

FgC2F-UDiff: Frequency-guided and Coarse-to-fine Unified Diffusion Model for Multi-modality Missing MRI Synthesis

Xiaojiao Xiao, Qinmin Vivian Hu, and Guanghui Wang, *Senior Member, IEEE*

Abstract—Multi-modality magnetic resonance imaging (MRI) is essential for the diagnosis and treatment of brain tumors. However, missing modalities are commonly observed due to limitations in scan time, scan corruption, artifacts, motion, and contrast agent intolerance. Synthesis of missing MRI has been a means to address the limitations of modality insufficiency in clinical practice and research. However, there are still some challenges, such as poor generalization, inaccurate non-linear mapping, and slow processing speeds. To address the aforementioned issues, we propose a novel unified synthesis model, the Frequency-guided and Coarse-to-fine Unified Diffusion Model (FgC2F-UDiff), designed for multiple inputs and outputs. Specifically, the Coarse-to-fine Unified Network (CUN) fully exploits the iterative denoising properties of diffusion models, from global to detail, by dividing the denoising process into two stages—coarse and fine—to enhance the fidelity of synthesized images. Secondly, the Frequency-guided Collaborative Strategy (FCS) harnesses appropriate frequency information as prior knowledge to guide the learning of a unified, highly non-linear mapping. Thirdly, the Specific-acceleration Hybrid Mechanism (SHM) integrates specific mechanisms to accelerate the diffusion model and enhance the feasibility of many-to-many synthesis. Extensive experimental evaluations have demonstrated that our proposed FgC2F-UDiff model achieves superior performance on two datasets, validated through a comprehensive assessment that includes both qualitative observations and quantitative metrics, such as PSNR SSIM, LPIPS, and FID. The source code is available at <https://github.com/xiaojiao929/FgC2F-UDiff>.

Index Terms—Diffusion model, Frequency, Synthesis, Multi-modality.

I. INTRODUCTION

MULTI-MODALITY magnetic resonance imaging (MRI), encompassing T1, T2, FLAIR, and T1 contrast-enhanced (T1ce) sequences, is indispensable for the diagnosis, monitoring, and treatment of brain tumors, as it provides complementary information about tissue views and spatial details [1], [2], [3]. As shown in Fig.1, T1 images provide anatomical structure, FLAIR highlights the entire tumor region, T2 delineates a clear outline of the tumor edema area, and T1ce depicts a clear boundary of enhanced areas. However, the absence of modalities across different clinical centers is unavoidable [4]. Furthermore,

This work is partly supported by the Natural Sciences and Engineering Research Council of Canada (NSERC) and TMU FOS Postdoctoral Fellowship.

X. Xiao and Q. Hu are with the Department of Computer Science, Toronto Metropolitan University, Toronto, Canada (e-mail: xiaojiao@torontomu.ca, vivian@torontomu.ca)

Guanghui Wang with the Department of Computer Science, Toronto Metropolitan University, Toronto, Canada (e-mail: wangcs@torontomu.ca).

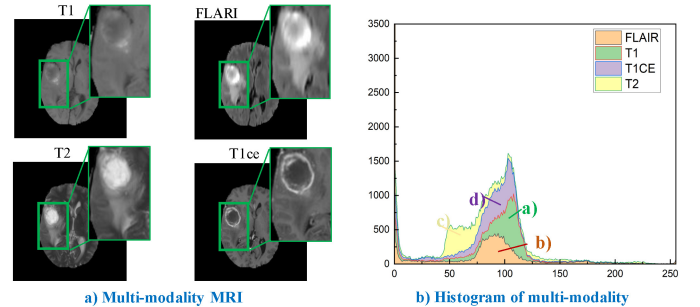


Fig. 1. Different image modalities provide different types of tissue contrast views and spatial resolution, which has variety in the histogram.

modalities are often missing due to limitations in scan time, scan corruption, artifacts, motion, and contrast agent intolerance [5], [6]. These factors restrict the ubiquity of multi-modality, adversely affecting various downstream tasks such as segmentation, detection, and quantification. Consequently, the synthesis of missing multi-modality MRI has garnered increasing research interest as a means to overcome the limitations of modality insufficiency in clinical practice and research.

In recent years, numerous studies have provided evidence of the effectiveness of deep learning in synthesizing missing modality [7], [8], [9], [10], [11], [12], [13], [14], [15]. Task-specific models, including one-to-one and many-to-one, focus on learning individual non-linear mappings from source to target imaging modalities, neglecting unique information present in shared features. These limitations can compromise synthesis performance and fail to meet the clinical needs for multi-modality applications. GAN-based methods utilize adversarial losses to better capture organizational structure and further enhance synthesis quality. Moreover, some unified method of GAN-based achieves many-to-many synthesis [16], [17], [18] to address the limitation of one modality synthesis. However, GAN-based methods suffer from mode collapse, non-convergence, instability, and high sensibility to hyperparameters [19], [20].

Denoising diffusion probabilistic models (DDPM) may stably generate high-quality dependable images to improve MRI synthesis as a promising alternative to GAN [21]. Diffusion models show superiority in a wide variety of areas, ranging from generative modeling tasks (e.g., image generation [22], image super-resolution [23], [24], [25]) to discriminative tasks

(e.g., image segmentation [26], [27], classification [28], [29], and anomaly detection [30], [31]) while reducing risk of modality collapse. The work [32] proposed a conditioned latent diffusion model for many-to-one synthesis, which preserves anatomical structure with accelerated sampling. However, only learning the mapping of a single task (i.e., many to one) is unable to meet clinical needs for multi-modality. And, the long time steps ($t=1000$) are unable to solve the slow sampling problem of DDPM [33]. Besides, significant challenges persist that limit the widespread application of existing methods: 1) Lack of an effective unified model to synthesize high-fidelity images from many to many. The existing methods suffer from poor generalization ability for atypical anatomical features, which limits the fidelity of synthesis, especially details. 2) Lack of effective strategy to guide the learning of unified highly non-linear mapping between multi-modality MRI. During the iterative denoising process, the randomness of noise alters the original distribution of the target image [34]. During the iterative denoising process, the randomness of the noise changes the original distribution of the target image. Therefore, the learned nonlinear mapping relationships are unable to accurately reflect the consistency of anatomical structures. 3) Lack of feasibility of many-to-many because of the slow sampling problem of DDPM.

Motivation: To address the above-mentioned challenges, we propose a unified diffusion model, the Frequency-guided and Coarse-to-fine Unified Diffusion Model (FgC2F-UDiff) to synthesize the missing modality for multiple inputs and outputs. Our key hypothesis is to decompose the frequency domain into low-high-frequency as guidance information in coarse-to-fine stages based on the iterative denoising properties of diffusion models. Specifically, FgC2F-UDiff relies primarily on low-frequency with global anatomical features to guide coarse denoising and subsequently performs fine denoising progress guided by high-frequency with texture and detail. FgC2F-UDiff enables the effective synthesis of target modalities with high fidelity from any combination of modalities.

Coarse-to-fine Unified Network (CUN) fully utilizes iterative denoising properties of diffusion model (global-to-detail), novelty dividing the denoising process into two stages (i.e., coarse-to-fine) to improve the fidelity of denoised images. Moreover, CUN based on the diffusion model provides an effective unified program for cross-modality missing synthesis of multiple inputs and outputs.

Frequency-guided Collaborative Strategy (FCS) guides the learning of accurate non-linear mapping by enhancing the diversity of prior knowledge. Specifically, inspired by image signals, we incorporate low-frequency with global anatomical features in the early stages of coarse denoising and introduce high-frequency information with local fine-grained (i.e., texture and details) in the later stages of fine denoising. Therefore, frequency domain information of different granularity is decomposed and used to guide denoising collaboratively as prior knowledge, so the synthesized image has diverse features and similar edge and detail characteristics to the real image. At the same time, novel strategies were designed to dynamically search for appropriate frequency domain information to maximize the information of available modalities.

Specific-acceleration Hybrid Mechanism (SHM) is designed for specific tasks to accelerate the diffusion model and improve the fidelity of synthesized images. First, the curriculum learning (CL) mechanism is employed to simulate the easy-to-hard learning of missing modalities. Second, the network divided into two coarse-to-fine phases fits the iterative denoising properties of diffusion models, thus accelerating the synthesis process from the whole image to the details. Finally, dynamically selecting constraint conditions of frequency ensures the maximization of information from available modalities, guaranteeing the learning of non-linear mapping of image texture and fine-detail structures. Consequently, the trained model is adaptable to any number of original modalities and exhibits increased robustness in specific complex regions of images, enhancing the feasibility and synthesis performance of many-to-many FgC2F-UDiff.

Our contributions include the following:

- To the best of our knowledge, this is the first work to introduce a unified diffusion model guided by a frequency domain, which provides an effective cross-modality synthesis mechanism for multiple inputs and outputs.
- We propose an innovative frequency-domain-guided coarse-to-fine network that effectively incorporates the iterative denoising characteristics of the diffusion model. This approach strategically shifts guidance across the appropriate frequency domains from coarse to fine, enhancing the fidelity of synthesized images.
- We propose an efficient mechanism, SHM, which intelligently blends specific mechanisms to accelerate the diffusion model and improve the feasibility of many-to-many.

II. RELATED WORK

Multi-domain synthesis of medical images provides a promising solution to address the limitations of modality insufficiency, which has attracted significant interest and gained popularity in recent years. Many research works and various technologies have been presented in the multi-domain synthesis of medical images. This section briefly reviews known synthesis methods by categorizing them into task-specific models (i.e., one-to-one and many-to-one) and unified models (i.e., many-to-many).

A. Task-specific model for missing image synthesis

a) One-to-one: Earlier one-to-one studies have proposed patch-based regression [35], [36], [37], sparse dictionary representation [38], [39], and atlas [40], [41]. However, handcrafted features constrain the performance and development of these traditional methods. To improve the automatic extract feature, deep learning (DL) has been employed in cross-modality synthesis [42], [43], [44]. For instance, the work of [42] developed a patch-based location-sensitive deep network (LSDN), which combines intensity and spatial information for synthesizing T2 MRI from T1 MRI and vice versa. The work of [43] proposed a deep encoder-decoder image synthesizer (DEDIS) for whole image synthesis. Despite yielding enhancements, CNN-based has the drawback of losing detailed structural

information [7]. GAN-based achieved great success with the development of deep learning techniques [11], [9], [10], [45], [46]. For instance, the work of [9] designed CoCa-GAN for synthesizing MRI data (i.e., T2, FLAIR, and T1ce) from T1, which utilizes adversarial learning and context-aware learning to learn common feature spaces. The work of [11] proposed a unified GAN, which learns the modality-invariant features by modality translation for segmentation tasks.

b) **Many-to-one**: Earlier many-to-one studies have proposed patch-based regression [47], [39]. Later studies also used DL-based, which also achieved great success [48], [49], [50]. For instance, the work of [50] designed a Hybrid-fusion Network (Hi-Net) for multi-modal MR image synthesis, which employed a fusion network to learn the common latent representation of multi-modality data. The work of [49] proposed a multi-modality synthesis framework, which fused disentangled content code from each modality into a shared representation via gated feature fusion. Recently, GAN-based methods were demonstrated to outperform other DL-based methods in many-to-one tasks [10], [51], [52], [14], [53], [54], [55]. For instance, the work of [10] presents EaGANs to generate T2 and FLAIR images from T1. The EaGAN captured the edges of key texture information, and two GAN variants are proposed to integrate the edge information through different learning strategies. Lee et. al[51] proposed a CollaGAN framework for missing image data imputation, which converts the image imputation problem to multi-domain images-to-image translation tasks. The work of [12] generated any missing modality in a single unified Auto-GAN model, which performs self-supervised learning to learn multi-facet information, further guaranteeing its generalizability.

However, when an insufficient number of modalities are available, especially when many modalities (e.g., three modalities) do not exist, applying the above strategies can not necessarily ensure that the lost data is recovered since there are not enough features to reconstruct the missing data [56].

B. Unified model for missing image synthesis

Unified synthesis methods take multiple inputs and generate multiple inputs and outputs, which is a relatively new study in synthesis tasks. Several studies have attempted to propose a unified model on many-to-many synthesis that can exploit all available data. For instance, the work of [16] proposed a multi-modality generative adversarial network (MM-GAN), which was one of the first to propose a multi-input and multi-output architecture that generalizes to any combination of available and missing modalities. The work of [17] proposed an adversarial model with a residual vision transformers (ResViT) generator to translate between multi-modal imaging data. The work of [18] exploits the commonality information of available modalities for unified multi-modal image synthesis based on GAN. However, the above methods are all based on GAN [57], which has some common issues while training GAN, such as mode collapse, non-convergence, instability, and high sensibility to hyperparameters, thus limiting the fidelity and diversity of synthesized images [19], [20]. The work of [58] proposed unified Multi-modal Modality-masked Diffusion

Network (M2DN), tackling multi-modal synthesis from the perspective of “progressive whole-modality inpainting”, instead of “cross-modal translation”. However, it only takes the available modes as conditions and does not take into account the denoising iterative properties of the diffusion model.

III. METHODOLOGY

As shown in Fig.2, the proposed FgC2F-UDiff integrates any available number of source image modalities (i.e., T1, T2, FLAIR, and T1ce) for coarse-to-fine synthesizing missing target modalities $X_0^m \in \mathbb{R}^{H \times W}$, where H and W represent the height and width, respectively. The FgC2F-UDiff works via a forward diffusion process in Section III-A and a coarse-to-fine reverse denoising process. Specifically, the coarse-to-fine unified network (CUN) divides the denoising process into two stages of coarse-to-fine for improving the fidelity of synthesizing image in Section III-B. Among them, the frequency-guided collaborative strategy (FCS) decomposes the frequency information to guide the learning non-linear mapping of many-to-many. It utilizes low-frequency and high-frequency to enhance the realism of the synthesized image structures in Section III-C. The entire network benefits from a specific-acceleration hybrid mechanism (SHM) to accelerate the time steps to improve the availability of many-to-many FgC2F-UDiff in Section III-D.

A. Forward diffusion process

The forward diffusion process of FgC2F-UDiff is defined as a Markov chain as DDPM [21], which maps between source samples and pure noise samples. Formally, given the multi-modality of S_0^m , $m \in \{T1, T2, FLAIR, T1ce\}$, among missing-images are $X_0^m \in S_0^m$, $X_0^m \sim \mathbb{Q}(X_0^m)$. The FgC2F-UDiff gradually adds Gaussian noise to X_0^m and obtains a pure noise sample X_T^m with time steps $T \in \{1, 2, \dots, t+1, \dots, T-1, T\}$. At the time step t , the noisy X_t^m can be formulated as:

$$X_t^m = \sqrt{1 - \beta_t} X_{t-1}^m + \beta_t \epsilon_{t-1} \quad (1)$$

where $\beta_t \sim (0, I)$ is the variance of the Gaussian noise added at time step T , $\epsilon \sim \mathcal{N}(0, 1)$ is Gaussian distribution noise. Thus, the forward diffusion process can be formulated as:

$$\mathbb{Q}(X_t^m | X_{t-1}^m) := \mathcal{N}(X_t^m; \sqrt{1 - \beta_t} X_{t-1}^m, \beta_t I) \quad (2)$$

where I denotes the standard normal distribution. Using the notation $a_t := 1 - \beta_t$ and $\bar{a}_t = \prod_{s=1}^t a_s$, the forward process admits sampling X_t^m at an arbitrary timestep t can be deduced by Eq.1 and Eq.2:

$$\mathbb{Q}(X_t^m | X_0^m) = \mathcal{N}(X_t^m; \sqrt{\bar{a}_t} X_0^m, (1 - \bar{a}_t) I) \quad (3)$$

B. Coarse-to-fine Unified Network (CUN)

To improve the performance of the many-to-many synthesis, FgC2F-UDiff designed a coarse-to-fine collaborative reverse diffusion process as shown in Fig.2. Specifically, in the light of the iterative denoising properties of diffusion models (global-to-detail) in the diffusion model, our reverse diffusion process is divided into two phases: a coarse denoising process from

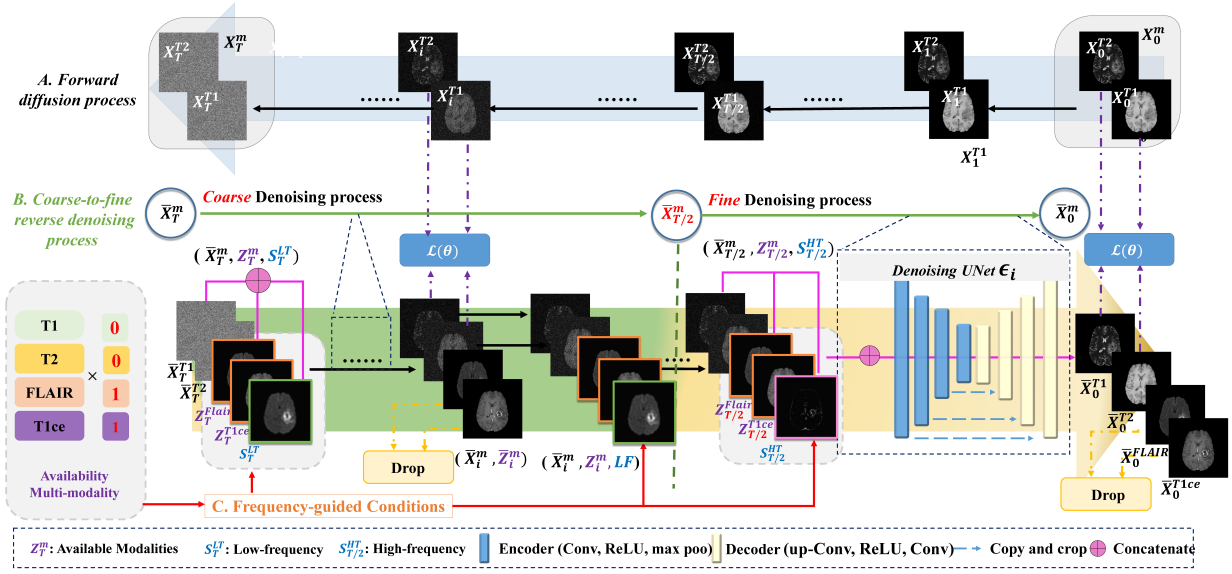


Fig. 2. Overview of the proposed FgC2F-UDiff, which cross-modality synthesizes missing modalities from multiple inputs and outputs. It includes forward diffusion progress and coarse-to-fine reverse denoising progress. The FgC2F-UDiff decomposes the frequency domain into low-high-frequency as guidance information in coarse-to-fine stages based on the iterative denoising properties of diffusion models.

T to $\frac{T}{2} - 1$, followed by a fine denoising process from $\frac{T}{2}$ to 0. Simultaneously, FgC2F-UDiff leverages the characteristics of image signals (high- and low-frequency), by decomposing the frequency domain to guide the synthesis in a staged manner. Additionally, FgC2F-UDiff fully leverages all available modalities as conditions to enhance and expedite the denoising process. As a result, our FgC2F-UDiff possesses the capability to cross-modality synthesis mechanism for multiple inputs and outputs.

Reverse denoising diffusion process. Since the reverse of the forward process is intractable, DDPM learns parameterized Gaussian transitions. Given X_t^m and corresponding conditional of all conditions C_t^m , in each time step of the reverse process, the denoising operation is performed on the noisy multi-channel image (X_t^m, C_t^m) to obtain the previous image X_{t-1}^m . The probability distribution of X_{t-1}^m under the condition X_t^m can be formulated as:

$$P(X_{t-1}^m | X_t^m, C_t^m) := \mathcal{N}(X_{t-1}^m; \mu_\theta(X_t^m, t, C_t^m), \sigma_\theta(X_t^m, t, C_t^m)I) \quad (4)$$

where σ_θ is the variance of conditional distribution $P(X_{t-1}^m | X_t^m, C_t^m)$, which can be formulated as:

$$\sigma_\theta = \frac{1 - \bar{a}_{t-1}}{1 - \bar{a}_t} \beta_t \quad (5)$$

where $\beta_t = 1 - \bar{a}_t$. The generative process is expressed as:

$$X_{t-1}^m = \frac{1}{\sqrt{\bar{a}_t}} \left(X_t^m - \frac{\beta_t}{\sqrt{1 - \bar{a}_t}} \epsilon_\theta(X_t^m, t, C_t^m) + \sigma_\theta(X_t^m, t, C_t^m) Z, Z \sim \mathcal{N}(0, 1) \right) \quad (6)$$

where ϵ_θ represents noise approximation.

Denoising models based on the UNet [59], which is widely used in a diverse range of segmentation and synthesis tasks due to the U-shaped symmetrical structure and skip connection between the encoder and decoder. The architecture is shown

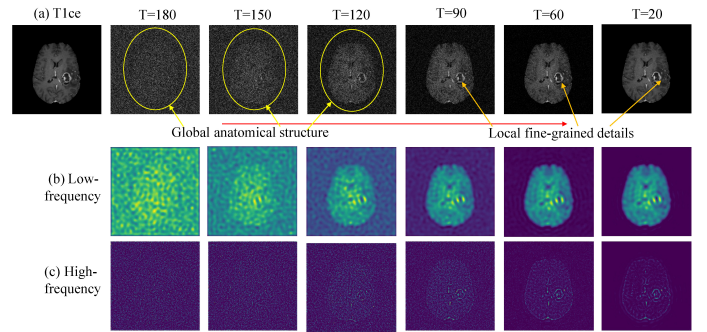


Fig. 3. Visualize analyzing and visualizing the denoising synthesis images, which shows the significant properties of iterative denoising. Specifically, (a) shows denoised synthesis images corresponding to different time steps T . (b) shows the low frequency. (c) shows the high frequency.

in Fig.2. FgC2F-UDiff is trained to predict a denoised variant of their input X_T^m , where X_t^m is a noisy version of the input X_0^m . Specifically, the network is designed with a connection input $(C_T \in \mathbb{R}^{(H \times W) \times N_c})$ with 5-channels, where channel $N_c = 0, 1, 2, 3$ and 4 corresponds to $T1, T2, FLAIR, T1ce$ and frequency-guided image, respectively. The input image C_T connects all missing-images X_T^m , and conditional images C_T^m (Z_T^m, S_T^{LF} or $S_{T/2}^{HF}$). In the frequency-guided conditions module (as shown in Fig.2.C), all modalities S_0^m have been sorted according to $T1, T2, FLAIR, T1ce$ to obtain the corresponding 4-digit numbers, the numeral “1” signifies the presence of a particular module, while “0” indicates its absence. After S_0^m multiplied by the corresponding digit number, the channels corresponding to each missing modality are inputted with noisy images of X_T^m , the channels corresponding to each condition modality are inputted with source images as Z_T^m , and the channel corresponding to conditional frequency is inputted with S_T^{LF} during coarse denoising process, while

S^{HF} during fine denoising process. For instance, if sequences $T1$ and $T2$ are missing, channels $N_c=0$ and $N_c=1$ are fed with noisy X_T^{T1} and X_T^{T2} , respectively. $N_c=2$ and $N_c=3$ are fed with original images Z_T^{FLAIR} and Z_T^{T1ce} , respectively. And, the last $N_c=4$ is fed with corresponding frequency image S^{LF} . After one time step, UNet outputs 4 channels corresponding to 4 modalities. FgC2F-UDiff calculates the difference between the output \bar{X}_{T-1}^m and X_{T-1}^m to train the network, while copying the output \bar{X}_{T-1}^m to enter the next iteration. To ensure effective guidance of available images, the synthesized \bar{Z}_{T-1}^m is dropped. The encoder is a traditional stack of 3×3 convolution and 2×2 max pooling layers. And, the symmetric decoder is a traditional stack of 2×2 up-conv, copy, and crop, and 3×3 convolution. To modify network performance, each convolutional layer follows a BN layer and a ReLU layer.

The FgC2F-UDiff is trained to synthesize the target modality by predicting the involved noise ϵ_θ under the guidance of the C_t^m , which is formulated below:

$$L_{FgC2F-UDiff} = E_{x_t, \epsilon \sim N(0, I), t} \|\epsilon - \epsilon_\theta(X_t^m, t, C_t^m)\|_2^2 \quad (7)$$

Discussion: What is the iterative denoising properties of diffusion models (global-to-detail) in the FgC2F-UDiff.

After analyzing and visualizing the denoising synthesis images, as shown in Fig.3, we can observe significant properties with iterative denoising. Specifically, (a) shows denoised synthesis images corresponding to different time steps T . In the early stage of denoising, it initially formed the global anatomical structure of the brain (as shown in the yellow circle). In the later stage, we are gradually synthesizing local fine-grained details such as the edges and texture of the tumor (as indicated by the orange arrow). These discrepancies gradually diminish from the global to the local details as the iteration of the denoising process T . (b) and (c) display the low-frequency and high-frequency images corresponding to the denoising image. In low-frequency images during the initial denoising phase (T to $\frac{T}{2} - 1$), the images exhibit clear and diverse global features, closely related to the anatomical structure of the brain. In high-frequency information images during the initial denoising phase ($\frac{T}{2}$ to 0), fine-grained features become increasingly evident and sharp, particularly around the edges of the brain and tumor. There is a higher demand for detailed features in later synthesis. Therefore, the proposed CUN fully considers the iterative denoising properties of diffusion models and the character of the frequency domain, divides denoising into two phases, and adds corresponding low-frequency and high-frequency information to guide the denoising process.

C. Frequency-guided Collaborative Strategy (FCS)

To learn the high non-linear mapping of many-to-many, especially anatomical structures, we designed an FCS strategy guided by frequency. Specifically, Fig.3 has verified the iterative denoising properties of diffusion models from global to local. So, FgC2F-UDiff incorporates low-frequency with global anatomical features in the early stages of coarse denoising from T to $\frac{T}{2} - 1$. And introducing high-frequency information with fine-grained (i.e., texture and details) as prior

knowledge in the later stages of fine denoising from $\frac{T}{2}$ to 0. The coarse and fine denoising stages work collaboratively in two phases to ensure the learning of the unified distribution of data and the production of high-quality synthesized images.

To find the most suitable frequency domain information as prior guidance knowledge, we have designed two strategies to dynamically search frequency information. Our FCS dynamically selects frequency information tailored to the available modalities at each stage of the different subject, employing a left-to-right scan for coarse structural low-frequency guidance and a right-to-left scan for fine detail high-frequency enhancement. This approach ensures optimal guidance for the diffusion process, adapting to the unique needs of each subject. Specifically, all modality S_0^m have been sorted according to $T1$, $T2$, $FLAIR$, $T1ce$ to obtain the corresponding 4-digit numbers, the numeral “1” signifies the presence of a particular module, while “0” indicates its absence. As shown in Fig.2.B, the missing modalities, when multiplied by zero, do not contribute any information to the denoising process, effectively excluding them from the calculation. Conversely, the available modalities are multiplied by one, preserving their original image data intact for further processing. Then, the “left to right” strategy extracts low-frequency by searching the corresponding modality of the first available image from left to right and defined as S^{LF} . And, the “right to left” strategy to extract high-frequency by searching the corresponding modality of the first available image from right to left and defined as S^{HF} .

To filter the images into their respective frequency domains, we apply Gaussian low pass filters (GLPF) [60] and Gaussian high pass filters (GHPF) [61] to process S^{LF} and S^{HF} , respectively. Specifically, we set that Gaussian kernel as:

$$\kappa_\sigma [i, j] = \frac{1}{2\pi\sigma^2} e^{-\frac{1}{2}(\frac{i^2+j^2}{\sigma^2})} \quad (8)$$

where $[i, j]$ is the original point, and σ is used to measure the width of the Gaussian curve. After the GLPF filter on S^{LF} , we obtain the low frequency image (LF):

$$LF [i, j] = \sum_m \sum_n \kappa [m, n] \cdot S^{LF} [i + m, j + n] \quad (9)$$

where m, n represents the index of GLPF. Then, the high frequency image (HF) of S^{HF} is expressed as:

$$HF [i, j] = 1 - \sum_m \sum_n \kappa [m, n] \cdot S^{HF} [i + m, j + n] \quad (10)$$

Discussion: Why design different strategies to dynamically search frequency domain information?

As shown in Fig.4, inspired by the character of image signals, the image can usually be decomposed into high-frequency sub-band and low-frequency sub-band. The high-frequency sub-bands contain more details and edge information, whereas the low-frequency sub-band contains the contour and structure information of images. And, the different modalities generated through scanning parameters usually provide different information [2]. Such as T1 brain images delineated

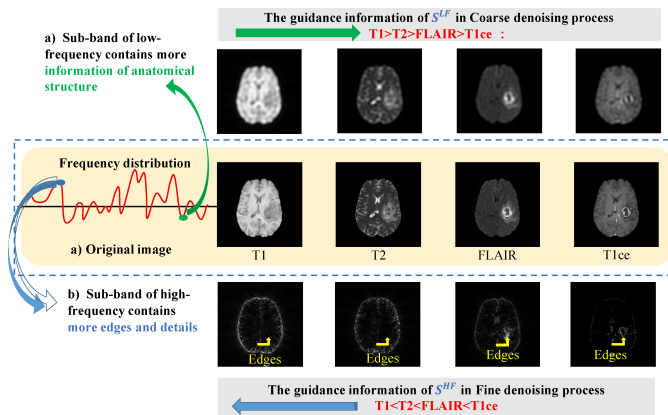


Fig. 4. The image can usually be decomposed into a high-frequency sub-band with edges and details and a low-frequency sub-band with anatomical structure.

low-frequency with anatomical structure, while T1ce provides a high-frequency with a clear boundary between the enhanced areas. So, sorting the four modalities according to the amount of low-frequency information they contain can be obtained: $T1 > T2 > FLAIR > T1ce$, whereas sorting the four modalities according to the amount of high-frequency information they contain can be obtained: $T1 < T2 < FLAIR < T1ce$. Therefore, we novelty designed two different selection methods according to the character of modality, that is, from "left to right" to find the best low-frequency information and from "right to left" to find the best high-frequency information. Through such a dynamic search strategy, the best guidance can be provided for the diffusion process.

D. Specific-acceleration Hybrid Mechanism (SHM)

Our SHM is designed to accelerate the sampling speed of the diffusion model to improve the feasibility of many-to-many FgC2F-UDiff. Specifically, first, the curriculum learning (CL) mechanism is employed to simulate the easy-to-hard learning of missing modalities. Second, the coarse-to-fine denoising process guided by frequency fits the iterative denoising properties of diffusion models to accelerate the synthesis process from the whole image to the details. Third, dynamically selecting high-frequency and low-frequency constraint conditions ensures the maximization of information from available modalities.

Curriculum learning (CL) mechanism. To ensure the adaptability of our FgC2F-UDiff for the synthesis of missing data across diverse inputs, we employ curriculum learning (CL) [62] as a rationality-enhancing training strategy. Due to the varying degrees of data loss and difficulty in obtaining complete modalities, it is imperative to devise a training strategy that effectively leverages available data and expedites model convergence, ultimately yielding higher-quality synthesized results. CL shares some similarities with boosting algorithms, where the focus is gradually shifted towards more challenging examples. However, unlike a uniform distribution of training data, CL begins by emphasizing easier examples and progressively introduces more complex instances as the

training process unfolds. In the context of CL-based training within our FgC2F-UDiff framework, we categorize the missing sequences into different difficulty levels. Specifically, we designate the task of synthesizing one missing sequence as "easy", followed by tackling the challenge of synthesizing two missing sequences categorized as "moderate", and finally, addressing the most demanding scenario of synthesizing all three missing sequences, denoted as "hard". This tiered approach to CL ensures that the model is systematically exposed to increasingly complex situations, enabling it to learn and adapt effectively across a range of input conditions.

IV. DATASET AND EVALUATION METRICS

A. Dataset

FgC2F-UDiff framework evaluated the performance on BraTS 2021 and IXI brain image datasets.

a) *Brain Tumor Segmentation Challenge 2021 (BraTS 2021)*: The BraTS 2021 [63], [1], [64] contains 1,251 cases, consisting of four different MRI sequences per case (i.e., T1, T2, FLAIR, and T1ce), acquired with different protocols and various scanners from multiple institutions. Standardized preprocessing has been applied to all the sequences. Specifically, the dimension of each data is resampled to $240 \times 240 \times 150$, and the intensity is normalized to the range $[-1, 1]$. More details about the preprocessing information can be found in the original publication [1].

b) *Information Extraction from Images (IXI)*: The IXI dataset [65] contains nearly 600 MRIs from normal and healthy subjects, consisting of three different MRI sequences (i.e., T1, T2, and PD-weighted). The images were acquired with the following parameters (T1 image: TE = 4.603 ms, TR = 9.813 ms, spatial resolution = $0.94 \times 0.94 \times 1.2 \text{ mm}^3$, matrix size = $256 \times 256 \times 150$. T2 image: TE = 100 ms, TR = 8178.34 ms, spatial resolution = $0.94 \times 0.94 \times 1.2 \text{ mm}^3$, matrix size = $256 \times 256 \times 150$. PD-weighted image: TE = 8 ms, TR = 8178.34 ms, spatial resolution = $0.94 \times 0.94 \times 1.2 \text{ mm}^3$, matrix size = $256 \times 256 \times 150$). Note that the multi-contrast images in this dataset were unregistered. Therefore, T2 and PD-weighted images were spatially registered onto T1-weighted images before modeling by rigid transformation. Registration was performed via an affine transformation in FSL [66] based on mutual information.

B. Implementation details

We employed five-fold cross-validation to train and test the grading. For each cross-validation split, the dataset was divided into a training/validation/testing as 7:1:2. Our network was implemented on Ubuntu 20.04 platform, using Python v3.6 and PyTorch v0.4.0, and was run on 2 NVIDIA GTX 3090Ti GPUs with 24 GB memory. FgC2F-UDiff are optimized using Adam optimizer [67] with a learning rate of 0.0001. Following the noise schedules of DDPM [21] and set the value of time steps T to 200. Following the work of [61], the kernel size κ in Eq.(8) was set to 21.

C. Evaluation index

The performance of FgC2F-UDiff is evaluated by four standard measures, including peak signal-to-noise ratio (PSNR), structural similarity index (SSIM) [68], learned perceptual image patch similarity (LPIPS), and fréchet inception distance (FID), which reflect the quality and variety of synthetic images. The significance of performance differences was evaluated with signed-rank tests ($p < 0.05$). The evaluation criteria are defined as:

$$PSNR = 10 \cdot \log_{10} \left(\frac{MAX_I^2}{\sqrt{MSE}} \right) \quad (11)$$

$$SSIM = \frac{(2Avg_x Avg_y + C1)(2\delta_{xy} + C2)}{(Avg_x^2 + Avg_y^2 + C1)(\delta_x^2 + \delta_y^2 + C2)} \quad (12)$$

where MAX_I represents the max value among the pixels in brain MRI with a size of $m \times n$. x and y are two images to be compared. Avg_x and Avg_y are the average pixel value of x and y , respectively. δ_x^2 and δ_y^2 are the variance of x and y , respectively. And δ_{xy} is the covariance of x and y . $C1 = (k_1 L)^2$, and $C2 = (k_2 L)^2$ are two constants, avoiding division by zero. L is the range of pixel values, $k_1 = 0.01$ and $k_2 = 0.03$ are the default.

$$LPIPS(x, y) = \sum_l w_l \cdot \frac{1}{H_l W_l} \sum_{h,w} \|\phi_l(x)_{h,w} - \phi_l(y)_{h,w}\|_2^2 \quad (13)$$

where $\phi_l(x)$ and $\phi_l(y)$ are the feature activations at layer l of the network for images x and y , respectively. H_l and W_l are the height and width of the feature maps at layer l . And, w_l are the learned weights for the features at layer l , optimized to align with human perceptual difference. $\|\cdot\|_2$ denotes the Euclidean distance.

$$FID(p, q) = \|\mu_p - \mu_q\|^2 + \text{Tr}(\Sigma_p + \Sigma_q - 2(\Sigma_p \Sigma_q)^{1/2}) \quad (14)$$

where p and q represent the distributions of features extracted from the real and generated images, respectively. And, μ_p, μ_q are the mean vectors of the features from distributions p and q . Σ_p, Σ_q are the covariance matrices of the features from distributions p and q . Tr denotes the trace of a matrix, which is the sum of the elements on the main diagonal. $\Sigma_p \Sigma_q^{1/2}$ represents the square root of the product of the covariance matrices, used to calculate the similarity between the two distributions. In this study, the FID was computed using a singular comprehensive evaluation of the model-generated images against the reference dataset.

D. Comparison settings

To demonstrate the superiority of our proposed framework, FgC2F-UDiff is compared with other methods on two datasets. The baseline methods include the task-specific models (pix2pix [7], pGAN [69], LDM [70], and CoLa-Diff [32]) and unified models (MM-GAN [16], ResVit [17], and UniGAN [18]). The hyperparameters of each competing method were optimized via identical cross-validation procedures.

V. EXPERIMENTS

The experiment results show that FgC2F-UDiff achieves high performance in both task-specific and unified synthesis regarding PSNR, SSIM, LPIPS, and FID. A set of experiments were performed to evaluate the performance of FgC2F-UDiff, including (1) synthesis results of the proposed FgC2F-UDiff in Section V-A; (2) performance comparison of task-specific synthesis with state-of-the-art (SOTA) methods in Section V-B; (3) performance comparison of unified synthesis models with state-of-the-art (SOTA) methods in Section V-C; (4) ablation studies of FgC2F-UDiff in Section V-D; and (5) analysis of SHM in Section V-E

A. Synthesis results of the proposed method

The visual qualitative results of FgC2F-UDiff are shown in Fig.5. The four-bit digits in the figures indicate the availability conditions of T1, T2, FLAIR, and T1ce modalities. The digit “1” signifies the availability of a particular modality, while “0” indicates its absence. Specifically, for T1 sequences, the synthetic results derived from multiple sequences (i.e., 0111) yield the most faithful quality with minimal error compared to the ground truth. This contrasted starkly with single-sequence inputs (i.e., 0100, 0010 and 0001). These multi-sequence inputs yielded synthetic results with diminished noise levels and sharply defined boundaries between white and gray matter regions. It is noteworthy that a clear anatomical structure, with less information about the tumor enhancement area, characterizes the T1 image. Therefore, it is easy to lose tumor enhancement areas without the guidance of complementary information from other modalities, such as 0110. Because T1ce modality provides a clear boundary between the regions enhanced around the tumor. Visual results eloquently establish the indispensability of integrating complementary information from diverse modalities to achieve precision in synthesizing tumor regions with accurate shapes and realistic textures. The synthetic sequence quality exhibits significant improvement with the increased number of available input sequences. This improvement is attributed to the inherent diversity and complementary nature of information encapsulated within different modalities. Meanwhile, the differential outcomes obtained across these scenarios vividly illustrate different modalities’ varying contributions to the target sequence synthesis process. Furthermore, these results validate our model’s ability to generalize effectively when confronted with varying quantities of available modalities.

The quantitative results of FgC2F-UDiff are summarized in Table.I and Table.II. The values in the tables represent the average results of using the input modality to synthesize all other target modalities. Specifically, On the BraTS dataset (as shown in Table.I), the first column (number 1) indicates the input modality (e.g., T1). The PSNR of 26.13 dB and SSIM of 0.887 values represent the average results of using the input modality to synthesize all other target modalities (e.g., T2, FLAIR, and T1ce). Compared with it, the combination of T1 and T2 (number 6) further improves the synthesis quality by 1.46 dB and 0.032 in terms of PSNR and SSIM, respectively. In addition, the combination of T2 and FLAIR information

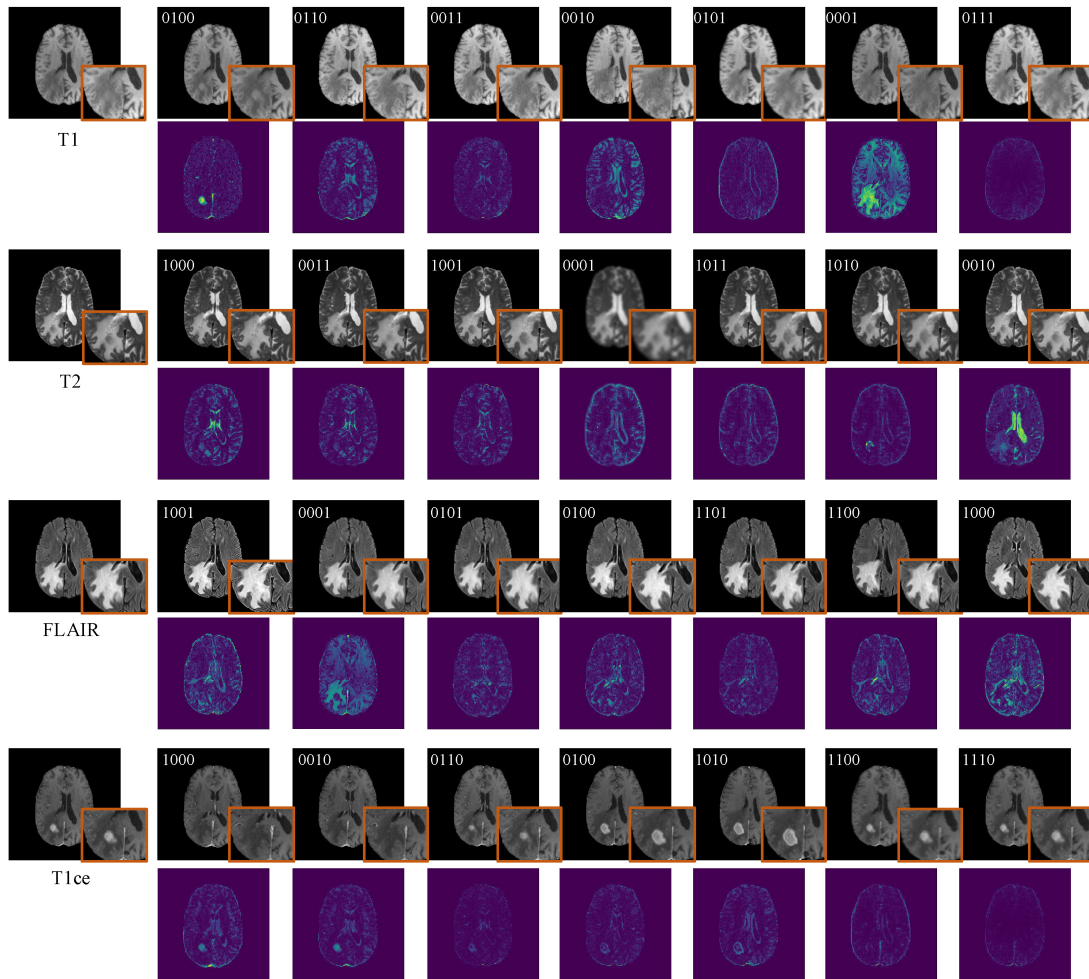


Fig. 5. Illustrative instances of synthetic images generated by our FgC2F-UDiff on the BraTS Dataset. Each row shows composite diagrams portraying distinct modes and error maps juxtaposed with the corresponding ground truth. The enlarged orange squares represent selected regions with notable disparities, providing enhanced insights into texture, edge enhancement, and shape characteristics.

(number 11) dramatically improves the synthesis quality by 2.19 *dB* and 0.055 for PSNR and SSIM, respectively. On IXI dataset (as shown in Table.II), FgC2F-UDiff achieves 34.00 *dB* in PSNR and 0.967 in SSIM using single T1 (number 1). Compared with it, the combination of T1 and T2 (number 4) improves the synthesis quality by 0.72 *dB* and 0.005 for PSNR and SSIM, respectively. The quantitative results indicate that the integration of the maximum number of accessible modalities attains optimal performance. This consistency between quantitative and qualitative findings underscores the significance of leveraging multiple modalities for enhanced synthesis outcomes. Besides, the synthesis results obtained in our experiments exhibit substantial variations across different datasets. These disparities are primarily attributed to two key factors. First, the inherent divergence in imaging principles gives rise to fundamental differences in the information encapsulated within the images themselves. Second, the complementary information embedded within different combinations of modalities leads to divergent guidance for the synthesis process. Consequently, the differential contributions of these modalities result in pronounced disparities in the quality and fidelity of the synthesized outputs.

TABLE I
QUANTITATIVE RESULTS OF OUR METHOD ON THE BRATS DATASET.
PSNR AND SSIM ARE REPORTED VALUES ARE MEAN \pm STD.

Number	Available modalities				Results	
	T1	FLAIR	T2	T1ce	PSNR(<i>dB</i>)	SSIM
1	✓				26.13 \pm 1.32	0.887 \pm 0.014
2		✓			25.96 \pm 1.51	0.881 \pm 0.015
3			✓		26.79 \pm 1.44	0.896 \pm 0.012
4				✓	27.12 \pm 1.09	0.903 \pm 0.017
5	✓	✓			27.38 \pm 1.32	0.917 \pm 0.015
6	✓		✓		27.59 \pm 1.24	0.919 \pm 0.009
7	✓			✓	27.93 \pm 1.08	0.923 \pm 0.010
8		✓	✓		27.45 \pm 1.26	0.915 \pm 0.021
9		✓		✓	28.28 \pm 1.17	0.933 \pm 0.012
10			✓	✓	28.66 \pm 1.05	0.943 \pm 0.007
11	✓	✓	✓		28.52 \pm 1.38	0.942 \pm 0.011
12	✓	✓		✓	28.75 \pm 1.29	0.946 \pm 0.008
13		✓	✓	✓	29.43 \pm 1.47	0.951 \pm 0.003

B. Synthesis results of task-specific comparison with state-of-the-art

To evaluate the synthesis performance in task-specific (one-to-one and many-to-one), task-specific FgC2F-UDiff was compared with the other seven SOTA methods (pix2pix, pGAN,

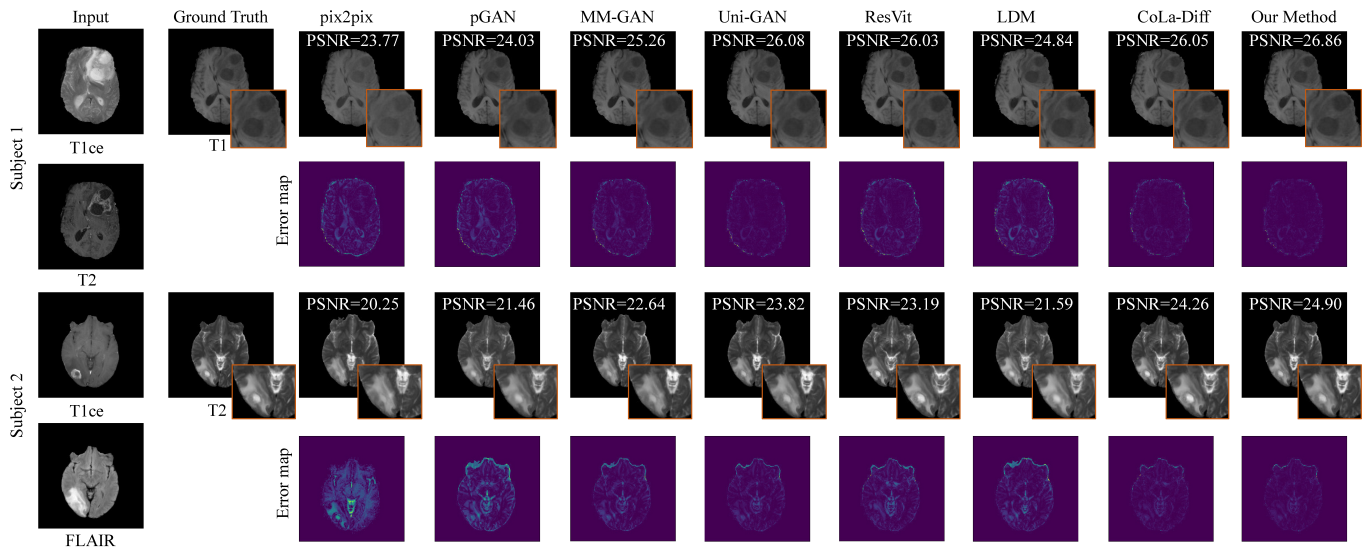


Fig. 6. Illustrative instances of synthetic images were demonstrated on the BraTS dataset for two representative tasks. Synthesized images from all competing methods are shown along with the source and reference target images. Partial enlargement and error plots can more intuitively observe the differences between the synthesized image and the ground truth, thereby reflecting the quality of the synthesis.

TABLE II
QUANTITATIVE RESULTS OF OUR METHOD ON THE IXI DATASET. PSNR AND SSIM ARE REPORTED VALUES ARE MEAN \pm STD.

Number	Modalities			Results	
	T1	T2	PD	PSNR(dB)	SSIM
1	✓			34.00 \pm 1.59	0.967 \pm 0.011
2		✓		34.57 \pm 1.72	0.971 \pm 0.012
3			✓	35.08 \pm 1.67	0.974 \pm 0.008
4	✓	✓		34.72 \pm 1.38	0.972 \pm 0.007
5	✓		✓	35.24 \pm 1.53	0.976 \pm 0.011
6		✓	✓	35.98 \pm 1.49	0.978 \pm 0.008

MM-GAN, Uni-GAN, ResVit, LDM, and CoLa-Diff) on the BraTS and IXI datasets. The quantitative results are shown in Table.III and Table.IV, which indicated that our proposed FgC2F-UDiff achieves the best performance in both task-specific (one-to-one and many-to-one) in terms of PSNR, SSIM, LPIPS, and FID metrics ($p < 0.05$). Specifically, on the BraTS dataset, one-to-one tasks of T1 \rightarrow T1ce; T1ce \rightarrow T1, many-to-one tasks of T1,T2 \rightarrow T1ce; and T1, FLAIR \rightarrow T1ce were considered, as shown in Table.III. In the one-to-one task of T1 \rightarrow T1ce, our method outperforms pix2pix, pGAN, MM-GAN, Uni-GAN, ResVit, LDM, and CoLa-Diff by a margin of 4.23 dB , 3.52 dB , 2.19 dB , 1.11 dB , 1.46 dB , 2.53 dB , and 0.85 dB , respectively. On the IXI dataset, one-to-one tasks of T1 \rightarrow PD; PD \rightarrow T1, many-to-one tasks of T1,T2 \rightarrow PD; and T1, PD \rightarrow T2 were considered, as shown in Table.IV. Our proposed FgC2F-UDiff achieves the best performance in both one-to-one tasks and many-to-one tasks in terms of PSNR, SSIM, LPIPS, and FID metrics ($p < 0.05$). Specifically, in the one-to-one task of T1 \rightarrow PD, our method outperforms pix2pix, pGAN, MM-GAN, Uni-GAN, ResVit, LDM and CoLa-Diff by a margin of 3.59 dB , 2.9 dB , 1.98 dB , 0.63 dB , 0.91 dB , 1.29 dB , and 0.64 dB , respectively. All these quantitative results prove that our FgC2F-UDiff is superior to other methods in the medical image synthesis task.

The visualized comparison results are shown in Fig.6 and Fig.7, which indicated that FgC2F-UDiff gains the best-synthesized performance compared with other SOTA methods on both datasets. Specifically, on the BraTS dataset, we selected two representative synthesis tasks: (1) T1ce, T2 \rightarrow T1 and (2) T1ce, FLAIR \rightarrow T2. The synthesized results for these tasks are displayed in Fig.6. To facilitate a more intuitive examination, we zoomed in on key regions within the synthesized results, which exhibit relatively conspicuous differences from the ground truth. Upon close observation of Subject 2, the disparities in synthesis primarily manifested in the tumor region. In comparison to the SOTA approaches, our proposed method closely approximated the ground truth. Other more advanced methods (i.e., LDM model and ResVit) have also achieved good synthesis results. Methods based on GANs (i.e., Uni-GAN, MM-GAN, and pGAN) exhibited some errors in fine details, while pix2pix nearly failed to capture the tumor-enhanced area. On the IXI dataset, we selected two representative synthesis tasks: (1) PD \rightarrow T2 and (2) T1 \rightarrow PD, and the synthesized results are shown in Fig.7. We magnified pivotal regions within the synthesized results to facilitate a more intuitive assessment that displayed discernible deviations from the corresponding ground truth. The results show that the synthesis results generated by our proposed FgC2F-UDiff model stand out notably. They exhibit superior quality are characterized by lower noise and retain better detail and structural information. The consistency observed between the quantitative and qualitative findings further demonstrates the superior synthesis performance of our proposed FgC2F-UDiff model. Compared to other state-of-the-art (SOTA) methods, FgC2F-UDiff consistently generates higher-quality results characterized by capturing fine details and preserving certain structural information with lower noise and clearer texture details, edges, and shapes.

TABLE III

QUANTITATIVE COMPARISON WITH SOTA METHODS ON BRATS DATASET. PSNR(dB), SSIM, LPIPS($\times 10^{-2}$), AND FID ARE LISTED AND REPORTED VALUES ARE MEAN \pm STD. THE **BOLDFACE** INDICATES THE TOP-PERFORMING MODEL FOR EACH TASK.

	T1 \rightarrow T1ce				T1ce \rightarrow T1				T1,T2 \rightarrow T1ce				T1,FLAIR \rightarrow T1ce			
	PSNR	SSIM	LPIPS	FID	PSNR	SSIM	LPIPS	FID	PSNR	SSIM	LPIPS	FID	PSNR	SSIM	LPIPS	FID
pix2pix	21.55	0.852	25.07	31.67	23.63	0.882	21.68	27.54	23.57	0.873	21.39	27.96	24.07	0.915	20.57	24.96
	± 1.28	± 0.017	± 1.53		± 1.46	± 0.013	± 1.35		± 1.54	± 0.013	± 1.27		± 1.52	± 0.013	± 1.31	
pGAN	22.26	0.857	23.76	29.56	24.19	0.887	20.17	25.37	24.23	0.875	20.37	25.72	24.45	0.919	19.16	24.03
	± 1.19	± 0.013	± 1.39		± 1.29	± 0.016	± 1.28		± 1.27	± 0.015	± 1.42		± 1.37	± 0.016	± 1.07	
MM-GAN	23.59	0.867	21.85	27.68	24.92	0.895	19.62	24.14	25.08	0.883	18.69	24.96	25.83	0.926	17.62	21.64
	± 1.54	± 0.014	± 1.94		± 1.32	± 0.014	± 0.99		± 1.62	± 0.012	± 1.17		± 1.63	± 0.013	± 1.15	
Uni-GAN	24.67	0.873	20.72	25.94	25.98	0.908	17.48	23.41	25.67	0.886	18.02	23.91	26.74	0.935	16.36	20.78
	± 1.37	± 0.013	± 0.78		± 1.78	± 0.016	± 1.06		± 1.28	± 0.014	± 1.20		± 1.32	± 0.016	± 1.72	
ResVit	24.32	0.871	21.38	26.51	25.51	0.903	18.08	23.97	25.77	0.889	17.73	23.57	26.58	0.933	16.94	21.96
	± 1.62	± 0.009	± 1.14		± 1.38	± 0.015	± 1.25		± 1.64	± 0.015	± 0.92		± 1.48	± 0.014	± 1.47	
LDM	23.25	0.866	22.19	28.53	25.08	0.896	18.96	24.03	24.74	0.880	19.11	24.76	25.19	0.922	18.54	22.59
	± 1.08	± 0.010	± 1.42		± 1.34	± 0.010	± 1.53		± 1.42	± 0.012	± 1.17		± 1.52	± 0.012	± 1.26	
CoLa-Diff	24.93	0.880	19.43	25.16	26.12	0.907	16.97	23.17	25.83	0.890	17.45	23.51	26.93	0.926	15.91	19.25
	± 1.52	± 0.010	± 0.97		± 1.62	± 0.011	± 1.37		± 1.52	± 0.015	± 1.35		± 1.09	± 0.011	± 1.42	
Our method	25.78	0.884	18.23	23.76	26.65	0.917	16.48	22.35	26.62	0.897	16.26	21.85	27.64	0.942	15.49	17.55
	± 1.63	± 0.015	± 1.06		± 1.30	± 0.012	± 1.42		± 1.37	± 0.011	± 1.28		± 1.33	± 0.013	± 1.08	

TABLE IV

QUANTITATIVE COMPARISON WITH SOTA METHODS ON IXI DATASET. PSNR(dB), SSIM, LPIPS($\times 10^{-2}$), AND FID ARE LISTED AND REPORTED VALUES ARE MEAN \pm STD. THE **BOLDFACE** INDICATES THE TOP-PERFORMING MODEL FOR EACH TASK.

	T1 \rightarrow PD				PD \rightarrow T1				T1,T2 \rightarrow PD				T1,PD \rightarrow T2			
	PSNR	SSIM	LPIPS	FID	PSNR	SSIM	LPIPS	FID	PSNR	SSIM	LPIPS	FID	PSNR	SSIM	LPIPS	FID
pix2pix	30.58	0.961	17.19	27.68	30.62	0.955	17.48	28.17	31.89	0.962	18.65	27.65	32.27	0.963	16.14	26.48
	± 1.75	± 0.016	± 2.08		± 1.47	± 0.015	± 1.69		± 1.63	± 0.017	± 1.87		± 1.39	± 0.016	± 1.74	
pGAN	31.27	0.964	16.25	25.26	32.02	0.965	16.35	25.74	32.00	0.966	17.96	24.76	32.48	0.968	15.47	23.59
	± 1.62	± 0.019	± 1.66		± 1.38	± 0.013	± 1.48		± 1.28	± 0.016	± 1.38		± 1.62	± 0.016	± 1.95	
MM-GAN	32.19	0.966	15.63	22.14	32.64	0.968	15.79	23.74	33.02	0.969	16.38	23.16	33.22	0.967	14.34	21.65
	± 1.53	± 0.013	± 1.48		± 1.53	± 0.016	± 1.57		± 1.63	± 0.014	± 1.56		± 1.38	± 0.014	± 1.68	
Uni-GAN	33.54	0.973	14.25	18.96	34.68	0.974	13.08	17.95	34.23	0.973	13.79	18.38	35.81	0.973	12.18	16.74
	± 1.96	± 0.015	± 1.62		± 1.29	± 0.012	± 1.43		± 1.57	± 0.013	± 1.47		± 1.48	± 0.015	± 1.37	
ResVit	33.26	0.972	15.37	20.79	34.51	0.972	13.57	19.42	33.95	0.975	14.53	19.52	35.63	0.977	12.67	17.32
	± 1.73	± 0.013	± 1.44		± 1.42	± 0.010	± 1.33		± 1.73	± 0.013	± 1.38		± 1.54	± 0.013	± 1.22	
LDM	32.88	0.971	13.98	18.34	33.53	0.970	14.92	21.95	33.43	0.973	15.72	22.51	33.99	0.972	13.79	20.78
	± 1.48	± 0.012	± 1.65		± 1.34	± 0.009	± 1.76		± 1.56	± 0.015	± 1.19		± 1.18	± 0.016	± 1.43	
CoLa-Diff	33.53	0.972	14.52	19.14	34.27	0.976	14.08	18.42	33.82	0.973	14.92	20.47	35.29	0.975	13.09	18.53
	± 1.42	± 0.015	± 1.59		± 1.08	± 0.012	± 1.37		± 1.27	± 0.009	± 1.32		± 1.52	± 0.013	± 1.49	
Our method	34.17	0.973	13.56	17.93	35.23	0.981	12.16	15.38	34.72	0.977	13.17	17.46	36.24	0.982	11.65	15.89
	± 1.58	± 0.014	± 1.37		± 1.27	± 0.014	± 1.58		± 1.46	± 0.010	± 1.49		± 1.43	± 0.004	± 1.42	

C. Synthesis results of unified model comparison with state-of-the-art

To evaluate the synthesis performance in unified synthesis models, unified FgC2F-UDiff was compared with the MM-GAN, ResVit, and Uni-GAN on many-to-one tasks of BraTS. Task-specific models are trained and tested to perform a single synthesis task to improve performance, but a separate model has to be built for each task. So, we demonstrate FgC2F-UDiff in learning unified synthesis models for multi-modality MRI. The quantitative results are shown in Table.V and Table.VI, which indicated that our proposed unified FgC2F-UDiff achieves the best performance on a many-to-one task in terms of PSNR, SSIM, LPIPS, and FID metrics ($p < 0.05$). Specifically, on the BraTS dataset, many-to-one tasks of T1,T2 \rightarrow T1ce and T1, FLAIR \rightarrow T1ce were considered. As shown in Table.V, our proposed FgC2F-UDiff achieves the best performance in many-to-one tasks in terms of PSNR, SSIM, LPIPS, and FID metrics ($p < 0.05$). On the many-to-one task of T1, T2 \rightarrow T1ce, our method outperforms MM-GAN, Uni-GAN, and ResVit by a margin of 1.55dB, 0.77dB, and

0.91dB, respectively. On the IXI dataset, many-to-one tasks of T1,T2 \rightarrow PD and T1, PD \rightarrow T2 were considered. As shown in Table.VI, our proposed FgC2F-UDiff achieves the best performance in many-to-one tasks in terms of PSNR, SSIM, LPIPS, and FID metrics ($p < 0.05$). Specifically, in the many-to-one task of T1, T2 \rightarrow PD, our method outperforms MM-GAN, Uni-GAN, and ResVit by a margin of 1.01dB, 0.60dB, and 0.93dB, respectively.

The visualized comparison results are shown in Fig.8, we report two representative synthesis tasks of T1 \rightarrow T1ce on BraTS and PD, T2 \rightarrow T1 on IXI, which indicated that FgC2F-UDiff gains the best-synthesized performance compared with other SOTA methods on both datasets. Our method demonstrates the synthesis of target images characterized by lower noise and clearer texture details, edges, and shapes when compared to baseline models. These results indicate the proficient consolidation of models for diverse source-target configurations by the unified FgC2F-UDiff model. Additionally, the consistency observed between quantitative and qualitative research findings further validates the superior synthesis performance of our proposed FgC2F-UDiff model in learning unified synthesis

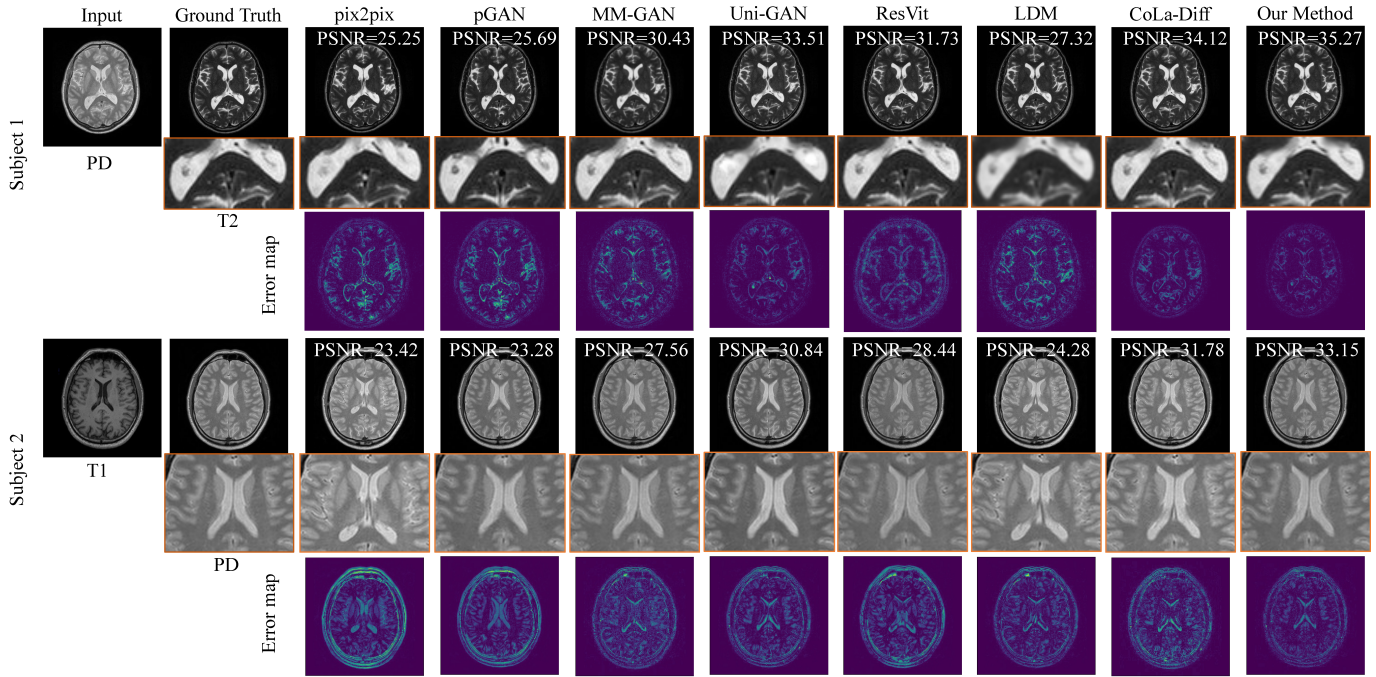


Fig. 7. Illustrative instances of synthetic images were demonstrated on the IXI dataset for two representative synthesis tasks. Synthesized images from all competing methods are shown along with the source and reference target images. Partial enlargement and error plots can more intuitively observe the differences between the synthesized image and the ground truth, thereby reflecting the quality of the synthesis.

TABLE V

QUANTITATIVE COMPARISON WITH SOTA METHODS. PSNR(dB), SSIM, LPIPS($\times 10^{-2}$), AND FID ARE LISTED AND REPORTED VALUES ARE MEAN \pm STD. THE **BOLDFACE** INDICATES THE TOP-PERFORMING MODEL FOR EACH TASK.

	T1,T2 \rightarrow T1ce				T1, FLAIR \rightarrow T1ce			
	PSNR	SSIM	LPIPS	FID	PSNR	SSIM	LPIPS	FID
MM-GAN	24.79	0.880	19.58	24.56	26.67	0.923	18.46	21.93
	± 1.43	± 0.011	± 1.38		± 1.38	± 0.024	± 1.42	
Uni-GAN	25.57	0.886	18.13	21.77	26.54	0.930	16.97	20.07
	± 1.52	± 0.014	± 1.28		± 1.23	± 0.028	± 1.08	
ResVit	25.43	0.885	18.49	22.34	26.34	0.929	17.28	20.76
	± 1.29	± 0.009	± 1.73		± 1.45	± 0.031	± 1.19	
Our method	26.34	0.894	17.45	20.85	27.19	0.937	16.14	18.74
	± 1.07	± 0.013	± 1.32		± 1.16	± 0.017	± 1.54	

TABLE VI

QUANTITATIVE COMPARISON WITH SOTA METHODS. PSNR(dB), SSIM, LPIPS($\times 10^{-2}$), AND FID ARE LISTED AND REPORTED VALUES ARE MEAN \pm STD. THE **BOLDFACE** INDICATES THE TOP-PERFORMING MODEL FOR EACH TASK.

	T1,T2 \rightarrow PD				T1, PD \rightarrow T2			
	PSNR	SSIM	LPIPS	FID	PSNR	SSIM	LPIPS	FID
MM-GAN	31.48	0.958	17.25	21.72	31.98	0.963	15.93	20.63
	± 1.37	± 0.007	± 1.55		± 1.53	± 0.011	± 1.67	
Uni-GAN	31.89	0.968	16.28	19.72	32.90	0.971	14.97	17.44
	± 1.77	± 0.006	± 1.36		± 1.27	± 0.008	± 1.26	
ResVit	31.56	0.963	16.93	20.35	32.17	0.971	15.42	18.35
	± 1.53	± 0.007	1.28		± 1.48	± 0.005	± 1.22	
Our method	32.49	0.971	15.73	18.13	33.77	0.973	14.52	16.38
	± 1.69	± 0.005	± 1.43		± 1.61	± 0.007	± 1.21	

models.

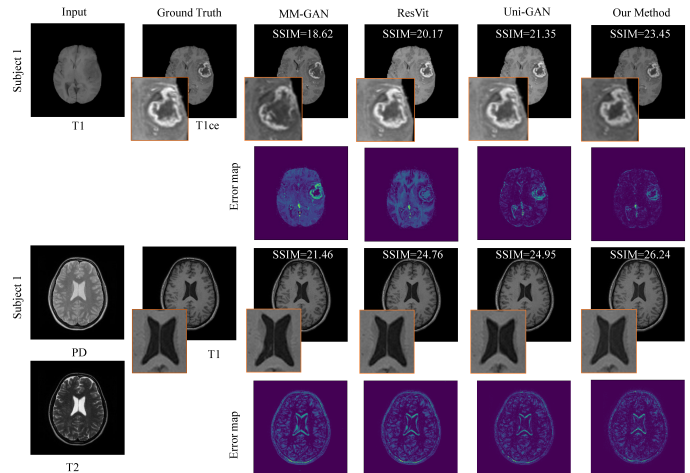


Fig. 8. Illustrative instances of synthetic images were demonstrated on the two datasets in learning unified synthesis models. Synthesized images from all competing methods are shown along with the source and reference target images.

D. Ablation studies show improvements of innovation

To validate the contributions of CUN, FCS, and SHM for cross-modality synthesis, we performed the comparison among our FgC2F-UDiff, the FgC2F-UDiff without CUN, the FgC2F-UDiff without FCS, and the FgC2F-UDiff without SHM. The structure of these ablation studies. Specifically, 1) To verify the contribution of our proposed CUN, we remove the module of frequency-guided conditions, leaving only the original diffusion model for synthesis tasks (Ours w/o CUN), as shown in Fig.V-D(b). All variants have the same network

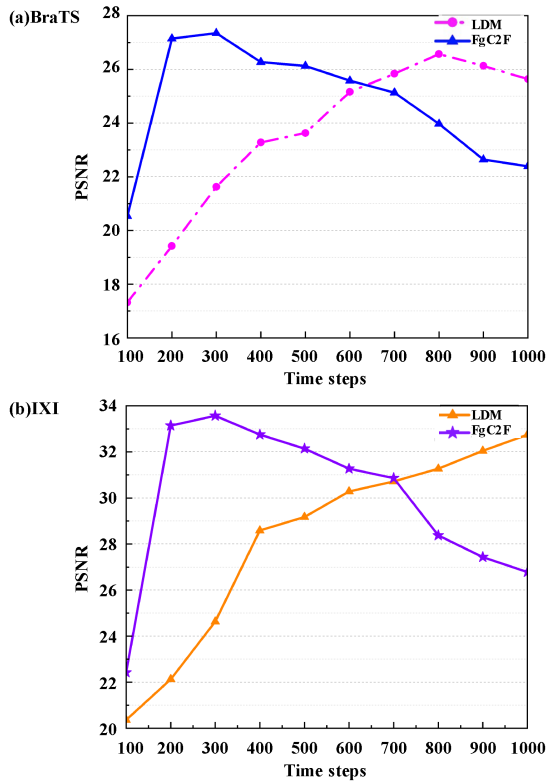


Fig. 9. The experimental results of our method compared with LDM at different time steps. The upper plot shows PSNR values of LDM and FgC2F-UDiff under different time steps on BraTS. The lower plot shows PSNR values of LDM and FgC2F-UDiff at different time steps on IXI.

structure except for the CUN. 2) To further examine the contribution of learning non-linear mapping guided by the proposed FCS, we designed two ablation (as shown in Fig.V-D(c)): one only utilizes all available modalities to obtain modality-class features (Ours w/o *LF* or *HF*), while the other solely employs frequency information for generating frequency-class guided features (Ours w/o Z_T^m). 3) To validate the contribution of SHM, we conducted the following experiments: (a) Our approach without CL (Ours w/o CL). We employ CL during training to facilitate network progression from easy to challenging tasks, thereby expediting both learning and the quality of synthesized samples. (b) Our approach without frequency-guided conditions (Ours w/o CUN). We leverage the frequency to guide the diffusion models to fit temporal regularities inherent in diffusion, thereby accelerating the coarse-to-fine synthesis. (c) Our approach without dynamically selecting frequency-information (Ours w/o dsf). To maintain fairness in the experiment, we use high-frequency and low-frequency information from the same available image as conditions. We utilize a dynamic strategy to retain specific features, elevating the synthesis quality of complex regions within images.

The quantitative analysis results of the ablation study are presented in Table.VII, which demonstrated that every part of FgC2F-UDiff contributes to the cross-modality synthesis. We report the average numerical results of each method across 14 input scenarios on the BraTS dataset. Especially for the CUN and Z_T^m all bring the cross-modality synthesis’s clear

TABLE VII
QUANTITATIVE COMPARISON OF ABLATION STUDY. PSNR(*dB*), SSIM, LPIPS($\times 10^{-2}$), AND FID ARE LISTED AND REPORTED VALUES ARE MEAN \pm STD. THE **BOLDFACE** INDICATES THE TOP-PERFORMING MODEL FOR EACH TASK.

	T1,T2 \rightarrow T1ce				T1, FLAIR \rightarrow T1ce			
	PSNR	SSIM	LPIPS	FID	PSNR	SSIM	LPIPS	FID
W/O CUN	20.58 ± 1.53	0.873 ± 0.009	22.42 ± 1.85	31.17	21.21 ± 1.34	0.907 ± 0.011	21.04 ± 1.46	28.86
W/O Z_T^m	22.64 ± 0.94	0.880 ± 0.013	18.28 ± 1.37	22.18	23.57 ± 1.27	0.919 ± 0.009	19.43 ± 1.53	25.14
W/O <i>LF</i> or <i>HF</i>	25.35 ± 1.27	0.889 ± 0.011	18.28 ± 1.36	22.14	26.11 ± 1.19	0.934 ± 0.014	17.25 ± 1.18	20.97
W/O CL	24.52 ± 1.17	0.883 ± 0.008	18.02 ± 1.52	23.86	25.15 ± 1.45	0.929 ± 0.012	18.63 ± 1.08	22.42
W/O dsf	25.87 ± 1.32	0.890 ± 0.013	17.83 ± 1.22	21.23	26.63 ± 1.16	0.932 ± 0.017	16.87 ± 1.06	20.35
Our method	26.34 ± 1.07	0.894 ± 0.013	17.45 ± 1.32	20.85	27.19 ± 1.16	0.937 ± 0.017	16.14 ± 1.54	18.74

performance gain. Specifically, for T1,T2 \rightarrow T1ce, the PSNR decreased from 26.34 ± 1.07 to 20.58 ± 1.53 when moving the CUN and decreased from 26.34 ± 1.07 to 22.64 ± 0.94 when moving the Z_T^m .

E. Analysis of SHM

To verify that our proposed acceleration mechanism can reduce the denoising time steps and inference time of our method without reducing the synthesis quality, we analyzed the impact of the SHM on our FgC2F-UDiff. In Fig.9, a comparative analysis of our approach, which integrates the LDM with SHM, reveals notable advancements in network training efficiency and a reduction in denoising steps. Specifically, when examining the BraTS dataset, it becomes evident that utilizing our method with a T of 200 or 300 surpasses the performance of the standalone LDM with a T set of 800. Likewise, when evaluating the IXI dataset, our model exhibits superior results with T values of 200 or 300 compared to the LDM with $T=1000$. Moreover, the decline in performance metrics beyond 300-time steps appears to be due to ‘over-diffusion’ [71], where the model continues to apply noise beyond the optimal range for image synthesis, resulting in unnecessary alterations that degrade image quality. Our experimental results clearly demonstrate that our method, combined with the proposed SHM, significantly accelerates the sampling speed, effectively addressing the sluggishness typically associated with conventional diffusion models. In addition, this discrepancy in performance can be attributed to several factors between the different datasets. The BraTS dataset boasts a greater number of modalities, thereby offering a richer array of diverse features for the model to leverage. Moreover, the increased sample size in BraTS contributes to heightened sample diversity, further enhancing the performance of FgC2F-UDiff.

Furthermore, our empirical analysis conclusively demonstrates that our proposed method significantly reduces inference times compared to the LDM. Utilizing 200 reverse diffusion steps, our method achieves an inference time of approximately 1.5 to 2 seconds per image in BraTS. In stark contrast, LDM, requiring 800 timesteps, has an inference time

ranging from 3.0 to 4.5 seconds per image under identical hardware conditions. This substantial reduction in inference time highlights the efficiency and practicality of our approach in real-world applications where rapid image synthesis is critical.

VI. CONCLUSION

This paper has presented a unified network for multi-modality missing MRI synthesis using a Frequency-guided and Coarse-to-fine Unified Diffusion Model (FgC2F-UDiff) from multiple inputs and outputs. CUN network has been introduced to leverage iterative denoising properties of the diffusion model to improve the fidelity of synthesizing images. In addition, an FCS strategy has been designed to utilize the frequency information to guide coarse-to-fine synthesis. The SHM further accelerates the diffusion process by intelligently integrating specific mechanisms, enhancing the efficiency and practicality of FgC2F-UDiff. Extensive experimental evaluations on two medical image synthesis datasets validate the effectiveness of our approach. This study provides a new perspective for addressing the missing modality issue in current technologies.

REFERENCES

- [1] B. H. Menze, A. Jakab, S. Bauer, J. Kalpathy-Cramer, K. Farahani, J. Kirby, Y. Burren, N. Porz, J. Slotboom, R. Wiest *et al.*, “The multimodal brain tumor image segmentation benchmark (brats),” *IEEE transactions on medical imaging*, vol. 34, no. 10, pp. 1993–2024, 2014.
- [2] U. Bagci, J. K. Udupa, N. Mendhiratta, B. Foster, Z. Xu, J. Yao, X. Chen, and D. J. Mollura, “Joint segmentation of anatomical and functional images: Applications in quantification of lesions from pet, pet-ct, mri-pet, and mri-pet-ct images,” *Medical image analysis*, vol. 17, no. 8, pp. 929–945, 2013.
- [3] X. Xiao, Q. V. Hu, and G. Wang, “Edge-aware multi-task network for integrating quantification segmentation and uncertainty prediction of liver tumor on multi-modality non-contrast mri,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2023, pp. 652–661.
- [4] R. M. Kronberg, D. Meskelevicius, M. Sabel, M. Kollmann, C. Rubbert, and I. Fischer, “Optimal acquisition sequence for ai-assisted brain tumor segmentation under the constraint of largest information gain per additional mri sequence,” *Neuroscience Informatics*, p. 100053, 2022.
- [5] A. Chartsias, T. Joyce, M. V. Giuffrida, and S. A. Tsaftaris, “Multi-modal mr synthesis via modality-invariant latent representation,” *IEEE transactions on medical imaging*, vol. 37, no. 3, pp. 803–814, 2017.
- [6] T. Varsavsky, Z. Eaton-Rosen, C. H. Sudre, P. Nachev, and M. J. Cardoso, “Pimms: permutation invariant multi-modal segmentation,” in *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: 4th International Workshop, DLMIA 2018, and 8th International Workshop, ML-CDS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, Proceedings 4*. Springer, 2018, pp. 201–209.
- [7] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, “Image-to-image translation with conditional adversarial networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.
- [8] J. Zhao *et al.*, “Tripartite-gan: Synthesizing liver contrast-enhanced mri to improve tumor detection,” *Medical Image Analysis*, vol. 63, p. 101667, 2020.
- [9] P. Huang, D. Li, Z. Jiao, D. Wei, G. Li, Q. Wang, H. Zhang, and D. Shen, “Coca-gan: common-feature-learning-based context-aware generative adversarial network for glioma grading,” in *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part III 22*. Springer, 2019, pp. 155–163.
- [10] B. Yu, L. Zhou, L. Wang, Y. Shi, J. Fripp, and P. Bourgeat, “Ea-gans: edge-aware generative adversarial networks for cross-modality mr image synthesis,” *IEEE transactions on medical imaging*, vol. 38, no. 7, pp. 1750–1762, 2019.
- [11] W. Yuan, J. Wei, J. Wang, Q. Ma, and T. Tasdizen, “Unified generative adversarial networks for multimodal segmentation from unpaired 3d medical images,” *Medical image analysis*, vol. 64, p. 101731, 2020.
- [12] B. Cao, H. Zhang, N. Wang, X. Gao, and D. Shen, “Auto-gan: self-supervised collaborative learning for medical image synthesis,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 07, 2020, pp. 10486–10493.
- [13] S. Wang, Z. Zhang, H. Yan, M. Xu, and G. Wang, “Mix-domain contrastive learning for unpaired h&e-to-ihc stain translation,” *arXiv preprint arXiv:2406.11799*, 2024.
- [14] M. Yurt, S. U. Dar, A. Erdem, E. Erdem, K. K. Oguz, and T. Çukur, “mustgan: multi-stream generative adversarial networks for mr image synthesis,” *Medical image analysis*, vol. 70, p. 101944, 2021.
- [15] W. Xu, C. Long, Y. Nie, and G. Wang, “Disentangled representation learning for controllable person image generation,” *IEEE Transactions on Multimedia*, 2024.
- [16] A. Sharma and G. Hamarneh, “Missing mri pulse sequence synthesis using multi-modal generative adversarial network,” *IEEE transactions on medical imaging*, vol. 39, no. 4, pp. 1170–1183, 2019.
- [17] O. Dalmaz, M. Yurt, and T. Çukur, “Resvit: residual vision transformers for multimodal medical image synthesis,” *IEEE Transactions on Medical Imaging*, vol. 41, no. 10, pp. 2598–2614, 2022.
- [18] Y. Zhang, C. Peng, Q. Wang, D. Song, K. Li, and S. K. Zhou, “Unified multi-modal image synthesis for missing modality imputation,” *arXiv preprint arXiv:2304.05340*, 2023.
- [19] D. Bau, J.-Y. Zhu, J. Wulff, W. Peebles, H. Strobelt, B. Zhou, and A. Torralba, “Seeing what a gan cannot generate,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4502–4511.
- [20] P. Dhariwal and A. Nichol, “Diffusion models beat gans on image synthesis,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 8780–8794, 2021.
- [21] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.
- [22] J. Zhao and S. Li, “Center-to-edge denoising diffusion probabilistic models with cross-domain attention for undersampled mri reconstruction,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer Nature Switzerland, 2024.
- [23] H. Li, Y. Yang, M. Chang, S. Chen, H. Feng, Z. Xu, Q. Li, and Y. Chen, “Srdiff: Single image super-resolution with diffusion probabilistic models,” *Neurocomputing*, vol. 479, pp. 47–59, 2022.
- [24] S. Gao, X. Liu, B. Zeng, S. Xu, Y. Li, X. Luo, J. Liu, X. Zhen, and B. Zhang, “Implicit diffusion models for continuous super-resolution,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 10021–10030.
- [25] W. Chao, F. Duan, X. Wang, Y. Wang, and G. Wang, “Lfsrdiff: Light field image super-resolution via diffusion models,” *arXiv preprint arXiv:2311.16517*, 2023.
- [26] A. Rahman, J. M. J. Valanarasu, I. Hacihaliloglu, and V. M. Patel, “Ambiguous medical image segmentation using diffusion models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 11536–11546.
- [27] K. Patel, F. Li, and G. Wang, “Multi-layer dense attention decoder for polyp segmentation,” *arXiv preprint arXiv:2403.18180*, 2024.
- [28] J. Zhao *et al.*, “United adversarial learning for liver tumor segmentation and detection of multi-modality non-contrast mri,” *Medical Image Analysis*, vol. 73, p. 102154, 2021.
- [29] T. Vyas, M. Chowdhury, X. Xiao, M. Claeys, G. Ong, and G. Wang, “Predicting mitral valve mteer surgery outcomes using machine learning and deep learning techniques,” in *Proceedings of the 2024 9th International Conference on Mathematics and Artificial Intelligence*, 2024, pp. 24–28.
- [30] J. Wolleb, R. Sandkühler, F. Bieder, P. Valmaggia, and P. C. Cattin, “Diffusion models for implicit image segmentation ensembles,” in *International Conference on Medical Imaging with Deep Learning*. PMLR, 2022, pp. 1336–1348.
- [31] X. Xiao, J. Zhao, and S. Li, “Task relevance driven adversarial learning for simultaneous detection, size grading, and quantification of hepatocellular carcinoma via integrating multi-modality mri,” *Medical Image Analysis*, vol. 81, p. 102554, 2022.
- [32] L. Jiang, Y. Mao, X. Wang, X. Chen, and C. Li, “Cola-diff: Conditional latent diffusion model for multi-modal mri synthesis,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2023, pp. 398–408.

- [33] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, "Diffusion models in vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [34] Q. Lyu and G. Wang, "Conversion between ct and mri images using diffusion and score-matching models," *arXiv preprint arXiv:2209.12104*, 2022.
- [35] A. Torrado-Carvajal, J. L. Herraiz, E. Alcain, A. S. Montemayor, L. Garcia-Canamaque, J. A. Hernandez-Tamames, Y. Rozenholc, and N. Malpica, "Fast patch-based pseudo-ct synthesis from t1-weighted mr images for pet/mr attenuation correction in brain studies," *Journal of Nuclear Medicine*, vol. 57, no. 1, pp. 136–143, 2016.
- [36] S. Roy, Y.-Y. Chou, A. Jog, J. A. Butman, and D. L. Pham, "Patch based synthesis of whole head mr images: Application to epi distortion correction," in *Simulation and Synthesis in Medical Imaging: First International Workshop, SASHIMI 2016, Held in Conjunction with MICCAI 2016, Athens, Greece, October 21, 2016, Proceedings 1*. Springer, 2016, pp. 146–156.
- [37] S. Roy, A. Carass, and J. L. Prince, "Magnetic resonance image example-based contrast synthesis," *IEEE transactions on medical imaging*, vol. 32, no. 12, pp. 2348–2363, 2013.
- [38] C. Bowles, C. Qin, C. Ledig, R. Guerrero, R. Gunn, A. Hammers, E. Sakka, D. A. Dickie, M. V. Hernández, N. Royle *et al.*, "Pseudo-healthy image synthesis for white matter lesion segmentation," in *Simulation and Synthesis in Medical Imaging: First International Workshop, SASHIMI 2016, Held in Conjunction with MICCAI 2016, Athens, Greece, October 21, 2016, Proceedings 1*. Springer, 2016, pp. 87–96.
- [39] A. Jog, A. Carass, S. Roy, D. L. Pham, and J. L. Prince, "Random forest regression for magnetic resonance image synthesis," *Medical image analysis*, vol. 35, pp. 475–488, 2017.
- [40] —, "Mr image synthesis by contrast learning on neighborhood ensembles," *Medical image analysis*, vol. 24, no. 1, pp. 63–76, 2015.
- [41] D. H. Ye, D. Zikic, B. Glocker, A. Criminisi, and E. Konukoglu, "Modality propagation: coherent synthesis of subject-specific scans with data-driven regularization," in *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2013: 16th International Conference, Nagoya, Japan, September 22–26, 2013, Proceedings, Part I 16*. Springer, 2013, pp. 606–613.
- [42] H. Van Nguyen, K. Zhou, and R. Vemulapalli, "Cross-domain synthesis of medical images using efficient location-sensitive deep network," in *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part I 18*. Springer, 2015, pp. 677–684.
- [43] V. Sevetlidis, M. V. Giuffrida, and S. A. Tsiftaris, "Whole image synthesis using a deep encoder-decoder network," in *Simulation and Synthesis in Medical Imaging: First International Workshop, SASHIMI 2016, Held in Conjunction with MICCAI 2016, Athens, Greece, October 21, 2016, Proceedings 1*. Springer, 2016, pp. 127–137.
- [44] R. Li, W. Zhang, H.-I. Suk, L. Wang, J. Li, D. Shen, and S. Ji, "Deep learning based imaging data completion for improved brain disease diagnosis," in *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2014: 17th International Conference, Boston, MA, USA, September 14–18, 2014, Proceedings, Part III 17*. Springer, 2014, pp. 305–312.
- [45] B. Peng, B. Liu, Y. Bin, L. Shen, and J. Lei, "Multi-modality mr image synthesis via confidence-guided aggregation and cross-modality refinement," *IEEE Journal of Biomedical and Health Informatics*, vol. 26, no. 1, pp. 27–35, 2021.
- [46] M. Yurt, O. Dalmaz, S. Dar, M. Ozbey, B. Tinaz, K. Oguz, and T. Çukur, "Semi-supervised learning of mri synthesis without fully-sampled ground truths," *IEEE Transactions on Medical Imaging*, vol. 41, no. 12, pp. 3895–3906, 2022.
- [47] A. Jog, A. Carass, D. L. Pham, and J. L. Prince, "Random forest flair reconstruction from t1, t2, and p d-weighted mri," in *2014 IEEE 11th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2014, pp. 1079–1082.
- [48] R. Mehta and T. Arbel, "Rs-net: Regression-segmentation 3d cnn for synthesis of full resolution missing brain mri in the presence of tumours," in *Simulation and Synthesis in Medical Imaging: Third International Workshop, SASHIMI 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Proceedings 3*. Springer, 2018, pp. 119–129.
- [49] C. Chen, Q. Dou, Y. Jin, H. Chen, J. Qin, and P.-A. Heng, "Robust multimodal brain tumor segmentation via feature disentanglement and gated fusion," in *Medical Image Computing and Computer Assisted Intervention—MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part III 22*. Springer, 2019, pp. 447–456.
- [50] T. Zhou, H. Fu, G. Chen, J. Shen, and L. Shao, "Hi-net: hybrid-fusion network for multi-modal mr image synthesis," *IEEE transactions on medical imaging*, vol. 39, no. 9, pp. 2772–2781, 2020.
- [51] D. Lee, J. Kim, W.-J. Moon, and J. C. Ye, "Collagan: Collaborative gan for missing image data imputation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 2487–2496.
- [52] D. Lee, W.-J. Moon, and J. C. Ye, "Assessing the importance of magnetic resonance contrasts using collaborative generative adversarial networks," *Nature Machine Intelligence*, vol. 2, no. 1, pp. 34–42, 2020.
- [53] W. Xu and G. Wang, "A domain gap aware generative adversarial network for multi-domain image translation," *IEEE Transactions on Image Processing*, vol. 31, pp. 72–84, 2021.
- [54] B. Zhan, D. Li, X. Wu, J. Zhou, and Y. Wang, "Multi-modal mri image synthesis via gan with multi-scale gate merge," *IEEE Journal of Biomedical and Health Informatics*, vol. 26, no. 1, pp. 17–26, 2021.
- [55] M. Yurt, M. Özbey, S. U. Dar, B. Tinaz, K. K. Oguz, and T. Çukur, "Progressively volumetrized deep generative models for data-efficient contextual learning of mr image recovery," *Medical Image Analysis*, vol. 78, p. 102429, 2022.
- [56] Y. Ding, X. Yu, and Y. Yang, "Rfnet: Region-aware fusion network for incomplete multi-modal brain tumor segmentation," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 3975–3984.
- [57] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.
- [58] X. Meng, K. Sun, J. Xu, X. He, and D. Shen, "Multi-modal modality-masked diffusion network for brain mri synthesis with random modality missing," *IEEE Transactions on Medical Imaging*, 2024.
- [59] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18*. Springer, 2015, pp. 234–241.
- [60] M. Heideman, D. Johnson, and C. Burrus, "Gauss and the history of the fast fourier transform," *IEEE Assp Magazine*, vol. 1, no. 4, pp. 14–21, 1984.
- [61] R. A. Schowengerdt, "Reconstruction of multispatial, multispectral image data using spatial frequency content," *Photogrammetric Engineering and Remote Sensing*, vol. 46, no. 10, pp. 1325–1334, 1980.
- [62] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proceedings of the 26th annual international conference on machine learning*, 2009, pp. 41–48.
- [63] U. Baid, S. Ghodasara, S. Mohan, M. Bilello, E. Calabrese, E. Colak, K. Farahani, J. Kalpathy-Cramer, F. C. Kitamura, S. Pati *et al.*, "The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification," *arXiv preprint arXiv:2107.02314*, 2021.
- [64] S. Bakas, H. Akbari, A. Sotiras, M. Bilello, M. Rozycki, J. S. Kirby, J. B. Freymann, K. Farahani, and C. Davatzikos, "Advancing the cancer genome atlas glioma mri collections with expert segmentation labels and radiomic features," *Scientific data*, vol. 4, no. 1, pp. 1–13, 2017.
- [65] IXI. Information extraction from images. [Online]. Available: www.brain-development.org
- [66] M. Jenkinson and S. Smith, "A global optimisation method for robust affine registration of brain images," *Medical image analysis*, vol. 5, no. 2, pp. 143–156, 2001.
- [67] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [68] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [69] S. U. Dar, M. Yurt, L. Karacan, A. Erdem, E. Erdem, and T. Cukur, "Image synthesis in multi-contrast mri with conditional generative adversarial networks," *IEEE transactions on medical imaging*, vol. 38, no. 10, pp. 2375–2388, 2019.
- [70] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10684–10695.
- [71] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," *arXiv preprint arXiv:2010.02502*, 2020.