# MedM2G: Unifying Medical Multi-Modal Generation via Cross-Guided Diffusion with Visual Invariant

Chenlu Zhan[1,2]    Yu Lin[2]    Gaoang Wang[1,2 (✉)]    Hongwei Wang[1,2(✉)]    Jian Wu[3]

[1] College of Computer Science and Technology, Zhejiang University
[2] ZIU-UIUC Institute, Zhejiang University
[3] Second Affiliated Hospital School of Medicine, and School of Public Health, Zhejiang University

{chenlu.22, yulin, gaoangwang, hongweiwang}@intl.zju.edu.cn
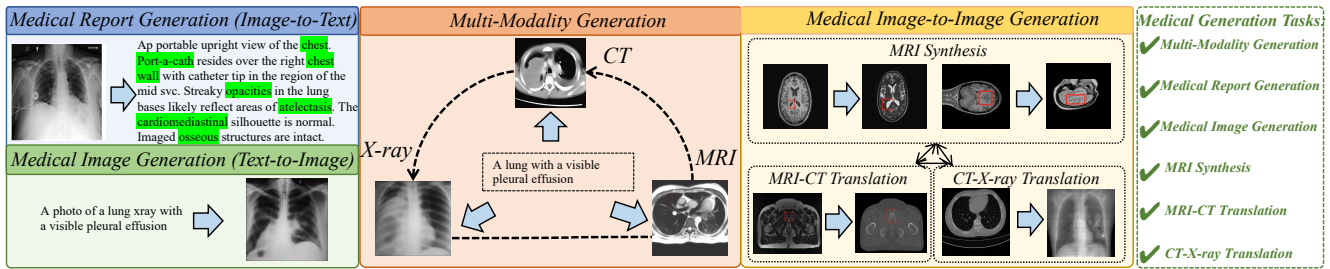
wujian2000@zju.edu.cn

Figure 1. **Our MedM2G on multiple medical generative tasks.** By effectively extracting clinical visual knowledge of multiple medical modalities and adopting the latent multi-flow cross-guided diffusion process, MedM2G has the capability of the unified medical image-to-text, text-to-image diffusion, as well as the unified generation of medical modalities (CT, MRI, X-ray).

## Abstract

*Medical generative models, acknowledged for their high-quality sample generation ability, have accelerated the fast growth of medical applications. However, recent works concentrate on separate medical generation models for distinct medical tasks and are restricted to inadequate medical multi-modal knowledge, constraining medical comprehensive diagnosis. In this paper, we propose **MedM2G**, a **Med**ical **M**ulti-**M**odal **G**enerative framework, with the key innovation to align, extract, and generate medical multi-modal within a unified model. Extending beyond single or two medical modalities, we efficiently align medical multi-modal through the central alignment approach in the unified space. Significantly, our framework extracts valuable clinical knowledge by preserving the medical visual invariant of each imaging modal, thereby enhancing specific medical information for multi-modal generation. By conditioning the adaptive cross-guided parameters into the multi-flow diffusion framework, our model promotes flexible interactions among medical multi-modal for generation. MedM2G is the first medical generative model that unifies medical generation tasks of text-to-image, image-to-text, and unified generation of medical modalities (CT, MRI, X-ray). It performs 5 medical generation tasks across 10 datasets, consistently outperforming various state-of-the-art works.*

Corresponding Authors[(✉)]: Hongwei Wang and Gaoang Wang.

## 1. Introduction

Recently various advanced medical generative works based on denoising diffusion models [17, 40, 41, 43] have significantly improved the efficiency of medical diagnostics tasks, such as medical text-to-image [25, 57], image-to-text generation tasks [58, 59], MRI-CT transaction task [6, 36], MRI synthesis task [22, 41, 60]. The generation of medical modality concentrates on capturing the distinctive specific medical knowledge of each modal and extends to corresponding medical applications.

However, most of these medical generative models [36, 46, 54, 57] rely on distinct single-flow pipelines for specialized generative tasks [22, 41] with cumbersome and slow processes. In real-world medical scenarios that demand the integration of multiple medical modalities for analysis, this generative approach faces substantial limitations in its ex-

tension. Besides, recent advanced multi-modal generative works [13, 52, 56] face challenges in extracting specific medical knowledge and leveraging limited medical paired data to attain cross-modal generation capabilities. These motivate us to construct a unified medical generative model capable of handling tasks of multiple medical modalities. There still exist some non-trivial challenges, as follows: (1) The substantial disparities among multiple medical modalities pose significant challenges in achieving alignment and come with expensive costs. (2) Distinct from images in the general domain, medical imaging modalities (CT, MRI, X-ray) each possess their specific clinical properties. The conventional unified alignment methods [13, 52, 56] often lead to a mixing. (3) Unlike the general multi-modal generative models [52, 56] pre-trained with large well-matched cross-modal databases, the lack of medical cross-modal paired training datasets poses difficulty in retraining generative capabilities of medical multi-modal.

To address the above challenges, we propose **MedM2G**, a unified **Med**ical **M**ulti-**M**odal **G**enerative Model that innovates to align, extract, and generate multiple medical modalities in a unified model, as shown in Fig. 1. MedM2G enables medical multi-modal generation by interacting with multiple diffusion models. The primary motivation is to address the following issues: 1) MedM2G can generate paired data for arbitrary modalities. We leverage the data generated to pre-train and improve the performance of downstream tasks (classification, segmentation, detection, translation). 2) MedM2G can compensate for scarce medical modals by generation. 3) MedM2G can fuse and generate multi-modal for medical comprehensive analysis. 4) MedM2G can handle multiple tasks within a unified model and achieves SOTA results. Specifically, extending to align multiple medical modalities with efficient cost, we first propose the central alignment efficiently adopted in the input and output sharing space, which simply aligns the embedding of each modality with the text embedding, resulting in the alignment across all modalities (Section 3.2). Significantly, with the innovation to maintain the specific medical knowledge of three medical imaging modalities unique to the cross-modal concept generation, we propose the medical visual invariant preservation by minimizing the off-diagonal elements of the two augmented views for better extraction (Section 3.3). Moreover, boosting the interaction of medical cross-modal is crucial, we hence condition the adaptive representation and a shareable cross-attention sublayer into each cross-modal diffuser (Section 3.4). Combined with the proposed multi-flow training strategy (Section 3.5), our model can seamlessly handle multiple medical generation tasks without cross-modal paired datasets. We conduct extensive experiments on 5 medical multi-modal generation tasks across corresponding 10 datasets. Comprehensive experiments validate the effectiveness and ef-

ficiency of MedM2G in its capacity to align, extract and generate multiple medical modalities. Our contributions are summarized as follows:

- We propose MedM2G, the first unified medical multi-flow generative framework capable of aligning, extracting and generating multiple medical modalities.
- We present the multi-flow cross-guided diffusion strategy with the adaptive parameter as the condition for efficient medical multi-modal generation, cooperating with the medical visual invariant preservation to maintain specific medical knowledge.
- MedM2G attains state-of-the-art results on 5 medical multi-modal generation tasks with 10 corresponding benchmarks, illustrating the novel capacity of multi-modal medical generation.

## 2. Related Work

### 2.1. Diffusion Model

Diffusion models (DM) [17, 41, 48–50] acquire the data distribution by outlining the forward diffusion phase and reverse this diffusion process by recovering noise-free data from noisy data samples. For recent diffusion works [17, 34, 41, 48], some models [17, 34] generate high-quality samples through the correlation of the adjacent pixels and the others [40, 41, 50] try to construct latent semantic space for improving efficiency. DDP [48] acquires the capability to learn an inverse diffusion procedure that transforms the input image into a latent space and utilizes a decoder to map these latent variables back to an output image that reconstructs the data's structure. DDPM [17] utilizes the diffusion process, optimizing a weighted variational bound that is constructed through an innovative connection between probabilistic diffusion models and denoising score matching using Langevin dynamics. DDIM [49] introduces an implicit diffusion procedure that yields deterministic samples originating from latent variables with minimal expense and superior quality. Another works [18, 44] introduce an adaptable learning approach that enables the gradual adjustment of noise parameters to achieve superior quality and speed. LDM [41] employs VAE for embedding inputs into a latent space to reduce modeling dimensions and enhance efficiency. These works are primarily centered on enhancing single-flow diffusion pipelines, lacking the capability to handle the multi-flow generation in a unified model. To overcome this, some multi-modal generative works [13, 52, 56] are effective in handling multiple modalities in the general domain but are constrained to the large distinction of medical modalities and absent well-paired datasets. There still remains a challenge in the effective extraction of medical information while aligning multiple modalities in a unified space.
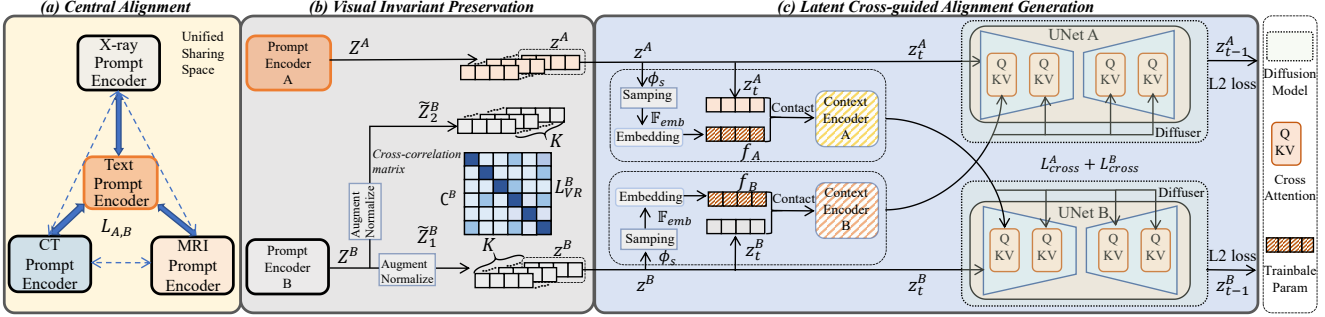
Figure 2. **The network structure of MedM2G.** (a) The multiple medical modalities are embedded into a unified sharing space and present the text as the central modality to efficiently align the other modalities. (b) To maintain the clinic knowledge, we minimize the off-diagonal elements of the cross-correlation matrix of the two augmented image views. (c) We directly condition the representation as the trainable adaptation to capture the semantic knowledge for the generation and adopt the cross-attention sub-layer of one modality to align another.

## 2.2. Medical Generative Modeling

Recently, there has been a remarkable surge in the use of diffusion-based methods [14, 22, 32, 36] in the medical imaging community, which encompass various medical generative tasks, such as medical image-to-text generation tasks [9, 54], text-to-image tasks [24, 25, 46], and the medical image-to-image tasks (e.g. MRI-CT [6, 19, 36, 45], MRI synthesis [22, 41, 47, 60, 63], Xray-CT [8, 12, 33]). For the single-modal translation, CoLa-Diff [22] introduces brain region masks as the dense distribution priors into diffusion guidance. GoentGen [3] devise a pre-trained latent diffusion model to address the substantial natural medical distributional discrepancy. For the multi-modal generation tasks, SynDiff [36] utilizes a conditional diffusion procedure to gradually transform noise and source images into the target image, achieving high-fidelity synthesis. MT-Diffusion [32] proposes denoising diffusion probabilistic and score-matching models for generating high-quality CT images. BrainGen [14] adopts a fast diffusion prior coupled with an adversarial mapping process to enable efficient image generation. These works are devised for the conversions between a single modality or two modalities, which motivates us to exploit a unified generative diffusion model for aligning and generating multiple medical modalities.

## 3. Methodology

In this section, we propose MedM2G, a unified medical generative model capable of aligning and generating multiple medical modalities. Fig. 2 illustrates the main structure which consists of (a) the central alignment strategy (b) the medical visual invariance preservation (c) the latent cross-guided diffusion process with multi-flow training structures.

### 3.1. Preliminary: Latent Diffusion Model

We base our diffusion model on LDM [41] which consists of a forward process and a reverse process. LDM diffuses the latent variable $z$ across multiple time steps $t$ following a variance schedule $\beta_t$ and reconstructs $z_t$ from the noise of the $t$-step through the UNet $\epsilon_\theta$ parameterized by $\theta$. These processes can be parameterized as:

$$q\left(z_t \mid z_{t-1}\right) = \mathcal{N}\left(z_t; \sqrt{1-\beta_t}z_{t-1}, \beta_t\boldsymbol{I}\right)$$
$$p\left(z_{t-1} \mid z_t\right) = \mathcal{N}\left(z_{t-1}; \frac{1}{\sqrt{\alpha_t}}\left(z_t - \frac{\beta_t}{\sqrt{\sigma_t}}\epsilon_\theta\right), \beta_t\boldsymbol{I}\right) \quad (1)$$

where $\beta_t$ is a series of hyper-parameters. $z_t = \alpha_t z + \sigma_t\epsilon$, $\alpha_t = 1 - \beta_t$ and $\sigma_t = 1 - \prod_{s=1}^{t}\alpha_s$. The training objective of the denoising process can be defined as:

$$\mathcal{L} = \mathbb{E}_{\boldsymbol{z},\boldsymbol{\epsilon}\sim\mathcal{N}(0,I),t}\left\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta\left(\boldsymbol{z}_t, t, C(y)\right)\right\|_2^2 \quad (2)$$

where $y$ is the variable for generations; $C$ is a prompt encoder which embeds $y$ into the encoder and controls the $C(y)$ through the cross-attention layers in the UNet $\epsilon_\theta$.

**Our Works** MedM2G extends to unify multiple medical modalities generation tasks in three steps: Align, Extract, Generate. (1) MedM2G first efficiently aligns multiple medical modalities in a unified space with the central alignment tackle with the limited paired dataset (Section 3.2, 3.5). (2) Notably, we extract effective clinic knowledge of each modal through the medical invariance for generation (Section 3.3). (3) For multi-modal generations, we proposed the cross-guided alignment diffusion with trainable adaptative parameters to further enhance the interaction of multi-modal (Section 3.4).

### 3.2. Unified Central Alignment

To facilitate our model with the capability to align and integrate multiple medical modalities (text, CT, MRI, X-ray), we initially aligned the four prompt encoders ($C_M$: $C_T$, $C_{MRI}$, $C_{CT}$, $C_{Xray}$) into a unified sharing space. However, optimizing multiple encoders in a pairwise fashion imposes a substantial computational burden, demanding $\mathcal{O}\left(n^2\right)$ pairs. Furthermore, there is a lack of well-matched

medical multi-modal data pairs for training the cross-modal frameworks, such as Xray-MRI data pairs.

**Central Alignment** To address the above two challenges, as shown in Fig. 2 (a), we developed a "Central Alignment" method to effectively align multiple modalities with $\mathcal{O}(n)$ pairs. Since the text mode is present in most medical cross-modal paired data, we first choose the text model $T$ as the central to align the other three medical imaging modalities, which are denoted as $M$. Afterward, we proceed with pairwise alignment between the remaining modalities. Given a medical feature of A modality $x_i^A$ and the feature of other modalities $x_i^B$, the embeddings $z_i^A = C_T(x_i^A)$ and $z_i^B = C_B(x_i^B)$ are aligned through the InfoNCE contrastive loss [39]:

$$\mathcal{L}_{A,B} = -\log \frac{\exp(z_i^{A\top} z_i^B/\tau)}{\exp(z_i^{A\top} z_i^B/\tau) + \sum_{j \neq i} \exp(z_i^{A\top} z_j^B/\tau)} \quad (3)$$

where $\tau$ is the scalar temperature regulating the softness of the softmax distribution, and $j$ refers to negative samples. We adopt the symmetric loss $L_{A,B} + L_{A,B}$ to make the embeddings $q_i^A$ and $k_i^B$ closer to align the dual modalities.

**Alignment of Modality Pairs** Taking text-Xray pairs as an example, based on a symmetric loss, we train text and CT prompt encoders $C_t$, $C_{Xray}$ on the text-Xray paired dataset and freeze the weights of the other encoders. The remaining encoders are also aligned in the same sharing space as the text modality. Afterward, the other paired modalities (except text) are trained on existing paired data using the alignment method described in Section 3.2, freezing the parameters of other modality encoders. This alignment approach results in a spontaneous and efficient alignment with limited paried data across all modalities. Notably, medical multi-modal (CT, MRI, X-Ray) without well-paired data can also be aligned implicitly within the same space, providing the capability for a versatile generation.

### 3.3. Medical Visual Invariant Preservation

In order to maintain the valuable clinical information of the three medical imaging modalities, we designed a medical visual invariant preservation method to extract high-quality medical feature representations in Fig. 2 (b). For each medical imaging modality $M$, given the dataset $D'$ consists of the medical images $X^M$, we start by generating two augmented views $X_1^M$ and $X_2^M$ of each medical imaging feature and feed them into the encoder for obtaining the augmented embeddings $\{Z_1^M, Z_2^M\} \in \mathbb{R}^{N \times d}$, where $N$ is the batch size and $d$ is the feature dimension. Then we retrain the $\{\widetilde{Z}_1^M, \widetilde{Z}_2^M\}$ by normalizing the augmented embeddings along the batch $K$ dimension. The feature dimension of normalized $\widetilde{Z}^M$ has a zero-mean and $1/\sqrt{K}$ standard deviation distribution. Next, we compute their cross-correlation $\mathcal{C}^M = \widetilde{Z}_1^{M\top} \widetilde{Z}_2^M$.
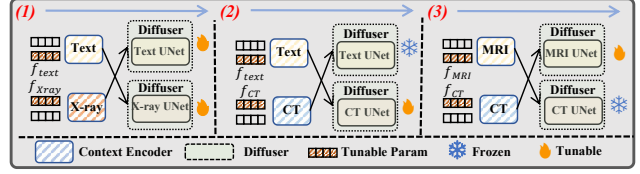


Figure 3. The multi-flow training strategy through 3 rounds of paired training for the multi-modal generation with the central alignment.

The objective $L_{VR}^M$ of the visual invariant preservation is to minimize the off-diagonal elements of the cross-correlation matrix $\mathcal{C}_{ij}^M$ while maximizing the diagonal elements, which is defined:

$$\mathcal{L}_{VR}^M = \frac{1}{D'}\{\sum_j^{D'} \left(1 - \mathcal{C}_{ii}^M\right)^2 + \lambda_1 \sum_j^{D'} \mathcal{C}_{ij}^{M^2}\} \quad (4)$$

where $\lambda_1$ is a non-negative balancing hyperparameter, which follows the default setting in Barlow twins [61].

In this way, the multiple medical modalities are all aligned in a unified sharing space. One may notice that the medical clinical knowledge of each imaging modality is also well-maintained by preserving the visual invariant. It should be noted that the VI module is optimized as a negative-free objective (Barlow twins [61]) instead of general positive-negative loss, which aims to disentangle the latent space feature-wisely.

### 3.4. Latent Cross-guided Alignment Generation

As shown in Fig. 2 (c), we established the latent cross-guided alignment generation structure, which is devised to acquire adaptive interaction information among different modalities for medical multi-modal generation.

For the medical single-modal generation, we first train the individual LDM [41] of the medical text, CT, MRI, and X-ray modalities, the detailed introduction is conducted in Appendix A. These diffusion models subsequently train for medical multi-modal generation through the proposed cross-guided alignment generation method.

**Guided Adaptation.** In order to fully promote the interaction of medical multi-modal, we invert the modality representation of $B$ as the continuously guided trainable adaptation $f_B$ to capture the valuable clinic knowledge unique to the cross-modal concept generation. Following the Textual Inversion method [11], we initialize $f_B$ as a set of context parameters sampled randomly from $z_B = C_B(x_B)$ of modality $B$, with the size the same as the representation of the cross-modality $A$ through the embedding layer $\mathbb{F}_{emb}$:

$$f_B = \mathbb{F}_{emb}(\phi_s(z_B)) \quad (5)$$

where the $\phi_s$ is the sampling strategy. The trainable parameters $f_B$ are then integrated into the modality generation

process of $A$ and assist in aligning the medical cross-modal representation within the unified latent space by directly optimizing the aforementioned loss function in Eq. 2.

**Cross-condition.** Specifically, based on the LDM [41] illustrated in Section 3.1, our cross-modal diffusion model aims to make a condition of the modality $A$ and $B$. We denote the latent variables for modality $A$ and $B$ at diffusion step $t$ as $z_t^A$ and $z_t^B$, respectively. We first project $z_t^B$ and adaptive parameter $f_B$ into a shared latent space of another modality through a context encoder $V_B$ and then adopt the cross-attention sublayer of the UNet for modality A to align the $V_B([z_t^B, f_B])$. The context encoder is devised to embed the latent variable into a unified sharing latent space.

Finally, the training objective for our diffusion model of modality $A$ can be formalized as:

$$\mathcal{L}_{Cross}^A = \mathbb{E}_{z,\epsilon,t,f_B} \parallel \epsilon - \epsilon_{\theta_c}(z_t^A, t, V_B([z_t^B, f_B]) \parallel_2^2 \quad (6)$$

where $\theta_c$ is the weights of the cross-attention layers, $[\cdot, \cdot]$ is the concatenation. We denote the multiple generations of modalities $A$ and $B$ as the $L_{Cross}^A + L_{Cross}^B$.

### 3.5. Multi-flow Training Strategy

The multi-flow training strategy enables the model capable of the medical multi-modal generation abilities in the absence of well-paired data, with a linear procedure through the central alignment in Section 3.2. Our pipeline consists of diffusion models with the VI module for multi-flow training. We start by adopting the pre-trained diffusion models for each medical modality. Then, these diffusion models effectively engage in joint multi-modal generation through 3 rounds of paired training (Text-Xray, Text-CT, CT-MRI) with "Cross-guided Alignment". As shown in Fig. 3, we first train the context encoder $V_T$, $V_{Xray}$ and the cross-attention sub-layer weights of the text and X-ray diffusers on the text-Xray paired dataset. Then we freeze the trainable parameter of the text diffuser and train the context encoder $V_{CT}$ and cross-attention sub-layer weights of the CT diffuser on the text-CT paired datasets. At last, we freeze the trainable parameter of the CT diffuser and train the context encoder $V_{MRI}$ and cross-attention sub-layer weight of the MRI diffuser on the MRI-CT paired datasets. In this multi-flow training procedure, our proposed unified diffusion model can deal with multiple medical generation tasks (Section 5) with merely three medical paired datasets.

### 4. Datasets and Implementation Details

**Datasets** We pre-train our unified diffusion model with the MIMIC-CXR [23], MedICat [51], and Brain tumor MRI, and CT scan [2] datasets for the central alignment. MIMIC-CXR [23] comprises a substantial collection of X-ray data, encompassing $377,100$ chest radiology images and $227,835$ corresponding patient reports. MedI-

Cat [51] is a dataset of contextual medical images, comprising $217,000$ images sourced from $131,000$ freely accessible biomedical papers. Brain tumor MRI and CT scan dataset [2] contains $4,500$ 2D MRI-CT slices. We adhere to the official data partitioning guidelines and filter the paired datasets for aligning different modalities. Detailed introductions of different pre-train tasks with corresponding datasets are in Appendix B. To assess our model's ability to align and generate medical multiple modalities, we conducted evaluations across 10 datasets, spanning 5 medical text-to-image, image-to-text, image-to-image, and multi-modality generation tasks. The experimental setup, i.e. setting of over $40$ parameters and train/fine-tuning processes, is detailed in Appendix C, D, L.

**Medical Multi-modality Generation Tasks** We conduct experiments of MRI synthesis task across BraTS 2020 [2] and IXI [21] datasets. For MRI-CT translation tasks, we train and evaluate on the Gold Atlas male pelvis datasets [35]. We also conduct chest X-ray generation tasks on MIMIC-CXR [23] and Chest X-ray [7] datasets.

**Medical Text-Image Generation Tasks** We evaluate the medical report generation task on MIMIX-CXR [23] and IU X-ray [9] and fine-tune the medical image generation task on Chest X-ray [54], SLIVER07 [16], ACDC [1] datasets. We all follow the official data splits and details of the fine-tuning tasks with datasets can be found in Appendix C.

**Implementation Details** We train the MedM2D with 3 settings on the 6 NVIDIA 3090 GPUs: medical text-Xray, text-MRI, MRI-CT. These training pairs are devised for various downstream tasks. In the course of training, we maintained diffusion settings in close proximity to LDM [41], *i.e.* the diffusion steps of different diffusion models set to 1000 and adopt the Linear noise schedule, the $\beta_0$ and $\beta_T$ are $0.00085$ and $0.0120$ respectively. The learning rates are set to $2e-5$ for medical image LDM and are set to $5e-5$ for text LDM. The weights of medical image diffusion models are initialized from Stable Diffusion-1.5 [41] and the weights of medical text diffusion model with OPTIMUS [28]-BERT [10] and GPT-2 [38] VAE are initialized from Versatile Diffusion [56]. The batch size is 256 for image modalities and 1024 for text training. We also embrace the DDIM [49] sampler for the sampling strategy and set 50 sampling steps, the $\eta$ and the guidance scale set to 1.0 and 2.0. For the diffusion models, the $z$-shape of the medical image and text diffusion models is set to $4 \times 64 \times 64$ and $768 \times 1 \times 1$ respectively. The depth of image and text LDM are 4 and 2. For the cross-attention guided layers in the diffusion modules, we adopt Adam [26] optimizer whose learning rate and weight decay are $1e-5$ and $1e-4$ respectively. Due to space limits, we conduct detailed diffusion model structure hyperparameters and configurations in Appendix D.

| Methods | IU X-Ray(mean±std) | | | | | MIMIC-CXR(mean±std) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 | ROUGE-L | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 | ROUGE-L |
| R2Gen [4] | 0.470 | 0.304 | 0.219 | 0.165 | 0.371 | 0.353 | 0.218 | 0.145 | 0.103 | 0.277 |
| R2GenCMN [5] | 0.475 | 0.309 | 0.222 | 0.170 | 0.375 | 0.353 | 0.218 | 0.148 | 0.106 | 0.278 |
| PPKED [30] | 0.483 | 0.315 | 0.224 | 0.168 | 0.376 | 0.360 | 0.224 | 0.149 | 0.106 | 0.284 |
| AlignTrans [59] | 0.484 | 0.313 | 0.225 | 0.173 | 0.379 | 0.378 | 0.235 | 0.156 | 0.112 | 0.283 |
| Clinical-BERT [58] | 0.495 | 0.330 | 0.231 | 0.170 | 0.376 | 0.383 | 0.230 | 0.151 | 0.106 | 0.275 |
| METransformer [55] | 0.483 | 0.322 | 0.228 | 0.172 | 0.380 | 0.386 | 0.250 | 0.169 | 0.124 | 0.291 |
| COMG [53] | **0.536** | **0.378** | 0.275 | 0.206 | 0.383 | 0.363 | 0.235 | 0.167 | 0.124 | 0.290 |
| Kiut [20] | 0.525 | 0.360 | 0.251 | 0.185 | 0.409 | 0.393 | 0.243 | 0.159 | 0.113 | 0.285 |
| **Ours** | $0.533_{\pm0.009}$ | $0.369_{\pm0.010}$ | $\mathbf{0.278}_{\pm0.011}$ | $\mathbf{0.212}_{\pm0.009}$ | $\mathbf{0.416}_{\pm0.008}$ | $\mathbf{0.412}_{\pm0.007}$ | $\mathbf{0.260}_{\pm0.009}$ | $\mathbf{0.179}_{\pm0.011}$ | $\mathbf{0.142}_{\pm0.010}$ | $\mathbf{0.309}_{\pm0.009}$ |

Table 1. The comparisons between MedM2G and medical report generation methods on IU X-Ray and MIMIC-CXR datasets.

| Methods | BraST | | | | | | | | IXI | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | T2+T1ce+FLAIR→T1 | | T1+T1ce+FLAIR→T2 | | | | | | T2+PD →T1 | | T1+PD →T2 | | | | | |
| | PSNR | SSIM | PSNR | SSIM | | | | | PSNR | SSIM | PSNR | SSIM | | | | |
| MM-GAN [47] | $25.78_{\pm2.16}$ | $90.67_{\pm1.45}$ | $26.11_{\pm1.62}$ | $90.58_{\pm1.39}$ | | | | | $27.32_{\pm1.70}$ | $92.35_{\pm1.58}$ | $30.87_{\pm1.75}$ | $94.68_{\pm1.42}$ | | | | |
| Hi-Net [63] | $27.42_{\pm2.58}$ | $93.46_{\pm1.75}$ | $25.64_{\pm2.01}$ | $92.59_{\pm1.42}$ | | | | | $28.89_{\pm1.43}$ | $93.78_{\pm1.31}$ | $32.58_{\pm1.85}$ | $96.54_{\pm1.74}$ | | | | |
| ProvoGAN [60] | $27.79_{\pm4.42}$ | $93.51_{\pm3.16}$ | $26.72_{\pm2.87}$ | $92.98_{\pm3.91}$ | | | | | $24.21_{\pm2.63}$ | $90.46_{\pm3.58}$ | $29.19_{\pm3.04}$ | $94.08_{\pm3.87}$ | | | | |
| LDM [41] | $24.55_{\pm2.62}$ | $88.34_{\pm2.51}$ | $24.79_{\pm2.67}$ | $88.47_{\pm2.60}$ | | | | | $24.19_{\pm2.51}$ | $88.75_{\pm2.47}$ | $27.04_{\pm2.31}$ | $91.23_{\pm2.24}$ | | | | |
| CoLa-Diff [22] | $28.26_{\pm3.13}$ | $93.65_{\pm3.02}$ | $28.33_{\pm2.27}$ | $93.80_{\pm2.75}$ | | | | | $30.21_{\pm2.38}$ | $94.49_{\pm2.15}$ | $32.86_{\pm2.83}$ | $96.57_{\pm2.27}$ | | | | |
| **Ours** | $\mathbf{29.89}_{\pm2.26}$ | $\mathbf{95.36}_{\pm1.43}$ | $\mathbf{30.51}_{\pm2.02}$ | $\mathbf{96.60}_{\pm1.66}$ | | | | | $\mathbf{32.45}_{\pm2.87}$ | $\mathbf{97.64}_{\pm1.88}$ | $\mathbf{34.81}_{\pm1.78}$ | $\mathbf{98.23}_{\pm1.66}$ | | | | |

Table 2. The comparisons between our model MedM2G and advanced MRI synthesis models on BraST and IXI datasets. Different MRI modalities: T1, T2, T1ce, FLAIR, PD-weighted.

| Method | ChestXray14 | | | ACDC | | | SLIVER07 | | | MIMIC-CXR | | OpenI | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR | SSIM | NIQE↓ | PSNR | SSIM | NIQE | PSNR | SSIM | NIQE | Fid↓ | NIQE | Fid | NIQE |
| StyleGAN | 20.13 | 88.47 | 8.41 | 24.69 | 90.13 | 9.22 | 23.19 | 89.15 | 7.33 | 19.23 | 5.14 | 22.91 | 7.45 |
| GCDP* | 24.51 | 88.69 | 8.02 | 28.14 | 90.69 | 7.92 | 31.43 | 86.75 | 7.18 | 13.23 | 4.82 | 15.72 | 6.58 |
| GLIGEN* | 32.12 | 88.95 | 7.61 | 33.27 | 91.81 | 8.02 | 32.89 | 88.41 | 6.61 | 12.49 | 4.26 | 13.17 | 6.22 |
| RoentGen | 33.24 | 90.25 | 6.33 | 34.91 | 93.27 | 6.82 | 34.25 | 89.96 | 6.22 | 9.54 | 3.88 | 6.56 | 4.90 |
| UniXGen | 34.75 | 91.86 | 5.05 | 36.45 | 94.52 | 5.62 | 35.66 | 91.42 | 5.14 | 6.72 | 3.71 | 11.98 | 4.66 |
| LLM-CXR | 35.92 | 93.56 | 3.81 | 37.89 | 95.68 | 4.42 | 36.94 | 92.89 | 4.59 | 2.18 | 3.60 | 1.66 | 3.82 |
| AdaMatch-Cyclic | 36.82 | 94.91 | 3.77 | 39.32 | 96.74 | 3.22 | 38.25 | 94.27 | 3.69 | 1.09 | 3.39 | 1.59 | 3.30 |
| **Ours** | 40.16 | 98.27 | 2.49 | 42.48 | 98.92 | 2.02 | 39.51 | 95.68 | 2.31 | 0.48 | 2.91 | 0.92 | 2.66 |

Table 3. More eval metrics of medical image generation. Brown: medical-domain models.

| Method | MIMIC-CXR | | Chest X-ray |
|---|---|---|---|
| | FID(↓) | MS-SSIM(↑) | FID(↓) |
| Original SD [41] | 52.7 | $0.09_{\pm0.05}$ | 78.86 |
| DreamBooth SD [42] | 18.6 | $0.28_{\pm0.07}$ | 69.14 |
| SD-RadBERT [57] | 6.0 | $0.26_{\pm0.12}$ | 45.28 |
| UniXGen [27] | 2.5 | $0.29_{\pm0.06}$ | 19.99 |
| Ours | **1.7** | $\mathbf{0.38}_{\pm0.07}$ | **9.76** |

Table 5. The comparisons on the chest X-ray generation task across the MIMIC-CXR and chest X-ray datasets. MS-SSIM: Multi-scale structural similarity index measure.

| Method | Dataset FID (↓) | | |
|---|---|---|---|
| | ChestXray14 | ACDC | SLIVER07 |
| Progressive Growing GAN [46] | 8.02 | - | - |
| StyleGAN [24] | 3.52 | 24.74 | 29.06 |
| StyleGAN2-ADA [25] | - | 21.17 | 10.78 |
| GCDP* [37] | 2.89 | 21.32 | 9.56 |
| GLIGEN* [29] | 2.66 | 20.19 | 8.45 |
| **Ours** | **1.84** | **15.89** | **6.89** |

Table 4. The comparisons of medical image generation across ChestX-ray, ACDC, and SLIVER07 datasets. *: Re-implement on the same pre-train datasets.

# 5. Experiments and Results

To demonstrate the outperformance of MedM2G, we conduct abundant experiments on 5 medical image-to-image generation tasks of MRI (Table. 2), CT (Table. 6), X-ray (Table. 5) and multiple report generation task (Table 1) and medical image generation task (Table 4) over 10 datasets. We also provide quantitative assessments (Fig. 4 and 5) on fine-tuning datasets and the unified medical multi-modal generation capability in Fig 4 (c). The ablation studies are conducted in Table 7 and the comparison between multi-modal generative models is in Table 9 and Fig. 7.

## 5.1. Comparison with State-of-the-art Methods

**Medical Image-to-Report Generation** As shown in Table 1, for the medical image-to-text generation task, we utilize the IU X-ray [9] and the MIMIC-CXR [23] to assess the resemblance scores between the generated reports and the annotated ones. It can be illustrated that our model is superior to the advanced GAN-based works [4, 5], as well as the well-trained Med-VLP works [20, 30, 53, 55, 58, 59], achieving 0.416 and 0.309 of ROUGE-L on two datasets, respectively. The substantial enhancement underscores the effectiveness of the multi-flow cross-guided diffusion process in modalities alignments.

Besides, we also show the visualization samples for qualitative analysis in Fig. 4 (a). Compared with the SOTA model Kiut [20] which devised medical domain-specific knowledge into training, our model has outperformance in generating more accurate and semantic reports. MedM2G aims to facilitate interaction among multiple modalities for broad generative capability. The majority of MeSH terms are correctly predicted (indicated in green), including terms
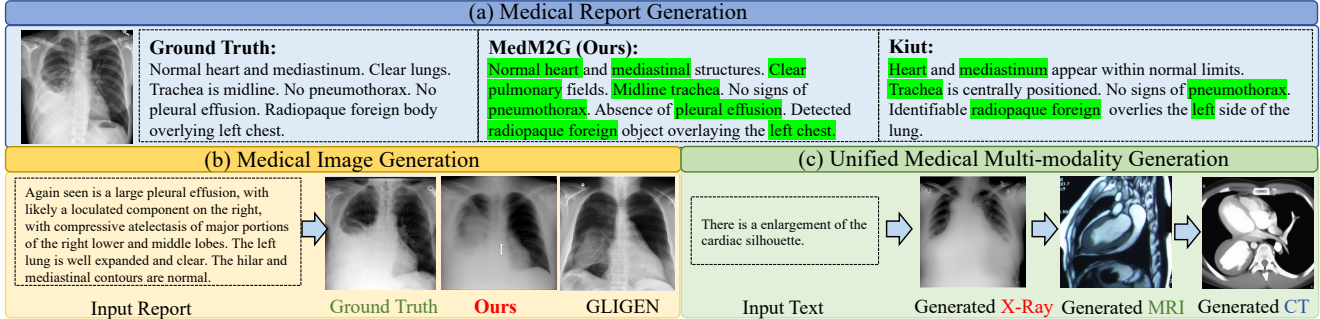
Figure 4. The qualitative analysis of (a) Medical report generation task (b) Medical text-image generation task (c) Unified medical multi-modality generation. The indication in green: the correctly predicted MeSH terms.
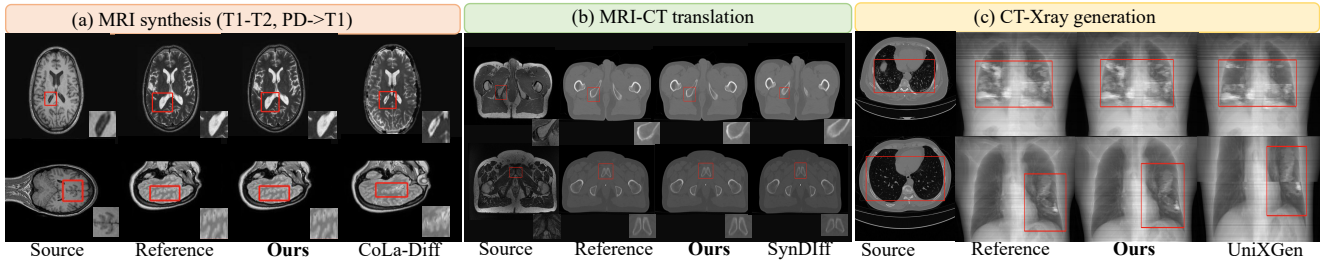


Figure 5. Multiple medical modalities generation tasks by MedM2D. (a) MRI synthesis task on IXI dataset. (b) MRI-CT transition task on Pelvi dataset. (c) CT-Xray generation task on Chestxray dataset..

| Methods | T2→CT | | T1→CT | | acc T2→CT | | acc T1→CT | |
|---|---|---|---|---|---|---|---|---|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| UNIT-DDPM [45] | $21.03_{\pm0.72}$ | $80.23_{\pm2.69}$ | $20.26_{\pm1.17}$ | $76.79_{\pm1.37}$ | $21.89_{\pm0.77}$ | $77.69_{\pm3.06}$ | $21.45_{\pm0.23}$ | $77.10_{\pm2.83}$ |
| DDPM [17] | $21.49_{\pm0.19}$ | $83.24_{\pm2.62}$ | $21.10_{\pm2.41}$ | $73.58_{\pm7.17}$ | $24.35_{\pm0.47}$ | $83.25_{\pm1.70}$ | $24.62_{\pm0.59}$ | $83.04_{\pm2.40}$ |
| SAGAN [62] | $22.90_{\pm0.33}$ | $67.77_{\pm0.86}$ | $23.89_{\pm1.02}$ | $77.05_{\pm2.87}$ | $19.61_{\pm0.78}$ | $61.92_{\pm0.32}$ | $23.28_{\pm0.96}$ | $70.02_{\pm2.85}$ |
| AttGAN [15] | $23.81_{\pm0.18}$ | $74.35_{\pm0.84}$ | $24.76_{\pm1.06}$ | $82.48_{\pm2.49}$ | $23.91_{\pm0.29}$ | $76.47_{\pm0.66}$ | $21.34_{\pm0.51}$ | $67.24_{\pm1.52}$ |
| MUNIT [19] | $24.66_{\pm1.05}$ | $77.42_{\pm2.17}$ | $24.76_{\pm0.62}$ | $79.81_{\pm1.20}$ | $23.44_{\pm0.77}$ | $77.88_{\pm2.04}$ | $24.42_{\pm0.34}$ | $79.64_{\pm1.05}$ |
| UNIT [31] | $25.07_{\pm0.49}$ | $86.40_{\pm2.71}$ | $25.04_{\pm0.39}$ | $82.62_{\pm1.52}$ | $25.20_{\pm0.37}$ | $84.83_{\pm1.43}$ | $24.92_{\pm0.39}$ | $81.44_{\pm1.13}$ |
| cGAN [6] | $26.10_{\pm0.17}$ | $84.91_{\pm1.84}$ | $24.11_{\pm1.00}$ | $77.81_{\pm1.84}$ | $21.24_{\pm0.51}$ | $69.62_{\pm0.85}$ | $20.35_{\pm0.32}$ | $64.73_{\pm1.47}$ |
| SynDiff [36] | $26.86_{\pm0.51}$ | $87.94_{\pm2.53}$ | $25.16_{\pm1.53}$ | $86.02_{\pm2.05}$ | $26.71_{\pm0.63}$ | $87.32_{\pm2.84}$ | $25.47_{\pm1.09}$ | $85.00_{\pm2.10}$ |
| Ours | $\mathbf{27.45}_{\pm0.64}$ | $\mathbf{89.23}_{\pm1.25}$ | $\mathbf{26.08}_{\pm0.68}$ | $\mathbf{88.34}_{\pm1.34}$ | $\mathbf{28.13}_{\pm0.46}$ | $\mathbf{89.99}_{\pm1.78}$ | $\mathbf{26.94}_{\pm0.24}$ | $\mathbf{87.23}_{\pm2.05}$ |

Table 6. The comparisons of MRI-CT translation tasks across Pelvic dataset. acc: accelerated tasks.

such as "mediastinum" and "pleura effusion". More qualitative analysis samples can be found in Appendix E.

**Medical Text-to-Image Generation** In Table 4, we take comparisons on Chest X-ray14 [54], ACDC [1] and SLIVER07 [16] datasets to quantify the generated images by assessing the similarity (FID) between their feature distribution and that of real images.

Compared with the advanced generative adversarial networks [24, 25, 46] and the text-to-image diffusion works [29, 37] which have the capability of generating high-resolution medical images, our proposed model can considerably decrease the FID of these SOTA works by 0.82, 4.30, 1.56 on above 3 datasets respectively. Besides, we employed more reasonable evaluation metrics (PSNR, SSIM, NIQE) on more relevant SOTA medical models in Tab. 3 to validate the outperformance. Overall, ours achieved SOTA

results on 5 evaluation metrics. This demonstrates the superior generative ability of MedM2G in medical bidirectional text-image generation. We also show the qualitative analysis in Fig. 4 (b). In comparison to the SOTA model GLINGEN [25], our model excels in its ability to precisely and semantically generate depictions of critical pathological regions based on input medical reports. More comparisons between MedM2G and advanced text-image generative models are in Appendix F. **Medical MRI Systhesis** As shown in Table 2, we conducted MRI synthesis tasks on four modalities on the IXI [21] and BraST [2] datasets. We designate one target modality while using the remaining modalities as conditioning factors. It illustrates that our model outperforms the advanced GAN-based generative models [47, 60, 63]. Moreover, compared with the preeminent diffusion works [22, 41] which devises the condi-

| CA | LCAG | VI | MIMIC-CXR | | | ACDC | MIMIC-CXR (X-Ray generation) | BraTS2020 (T2+T1→PD) | | Pelvic (T2→CT) | | Pre-train(h) time | Add parameter |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | BLEU-1 | BLEU-4 | ROUGE_L | FID(↓) | FID(↓) | PSNR | SSIM | PSNR | SSIM | /epoch | /M |
| × | × | × | $0.365_{\pm0.012}$ | $0.100_{\pm0.011}$ | $0.261_{\pm0.013}$ | 26.02 | 8.3 | $31.02_{\pm2.46}$ | $95.77_{\pm1.76}$ | $24.56_{\pm0.78}$ | $88.05_{\pm2.13}$ | 0.6 | - |
| ✓ | × | × | $0.385_{\pm0.009}$ | $0.108_{\pm0.007}$ | $0.274_{\pm0.012}$ | 21.02 | 4.5 | $32.56_{\pm1.53}$ | $97.23_{\pm2.32}$ | $26.99_{\pm0.23}$ | $88.45_{\pm2.67}$ | 1.2 | 32.5 |
| × | ✓ | × | $0.379_{\pm0.008}$ | $0.105_{\pm0.013}$ | $0.271_{\pm0.014}$ | 21.43 | 4.3 | $32.63_{\pm1.89}$ | $97.31_{\pm2.11}$ | $27.09_{\pm0.45}$ | $88.54_{\pm2.98}$ | 3.2 | 128.3 |
| × | × | ✓ | $0.375_{\pm0.011}$ | $0.106_{\pm0.014}$ | $0.272_{\pm0.008}$ | 23.45 | 3.9 | $33.12_{\pm3.45}$ | $97.39_{\pm2.12}$ | $27.11_{\pm0.38}$ | $88.67_{\pm2.38}$ | 0.8 | 13.7 |
| ✓ | × | ✓ | $0.385_{\pm0.014}$ | $0.110_{\pm0.015}$ | $0.278_{\pm0.009}$ | 21.67 | 3.0 | $33.89_{\pm2.24}$ | $97.74_{\pm2.98}$ | $27.40_{\pm0.78}$ | $89.08_{\pm1.45}$ | 0.5 | 40.3 |
| × | ✓ | ✓ | $0.389_{\pm0.006}$ | $0.113_{\pm0.012}$ | $0.281_{\pm0.008}$ | 21.34 | 3.2 | $33.68_{\pm3.03}$ | $97.67_{\pm2.55}$ | $27.36_{\pm0.44}$ | $88.99_{\pm0.97}$ | 3.6 | 140.2 |
| ✓ | ✓ | × | $0.392_{\pm0.008}$ | $0.118_{\pm0.008}$ | $0.287_{\pm0.011}$ | 20.56 | 3.6 | $33.21_{\pm2.33}$ | $97.45_{\pm1.98}$ | $27.24_{\pm0.47}$ | $88.78_{\pm2.91}$ | 1.6 | 88.6 |
| ✓ | ✓ | ✓ | $\mathbf{0.412}_{\pm0.007}$ | $\mathbf{0.142}_{\pm0.010}$ | $\mathbf{0.309}_{\pm0.009}$ | **15.89** | **1.7** | $\mathbf{34.12}_{\pm1.98}$ | $\mathbf{97.88}_{\pm1.89}$ | $\mathbf{27.45}_{\pm0.19}$ | $\mathbf{89.23}_{\pm1.54}$ | 1.8 | 96.6 |

Table 7. Ablation study on the MIMIC-CXR(test set), ACDC, BraTS2020, and the Pelvis datasets. "CA" represents the central alignment strategy. "LCGA" is the Latent Cross-guided Alignment Generation procedure, and "VI" represents the medical visual invariant.
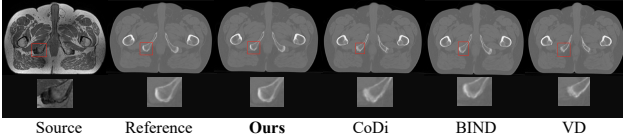


Figure 6. The comparison of MRI-CT translation task across Pelvis dataset between the multi-modal generative models.



Figure 7. The t-SNE of multi-modal with/without VI module.

| Method | CheXpert(AUC) (Classification) | | SIIM(Dice) (Segmentation) | | RSNA (mAP) (Object Detection) | | MimicCXR (generation) | BraTS2020 (T2+T1→PD) | |
|---|---|---|---|---|---|---|---|---|---|
| +Data | 1% | 100% | 1% | 100% | 1% | 100% | ROUGE-L↑ | PSNR | SSIM |
| MGCA | 87.6 | 88.2 | 49.7 | 64.2 | 12.9 | 16.8 | / | / | / |
| +UniXGen | 85.3(-2.3) | 86.2(-2.0) | 47.6(-2.1) | 62.7(-1.5) | 11.2(-1.7) | 14.3(-2.5) | / | / | / |
| +Ours | **89.4**(+1.8) | **90.3**(+2.1) | **51.6**(+1.9) | **66.5**(+2.3) | **14.9**(+2.0) | **19.1**(+2.3) | / | / | / |
| Baseline | 89.5 | 90.4 | 57.8 | 65.5 | 15.9 | 27.4 | 30.9 | 34.1 | 97.9 |
| +UniXGen | 86.8(-2.7) | 87.9(-2.5) | 54.9(-2.9) | 64.5(-1.0) | 13.8(-2.1) | 24.8(-2.6) | 29.7(-1.2) | 32.5(-1.6) | 96.6(-1.3) |
| +Ours | **91.9**(+2.4) | **92.7**(+2.3) | **60.1**(+2.3) | **68.2**(+2.7) | **18.2**(+2.3) | **30.3**(+2.9) | **34.1**(+3.2) | **36.9**(+2.8) | **99.2**(+1.4) |

Table 8. Comparison of our baseline and SOTA medical vision-language pre-train model MGCA after adding(+) the data generated by us and by SOTA medical generative model UniXGen.

| Methods | Pre-train samples | MIMIC-CXR FID | BraTS | | Pelvic | |
|---|---|---|---|---|---|---|
| | | | PSNR | SSIM | PSNR | SSIM |
| VD* [56] | 700M | 12.7 | $28.97_{\pm2.12}$ | $78.45_{\pm2.33}$ | $17.87_{\pm0.98}$ | $71.43_{\pm1.34}$ |
| BIND* [13] | 2270K | 14.6 | $27.66_{\pm1.45}$ | $71.12_{\pm1.87}$ | $15.43_{\pm1.13}$ | $65.38_{\pm1.88}$ |
| CoDi* [52] | 512M | 10.9 | $29.12_{\pm2.11}$ | $80.68_{\pm1.86}$ | $19.12_{\pm0.88}$ | $73.23_{\pm1.22}$ |
| **Ours** | **598K** | **2.7** | $\mathbf{34.12}_{\pm1.98}$ | $\mathbf{97.88}_{\pm1.89}$ | $\mathbf{27.45}_{\pm0.19}$ | $\mathbf{89.23}_{\pm1.54}$ |

Table 9. The comparisons between the multi-modal generative models. *: Re-implement results on medical downstream tasks.

tional latent diffusion for more effective MRI synthesis, our model considerably exceeds by up to 1.63 dB on PSNR of T2+T1ce+FLAIR-T1 in BraST dataset. A possible explanation could be that the unified central alignment with the medical visual invariants promotes the medical knowledge alignment of multiple modalities to synthesize accurate and high-quality MRI. As shown in Fig. 5 (a), we showcase the high-quality MRI generated by our model on the IXI dataset [21]. Compared to the advanced generative model CoLa-Diff [22], our model exhibits outperformance in generating intricate brain sulci and tumor boundaries while effectively preserving anatomical structure. We present more comparisons in Appendix G.

**Medical MRI-CT Translation** We compare MedM2G with SOTA generative works in Table 6, including the diffusion-based models [17, 45] and the GAN-based works [6, 19, 31], attention-GAN-based works [15, 62]. Our proposed model yields the best performance on all four MRI-CT modality translation tasks (p< 0.05). Besides, we observe that MedM2G outperforms the SOTA work Syn-Diff [36] by 0.59dB, 0.92dB, 1.42dB,1.47dB PSNR of all tasks on average respectively. It illustrates that our model is
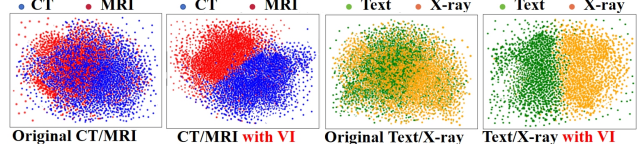
superior in generating more accurate and high-anatomical fidelity CT scans with unified cross-guided alignment diffusion and visual invariant preservation. Fig 5 (b) showcases MedM2M's proficiency in generating CT from MRI in the Pelvi dataset. When compared to SynDiff [36], an advanced medical image translation model, our MedM2M consistently excels in lower artifact levels and more accurate estimation of anatomical structures around diagnostically significant areas. This highlights the unified performance of our model in medical multi-modal generation.

**Chest X-ray Generation** As shown in Table 5, we conduct the fidelity and diversity through FID and MS-SSIM metrics on the chest X-ray generation task over the MIMIC-CXR [23] and Chest X-ray [7] datasets. MedM2G outperforms all the SOTA works [27, 41, 42, 57] which all pre-train with the large-scale clinic text datasets, achieving 1.7 FID and 0.38 MS-SSIM on average. Benefiting from the multi-flow cross-guided diffusion process and the medical visual invariants, our model has a significant advantage in generating higher fidelity and diversity of X-rays. Likewise, as shown in Fig. 5 (c), we showcase MedM2M's superior generative capabilities in precisely generating the contours of both lungs, the heart, and the trachea, along with corresponding anomalous chest regions of nodules.

**Unified Multi-modality Joint Generation** To demonstrate the unified generation ability of medical multi-modal within a diffusion model, we also illustrate the high-quality medical multi-modal generation results in Fig. 4 (c). It becomes evident that, based on the provided medical description, our model can simultaneously generate multiple modalities of MRI, CT, and X-ray (columns 2-4). The generated medical images across three modalities accurately pinpoint the medical abnormalities regions, as exemplified by " degenerative changes" in the first line. We also provide more joint multi-modal generation samples by MedM2G in Appendix H. Notably, MedM2G is the first medical generative model that not only performs generations between the text and images, but also acts as a bridge for medical multi-modality generation between MRI, CT, and X-ray. Different modalities may contain complementary information. Note that our "Cross-guided Alignment" is trained on well-paired open-source data, ensuring the absence of conflicting information. Experiments illustrated that no complementary arises.

## 5.2. Ablation Study

As shown in Table 7, We conduct ablation studies to validate the efficacy of the proposed methods. We take the LDM [41] model pre-trained with MIMIC-CXR datasets as the baseline, as shown in row 1 of Table 7.

**Multi-flow Central Alignment** From the comparison of the rows 1 and 2, we can obverse that the central alignment strategy effectively obtains the 0.013 improvements of ROUGE-L and 5.0 decrease on ACDC datasets, illustrating that the central alignment efficiently benefits the alignment of various medical modalities with linearly increased computation costs. Moreover, we also conduct the ablation studies of the multi-flow training strategy in Appendix I and demonstrate that it can merge medical multi-modal on a deeper alignment with efficient computational costs.

**Cross-guided Diffusion** Besides, the comparison between the rows 1 and 3 in Table 7 reveals that the latent cross-guided diffusion process also decreases the FID of chest x-ray generation task by 4.0, which effectively promotes the interaction of multiple modalities with cost-efficient.

**Medical Visual Invariant** As shown in row 1 and 4 in Table 7, the model equipped with the visual invariants improves the PSNR of BraTS2020 by 2.10. VI excels in capturing intricate clinic structural information, especially when medical multi-modal coexist, effectively preserving clinical knowledge. We also provide a visualization of embedding through t-SNE in Fig. 7, which reveals the two modalities exhibit confusion within a unified space. In contrast, the imaging modality combined **with the VI module** exhibits **cohesion** to preserve its own visual information. When aligning multi-modal to the same space through CA&LCAG modules with multi-flow training, confusion occurs among modalities. Hence, VI is crucial to preserve the features of each imaging modality and achieve obvious improvements. This strategy enhances the medical valuable imaging information of each modal, prompting high-quality medical image generation. When added independently, VI acts as an augmentation without significant improvement.

**Comparison with Multi-modal Generative Model** As shown in Table 9 and Fig. 7, we compare our model with the SOTA multi-modal generative models [13, 52, 56] under the same settings. In scenarios with comparatively smaller training resources, our multi-modal generative model demonstrates a distinct advantage in medical downstream tasks and generates more high-quality medical images in granularity and lower artifact levels, attributable to the proposed medical visual invariance and cross-guided diffusion. More comparison between MedM2G and multi-modal generative model can be found in Appendix J.

**Computation Cost** We list the pre-train time and parameter cost in Table 7. The comparison proves that central alignment effectively minimizes the computational expense for pairwise alignment across multiple modalities, maintaining a cost-efficient linear increase. More computation costs of different training settings can be found in Appendix K. **Pre-train with Generated Data** In Tab. 8, we utilized the data generated by us to pre-train and significantly benefit the downstream medical imaging and translation tasks.

## 6. Conclusion

In this paper, we introduce MedM2G, the first medical generative model to align, extract and generate medical multi-modal within a unified model. The key innovation concentrates on the effective clinical knowledge extraction of each medical modality through the proposed visual invariant preservation, as well as the proposed latent multi-flow cross-guided diffusion framework to efficiently enhance the cross-modal interaction for multi-modal generation. MedM2G achieves superior results across 5 medical generation tasks on 10 datasets. Codes will be released.

## 7. Limitation

Although MedM2G achieved excellent performance in multiple medical generation tasks, we also considered the potential limitations, including: 1) Fake information. Malicious actors could exploit the powerful medical modal generation ability of MedM2G to fabricate false medical information. 2) Comprehensive medical information. For some diseases that do not have multimodal features clinically, the modal generated by ours can only be used as auxiliary information, which needs to be comprehensively analyzed.

## 8. Acknowledgements

# References

[1] Olivier Bernard, Alain Lalande, Clement Zotti, Frederick Cervenansky, Xin Yang, Pheng-Ann Heng, Irem Cetin, Karim Lekadir, Oscar Camara, Miguel Angel Gonzalez Ballester, et al. Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: is the problem solved? *IEEE transactions on medical imaging*, 37 (11):2514–2525, 2018. 5, 7

[2] Brain tumor MRI and CT scan. https://www.kaggle.com/datasets/chenghanpu/brain-tumor-mri-and-ct-scan, 2022. Accessed: 2022-03-18. 5, 7

[3] 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. *arXiv preprint arXiv:2211.12737*, 2022. 3

[4] Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. Generating radiology reports via memory-driven transformer. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1439–1449, 2020. 6

[5] Zhihong Chen, Yaling Shen, Yan Song, and Xiang Wan. Cross-modal memory networks for radiology report generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5904–5914, 2021. 6

[6] Grigorios G Chrysos, Jean Kossaifi, and Stefanos Zafeiriou. Robust conditional generative adversarial networks. *arXiv preprint arXiv:1805.08657*, 2018. 1, 3, 7, 8

[7] Joseph Paul Cohen, Joseph D Viviano, Paul Bertin, Paul Morrison, Parsa Torabian, Matteo Guarrera, Matthew P Lungren, Akshay Chaudhari, Rupert Brooks, Mohammad Hashir, et al. Torchxrayvision: A library of chest x-ray datasets and models. In *International Conference on Medical Imaging with Deep Learning*, pages 231–249. PMLR, 2022. 5, 8

[8] Abril Corona-Figueroa, Jonathan Frawley, Sam Bond-Taylor, Sarath Bethapudi, Hubert PH Shum, and Chris G Willcocks. Mednerf: Medical neural radiance fields for reconstructing 3d-aware ct-projections from a single x-ray. In *2022 44th Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pages 3843–3848. IEEE, 2022. 3

[9] Dina Demner-Fushman, Marc D Kohli, Marc B Rosenman, Sonya E Shooshan, Laritza Rodriguez, Sameer Antani, George R Thoma, and Clement J McDonald. Preparing a collection of radiology examinations for distribution and retrieval. *Journal of the American Medical Informatics Association*, pages 304–310, 2016. 3, 5, 6

[10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 5

[11] 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. *arXiv preprint arXiv:2208.01618*, 2022. 4

[12] Rongjun Ge, Yuting He, Cong Xia, Chenchu Xu, Weiya Sun, Guanyu Yang, Junru Li, Zhihua Wang, Hailing Yu, Daoqiang Zhang, et al. X-ctrsnet: 3d cervical vertebra ct reconstruction and segmentation directly from 2d x-ray images. *Knowledge-Based Systems*, 236:107680, 2022. 3

[13] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15180–15190, 2023. 2, 8, 9

[14] Alper Güngör, Salman UH Dar, Şaban Öztürk, Yilmaz Korkmaz, Hasan A Bedel, Gokberk Elmas, Muzaffer Ozbey, and Tolga Çukur. Adaptive diffusion priors for accelerated mri reconstruction. *Medical Image Analysis*, page 102872, 2023. 3

[15] Zhenliang He, Wangmeng Zuo, Meina Kan, Shiguang Shan, and Xilin Chen. Attgan: Facial attribute editing by only changing what you want. *IEEE transactions on image processing*, 28(11):5464–5478, 2019. 7, 8

[16] Tobias Heimann, Bram Van Ginneken, Martin A Styner, Yulia Arzhaeva, Volker Aurich, Christian Bauer, Andreas Beck, Christoph Becker, Reinhard Beichel, György Bekes, et al. Comparison and evaluation of methods for liver segmentation from ct datasets. *IEEE transactions on medical imaging*, 28(8):1251–1265, 2009. 5, 7

[17] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1, 2, 7, 8

[18] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *The Journal of Machine Learning Research*, 23(1):2249–2281, 2022. 2

[19] Xun Huang, Ming-Yu Liu, Serge J. Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. *CoRR*, abs/1804.04732, 2018. 3, 7, 8

[20] Zhongzhen Huang, Xiaofan Zhang, and Shaoting Zhang. Kiut: Knowledge-injected u-transformer for radiology report generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19809–19818, 2023. 6

[21] IXI dataset. https://brain-development.org/ixi-dataset/, 2023. Accessed: 2023-02-14. 5, 7, 8

[22] Lan Jiang, Ye Mao, Xiangfeng Wang, Xi Chen, and Chao Li. Cola-diff: Conditional latent diffusion model for multi-modal mri synthesis. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 398–408. Springer, 2023. 1, 3, 6, 7, 8

[23] Alistair EW Johnson, Tom J Pollard, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Yifan Peng,

Zhiyong Lu, Roger G Mark, Seth J Berkowitz, and Steven Horng. Mimic-cxr-jpg, a large publicly available database of labeled chest radiographs, 2019. 5, 6, 8

[24] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 3, 6, 7

[25] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*, 34:852–863, 2021. 1, 3, 6, 7

[26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5

[27] Hyungyung Lee, Wonjae Kim, Jin-Hwa Kim, Tackeun Kim, Jihang Kim, Leonard Sunwoo, and Edward Choi. Unified chest x-ray and radiology report generation model with multi-view chest x-rays. *arXiv preprint arXiv:2302.12172*, 2023. 6, 8

[28] Chunyuan Li, Xiang Gao, Yuan Li, Baolin Peng, Xiujun Li, Yizhe Zhang, and Jianfeng Gao. Optimus: Organizing sentences via pre-trained modeling of a latent space, 2020. 5

[29] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. *CVPR*, 2023. 6, 7

[30] Fenglin Liu, Xian Wu, Shen Ge, Wei Fan, and Yuexian Zou. Exploring and distilling posterior and prior knowledge for radiology report generation. In *CVPR*, pages 13753–13762, 2021. 6

[31] Ming-Yu Liu, Thomas M. Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. *CoRR*, abs/1703.00848, 2017. 7, 8

[32] Qing Lyu and Ge Wang. Conversion between ct and mri images using diffusion and score-matching models. *arXiv preprint arXiv:2209.12104*, 2022. 3

[33] Payal Maken and Abhishek Gupta. 2d-to-3d: A review for computational 3d image reconstruction from x-ray images. *Archives of Computational Methods in Engineering*, 30(1): 85–114, 2023. 3

[34] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 2

[35] Tufve Nyholm, Stina Svensson, Sebastian Andersson, Joakim Jonsson, Maja Sohlin, Christian Gustafsson, Elisabeth Kjellén, Karin Söderström, Per Albertsson, Lennart Blomqvist, et al. Mr and ct data with multiobserver delineations of organs in the pelvic area—part of the gold atlas project. *Medical physics*, 45(3):1295–1300, 2018. 5

[36] Muzaffer Özbey, Onat Dalmaz, Salman UH Dar, Hasan A Bedel, Şaban Öztürk, Alper Güngör, and Tolga Çukur. Unsupervised medical image translation with adversarial diffusion models. *IEEE Transactions on Medical Imaging*, 2023. 1, 3, 7, 8

[37] Minho Park, Jooyeol Yun, Seunghwan Choi, and Jaegul Choo. Learning to generate semantic layouts for higher text-image correspondence in text-to-image synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7591–7600, 2023. 6, 7

[38] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 5

[39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 4

[40] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1 (2):3, 2022. 1, 2

[41] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 2, 3, 4, 5, 6, 7, 8, 9

[42] 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 6, 8

[43] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 1

[44] Robin San-Roman, Eliya Nachmani, and Lior Wolf. Noise estimation for generative diffusion models. *arXiv preprint arXiv:2104.02600*, 2021. 2

[45] Hiroshi Sasaki, Chris G Willcocks, and Toby P Breckon. Unit-ddpm: Unpaired image translation with denoising diffusion probabilistic models. *arXiv preprint arXiv:2104.05358*, 2021. 3, 7, 8

[46] Bradley Segal, David M Rubin, Grace Rubin, and Adam Pantanowitz. Evaluating the clinical realism of synthetic chest x-rays generated using progressively growing gans. *SN Computer Science*, 2(4):321, 2021. 1, 3, 6, 7

[47] Anmol Sharma and Ghassan Hamarneh. Missing mri pulse sequence synthesis using multi-modal generative adversarial network. *IEEE transactions on medical imaging*, 39(4): 1170–1183, 2019. 3, 6, 7

[48] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015. 2

[49] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 2, 5

[50] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 2

[51] Sanjay Subramanian, Lucy Lu Wang, Sachin Mehta, Ben Bogin, Madeleine van Zuylen, Sravanthi Parasa, Sameer Singh, Matt Gardner, and Hannaneh Hajishirzi. Medicat: A dataset of medical images, captions, and textual references. *CoRR*, abs/2010.06000, 2020. 5

[52] Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. Any-to-any generation via composable diffusion. *arXiv preprint arXiv:2305.11846*, 2023. 2, 8, 9

[53] Gu Tiancheng, Liu Dongnan, Li Zhiyuan, and Cai Weidong. Complex organ mask guided radiology report generation. *arXiv preprint arXiv:2311.02329*, 2023. 6

[54] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2097–2106, 2017. 1, 3, 5, 7

[55] Zhanyu Wang, Lingqiao Liu, Lei Wang, and Luping Zhou. Metransformer: Radiology report generation by transformer with multiple learnable expert tokens. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11558–11567, 2023. 6

[56] Xingqian Xu, Zhangyang Wang, Gong Zhang, Kai Wang, and Humphrey Shi. Versatile diffusion: Text, images and variations all in one diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7754–7765, 2023. 2, 5, 8, 9

[57] An Yan, Julian McAuley, Xing Lu, Jiang Du, Eric Y Chang, Amilcare Gentili, and Chun-Nan Hsu. Radbert: Adapting transformer-based language models to radiology. *Radiology: Artificial Intelligence*, 4(4):e210258, 2022. 1, 6, 8

[58] Bin Yan and Mingtao Pei. Clinical-bert: Vision-language pre-training for radiograph diagnosis and reports generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2982–2990, 2022. 1, 6

[59] Di You, Fenglin Liu, Shen Ge, Xiaoxia Xie, Jing Zhang, and Xian Wu. Aligntransformer: Hierarchical alignment of visual regions and disease tags for medical report generation. In *MICCAI*, 2021. 1, 6

[60] Mahmut Yurt, Muzaffer Özbey, Salman UH Dar, Berk Tinaz, Kader K Oguz, and Tolga Çukur. Progressively volumetrized deep generative models for data-efficient contextual learning of mr image recovery. *Medical Image Analysis*, 78:102429, 2022. 1, 3, 6, 7

[61] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *International Conference on Machine Learning*, pages 12310–12320. PMLR, 2021. 4

[62] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *International conference on machine learning*, pages 7354–7363. PMLR, 2019. 7, 8

[63] Tao Zhou, Huazhu Fu, Geng Chen, Jianbing Shen, and Ling Shao. Hi-net: hybrid-fusion network for multi-modal mr image synthesis. *IEEE transactions on medical imaging*, 39(9): 2772–2781, 2020. 3, 6, 7