdot W} \sum_{i=1}^H \sum_{j=1}^W |x_{ij} - y_{ij}|. \quad (6)$$

197 To balance the MAE we also used the SSIM [52]. The SSIM is a commonly used metric for measuring
 198 the structural similarity between two images, it is similar to human perception of image quality. It is
 199 defined as:

$$SSIM(x, y) = \frac{\mu_x \mu_y + C_1}{\mu_x^2 + \mu_y^2 + C_1} \cdot \frac{2 \cdot \sigma_{xy} + C_2}{\sigma_x^2 + \sigma_y^2 + C_2}, \quad (7)$$

200 where μ_x, μ_y and σ_x, σ_y are the means and standard deviations of images x and y respectively, σ_{xy}
 201 means the covariance of x and y . C_1 is $(k_1 L)^2$ and C_2 is $(k_2 L)$, where $k_1 = 0.01$, $k_2 = 0.01$ and
 202 L is the largest pixel of the image x . Also we used the peak signal-to-noise ratio (PSNR) [53] to
 203 evaluate the similarity between the original image and the generated image. PSNR measures the ratio
 204 between the maximum possible value of the original image and the power of distorting noise:

$$PSNR = 20 \cdot \log_{10} \left(\frac{\text{MAX}(image)}{\sqrt{MSE}} \right), \quad (8)$$

205 where $\text{MAX}(x)$ denotes the maximum possible pixel value of the image.

206 4.4 Quantitative Results

207 All values of the MAE metric are specified for images whose values range from -1000 to 1000 HU,
 208 which corresponds to the values of real images. The SSIM values are calculated by translating values
 209 from 0 to 1, the PSNR metric is obtained taking into account that the values take values from 0 to
 210 2000, and the maximum value is 2000. All models were compared by translation based on a one-step
 solution using the Euler method.

Name	Test			Hold-out			Time	Params
	MAE↓	SSIM↑	PSNR↑	MAE↓	SSIM↑	PSNR↑		
UNet	8.132	0.976	37.629	7.051	0.979	39.220	0.015	104.3
SegResNet	7.177	0.979	38.103	6.146	0.983	40.108	0.056	214.1
SwinUNETR	7.465	0.978	38.046	6.485	<u>0.982</u>	39.865	0.087	120.1
UMambaBot	<u>7.220</u>	0.979	38.103	<u>6.201</u>	<u>0.982</u>	<u>39.955</u>	<u>0.077</u>	141.6

Table 1: Comparison of neural network regression performance on test and hold-out datasets using
 MAE, SSIM, PSNR, and inference time in seconds. Best values in **bold**, second-best values are
underlined.

211
 212 The regression-based method produces good metrics for MAE, but you can see that SSIM is signifi-
 213 cantly inferior to the FlowMatching method. This loss can also be observed visually, the image is
 214 blurred.

Split	MAE↓	SSIM ↑	PSNR ↑	Time (seconds) ↓
Test	17.822	0.935	30.192	44.2
Hold-out	18.203	0.934	30.051	44.2

Table 2: Evaluate DDPM performance using MAE, SSIM, and PSNR on test and hold-out datasets.

215 The diffusion generation method produces good metrics, but it lags far behind other approaches.
 216 For the task of translating Contrast to Native, the best metrics are provided by the SwinUNETR
 217 architecture. To compare different architectures, we used the simplest method of solving the equation
 218 - the one-step Euler method. We conducted a more thorough analysis in terms of selecting a solver
 219 and choosing the number of steps, using the TimeResNet example, you can find the results in B.
 220 Careful selection of the solver leads to a significant improvement in the metrics of our proposed
 221 model compared to the values of the architecture metrics in the table.

Name	Test			Hold-out			Time	Params
	MAE \downarrow	SSIM \uparrow	PSNR \uparrow	MAE \downarrow	SSIM \uparrow	PSNR \uparrow	Seconds \downarrow	Millions \downarrow
UMambaBot	<u>6.207</u>	0.992	38.623	6.171	0.992	38.545	<u>0.077</u>	141.6
SegResNet	6.326	0.991	<u>38.783</u>	6.494	0.991	<u>38.572</u>	0.056	214.1
SwinUNETR	6.235	0.992	38.942	<u>6.229</u>	0.992	38.933	0.087	120.1
TimeResNet (Our)	6.513	0.990	38.017	6.363	0.991	38.182	0.105	<u>124.7</u>

Table 3: Comparison of neural network Flow Matching performance (from Contrast to Native) using MAE, SSIM, PSNR, and inference time in seconds on test and hold-out datasets. Best values in **bold**, second-best values are underlined.

222 4.5 Results on subtask segmentation

223 In order to test the possibility of training real models on our generated images, we conducted the
 224 following experiment. We took a hold-out dataset that was not used in training the model, marked the
 225 aorta on Native images, and marked the aorta on contrasting images. Next, the contrasting images
 226 were converted to native images using the proposed model in the article. We used the nnUNetV2
 227 framework to train the segmentation model. We trained on all 20 images, in order to validate our
 228 models, we marked the aorta on an additional 20 native images, on which our models failed. Dice
 229 of the model that was trained on real Native images = 0.966, Dice of the model that was trained on
 230 translated Native images = 0.926. These results confirm the huge opportunity for training models on
 231 the generated ones, they show metrics comparable with real images.

232 4.6 Discussion and Limitations

233 Our proposed method and models yield substantial improvements in CT image translation—thereby
 234 expanding the available training data and markedly enhancing performance across diverse do-
 235 mains—they also serve effectively as a data-augmentation tool for training robust models capable
 236 of operating in dual-domain settings. It is imperative to define the limits of applicability and to
 237 caution that radiologists should not draw clinical conclusions based solely on a single series of
 238 images synthesized by these models, whether generating contrast-enhanced scans from native data
 239 or vice versa. Our future work will focus on addressing challenges associated with extending these
 240 techniques to volumetric imaging and on reducing the requisite GPU memory footprint.

241 5 Conclusion

242 In this paper, we propose a novel method in medical imaging that allows us to solve two tasks at
 243 once, translation from Contrast to Native and translation from Native to Contrast, which reduces
 244 the necessary computational resources to train a model for these tasks, and also demonstrates an
 245 unprecedented generation rate compared to diffusion models. We offer an architecture that reliably
 246 solves these problems and generate stable images. These developments are expected to enhance the
 247 potential accessibility of data across diverse modalities for downstream tasks, as well as to support
 248 additional augmentation of medical datasets for foundation models.

249 6 Data availability

250 Due to confidentiality, data collected for the study are not publicly available for download, however
 251 the corresponding authors can be contacted for academic purposes. Tools for deep learning are
 252 indicated in the methods section.

253 References

- 254 [1] Scott Mayer McKinney, Marcin Sieniek, Varun Godbole, Jonathan Godwin, Natasha Antropova,
 255 Hutan Ashrafian, Trevor Back, Mary Chesus, Greg S Corrado, Ara Darzi, et al. International
 256 evaluation of an ai system for breast cancer screening. *Nature*, 577(7788):89–94, 2020.

- 257 [2] Diego Ardila, Atilla P Kiraly, Sujeeth Bharadwaj, Bokyung Choi, Joshua J Reicher, Lily
258 Peng, Daniel Tse, Mozziyar Etemadi, Wenxing Ye, Greg Corrado, et al. End-to-end lung
259 cancer screening with three-dimensional deep learning on low-dose chest computed tomography.
260 *Nature medicine*, 25(6):954–961, 2019.
- 261 [3] Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and
262 Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks.
263 *nature*, 542(7639):115–118, 2017.
- 264 [4] Varun Gulshan, Lily Peng, Marc Coram, Martin C Stumpe, Derek Wu, Arunachalam
265 Narayanaswamy, Subhashini Venugopalan, Kasumi Widner, Tom Madams, Jorge Cuadros,
266 et al. Development and validation of a deep learning algorithm for detection of diabetic
267 retinopathy in retinal fundus photographs. *jama*, 316(22):2402–2410, 2016.
- 268 [5] Eric J Topol. High-performance medicine: the convergence of human and artificial intelligence.
269 *Nature medicine*, 25(1):44–56, 2019.
- 270 [6] Karin Dembrower, Alessio Crippa, Eugenia Colón, Martin Eklund, and Fredrik Strand. Artificial
271 intelligence for breast cancer detection in screening mammography in sweden: a prospective,
272 population-based, paired-reader, non-inferiority study. *The Lancet Digital Health*, 5(10):e703–
273 e711, 2023.
- 274 [7] Rebecca Smith-Bindman, Philip W. Chu, Hana Azman Firdaus, Carly Stewart, Matthew
275 Malekhedayat, Susan Alber, Wesley E. Bolch, Malini Mahendra, Amy Berrington de González,
276 and Diana L. Miglioretti. Projected lifetime cancer risks from current computed tomography
277 imaging. *JAMA Internal Medicine*, 04 2025. ISSN 2168-6106. doi: 10.1001/jamainternmed.
278 2025.0505. URL <https://doi.org/10.1001/jamainternmed.2025.0505>.
- 279 [8] Mayur K Virarkar, Sai Swarupa R Vulasala, Anjali Verma Gupta, DheerajReddy Gopireddy,
280 Sindhu Kumar, Mauricio Hernandez, Chandana Lall, and Priya Bhosale. Virtual non-contrast
281 imaging in the abdomen and the pelvis: An overview. *Seminars in Ultrasound, CT and MRI*,
282 43(4):293–310, 2022. ISSN 0887-2171. doi: <https://doi.org/10.1053/j.sult.2022.03.004>. URL
283 <https://www.sciencedirect.com/science/article/pii/S0887217122000312>. Dual
284 source CT: Applications, Technology.
- 285 [9] M Jorge Cardoso, Wenqi Li, Richard Brown, Nic Ma, Eric Kerfoot, Yiheng Wang, Benjamin
286 Murrey, Andriy Myronenko, Can Zhao, Dong Yang, et al. Monai: An open-source framework
287 for deep learning in healthcare. *arXiv preprint arXiv:2211.02701*, 2022.
- 288 [10] T.F. Cootes, G.J. Edwards, and C.J. Taylor. Active appearance models. *IEEE Transactions on
289 Pattern Analysis and Machine Intelligence*, 23(6):681–685, 2001. doi: 10.1109/34.927467.
- 290 [11] Berengère Aubert-Broche, Mark Griffin, G Bruce Pike, Alan C Evans, and D Louis Collins.
291 Twenty new digital brain phantoms for creation of validation image data bases. *IEEE transactions on medical imaging*, 25(11):1410–1416, 2006.
- 292 [12] Marcel Prastawa, Elizabeth Bullitt, and Guido Gerig. Simulation of brain tumors in mr
293 images for evaluation of segmentation efficacy. *Medical Image Analysis*, 13(2):297–311,
294 2009. ISSN 1361-8415. doi: <https://doi.org/10.1016/j.media.2008.11.002>. URL
295 <https://www.sciencedirect.com/science/article/pii/S1361841508001357>. Includes Spe-
296 cial Section on Functional Imaging and Modelling of the Heart.
- 297 [13] W Paul Segars, G Sturgeon, S Mendonca, Jason Grimes, and Benjamin MW Tsui. 4d xcat
phantom for multimodality imaging research. *Medical physics*, 37(9):4902–4915, 2010.
- 298 [14] Sébastien Jan, Giovanni Santin, Daniel Strul, Steven Staelens, K Assié, Damien Autret, Stéphane
299 Avner, Remi Barbier, Manuel Bardès, PM Bloomfield, et al. Gate: a simulation toolkit for pet
and spect. *Physics in Medicine & Biology*, 49(19):4543, 2004.
- 300 [15] Shengye Hu, Baiying Lei, Shuqiang Wang, Yong Wang, Zhiguang Feng, and Yanyan Shen.
301 Bidirectional mapping generative adversarial networks for brain mr to pet synthesis. *IEEE
302 Transactions on Medical Imaging*, 41(1):145–157, 2021.

- 306 [16] Wen Wei, Emilie Poirion, Benedetta Bodini, Stanley Durrleman, Nicholas Ayache, Bruno
307 Stankoff, and Olivier Colliot. Predicting pet-derived demyelination from multimodal mri
308 using sketcher-refiner adversarial training for multiple sclerosis. *Medical Image Analysis*, 58:
309 101546, 2019. ISSN 1361-8415. doi: <https://doi.org/10.1016/j.media.2019.101546>. URL
310 <https://www.sciencedirect.com/science/article/pii/S1361841519300817>.
- 311 [17] Haoyu Lan, the Alzheimer Disease Neuroimaging Initiative, Arthur W. Toga, and Farshid
312 Sepehrband. Three-dimensional self-attention conditional gan with spectral normalization for
313 multimodal neuroimaging synthesis. *Magnetic Resonance in Medicine*, 86(3):1718–1733, 2021.
314 doi: <https://doi.org/10.1002/mrm.28819>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/mrm.28819>.
- 316 [18] Wanyun Lin, Weiming Lin, Gang Chen, Hejun Zhang, Qinquan Gao, Yechong Huang, Tong
317 Tong, and Min Du. Bidirectional mapping of brain mri and pet with 3d reversible gan for
318 the diagnosis of alzheimer’s disease. *Frontiers in Neuroscience*, 15, 2021. URL <https://api.semanticscholar.org/CorpusID:233186986>.
- 320 [19] Apoorva Sikka, Jitender Singh Virk, Deepti R Bathula, et al. Mri to pet cross-modality transla-
321 tion using globally and locally aware gan (gla-gan) for multi-modal diagnosis of alzheimer’s
322 disease. *arXiv preprint arXiv:2108.02160*, 2021.
- 323 [20] Jin Zhang, Xiaohai He, Linbo Qing, Feng Gao, and Bin Wang. Bpgan: Brain pet synthesis
324 from mri using generative adversarial network for multi-modal alzheimer’s disease diagno-
325 sis. *Computer Methods and Programs in Biomedicine*, 217:106676, 2022. ISSN 0169-2607.
326 doi: <https://doi.org/10.1016/j.cmpb.2022.106676>. URL <https://www.sciencedirect.com/science/article/pii/S016926072200061X>.
- 328 [21] Farideh Bazangani, Frédéric J. P. Richard, Badih Ghattas, and Eric Guedj. Fdg-pet to t1 weighted
329 mri translation with 3d elicit generative adversarial network (e-gan). *Sensors*, 22(12), 2022.
330 ISSN 1424-8220. doi: 10.3390/s22124640. URL <https://www.mdpi.com/1424-8220/22/12/4640>.
- 332 [22] Yan Wang, Yanmei Luo, Chen Zu, Bo Zhan, Zhengyang Jiao, Xi Wu, Jiliu Zhou, Dinggang
333 Shen, and Luping Zhou. 3d multi-modality transformer-gan for high-quality pet reconstruction.
334 *Medical Image Analysis*, 91:102983, 2024. ISSN 1361-8415. doi: <https://doi.org/10.1016/j.media.2023.102983>. URL <https://www.sciencedirect.com/science/article/pii/S1361841523002438>.
- 337 [23] Cheng Peng, Pengfei Guo, S Kevin Zhou, Vishal M Patel, and Rama Chellappa. Towards
338 performant and reliable undersampled mr reconstruction via diffusion model sampling. In
339 *International Conference on Medical Image Computing and Computer-Assisted Intervention*,
340 pages 623–633. Springer, 2022.
- 341 [24] Yutong Xie and Quanzheng Li. Measurement-conditioned denoising diffusion probabilistic
342 model for under-sampled medical image reconstruction. In *International Conference on Medical
343 Image Computing and Computer-Assisted Intervention*, pages 655–664. Springer, 2022.
- 344 [25] Chi-Hieu Pham, Carlos Tor-Díez, Hélène Meunier, Nathalie Bednarek, Ronan Fablet, Nicolas
345 Passat, and François Rousseau. Multiscale brain mri super-resolution using deep 3d convolu-
346 tional networks. *Computerized Medical Imaging and Graphics*, 77:101647, 2019.
- 347 [26] Cycle Learning Ensemble. Ct super-resolution gan constrained by the identical, residual, and
348 cycle learning ensemble.
- 349 [27] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
350 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF
351 conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- 352 [28] Alicia Durrer, Julia Wolleb, Florentin Bieder, Tim Sinnecker, Matthias Weigel, Robin Sandküh-
353 ler, Cristina Granziera, Özgür Yaldizli, and Philippe C Cattin. Diffusion models for contrast
354 harmonization of magnetic resonance images. *arXiv preprint arXiv:2303.08189*, 2023.

- 355 [29] Robert Graf, Joachim Schmitt, Sarah Schlaeger, Hendrik Kristian Möller, Vasiliki Sideri-
 356 Lampretsa, Anjany Sekuboyina, Sandro Manuel Krieg, Benedikt Wiestler, Bjoern Menze,
 357 Daniel Rueckert, et al. Denoising diffusion-based mri to ct image translation enables automated
 358 spinal segmentation. *European Radiology Experimental*, 7(1):70, 2023.
- 359 [30] Shaoyan Pan, Zach Eidex, Mojtaba Safari, Richard Qiu, and Xiaofeng Yang. Cycle-guided
 360 denoising diffusion probability model for 3d cross-modality mri synthesis. In *Medical Imaging*
 361 *2025: Clinical and Biomedical Imaging*, volume 13410, pages 515–522. SPIE, 2025.
- 362 [31] Lingting Zhu, Zeyue Xue, Zhenchao Jin, Xian Liu, Jingzhen He, Ziwei Liu, and Lequan Yu.
 363 Make-a-volume: Leveraging latent diffusion models for cross-modality 3d brain mri synthesis.
 364 In *International Conference on Medical Image Computing and Computer-Assisted Intervention*,
 365 pages 592–601. Springer, 2023.
- 366 [32] Reza Kalantar, Sumeet Hindocha, Benjamin Hunter, Bhupinder Sharma, Nasir Khan, Dow-Mu
 367 Koh, Merina Ahmed, Eric Aboagye, Richard Lee, and Matthew Blackledge. Non-contrast ct
 368 synthesis using patch-based cycle-consistent generative adversarial network (cycle-gan) for
 369 radiomics and deep learning in the era of covid-19. *Scientific Reports*, 13, 06 2023. doi:
 370 10.1038/s41598-023-36712-1.
- 371 [33] Jonghun Kim and Hyunjin Park. Adaptive latent diffusion model for 3d medical image to image
 372 translation: Multi-modal magnetic resonance imaging study. In *Proceedings of the IEEE/CVF*
 373 *Winter Conference on Applications of Computer Vision*, pages 7604–7613, 2024.
- 374 [34] Lingting Zhu, Noel Codella, Dongdong Chen, Zhenchao Jin, Lu Yuan, and Lequan Yu. Genera-
 375 tive enhancement for 3d medical images. *arXiv preprint arXiv:2403.12852*, 2024.
- 376 [35] Zolnamar Dorjsembe, Hsing-Kuo Pao, Sodtavilan Odonchimed, and Furen Xiao. Conditional
 377 diffusion models for semantic 3d brain mri synthesis. *IEEE Journal of Biomedical and Health*
 378 *Informatics*, 2024.
- 379 [36] Shaoyan Pan, Elham Abouei, Jacob Wynne, Chih-Wei Chang, Tonghe Wang, Richard LJ Qiu,
 380 Yuheng Li, Junbo Peng, Justin Roper, Pretesh Patel, et al. Synthetic ct generation from mri
 381 using 3d transformer-based denoising diffusion model. *Medical Physics*, 51(4):2538–2548,
 382 2024.
- 383 [37] Yitong Li, Igor Yakushev, Dennis M Hedderich, and Christian Wachinger. Pasta: P athology-a
 384 ware mri to pet cro s s-modal t r a nslation with diffusion models. In *International Conference*
 385 *on Medical Image Computing and Computer-Assisted Intervention*, pages 529–540. Springer,
 386 2024.
- 387 [38] Pengfei Guo, Can Zhao, Dong Yang, Ziyue Xu, Vishwesh Nath, Yucheng Tang, Benjamin
 388 Simon, Mason Belue, Stephanie Harmon, Baris Turkbey, et al. Maisi: Medical ai for synthetic
 389 imaging. In *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*,
 390 pages 4430–4441. IEEE, 2025.
- 391 [39] Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic
 392 interpolants. *arXiv preprint arXiv:2209.15571*, 2022.
- 393 [40] Marin Biloš, Johanna Sommer, Syama Sundar Rangapuram, Tim Januschowski, and Stephan
 394 Günnemann. Neural flows: Efficient alternative to neural odes. *Advances in neural information*
 395 *processing systems*, 34:21325–21337, 2021.
- 396 [41] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow
 397 matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- 398 [42] Alexander Tong, Kilian Fatras, Nikolay Malkin, Guillaume Huguet, Yanlei Zhang, Jarrid Rector-
 399 Brooks, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative
 400 models with minibatch optimal transport. *arXiv preprint arXiv:2302.00482*, 2023.
- 401 [43] Milad Yazdani, Yasamin Medghalchi, Pooria Ashrafian, Ilker Hacihaliloglu, and Dena Shahriari.
 402 Flow matching for medical image synthesis: Bridging the gap between speed and quality. *arXiv*
 403 *preprint arXiv:2503.00266*, 2025.

- 404 [44] Hadrien Reynaud, Alberto Gomez, Paul Leeson, Qingjie Meng, and Bernhard Kainz. Echoflow:
 405 A foundation model for cardiac ultrasound image and video generation. *arXiv preprint*
 406 *arXiv:2503.22357*, 2025.
- 407 [45] Wenzuan Li, Chongyu Qu, Xiaoxi Chen, Pedro RAS Bassi, Yijia Shi, Yuxiang Lai, Qian Yu,
 408 Huimin Xue, Yixiong Chen, Xiaorui Lin, et al. Abdomenatlas: A large-scale, detailed-annotated,
 409 & multi-center dataset for efficient transfer learning and open algorithmic benchmarking.
 410 *Medical Image Analysis*, 97:103285, 2024.
- 411 [46] Jungye Kim, Jimin Lee, Bitbyeo Kim, Sangwook Kim, Hyeongmin Jin, and Seongmoon Jung.
 412 Generation of deep learning based virtual non-contrast ct using dual-layer dual-energy ct and
 413 its application to planning ct for radiotherapy. *PLOS ONE*, 19, 12 2024. doi: 10.1371/journal.
 414 *pone.0316099*.
- 415 [47] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*
 416 *in neural information processing systems*, 33:6840–6851, 2020.
- 417 [48] Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning
 418 Wang. Swin-unet: Unet-like pure transformer for medical image segmentation. In *European*
 419 *conference on computer vision*, pages 205–218. Springer, 2022.
- 420 [49] Baosong Yang, Longyue Wang, Derek F. Wong, Lidia S. Chao, and Zhaopeng Tu. Convolutional
 421 self-attention networks. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for*
 422 *Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*,
 423 pages 4040–4045, Minneapolis, Minnesota, June 2019. Association for Computational
 424 Linguistics. doi: 10.18653/v1/N19-1407. URL <https://aclanthology.org/N19-1407/>.
- 426 [50] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference*
 427 *on computer vision (ECCV)*, pages 3–19, 2018.
- 428 [51] A Paszke. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint*
 429 *arXiv:1912.01703*, 2019.
- 430 [52] Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from
 431 error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612,
 432 2004. doi: 10.1109/TIP.2003.819861.
- 433 [53] Jari Korhonen and Junyoung You. Peak signal-to-noise ratio revisited: Is simple beautiful? In
 434 *2012 Fourth international workshop on quality of multimedia experience*, pages 37–38. IEEE,
 435 2012.
- 436 [54] L. Euler. *Institutionum calculi integralis*. Impensis Academiae Imperialis Scientiarum, 1769.
 437 URL <https://books.google.ru/books?id=XmsCrgEACAAJ>.
- 438 [55] Mark Lotkin. A note on the midpoint method of integration. *J. ACM*, 3(3):208–211, July 1956.
 439 ISSN 0004-5411. doi: 10.1145/320831.320840. URL <https://doi.org/10.1145/320831.320840>.
- 441 [56] C. Runge. Ueber die numerische auflösung von differentialgleichungen. *Mathematische*
 442 *Annalen*, 46:167–178, 1895. URL <http://eudml.org/doc/157756>.
- 443 [57] W. Kutta. *Beitrag zur näherungsweisen Integration totaler Differentialgleichungen*. Teubner,
 444 1901. URL <https://books.google.ru/books?id=K5e6kQEACAAJ>.
- 445 [58] Francis Bashforth and John Couch Adams. *An Attempt to Test the Theories of Capillary Action by*
 446 *Comparing the Theoretical and Measured Forms of Drops of Fluid*. Cambridge University Press,
 447 Cambridge, 1883. URL <https://archive.org/details/attempttest00bashrich>.
- 448 [59] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 449 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
 450 *processing systems*, 30, 2017.

451 The additional material is organized as follows: section A will contain details about metrics and
 452 methods of model inference on a test sample, section B will provide details about the selection
 453 of hyperparameters in training, section C will provide visual diagrams of training and operation
 454 of various parts of the neural network, and section D will provide additional examples of image
 455 generation.

456 A Additional quantitative results

457 A.1 Details

458 We measured the ability to convert images with contrast to native images using the models we trained
 459 to test their ability to do this. They showed good results. This gives us the opportunity to use the
 460 same neural network, but for two types of tasks, it allows us to reduce the time for training two neural
 networks in the case of UNet Regression, DDPM.

Name	Test			Hold-out			Time (s)	Params
	MAE↓	SSIM↑	PSNR↑	MAE↓	SSIM↑	PSNR↑	Seconds↓	Millions ↓
UMambaBot	<u>6.996</u>	0.989	<u>39.095</u>	<u>6.882</u>	<u>0.989</u>	<u>38.933</u>	<u>0.077</u>	141.1
SegResNet	6.431	0.989	39.544	6.426	0.990	39.369	0.056	214.1
SwinUNETR	7.593	0.987	38.692	7.604	0.987	38.524	0.087	120.1
TimeResNet (Our)	7.329	0.988	38.947	7.253	0.988	38.777	0.105	<u>124.7</u>

Table 4: Comparison of neural network Flow Matching performance (from Native to Contrast) using MAE, SSIM, PSNR, and inference time in seconds on test and hold-out datasets. Best values in **bold**, second-best values are underlined.

461
 462 In the task of translating from Native to Contrast, the SegResNet architecture wins by a large margin.

463 A.2 ODE-Solver Ablation

464 Table A.2 benchmarks five numerical solvers (Euler [54], Midpoint [55], RK2 [56], RK3, RK4 [57],
 465 Adams [58]) Exponential Integrator) for the best FM model over step budgets {1, 3, 5, 10}.

466 We can numerically solve the neural ordinary differential equation using various methods to figure out
 467 which method works best and what is the optimal number of steps to choose. With the proper selection
 468 of a solver to solve the problem, our model wins by a significant margin over other architectures.

Solver	Test			Hold-out			Time	Complexity
	MAE↓	SSIM↑	PSNR↑	MAE↓	SSIM↑	PSNR↑	Seconds ↓	↓
Euler (1)	6.513	0.990	38.017	6.362	0.991	38.182	0.105	1
Euler (2)	<u>5.591</u>	<u>0.995</u>	<u>39.039</u>	<u>5.525</u>	<u>0.995</u>	<u>39.017</u>	<u>0.210</u>	2
Euler (3)	<u>5.529</u>	0.996	38.941	<u>5.472</u>	0.995	38.994	0.313	3
Euler (5)	5.565	0.995	38.899	5.526	0.995	38.951	0.521	5
RK2 (1)	5.399	0.996	40.017	5.436	0.996	39.776	0.209	2
RK3 (1)	5.553	0.995	39.193	5.553	0.995	39.122	0.313	3
RK4 (1)	5.674	0.995	<u>39.408</u>	5.716	0.995	<u>39.257</u>	0.419	4
Mid Point (1)	5.601	0.995	39.106	5.611	0.995	39.047	<u>0.210</u>	2
Mid Point (3)	5.744	0.995	38.882	5.759	0.995	38.923	0.628	6

Table 5: Comparison of different solvers (from Contrast to Native) for TimeResNet (our) using MAE, SSIM, PSNR, and inference time in seconds on test and hold-out datasets. Best values in **bold**, second-best values are underlined.

469 By complexity here, we mean the number of calls to the neural network to calculate the gradient
 470 when generating a single image.

471 **B Additional implementation details**

472 **B.1 2D Processing Justification**

473 Due to the limited computing resources available, we decided to solve this problem in 2D, since 3D
474 requires much more computing resources and a lot of data, but a full-fledged volume does not fit
475 on one GPU, so it would require a latent-space solution, which still demands significant resources.
476 Moreover, clinicians typically review studies in the axial plane (though not always), making 2D axial
477 slices a natural choice. We therefore restrict our method to 2D axial processing.

478 **B.2 Image Preprocessing and Resizing**

479 When down-sampling inputs, one can either resample to a uniform spacing or to a fixed pixel size.
480 We chose a fixed image size of 512×512 , since nearly all DICOM slices in our dataset already
481 conform to this dimension. Avoiding on-the-fly resizing during inference eliminates an extra source of
482 error. Additionally, variations in original spacing act as implicit data augmentation, improving model
483 robustness. We acknowledge that 2D methods may lack inter-slice context and can yield inconsistent
484 predictions across different planes (e.g., axial vs. sagittal), as seen in diffusion experiments; however,
485 our approach demonstrates stable and consistent contrast-to-native predictions slice-by-slice within
486 the same study.

487 **B.3 Data Normalization and Clipping**

488 To harmonize input intensities across scanners, we clip CT Hounsfield Units to $[-1000, 1000]$, then
489 linearly rescale to $[-1, 1]$. This range covers the vast majority of soft-tissue contrasts while ensuring
490 numerical stability during training.

491 **B.4 Data Augmentation and Optimization**

492 We apply geometric augmentations via the MONAI framework:

- 493
 - Random flips along x and y axes (50% each),
 - Random 90° rotations (up to 3 turns, 50% probability),
 - Random affine transforms (70% probability) with: rotation up to $\pm 45^\circ$, translation up to
496 ± 102.4 pixels, shear up to $\pm 10^\circ$.

497 We train using Adam with hyperparameters: learning rate = $2 \cdot 10^{-5}$, weight decay = 0, $\beta_1 =$
498 0.9, $\beta_2 = 0.999$. Batch size was 2.

499 **B.5 Training Data Selection and Dual-Energy Exclusion**

500 To minimize misregistration due to patient motion between acquisitions, we remove series where we
501 have different origins and spacing between two series. Although dual-energy CT could in principle
502 provide “native” images directly, such data are scarce, and the vendor-reconstructed “native” series
503 result from proprietary algorithms that may not faithfully represent true tissue contrast. Therefore,
504 we exclude dual-energy studies to ensure the network learns genuine native-to-contrast mappings.

505 **B.6 Image Registration**

506 In order to improve the quality and purity of the data, we registered images with contrast to native
507 images. We employ the ANTs library with the “deformable SyN only” transform, providing a
508 body-mask to focus optimization on the anatomy. This configuration yielded the best mean absolute
509 error (MAE) and mutual information (MI) improvements, and was faster than alternative ANTs
510 transforms.

511 **C Additional visual results and discussion**

512 The generation examples are presented in 3 projections, and each of the projections shows the
513 difference between the original native image. The picture was clipped into the abdominal window for

514 a clearer understanding of the difference, since in other windows or without an image clip at all, the difference is much less noticeable, and we are primarily interested in the abdominal window.

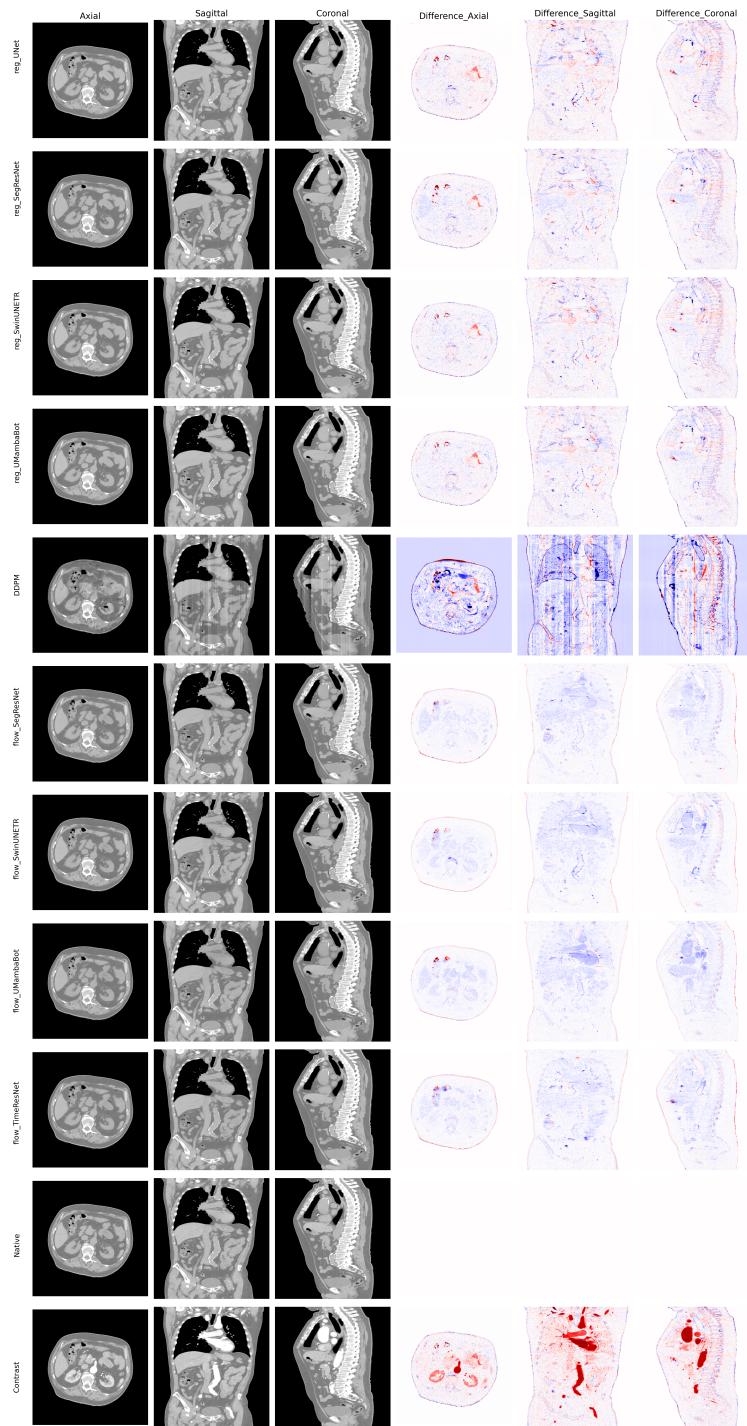


Figure 10: Compare different methods and models for translation from Contrast to Native image. The original Contrast and Native are shown below

516 **D Architecture details**

517 **D.1 Time embedding mechanism**

518 One nontrivial challenge is how to encode temporal information so that the model fully captures
 519 context and can accurately predict the vector field between two images at a given time point. We
 520 evaluated several strategies and found the following approach to be most effective, yielding substantial
 521 performance gains. Because the model tends to “forget” the temporal signal, we inject temporal
 522 embeddings into every block of the encoder. In addition, we introduce a lightweight linear projection
 523 layer in each block, enabling it to learn the specific dependencies it requires without unduly increasing
 524 the overall parameter count. This linear layer operates analogously to an instance-normalization
 525 layer, allowing independent modulation of each channel’s activations. The temporal embeddings
 526 themselves are constructed using the sinusoidal scheme described in [59]. To capture nonlinear
 527 temporal interactions, we further process these sinusoidal embeddings through a small multilayer
 528 perceptron, and the resulting feature vector is added to the input of each ResNet block. In the
 529 schematic below A, B denote trainable vectors of length C , which are applied element-wise to the
 530 output of the SiLU activation. For convenience, one linear layer of length $2 \cdot C$ is trained, which is
 531 then divided into two vectors – A, B . If the output of the nonlinear layer, h , has dimensions $(B \times C \times$
 532 $H \times W)$, then the resulting tensor likewise has dimensions $(B \times C \times H \times W)$, with the multiplication
 533 broadcast along the channel dimension.

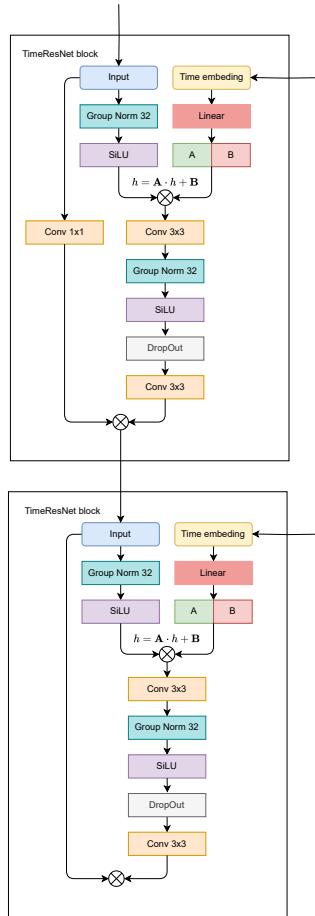


Figure 11: Time embedding mechanism.

534 **NeurIPS Paper Checklist**

535 **1. Claims**

536 Question: Do the main claims made in the abstract and introduction accurately reflect the
537 paper's contributions and scope?

538 Answer: [Yes]

539 Justification: The main claims made in the abstract and introduction accurately reflect the
540 paper's contributions and scope.

541 Guidelines:

- 542 • The answer NA means that the abstract and introduction do not include the claims
543 made in the paper.
- 544 • The abstract and/or introduction should clearly state the claims made, including the
545 contributions made in the paper and important assumptions and limitations. A No or
546 NA answer to this question will not be perceived well by the reviewers.
- 547 • The claims made should match theoretical and experimental results, and reflect how
548 much the results can be expected to generalize to other settings.
- 549 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
550 are not attained by the paper.

551 **2. Limitations**

552 Question: Does the paper discuss the limitations of the work performed by the authors?

553 Answer: [Yes]

554 Justification: We have discussed the limitations of the work in Section 4.6

555 Guidelines:

- 556 • The answer NA means that the paper has no limitation while the answer No means that
557 the paper has limitations, but those are not discussed in the paper.
- 558 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 559 • The paper should point out any strong assumptions and how robust the results are to
560 violations of these assumptions (e.g., independence assumptions, noiseless settings,
561 model well-specification, asymptotic approximations only holding locally). The authors
562 should reflect on how these assumptions might be violated in practice and what the
563 implications would be.
- 564 • The authors should reflect on the scope of the claims made, e.g., if the approach was
565 only tested on a few datasets or with a few runs. In general, empirical results often
566 depend on implicit assumptions, which should be articulated.
- 567 • The authors should reflect on the factors that influence the performance of the approach.
568 For example, a facial recognition algorithm may perform poorly when image resolution
569 is low or images are taken in low lighting. Or a speech-to-text system might not be
570 used reliably to provide closed captions for online lectures because it fails to handle
571 technical jargon.
- 572 • The authors should discuss the computational efficiency of the proposed algorithms
573 and how they scale with dataset size.
- 574 • If applicable, the authors should discuss possible limitations of their approach to
575 address problems of privacy and fairness.
- 576 • While the authors might fear that complete honesty about limitations might be used by
577 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
578 limitations that aren't acknowledged in the paper. The authors should use their best
579 judgment and recognize that individual actions in favor of transparency play an impor-
580 tant role in developing norms that preserve the integrity of the community. Reviewers
581 will be specifically instructed to not penalize honesty concerning limitations.

582 **3. Theory Assumptions and Proofs**

583 Question: For each theoretical result, does the paper provide the full set of assumptions and
584 a complete (and correct) proof?

585 Answer: [NA]

586 Justification: The paper does not include theoretical results.

587 Guidelines:

- 588 • The answer NA means that the paper does not include theoretical results.
- 589 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 590 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 591 • The proofs can either appear in the main paper or the supplemental material, but if
- 592 they appear in the supplemental material, the authors are encouraged to provide a short
- 593 proof sketch to provide intuition.
- 594 • Inversely, any informal proof provided in the core of the paper should be complemented
- 595 by formal proofs provided in appendix or supplemental material.
- 596 • Theorems and Lemmas that the proof relies upon should be properly referenced.

597 4. Experimental Result Reproducibility

599 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
600 perimental results of the paper to the extent that it affects the main claims and/or conclusions
601 of the paper (regardless of whether the code and data are provided or not)?

602 Answer: [Yes]

603 Justification: All details of the proposed method are included in the paper.

604 Guidelines:

- 605 • The answer NA means that the paper does not include experiments.
- 606 • If the paper includes experiments, a No answer to this question will not be perceived
- 607 well by the reviewers: Making the paper reproducible is important, regardless of
- 608 whether the code and data are provided or not.
- 609 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 610 to make their results reproducible or verifiable.
- 611 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 612 For example, if the contribution is a novel architecture, describing the architecture fully
- 613 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 614 be necessary to either make it possible for others to replicate the model with the same
- 615 dataset, or provide access to the model. In general, releasing code and data is often
- 616 one good way to accomplish this, but reproducibility can also be provided via detailed
- 617 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 618 of a large language model), releasing of a model checkpoint, or other means that are
- 619 appropriate to the research performed.
- 620 • While NeurIPS does not require releasing code, the conference does require all submis-
- 621 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 622 nature of the contribution. For example
 - 623 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
 - 624 to reproduce that algorithm.
 - 625 (b) If the contribution is primarily a new model architecture, the paper should describe
 - 626 the architecture clearly and fully.
 - 627 (c) If the contribution is a new model (e.g., a large language model), then there should
 - 628 either be a way to access this model for reproducing the results or a way to reproduce
 - 629 the model (e.g., with an open-source dataset or instructions for how to construct
 - 630 the dataset).
 - 631 (d) We recognize that reproducibility may be tricky in some cases, in which case
 - 632 authors are welcome to describe the particular way they provide for reproducibility.
 - 633 In the case of closed-source models, it may be that access to the model is limited in
 - 634 some way (e.g., to registered users), but it should be possible for other researchers
 - 635 to have some path to reproducing or verifying the results.

636 5. Open access to data and code

637 Question: Does the paper provide open access to the data and code, with sufficient instruc-

638 tions to faithfully reproduce the main experimental results, as described in supplemental

639 material?

640 Answer: [Yes]

641 Justification: All the code and model weights will be released upon paper acceptance, data
642 will be available on request.

643 Guidelines:

- 644 • The answer NA means that paper does not include experiments requiring code.
- 645 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 646 • While we encourage the release of code and data, we understand that this might not be
647 possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
648 including code, unless this is central to the contribution (e.g., for a new open-source
649 benchmark).
- 650 • The instructions should contain the exact command and environment needed to run to
651 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 652 • The authors should provide instructions on data access and preparation, including how
653 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 654 • The authors should provide scripts to reproduce all experimental results for the new
655 proposed method and baselines. If only a subset of experiments are reproducible, they
656 should state which ones are omitted from the script and why.
- 657 • At submission time, to preserve anonymity, the authors should release anonymized
658 versions (if applicable).
- 659 • Providing as much information as possible in supplemental material (appended to the
660 paper) is recommended, but including URLs to data and code is permitted.

663 6. Experimental Setting/Details

664 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
665 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
666 results?

667 Answer: [Yes]

668 Justification: All details included in the paper.

669 Guidelines:

- 670 • The answer NA means that the paper does not include experiments.
- 671 • The experimental setting should be presented in the core of the paper to a level of detail
672 that is necessary to appreciate the results and make sense of them.
- 673 • The full details can be provided either with the code, in appendix, or as supplemental
674 material.

675 7. Experiment Statistical Significance

676 Question: Does the paper report error bars suitably and correctly defined or other appropriate
677 information about the statistical significance of the experiments?

678 Answer: [No]

679 Justification: NA

680 Guidelines:

- 681 • The answer NA means that the paper does not include experiments.
- 682 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
683 dence intervals, or statistical significance tests, at least for the experiments that support
684 the main claims of the paper.
- 685 • The factors of variability that the error bars are capturing should be clearly stated (for
686 example, train/test split, initialization, random drawing of some parameter, or overall
687 run with given experimental conditions).
- 688 • The method for calculating the error bars should be explained (closed form formula,
689 call to a library function, bootstrap, etc.)
- 690 • The assumptions made should be given (e.g., Normally distributed errors).

- 691 • It should be clear whether the error bar is the standard deviation or the standard error
692 of the mean.
693 • It is OK to report 1-sigma error bars, but one should state it. The authors should
694 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
695 of Normality of errors is not verified.
696 • For asymmetric distributions, the authors should be careful not to show in tables or
697 figures symmetric error bars that would yield results that are out of range (e.g. negative
698 error rates).
699 • If error bars are reported in tables or plots, The authors should explain in the text how
700 they were calculated and reference the corresponding figures or tables in the text.

701 **8. Experiments Compute Resources**

702 Question: For each experiment, does the paper provide sufficient information on the com-
703 puter resources (type of compute workers, memory, time of execution) needed to reproduce
704 the experiments?

705 Answer: [Yes]

706 Justification: All the experiments are conducted on 4 RTX 3090 24G GPUs.

707 Guidelines:

- 708 • The answer NA means that the paper does not include experiments.
709 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
710 or cloud provider, including relevant memory and storage.
711 • The paper should provide the amount of compute required for each of the individual
712 experimental runs as well as estimate the total compute.
713 • The paper should disclose whether the full research project required more compute
714 than the experiments reported in the paper (e.g., preliminary or failed experiments that
715 didn't make it into the paper).

716 **9. Code Of Ethics**

717 Question: Does the research conducted in the paper conform, in every respect, with the
718 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

719 Answer: [Yes]

720 Justification: Research conducted in the paper conforms with the NeurIPS Code of Ethics.

721 Guidelines:

- 722 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
723 • If the authors answer No, they should explain the special circumstances that require a
724 deviation from the Code of Ethics.
725 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
726 eration due to laws or regulations in their jurisdiction).

727 **10. Broader Impacts**

728 Question: Does the paper discuss both potential positive societal impacts and negative
729 societal impacts of the work performed?

730 Answer: [Yes]

731 Justification: We contribute translation model and method, which can benefit numerous
732 clinical study and applications. All limitations discussed in the Section 4.6.

733 Guidelines:

- 734 • The answer NA means that there is no societal impact of the work performed.
735 • If the authors answer NA or No, they should explain why their work has no societal
736 impact or why the paper does not address societal impact.
737 • Examples of negative societal impacts include potential malicious or unintended uses
738 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
739 (e.g., deployment of technologies that could make decisions that unfairly impact specific
740 groups), privacy considerations, and security considerations.

- 741 • The conference expects that many papers will be foundational research and not tied
 742 to particular applications, let alone deployments. However, if there is a direct path to
 743 any negative applications, the authors should point it out. For example, it is legitimate
 744 to point out that an improvement in the quality of generative models could be used to
 745 generate deepfakes for disinformation. On the other hand, it is not needed to point out
 746 that a generic algorithm for optimizing neural networks could enable people to train
 747 models that generate Deepfakes faster.
- 748 • The authors should consider possible harms that could arise when the technology is
 749 being used as intended and functioning correctly, harms that could arise when the
 750 technology is being used as intended but gives incorrect results, and harms following
 751 from (intentional or unintentional) misuse of the technology.
- 752 • If there are negative societal impacts, the authors could also discuss possible mitigation
 753 strategies (e.g., gated release of models, providing defenses in addition to attacks,
 754 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
 755 feedback over time, improving the efficiency and accessibility of ML).

756 11. Safeguards

757 Question: Does the paper describe safeguards that have been put in place for responsible
 758 release of data or models that have a high risk for misuse (e.g., pretrained language models,
 759 image generators, or scraped datasets)?

760 Answer: [No]

761 Justification: No such risks.

762 Guidelines:

- 763 • The answer NA means that the paper poses no such risks.
- 764 • Released models that have a high risk for misuse or dual-use should be released with
 765 necessary safeguards to allow for controlled use of the model, for example by requiring
 766 that users adhere to usage guidelines or restrictions to access the model or implementing
 767 safety filters.
- 768 • Datasets that have been scraped from the Internet could pose safety risks. The authors
 769 should describe how they avoided releasing unsafe images.
- 770 • We recognize that providing effective safeguards is challenging, and many papers do
 771 not require this, but we encourage authors to take this into account and make a best
 772 faith effort.

773 12. Licenses for existing assets

774 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
 775 the paper, properly credited and are the license and terms of use explicitly mentioned and
 776 properly respected?

777 Answer: [Yes]

778 Justification: Yes.

779 Guidelines:

- 780 • The answer NA means that the paper does not use existing assets.
- 781 • The authors should cite the original paper that produced the code package or dataset.
- 782 • The authors should state which version of the asset is used and, if possible, include a
 783 URL.
- 784 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 785 • For scraped data from a particular source (e.g., website), the copyright and terms of
 786 service of that source should be provided.
- 787 • If assets are released, the license, copyright information, and terms of use in the
 788 package should be provided. For popular datasets, paperswithcode.com/datasets
 789 has curated licenses for some datasets. Their licensing guide can help determine the
 790 license of a dataset.
- 791 • For existing datasets that are re-packaged, both the original license and the license of
 792 the derived asset (if it has changed) should be provided.

- 793 • If this information is not available online, the authors are encouraged to reach out to
794 the asset's creators.

795 **13. New Assets**

796 Question: Are new assets introduced in the paper well documented and is the documentation
797 provided alongside the assets?

798 Answer: [Yes]

799 Justification: The trained model will be released after reviewing.

800 Guidelines:

- 801 • The answer NA means that the paper does not release new assets.
802 • Researchers should communicate the details of the dataset/code/model as part of their
803 submissions via structured templates. This includes details about training, license,
804 limitations, etc.
805 • The paper should discuss whether and how consent was obtained from people whose
806 asset is used.
807 • At submission time, remember to anonymize your assets (if applicable). You can either
808 create an anonymized URL or include an anonymized zip file.

809 **14. Crowdsourcing and Research with Human Subjects**

810 Question: For crowdsourcing experiments and research with human subjects, does the paper
811 include the full text of instructions given to participants and screenshots, if applicable, as
812 well as details about compensation (if any)?

813 Answer: [NA]

814 Justification: The paper does not involve crowdsourcing nor research with human subjects

815 Guidelines:

- 816 • The answer NA means that the paper does not involve crowdsourcing nor research with
817 human subjects.
818 • Including this information in the supplemental material is fine, but if the main contribu-
819 tion of the paper involves human subjects, then as much detail as possible should be
820 included in the main paper.
821 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
822 or other labor should be paid at least the minimum wage in the country of the data
823 collector.

824 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
825 Subjects**

826 Question: Does the paper describe potential risks incurred by study participants, whether
827 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
828 approvals (or an equivalent approval/review based on the requirements of your country or
829 institution) were obtained?

830 Answer: [NA]

831 Justification: The paper does not involve crowdsourcing nor research with human subjects.

832 Guidelines:

- 833 • The answer NA means that the paper does not involve crowdsourcing nor research with
834 human subjects.
835 • Depending on the country in which research is conducted, IRB approval (or equivalent)
836 may be required for any human subjects research. If you obtained IRB approval, you
837 should clearly state this in the paper.
838 • We recognize that the procedures for this may vary significantly between institutions
839 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
840 guidelines for their institution.
841 • For initial submissions, do not include any information that would break anonymity (if
842 applicable), such as the institution conducting the review.