# DIFFUSION MODELS IN MEDICAL IMAGING: A COMPREHENSIVE SURVEY

**Amirhossein Kazerouni**
School of Electrical Engineering
Iran University of Science and Technology
Tehran, Iran
amirhossein477@gmail.com

**Ehsan Khodapanah Aghdam**
Department of Electrical Engineering
Shahid Beheshti University
Tehran, Iran
ehsan.khpaghdam@gmail.com

**Moein Heidari**
School of Electrical Engineering
Iran University of Science and Technology
Tehran, Iran
moein_heidari@elec.iust.ac.ir

**Reza Azad**
Institute of Imaging and Computer Vision
RWTH Aachen University
Aachen, Germany
azad@lfb.rwth-aachen.de

**Mohsen Fayyaz**
Microsoft
Berlin, Germany
mohsenfayyaz@microsoft.com

**Ilker Hacihaliloglu**
Department of Radiology
Department of Medicine
University of British Columbia
British Columbia, Canada
ilker.hacihaliloglu@ubc.ca

**Dorit Merhof**
Faculty of Informatics and Data Science
University of Regensburg
Regensburg, Germany
dorit.merhof@ur.de

## ABSTRACT

Denoising diffusion models, a class of generative models, have garnered immense interest lately in various deep-learning problems. A diffusion probabilistic model defines a forward diffusion stage where the input data is gradually perturbed over several steps by adding Gaussian noise and then learns to reverse the diffusion process to retrieve the desired noise-free data from noisy data samples. Diffusion models are widely appreciated for their strong mode coverage and quality of the generated samples in spite of their known computational burdens. Capitalizing on the advances in computer vision, the field of medical imaging has also observed a growing interest in diffusion models. With the aim of helping the researcher navigate this profusion, this survey intends to provide a comprehensive overview of diffusion models in the discipline of medical imaging. Specifically, we start with an introduction to the solid theoretical foundation and fundamental concepts behind diffusion models and the three generic diffusion modeling frameworks, namely, diffusion probabilistic models, noise-conditioned score networks, and stochastic differential equations. Then, we provide a systematic taxonomy of diffusion models in the medical domain and propose a multi-perspective categorization based on their application, imaging modality, organ of interest, and algorithms. To this end, we cover extensive applications of diffusion models in the medical domain, including image-to-image translation, reconstruction, registration, classification, segmentation, denoising, 2/3D generation, anomaly detection, and other medically-related challenges. Furthermore, we emphasize the practical use case of some selected approaches, and then we discuss the limitations of the diffusion models in the medical domain and propose several directions to fulfill the demands of this field. Finally, we gather the overviewed studies with their available open-source implementations at our GitHub[1]. We aim to update the relevant latest papers within it regularly.

**Keywords** Generative models · Diffusion models · Denoising diffusion models · Noise conditioned score networks · Score-based models · Medical imaging · Medical applications · Survey
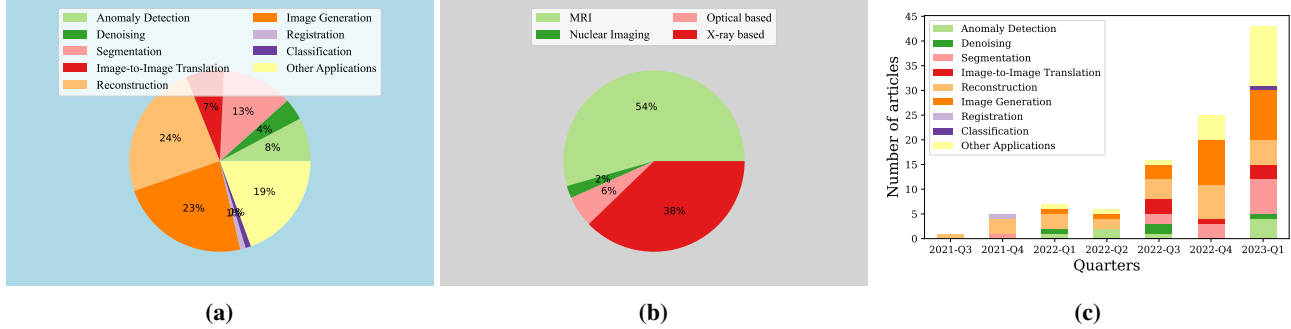
---

[1] https://github.com/amirhossein-kz/Awesome-Diffusion-Models-in-Medical-Imaging

Figure 1: The diagram **(a)** shows the relative proportion of published papers categorized according to their application and **(b)** according to their imaging modalities. **(c)** indicates the number of diffusion-based research papers published in the medical domain. The growth rate per year reveals the importance of diffusion models for future work. It is worth mentioning that the overall number of papers is 103.

# 1 Introduction

Generative modeling using neural networks has been a leading force in the past decade of deep learning. Since their emergence, generative models have made a tremendous impact in various domains ranging from images [1, 2], audio [3, 4], to text [5], and point clouds [6]. From a probabilistic modeling viewpoint, the key defining characteristic of a generative model is that it is trained in such a way so that its samples $\tilde{x} \sim p_\theta(\tilde{x})$ come from the same distribution as the training data distribution, $x \sim p_d(x)$. The energy-based models (EBMs) achieve this by defining an unnormalized probability density over a state space; however, these methods require Markov Chain Monte Carlo (MCMC) sampling during both training and inference, which is a slow iterative process [7]. In the past few years, due to the progressions in general deep learning architectures, there has been a resurgence of interest in generative models, unveiling improved visual fidelity and sampling speed. Specifically, generative adversarial networks (GANs) [8], variational autoencoders (VAEs) [9], and normalizing flows [10] have emerged. Apart from these, generative models based on diffusion processes offer an alternative to existing VAEs, EBMs, GANs, and normalizing flows, which do not require the alignment of posterior distributions, the estimation of intractable partition functions, the introduction of additional discriminator networks or the placement of network constraints respectively. To date, diffusion models have been found to be useful in a wide variety of areas, ranging from generative modeling tasks such as image generation [11], image super-resolution [12], image inpainting [13] to discriminative tasks such as image segmentation [14], classification [15], and anomaly detection [16]. Recently, the medical imaging community has witnessed exponential growth in the number of diffusion-based techniques (see Figure 1). As shown in Figure 1, a wealth of research is dedicated to the applications of diffusion models in diverse medical imaging scenarios. Since diffusion models have recently received significant attention from the research community, the literature is experiencing a large influx of contributions in this direction. Thus, a survey of the existing literature is beneficial for the community and timely. To this end, this survey sets out to provide a comprehensive overview of the recent advances made and provides a holistic overview of this class of models in medical imaging. A thorough search of the relevant literature revealed that we are the first to cover the diffusion-based models exploited in the medical domain. We hope this work will point out new paths, provide a road map for researchers, and inspire further interest in the vision community to leverage the potential of diffusion models in the medical domain. Our major contributions include:

• This is the first survey paper that comprehensively covers applications of diffusion models in the medical imaging domain. Specifically, we present a comprehensive overview of all available relevant papers (until October 2022) as well as showcase some of the latest techniques through April 2023.

• We devise a multi-perspective categorization of diffusion models in the medical community, providing a systematical taxonomy of research in diffusion models and their applications. We divide the existing diffusion models into two categories: variational-based models and score-based models. Moreover, we group the applications of diffusion models into nine categories: image-to-image translation, reconstruction, registration, classification, segmentation, denoising, image generation, anomaly detection, and other applications.

• We do not restrict our attention to application and provide a new taxonomy (see Figure 5) where each paper is broadly classified according to the proposed algorithm along with the organ concerned and imaging modality, respectively.

• Finally, we discuss the challenges and open issues and identify the new trends raising open questions about the future development of diffusion models in the medical domain in both algorithms and applications.

***Motivation and uniqueness of this survey***. Generative approaches have undergone significant advances in medical imaging over the past few decades. Therefore, there have been numerous survey papers published on deep generative models for medical imaging [17, 18, 19]. Some of these papers only focus on a specific application, while others concentrate on a specific image modality. There have also been review articles on diffusion models surfacing recently for computer vision tasks [20, 21, 22]. Although reviews had already been released before this area is fully developed, a lot of advances in the medical field have come out since then. On the other hand, none of these surveys focuses on the applications of diffusion models in medical imaging, which is the central aspect in pushing this research direction forward. Hence, these surveys leave a clear open gap. Besides, we believe that the medical community can leverage insights from successful products of diffusion models in vision by a retrospective on the past and future research directions of diffusion models provided in our survey. Moreover, diffusion models have demonstrated their potential in generating synthetic data and can serve as an effective supplement to existing real data, as well as a generative prior in biomedical inverse imaging problems (as discussed in Section 3). Lastly, we think our survey can help medical researchers (e.g., radiologists) and guide them toward exploiting up-to-date methodologies in their fields. To this end, in this survey, we devise a multi-perspective vision of diffusion models in which we discuss existing literature based on their applications in the medical domain. Nonetheless, we do not restrict our interest to the applications but describe the underlying working principles, the organ, and the imaging modality of the proposed method. We further discuss how this additional information can help researchers attempt to consolidate the literature across the spectrum. A brief outlook of our paper is depicted in Figure 5.

***Search Strategy***. We searched DBLP, Google Scholar, and Arxiv Sanity Preserver with customized search queries, as they allow for customized search queries and provide lists for all scholarly publications: peer-reviewed journal papers or papers published in the proceedings of conferences or workshops, non-peer-reviewed papers, and preprints. Our search query was (diffusion* deep | medical | imaging*) (denoising | medical*) (diffusion* | medical* | probabilistic* | model*) (score* | diffusion* | model* | medical*). We filtered our search results to remove false positives and included only papers related to diffusion probabilistic models (e.g., we had many false-positive search results about diffusion Magnetic Resonance Imaging (MRI) models). Notably, we selected the papers for detailed examination based on a careful evaluation of their novelty, contribution, significance, and if being the first introduced paper in medical imaging. After applying these criteria, we selected two or three of the highest-ranked papers to examine in more detail. We acknowledge that there may be other important papers in the field that we did not discuss in our review, but we aimed to provide a comprehensive overview of the most important and impactful papers.

***Paper Organization***. In Section 2, we present a detailed overview of the concepts and theoretical foundations behind diffusion models, covering two perspectives of diffusion models. In Section 3, we will delve into the significance of employing generative models, specifically diffusion models, in clinical settings and discuss the benefits they offer. Section 4 comprehensively covers the applications of diffusion models in several medical imaging tasks, as shown in Figure 5, and finally provides a comparative task-specific overview of different literature work. We conclude this survey by pinpointing future directions and open challenges facing diffusion models in the medical imaging domain in Section 5.

## 2 Theory

Diffusion models are a cutting-edge class of generative models that have been demonstrated to be highly effective in learning complex data distributions. They are a relatively new addition to the generative learning landscape but have shown to be useful in various applications. In this section, we take an in-depth look at the theory of diffusion models. We begin by discussing the position of diffusion models within the broader generative learning landscape and provide a new perspective on how they compare to other generative models. We further classify diffusion models into two main perspectives: the **Variational Perspective** and the **Score Perspective**. We delve into their details and highlight the specific models that fall under them, such as DDPMs in the Variational Perspective and NCSNs and SDEs in the Score Perspective. Ultimately, we provide a comprehensive understanding of the underlying theory behind these methods.

### 2.1 Where do diffusion models fit the generative learning landscape?

Following the remarkable surge of available datasets, as well as advances in general deep learning architectures, there has been a revolutionary paradigm shift in generative modeling. Specifically, the three mainstream generative frameworks include, namely, GANs [8], VAEs [9, 24], and normalizing flows [10] (see Figure 2). Generative models typically entail key requirements to be adopted in real-world problems. These requirements include (i) high-quality sampling, (ii) mode coverage and sample diversity, and (iii) fast execution time and computationally inexpensive sampling [26] (see Figure 3).

Generative models often make accommodations between these criteria. Specifically, GANs are capable of generating high-quality samples rapidly, but they have poor mode coverage and are prone to lack sampling diversity. Conversely, VAEs and normalizing flows suffer from the intrinsic property of low sample quality despite being witnessed in covering data modes. GANs consist of two models: a generator and a critic (discriminator), which compete with each other while

**(a)** GAN



**(b)** Energy-based Models



**(c)** VAE



**(d)** Flow-based models
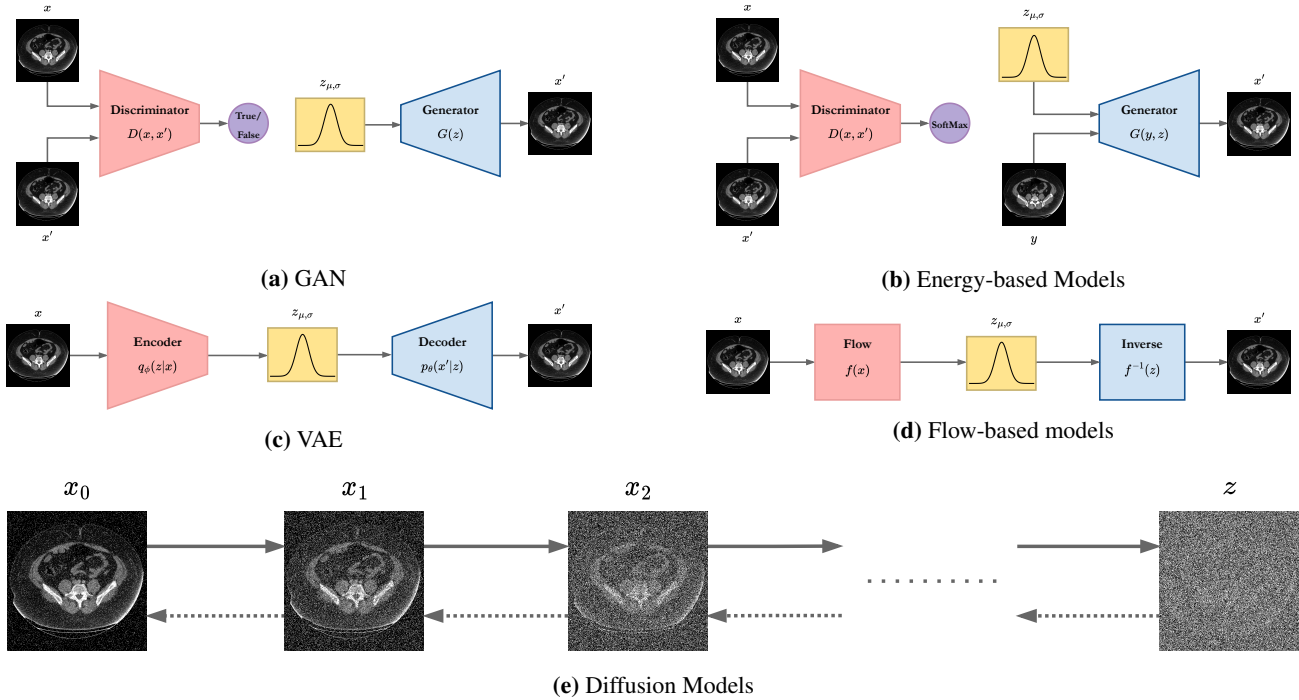


**(e)** Diffusion Models

Figure 2: This figure showcases different generative models and provides an overview of their underlying principles. (a) General Adversarial Network (GAN) [8] is an end-to-end pipeline that trains the generator in an adversarial manner to generate samples that the discriminator is capable of distinguishing from the real data sample. (b) Energy-based Model (EBM) [23], also known as non-normalized probabilistic models, trains in the same way as GANs with two major modifications. First, the discriminator learns a proper energy-based function that maps the data sample to a distribution space. Second, the generator utilizes a prior input to enhance the sample generation performance. (c) Variational AutoEncoder (VAE) [24] is a standalone network that follows a projection from a data sample to a low-dimensional latent space by the encoder and generates by sampling from it via a decoder path. (d) Normalizing flow (NF) [25] utilizes an invertible flow function to transform input to latent space and generate samples with the inverse flow function. (e) Diffusion Models intermingle the noise with the input in successive steps until it becomes a noise distribution before applying a reverse process to neutralize the noise addition in each step in the sampling procedure.

simultaneously making each other stronger. The generator tries to capture the distribution of true examples, while the discriminator, which is typically a binary classifier, estimates the probability of a given sample coming from the real dataset. It works as a critic and is optimized to recognize the synthetic samples from the real ones. A common concern with GANs is their training dynamics which have been recognized as being unstable, resulting in deficiencies such as mode collapse, vanishing gradients, and convergence [27]. Therefore, an immense interest has also influenced the research direction of GANs to propose more efficient variants [28, 29]. VAEs optimize the log-likelihood of the data by maximizing the evidence lower bound (ELBO). Despite the remarkable achievements, the behavior of VAEs is still far from satisfactory due to some theoretical and practical challenges such as balancing issue [30] and variable collapse phenomenon [31]. A flow-based generative model is constructed by a sequence of invertible transformations. Specifically, a normalizing flow transforms a simple distribution into a complex one by applying a sequence of invertible transformation functions where one can obtain the desired probability distribution for the final target variable using a change of variables theorem. Unlike GANs and VAEs, these models explicitly learn the data distribution; therefore, their loss function is simply the negative log-likelihood [32]. Despite being feasibly designed, these generative models have their specific drawbacks. Since the Likelihood-based method has to construct a normalized probability model, a specific type of architecture must be used (Autoregressive Model, Flow Model), or in the case of VAE, an alternative Loss such as ELBO is not calculated directly for the generated probability distribution. In contrast, the learning process of GANs is inherently unstable due to the nature of the adversarial loss of the GAN. Recently, diffusion models [33, 34] have emerged as powerful generative models, showcasing one of the leading topics in computer vision so that researchers and practitioners alike may find it challenging to keep pace with the rate of innovation.

Diffusion models are a powerful class of probabilistic generative models that are used to learn complex data distributions. These models accomplish this by utilizing two key stages: the forward diffusion process and the reverse diffusion process.

The forward diffusion process adds noise to the input data, gradually increasing the noise level until the data is transformed into pure Gaussian noise. This process systematically perturbs the structure of the data distribution. The reverse diffusion process, also known as denoising, is then applied to recover the original structure of the data from the perturbed data distribution. This process effectively undoes the degradation caused by the forward diffusion process. The result is a highly flexible and tractable generative model that can accurately model complex data distributions from random noise.

## 2.2 Variational Perspective

The Variational Perspective category includes models that use variational inference to approximate the target distribution, generally by minimizing the Kullback-Leibler divergence between the approximate and target distributions. Denoising Diffusion Probabilistic Models (DDPMs) [33, 34] are an example of this type of model, as they use a variational inference approach to estimate the parameters of a diffusion process.

### 2.2.1 Denoising Diffusion Probabilistic Models (DDPMs)

**Forward Process.** DDPM defines the forward diffusion process as a Markov Chain where Gaussian noise is added in successive steps to obtain a set of noisy samples. Consider $q(x_0)$ as the uncorrupted (original) data distribution. Given a data sample $x_0 \sim q(x_0)$, a forward noising process $p$ which produces latent $x_1$ through $x_T$ by adding Gaussian noise at time $t$ is defined as follows:

$$q(x_t \mid x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{1-\beta_t} \cdot x_{t-1}, \beta_t \cdot \mathbf{I}\right), \forall t \in \{1, \ldots, T\}, \tag{1}$$

where $T$ and $\beta_1, \ldots, \beta_T \in [0, 1)$ represent the number of diffusion steps and the variance schedule across diffusion steps, respectively. $\mathbf{I}$ is the identity matrix and $\mathcal{N}(x; \mu, \sigma)$ represents the normal distribution of mean $\mu$ and covariance $\sigma$. Considering $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=0}^{t} \alpha_s$, one can directly sample an arbitrary step of the noised latent conditioned on the input $x_0$ as follows:
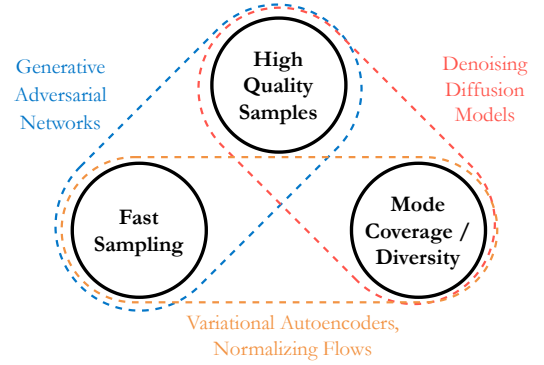


Figure 3: Generative learning trilemma [26]. Despite the ability of GANs to quickly generate high-fidelity samples, their mode coverage is limited. In addition, VAEs and normalizing flows have been revealed to have a great deal of diversity; however, they generally have poor sampling quality. Diffusion models have emerged to compensate for the deficiency of VAEs and GANs by showing adequate mode coverage and high-quality sampling. Nevertheless, due to their iterative nature, which causes a slow sampling process, they are practically expensive and require more improvement.

$$q(\mathbf{x}_t \mid \mathbf{x}_0) = N\left(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}\right) \tag{2}$$

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_l}\epsilon. \tag{3}$$

**Reverse Process.** Leveraging the above definitions, we can approximate a reverse process to get a sample from $q(x_0)$. To this end, we can parameterize this reverse process by starting at $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$ as follows:

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T)\prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) \tag{4}$$

$$p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t)). \tag{5}$$

To train this model such that $p(x_0)$ learns the true data distribution $q(x_0)$, we can optimize the following variational bound on negative log-likelihood:

$$\mathbb{E}\left[-\log p_\theta(\mathbf{x}_0)\right] \leq \mathbb{B}_q\left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)}\right]$$

$$= \mathbb{E}_q\left[-\log p(\mathbf{x}_T) - \sum_{t \geq 1} \log \frac{p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_t \mid \mathbf{x}_{t-1})}\right] \tag{6}$$

$$= -L_{\text{VL.B}}.$$

Ho et al. [34] found it better not to directly parameterize $\mu_\theta(x_t, t)$ as a neural network, but instead to train a model $\epsilon_\theta(x_t, t)$ to predict $\epsilon$. Hence, by reparameterizing Equation (6), they proposed a simplified objective as follows:

$$L_{\text{simple}} = E_{t,x_0,\epsilon} \left[ \| \epsilon - \epsilon_\theta(x_t, t) \|^2 \right], \tag{7}$$

where the authors draw a connection between the loss in Equation (6) to generative score networks in Song et al. [35].

## 2.3 Score Perspective

Score Perspective models rely on a maximum likelihood-based estimation approach, using the score function of the log-likelihood of the data to estimate the parameters of the diffusion process. Noise-conditioned Score Networks (NCSNs) [35] and Stochastic Differential Equations (SDEs) [36] are both subcategories that fall into this category. NCSNs focus on estimating the derivative of the log density function of the perturbed data distribution at different noise levels, while SDEs are a generalization of previous approaches and encompass both DDPMs and NCSNs characteristics. We hereinafter elaborate on the details of each subcategory.

### 2.3.1 Noise Conditioned Score Networks (NCSNs)

The score function of some data distribution $p(x)$ is defined as the gradient of the log density with respect to the input, $\nabla_x \log p(x)$. To estimate this score function, one can train a shared neural network with score matching. Specifically, the score network $s_\theta$ is a neural network parameterized by $\theta$, which is trained to approximate the score of $p(x)$ ($s_\theta(x) \approx \nabla_x \log p(x)$) by minimizing the following objective:

$$\mathbb{E}_{x \sim p(x)} \| s_\theta(x) - \nabla_x \log p(x) \|_2^2. \tag{8}$$

However, due to the computational burden of calculating $\nabla_x \log p(x)$, score matching is not scalable to deep networks and high dimensional data. To mitigate this problem, the authors of [35] propose to exploit denoising score matching [37] and sliced score matching [38]. Moreover, Song et al. [35] highlight major challenges that prevent a naive application of score-based generative modeling in real data. The key challenge is the fact that the estimated score functions are inaccurate in low-density regions since data in the real world tend to concentrate on low-dimensional manifolds embedded in a high-dimensional space (the manifold hypothesis). The authors demonstrate that these problems can be addressed by perturbing the data with Gaussian noise at different scales, as it makes the data distribution more amenable to score-based generative modeling. They propose to estimate the score corresponding to all noise levels by training a single noise-conditioned score network (NCSN). They derive $\nabla_x \log(p_{\sigma_t}(x))$ as $\nabla_{x_t} \log p_{\sigma_t}(x_t \mid x) = -\frac{x_t - x}{\sigma_t}$ by choosing the noise distribution to be $p_{\sigma_t}(x_t \mid x) = \mathcal{N}(x_t; x, \sigma_t^2 \cdot \mathbf{I})$ where $x_t$ is a noised version of $x$. Thus, for a given sequence of Gaussian noise scales $\sigma_1 < \sigma_2 < \cdots < \sigma_T$, Equation (8) can be written as:

$$\frac{1}{T} \sum_{t=1}^{T} \lambda(\sigma_t) \, \mathbb{E}_{p(x)} \mathbb{E}_{x_t \sim p_{\sigma_t}(x_t \mid x)} \left\| s_\theta(x_t, \sigma_t) + \frac{x_t - x}{\sigma_t} \right\|_2^2, \tag{9}$$

where $\lambda(\sigma_t)$ is a weighting function. The inference is done using an iterative procedure called "Langevin dynamics" [39, 40]. Langevin dynamics design an MCMC procedure to sample from a distribution $p(\mathbf{x})$ using only its score function $\nabla_\mathbf{x} \log p(\mathbf{x})$. Specifically, to move from a random sample $\mathbf{x}_0 \sim \pi(\mathbf{x})$ towards samples from $p(\mathbf{x})$, it iterates the following:

$$x_i = x_{i-1} + \frac{\gamma}{2} \nabla_x \log p(x) + \sqrt{\gamma} \cdot \omega_i, \tag{10}$$

where $\omega_i \sim \mathcal{N}(0, \mathbf{I})$, and $i \in \{1, \ldots, N\}$. When $\gamma \to 0$ and $N \to \infty$, $\mathbf{x}_i$ samples obtained from this procedure converge to a sample from $p(\mathbf{x})$. The authors of [35] propose a modification of this algorithm nomenclature as the annealed Langevin dynamics algorithm since the noise scale $\sigma_i$ decreases (anneals) gradually over time to mitigate some pitfalls and failure modes of score matching [41].

### 2.3.2 Stochastic Differential Equations (SDEs)

Similar to the aforementioned two approaches, score-based generative models (SGMs) [36] transform the data distribution $q(x_0)$ into noise. However, by generalizing the number of noise scales to infinity, one can view the previous probabilistic models as a discretization of an SGM. We know that many stochastic processes, such as the diffusion process, are the solution to a stochastic differential equation (SDE) in the following form:

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w}, \tag{11}$$

where $\mathbf{f}(.,t)$ is the drift coefficient of the SDE, $g(t)$ is the diffusion coefficient, and w represents the standard Brownian motion. Let $\mathbf{x}_0$ be the uncorrupted data sample, and $\mathbf{x}_T$ denote the perturbed data approximating standard Gaussian distribution. For the given forward SDE, there exists a reverse time SDE running backward where, by starting with a sample from $p_T$ and reversing this diffusion SDE, we can obtain samples from our data distribution $p_0$. The reverse-time SDE is given as:

$$d\mathbf{x} = \left[\mathbf{f}(\mathbf{x},t) - g^2(t)\nabla_x \log p_t(x)\right] dt + g(t)d\bar{\mathbf{w}}, \tag{12}$$

where $dt$ is the infinitesimal negative time step, and $\bar{w}$ is the Brownian motion running backward. In order to numerically solve the reverse-time SDE, one can train a neural network to approximate the actual score function via score matching [35, 36] to estimate $s_\theta(\boldsymbol{x},t) \simeq \nabla_{\boldsymbol{x}} \log p_t(\boldsymbol{x})$ (denoted red in Equation (12)). This score model is trained using the following objective:

$$\mathcal{L}(\theta) = \mathbb{E}_{\mathbf{x}(t)\sim p(\mathbf{x}(t)|\mathbf{x}(0)),\mathbf{x}(0)\sim P_{\text{data}}} \left[\frac{\lambda(t)}{2} \left\|s_\theta(\mathbf{x}(t),t) - \nabla_{\mathbf{x}(t)} \log p_t(\mathbf{x}(t) \mid \mathbf{x}(0))\right\|_2^2\right], \tag{13}$$

where $\lambda$ is a weighting function, and $t \sim \mathcal{U}([0,T])$. Notably, $\nabla_{\boldsymbol{x}} \log p_t(\boldsymbol{x})$ is replaced with $\nabla_{\boldsymbol{x}} \log p_{0t}(\boldsymbol{x}(t) \mid \boldsymbol{x}(0))$ to circumvent technical difficulties.

The sampling process of SDEs can be accomplished by applying any numerical method to Equation (12). Three commonly used techniques are discussed in detail below.

1. **Euler-Maruyama (EM) method:** Using a simple discretization technique in which $dt$ is replaced with $\Delta t$ and $d\bar{w}$ with Gaussian noise $z \sim \mathcal{N}(0, \Delta t \cdot I)$, Equation (12) can be solved.

2. **Prediction-Correction (PC) method:** In this method, the prediction and correction process takes place in a nested loop, in which the prior data is first predicted and then corrected in several steps. The predictor can be solved using EM. Since the corrector can be any score-based Markov Chain Monte Carlo (MCMC) method, including annealed Langevin dynamics, it can be solved utilizing Langevin dynamics in Equation (10).

3. **Probability Flow ODE (ODE) method:** SDE equations in Equation (11) can be written as ODE equations as follows:

$$d\mathbf{x} = \left[\mathbf{f}(\mathbf{x},t) - \frac{1}{2}g^2(t)\nabla_x \log p_t(x)\right] dt. \tag{14}$$

Hence, by solving the ODE problem, $x_0$ can be found. However, while ODE is a quick solver, it lacks a stochastic term to correct errors, resulting in slightly diminished performance.

# 3 Clinical Importance

Generative models have significantly impacted the field of medical imaging, where there is a strong need for tools to improve the routines of clinicians and patients. Concretely, the complexity of data collection procedures, the lack of experts, privacy concerns, and the compulsory requirement of authorization from patients create a major bottleneck in the annotation process in medical imaging. This is where generative models become advantageous [42]. Several perspectives have driven our interest in generative diffusion models for medical imaging. In the medical field, many datasets suffer from severe class imbalance due to the rare nature of some pathologies. Diffusion models can alleviate this restriction by generating diverse realistic-looking images to leverage in the medical field. Furthermore, generating synthetic medical images has substantial educational value. With its ability to produce a limitless source of unique instances of different medical imaging modalities, diffusion models can satisfy educational demands by constructing distinct synthetic samples for teaching and practice. Additionally, these artificial images can mitigate data security concerns associated with using patient data in public settings. These artificial images can also solve a particular significant difficulty in training deep neural networks for medical applications. Generally, the annotation of medical images is a lengthy and costly process that necessitates the assistance of an expert. Hence, using diffusion models to generate synthetic samples can alleviate the problem of medical data scarcity to a great extent. A case study using [43] to generate histopathology images with rare cancer subtypes is described in Figure 4.

While independent use of synthetic data generated by generative models is still in its early stages, studies have shown promising results in utilizing them in real-world scenarios. Studies such as Goncalves et al. [44] have evaluated different generative methods for creating synthetic electronic health records and found that some have the potential to be useful in practice as they produce synthetic samples that have similar statistical properties to real data without compromising patient privacy. In another study, Chen et al. [45] found that using both synthetic and real data to train classifiers for histology images improves performance compared to using only real data. Furthermore, studies conducted by Akrout et al. [46] have demonstrated that the utilization of synthetic images generated by diffusion models improves the accuracy of skin classifiers, and models that are trained using a combination of synthetic and real data perform better than those trained using only one
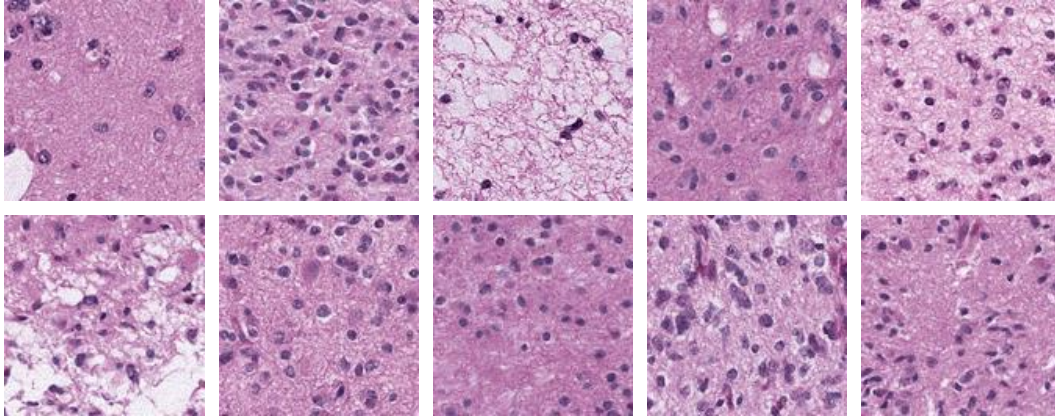
Figure 4: Ten synthetic histopathology images generated by MFDPM [43].

data source. Moghadam et al. [43] conducted a study to assess the morphological properties of synthetic and actual images by administering a survey to two pathologists with different levels of expertise. The results revealed that the pathologists could not distinguish real from synthetic images generated by the diffusion model, and the majority of the small percentage they correctly identified had lower confidence levels. Overall, the research findings indicate that the synthetic images are convincingly similar to real images and can be effective in training models for medical research. In addition, incorporating both synthetic and real data has the potential to improve performance in a variety of applications, as synthetic data serves as a powerful augmentation to real data.

The use of generative models, particularly diffusion models, as a generative prior in biomedical inverse imaging problems is a recent development in the field. In inverse imaging problems, the goal is to infer the underlying physical properties of an object or system from observations or measurements. Conventional methods for inverse imaging problems typically use model-based priors, which are assumptions about the properties of the object or system that are used to guide the reconstruction process. Generative models, specifically diffusion models, can be used as an alternative to these model-based priors since they can provide a more accurate representation of the data distribution [47, 48]. This is because generative models are trained on real data, which means that they can learn the complex patterns and structures present in the data. As a result, they can provide a more accurate prior for the reconstruction process, leading to more accurate reconstructions of the object or system. There are other reasons why diffusion models are proliferating in the field of biomedical inverse imaging problems. One reason is that they can be used to handle high-dimensional and complex data, such as medical images, which are difficult to model using conventional methods. Furthermore, diffusion models can also provide a more efficient and accurate way to infer the underlying physical properties of the object or system and can handle uncertainty and noise in the measurements or observations.

In conclusion, diffusion models have proven to be a valuable and versatile tool that can be utilized in clinical settings and address a wide range of imaging challenges, and it is expected that their use will continue to expand in the future, providing new opportunities for medical imaging and research.

# 4 Diffusion Models in Action

Providing a taxonomy for diffusion models more or less follows the same route as other techniques for medical imaging analysis. We provide, however, detailed additional information for each sub-category paper in Figure 5. In this section, we explore diffusion-based methods, which are proposed to solve any disentanglement from the medical imaging analysis in **seven** application categories, as in Figure 5: **(I)** Image-to-Image Translation, **(II)** Image Reconstruction, **(III)** Image Registration, **(IV)** Image Classification, **(V)** Image Segmentation, **(VI)** Image Denoising, **(VII)** Image Generation, **(VIII)** Anomaly Detection, and **(IX)** multi-disciplinary applications, named Other Applications. Figure 5 represents a collection of numerous studies for each category with extensive information on each study, such as the modality of study, the organ of interest, and the specific algorithm utilized in the reverse process of the diffusion model for study. Finally, in Section 4.10, we discuss the overall algorithms used in the studies and try to shed light on the main novelty and contribution of the papers in Table 1.
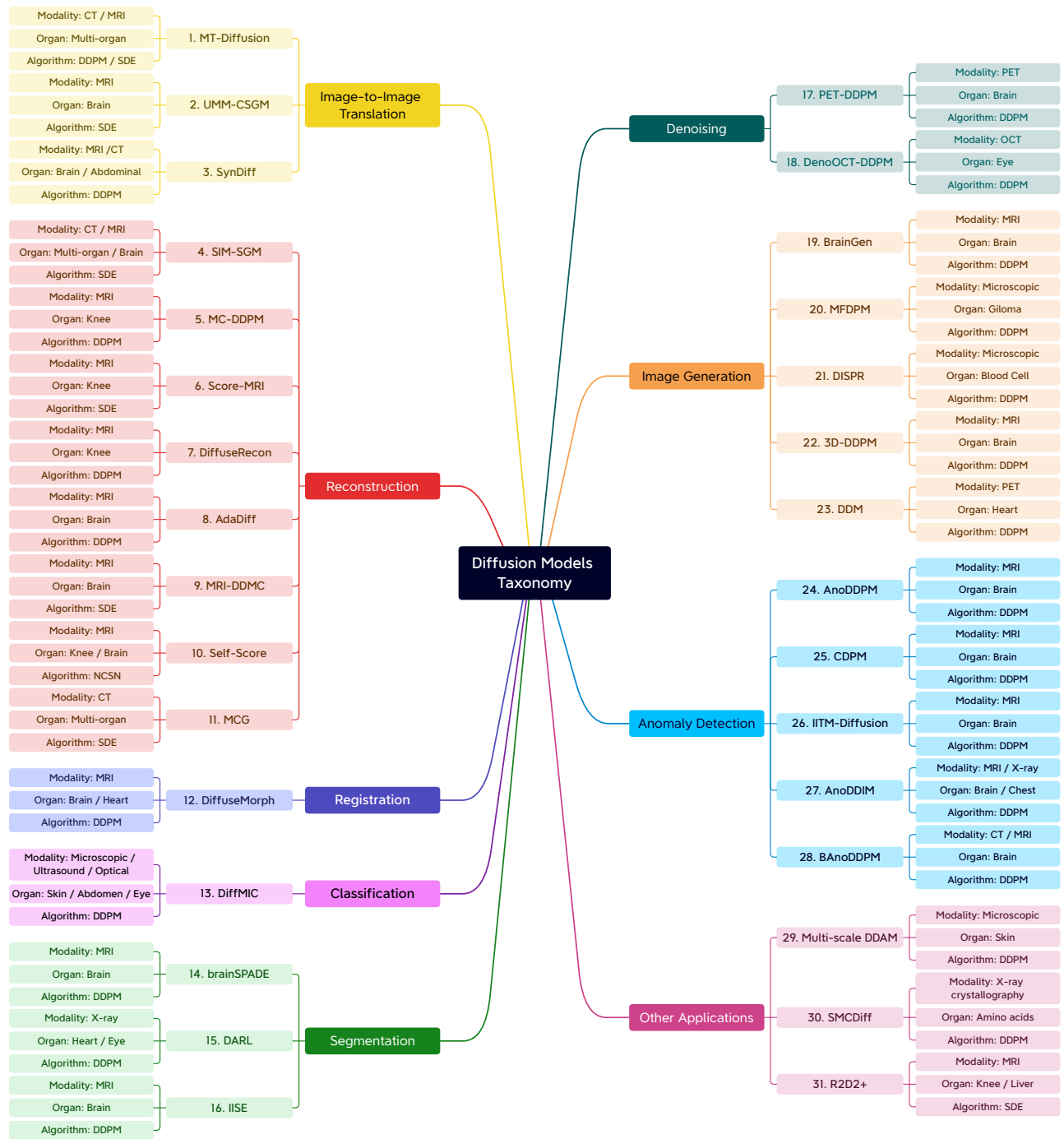
Figure 5: The proposed taxonomy for diffusion-based medical imaging research is built on nine sub-fields: I) Image-to-Image Translation, II) Image Reconstruction, III) Image Registration, IV) Image Classification, V) Image Segmentation, VI) Image Denoising, VII) Image Generation, VIII) Anomaly Detection, and IX) Multi-disciplinary applications, named Other Applications. For the sake of brevity, we utilize the prefix numbers in the paper's name in ascending order and denote the reference for each study as follows: 1. [49], 2. [50], 3. [51], 4. [52], 5. [53], 6. [48], 7. [54], 8. [55], 9. [56], 10. [57], 11. [58], 12. [59], 13. [60], 14. [61], 15. [62], 16. [63], 17. [64], 18. [65], 19. [66], 20. [43], 21. [67], 22. [68], 23. [69], 24. [70], 25. [71], 26. [72], 27. [16], 28. [73], 29. [74], 30. [75], 31. [76].

## 4.1 Image-to-Image Translation

Acquiring multi-modality images for diagnosis and therapy is often crucial. Also, we may miss modalities in some conditions. Diffusion models have shown favorable results for generating missing modalities utilizing cross-modalities and producing ones using other modality types, e.g., translating from MRI to Computed Tomography (CT).

CT and MRI are two of the most prevalent imaging types. CT, however, is limited in displaying the intricacies of the images for soft tissue injuries. Hence, a subsequent MRI may be needed for a conclusive diagnosis after receiving the initial CT results. Nevertheless, in addition to being time-consuming and costly, this process may also cause misalignment between MRI and CT images. To this end, Lyu et al. [49] take advantage of the recently introduced DDPMs [77, 34] and score-based diffusion models [36] in solving the translation problem between two modalities, i.e., from MRI to CT. In particular, they present conditional DDPM and conditional SDE, in which their reverse process is conditioned on T2w MRI images and conducts comprehensive experiments. The authors adopt the DDPM and SDE with three different sampling methods (EM, PC, and ODE) that are explained in Section 2.3.2 and compare their results with the existing GAN-based [78] and CNN-based [79] methods. Their extensive experiments on the Gold Atlas male pelvis dataset [80] demonstrate that diffusion models outperform both CNN and GAN-based methods in terms of Structural Similarity Index Measure (SSIM) and Peak Signal-to-Noise Ratio (PSNR). Additionally, they employ the Monte Carlo (MC) method to investigate the uncertainties of diffusion models; in this technique, the model outputs ten times, and the average yields the final result. Qualitative results are depicted in Figure 6.
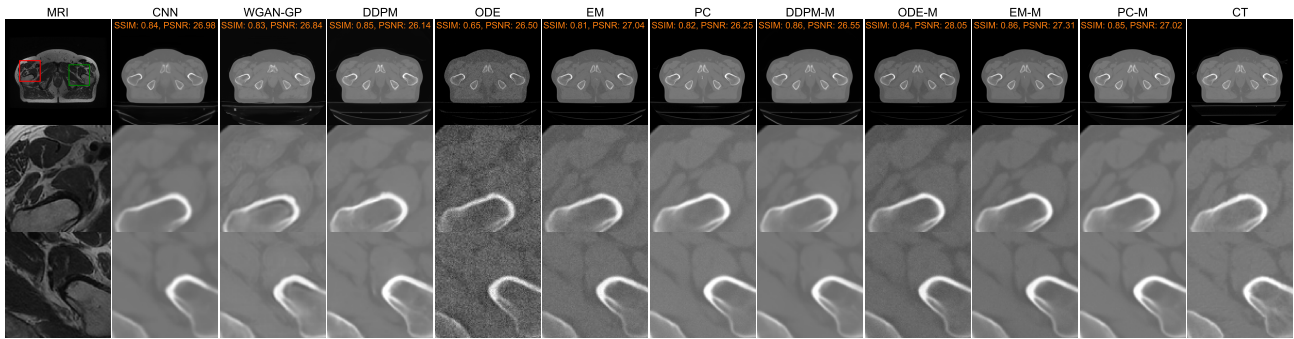


Figure 6: Visual and quantitative comparison of different methods of translating an MR image to a CT image conducted by [49]. The second and third rows indicate zoomed regions in the original image, delineated by red and green boxes. Results after applying the Monte Carlo method (taking an average from ten outputs) using DDPM, ODE, EM, and ODE sampling methods are displayed with DDPM-M, ODE-M, EM-M, and PC-M, respectively.

To cope with the missing modality issue, Meng et al. [50] propose a unified multi-modal conditional score-based generative approach (UMM-CSGM), which synthesizes the missing modality based on all remaining modalities as conditions. The proposed model is a conditional format of SDE [36], in which it employs only a score-based network to learn different cross-modal conditional distributions. Experiments on the BraTS19 dataset [81, 82, 83], which contains four MRI modalities for each subject, show that the UMM-CSGM is capable of generating missing-modality images with higher fidelity and structural information of the brain tissue compared to SOTA methods [84, 85, 86, 87, 88].

For the task of translating medical images, diffusion models inherently lack the ability to maintain the structural information accurately, which is due to the fact that the structured details of the source domain images are lost during the forward diffusion process and cannot be completely recovered through learned reverse diffusion, although it would be crucial to preserve the integrity of anatomical structures in medical images. To mitigate the aforementioned problem, Li et al. [89] introduce a novel approach for structure-preserving image translation using frequency-domain filters, the Frequency-Guided Diffusion Model (FGDM). The proposed FGDM architecture enables zero-shot learning and can be exclusively trained on target domain data. Furthermore, their model does not require exposure to source domain data for training and can be directly deployed for source-to-target domain translation. The proposed method shows significant advantages in zero-shot medical image translation over the baseline SOTA methods.

## 4.2 Reconstruction

Medical image reconstruction plays a critical role in medical imaging. Its main goal is to produce high-quality medical images for clinical use while minimizing costs and patient risk [90, 91]. Medical imaging modalities such as CT and MRI are medicine's most popular imaging tools. However, their physics restricts their efficacy, directly affects their performance, and degrades their desired results. The high resolution and complete result acquisition from the subject requires a higher
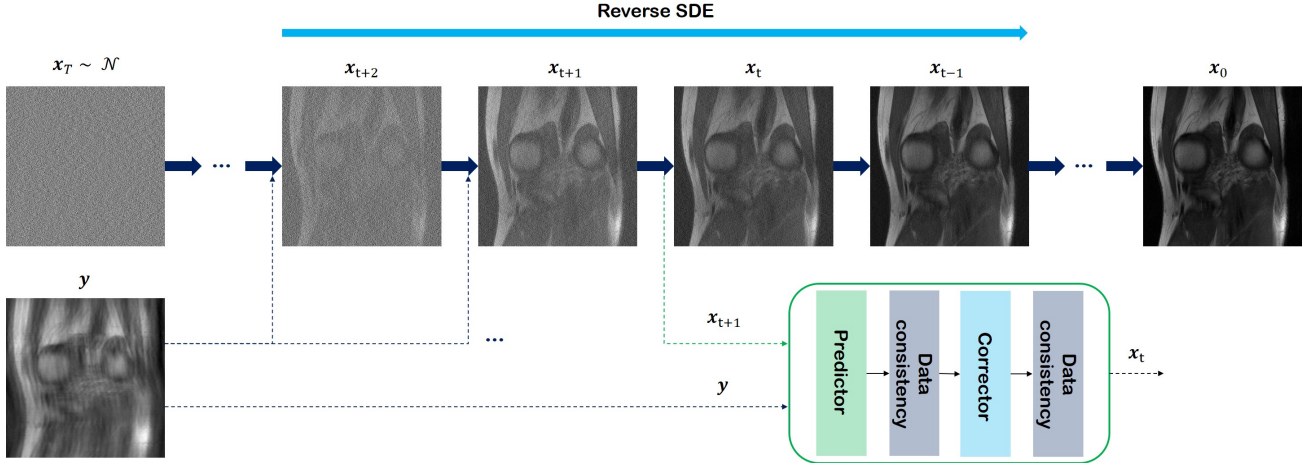
Figure 7: An overview of Score-MRI [48]. $x_T$ is derived from the pre-trained prior distribution, and $x_0$ is retrieved by sequentially applying the predictor, data consistency, corrector, and data consistency steps, given the measurement through the reverse SDE procedure.

radiation dose and a relatively longer resting time in the tube, which are only partially applicable due to the health precautions and the patient's relentless. Therefore, faster acquisition speed in medical imaging techniques like CT, Positron Emission Tomography (PET), and MRI is essential for reducing exam time, improving access to imaging services and reducing waiting times, and, more importantly, producing accurate images, particularly in dynamic studies that require fast imaging sequences. Accordingly, their radiation exposure reduces from the standard dose, or the imaging process is done in an under-sampled or sparse-view manner [92, 93, 94, 95]. To diminish these drawbacks, e.g., low Signal-to-Noise Ratio (SNR) and Contrast-to-Noise Ratio (CNR), medical image reconstruction must overcome the challenges mentioned and solve this ill-posed inversion problem [96]. This section overviews the diffusion-based paradigms for medical image reconstruction and enhancement.

MRI is a popular non-invasive imaging utility in medical diagnosis treatment, but due to its innate physics, it is a time-consuming process of imaging sessions in which the movement of patients results in various artifacts in images. Therefore, to decrease bedtime and to accelerate the reverse problem-solving from the spatial domain (or *k-space*) to image level, miscellaneous solutions are provided in the supervised-learning concept. However, these methods are not robust to distribution changes or drifts in their train/test sets. Jalal et al. [47] proposed the first study in the MRI reconstruction domain via Compressed Sensing with Generative Models (CSGM). To this end, CSGM trains the score-based generative models [41] on MRI images to utilize as prior information for the inversion pathway in reconstructing realistic MRI data from under-sampled MRI in posterior sampling scheme with Langevin dynamics [97]. CSGM [47] demonstrated its better performance over fastMRI [98] and Stanford MRI [99] datasets with SSIM and PSNR metrics in comparison with end-to-end supervised-learning paradigms.

Chung et al. [48] propose a score-based diffusion framework that solves the inverse problem for image reconstruction from accelerated MRI scans, as depicted in Figure 7. In the first step, a single continuous time-dependent score function with denoising score matching is trained only with magnitude images. Then, in the reverse SDE process, the Variance Exploding (VE)-SDE [36] is exploited to sample from the pretrained score model distribution, conditioned on the measurement. Afterward, the image is first split into real and imaginary components at each step. Each part is fed into the predictor, followed by data consistency mapping to reconstruct the image. The obtained image is split again, and the correcter and the data consistency mapping are applied to each part to compensate for errors during the diffusion and reconstruct the improved image, respectively. Results demonstrate that the proposed model outperforms the previous SOTA methods [100, 101, 98] and can even reconstruct the data, which is considerably outside the training distribution with high fidelity, e.g., reconstructing anatomy not seen during training. In addition, the proposed framework has shown to be very effective for reconstructing the image when multiple coils exist. For the aforementioned problem, they present two approaches: (1) they reconstruct each coil image in parallel; (2) they take into account the correlation between the coil images by injecting the dependency between them at each given step during reverse SDE. Then, the final image is acquired by taking the sum-of-root-sum-of-squares (SSoS) [102] of each reconstructed coil image. Although these two techniques have shown great results qualitatively and practically, they are time-consuming.

Liu et al. [103] addressed the limited-angle CT reconstruction with a model-based DDPM paradigm called DOLCE. Based on the Fourier slice theorem, the conventional algorithm for mapping CT images from sinograms is Filtered Back Projection

(FBP) [104]. Therefore, the limited-angle measurements can lead to the loss of Fourier measurements and, thus, degraded reconstruction results. However, due to the ill-posed nature of the reconstruction framework, the DDPM can not be utilized directly. DOLCE [103] incorporates the output of FBP on the limited sinograms as prior information to condition the diffusion model. Further, due to the consistency conditions presented by sinograms, DOLCE imposes a consistency term within an additional step in the denoising iteration under $\ell_2$-norm loss in the inference step. The results of the Kidney CT (C4KC-KiTS) dataset [105] regarding SSIM and PSNR metrics demonstrate that DOLCE is effective in producing sharp CT images.

### 4.3 Registration

Deformable image registration is a critical medical image analysis technique focused on identifying non-rigid relationships between a pair of moving and fixed images. It plays a vital role when image shapes change due to factors such as subject, scanning time, and imaging modality. Traditional registration algorithms can be computationally expensive, while deep-learning methods are faster but still struggle with realistic continuous deformations.

To overcome these limitations, Kim et al. [59] introduce a novel diffusion-based method called DiffuseMorph. DiffuseMorph has two main networks - a diffusion network and a deformation network - both trained in an end-to-end fashion. The diffusion network scores the deformation between the moving and fixed images, while the deformation network uses this information to estimate the deformation field. The information used contains spatial information, allowing the generation of deformation fields along a continuous trajectory from the moving to the fixed image. The moving image is then transformed into a deformed image using the generated deformation fields and the spatial transformation layer (STL) [106]. In the inference phase, the model can provide both image registration and generation tasks. The experimental results affirm the high accuracy of the proposed method in registering both 2D facial expressions [107] and 3D medical images [108, 109].

### 4.4 Classification

The classification task is of great importance in medical image analysis, as it allows for accurately identifying and characterizing different structures and anomalies within medical images. Its ability to aid medical professionals in interpreting large amounts of complex data has the potential to revolutionize the healthcare industry [110]. Despite this potential, adopting diffusion models to enhance classification results remains a significant challenge that needs to be addressed further.

DiffMIC [60] presents a new approach for classifying different medical image modalities using diffusion models, as shown in Figure 8. It first encodes the input image into a feature embedding space and uses a Dual-granularity Conditional Guidance (DCG) model to capture global and local prior information. The ground truth and two priors are then diffused to generate three noisy variables, which are concatenated with their corresponding priors and projected to a latent space to obtain three feature embeddings. The denoising U-Net integrates these embeddings with the image feature embedding and predicts the noise distribution for each embedding. Next, the feature embeddings obtained are projected back to their original dimensions. In order to estimate the amount of noise added to both the global and local priors, as well as the ground truth, DiffMIC utilizes maximum-mean discrepancy (MMD) [111, 112] regularization loss and mean squared error (MSE) loss, respectively. In the inference stage, the DCG model is used to acquire dual priors from the input image, and the final prediction is denoised iteratively using the trained U-Net conditioned by the dual priors and the image feature embedding. Overall, DiffMIC presents a promising approach for accurately classifying medical images using diffusion models for the three considered tasks: placental maturity grading using ultrasound images, skin lesion classification using dermatoscopic images, and diabetic retinopathy grading using fundus images.

### 4.5 Segmentation

Image segmentation is a vital task in computer vision, which investigates simplifying the complexity of the image by decomposing an image into multiple meaningful image segments [113, 114]. Specifically, it facilitates medical analysis by providing beneficial information about anatomy-related areas. However, deep learning models often require vast amounts of diverse pixel-annotated training data in order to produce generalizable results [115, 116]. Nonetheless, the number of images and labels accessible for medical image segmentation is restricted due to the time, cost, and expertise required [117, 118, 119]. To this end, diffusion models have emerged as a promising approach in image segmentation research by synthesizing the labeled data and obviating the necessity for pixel-annotated data.

brainSPADE [61] proposes a generative model for synthesizing labeled brain MRI images that can be used for training segmentation models. brainSPADE is composed of a label generator and an image generator sub-model. The former is responsible for creating synthetic segmentation maps, and the latter for synthesizing images based on generated labels. In the label generator, the input segmentation map is first encoded during training using a spatial VAE encoder and builds a
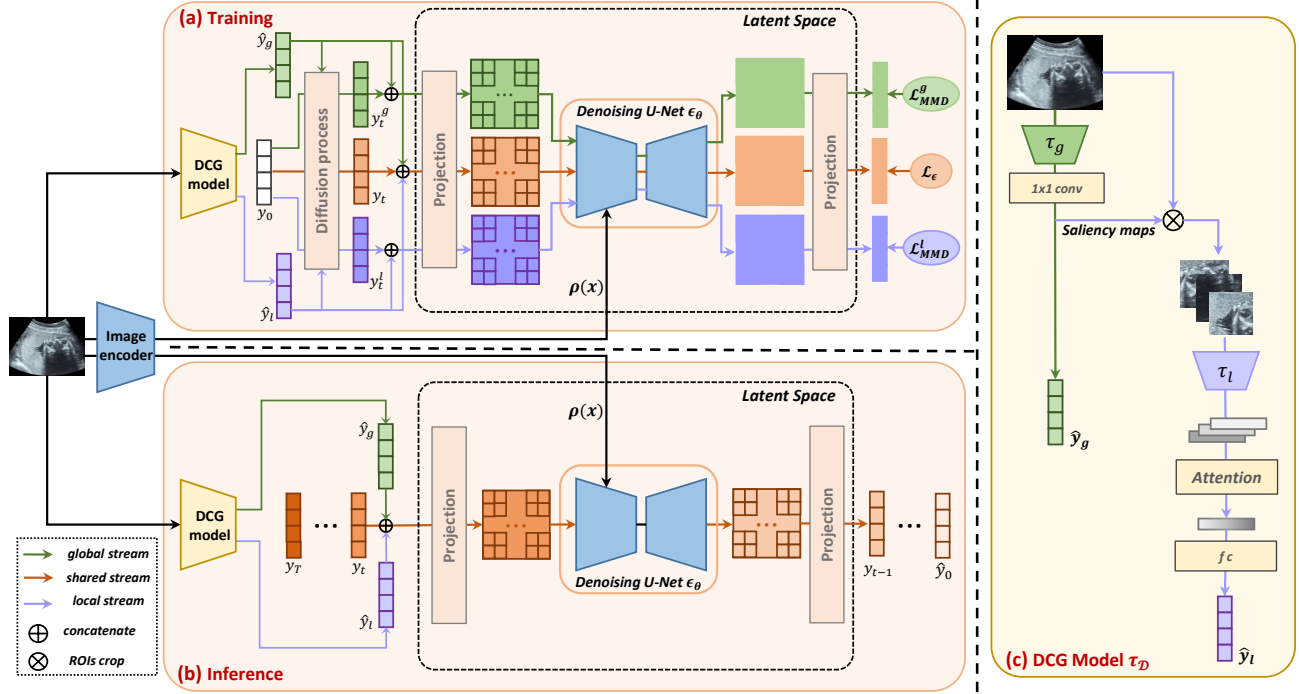
Figure 8: The image depicts the two stages of DiffMIC [60] - the training (a) and inference (b) stages - as well as the Dual-granularity Conditional Guidance (DCG) model shown in (c). The DCG model generates global and local priors from the raw image and Regions of Interest (ROIs), which serve to guide the diffusion process.

latent space. The compressed latent code is then diffused and denoised via LDMs [120] and produces an efficient latent space in which imperceptible details are ignored, and semantic information is highlighted more. A spatial VAE decoder then constructs the artificial segmentation map via the latent space. In the image generator, Fernandez et al. [61] take advantage of SPADE [121], a VAE-GAN model, to build a style latent space from the input arbitrary style and use it with the artificial segmentation map to decode the output image. nnU-Net [122] was leveraged to examine the performance. Findings show that the model achieves comparable results when trained on synthetic data compared to that trained on factual data, and their combination significantly improves the model result.

Kim et al. [62] propose a novel diffusion adversarial representation learning (DARL) model for self-supervised vessel segmentation, aiming to diagnose vascular diseases. There are two main modules in the proposed DARL model: a diffusion module, which learns background image distribution, and a generation module, which generates vessel segmentation masks or synthetic angiograms using a switchable SPADE algorithm [121]. Figure 9 illustrates two ways in which this method can be applied. In path (A), a real noisy angiography image $x_{t_a}^a$ is input into the model to produce a segmentation mask $\hat{s}^v$, and the SPADE switch is off. In path (B), a real noisy background image $x_{t_b}^b$ is fed into the model, and the SPADE becomes active and receives a vessel-like fractal mask, generating a synthetic angiography image $\hat{x}^a$. Then, by giving the generated synthetic angiography images into the path (A), a cycle is formed, which helps in learning the vessel information. In addition, during inference, path (A) is performed at one step, where the model produces the mask by only inputting the noisy angiography image into the model. Results verify the generalization, robustness, and superiority of the proposed method compared to SOTA un/self-supervised learning approaches.

In addition to the studies mentioned earlier, Rahman et al. [123] introduced the CIMD framework, a single probabilistic diffusion-based model, to address the ambiguous medical image segmentation task. Deterministic medical image segmentation frameworks such as [119, 124] produce a pixel-wise uncertainty, but the results are not consistent [125]. Also, from a medical aspect, medical image segmentation should not be considered just as a pixel-wise task. In clinical practice, analyzing organs or other structures from medical images is not a deterministic pixel-wise process but underlies the assessment of the whole image or, on a smaller scale, assessing the neighboring pixels' diversity. The stochastic sampling step in the diffusion model can produce diverse and multiple masks. During its training step, the CIMD [123] utilizes the noisy segmentation ground-truth masks concatenated to the original image to hinder the conventional usage of the diffusion process in the segmentation task from producing more resilient results, rather than arbitrary masks. CIMD outperforms the
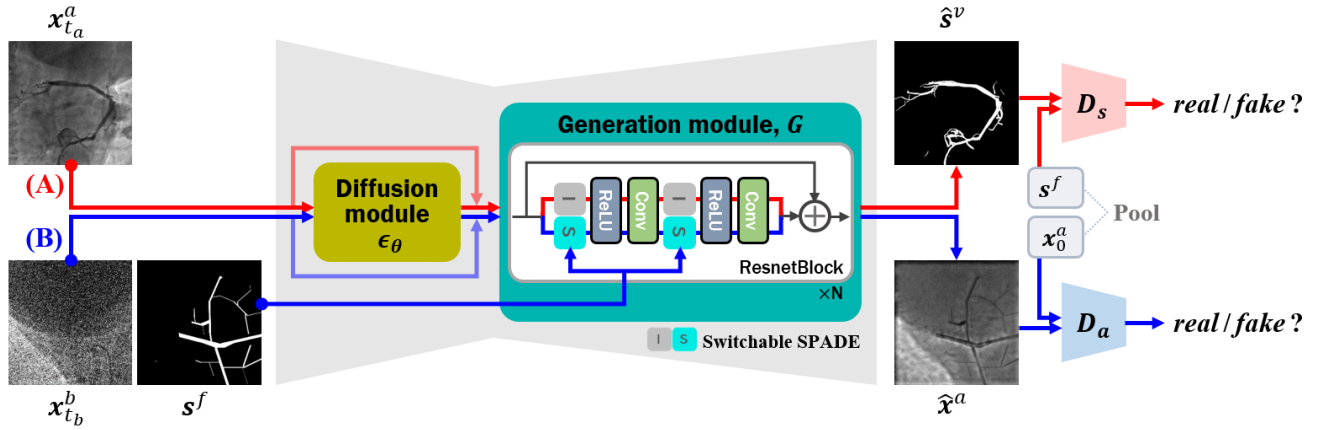
Figure 9: An overview of DARL [62]. Path (A) involves feeding a real noisy angiography image through the model to generate a segmentation map. Path (B) incorporates passing a noisy background image alongside a vessel-like fractal mask through the model to synthesize a synthetic angiography image.

probabilistic U-Net model [126] in terms of Collective Insight Score, which Rahman et al. [123] introduced in their work, which investigates three datasets (one private and two publicly available ones) with different modalities [127, 128].

Bieder et al. [129] present a memory-efficient patch-based diffusion model called PatchDDM that can be applied to large 3D volumes, making it suitable for medical tasks. The authors evaluate PatchDDM on the tumor segmentation task of the BraTS2020 [83, 82] dataset and demonstrate that it can generate meaningful three-dimensional segmentation while requiring less computational resources than traditional diffusion models.

## 4.6 Denoising

The major challenge in medical imaging is obtaining an image without losing important information. The images obtained may be corrupted by noise or artifacts during the acquisition and/or further processing stages [130, 131]. Noise reduces the image quality and is especially significant when the imaged objects are small and have relatively low contrast [132]. Due to the nature of generative models, diffusion models are convenient for diverse denoising problems [133, 134]. In this section, we will explore the contribution of diffusion models to this task.

Hu et al. [65] utilized a DDPM [34] to despeckle Optical Coherence Tomography (OCT) volumetric retina data in an unsupervised manner, denoted as DenoOCT-DDPM. OCT imaging utility suffers from limited spatial-frequency bandwidth, which leads to the resulting images containing speckle noise. Speckle noise hinders the ophthalmologist's diagnosis and can severely affect the visibility of the tissue. The classic methods, such as averaging multiple b-scans at the same location, have extreme drawbacks, such as prolonged acquisition time and registration artifacts. Due to the multiplicative properties of speckle noise, these methods enrich rather than reduce the noise. Deep-based models perform outstandingly. This performance, however, depends on the availability of noise-free images, which is a rare and costly process to acquire. To this end, DenoOCT-DDPM [65] utilizes DDPM's feasibility in noise patterns rather than real-data pattern. Therefore, they use a self-fusion [135] as a preprocessing step to feed the DDPM with a clear reference image and train the parameterized Markov chain (see Figure 10a). Their investigation demonstrated the SOTA results over the Pseudo-Modality Fusion Network (PMFN) [136], which uses information from the single-frame noisy b-scan and a pseudo-modality that is created with the aid of the self-fusion [135] method, regarding Signal-to-Noise Ratio (SNR) metric. The qualitative results over PMNF depicted in Figure 10b (represented in multiple acquisition SNRs) endorse the ability of diffusion models to remove the speckle noise while preserving the fine-grained features like small vessels.

PET is a non-invasive imaging utility that plays a crucial role in cancer screening and diagnosis. However, as with OCT devices, PETs suffer from low SNR and resolution due to the low beam count radiation to patients. Deep learning methods over PET image denoising have advanced, but over-smoothing is a prominent drawback of CNN-based approaches. Therefore, Conditional Generative Adversarial Networks (CGANs) [137, 138] neutralize the mentioned deficiency but still depend on training and test set distributions. Gong et al. [64] proposed the DDPM-based framework for PET denoising in collaboration with an assistive modality embedding as prior information to DDPM formulation, namely PET-DDPM. Gong et al. used $^{18}$F-FDG and $^{18}$F-MK-6240 datasets for PET and MR modalities, respectively. PET-DDM is a multi-disciplinary study investigating the excessive modalities of collaboration in learning noise distribution through PET images. This intuition follows the original paper in a generative paradigm [11] with a guided classifier to converge the learned distribution
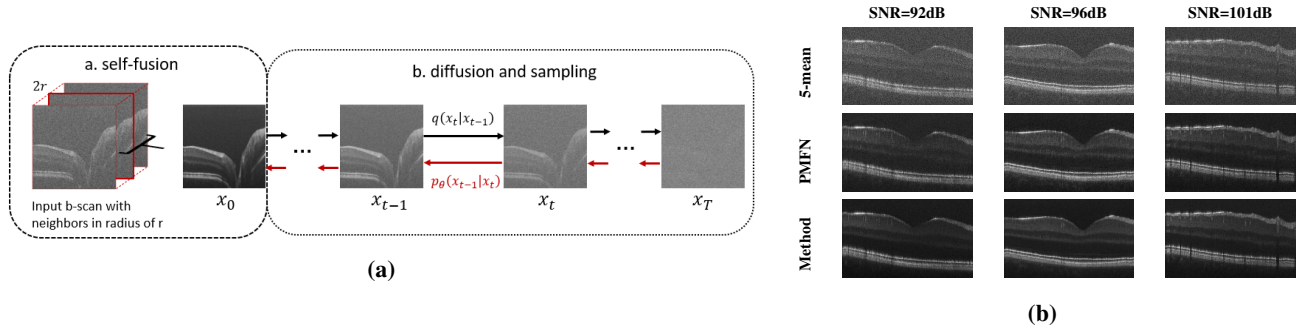
Figure 10: (Left) General pipeline of DenoOCT-DDPM [65]. To feed the diffusion model more low-noise reference images, DenoOCT-DDPM applies a self-fusion [135] on the red canvas b-scan in (a) with neighboring b-scans for higher SNR b-scan. (b) indicates a straightforward DDPM scheme in an unsupervised manner to learn speckle noise distribution in the red arrow sampling stream. (Right) Ultimate visual comparison of DenoOCT-DDPM [65] for denoising b-scans of volumetric OCT images demolished by speckle noise. DenoOCT-DDPM compared itself with PMFN [136] and the average of 5 successive b-scans as the ground truth. PMFN results in good feature preservation in diverse retinal layers, but it can easily over-smooth small vessel regions. In contrast, the DenoOCT-DDPM retinal layers are more homogenous than PMFN in diverse acquisition SNRs.
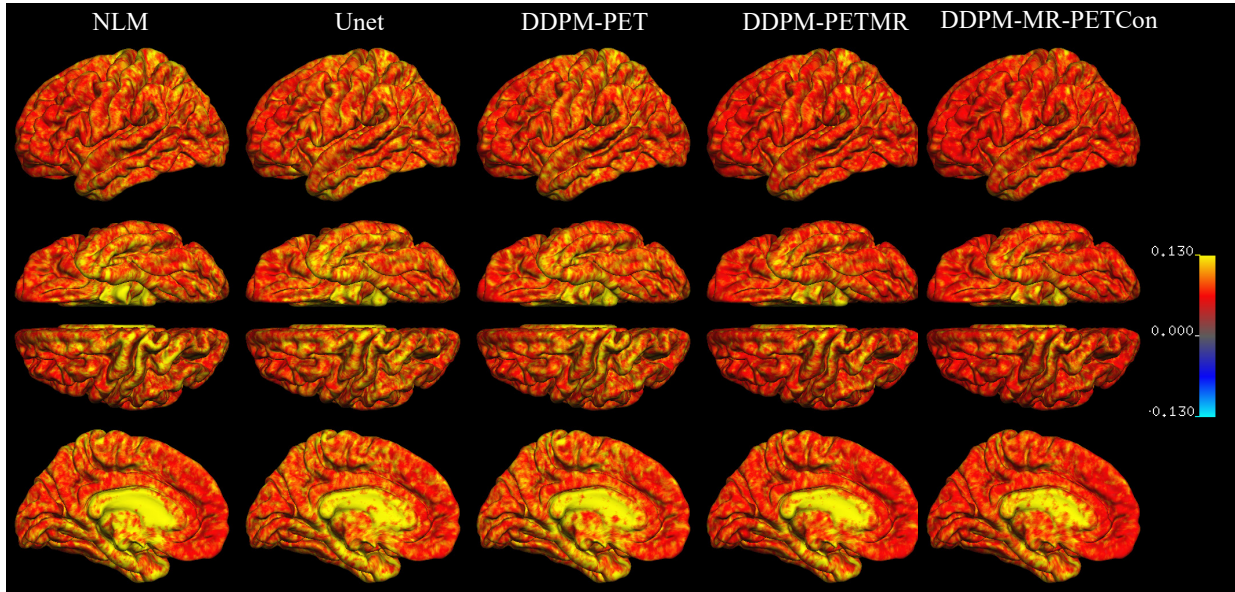


Figure 11: Comparison by different methods of surface error mapping of the left hemisphere from 20 $^{18}$F-MK-6240 test dataset [64]. The results confirm that DDPM-MR-PETCon has the lowest error, followed by the DDPM-PETMR method. DDPM-MR-PETCon is short for using the MR image as network input while using the PET image as a data consistency item, and DDPM-PETMR denotes that PET and MR images are used as network input.

to desired distribution. As qualitatively illustrated in Figure 11, PET-DDM produced SOTA results compared with U-Net [79] based denoising network in terms of PSNR and SSIM.

Diffusion MRI is an important modality for studying oncologic and neurologic biomarkers but suffers from long acquisition times and low SNR. Xiang et al. [139] explore these deficiencies by imposing the self-supervised statistic-based denoising strategy into diffusion models and by performing denoising through the conditional generation process. To this end, their approach consists of three main stages. First, they learn an initial noise distribution in a self-supervised manner. Next, they estimate a noise model with $\mu = 0$ and $\sigma$ (a Gaussian distribution) from the learned noise distribution in the previous step. They apply a $p-$norm to minimize the distance between the $\sigma$ and diffusion sampling noise $\beta$. Ultimately, they train another diffusion model to produce clean images in an unsupervised manner. They elaborated their findings on one private and three public datasets (Sherbrooke 3-Shell [140], Stanford HARDI [141], and Parkinson's Progression Markers Initiative (PPMI) [142]) and reported superior denoising performances.
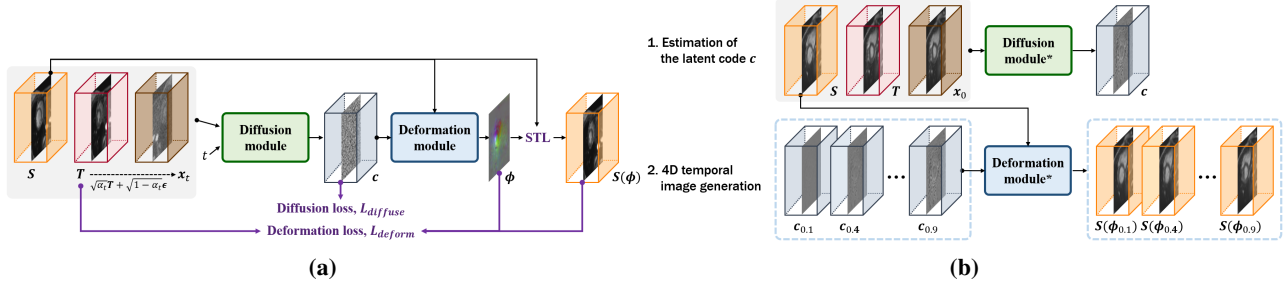
15

Figure 12: **(a)** demonstrates the DDM [69] training phase and **(b)** the inference phase.
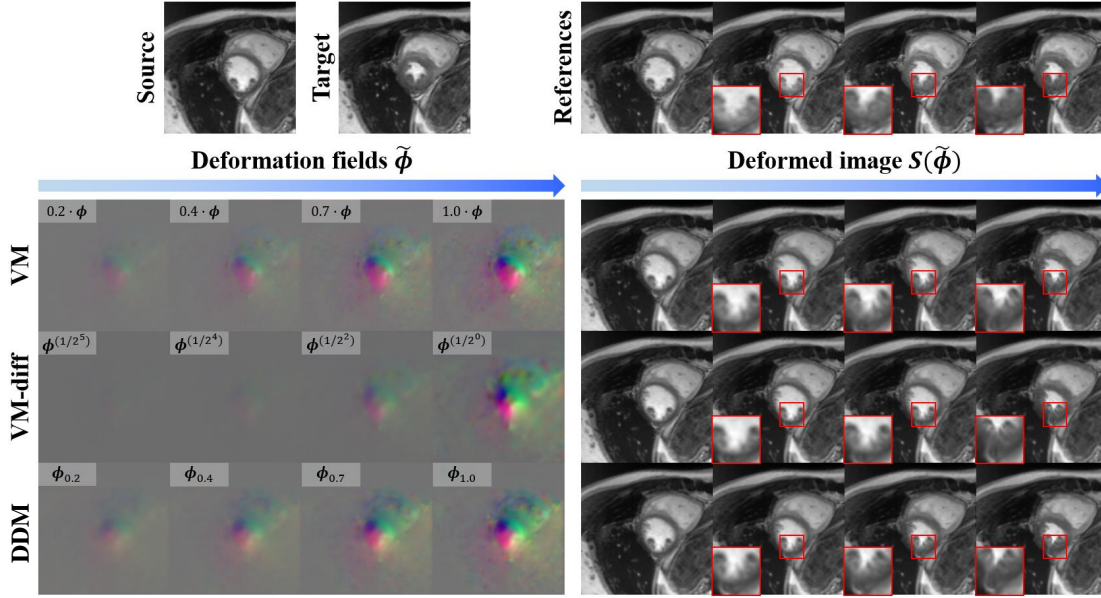


Figure 13: Visual comparison of DDM [69], VM [143], and VM-Diff [144] for generating temporal cardiac images. The deformed intermediate frames $S(\tilde{\phi})$ (right) are constructed using the source and target, and produce the deformation fields $\tilde{\phi}$ (left).

## 4.7 Image Generation

Image generation is one of the primary objectives of diffusion models, which has been widely applied in a variety of styles, including generating synthetic 2D/3D medical images [66, 43, 68, 69], reconstructing 3D cell from 2D cell images [67], etc. This section will outline the diffusion-based approaches for medical image generation.

Using 4D imaging to follow anatomical changes is one of the methods used in medicine to track 3D volumes over time to detect anomalies and disease progression. Such 4D images are primarily obtained with MRI, but this process is relatively time-consuming. Kim et al. [69] recently proposed the Diffusion Deformable Model (DDM), which takes source and target images and generates intermediate temporal frames along the continuous trajectory. This approach comprises two main modules: (i) a denoising diffusion probabilistic model (DDPM) module and (ii) a deformation module. In the DDPM module, a latent code is constructed by learning the source and target images, and in the deformation module, the acquired latent code and the source image are used to render the deformed image. In the training phase, as shown in Figure 12a, the diffusion model, derived from [34], takes source, target, and perturbed target images and outputs a latent code. The learned latent code along the source image is fed into the deformation module, adopted from [143], and creates deformation fields. Then, the spatial transformation layer (STL) [106], with tri-linear interpolation, is employed to warp the source volume using the deformation fields in order to build the deformed source image. Afterward, inference begins with the diffusion module providing the latent code, which contains spatial information from the source toward the target (see Figure 12b). Then, deformed intermediate frames are generated using the deformation module by scaling the latent code with a factor, which is an element of [0, 1]. Additionally, qualitative results are illustrated in Figure 13.

16

Packhauser et al. [145] utilize a latent diffusion model [120] to produce high-quality class-conditional chest X-ray images while proposing a sampling strategy to maintain sensitive biometric information's privacy during the generation process. To assess the potential utility of the generated dataset, the images are evaluated on a thoracic abnormality classification task, and the results show that the proposed approach outperforms GAN-based methods.

Histopathology involves the study of tissues and cells at the microscopic level in order to diagnose diseases and cancer [146]. Histology images, however, are rare for some cancer subtypes, thereby increasing the significance of generative models to fill the void. To this end, Moghadam et al. [43] investigate adopting DDPMs for generating histopathology images for the first time. Specifically, they exploit the DDPMs with genotype guidance to synthesize images containing various morphological and genomic information. To tackle this data consistency problem and enforce the model to focus more on morphological patterns, they first feed input images into a color normalization module [147] to unify the domain of all images. In addition, they apply a morphology levels prioritization module [148] that designates higher weight values to the loss at earlier levels to emphasize perceptual information and lower weights to the loss at later levels, resulting in higher fidelity samples. Experiments on the Cancer Genome Atlas (TCGA) dataset [149] exhibit the superiority of the proposed method compared to GAN-based approaches [150].

In the case of diffusion-based MRI synthesis models, it is common to employ a unimodal approach. However, since they rely on the original image domain, these models often suffer from high memory demands and are less practical for multi-modal synthesis purposes. To mitigate this problem, Jiang et al. [151] propose the first diffusion-based multi-modality MRI synthesis model, namely the Conditioned Latent Diffusion Model (CoLa-Diff). Specifically, they propose an architecture designed to reduce memory consumption by operating in the latent space. In order to address potential issues with compression and noise present in the latent space, they utilize a cooperative filtering approach inspired by collaborative filtering techniques. Moreover, to ensure the preservation of anatomical structures, they consider the inclusion of brain region masks as priors for density distributions to guide the diffusion process. Additionally, they implement an auto-weight adaptation technique to leverage multi-modal information effectively. Their experiments demonstrate that the proposed method outperforms other SOTA MRI synthesis methods, indicating that CoLa-Diff has significant promise as an effective tool for facilitating multi-modal MRI synthesis.

## 4.8 Anomaly Detection

Medical anomaly detection is an important topic in computer vision, aiming to highlight the anomalous regions of the image [152, 153, 154, 155, 156]. Generative models have dramatically shaped queries on anomaly detection in recent years and have shown promising results. Accordingly, we explore diffusion models as dominant generative models in anomaly detection in the following section.

Wolleb et al. [16] introduce a weakly supervised learning method based on Denoising Diffusion Implicit Models (DDIMs) [157] for medical anomaly detection. Given an input image of a healthy or diseased subject, image-to-image translation first performs such that the objective is to translate the input image into the healthy one. Then, the anomaly regions are identified by subtracting the output image from the input. This process begins by encoding an input image into a noisy image with reversed DDIM sampling. Then, the denoising process is guided through a binary classifier trained beforehand on the healthy and diseased images to produce the healthy image. Finally, the anomaly map is calculated by taking the difference between the output and input. Results on BraTS2020 [81, 82, 83] and CheXpert [158] datasets demonstrate the superiority of the proposed approach compared to both VAE [159] and GAN [160] models.

Wyatt et al. [70] in AnoDDPM train a DDPM only on healthy medical samples. The anomaly image is then rendered by computing the difference between the output and input images. They also show that leveraging Simplex noise over Gaussian noise significantly enhances the performance.

In contrast, CDPM [71] demonstrates that training the diffusion probabilistic models only on healthy data generates poor segmentation performance. Thus, CDPM presents a counterfactual diffusion probabilistic model for generating healthy counterfactuals from factual input images. As illustrated in Figure 14, the input image is initially encoded into a latent space by iteratively applying diffusion models using an unconditional model. Then, the decoding step is accomplished by reversing the diffusion process. Using implicit guidance [161], the latent is decoded into a counterfactual by conditioning it on a healthy state and $\emptyset$. Inspired by [162, 163, 120], Sanchez et al. [71] then enhance the conditioning process by incorporating conditional attention into the U-Net backbone. As a final step, a dynamic normalization technique is applied during inference to avoid saturation in latent space pixels, caused by the guided iterative process that may change the image statistics. Eventually, the location of abnormality is determined by subtracting the input image from the generated healthy counterfactual.

Pinaya et al. [73] propose a fast DDPM-based approach for detecting and segmenting anomalous regions in the brain MR images (see Figure 15). This method follows the strategy of generating a healthy sample and delineating the anomaly segmentation map by subtracting it from the input image. To this end, VQ-VAE [164] is first adopted following [120],
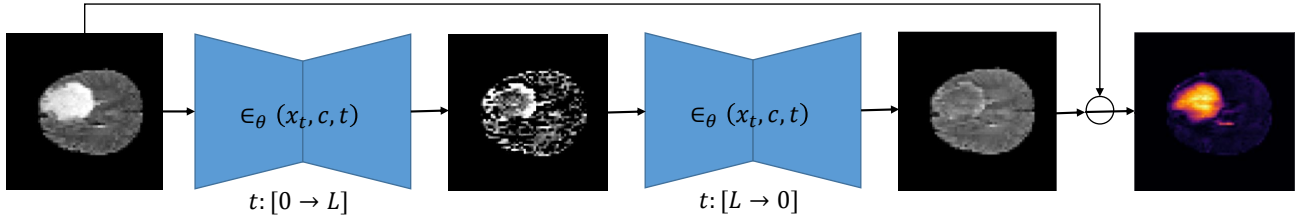
Figure 14: An overview of CDPM [71]. Iteratively applying diffusion models using an unconditional model $(c = \emptyset)$ encodes the input image into a latent space. Then, reversing the diffusion process from the latent space decodes a healthy state image. The decoding process is guided by conditioning it on the healthy state and $\emptyset$. The anomaly heatmap is generated by subtracting the input image from the generated counterfactual.

which encodes the input images into a compact latent representation and provides the quantized latent representation from input images utilizing a codebook. The DDPM then uses the acquired latent space and learns the distribution of the latent representation of the healthy samples. A binary mask indicating the location of the anomaly is constructed by applying a pre-calculated threshold on the average of intermediate samples of the reverse process, which contain less noisy and more distinct values. Using the middle step as the starting point for the reverse process, they denoise the anomalous areas of the image and preserve the rest using the obtained mask, thereby removing the lesion from the sample. Eventually, upon decoding the sample, a healthy image is produced.
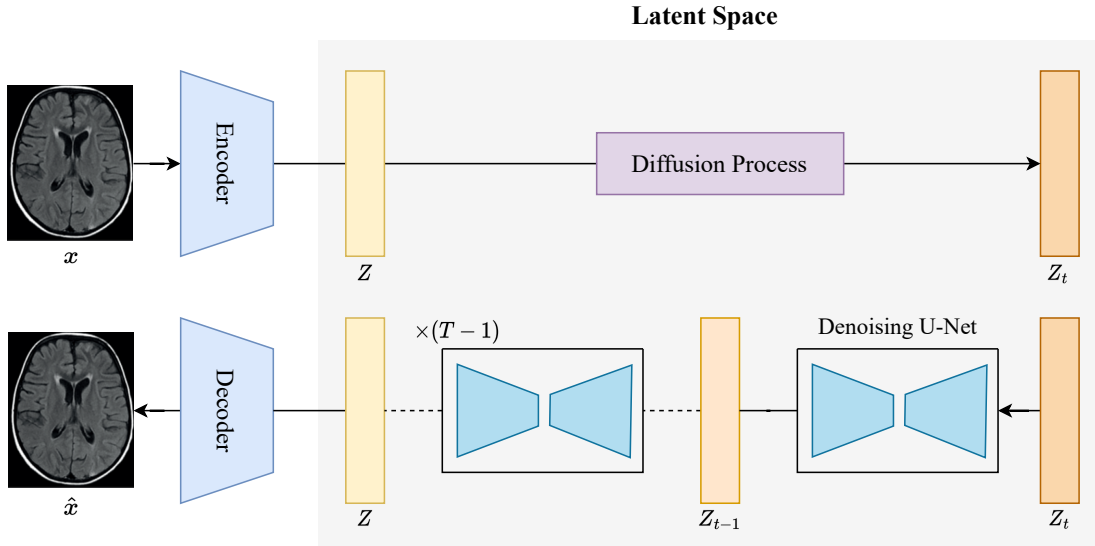


Figure 15: An overview of BAnoDDPM [73]. An autoencoder compresses the input image into a latent code, further enhanced by applying diffusion and reverse processes, and decodes into the pixel space.

In recent work by Behrendt et al. [165], the generation task of diffusion models is reframed as an estimation of healthy brain anatomy based on patches, utilizing spatial context to guide and enhance the reconstruction process. Specifically, they demonstrate that applying noise to the entire image simultaneously can pose challenges in achieving accurate reconstruction of the intricate structure of the brain. Therefore, patch-based Denoising Diffusion Probabilistic Models (pDDPMs) are introduced for Unsupervised Anomaly Detection (UAD) in brain MRI. In the proposed pDDPMs, they perform the forward diffusion process on a localized patch of the input image while employing the entire, partially noised image in the backward diffusion process to recover the noised patch. During inference, their trained pDDPM sequentially operates on sliding patches of the input image, alternately applying noise and denoising operations before stitching together the resultant denoised patches to yield the final denoised output image. Experiments on the public BraTS21 [166] and MSLUB [167] datasets verify that their approach performs superior or on par with most contemporary works in terms of UAD performance.
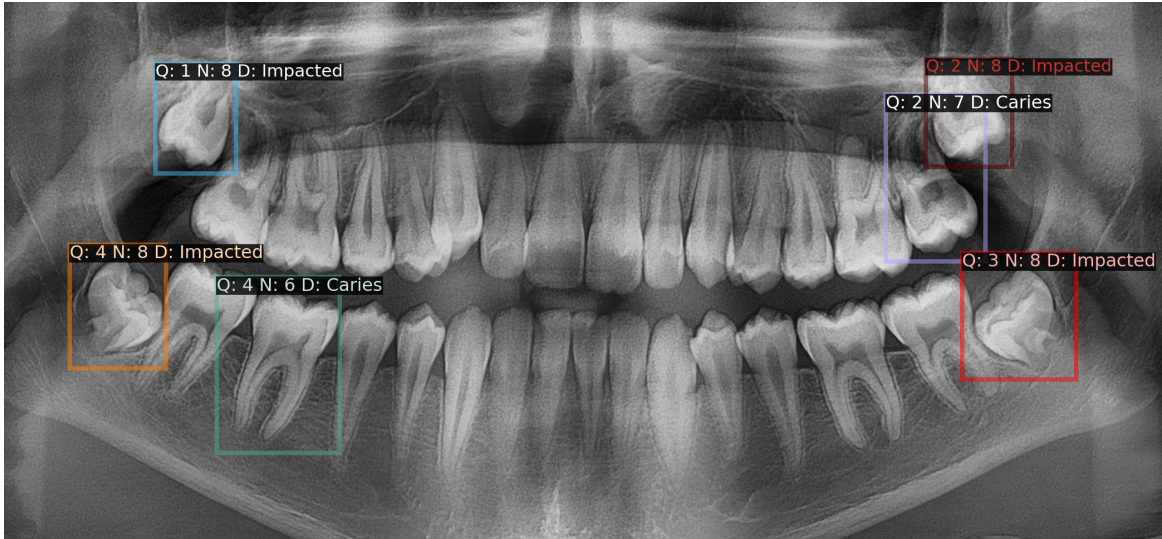
Figure 16: The final result of the HierarchicalDet [174] model displays bounding boxes around unhealthy teeth, along with the predicted quadrant (Q), enumeration (N), and diagnosis (D) labels.

## 4.9 Other Applications and Multiple Tasks

Based on Figure 5, there are still studies that could not be assigned to a particular category, and the use of diffusion is not limited to those nine categories. Gong et al. [168] present an innovative semi-supervised learning framework that utilizes diffusion models to accurately quantify the brain midline shift observed in head CT images. Keicher et al. [169] introduced a new method for grading vertebral fractures using a Diffusion Autoencoder (DAE) as an unsupervised, generative feature extractor. Also, diffusion models are not limited to vision-related tasks in the medical domain and can also foster research innovations in biology; e.g. the platforms presented in [75, 170] can be used for designing drugs and vaccines. In the following paragraphs, we explore some of the recent diffusion-based approaches used in multi-task learning and their distinctive uses in the medical domain.

Un/self-supervised learning is an ideal alternate approach in medical image denoising, where accessing paired clean and noisy images is difficult to achieve [171, 172]. Conventional deep-based networks utilize Minimum Mean Square Error (MMSE) estimates, which lead to unsatisfactory and blurred images due to the distribution change in train/test data or the preliminary assumption of Gaussian noise is at odds with the data's actual distribution. Chung et al. [76] proposed a multi-successive paradigm for MRI image denoising and super-resolution, namely R2D2+, with the SDE [36] algorithm to tackle the mentioned deficiencies. Diffusion generative models are robust to any distribution change over the data and produce more realistic data [11]. Despite the advantages of diffusion models, they are very time-consuming. To this end, Chung et al. [76] do not start the reverse diffusion process from the pure noise but start from the initial noisy image. R2D2+ [76] solves a reverse time SDE procedure with a non-parametric estimation method based on eigenvalue analysis of covariance matrix rather than the conventional numerical methods [36]. To restrain structure alleviation through the process, R2D2+ uses a low-frequency regularization to hamper any change in the low-frequency counterparts of the image. R2D2+ utilizes the same network for the super-resolution task after the denoising step. The overall results over the single coiled fastMRI [98] knee dataset and private liver MRI dataset indicate the superiority of this approach over conventional SOTA un/self-supervised learning schemes in terms of SNR and Contrast-to-Noise Ratio (CNR) metrics.

Accurate models for diagnosing skin cancer are crucial for early detection and treatment. Current computer-aided systems use deep learning, but recent research has shown that these models are highly vulnerable to attacks that subtly alter images, causing them to misclassify skin lesions. To address this problem, a new defense method [74] is proposed that can reverse these distortions by using a multiscale image pyramid and injecting Gaussian noise at each scale to neutralize the effects of adversarial perturbations. A denoising mechanism is then employed to remove added noise and aggregate information from neighboring scales. By repeating this process, images become resistant to noise and achieve actual probability value. The final step involves fusing sub-images at different scales to produce a reversed image. Experimental results on the ISIC 2019 dataset [173] demonstrate the superiority of the proposed method in defending against different attacks.

Hamamci et al. [174] propose a novel method for detecting abnormal teeth in panoramic X-rays using a hierarchical multi-label approach. They employ the DiffusionDet model [175], which utilizes a denoising diffusion process to predict objects and their categories from noisy boxes. In the first stage of their approach, the model is trained to predict quadrants

and their corresponding bounding box coordinates. The researchers input a raw image into the encoder to create high-level features, which are then used in the decoder's denoising step to refine the bounding boxes. The second stage involves predicting the tooth number, quadrant number, and bounding box. However, instead of refining the complete noisy boxes, the researchers propose manipulating bounding boxes by concatenating inferred boxes from the previous stage with the noisy boxes. The encoder and decoder weights are transferred from the previous stage to this stage. Finally, the third stage uses a similar approach to detect abnormal teeth with their quadrant-enumeration-diagnosis label. The output of the final model is illustrated in Figure 16. This hierarchical approach can handle partially labeled data and capture the full complexity of the underlying data. Additionally, the researchers present a new public dataset with three distinct data types: 1) for quadrant detection, 2) for tooth detection with both quadrant and tooth enumeration classifications, and 3) for diseased tooth detection with quadrant, tooth enumeration, and diagnosis classifications.

## 4.10 Comparative Overview

Table 1 comprehensively categorizes the reviewed diffusion model papers according to which algorithm they directly used or inspired: (1) DDPMs, (2) NCSNs, and (3) SDEs. In addition, Table 1 highlights the key concepts and objectives of each algorithm and represents the practical use cases that can be investigated and utilized in future research based on reviewed papers.

It is evident that conditioning the reverse diffusion process is one of the most studied methods for obtaining the desired output. This guiding process can be done using different constraint types. In [49, 53, 67], they control the reverse process by applying conditions using images; in particular, Lyu et al. [49] condition the DDPM and SDE utilizing T2w MR images to obtain CT images, Xie et al. [53] propose measurement-conditioned DDPM constituted for under-sampled medical image reconstruction, and Waibel et al. [67] constrain the 3D model using 2D microscopy images for generating 3D single cell shapes. Moreover, BrainGen [55] produces realistic examples of brain scans, which are conditioned on meta-data such as age, gender, ventricular volume, and brain volume relative to intracranial volume. In addition, the use of a classifier and implicit guidance methods in [16] and [71] have been investigated thoroughly. In this way, the distribution is shifted in a manner that is more likely to reach the expected outcome.

Some of the primary concerns and limitations of diffusion models are their slow speed and required computational cost. Several methods have been developed to address these drawbacks. Training-free Denoising Diffusion Implicit Model (DDIM) [157] is one of the advancements designed to accelerate the sampling process. DDIM extends the DDPM by substituting the Markovian process with the non-Markovian one, resulting in a faster sampling procedure with negligible quality degradation. [176, 177] propose a suitable initialization instead of a random Gaussian noise in the reverse process, causing a significant acceleration. Specifically, Chung et al. [177] prove that, based on stochastic contraction theory, beginning the reversion, for example, after one step prediction of the pre-trained neural network, fastens the reverse diffusion and reduces the number of reverse samplings. Dar et al. [55] also verify that adversarial learning can boost reverse diffusion speed by two orders of magnitude.

Several methods have also been worked on to enhance the output quality of diffusion models. AnoDDPM [70] substantiates that generalization to other types of noise distributions can enhance task-specific quality. They ascertain that in the case of anomaly detection, Simplex noise shows superiority over Gaussian noise. Furthermore, Cao et al. [176] corroborate that operating the diffusion process only in the high-frequency part of the image improves the stability and quality of the MRI reconstruction.

Despite the mentioned improvements in diffusion models, there remains a need to investigate why diffusion models have become popular in medical imaging and why some tasks are more successful in adopting diffusion models. Primarily, diffusion models have been increasingly used in medical imaging due to their effectiveness, ease of implementation, and high quality of output. In medical imaging, it is crucial to have high-resolution images that provide accurate local information for disease detection. Diffusion models have been able to achieve this, leading to their growing popularity in this field.

As seen in Figure 17 and Figure 5, some specific applications, such as image reconstruction, denoising, and generation, have received more attention compared to other tasks like segmentation, text-to-image translation, registration, etc. This is largely due to the compatibility of the diffusion model theory with the objectives of reconstruction and denoising tasks. The process in diffusion models involves adding noise to data and then denoising it until the original data is reconstructed, making it easier to implement these two tasks within the framework of diffusion models. Moreover, diffusion models have the ability to capture the underlying physical processes involved in these tasks and effectively model the complex interactions between signals and noise in data, resulting in more accurate image reconstruction. In addition, diffusion models are a class of generative models that operate on probabilistic distributions and can be conditioned to create synthetic data with a high degree of diversity and quality, which is why image generation has also been a popular application from the start.

While diffusion models have the potential to be applied to different tasks, they may require further modifications to be adapted to other specific tasks. For example, text-to-image translation requires an auxiliary network with strong text

encoding capabilities. In the case of object detection, recent vision-based work has demonstrated the potential of applying diffusion models to the object detection task by progressively refining randomly generated boxes to produce the final output results [175]. Hence, while initial works tend to focus on image generation, reconstruction, and denoising tasks, it is expected that, over time, more research will emerge addressing a wider range of tasks, as an examination of Figure 17 reveals the promising future of this field in the academic environment.

Table 1: Overview of the reviewed diffusion models in medical imaging based on their algorithm choice presented in our taxonomy, Figure 5. The symbol * indicates that the mentioned paper explores both DDPM and SDE algorithms.

| Algorithm | Networks | Core Ideas | Practical Use Cases |
|---|---|---|---|
| Denoising Diffusion Probabilistic Models (DDPMs) | [1]AnoDDPM [70]<br>[2]CDPM [71]<br>[3]IITM-Diffusion [72]<br>[4]AnoDDIM [16]<br>[5]PET-DDPM [64]<br>[6]DenoOCT-DDPM [65]<br>[7]brainSPADE [61]<br>[8]DARL [69]<br>[9]IISE [63]<br>[10]*MT-Diffusion [49]<br>[11]SynDiff [51]<br>[12]MC-DDPM [53]<br>[13]DiffuseRecon [54]<br>[14]AdaDiff [66]<br>[15]BrainGen [55]<br>[16]MFDPM [43]<br>[17]DISPR [67]<br>[18]3D-DDPM [68]<br>[19]DDM [69]<br>[20]Multi-scale DDAM [74]<br>[21]SMCDiff [75]<br>[22]BAnoDDPM [73]<br>[23]DiffuseMorph [59]<br>[24]20x-DenoDDPM [133]<br>[25]DOLCE [103]<br>[26]DiffMIC [60]<br>[27]CIMD [123]<br>[28]PatchDDM [129]<br>[29]DDM$^2$ [139]<br>[30]CoLa-Diff [151]<br>[31]pDDPMs [165]<br>[32]MLS-DDPM [168]<br>[33]X-ray LDM [145]<br>[34]DAE [169]<br>[35]HierarchicalDet [174] | In DDPMs [34], the forward diffusion process is represented as a Markov chain in which Gaussian noise is gradually added to the data. Data generation is then accomplished using the attained pure random noise and begins iterative denoising through a parametrized reverse process. Unlike VAE, where both the encoder and decoder are trained, only a single network is trained during the reverse process, and the forward process is considered fixed. An objective function of DDPMs is to simulate noise, which means that given a noisy input image, the neural network will produce the distribution modeled as normal distribution and indicate where the noise originates. | • Generalization to other types of noise distributions [70]<br>• Facilitating diffusion process via implicit guidance [71]<br>• Acceleration improvement using DDIM sampling [72, 169, 16]<br>• Guiding diffusion process via classifier guidance [16]<br>• Conditional DDPMs [64, 49, 53, 67, 133, 75]<br>• Generating synthetic segmentation datasets [61, 62]<br>• Cross-modality translation [51]<br>• Multi-modal conversion [49, 51]<br>• Exploiting K-space parameter-free guidance [54]<br>• Accelerate MC sampling using coarse-to-fine sampling [54]<br>• Adversarial learning in the reverse diffusion process [55, 51]<br>• 3D reconstruction from 2D images [67]<br>• Conditioning on medical meta-data [66]<br>• Using LDM to enhance the training and sampling efficiency [66, 151, 73, 145]<br>• DDPMs for histopathology images [43]<br>• 3D medical image generation [68]<br>• Using deformation fields for medical image generation [59, 69]<br>• DDPMs in skin image adversarial attacks [74]<br><br>• CT image reconstruction using limited-angle sinograms [103]<br>• Improving efficiency using a patch-based strategy [129, 165]<br>• Denoising diffusion MRI [139] |
| Noise Conditioned Score Networks (NCSNs) | [1]CSGM-MRI-Langevin [47]<br>[2]Self-Score [57] | In this algorithm [35], creating samples requires solving the Langevin dynamics equation. However, this equation mandates the solution of the gradient of the log density w.r.t. the input, $\nabla_x \log p(x)$, which is unknown and intractable. NCSN formulates the forward diffusion process by disturbing the data with Gaussian noise at different scales. Through this approach, this equation can be solved by training a single score network conditioned on the noise level and modeling the scores at all noise levels. Therefore, using the annealed Langevin dynamics algorithm and estimated score function at each scale, data can be generated. | • Using Langevin dynamics with different random initializations to get multiple reconstructions [47]<br>• Using conditioned Langevin Markov chain Monte Carlo (MCMC) sampling [57]<br>• Utilizing K-space data [57] |
| Stochastic Differential Equations (SDEs) | [1]*MT-Diffusion [49]<br>[2]UMM-CSGM [50]<br>[3]SIM-SGM [52]<br>[4]Score-MRI [48]<br>[5]MRI-DDMC [56]<br>[6]MCG [58]<br>[7]HFS-SDE [176]<br>[8]Diffuse-Faster [177]<br>[9]HKGM [178]<br>[10]WKGM [179]<br>[11]R2D2+ [76] | As in the cases of DDPMs and NCSNs, SDEs [36] follow a similar approach in the forward diffusion process, in which the input data is perturbed consecutively utilizing Gaussian noise. In contrast, previous probabilistic models can be viewed as discretizations of score-based models by extending the number of noise scales to infinity and treating them as continuous.<br>Having operated backward in time, the reverted SDE Process solves the reverse-time Stochastic Differential Equation and recovers the data. Even so, each time step of this process requires the score function of the density. SDEs are, therefore, intended to learn the actual score function via a neural network and construct samples using numerical SDE solvers.<br>Any numerical method can be adopted to solve the reverse-time SDE equation. In particular, three notable sampling methods are named and described in Section 2.3.2: Euler-Maruyama (EM) method, Prediction-Correction (PC), and Probability Flow ODE. EM follows a simple strategy, PC generates more high-fidelity samples, and Probability Flow ODE is fast and efficient. | • SDEs achieve better results in multi-modal conversion [49]<br>• Cross-modality translation [50]<br>• Conditional SDEs [50, 54, 58]<br>• Solving linear inverse problems [52]<br>• Using K-space data [178, 179]<br>• Using manifold constraints [58]<br>• Good initialization in the reverse process instead of random Gaussian noise leads to a faster convergence [177, 176]<br>• Diffusion only in high-frequency space increases the stability and quality [176] |

# 5  Future Direction and Open Challenges

Diffusion models have emerged as a popular topic in the medical vision and medical-biology fields, as evidenced by the upward trend shown in Figure 1. One of the primary advantages of diffusion models in medical imaging is that they do not require labeled data, making them a strong candidate for many medical applications. In addition, Diffusion models have gained popularity due to their impressive performance in foundational models like large text-to-image models [162]. These models remain appealing for modeling images and other data distributions for several reasons, including their ability to
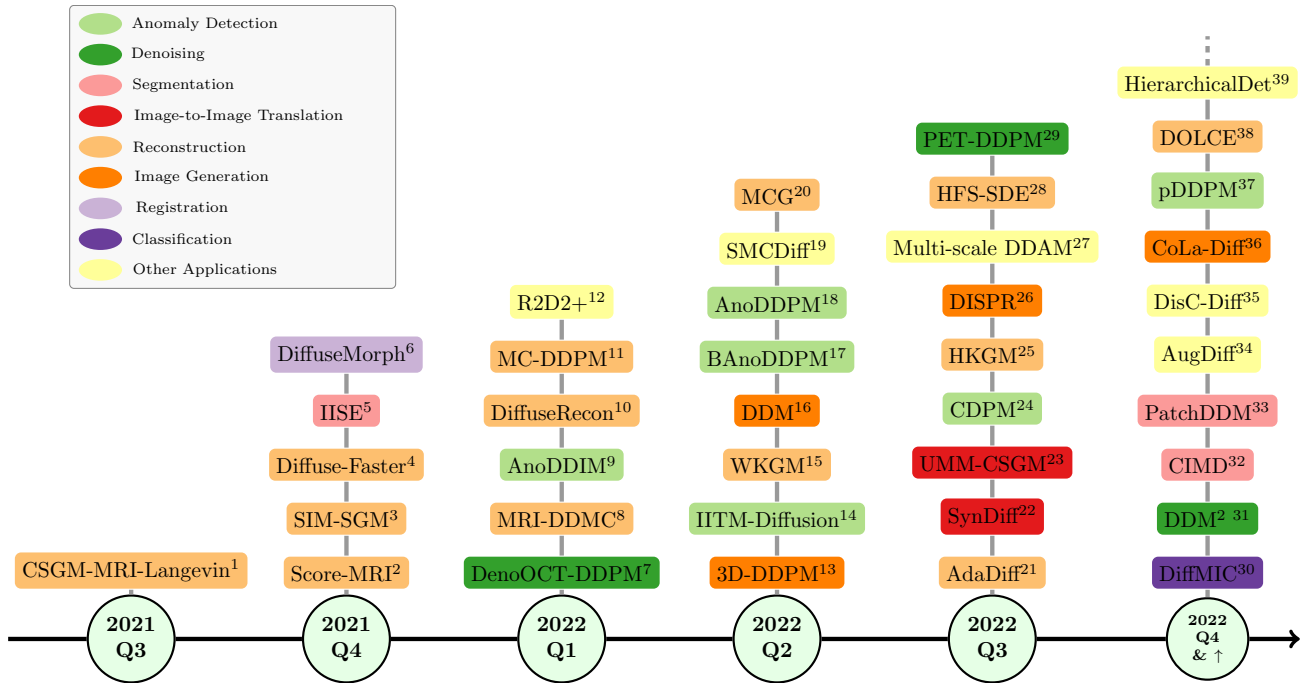
Figure 17: Diffusion models timeline from 2021 to April 2023 through the first paper in the medical field. The superscripts in ascending order represent 1.[47], 2.[48], 3.[52], 4.[177], 5.[63], 6.[59], 7.[65], 8.[56], 9.[16], 10.[54], 11.[53], 12.[76], 13.[68], 14.[72], 15.[179], 16.[69], 17.[73], 18.[70], 19.[75], 20.[58], 21.[55], 22.[51], 23.[50], 24.[71], 25.[178], 26.[67], 27.[74], 28.[176], 29.[64], 30.[60], 31.[139], 32.[123], 33.[129], 34.[180], 35.[181], 36.[151], 37.[165], 38.[103], 39.[174], respectively.

effectively represent high-dimensional data like images due to their inductive bias. These models are trained on a reweighted variational lower bound that emphasizes the global consistency and dominant patterns of images while giving less attention to less noticeable details, making them an excellent inductive bias for spatial data. However, Diffusion models, compared to other generative models, have some limitations that should be taken into account when considering their applications. These limitations include a slower generation process compared to some other generative models, limited applicability to certain data types (e.g., audio, text, or structured data), lower likelihood compared to other models, and an inability to perform dimensionality reduction [20]. However, these limitations do not diminish the unique strengths and advantages of diffusion models in generating high-quality images and their ability to work without a pair of labeled or unlabeled data but rather present open challenges for further research and improvements.

This paper aims to provide a comprehensive review of the latest medical research papers utilizing diffusion models. We categorized the studies based on the taxonomy proposed in Figure 5 to demonstrate the potential of diffusion models. Through this review, we hope to highlight the power of diffusion models and shed light on their importance in advancing the capabilities of medical imaging techniques. This section identifies areas for future investigation, emphasizing the need for continued research in this exciting and rapidly evolving field.

**Exploring more diverse medical imaging modalities:** Due to the nature of diffusion models, they are a strong candidate for exploring diverse modalities for distinct downstream tasks. According to Figure 1b, most of the published studies utilize CT and MRI as modalities. Nevertheless, other modalities may also benefit from the capabilities of diffusion models. For instance, ultrasound imaging techniques may suffer from a non-ideal Point Spread Function (PSF) of the imaging system as well as intrinsic physical limitations [182]. To this end, several studies explored the impact of various generative pipelines on ultrasound data in image quality enhancement and denoising [182, 183, 184], resulting in improved image quality.

**Representation space:** VAEs and GANs are designed to preserve and learn meaningful representations of data in their latent space. Nevertheless, diffusion models have been proven to be less successful in creating semantically meaningful data representations in their latent space [21, 185]. Hence, the lack of semantically expressive data representations in the latent space of diffusion models poses a significant hindrance in performing tasks that involve the manipulation of data based on semantic representations. This is probably because diffusion models mainly destruct information in the latent variables during the diffusion process, resulting in less meaningful representation space. Therefore, there is a need to develop models that can learn semantically meaningful representations, as they enable better image reconstructions and semantic

interpolations [185]. For example, Abstreiter et al. [186] have proposed a novel diffusion-based representation learning approach based on conditional denoising score matching. Specifically, they introduce an additional trainable encoder and condition the score estimator on the encoder output, which is the latent representation of the clean data. This produces interpretable features in the latent space and enables altering the encoded features' granularity without requiring architecture modifications or data manipulations. Therefore, the lack of proper representation space of diffusion models presents an open challenge for researchers to work on.

**Architecture design:** In the context of diffusion models, the network structure is a critical design choice that directly affects their ability to learn complex data relationships and produce high-quality results on large and diverse datasets [11]. Most diffusion models currently utilize CNN-based architectures with a global attention layer, but recent studies have explored the use of transformer models [187]. Compared to CNNs, transformers offer several advantages, including the ability to model non-local interactions and capture long-range dependencies within the data. These models have also shown promising results in natural language processing tasks, indicating their potential for modeling sequential or spatiotemporal data. Despite this potential, the use of transformers in diffusion models is still in its early stages, and further research is necessary to fully understand their capabilities and limitations in this context. Specifically, most DDPM-based approaches follow the baseline of [34, 11, 188], while Score-based approaches follow the baseline of [36, 41]. As a result, there is a lack of research focused on improving the architecture of diffusion models for medical imaging, making it an open challenge for future research.

**Causal discovery, inference, and counterfactual generation of diffusion models:** Diffusion models can generally learn the underlying probability distribution of a dataset and be used to generate new data points that follow the same distribution. This makes them even more useful for complex causal inference, discovery, and counterfactual generation tasks. The benefit of using diffusion models for causal inference is their ability to handle missing data and their robustness to distributional shifts, allowing them to estimate the causal effects of interventions in real-world settings where data may be incomplete. For example, Sanchez et al. [189] propose a novel framework for estimating causal effects involving high-dimensional variables using diffusion models. In addition, causal discovery seeks to identify the underlying causal structure of a system without necessarily being given a specific intervention. Sanchez et al. [190] proposed a novel method using diffusion models for causal discovery via topological ordering based on diffusion models. Additionally, diffusion models can be used to generate counterfactuals, which explore hypothetical scenarios that assess the impact of an intervention that was or was not made. This is particularly useful in fields such as medicine and public policy, where conducting randomized controlled trials may be difficult or unethical. Overall, diffusion models provide a powerful tool for generating data and conducting causal inference, discovery, and counterfactual analysis in various fields.

**Privacy concerns:** The medical community is highly concerned about the privacy of medical data. AI image synthesis models are currently under scrutiny for potentially violating copyright laws and compromising the privacy of their training data. Carlini et al. [191] conducted extensive experiments to evaluate the privacy concerns associated with generative diffusion models. The results show that diffusion models tend to memorize individual images from their training data and reproduce them during generation. This would consequently enable adversaries to attack to extract training data. Furthermore, the study reveals that diffusion models are much less private compared to other generative models like GANs. As a result, new developments in privacy-preserving training are needed to handle these vulnerabilities, particularly in delicate fields like medicine.

**Federated learning and diffusion models:** Due to privacy concerns in medical imaging, which limit data integration, diffusion models and federated learning can create a profound and robust learning platform in the medical domain. Data is collected from various sources in privacy preservation smart healthcare systems and stored in decentralized locations [192]. Federated Learning, as it allows for training machine learning models on decentralized data without exposing sensitive information in conjunction with diffusion models, can capture the underlying distribution of the data across multiple participating devices. This intuition, using the diffusion models as a "generative prior," can help to mitigate the effects of data heterogeneity, reduce the risk of privacy leaks, and improve the quality of the learned models and their trustworthiness in fairness generalization. Furthermore, generative models (i.e., diffusion models) are trained to learn the underlying probability distribution of the data rather than memorizing the training data. The stability against perturbations is a desirable property of machine learning models, especially in safety-critical applications, such as medical diagnosis, where small changes in the input data can have significant consequences. Diffusion models can provide stability against perturbations by using regularization techniques, such as adding noise to the input data during training, which can help the model generalize better to unseen data [47, 48]. Excessive to this, federated learning models perform on a broadened range of sources with various data distributions. Therefore, a federated learning paradigm could achieve strong out-of-distribution (OOD) generalization capacity [57, 48]. In addition, due to the same privacy concerns, diffusion models can consolidate their steps in generating synthetic medical data for educational purposes.

**Reverse process using reinforcement learning:** The inverse problem-solving of the diffusion models could be performed by the reinforcement learning paradigm to estimate the best inversion path rather than solid mathematical solutions. In this process, reinforcement learning can be used to search for the optimal values of the diffusion model parameters that

maximize a given reward function. The reward function can be designed to evaluate how well the diffusion model fits the data and penalize deviations from the observed data where the traditional optimization methods, such as maximum likelihood estimation or Bayesian inference, become computationally expensive or intractable.

# 6    Conclusion

In this paper, we provided a survey of the literature on diffusion models with a focus on applications in the medical imaging field. Specifically, we investigated the applications of diffusion models in anomaly detection, medical image segmentation, denoising, classification, reconstruction, registration, generation, and other tasks. In particular, for each of these applications, we provided a taxonomy and high-level abstraction of the core techniques from various angles. Moreover, we characterized the existing models based on techniques where we identified three primary formulations of diffusion modeling based on: DDPMs, NCSNs, and SDEs. Finally, we outlined possible avenues for future research.

While our survey highlights the rapid growth of diffusion-based techniques in medical imaging, we also acknowledge that the field is still in its early stages and subject to change. As diffusion models continue to gain popularity and more research is conducted in this field, our survey serves as an important starting point and reference for researchers and practitioners looking to utilize these models in their work. We hope that this survey will inspire further interest and exploration of the potential of diffusion models in the medical domain. It is important to note that some of the papers cited in this survey are pre-prints. However, we made every effort to include only high-quality research from reputable sources, and we believe that the inclusion of pre-prints provides a comprehensive overview of the current state-of-the-art in this rapidly evolving field. Overall, we believe that our survey provides valuable insights into the use of diffusion models in medical imaging and highlights promising areas for future research.

# References

[1] Jianmin Bao, Dong Chen, Fang Wen, Houqiang Li, and Gang Hua. CVAE-GAN: fine-grained image generation through asymmetric training. In *Proceedings of the IEEE international conference on computer vision*, pages 2745–2754, 2017.

[2] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32, 2019.

[3] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. DiffWave: A versatile diffusion model for audio synthesis. In *International Conference on Learning Representations*, 2021.

[4] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. WaveNet: A generative model for raw audio. In *9th ISCA Speech Synthesis Workshop*, pages 125–125, 2016.

[5] Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori Hashimoto. Diffusion-LM improves controllable text generation. In *Advances in Neural Information Processing Systems*, 2022.

[6] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. PointFlow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4541–4550, 2019.

[7] Sam Bond-Taylor, Adam Leach, Yang Long, and Chris G Willcocks. Deep generative modelling: A comparative review of VAEs, GANs, normalizing flows, energy-based and autoregressive models. *IEEE transactions on pattern analysis and machine intelligence*, 2021.

[8] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.

[9] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*, pages 1278–1286. PMLR, 2014.

[10] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In *International Conference on Learning Representations*, 2017.

[11] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.

[12] Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yueting Chen. SRDiff: Single image super-resolution with diffusion probabilistic models. *Neurocomputing*, 479:47–59, 2022.

[13] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022.

[14] Tomer Amit, Eliya Nachmani, Tal Shaharbany, and Lior Wolf. SegDiff: Image segmentation with diffusion probabilistic models. *arXiv preprint arXiv:2112.00390*, 2021.

[15] Roland S. Zimmermann, Lukas Schott, Yang Song, Benjamin Adric Dunn, and David A. Klindt. Score-based generative classifiers. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.

[16] Julia Wolleb, Florentin Bieder, Robin Sandkühler, and Philippe C Cattin. Diffusion models for medical anomaly detection. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part VIII*, pages 35–45. Springer, 2022.

[17] Manal AlAmir and Manal AlGhamdi. The role of generative adversarial network in medical image analysis: an in-depth survey. *ACM Computing Surveys (CSUR)*, 2022.

[18] Hazrat Ali, Rafiul Biswas, Farida Ali, Uzair Shah, Asma Alamgir, Osama Mousa, and Zubair Shah. The role of generative adversarial networks in brain MRI: a scoping review. *Insights into Imaging*, 13(1):1–15, 2022.

[19] Yizhou Chen, Xu-Hua Yang, Zihan Wei, Ali Asghar Heidari, Nenggan Zheng, Zhicheng Li, Huiling Chen, Haigen Hu, Qianwei Zhou, and Qiu Guan. Generative adversarial networks in medical image augmentation: a review. *Computers in Biology and Medicine*, page 105382, 2022.

[20] Hanqun Cao, Cheng Tan, Zhangyang Gao, Guangyong Chen, Pheng-Ann Heng, and Stan Z Li. A survey on generative diffusion model. *arXiv preprint arXiv:2209.02646*, 2022.

[21] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Yingxia Shao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *arXiv preprint arXiv:2209.00796*, 2022.

[22] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in vision: A survey. *arXiv preprint arXiv:2209.04747*, 2022.

[23] Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, and F Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.

[24] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[25] George Papamakarios, Eric T Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *J. Mach. Learn. Res.*, 22(57):1–64, 2021.

[26] Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the generative learning trilemma with denoising diffusion GANs. In *International Conference on Learning Representations*, 2022.

[27] Maciej Wiatrak, Stefano V Albrecht, and Andrew Nystrom. Stabilizing generative adversarial networks: A survey. *arXiv preprint arXiv:1910.00927*, 2019.

[28] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *International Conference on Learning Representations*, 2018.

[29] Tanya Motwani and Manojkumar Parmar. A novel framework for selection of GANs for an application. *arXiv preprint arXiv:2002.08641*, 2020.

[30] Tim R Davidson, Luca Falorsi, Nicola De Cao, Thomas Kipf, and Jakub M Tomczak. Hyperspherical variational auto-encoders. *arXiv preprint arXiv:1804.00891*, 2018.

[31] Andrea Asperti. Variational autoencoders and the variable collapse phenomenon. *Sensors & Transducers*, 234(6):1–8, 2019.

[32] Lilian Weng. Flow-based deep generative models. *lilianweng.github.io*, 2018.

[33] 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.

[34] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.

[35] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019.

[36] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.

[37] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.

[38] Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach to density and score estimation. In *Uncertainty in Artificial Intelligence*, pages 574–584. PMLR, 2020.

[39] Giorgio Parisi. Correlation functions and computer simulations. *Nuclear Physics B*, 180(3):378–384, 1981.

[40] Ulf Grenander and Michael I Miller. Representations of knowledge in complex systems. *Journal of the Royal Statistical Society: Series B (Methodological)*, 56(4):549–581, 1994.

[41] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020.

[42] Yanbin Liu, Girish Dwivedi, Farid Boussaid, and Mohammed Bennamoun. 3d brain and heart volume generative models: A survey. *arXiv preprint arXiv:2210.05952*, 2022.

[43] Puria Azadi Moghadam, Sanne Van Dalen, Karina C Martin, Jochen Lennerz, Stephen Yip, Hossein Farahani, and Ali Bashashati. A morphology focused diffusion probabilistic model for synthesis of histopathology images. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2000–2009, 2023.

[44] Andre Goncalves, Priyadip Ray, Braden Soper, Jennifer Stevens, Linda Coyle, and Ana Paula Sales. Generation and evaluation of synthetic patient data. *BMC medical research methodology*, 20(1):1–40, 2020.

[45] Richard J Chen, Ming Y Lu, Tiffany Y Chen, Drew FK Williamson, and Faisal Mahmood. Synthetic data in machine learning for medicine and healthcare. *Nature Biomedical Engineering*, 5(6):493–497, 2021.

[46] Mohamed Akrout, Bálint Gyepesi, Péter Holló, Adrienn Poór, Blága Kincső, Stephen Solis, Katrina Cirone, Jeremy Kawahara, Dekker Slade, Latif Abid, et al. Diffusion-based data augmentation for skin disease classification: Impact across original medical datasets to fully synthetic images. *arXiv preprint arXiv:2301.04802*, 2023.

[47] Ajil Jalal, Marius Arvinte, Giannis Daras, Eric Price, Alexandros G Dimakis, and Jon Tamir. Robust compressed sensing MRI with deep generative priors. *Advances in Neural Information Processing Systems*, 34:14938–14954, 2021.

[48] Hyungjin Chung and Jong Chul Ye. Score-based diffusion models for accelerated MRI. *Medical Image Analysis*, page 102479, 2022.

[49] Qing Lyu and Ge Wang. Conversion between CT and MRI images using diffusion and score-matching models. *arXiv preprint arXiv:2209.12104*, 2022.

[50] Xiangxi Meng, Yuning Gu, Yongsheng Pan, Nizhuan Wang, Peng Xue, Mengkang Lu, Xuming He, Yiqiang Zhan, and Dinggang Shen. A novel unified conditional score-based generative framework for multi-modal medical image completion. *arXiv preprint arXiv:2207.03430*, 2022.

[51] Muzaffer Özbey, Salman UH Dar, Hasan A Bedel, Onat Dalmaz, Şaban Öztürk, Alper Güngör, and Tolga Çukur. Unsupervised medical image translation with adversarial diffusion models. *arXiv preprint arXiv:2207.08208*, 2022.

[52] Yang Song, Liyue Shen, Lei Xing, and Stefano Ermon. Solving inverse problems in medical imaging with score-based generative models. In *International Conference on Learning Representations*, 2022.

[53] Yutong Xie and Quanzheng Li. Measurement-conditioned denoising diffusion probabilistic model for under-sampled medical image reconstruction. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part VI*, pages 655–664. Springer, 2022.

[54] Cheng Peng, Pengfei Guo, S Kevin Zhou, Vishal M Patel, and Rama Chellappa. Towards performant and reliable undersampled MR reconstruction via diffusion model sampling. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 623–633. Springer, 2022.

[55] Salman UH Dar, Şaban Öztürk, Yilmaz Korkmaz, Gokberk Elmas, Muzaffer Özbey, Alper Güngör, and Tolga Çukur. Adaptive diffusion priors for accelerated MRI reconstruction. *arXiv preprint arXiv:2207.05876*, 2022.

[56] Guanxiong Luo, Martin Heide, and Martin Uecker. MRI reconstruction via data driven markov chain with joint uncertainty estimation. *arXiv preprint arXiv:2202.01479*, 2022.

[57] Zhuo-Xu Cui, Chentao Cao, Shaonan Liu, Qingyong Zhu, Jing Cheng, Haifeng Wang, Yanjie Zhu, and Dong Liang. Self-score: Self-supervised learning on score-based models for MRI reconstruction. *arXiv preprint arXiv:2209.00835*, 2022.

[58] Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Improving diffusion models for inverse problems using manifold constraints. In *Advances in Neural Information Processing Systems*, 2022.

[59] Boah Kim, Inhwa Han, and Jong Chul Ye. DiffuseMorph: Unsupervised deformable image registration using diffusion model. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXI*, pages 347–364. Springer, 2022.

[60] Yijun Yang, Huazhu Fu, Angelica Aviles-Rivero, Carola-Bibiane Schönlieb, and Lei Zhu. Diffmic: Dual-guidance diffusion network for medical image classification. *arXiv preprint arXiv:2303.10610*, 2023.

[61] Virginia Fernandez, Walter Hugo Lopez Pinaya, Pedro Borges, Petru-Daniel Tudosiu, Mark S Graham, Tom Vercauteren, and M Jorge Cardoso. Can segmentation models be trained with fully synthetically generated data? In *International Workshop on Simulation and Synthesis in Medical Imaging*, pages 79–90. Springer, 2022.

[62] Boah Kim, Yujin Oh, and Jong Chul Ye. Diffusion adversarial representation learning for self-supervised vessel segmentation. In *The Eleventh International Conference on Learning Representations*, 2023.

[63] Julia Wolleb, Robin Sandkühler, Florentin Bieder, Philippe Valmaggia, and Philippe C Cattin. Diffusion models for implicit image segmentation ensembles. In *International Conference on Medical Imaging with Deep Learning*, pages 1336–1348. PMLR, 2022.

[64] Kuang Gong, Keith A Johnson, Georges El Fakhri, Quanzheng Li, and Tinsu Pan. PET image denoising based on denoising diffusion probabilistic models. *arXiv preprint arXiv:2209.06167*, 2022.

[65] Dewei Hu, Yuankai K Tao, and Ipek Oguz. Unsupervised denoising of retinal OCT with diffusion probabilistic model. In *Medical Imaging 2022: Image Processing*, volume 12032, pages 25–34. SPIE, 2022.

[66] Walter HL Pinaya, Petru-Daniel Tudosiu, Jessica Dafflon, Pedro F Da Costa, Virginia Fernandez, Parashkev Nachev, Sebastien Ourselin, and M Jorge Cardoso. Brain imaging generation with latent diffusion models. In *MICCAI Workshop on Deep Generative Models*, pages 117–126. Springer, 2022.

[67] Dominik JE Waibel, Ernst Röoell, Bastian Rieck, Raja Giryes, and Carsten Marr. A diffusion model predicts 3D shapes from 2D microscopy images. *arXiv preprint arXiv:2208.14125*, 2022.

[68] Zolnamar Dorjsembe, Sodtavilan Odonchimed, and Furen Xiao. Three-dimensional medical image synthesis with denoising diffusion probabilistic models. In *Medical Imaging with Deep Learning*, 2022.

[69] Boah Kim and Jong Chul Ye. Diffusion deformable model for 4d temporal medical image generation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 539–548. Springer, 2022.

[70] Julian Wyatt, Adam Leach, Sebastian M Schmon, and Chris G Willcocks. AnoDDPM: Anomaly detection with denoising diffusion probabilistic models using simplex noise. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 650–656, 2022.

[71] Pedro Sanchez, Antanas Kascenas, Xiao Liu, Alison Q O'Neil, and Sotirios A Tsaftaris. What is healthy? generative counterfactual diffusion for lesion localization. In *Deep Generative Models: Second MICCAI Workshop, DGM4MICCAI 2022, Held in Conjunction with MICCAI 2022, Singapore, September 22, 2022, Proceedings*, pages 34–44. Springer, 2022.

[72] Julia Wolleb, Robin Sandkühler, Florentin Bieder, and Philippe C Cattin. The swiss army knife for image-to-image translation: Multi-task diffusion models. *arXiv preprint arXiv:2204.02641*, 2022.

[73] Walter HL Pinaya, Mark S Graham, Robert Gray, Pedro F Da Costa, Petru-Daniel Tudosiu, Paul Wright, Yee H Mah, Andrew D MacKinnon, James T Teo, Rolf Jager, et al. Fast unsupervised brain anomaly detection and segmentation with diffusion models. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part VIII*, pages 705–714. Springer, 2022.

[74] Yongwei Wang, Yuan Li, and Zhiqi Shen. Fight fire with fire: Reversing skin adversarial examples by multiscale diffusive and denoising aggregation mechanism. *arXiv preprint arXiv:2208.10373*, 2022.

[75] Brian L. Trippe, Jason Yim, Doug Tischer, David Baker, Tamara Broderick, Regina Barzilay, and Tommi S. Jaakkola. Diffusion probabilistic modeling of protein backbones in 3d for the motif-scaffolding problem. In *International Conference on Learning Representations*, 2023.

[76] Hyungjin Chung, Eun Sun Lee, and Jong Chul Ye. MR image denoising and super-resolution using regularized reverse diffusion. *IEEE Transactions on Medical Imaging*, 2022.

[77] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.

[78] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein GANs. *Advances in neural information processing systems*, 30, 2017.

[79] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.

[80] 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.

[81] Spyridon Bakas, Hamed Akbari, Aristeidis Sotiras, Michel Bilello, Martin Rozycki, Justin S Kirby, John B Freymann, Keyvan Farahani, and Christos Davatzikos. Advancing the cancer genome atlas glioma MRI collections with expert segmentation labels and radiomic features. *Scientific data*, 4(1):1–13, 2017.

[82] Spyridon Bakas, Mauricio Reyes, Andras Jakab, Stefan Bauer, Markus Rempfler, Alessandro Crimi, Russell Takeshi Shinohara, Christoph Berger, Sung Min Ha, Martin Rozycki, et al. Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the BRATS challenge. *arXiv preprint arXiv:1811.02629*, 2018.

[83] Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree Kalpathy-Cramer, Keyvan Farahani, Justin Kirby, Yuliya Burren, Nicole Porz, Johannes Slotboom, Roland Wiest, et al. The multimodal brain tumor image segmentation benchmark (BRATS). *IEEE transactions on medical imaging*, 34(10):1993–2024, 2014.

[84] Pu Huang, Dengwang Li, Zhicheng Jiao, Dongming Wei, Guoshi Li, Qian Wang, Han Zhang, and Dinggang Shen. CoCa-GAN: common-feature-learning-based context-aware generative adversarial network for glioma grading. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 155–163. Springer, 2019.

[85] Agisilaos Chartsias, Thomas Joyce, Mario Valerio Giuffrida, and Sotirios A Tsaftaris. Multimodal MR synthesis via modality-invariant latent representation. *IEEE transactions on medical imaging*, 37(3):803–814, 2017.

[86] Xiaofeng Liu, Fangxu Xing, Georges El Fakhri, and Jonghye Woo. A unified conditional disentanglement framework for multimodal brain MR image translation. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 10–14. IEEE, 2021.

[87] 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.

[88] 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.

[89] Yunxiang Li, Hua-Chieh Shao, Xiao Liang, Liyuan Chen, Ruiqi Li, Steve Jiang, Jing Wang, and You Zhang. Zero-shot medical image translation via frequency-guided diffusion models. *arXiv preprint arXiv:2304.02742*, 2023.

[90] Brett Levac, Ajil Jalal, and Jonathan I Tamir. Accelerated motion correction for mri using score-based generative models. *arXiv preprint arXiv:2211.00199*, 2022.

[91] Chentao Cao, Zhuo-Xu Cui, Jing Cheng, Sen Jia, Hairong Zheng, Dong Liang, and Yanjie Zhu. Spirit-diffusion: Spirit-driven score-based generative modeling for vessel wall imaging. *arXiv preprint arXiv:2212.11274*, 2022.

[92] Chang Min Hyun, Hwa Pyung Kim, Sung Min Lee, Sungchul Lee, and Jin Keun Seo. Deep learning for undersampled MRI reconstruction. *Physics in Medicine & Biology*, 63(13):135007, 2018.

[93] Yilmaz Korkmaz, Salman UH Dar, Mahmut Yurt, Muzaffer Özbey, and Tolga Cukur. Unsupervised MRI reconstruction via zero-shot learned adversarial transformers. *IEEE Transactions on Medical Imaging*, 2022.

[94] Chun-Mei Feng, Yunlu Yan, Huazhu Fu, Li Chen, and Yong Xu. Task transformer network for joint MRI reconstruction and super-resolution. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 307–317. Springer, 2021.

[95] Yilmaz Korkmaz, Mahmut Yurt, Salman Ul Hassan Dar, Muzaffer Özbey, and Tolga Cukur. Deep MRI reconstruction with generative vision transformers. In *International Workshop on Machine Learning for Medical Image Reconstruction*, pages 54–64. Springer, 2021.

[96] Ritu Gothwal, Shailendra Tiwari, and Shivendra Shivani. Computational medical image reconstruction techniques: A comprehensive review. *Archives of Computational Methods in Engineering*, pages 1–28, 2022.

[97] Dominique Bakry and Michel Émery. Diffusions hypercontractives. In *Seminaire de probabilités XIX 1983/84*, pages 177–206. Springer, 1985.

[98] Jure Zbontar, Florian Knoll, Anuroop Sriram, Tullie Murrell, Zhengnan Huang, Matthew J Muckley, Aaron Defazio, Ruben Stern, Patricia Johnson, Mary Bruno, et al. fastMRI: An open dataset and benchmarks for accelerated MRI. *arXiv preprint arXiv:1811.08839*, 2018.

[99] Stanford University. Stanford MRI. http://mridata.org/. Online; accessed 6 October 2022.

[100] Kai Tobias Block, Martin Uecker, and Jens Frahm. Undersampled radial MRI with multiple coils. iterative image reconstruction using a total variation constraint. *Magnetic Resonance in Medicine: An Official Journal of the International Society for Magnetic Resonance in Medicine*, 57(6):1086–1098, 2007.

[101] Bo Zhou and S Kevin Zhou. DuDoRNet: learning a dual-domain recurrent network for fast MRI reconstruction with deep t1 prior. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4273–4282, 2020.

[102] Peter B Roemer, William A Edelstein, Cecil E Hayes, Steven P Souza, and Otward M Mueller. The NMR phased array. *Magnetic resonance in medicine*, 16(2):192–225, 1990.

[103] Jiaming Liu, Rushil Anirudh, Jayaraman J Thiagarajan, Stewart He, K Aditya Mohan, Ulugbek S Kamilov, and Hyojin Kim. DOLCE: A model-based probabilistic diffusion framework for limited-angle CT reconstruction. *arXiv preprint arXiv:2211.12340*, 2022.

[104] Avinash C Kak and Malcolm Slaney. *Principles of computerized tomographic imaging*. SIAM, 2001.

[105] Nicholas Heller, Fabian Isensee, Klaus H Maier-Hein, Xiaoshuai Hou, Chunmei Xie, Fengyi Li, Yang Nan, Guangrui Mu, Zhiyong Lin, Miofei Han, et al. The state of the art in kidney and kidney tumor segmentation in contrast-enhanced CT imaging: Results of the KiTS19 challenge. *Medical Image Analysis*, page 101821, 2020.

[106] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. *Advances in neural information processing systems*, 28, 2015.

[107] Oliver Langner, Ron Dotsch, Gijsbert Bijlstra, Daniel HJ Wigboldus, Skyler T Hawk, and AD Van Knippenberg. Presentation and validation of the radboud faces database. *Cognition and emotion*, 24(8):1377–1388, 2010.

[108] Pamela J LaMontagne, Tammie LS Benzinger, John C Morris, Sarah Keefe, Russ Hornbeck, Chengjie Xiong, Elizabeth Grant, Jason Hassenstab, Krista Moulder, Andrei G Vlassenko, et al. OASIS-3: longitudinal neuroimaging, clinical, and cognitive dataset for normal aging and alzheimer disease. *MedRxiv*, pages 2019–12, 2019.

[109] 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.

[110] Mohammad Shehab, Laith Abualigah, Qusai Shambour, Muhannad A Abu-Hashem, Mohd Khaled Yousef Shambour, Ahmed Izzat Alsalibi, and Amir H Gandomi. Machine learning in medical applications: A review of state-of-the-art methods. *Computers in Biology and Medicine*, 145:105458, 2022.

[111] Arthur Gretton, Karsten Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex Smola. A kernel method for the two-sample-problem. *Advances in neural information processing systems*, 19, 2006.

[112] Yujia Li, Kevin Swersky, and Rich Zemel. Generative moment matching networks. In *International conference on machine learning*, pages 1718–1727. PMLR, 2015.

[113] Reza Azad, Moein Heidari, Yuli Wu, and Dorit Merhof. Contextual attention network: Transformer meets U-Net. In *Machine Learning in Medical Imaging: 13th International Workshop, MLMI 2022, Held in Conjunction with MICCAI 2022, Singapore, September 18, 2022, Proceedings*, pages 377–386. Springer, 2022.

[114] Moein Heidari, Amirhossein Kazerouni, Milad Soltany, Reza Azad, Ehsan Khodapanah Aghdam, Julien Cohen-Adad, and Dorit Merhof. HiFormer: Hierarchical multi-scale representations using transformers for medical image segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 6202–6212, 2023.

[115] Reza Azad, Mohammad T Al-Antary, Moein Heidari, and Dorit Merhof. TransNorm: Transformer provides a strong spatial normalization mechanism for a deep segmentation model. *IEEE Access*, 10:108205–108215, 2022.

[116] Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L Yuille, and Yuyin Zhou. TransUNet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306*, 2021.

[117] Reza Azad, Moein Heidari, Moein Shariatnia, Ehsan Khodapanah Aghdam, Sanaz Karimijafarbigloo, Ehsan Adeli, and Dorit Merhof. TransDeepLab: Convolution-free transformer-based deeplab v3+ for medical image segmentation. In *International Workshop on PRedictive Intelligence In MEdicine*, pages 91–102. Springer, 2022.

[118] Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning Wang. Swin-unet: Unet-like pure transformer for medical image segmentation. In *Proceedings of the European Conference on Computer Vision Workshops(ECCVW)*, 2022.

[119] Ehsan Khodapanah Aghdam, Reza Azad, Maral Zarvani, and Dorit Merhof. Attention swin u-net: Cross-contextual attention mechanism for skin lesion segmentation. *arXiv preprint arXiv:2210.16898*, 2022.

[120] 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.

[121] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2337–2346, 2019.

[122] Fabian Isensee, Paul F Jaeger, Simon AA Kohl, Jens Petersen, and Klaus H Maier-Hein. nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods*, 18(2):203–211, 2021.

[123] Aimon Rahman, Jeya Maria Jose Valanarasu, Ilker Hacihaliloglu, and Vishal M Patel. Ambiguous medical image segmentation using diffusion models. *arXiv preprint arXiv:2304.04745*, 2023.

[124] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, pages 565–571. IEEE, 2016.

[125] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *Advances in neural information processing systems*, 30, 2017.

[126] Simon Kohl, Bernardino Romera-Paredes, Clemens Meyer, Jeffrey De Fauw, Joseph R Ledsam, Klaus Maier-Hein, SM Eslami, Danilo Jimenez Rezende, and Olaf Ronneberger. A probabilistic U-Net for segmentation of ambiguous images. *Advances in neural information processing systems*, 31, 2018.

[127] Samuel G Armato III, Geoffrey McLennan, Michael F McNitt-Gray, Charles R Meyer, David Yankelevitz, Denise R Aberle, Claudia I Henschke, Eric A Hoffman, Ella A Kazerooni, Heber MacMahon, et al. Lung image database consortium: developing a resource for the medical imaging research community. *Radiology*, 232(3):739–748, 2004.

[128] Aaron Carass, Snehashis Roy, Amod Jog, Jennifer L Cuzzocreo, Elizabeth Magrath, Adrian Gherman, Julia Button, James Nguyen, Ferran Prados, Carole H Sudre, et al. Longitudinal multiple sclerosis lesion segmentation: resource and challenge. *NeuroImage*, 148:77–102, 2017.

[129] Florentin Bieder, Julia Wolleb, Alicia Durrer, Robin Sandkuehler, and Philippe C. Cattin. Memory-efficient 3d denoising diffusion models for medical image processing. In *Medical Imaging with Deep Learning*, 2023.

[130] Zhicheng Zhang, Lequan Yu, Xiaokun Liang, Wei Zhao, and Lei Xing. TransCT: dual-path transformer for low dose computed tomography. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 55–64. Springer, 2021.

[131] Achleshwar Luthra, Harsh Sulakhe, Tanish Mittal, Abhishek Iyer, and Santosh Yadav. Eformer: Edge enhancement based transformer for medical image denoising. *arXiv preprint arXiv:2109.08044*, 2021.

[132] Salome Kazeminia, Christoph Baur, Arjan Kuijper, Bram van Ginneken, Nassir Navab, Shadi Albarqouni, and Anirban Mukhopadhyay. GANs for medical image analysis. *Artificial Intelligence in Medicine*, 109:101938, 2020.

[133] Wenjun Xia, Qing Lyu, and Ge Wang. Low-dose CT using denoising diffusion probabilistic model for 20× speedup. *arXiv preprint arXiv:2209.15136*, 2022.

[134] Qi Gao, Zilong Li, Junping Zhang, Yi Zhang, and Hongming Shan. CoreDiff: Contextual error-modulated generalized diffusion model for low-dose CT denoising and generalization. *arXiv preprint arXiv:2304.01814*, 2023.

[135] Ipek Oguz, Joseph D Malone, Yigit Atay, and Yuankai K Tao. Self-fusion for OCT noise reduction. In *Medical Imaging 2020: Image Processing*, volume 11313, pages 45–50. SPIE, 2020.

[136] Dewei Hu, Joseph D Malone, Yigit Atay, Yuankai K Tao, and Ipek Oguz. Retinal OCT denoising with pseudo-multimodal fusion network. In *International Workshop on Ophthalmic Medical Image Analysis*, pages 125–135. Springer, 2020.

[137] Long Zhou, Joshua D Schaefferkoetter, Ivan WK Tham, Gang Huang, and Jianhua Yan. Supervised learning with cyclegan for low-dose FDG PET image denoising. *Medical image analysis*, 65:101770, 2020.

[138] Tzu-An Song, Samadrita Roy Chowdhury, Fan Yang, and Joyita Dutta. PET image super-resolution using generative adversarial networks. *Neural Networks*, 125:83–91, 2020.

[139] Tiange Xiang, Mahmut Yurt, Ali B Syed, Kawin Setsompop, and Akshay Chaudhari. $DDM^2$: Self-supervised diffusion MRI denoising with generative diffusion models. In *The Eleventh International Conference on Learning Representations*, 2023.

[140] Eleftherios Garyfallidis, Matthew Brett, Bagrat Amirbekian, Ariel Rokem, Stefan Van Der Walt, Maxime Descoteaux, Ian Nimmo-Smith, and Dipy Contributors. Dipy, a library for the analysis of diffusion MRI data. *Frontiers in neuroinformatics*, 8:8, 2014.

[141] Ariel Rokem. Stanford HARDI surfaces. 10 2016.

[142] Kenneth Marek, Danna Jennings, Shirley Lasch, Andrew Siderowf, Caroline Tanner, Tanya Simuni, Chris Coffey, Karl Kieburtz, Emily Flagg, Sohini Chowdhury, Werner Poewe, Brit Mollenhauer, Paracelsus-Elena Klinik, Todd Sherer, Mark Frasier, Claire Meunier, Alice Rudolph, Cindy Casaceli, John Seibyl, Susan Mendick, Norbert Schuff, Ying Zhang, Arthur Toga, Karen Crawford, Alison Ansbach, Pasquale De Blasio, Michele Piovella, John Trojanowski, Les Shaw, Andrew Singleton, Keith Hawkins, Jamie Eberling, Deborah Brooks, David Russell, Laura Leary, Stewart Factor, Barbara Sommerfeld, Penelope Hogarth, Emily Pighetti, Karen Williams, David Standaert, Stephanie Guthrie, Robert Hauser, Holly Delgado, Joseph Jankovic, Christine Hunter, Matthew Stern, Baochan Tran, Jim Leverenz, Marne Baca, Sam Frank, Cathi-Ann Thomas, Irene Richard, Cheryl Deeley, Linda Rees, Fabienne Sprenger, Elisabeth Lang, Holly Shill, Sanja Obradov, Hubert Fernandez, Adrienna Winters, Daniela Berg, Katharina Gauss, Douglas Galasko, Deborah Fontaine, Zoltan Mari, Melissa Gerstenhaber, David Brooks, Sophie Malloy, Paolo Barone, Katia Longo, Tom Comery, Bernard Ravina, Igor Grachev, Kim Gallagher, Michelle Collins, Katherine L. Widnell, Suzanne Ostrowizki, Paulo Fontoura, Tony Ho, Johan Luthman, Marcel van der Brug, Alastair D. Reith, and Peggy Taylor. The parkinson progression marker initiative (ppmi). *Progress in Neurobiology*, 95(4):629–635, 2011. Biological Markers for Neurodegenerative Diseases.

[143] Guha Balakrishnan, Amy Zhao, Mert R Sabuncu, John Guttag, and Adrian V Dalca. An unsupervised learning model for deformable medical image registration. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9252–9260, 2018.

[144] Adrian V Dalca, Guha Balakrishnan, John Guttag, and Mert R Sabuncu. Unsupervised learning for fast probabilistic diffeomorphic registration. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 729–738. Springer, 2018.

[145] Kai Packhäuser, Lukas Folle, Florian Thamm, and Andreas Maier. Generation of anonymous chest radiographs using latent diffusion models for training thoracic abnormality classification systems. *arXiv preprint arXiv:2211.01323*, 2022.

[146] Cecilia Di Ruberto, Lorenzo Putzu, HR Arabnia, and T Quoc-Nam. A feature learning framework for histology images classification. *Emerging trends in applications and infrastructures for computational biology, bioinformatics, and systems biology: systems and applications*, pages 37–48, 2016.

[147] Abhishek Vahadane, Tingying Peng, Amit Sethi, Shadi Albarqouni, Lichao Wang, Maximilian Baust, Katja Steiger, Anna Melissa Schlitter, Irene Esposito, and Nassir Navab. Structure-preserving color normalization and sparse stain separation for histological images. *IEEE transactions on medical imaging*, 35(8):1962–1971, 2016.

[148] Jooyoung Choi, Jungbeom Lee, Chaehun Shin, Sungwon Kim, Hyunwoo Kim, and Sungroh Yoon. Perception prioritized training of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11472–11481, 2022.

[149] Robert L Grossman, Allison P Heath, Vincent Ferretti, Harold E Varmus, Douglas R Lowy, Warren A Kibbe, and Louis M Staudt. Toward a shared vision for cancer genomic data. *New England Journal of Medicine*, 375(12):1109–1112, 2016.

[150] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. In *International Conference on Learning Representations*, 2018.

[151] Lan Jiang, Ye Mao, Xi Chen, Xiangfeng Wang, and Chao Li. CoLa-Diff: Conditional latent diffusion model for multi-modal MRI synthesis. *arXiv preprint arXiv:2303.14081*, 2023.

[152] Liyang Chen, Zhiyuan You, Nian Zhang, Juntong Xi, and Xinyi Le. UTRAD: Anomaly detection and localization with u-transformer. *Neural Networks*, 147:53–62, 2022.

[153] Maximilian E Tschuchnig and Michael Gadermayr. Anomaly detection in medical imaging-a mini review. *Data Science–Analytics and Applications*, pages 33–38, 2022.

[154] Tharindu Fernando, Harshala Gammulle, Simon Denman, Sridha Sridharan, and Clinton Fookes. Deep learning for medical anomaly detection–a survey. *ACM Computing Surveys (CSUR)*, 54(7):1–37, 2021.

[155] Cosmin I Bercea, Benedikt Wiestler, Daniel Rueckert, and Julia A Schnabel. Reversing the abnormal: Pseudo-healthy generative networks for anomaly detection. *arXiv preprint arXiv:2303.08452*, 2023.

[156] Jian Shi, Pengyi Zhang, Ni Zhang, Hakim Ghazzai, and Yehia Massoud. Dissolving is amplifying: Towards fine-grained anomaly detection. *arXiv preprint arXiv:2302.14696*, 2023.

[157] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021.

[158] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 590–597, 2019.

[159] David Zimmerer, Simon Kohl, Jens Petersen, Fabian Isensee, and Klaus Maier-Hein. Context-encoding variational autoencoder for unsupervised anomaly detection, 2019.

[160] Md Mahfuzur Rahman Siddiquee, Zongwei Zhou, Nima Tajbakhsh, Ruibin Feng, Michael B Gotway, Yoshua Bengio, and Jianming Liang. Learning fixed points in generative adversarial networks: From image-to-image translation to disease detection and localization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 191–200, 2019.

[161] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.

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

[163] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob Mcgrew, Ilya Sutskever, and Mark Chen. GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*, pages 16784–16804. PMLR, 2022.

[164] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.

[165] Finn Behrendt, Debayan Bhattacharya, Julia Krüger, Roland Opfer, and Alexander Schlaefer. Patched diffusion models for unsupervised anomaly detection in brain MRI. In *Medical Imaging with Deep Learning*, 2023.

[166] Ujjwal Baid, Satyam Ghodasara, Suyash Mohan, Michel Bilello, Evan Calabrese, Errol Colak, Keyvan Farahani, Jayashree Kalpathy-Cramer, Felipe C Kitamura, Sarthak Pati, et al. The RSNA-ASNR-MICCAI BraTS 2021 benchmark on brain tumor segmentation and radiogenomic classification. *arXiv preprint arXiv:2107.02314*, 2021.

[167] Žiga Lesjak, Alfiia Galimzianova, Aleš Koren, Matej Lukin, Franjo Pernuš, Boštjan Likar, and Žiga Špiclin. A novel public MR image dataset of multiple sclerosis patients with lesion segmentations based on multi-rater consensus. *Neuroinformatics*, 16:51–63, 2018.

[168] Shizhan Gong, Cheng Chen, Yuqi Gong, Nga Yan Chan, Wenao Ma, Calvin Hoi-Kwan Mak, Jill Abrigo, and Qi Dou. Diffusion model based semi-supervised learning on brain hemorrhage images for efficient midline shift quantification. *arXiv preprint arXiv:2301.00409*, 2023.

[169] Matthias Keicher, Matan Atad, David Schinz, Alexandra S Gersing, Sarah C Foreman, Sophia S Goller, Juergen Weissinger, Jon Rischewski, Anna-Sophia Dietrich, Benedikt Wiestler, et al. Semantic latent space regression of diffusion autoencoders for vertebral fracture grading. *arXiv preprint arXiv:2303.12031*, 2023.

[170] Namrata Anand and Tudor Achim. Protein structure and sequence generation with equivariant denoising diffusion probabilistic models. *arXiv preprint arXiv:2205.15019*, 2022.

[171] Jaakko Lehtinen, Jacob Munkberg, Jon Hasselgren, Samuli Laine, Tero Karras, Miika Aittala, and Timo Aila. Noise2Noise: Learning image restoration without clean data. In *International Conference on Machine Learning*, pages 2965–2974. PMLR, 2018.

[172] Kwanyoung Kim and Jong Chul Ye. Noise2Score: tweedie's approach to self-supervised image denoising without clean images. *Advances in Neural Information Processing Systems*, 34:864–874, 2021.

[173] Marc Combalia, Noel CF Codella, Veronica Rotemberg, Brian Helba, Veronica Vilaplana, Ofer Reiter, Cristina Carrera, Alicia Barreiro, Allan C Halpern, Susana Puig, et al. Bcn20000: Dermoscopic lesions in the wild. *arXiv preprint arXiv:1908.02288*, 2019.

[174] Ibrahim Ethem Hamamci, Sezgin Er, Enis Simsar, Anjany Sekuboyina, Mustafa Gundogar, Bernd Stadlinger, Albert Mehl, and Bjoern Menze. Diffusion-based hierarchical multi-label object detection to analyze panoramic dental x-rays. *arXiv preprint arXiv:2303.06500*, 2023.

[175] Shoufa Chen, Peize Sun, Yibing Song, and Ping Luo. DiffusionDet: Diffusion model for object detection. *arXiv preprint arXiv:2211.09788*, 2022.

[176] Chentao Cao, Zhuo-Xu Cui, Shaonan Liu, Dong Liang, and Yanjie Zhu. High-frequency space diffusion models for accelerated MRI. *arXiv preprint arXiv:2208.05481*, 2022.

[177] Hyungjin Chung, Byeongsu Sim, and Jong Chul Ye. Come-closer-diffuse-faster: Accelerating conditional diffusion models for inverse problems through stochastic contraction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12413–12422, 2022.

[178] Hong Peng, Chen Jiang, Yu Guan, Jing Cheng, Minghui Zhang, Dong Liang, and Qiegen Liu. One-shot generative prior learned from hankel-k-space for parallel imaging reconstruction. *arXiv preprint arXiv:2208.07181*, 2022.

[179] Zongjiang Tu, Die Liu, Xiaoqing Wang, Chen Jiang, Minghui Zhang, Qiegen Liu, and Dong Liang. WKGM: Weight-k-space generative model for parallel imaging reconstruction. *arXiv preprint arXiv:2205.03883*, 2022.

[180] Zhuchen Shao, Liuxi Dai, Yifeng Wang, Haoqian Wang, and Yongbing Zhang. AugDiff: Diffusion based feature augmentation for multiple instance learning in whole slide image. *arXiv preprint arXiv:2303.06371*, 2023.

[181] Ye Mao, Lan Jiang, Xi Chen, and Chao Li. DisC-Diff: Disentangled conditional diffusion model for multi-contrast mri super-resolution. *arXiv preprint arXiv:2303.13933*, 2023.

[182] Sobhan Goudarzi and Hassan Rivaz. Deep ultrasound denoising without clean data. In *Medical Imaging 2023: Ultrasonic Imaging and Tomography*, volume 12470, pages 131–136. SPIE, 2023.

[183] Lun Zhang and Junhua Zhang. Ultrasound image denoising using generative adversarial networks with residual dense connectivity and weighted joint loss. *PeerJ Computer Science*, 8:e873, 2022.

[184] Vincent van de Schaft and Ruud JG van Sloun. Ultrasound speckle suppression and denoising using MRI-derived normalizing flow priors. *arXiv preprint arXiv:2112.13110*, 2021.

[185] Jeremias Traub. Representation learning with diffusion models. *arXiv preprint arXiv:2210.11058*, 2022.

[186] Korbinian Abstreiter, Sarthak Mittal, Stefan Bauer, Bernhard Schölkopf, and Arash Mehrjou. Diffusion-based representation learning. *arXiv preprint arXiv:2105.14257*, 2021.

[187] William Peebles and Saining Xie. Scalable diffusion models with transformers. *arXiv preprint arXiv:2212.09748*, 2022.

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

[189] Pedro Sanchez and Sotirios A. Tsaftaris. Diffusion causal models for counterfactual estimation. In *First Conference on Causal Learning and Reasoning*, 2022.

[190] Pedro Sanchez, Xiao Liu, Alison Q O'Neil, and Sotirios A. Tsaftaris. Diffusion models for causal discovery via topological ordering. In *The Eleventh International Conference on Learning Representations*, 2023.

[191] Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwag, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. *arXiv preprint arXiv:2301.13188*, 2023.

[192] Mansoor Ali, Faisal Naeem, Muhammad Tariq, and Geroges Kaddoum. Federated learning for privacy preservation in smart healthcare systems: A comprehensive survey. *IEEE journal of biomedical and health informatics*, 2022.