# Comprehensive Survey on Generative Diffusion Models: Foundations, Innovations, and Applications

SurveyForge

**Abstract**— Generative diffusion models have significantly advanced generative modeling by enabling high-quality data synthesis across varied domains such as image, text, and scientific applications. This comprehensive survey assesses the evolution, mechanisms, and theoretical foundations of these models, emphasizing their notable adaptability and integration with frameworks like GANs and VAEs for improved performance. Key aspects include the elaboration of forward and reverse diffusion processes, efficient sampling techniques, and architectural innovations that promote scalability and computational efficiency. Despite their strengths, challenges such as high computational costs, robustness, and ethical considerations remain. Advancing sampling methods and integrating reinforcement learning are pivotal for addressing these constraints while fostering more responsible deployment. This survey highlights the potential of diffusion models to enhance AI capabilities through cross-modal applications and proposes future research directions to maximize their impact across diverse fields. As diffusion models continue to evolve, they offer tremendous opportunities for bridging human and machine-driven creativity, underscoring the importance of interdisciplinary research to realize their full potential responsibly.

**Index Terms**—diffusion model evolution, sampling techniques innovation, cross-modal applications

✦

## 1 INTRODUCTION

GENERATIVE diffusion models have emerged as a formidable approach within the domain of generative modeling, showcasing remarkable potential in synthesizing high-quality data across various forms including images, text, and even complex scientific datasets. This subsection aims to elucidate the historical journey, foundational methodologies, and the overarching significance of generative diffusion models in modern artificial intelligence (AI).

The proliferation of diffusion models can be traced back to early interests in statistical physics and stochastic processes, particularly the use of Gaussian noise to progressively perturb data representations [1]. Historically, the methodology originated from a need to model high-dimensional data distribution spaces effectively. This motivation saw the advent of Denoising Diffusion Probabilistic Models (DDPM) and their various adaptations, which emphasized iterative noise reduction as a core operational principle [2], [3].

Fundamentally, the working mechanism of diffusion models encapsulates two critical phases: a forward diffusion process and a reverse diffusion process. These models essentially involve a gradual noising of input data, transforming it into a tractable Gaussian distribution, followed by a learned reverse diffusion process that aims to reconstruct the original data distribution from the noise. Mathematically, this can be expressed as a series of stochastic differential equations (SDEs) that describe the evolution of data over time [4]. The reverse process notably implements learned score functions to map noise back into coherent data forms, highlighting a key differentiator from traditional generative adversarial networks (GANs) which do not incorporate such explicit transformation steps.

In evaluating the strengths and limitations of diffusion models, several comparative analyses aid in discerning their unique contributions to generative modeling. Unlike GANs, diffusion models inherently exhibit a higher stabilization in training, circumventing issues such as mode collapse that often plague adversarial approaches [5]. Other attributes, such as architectural modularity and scalability, endow diffusion models with a degree of flexibility that is conducive to large-scale data generation and multifaceted applications [6]. However, the computational cost associated with iterative denoising processes remains a noted drawback, making efficient sampling and model optimization paramount in contemporary research [7].

The applications of generative diffusion models cut across a diverse range of domains. In computer vision, these models have set new benchmarks for image and video synthesis, enabling refined outputs that rival human creativity in detail and accuracy. The adaptability of diffusion models to tasks such as image inpainting or style transfer further consolidates their value proposition in visual applications [8]. Beyond vision, their utility extends to natural language processing, where diffusion frameworks have innovated text augmentation and interactive storytelling applications by integrating with large language models [9].

Moreover, the flexibility of diffusion models allows for their integration with other generative frameworks, such as VAEs and hybrid discriminative-generative models, to craft enhanced solutions capable of tackling increasingly complex datasets [10]. These synergies foster novel research directions, especially in fields where multi-modal data synthesis is increasingly relevant [11]. As foundational models in AI continue to evolve, diffusion models are expected to expand their dominance, redefining the paradigms of data synthesis and manipulation.

The significance of diffusion models is further underscored by their application in scientific and interdisciplinary domains. For instance, their adaptability to graph-based data has seen them applied to molecule and protein generation, where precision and structural fidelity are paramount [12]. Similarly, in domains like climate science and bioinformatics, diffusion models are instrumental in simulating complex dynamical systems that require an intricate understanding of temporal and spatial interdependencies [13].

Examining the challenges and future of diffusion models, several key areas emerge as focal points for ongoing and future research. Primarily, the quest for computational efficiency, specifically reducing the substantial resource requirements for training and inference, remains at the forefront [7]. Innovations such as parallel sampling and strategic architectural designs hold promise for addressing these limitations, paving the way for more accessible diffusion modeling frameworks [14].

In conclusion, the journey of generative diffusion models encapsulates a paradigm shift in how AI interprets and generates complex data structures. With continuing advancements in algorithmic efficiency and integration methodologies, these models are poised to contribute significantly to the future trajectories of AI, offering a holistic framework for creative and scientific exploration alike. The prospect of diffusion models continues to captivate scholars, with ongoing research exploring their potential to facilitate breakthroughs in unexplored domains.

## 2 THEORETICAL FOUNDATIONS

### 2.1 Mathematical Formulation of Diffusion Processes

The mathematical formulation of diffusion processes lies at the heart of understanding generative diffusion models, an advanced class of deep generative models renowned for their ability to generate high-quality, diverse data samples across various domains. Fundamentally, diffusion processes describe a framework where a data distribution is progressively transformed into noise, and subsequently, the noise is gradually reverted to reform the data distribution. This section explores the mathematical intricacies of these processes, focusing on their formulation, underlying stochastic processes, and the duality of forward and reverse diffusion mechanisms.

At the core of diffusion models are two key stochastic processes: the forward (diffusion) process and the reverse (denoising) process. The forward process can be understood as a Markov chain where noise is incrementally added to the data. Mathematically, the forward diffusion process typically begins with data sampled from a real-world distribution $p(x_0)$ and involves successive data corruption over time, resulting in a Gaussian distribution as the final state. The forward process can be denoted by $q(x_t|x_{t-1})$, which encapsulates Gaussian noise addition. Specifically, if we define a Gaussian perturbation as $\epsilon \sim \mathcal{N}(0, I)$, the forward step can be described as:

$$x_t = \sqrt{1 - \beta_t} x_{t-1} + \beta_t \epsilon, \qquad (1)$$

where $t$ signifies a discrete time step, and $\beta_t$ is the variance schedule controlling the noise level at each step.

In contrast, the reverse process aims to model the inverse of this diffusion by removing noise step-by-step to approximate the reconstruction of the original data. The reverse diffusion, often parameterized using a learned model, is typically expressed as $p_\theta(x_{t-1}|x_t)$, where the goal of the model $\theta$ is to approximate the reverse conditional distributions. The reverse process equations take the form:

$$x_{t-1} = \frac{1}{\sqrt{1 - \beta_t}} \left( x_t - \frac{\beta_t}{\sqrt{1 - \beta_t}} \epsilon_\theta(x_t, t) \right) + \sigma_t z, \qquad (2)$$

where $\epsilon_\theta(x_t, t)$ represents the model's predicted noise, and $\sigma_t$ is related to the variance of the reverse diffusion's noise term $z \sim \mathcal{N}(0, I)$ [2].

The reverse process heavily relies on accurately learning the score function—essentially, the gradient of the log probability of the data with respect to the noise level. Mathematical breakthroughs, such as Tweedie's formula, leverage this score-based approach, which is quintessential for score-based generative models [1]. The notion of score matching, specifically the minimization of the Fisher divergence, is pivotal here as it enables the estimation of such gradients efficiently [5].

While the general framework described above is common, the choice of noise schedules $\beta_t$ and the intricacies of modeling $\epsilon_\theta$ critically influence the model's performance and stability. Recent innovations have examined alterations in noise schedules, with some advocating for learnable variance schedules that adapt during training [2]. These adaptive schedules have been shown to reduce the number of forward passes required, thus improving computational efficiency without sacrificing sample quality.

Moreover, exploration into the underlying mechanics often involves stochastic differential equations (SDEs), which offer a continuous-time view of these discrete processes. The adoption of SDEs broadens the mathematical toolkit available for understanding diffusion processes within the larger framework of stochastic calculus [15]. This perspective facilitates novel advancements such as time-reversal symmetry, which further enhances reverse diffusion stability.

Despite these advancements, generative diffusion models are not without challenges. The high computational cost and slow generation speed due to iterative sampling processes are notable limitations [4]. Innovative approaches have been proposed to address these challenges, such as parallel sampling techniques and optimized architectures that diminish inference latency [16].

Looking ahead, continuous refinement in the mathematical understanding and algorithmic efficiency of diffusion processes is expected to play a pivotal role in broadening their applicability and overcoming current limitations. Future research directions include exploring diffusion models' integration with other generative frameworks for enhanced performance, addressing computational challenges, and expanding their application in new domains [3]. Through continued exploration and innovation, the mathematical depth of diffusion processes promises to unlock even greater potential for generative modeling in the future.

## 2.2 Probabilistic and Score-Based Frameworks

In the sphere of generative modeling, diffusion models stand out due to their foundation in both probabilistic and score-based frameworks. This subsection provides a detailed exploration of these approaches, emphasizing their crucial roles in bolstering model performance and reliability. At its core, the generative function of diffusion models relies on converting data into noise and reverting it back to data—a stochastic process fundamentally underpinned by probabilistic and score-based interpretations. Here, we delve into the theoretical bases, strengths, constraints, and key advancements that have driven the development of these frameworks.

Fundamentally, diffusion models utilize probabilistic methods via stochastic differential equations (SDEs) to describe the dynamic transformation of data distributions affected by noise. The process begins with forward diffusion, which gradually perturbs data with Gaussian noise, creating a trajectory from the original data distribution to a more tractable Gaussian prior. This stochastic journey is encapsulated in a forward SDE, generally expressed as $d\mathbf{x}_t = f(\mathbf{x}_t, t)dt + g(t)d\mathbf{w}_t$, where $\mathbf{x}_t$ symbolizes the data at time $t$, $f$ and $g$ represent drift and diffusion coefficients, and $\mathbf{w}_t$ is a Wiener process [17].

The reverse process—essential for generating data—is governed by the reverse-time SDE, showcasing the sophistication of probabilistic models in diffusion techniques. Its formulation [18], $d\mathbf{x}_t = dt + g(t)d\tilde{\mathbf{w}}_t$, employs the gradient field of the log-density, known as the "score," to guide the process back to the original data structure [17]. This score function, $\nabla \log p_t(\mathbf{x}_t)$, is vital for this reverse operation, highlighting the geometric features intrinsic to the data distribution.

Beyond SDE formulation, probabilistic interpretation impacts the training strategies of diffusion models. Score-based methodologies, like denoising score matching (DSM), allow the estimation of the score function by minimizing a loss that measures the difference between the true score and its approximated version. As detailed by Song et al., this approach effectively incorporates the objective within a variational perspective, where the score matching loss approximates a lower bound on the reverse diffusion's likelihood, thus providing a coherent training path [19].

Score-based approaches enhance the robustness of diffusion models by offering a means to estimate log-density gradients without explicitly constructing the data distribution—circumventing the high-dimensionality challenges common to other generative models like GANs and VAEs. However, these frameworks have limitations. Score estimation accuracy can suffer from model errors and the curse of dimensionality, adversely affecting generative performance in high-dimensional spaces. Continuous advancements in noise estimation techniques and refined architectural designs aim to address these challenges [20].

Probabilistically, the synergy between the score-matching framework and maximum likelihood estimation (MLE) enhances the generative capacity of diffusion models. Recent developments integrating maximum likelihood principles into score-based diffusion training have refined model precision and reliability by aligning likelihood estimations with training objectives [21]. These alignments have led to promising results, positioning score-based diffusion models as competitive with leading benchmarks in image synthesis and other domains.

The intersection of probabilistic and score-based frameworks necessitates examining convergence and stability within diffusion processes. Model convergence depends on the precision of score estimates and SDE discretization, both crucial to generative success. Recent research into advanced discretization and sampling techniques addresses numerical drift issues in the generative process, enhancing model fidelity despite complex manifold challenges [22]. This intersection not only emphasizes diffusion models' theoretical soundness but also showcases practical applications from integrating probabilistic insights with score-based methods.

Despite notable progress, challenges persist. The dependency on accurate score estimation remains an obstacle, prompting the exploration of innovative training methods such as advanced noise scheduling and variational inference [23]. Moreover, applying these methodologies to areas with unique structural demands like low-dimensional manifold learning and inverse problems invites further research [24], [25].

As diffusion models progress, the blend of probabilistic principles with score-based techniques is set to inspire new avenues in generative modeling. This development signals potential advancements in adaptive learning, precise likelihood estimations, and the integration of specialized domain knowledge. The amalgamation of probabilistic thoroughness with innovative score-based methods promises to address future challenges, advancing the generative capabilities of diffusion models across a diverse and intricate range of applications.

## 2.3 Theoretical Insights and Challenges

Understanding the theoretical underpinnings of generative diffusion models is crucial for advancing their development and application across various domains. This subsection delves into the theoretical insights and challenges that define the core mathematical and computational issues within diffusion models, with a focus on mathematical assumptions, convergence, and stability.

The foundation of generative diffusion models rests upon stochastic differential equations (SDEs) that facilitate the transformation of data distributions. These mathematical formulations are sensitive to the assumptions made about the nature of noise and the data distribution itself. For instance, many models assume that the diffusion process operates within a Gaussian noise framework, allowing for tractable analytic formulations and score-based estimations. However, this assumption can be restrictive, particularly when data distributions exhibit significant non-Gaussian features. The work by Song et al. in [17] explores the stochastic differential equation framework, highlighting the flexibility and limitations inherent in its application.

Further, diffusion models often assume that the diffusion process spans a continuous space, which can limit their applicability to data with intrinsic discrete structures, such as text or graph data. The transition from continuous to discrete spaces remains a significant hurdle, as demonstrated

by ongoing research aimed at adapting continuous diffusion frameworks for discrete data by leveraging techniques such as discrete time marking chains [26].

Convergence in diffusion models is defined by the system's ability to map noisy input through iterative steps back to a coherent, high-fidelity output that resembles the training data distribution. However, achieving reliable convergence is fraught with challenges. Current diffusion models benefit from insights into probability density function optimization, but ensuring convergence across varied datasets and conditions remains complex. Papers such as [27] illustrate the convergence properties tied to score-based models, emphasizing the challenges inherent in maintaining accurate score estimation throughout the diffusion process.

Different algorithms exhibit different convergence rates, and the choice of sampling technique and parameter settings—time schedules and noise scales, for instance—critically influences the trajectory of convergence. Studies like [28] investigate the analytical convergence guarantees for diffusion models, especially when optimizing variance and noise parameters to facilitate convergence.

Stability is another central theoretical consideration that influences the reliability and robustness of diffusion models. Instabilities can cause generative processes to diverge or produce suboptimal solutions, rendering the models ineffective. Various approaches have been introduced to regulate these instabilities, such as ensuring balanced model architectures and carefully tuned noise schedules [29].

The role of score matching techniques is pivotal here, as they govern how well the model can learn the required perturbations to reconstruct the distribution through sampling. Previous works, such as [30], demonstrate that minimizing distance metrics like the Wasserstein distance can effectively bolster model robustness and improve output stability.

Emerging trends in the field are shifting towards resolving these theoretical challenges through innovative approaches. Among these are attempts to enhance sampling methods for speed without sacrificing quality, as outlined in [31]. These approaches introduce adaptive steps and high-order solvers to hasten the sampling phase, which traditionally is computationally intensive.

Moreover, theoretical work is increasingly exploring hybrid models that aim to blend diffusion models with other generative strategies—an effort that may assist in addressing scale and dimensionality challenges. For instance, [32] introduces control variations to manage the diffusion process's complexity, helping focus on both speed and fidelity.

Despite these advancements, significant hurdles remain, particularly around integrating these frameworks into practical, real-world applications that demand high efficiency and quality simultaneously. To drive future research, the field must continue to innovate across theoretical and algorithmic dimensions, embracing cross-disciplinary insights that can address the deeply rooted challenges of convergence, stability, and computational feasibility.

In conclusion, while the theoretical landscape of generative diffusion models is well-grounded, the path to comprehensive solutions entails navigating rigorous mathematical landscapes and computational limitations. Ongoing and future research must bridge these gaps with innovative frameworks, uncovering comprehensive methodologies that enhance the utility and efficacy of diffusion models across countless applications.

## 2.4 Variational and Optimal Transport Perspectives

In the exploration of generative diffusion models, the variational and optimal transport perspectives offer profound insights into refining both the construction and optimization of these models. This subsection delves into these advanced theoretical paradigms, highlighting their contributions to enhancing model performance and inspiring new avenues for exploration in generative modeling.

Variational approaches leverage the principles of variational inference, renowned for approximating posterior distributions in complex probabilistic models with efficiency. This technique adeptly addresses the challenge of intractable likelihoods, a common barrier in generative modeling, by optimizing a surrogate evidence lower bound (ELBO). The merits of variational inference lie in its systematic handling of uncertainty and its capability to learn efficient representations within diffusion models. By minimizing the Kullback-Leibler divergence between the approximate and true posteriors, variational methods ensure close alignment with the optimal latent distribution [33].

Further integration of variational methods within the diffusion framework capitalizes on the stochastic processes intrinsic to these models, framing the generative process as a latent variable model parameterized by stochastic differential equations (SDEs). This integration proves advantageous for developing deep generative models with enhanced flexibility and robustness, facilitated through stochastic control perspectives [15], [33]. These integrations amplify the expressiveness and scalability of diffusion models, enabling them to adeptly handle diverse data complexities while upholding computational feasibility.

In parallel, optimal transport theory provides a geometric framework crucial for comprehending and addressing discrepancies between probabilistic distributions. This approach has emerged as a cornerstone in diffusion models to efficiently transform input distributions into target distributions. The core efficiency of this transformation is encapsulated by the minimization of transportation cost, often gauged by the Wasserstein distance [34]. This metric not only offers robustness in optimizing generative models but also aligns closely with empirical discrepancies encountered during training, thereby bolstering the fidelity and diversity of synthesized samples.

Embedding optimal transport theories into generative modeling results in several benefits. Firstly, it promotes effective handling of high-dimensional data distributions by optimizing pathways that cater to both global structure and local nuances. Secondly, it enhances comparability between model and data distributions, establishing a foundation for rigorous evaluation metrics and potential model improvements [35]. The Monge-Ampère equation, for instance, serves as a continuous gradient flow that matches the latent space with the data space efficiently, facilitating tractable likelihood estimation and scalable inference processes [35].

Despite their strengths, each approach comes with trade-offs. Variational methods, while computationally efficient and widely applicable, may encounter challenges in securing tight bounds for complex distributions or in navigating

model selection constraints. Conversely, optimal transport methods, although precise and theoretically robust, often demand substantial computational resources and sophisticated implementations, which can confine their application to large-scale or real-time scenarios.

Emerging trends focus on hybrid architectures that marry the strengths of variational inference and optimal transport, utilizing the former for model training and the latter for model evaluation. These hybrid models present promising solutions across various applications, enhancing the interpretability of generated structures and enabling more nuanced control over the generative process [36].

Future research directions highlight the development of algorithms to efficiently harness these theoretical frameworks. This involves optimizing the computational demands of optimal transport methods through parallelization and adaptive algorithms, as well as elevating the expressiveness of variational models with advanced neural architectures and training strategies. Further investigation into adapting these methodologies to emerging applications—such as multi-modal generative tasks or reinforcement learning-based generative policies, which demand high precision and adaptability—presents fertile ground for exploration [37].

In conclusion, variational and optimal transport perspectives not only deepen our theoretical grasp of diffusion models but also enrich the practical toolkit available for their development. Through continued refinement and exploration of these methodologies, researchers can unlock new potentials for diffusion models, extending their applicability and impact across the diverse field of artificial intelligence.

## 2.5 Extensions to Complex Domains

The integration of diffusion models into complex domains raises significant theoretical and practical challenges due to the inherent structural and dimensional intricacies associated with such domains. This subsection delves into the advanced methodologies devised to adapt generative diffusion models for high-dimensional spaces and structured data environments, such as graphs or manifolds. We aim to map the terrain of extending diffusion models to these complex settings, analyzing the strategies, strengths, limitations, and potential trajectories for future research.

Various techniques have emerged to tackle the high-dimensional complexity inherent in diffusion models, inspired by the need to accurately capture the manifold structure of the data without succumbing to the "curse of dimensionality". One promising direction is the incorporation of manifold learning techniques, which seek to embed the data in a lower-dimensional latent space that preserves its geometric properties. Riemannian diffusion models, for instance, extend classical diffusion processes to accommodate arbitrary manifold structures, harnessing the power of Riemannian geometry to maintain fidelity to the data's inherent structure [38]. The use of Riemannian divergence computation offers an enhanced variational framework for likelihood estimation in these spaces, allowing for both expressive power and theoretical rigor to be achieved in applications spanning from spherical to hyperboloidal manifolds.

In another vein, diffusion models have been successfully adapted for use on structured data, such as graphs, where the topology and connectivity impose unique constraints. These adaptations often necessitate novel formulations of the diffusion process to respect the graph's structure while maintaining efficient computation and sampling. The Dynamically Conditioned Optimal Transport through Simulation-Free Flows offers an example of harnessing the inherent geometry of conditional optimal transport to generate structurally coherent data in high-dimensional environments [39]. This approach highlights how advanced transport methodologies can effectively mediate between source and target distributions in structured domains.

Beyond manifold-based approaches, high-dimensional diffusion models have also benefited from insights offered by optimal transport theory, particularly in defining paths or flows that minimize specific cost functions while efficiently capturing data distributions. For example, the use of Monge-Ampère flow introduces a continuous-time gradient flow conceptually tied to optimal transport, revealing promising results in modeling data distribution through a compressible fluid dynamic system [35]. Such approaches not only extend the application of diffusion models in infinite dimensions but also enhance performance by adhering to the intrinsic geometry of the data.

However, these methodologies are not without their limitations. Issues of computational scalability and stability remain particularly pertinent in high-dimensional settings. The mathematical formulation of diffusion processes in these domains often involves solving complex differential equations, which pose significant computational demands. Techniques such as the Iterated Diffusion Bridge Mixture (IDBM) attempt to address these challenges by offering scalable solutions that leverage Schrödinger Bridge algorithms, thus demonstrating superior sample quality over larger discretization intervals [40].

The intersection of generative diffusion models with optimal transport also underscores a promising research trajectory. Innovations in unbalanced optimal transport frameworks provide enhanced robustness against outliers and superior optimization stability, offering scalability to diffusion processes initially constrained by the assumptions of classical optimal transport methods [41]. These advancements align with broader trends in variance reduction techniques and adaptive optimization, which collectively aim to refine the robustness and efficiency of generative models in capturing complex, high-dimensional distributions.

Another crucial aspect is the extension to constrained domains, where traditional assumptions of forward and backward diffusion processes may not hold. These settings often require novel analytical tools to define boundaries, constraints, and transitions effectively. The Reflected Schrödinger Bridge, for instance, offers an elegant solution by integrating boundary conditions into the stochastic differential equations that describe the diffusion process [42]. This method underscores the benefits of incorporating optimal transport constraints in maintaining tractable and effective models across varied domains.

Looking forward, the synthesis of these approaches suggests several emergent trends in the diffusion model landscape. There is a marked movement toward integrating con-

textual and domain-specific knowledge directly into the diffusion modeling framework. Advances in constraint-aware training protocols, new variational approaches for infinite-dimensional settings, and the development of physics-informed models reveal a growing appreciation for the multifaceted nature of complex domains, which demand nuanced approaches to model training and inference [43].

Moreover, there is substantial interest in improving the adaptive capabilities of diffusion models. The promise of scalable and robust methodologies lies in their potential to dynamically adjust to data properties, thus ensuring accurate and convincing generative performance across diverse application areas. As diffusion models continue to traverse the interface of statistical mechanics, optimal transport, and information theory, their adaptability to complex domains will likely enhance their efficacy in cutting-edge applications across scientific and industrial landscapes.

In conclusion, the extension of diffusion models into complex domains is an area ripe with challenges and opportunities. By drawing on theoretical insights from manifold learning, optimal transport, and constraint handling, researchers can push the boundaries of what diffusion models can achieve. These advancements not only contribute fundamental understanding but also pave the way for new applications and innovations in the rapidly evolving landscape of generative modeling.

## 3  ALGORITHMIC ENHANCEMENTS AND INNOVATIONS

### 3.1  Advances in Sampling Techniques

The pursuit of improving sampling techniques in generative diffusion models has gained significant traction in recent years, driven by the increasing demand for efficient and scalable solutions. This subsection critically examines the evolution of sampling methodologies, emphasizing innovations that streamline the generative processes by reducing computational burdens while maintaining or enhancing output quality.

Generative diffusion models, at their core, rely on sampling to iteratively reverse a stochastic process, traditionally implemented via stochastic differential equations (SDEs) or discrete-time Markov chains. These processes, however, can be computationally demanding due to their stepwise nature where each step incrementally refines the sample. The conventional approach requires numerous forward and backward passes to achieve high fidelity outputs, leading to high latency and resource consumption, as noted in [2].

One significant advancement in this area is the development of stochastic and exponential integrator sampling techniques. These methods aim to optimize the discretization steps within the diffusion process, often incorporating auxiliary variables to increase precision and speed [33]. By redefining the underlying mathematical framework for sampling, such integrators balance trade-offs between computation and accuracy. For example, stochastic integrators use tailored noise distributions to enhance variability control, providing robustness against numerical instabilities without heavily compromising on speed.

Beyond these integrators, multi-stage and parallel sampling frameworks have emerged as robust solutions to computational challenges. The integration of multi-stage frameworks involves multiple, succinct phases of sampling that effectively leverage hierarchical modeling approaches. This advancement condenses the sampling timeline and improves model scaling to larger datasets, as illustrated by the breakthroughs in parallel sampling, which splits the sampling process into concurrent tasks that drastically reduce overall inference time [16]. Such parallelization is crucial in modern GPU architectures, where data and model parallelism can be substantially utilized to enhance efficiency.

Furthermore, conditional and hybrid sampling methods have become pivotal in extending the flexibility of generative diffusion models. These approaches incorporate additional conditional variables into the sampling process, allowing models to adapt to a wide range of input conditions. Hybrid frameworks, such as those combining score-based and variational techniques, have demonstrated improved adaptability in modality-specific data synthesis, offering a nuanced balance between conditioning and generative capabilities [5]. The conditional techniques have been expanded with manifold-constrained methods, which ensure that the sampling remains within plausible data manifold bounds, addressing potential issues of mode collapse and improving sample diversity [44].

Each of these advancements brings unique strengths; however, they present challenges that warrant further exploration. For instance, while multi-stage and parallel sampling frameworks significantly reduce computational demands, they introduce complexities in managing synchronization and data consistency across parallel processes. Achieving a balance between precision and speed in stochastic integrator techniques remains an ongoing research focus, as the trade-offs can vary significantly depending on the model architecture and inference objectives [45].

Moreover, hybrid approaches, though flexible, require careful calibration to manage the additional conditional variables, which can lead to increased complexity in both training and inference phases. The integration of hybrid strategies with existing architectures also poses non-trivial integration challenges, particularly in tuning the interplay between continuous and discrete components within a unified probabilistic framework [46].

Looking forward, future research directions suggest a trajectory towards more adaptive and context-aware sampling strategies. Exploration of self-tuning algorithms that dynamically adjust sampling parameters in response to real-time data characteristics would represent a significant leap in efficiency and adaptability. There is also a potential in exploiting the intersection of reinforcement learning and diffusion processes to optimize sampling strategies through learned dynamic feedback mechanisms [47]. Integrating these advancements into scalable architectures will likely set the stage for next-generation adaptive diffusion models capable of operating across diverse domains with heightened efficiency and specificity.

In conclusion, the innovations in sampling techniques for generative diffusion models provide exciting avenues for both practical application and theoretical exploration. By addressing the dual imperatives of efficiency and sample quality, these advancements pave the way for diffusion models

to be more accessible, versatile, and powerful tools in the AI landscape. As these techniques continue to mature, they promise to extend the frontiers of what is achievable with generative modeling, particularly in areas requiring high-quality synthesis under strict computational constraints.

## 3.2  Training and Optimization Improvements

In recent years, generative diffusion models have witnessed impressive advancements, yet their training phase presents notable challenges related to convergence, stability, and computational efficiency. These challenges have spurred intensive research into innovative training methods and optimization strategies, aiming to ensure more effective learning. This subsection delves into the latest algorithmic developments addressing these challenges, offering a comparative analysis of methodologies, examining their strengths and limitations, and exploring future directions.

Central to the training of diffusion models is the accurate estimation of score functions, which are crucial for reversing the diffusion process and effectively synthesizing data from noise. An important improvement in training schemes is the use of curriculum and momentum-based learning strategies. These strategies adaptively adjust learning rates and training difficulty to improve convergence speeds and stability. By progressively increasing the complexity of training samples or dynamically adjusting the noise schedule, models can grasp easier patterns before tackling more difficult tasks, thereby minimizing computational overhead while enhancing model robustness.

Significant advancements also include parameter-efficient and transfer learning approaches. These involve fine-tuning pre-trained models to allow diffusion models to adapt across different tasks without incurring high resource costs. Modular customization strategies enable selective training of specific model components, reducing the need to train the entire model from scratch. This modular approach facilitates efficient domain adaptation, permitting the same model architecture to be tailored to various contexts by modifying only task-specific components [48].

Noise optimization and dynamic scheduling play pivotal roles in optimizing the training of diffusion models. Noise parameters, which determine data corruption during the forward diffusion process, are critical for model performance and efficiency. Recent research has focused on developing methods for step-by-step adjustment of these parameters to ensure optimal learning trajectories regardless of the number of denoising steps [20]. This dynamic scheduling enhances robustness to varying noise levels and converges to a more stable training regime, reducing performance degradation from suboptimal noise configurations.

A recent methodological exploration has been the incorporation of maximum likelihood training frameworks in conjunction with score-based techniques. Studies indicate that linking score matching loss with log-likelihood provides a robust mechanism to enhance likelihood estimation of the model's output [21]. This integration achieves a more consistent optimization criterion aligning with the generative task's objectives, potentially leading to improved generalization across diverse datasets and tasks.

The landscape of diffusion model training also benefits from advanced optimization frameworks, such as Stochastic Gradient Descent (SGD) with momentum adaptations and the Adam optimizer. These frameworks accommodate the peculiarities of diffusion processes, which include high-dimensional parameter spaces and the necessity for fine-grained control over training dynamics. By fine-tuning hyperparameters such as learning rates and decay schedules, these optimizers significantly stabilize training and accelerate convergence.

Despite these advancements, challenges remain in balancing training efficiency with model performance, especially considering the computational demands of high-dimensional generative tasks. Future directions include further refinement of transfer learning techniques, particularly concerning domain shift and adversarial robustness. Moreover, exploring novel optimization algorithms leveraging modern computing architectures, such as distributed and parallel training paradigms, promises to mitigate existing bottlenecks and enhance scalability.

Emerging trends suggest a growing interest in hybrid training frameworks integrating elements from different generative approaches, such as Variational Autoencoders (VAEs) and generative adversarial networks (GANs), to exploit their complementary strengths. This cross-breeding could lead to models that are both flexible and robust, overcoming limitations inherent to each standalone framework [19]. Additionally, adopting insights from differential equations and statistical mechanics, as evidenced in score-based diffusion models' application in broader scientific domains, may provide innovative pathways for refining training methodologies.

In conclusion, the training and optimization of diffusion models continue to evolve rapidly, driven by the complexity of underlying processes and the demands for high efficiency and scalability. A strategic focus on enhancing training stability, convergence, and adaptability will undoubtedly catalyze the next wave of innovations in this dynamic field. Integrating multidisciplinary insights is anticipated to pave the way for breakthroughs that could redefine model performance benchmarks and open up new vistas in generative modeling applications.

## 3.3  Architectural and Structural Innovations

In recent years, the architectural landscape of generative diffusion models has undergone significant evolution with innovations aimed at enhancing their generative capabilities and scalability. This subsection delves into these architectural advancements, focusing on how they mold the efficiency, flexibility, and applicability of diffusion models in various domains. By dissecting foundational components such as model modularity, integration with latent spaces, and the optimization of neural networks like U-Nets and Vision Transformers, we offer a comprehensive analysis of the state-of-the-art architectural strategies that bolster the performance of diffusion models.

The transformation of network architectures in generative diffusion models often begins with modular and scalable designs. The adoption of modular structures, akin to those used in other domains of machine learning, facilitates scalability and customization. A modular approach allows components of the diffusion model to operate independently and be individually optimized or replaced, leading to

enhanced performance across diverse tasks. This flexibility is critical for deploying models across varying environments and computational budgets. Such architectures, similar to stackable layers in convolutional neural networks, offer scalable solutions that cater to the increasing complexity and data dimensionalities faced by diffusion models.

Moreover, the integration of latent and hybrid diffusion frameworks provides a compelling avenue for enhancing both the expressiveness and efficiency of generative models. By incorporating latent variable models, such as Variational Autoencoders (VAEs), into the diffusion process, these models gain the ability to learn more nuanced feature representations. Latent spaces act as bottlenecks that can encode complex data distributions into simpler forms, facilitating efficient learning and sampling procedures [49]. This synergy allows diffusion models to handle data with intricate structures, enhancing their versatility across applications ranging from image synthesis to more structured domains like graph modeling [50].

Significant strides have also been made by optimizing core neural network architectures within diffusion models. Among them, U-Nets and Vision Transformers (ViTs) have emerged as pivotal structures. U-Nets, widely recognized for their success in image segmentation tasks, introduce an encoder-decoder architecture with skip connections that preserve spatial hierarchies essential for generating high-fidelity outputs. Their ability to integrate detail retention with efficient feature extraction makes them well-suited for diffusion models dealing with high-resolution image synthesis. Furthermore, enhancements to U-Nets, such as deeper architectures or pixel-shuffle layers, can further elevate the fidelity of generated samples [21].

Vision Transformers bring another layer of innovation, particularly in their capacity to model long-range dependencies through attention mechanisms. By processing data as sequences, ViTs can capture intricate patterns that might be overlooked by convolutional networks. The adaptation of ViTs in diffusion models underscores an architectural shift towards capturing global contextual information in generative tasks. This attention-based approach is particularly advantageous for tasks requiring integration of diverse data modalities or maintaining coherence across sequential outputs, such as video generation [49].

The union of these architectural paradigms amounts to a robust framework for diffusion models, yet it also introduces a set of challenges and trade-offs. Modularity offers flexibility but may lead to increased computational overhead and complexity in model management. The integration of latent spaces can enhance expressiveness but poses risks of added noise in the latent representations, potentially complicating the training dynamics. Similarly, while U-Nets and ViTs introduce substantial generative improvements, their computational demands require careful resource balancing, especially in resource-constrained environments.

Moreover, as diffusion models continue to integrate more complex architectures, it becomes increasingly imperative to address the emerging challenge of computational efficiency. The intricate architectures often necessitate high-capacity computational resources for training and inference, which can become a bottleneck for widespread deployment. Techniques such as distillation [51], pruning, and quantization are being explored to alleviate these computational burdens. These strategies aim to streamline models, reducing inference times and resource consumption without significant detriment to generative quality.

Emerging trends also signal a growing interest in the harmonization of different architectural innovations. For instance, exploring the fusion of U-Net architectures with transformer dynamics might offer a hybrid approach capable of leveraging the strengths of both paradigms — the detailed local processing of U-Nets combined with the global context capturing of transformers. Experimenting with such innovative combinations could serve as a pathway to overcoming the limitations typically associated with one singular approach.

Looking forward, the landscape of architectural design in diffusion models is poised to extend beyond traditional boundaries. The synthesis of current innovations with novel approaches like synthetic feature engineering, adaptive networks, and AI-driven optimization holds promise for more intelligent, adaptable diffusion models. As research progresses, fostering interdisciplinary dialogue and collaboration will be pivotal in shaping architectures that are not only technically sound but also ethically grounded and responsive to varied application contexts.

In conclusion, the architectural evolution within generative diffusion models illuminates a path of increased complexity and capability. While significant progress has been made, ongoing innovation and rigorous exploration of architectural frontiers remain crucial. This dynamic interplay of designs will continue to evolve as researchers strive to refine generative diffusion models' ability to create high-quality, contextually relevant data across an expanding array of domains [49], [50].

## 3.4 Improved Model Customization and Adaptability

In recent years, the versatility of generative diffusion models has expanded significantly, fueled by the increasing demand for customizable and adaptable solutions that can cater to diverse application requirements across various domains. Building on the architectural innovations discussed earlier, this section delves into advances that transform diffusion models into agile tools for domain-specific applications, highlighting methodologies that promote minimal reconfiguration while maintaining or enhancing performance.

A notable advancement in model customization is the development of orthogonal adaptation and modular customization techniques. These strategies enable diffusion models to seamlessly integrate auxiliary modules tailored to specific tasks, such as noise adaptation layers or context-specific embedding modules, thus enhancing their applicability across different contexts. Modular architectures provide the flexibility to add or remove components based on task requirements without altering the model's core structure [52]. This flexibility is crucial for adapting diffusion models to specific needs while preserving their foundational capabilities.

An illustrative example of modular innovation is the use of specialized plug-and-play modules, akin to the adaptor layers used in language models, which allow diffusion models to adjust their outputs based on varying input conditions

or desired specifications. This approach is exemplified by Multi-Architecture Multi-Expert diffusion models, which strategically employ distinct architectural designs for different tasks and time frames within the generative process to enhance adaptability and efficiency [53]. These architectures effectively balance computational demand and precision, adapting the operation mode to fulfill requirements ranging from light computational loads for faster, less resource-intensive tasks to highly detailed, resource-demanding ones.

Simultaneously, the integration of multi-modal and cross-modal techniques has emerged as a vital direction in improving the adaptability of diffusion models. The ability to handle inputs from disparate domains—such as text, audio, and images—within a unified framework is increasingly critical. Such frameworks facilitate the synthesis of heterogeneous data streams, enhancing the model's utility in complex scenarios like automated caption generation from image and video streams or the development of immersive mixed-reality environments. The versatility of these multi-modal frameworks is already evident in fields like bioinformatics and healthcare, where diffusion models are tailored to integrate and generate insights from multimodal datasets, improving the robustness and interpretability of predictive models [54].

Recent advancements in cross-modal and multi-modal generation have underscored the importance of shared latent spaces capable of encoding multiple data types. These frameworks typically employ joint embeddings to facilitate the seamless conversion of information across modalities while maintaining coherence and context. Techniques such as joint conditional modeling utilizing diffusion models leverage domain knowledge embedded within the latent spaces, thereby significantly reducing the gap between modality-specific training data contributions [55]. This approach empowers diffusion models to synthesize richer and more contextually relevant outputs, ultimately broadening their applicability to more complex and dynamic environments.

Furthermore, the increasing importance of domain-specific applications in scientific and interdisciplinary fields propels the demand for improved model customization. In such spheres, diffusion models show promise in specialized applications like climate simulation, bioinformatics, and other scientific domains, where they generate precise models capable of understanding and predicting complex phenomena. For instance, in climate modeling, diffusion models can be customized to simulate intricate environmental patterns over time and space, thus aiding the development of comprehensive predictive systems [56].

In these contexts, personalization is often supported through targeted training strategies and domain adaptation, where models are selectively fine-tuned on domain-specific datasets to optimize their performance for particular scientific or practical applications [57]. Domain adaptation often involves transfer learning methodologies that allow a base diffusion model, pre-trained on general data, to be recalibrated with specialized datasets. This approach not only minimizes the data and computational overhead associated with training models from scratch but also enhances the models' adaptability to the target domain.

As the field progresses, several challenges and future directions warrant consideration. The primary challenge remains the computational cost and complexity involved in customizing diffusion models, particularly as model dimensions and parameter complexities increase. Addressing this challenge requires developing streamlined adaptive techniques and efficient algorithms that can accommodate the growing demand for tailored applications without substantially increasing the computational burden [7].

In conclusion, ongoing innovations in generative diffusion models have markedly advanced their adaptability and customization across diverse fields and use cases. These advancements are crucial for the growing trend toward personalized and context-sensitive generative AI applications. As research proceeds, tackling current computational challenges and further enhancing model flexibility will be essential to realizing the full potential of diffusion models in dynamic and multifaceted environments. Scholars and practitioners must remain vigilant regarding emerging methodologies and frameworks designed to enhance diffusion model adaptability, leveraging these insights to drive further innovation and application in increasingly complex and specialized domains.

## 4 DIVERSE APPLICATIONS ACROSS DOMAINS

### 4.1 Image and Video Generation

The application of generative diffusion models in image and video generation stands as a leading frontier in the domain of artificial intelligence, characterized by their capability to generate high-fidelity, realistic visual content across various applications. This subsection examines the intricacies of diffusion models in synthesizing, enhancing, and manipulating images and videos, showcasing significant advancements and persistent challenges within the field.

Generative diffusion models, built upon principles of sequential data manipulation through noise addition and removal, have become pivotal to addressing complexities inherent in visual data generation tasks. The foundational mechanism involves gradually corrupting visual data with Gaussian noise—a forward diffusion process—followed by a reverse diffusion process that employs neural networks to denoise and reconstruct the data, thereby generating new samples [4]. This simplistic yet computationally intensive framework allows for the precise handling of visual features, resulting in outputs that are not only of high resolution but also rich in texture and color fidelity.

Recent advancements have focused on improving the efficiency and quality of image synthesis. The introduction of design spaces that distinguish specific training and sampling strategies has led to substantial improvements in both speed and fidelity of generated images [6]. By implementing optimized score networks and revising core algorithmic processes, researchers have achieved new benchmarks in performance metrics such as the Fréchet Inception Distance (FID), which measures similarity between generated images and real-world counterparts [2].

In video generation, one of the foremost challenges involves maintaining temporal coherence across frames while preserving spatial details. Diffusion models have been adeptly extended to manage temporal dynamics by aligning the generation of frames with corresponding score networks

trained on spatial-temporal datasets [58]. This allows for the production of videos where visual coherence is sustained throughout sequences, effectively bridging the gap between static image synthesis and dynamic video content creation.

A critical technique that emerges within this context is multi-speed diffusion. By adjusting noise schedules and leveraging frame-specific information, multi-speed diffusion provides a framework for faster convergence and smoother transitions between frames [5]. The ability to modulate the diffusion process according to the temporal nature of data showcases a strategic evolution from static image-focused methodologies toward more holistic, dynamic approaches.

Beyond synthesis, diffusion models have made considerable strides in enhancement and manipulation tasks, including image inpainting, super-resolution, and style transfer. For instance, when tasked with enhancing image resolution, diffusion models employ high-dimensional latent spaces to iteratively refine and upscale visual content—an application with significant implications for fields such as medical imaging and high-definition broadcasting [4]. Image manipulation tasks are particularly aided by the model's inherent flexibility and the increasing incorporation of guidance techniques that allow for conditional generation with predefined stylistic or structural constraints [59].

While the capabilities of diffusion models in generating and manipulating visual content are undeniable, several challenges persist. Chief among these is the high computational overhead associated with their iterative nature. Each forward and reverse pass demands substantial resources, which can be a barrier to real-time applications or deployments in low-resource environments [7]. Efforts to mitigate these constraints focus on sampling efficiency, with techniques such as stochastic sampling and parallel processing frameworks being actively explored to reduce operational costs [6].

Furthermore, as the models become more sophisticated, issues related to bias and fairness in the generated content have come to light. The models are invariably influenced by the datasets they are trained on, leading to potential propagation of biases present in training data [60]. Addressing these concerns necessitates rigorous evaluation metrics and ethical guidelines to ensure that the generated content meets diverse user requirements without perpetuating existing societal biases.

Looking towards the future, several promising directions for diffusion models in image and video generation are emerging. First, there is a growing interest in hybrid models that incorporate other generative frameworks and adversarial training techniques to enhance robustness and model diversity [8]. Second, the exploration of domain-specific models tailored to niche applications, such as scientific visualization or cultural preservation, underscores the adaptability and expanding applicability of diffusion models [61].

In conclusion, generative diffusion models signify a transformative leap in the field of image and video generation, balancing the creation of visually arresting content with the computational complexities posed by their underlying processes. As research continues to refine these models, addressing challenges of efficiency, bias, and scal-

ability, the potential for these models to redefine digital visual content generation is immense, heralding a new era of generative artificial intelligence where art and technology converge seamlessly.

## 4.2 Text and Natural Language Processing

In recent years, diffusion models have established a significant presence in the realm of natural language processing (NLP), marking a new era for text generation and enhancement. This subsection explores the adaptation of diffusion models to the intricacies of NLP tasks, delving into their integration with established language model architectures, such as transformers. We assess how these models enhance text generation capabilities, maintain coherence, improve stylistic attributes, and transform structured data into natural language text. Our goal is to provide a comprehensive overview of state-of-the-art approaches, identify emerging trends and challenges, and propose future research directions and applications.

Originally designed for continuous data generation, diffusion models have been adapted for the discrete nature of textual data, a transformation necessitated by the unique sequential dependencies of language. This adaptation is nontrivial, requiring innovations in how noise is introduced and removed from discrete elements like characters or tokens. Techniques such as Categorical SDEs with Simplex Diffusion [62] have been developed, diffusing categorical data over a probability simplex rather than through Gaussian processes to ensure semantic and syntactic fidelity.

A critical aspect of diffusion models' application in NLP is their integration with transformers, which enhances their expressiveness and flexibility, allowing them to leverage pre-learned language semantics. For instance, prompt-tuning latent diffusion models [63] provide valuable insights into guiding the diffusion process, aligning generated text with contextual and stylistic nuances for coherent and meaningful narratives.

Diffusion models show advantages over conventional models in text generation, particularly in coherence. Through iterative noise addition and removal—similar to methods used in physics-based modeling—language constructs can be refined. Neural Diffusion Processes [64], for example, incorporate exchangeability properties into model architectures, enabling text generation from learned distributions over finite marginals that closely resemble real textual data. This capability positions diffusion models to generate text that maintains coherence over long sequences, addressing challenges traditional sequence models like LSTMs face.

Beyond generation, diffusion models extend to the typographic and stylistic enhancement of text, marrying visual aesthetics with computational linguistics. This involves transforming text with features for font specification, alignment, and other stylistic applications, enhancing readability and engagement. Methods like Iterative $a$-(de)Blending [22] show how language structure can be stylized without sacrificing meaning or coherence.

In addition, diffusion models are valuable for converting structured data into text, crucial for tasks like report generation and summary writing. Approaches such as the PaDIS

model [65] highlight how patch-based diffusion models can learn efficient data priors, translating into NLP applications for generating narratives from structured inputs like tables or databases, ensuring coherence and completeness of the conveyed information.

Despite these advancements, challenges and limitations persist. The computational cost of diffusion processes, due to iterative noise operations, presents scalability issues for large-scale NLP applications. Adapting continuous processes to discrete data involves approximation methods that can introduce errors or biases. Theoretical work, such as Convergence of denoising diffusion models under the manifold hypothesis [66], provides foundational insights into managing these errors through convergence guarantees.

Ethical considerations remain a concern, as diffusion models can perpetuate biases from training datasets or generate misinformation. Efforts like EraseDiff [67] aim to mitigate these risks, but a comprehensive approach is needed as these models gain wider application.

Future research should focus on refining model architectures to enhance scalability and efficiency, potentially leveraging hybrid architectures combining energy-based models, GANs, and diffusion models [3]. Increasing interpretability and transparency can address ethical concerns. The integration with real-time text generation for interactive AI systems offers an exciting frontier in NLP research.

In conclusion, the integration of diffusion models into natural language processing signifies a paradigm shift with profound implications for text generation and enhancement. As research progresses, these models are poised to revolutionize machine understanding and generation of human language, promising more intuitive and impactful applications across domains like literature, education, and journalism.

### 4.3 Scientific and Medical Applications

Generative diffusion models have emerged as a transformational tool in scientific and medical domains, offering novel approaches to simulation, imaging, and problem-solving that extend beyond the capabilities of traditional methodologies. In this subsection, we explore the multifaceted applications of these models, examining their impact on medical imaging, scientific simulations, and multimodal data integration. We also assess their strengths, limitations, and future prospects within these high-stakes fields, providing a comprehensive analysis supported by academic literature.

At the forefront of medical applications, diffusion models are progressively being utilized for medical image synthesis and segmentation. The intrinsic ability of diffusion models to learn complex data distributions enables enhanced generation of realistic medical images, essential for training diagnostic algorithms and facilitating clinical decision-making. Research in accelerated MRI using score-based diffusion models exemplifies this application, addressing the challenge of reconstructing high-quality images from undersampled data [68]. This capability is paramount in clinical settings where rapid acquisition of high-resolution images is critical but often constrained by time and budgetary considerations.

Furthermore, the synthesis of synthetic medical images through diffusion models not only aids in augmenting training datasets but also in simulating various medical conditions for research purposes. This synthetic data generation potential is crucial in rare diseases where data scarcity impedes model training and validation. Moreover, diffusion models facilitate robust unsupervised anomaly detection by highlighting discrepancies between generated and actual images. For example, by generating baseline anatomical models, deviations can be assessed and potentially indicative of pathologies, thus enhancing diagnostic accuracy.

Scientific simulations are another domain where generative diffusion models have proven instrumental. They are particularly adept at creating simulations in climate modeling, material science, and other fields requiring the synthesis of complex temporal and spatial patterns. For instance, diffusion models enable the generation of realistic climate data sets under various hypothetical scenarios, thus assisting in disaster preparedness and environmental policymaking [1]. The ability to incorporate stochastic processes akin to real-world phenomena makes diffusion models a potent tool in scientific inquiry, offering nuanced insights into systems otherwise intractable by deterministic models.

In terms of complex problem-solving, diffusion models offer solutions to inverse problems frequently encountered in scientific imaging and medical diagnosis. By iteratively refining estimates of data parameters, diffusion models can reconstruct images from incomplete or noisy measurements, a task crucial in fields like computed tomography and MRI. Techniques such as those illustrated by Grad-TTS employ score-based frameworks to address inverse problems, proving beneficial in practical applications ranging from image inpainting to colorization tasks [69].

An area of burgeoning interest in the medical field is brain imaging and multimodal medical data synthesis. Diffusion models offer a framework for integrating data from multiple medical imaging modalities, such as MRI, CT, and PET scans, leading to richer, more comprehensive diagnostic information. By learning the joint distribution of these diverse data types, diffusion models can generate integrative diagnostics that enhance the interpretation and accuracy of medical assessments. This capability is invaluable in complex pathological conditions where a single imaging modality might fall short of providing a holistic view.

Despite their numerous advantages, generative diffusion models also face certain limitations and challenges. One significant concern is the computational cost associated with their iterative sampling processes, which can be computationally intensive, especially when scaling to high-resolution images or extensive data sets [31]. Addressing this challenge requires advancements in algorithmic efficiency and the development of more efficient solvers capable of preserving fidelity while enhancing the speed of convergence.

Moreover, ethical considerations such as data privacy and bias also loom large in deploying diffusion models, particularly in sensitive fields such as medical diagnostics. The potential for these models to memorize and inadvertently replicate training data underscores the importance of developing privacy-preserving techniques that safeguard patient confidentiality. Techniques like membership inference and anonymization must be refined to bolster trust and

compliance in medical applications [70].

Looking ahead, the integration of diffusion models with other advanced technologies, such as reinforcement learning and neural architecture search, holds promise for expanding their applicability and performance. Reinforcement learning, for example, can optimize the training process of diffusion models, tailoring their outputs to achieve desired objectives, whether in scientific simulations or medical diagnoses [71]. Furthermore, the application of diffusion models in high-dimensional and structured data spaces is expected to open new frontiers in complex domains requiring precision and tailored generative capabilities.

In summary, generative diffusion models are set to make substantial contributions to scientific and medical fields, fostering innovations in simulation, imaging, and analytical processes. By addressing current challenges and integrating emerging technologies, these models have the potential to transform how we approach and solve some of the most complex problems in these domains, paving the way for groundbreaking advancements in scientific research and healthcare.

## 4.4 Interdisciplinary and Emerging Domains

In recent years, generative diffusion models have increasingly found applications across interdisciplinary domains, demonstrating their versatility and innovative potential in addressing complex, multi-faceted challenges. This subsection delves into the diverse and emerging uses of diffusion models in fields that merge different data modalities or tackle novel interdisciplinary problems, offering a bridge from the groundbreaking successes in scientific and medical applications to wider explorations beyond traditional boundaries.

At the core of these interdisciplinary applications is the cross-modal generation capability of diffusion models, which excels in synthesizing coherent outputs from multiple forms of data. For instance, the integration of audio-visual information is paramount in applications such as generating video from textual descriptions, enabling the creation of new forms of multimedia content [72]. By embedding data from different modalities into a shared latent space, diffusion models facilitate seamless transformations, leveraging their intrinsic probabilistic nature to maintain consistency across outputs [73]. The aptitude for handling complex data forms not only extends the scope of generative tasks but also enhances the richness and diversity of synthesized outputs, marking a significant advancement over more traditional single-modal approaches.

In the realm of generative design and creativity, diffusion models are significantly impacting how artists and designers explore creative processes. They provide novel tools for digital art creation, allowing the fine-tuning of artistic elements through iterative processes informed by learned patterns across diverse styles and mediums [57]. By capturing the underlying statistical properties of various artistic styles, diffusion models can produce intermediate transformations that retain stylistic coherence while introducing novel elements, thus pushing creative boundaries. This capability transforms how digital content is conceived and produced, utilizing diffusion models' ability to model high-dimensional data spaces and translate conceptual inputs into visually coherent artistic expressions.

Beyond the arts and multimedia, diffusion models play an increasingly critical role in environmental monitoring and remote sensing, illustrating their adaptability in fields that require integration across temporal and spatial dimensions. The processing of remote sensing data is enhanced by diffusion models' adaptive nature, allowing them to deliver high-resolution environmental insights [74]. Their prowess in generating simulations and predictive models supports understanding environmental changes and informing decision-making in urban planning and climate impact assessments. By generating synthetic datasets that augment limited real-world data, diffusion models contribute to robust forecasting models, enhancing analytic precision and reliability.

The technical integration of diffusion models into these interdisciplinary domains often involves overcoming challenges related to computational complexity and data heterogeneity. The design and implementation of scalable architectures, as discussed in works such as [52], are essential to manage the extensive data flow and computational demands inherent in multimodal synthesis. Additionally, combining disparate data sources necessitates advanced algorithmic strategies, such as hierarchical or modular architectures, to effectively manage data inference and integration. These considerations are crucial in ensuring diffusion models' robust performance across applications, demanding profound depth and breadth in data modeling capabilities [75].

Emerging trends signal a growing interest in leveraging diffusion models' flexibility to support tasks requiring high adaptivity, such as personalized content generation and adaptive design frameworks. As these models evolve, enhanced learning efficiency and inferential accuracy are likely to drive further integration into domains like bioinformatics and computational biology, where precision and adaptability are critical [54]. These models are poised to address critical challenges in protein interaction modeling and drug discovery, showcasing their potential in solving highly specialized tasks with significant implications for healthcare and biological sciences.

Looking forward, the fusion of diffusion models with other generative frameworks, such as Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs), holds promise for creating richer, more robust generative systems that leverage the strengths of each approach [76]. This hybridization could enhance model capabilities, such as improving convergence rates and increasing sample diversity while mitigating individual weaknesses like mode collapse or sample variability. Additionally, enhancing the interpretability of diffusion models remains a key area of research, critical for fostering trust and facilitating adoption in sensitive domains like healthcare and environmental science [73].

In conclusion, the interdisciplinary applications of diffusion models represent a burgeoning field of research with extensive implications, seamlessly extending from their foundational roles in medical and scientific arenas. By successfully merging various data modalities and addressing complex challenges, these models are paving the way for

breakthroughs across a diverse array of sectors. Continued innovation in model design, optimization strategies, and domain-specific adaptations are essential for unlocking diffusion models' full potential, driving their advancement as integral tools in the interdisciplinary toolkit. As research progresses, exploring the ethical and societal impacts of these technologies will be vital to guide their responsible and beneficial deployment across various domains.

## 5 INTEGRATION WITH OTHER GENERATIVE FRAMEWORKS

### 5.1 Synergistic Integration with Adversarial and Autoencoder Models

The integration of diffusion models with Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs) offers a promising avenue for enhancing the capabilities and performance of generative models. This subsection explores the synergistic potential of these combinations, focusing on strategies that leverage their unique advantages while addressing inherent limitations, thereby advancing the state-of-the-art in generative modeling.

Diffusion models have emerged as a robust framework for generating high-quality samples through a progressive denoising process. Their strength lies in their ability to model complex data distributions without suffering from mode collapse, a common challenge in GANs. On the other hand, GANs excel in generating sharp and realistic images due to their adversarial training objectives [13]. However, GANs often struggle with training instability and mode dropping. VAEs bring to the table a principled framework for latent space representation and probabilistic inference, providing robust generation and reconstruction capabilities. Nevertheless, VAEs often produce blurrier outputs due to their reliance on Gaussian assumptions in the latent space [1].

A primary strategy for integrating diffusion models with GANs focuses on using diffusion processes to enhance GAN training dynamics. By incorporating the progressive noise-removal strategy of diffusion models into the GAN framework, one can achieve more stable, diverse, and high-quality image generation. This integration can mitigate mode collapse in GANs by leveraging the diversity-preserving properties of diffusion processes. For instance, diffusion components can be integrated as a regularizing mechanism in GANs, allowing the discriminator to better differentiate between diverse plausible samples and reducing the generator's tendency for mode collapse [6].

Moreover, combining VAEs with diffusion models can address some of the former's limitations in sample diversity and quality. VAEs can benefit from the superior reconstruction capabilities of diffusion models by adopting their denoising strategies to refine latent space sampling. This hybrid model can improve the sharpness of generated images and provide more expressive latent features, enhancing both reconstruction and generative tasks. The joint modeling of latent spaces using diffusion-enhanced VAEs also holds promise for better exploration of data manifolds, which can lead to superior data representation and generation [58].

The integration of these models can be mathematically described using variational inference techniques. Suppose $\mathbf{x}$ denotes the data, $\mathbf{z}$ the latent variable in VAEs, and $\mathbf{y}$ represents the noise in diffusion models. The goal is to optimize the joint likelihood $p(\mathbf{x}, \mathbf{z}, \mathbf{y})$ using an extended Evidence Lower Bound (ELBO) that incorporates terms from both VAE and diffusion model training. Such formulations enable the exploitation of structured priors and informed sampling strategies across diverse latent spaces, which in turn enhances the expressiveness of the generative model.

Emerging trends in this synergy include the development of training algorithms that exploit the strengths of individual components in the integrated model. For instance, curriculum learning techniques can be adapted to schedule the training phases of GAN, VAE, and diffusion components, thereby optimizing convergence dynamics and sample diversity [47]. Another innovative approach involves the use of pretrained diffusion models as plug-and-play modules for finetuning GANs or VAEs on specific datasets or tasks, by leveraging their superior generalization properties [59].

As we advance, the integration of GANs, VAEs, and diffusion models could lead to numerous practical implications, especially in applications demanding high fidelity and diverse outputs, such as image editing, scientific simulations, and content generation. However, addressing challenges related to computational complexity, training stability, and compatibility of different model architectures remains critical [4].

In conclusion, the synergistic integration of diffusion models with GANs and VAEs promises to revolutionize generative modeling by combining their complementary strengths. Moving forward, research should focus on developing efficient training paradigms, exploring novel integration techniques, and addressing inherent challenges of computational scalability and model integration. Such efforts will likely broaden the horizon of generative AI, unlocking unprecedented potential across diverse domains and applications.

### 5.2 Multimodal and Cross-Modal Generative Approaches

The integration of diffusion models with other generative architectures unlocks new potential for multimodal and cross-modal generation, enabling the creation of content that is not only diverse but also contextually enriched across multiple domains. This subsection delves into the mechanisms and innovations driving these integrations, offering a comparative analysis of methods while evaluating their strengths, weaknesses, emerging trends, and challenges.

Multimodal and cross-modal generative approaches involve the synthesis of outputs that integrate multiple types of data, such as combining text, images, and audio to produce rich and comprehensive representations. To facilitate this, diffusion models are often coupled with other generative frameworks like Variational Autoencoders (VAEs) and Generative Adversarial Networks (GANs), enhancing the model's capability to capture complex joint distributions [3].

The foundational mechanism of diffusion models operates through forward and reverse processes, whereby data is progressively noised and denoised. This aligns seamlessly

with score-based models that learn the gradient of the data distribution's log-density [17]. By integrating these diffusion processes with architectures representing different data modalities, models can learn joint representations capturing correlations across varied data types.

Cross-modal embeddings are pivotal in such frameworks. These embeddings establish a shared latent space where different modalities can align and interact, crucial for achieving coherent cross-modal synthesis. Utilizing approaches like contrastive learning, these embeddings can be effectively synchronized, enabling diffusion models to generate harmonized outputs consistent across data types. These techniques significantly boost the multimodal generation capabilities of diffusion models, leveraging their inherent ability to manage high-dimensional spaces [3].

The chief benefit of integrating diffusion models with other architectures lies in leveraging complementary strengths while mitigating individual weaknesses. GANs, for example, are known for producing high-quality outputs but often suffer from training instability and mode collapse. Diffusion models, with their progressive denoising approach, can provide stabilization, enhancing the quality of GAN-generated outputs in multimodal tasks [3]. Similarly, VAEs can capitalize on the refinement capabilities of diffusion models to enhance the fidelity and precision of generated outputs.

Technical challenges in these integrations center on efficiently managing and learning from the vast and heterogeneous data spaces inherent in multimodal tasks. For diffusion models, effectively capturing multimodal correlations during the reverse process is crucial, necessitating careful design of noise schedulers and score estimators to accommodate the characteristics of each modality. Insights from information theory and optimal transport can offer frameworks to better manage these challenges [77].

Emerging trends in multimodal diffusion frameworks include exploring latent variable models tailored to the unique attributes of multimodal datasets. The use of graph-based structural representations, which capture intricate relationships between modalities, is gaining traction. These representations allow diffusion models to effectively leverage spatial and temporal dependencies, enriching the generative process [27].

Furthermore, developing invertible architectures and efficient sampling strategies is critical to enhancing the practicality and scalability of diffusion models in multimodal settings. As discussed in [21], incorporating parallel sampling methods and scalability solutions is essential for these models to be viable in real-world scenarios where computational resources and time are constrained.

The advancements in multimodal and cross-modal generative approaches underscore the importance of robust, flexible frameworks capable of handling diverse and intricate data inputs. Although significant progress has been made, challenges remain in ensuring the scalability and efficiency of these integrated models, especially when dealing with the high-dimensional data spaces common in domains like audiovisual media and interactive applications.

Future directions are likely to focus on improving the efficiency of cross-modal embeddings, refining joint training strategies, and optimizing multimodal score estimation

techniques. Enhancements in computational techniques, such as distributed training and adaptive learning mechanisms, will be pivotal in advancing diffusion models' capabilities in these contexts [78]. Further exploration into leveraging domain-specific knowledge for modeling multimodal data could provide novel insights and practical solutions to current limitations.

In synthesizing insights from multiple studies and approaches, it is evident that the integration of diffusion models with other generative architectures holds substantial promise for advancing multimodal generative modeling. By continuing to refine these integrative frameworks, researchers can unlock new possibilities for creative and functional multimodal outputs, with far-reaching implications in diverse sectors such as entertainment, education, and scientific visualization.

## 5.3 Integrative Frameworks for Enhanced Generative Capabilities

In the landscape of generative modeling, the development of integrative frameworks that synergize diffusion models with other generative mechanisms is a burgeoning area of research, offering enhanced capabilities and novel applications across diverse domains. This subsection delves into these integrative frameworks, aiming to elucidate how diffusion models can be augmented via collaborative frameworks with other generative architectures, such as Variational Autoencoders (VAEs) and Generative Adversarial Networks (GANs). The goal is to explore these integrative strategies in terms of their technical methodologies, strengths, limitations, and prospective future directions while providing a comparative analysis and synthesis of the current state of this innovative field.

In essence, integrative frameworks strive to synthesize the mathematical elegance and practical efficiency of various generative modeling paradigms. Diffusion models traditionally follow a stochastic differential equation-based approach that methodically reverses noise into coherent structures, capitalizing on their robustness to generate high-fidelity outputs [17]. However, this methodology often incurs high computational overheads due to the sequence of iterative inversions required. By integrating with VAEs and GANs, these models can leverage the latent space exploration and adversarial training advantages, respectively.

One notable integration is the coupling of diffusion models with VAEs, which aims to optimize latent space utilization. This union facilitates efficient latent exploration by employing VAEs to encode data into a compact space while relying on diffusion models to decode and refine the samples progressively, enhancing output precision. This dual-process approach serves to mitigate the shortcomings of VAEs, such as an inability to capture fine-grained details, while simultaneously addressing the computational intensity typical of standalone diffusion models. Despite the clear benefits, challenges remain, particularly regarding the optimal balance between computational efficiency and the preservation of intricate data characteristics. The integration often demands substantial tuning of the latent structure and its interaction with the diffusion-based refinement process, which could potentially complicate model training and limit scalability.

A further integration explored is the amalgamation of diffusion models with GANs. This hybrid model seeks to take advantage of the adversarial approach's propensity for generating sharp and realistic outputs while employing diffusion processes to stabilize training and improve diversity [79]. GANs' capability to discern and emphasize minute details through discriminator feedback is augmented by the denoising mechanism in diffusion models, leading to outputs with enhanced realism and reduced artifacts. However, merging these paradigms involves intricate adjustment of the diffusion process's granularity to fit the adversarial game dynamically. Moreover, while diffusion models can enrich the diversity of GAN-generated content, the inherent mode-seeking behavior of GANs could occasionally counteract the diffusion model's aim for output variance maximization, thus requiring careful calibration [80].

Integration in hybrid frameworks extends into task-specific model enhancements and composability. There is an increasing trend in using diffusion models as foundational structures upon which task-oriented models are built, enabling specialized generators for particular applications like text-to-speech or conditional image synthesis. Such frameworks allow practitioners to design sophisticated generative pipelines that maintain general-purpose capabilities while being fine-tunable to address specific requirements. The success of these systems depends heavily on the seamless combination of different model architectures and their respective learning paradigms. Recent advances in task-specific integrations have shown that score-based methods are particularly adaptable, contributing to more tailored and accurate data synthesis in stressful domains, such as medical imaging and scientific simulations [68], [69].

From a broader perspective, these integrated frameworks represent a pivotal paradigm shift, bridging the theoretical constructs of diffusion processes with practical requirements across various domains. This synthesis not only fosters the enhancement of generation quality and adaptability but also invites new research into optimized algorithmic synthesis, improved training regimes, and enhanced model robustness.

Future directions in integrative frameworks point toward the development of more automated and adaptive systems that can dynamically select and integrate generative strategies based on task-specific requirements. Such adaptability would likely involve incorporating machine-learning techniques, such as reinforcement learning, to improve integration efficiency and reduce manual tuning. Moreover, there exists a burgeoning interest in exploring the convergence properties and theoretical guarantees of these hybrid models to ascertain their foundational stability and optimal performance across varying conditions [21].

In conclusion, the emerging field of integrative frameworks in generative modeling facilitates an expansion in both the capability and application scope of diffusion models. By effectively leveraging the complementary strengths of other generative architectures, researchers are charting a path toward more sophisticated, robust, and flexible generative systems. As the community continues to innovate and refine these integrative paradigms, the potential for more nuanced applications and deeper theoretical insights will likely propel the field into increasingly exciting areas of research and application.

## 5.4 Unified Generative Strategies using Diffusion Models

Within the realm of generative models, diffusion models have emerged as potent frameworks capable of generating high-quality synthetic data across diverse domains such as image, audio, and text. This subsection aims to articulate the unified generative strategies that incorporate diffusion models alongside other generative architectures, such as Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs), to establish comprehensive generative systems. By exploiting the synergistic potential of these models, this integrative approach seeks to push the boundaries of generative capabilities and addresses existing challenges.

The integration of diffusion models with other generative architectures leverages their complementary strengths. Diffusion models offer stability in training and robustness against mode collapse due to their progressive denoising process. On the other hand, GANs are distinguished by their ability to produce sharp and high-resolution images through adversarial training mechanisms [76]. In contrast, VAEs provide an efficient mapping of data to a latent space, facilitating data representation and enabling systematic exploration of this space, albeit often at the expense of generating fuzzier images. Unified strategies that capitalize on the strengths and ameliorate the weaknesses of each model can significantly enhance generative efficiency, diversity, and realism [73], [81].

A notable unified strategy involves the coupling of diffusion models with GANs to form hybrid architectures. These architectures leverage adversarial training for detailed synthesis while employing diffusion processes for structured noise reduction and denoising. This approach mitigates the challenges of mode collapse commonly associated with GANs, thereby enhancing the diversity and quality of generated outputs. Recent research elucidates that the diffusion component can regulate and guide the adversarial training process by providing structured noise patterns, fostering more stable learning dynamics and allowing for finer navigation of the generative process [11]. This synergy opens avenues for achieving outputs with both high fidelity and detail, broadening applicability across more demanding generative tasks.

Another promising integration involves the conjugation of diffusion models with VAEs. This hybridization exploits the efficient latent space representation characteristic of VAEs, combined with the robust generative capacity of diffusion models for data reconstruction and novel data generation. This coupling is effective in refining the latent space through enhanced variance estimation inherent in diffusion processes, aiding in generating more stable and consistent outputs [57]. Moreover, the architectural synergy between VAEs and diffusion models allows for the correction of modeling inaccuracies in the latent space via diffusion-induced stochastic sampling, ensuring more reliable generative outcomes that maintain high fidelity to the target distribution [82].

Emerging trends have seen the application of unified strategies in developing models that are robust and ca-

pable of operating efficiently across various domains and modalities. By aligning multimodal data such as text, image, and audio within shared latent spaces, these strategies enable comprehensive cross-modal generative capabilities—essential for tasks requiring contextual coherence and high fidelity across modes. This holistic approach is exemplified in models employing diffusion processes to facilitate seamless and logical transitions between diverse data types, thus enhancing the scope of multimodal synthesis [3].

Despite the potential of these unified strategies, challenges persist. The complexity involved in effectively marrying different generative frameworks without inducing instability or compromising computational efficiency is significant. Therefore, ongoing research focuses on refining integration methodologies to ensure scalable and resource-efficient implementations, especially when scaling to higher-dimensional data spaces and more intricate data distributions [83], [84]. Furthermore, while unified models enhance generative diversity and coherence, they emphasize addressing ethical concerns tied to data privacy and bias inherent in model training data [85], [86].

In conclusion, the integrative use of diffusion models with other generative frameworks represents a significant advancement in generative modeling. These unified strategies not only enhance the robustness and diversity of generative outputs but also broaden the applicability of these models across various domains. Future research endeavors will focus on refining these integrations through novel architectural designs and optimization algorithms that effectively balance computational demands with generative performance. With continued progress, these unified generative strategies hold the promise of significantly transforming generative modeling, providing tools capable of tackling increasingly complex and multi-faceted generative challenges.

## 6 EVALUATION METRICS AND BENCHMARKING

### 6.1 Performance Metrics for Generative Diffusion Models

The performance evaluation of generative diffusion models is a critical aspect of their development and application, aimed at ensuring high-quality, realistic, and diverse outputs. As these models grow in prominence across various domains, the need for robust and comprehensive metrics becomes increasingly essential. This subsection delves into the multi-faceted metrics that assess the efficacy of generative diffusion models, examining their strengths, limitations, and emerging trends, while also offering insights into future research directions.

To begin with, the core metrics for evaluating the *quality* of outputs generated by diffusion models primarily include the Fréchet Inception Distance (FID) and Inception Score (IS). The FID metric measures the distance between the feature representations of real and generated data, calculated using a pre-trained Inception network. This metric is valued for its sensitivity to both the quality and diversity of generated data, with a lower FID indicating closer alignment with the real distribution [6]. Meanwhile, the IS evaluates the entropy of the predicted class probabilities for generated images, rewarding models that produce highly diverse images among distinct classes while maintaining coherence within each class [2]. However, while IS provides insight into intra-class variety, it may overlook certain quality nuances when compared to FID. Despite this, both metrics play complementary roles in assessing generative performance and are widely adopted for benchmarking in various applications [12].

When considering the *realism* of generated outputs, newer metrics such as the Image Realism Score (IRS) have emerged, which aim to directly quantify the perceptual quality of individual images. This kind of assessment is crucial for applications like image and video generation, where maintaining the visual realism is often paramount [4]. Unlike IS and FID, the IRS focuses on perceptual fidelity rather than statistical comparison to large datasets, thus offering an alternative lens through which to evaluate model outputs. The perceived realism, especially in applications like high-fidelity image synthesis and video generation, demands metrics that consider not only pixel-level accuracy but also the semantic integrity of the images produced [13].

In terms of *diversity*, it is critical to avoid issues such as mode collapse, where the model generates a limited variety of outputs. To combat this, metrics that assess the variability and breadth of generative outputs are crucial. For instance, a set of metrics inspired by coverage and precision in point cloud generation can offer more detailed insights into the diversity and completeness of outputs by distinguishing between underrepresented regions in the generated data and those that the model captures accurately [1]. By analyzing these two dimensions, researchers and practitioners can better understand a model's capacity to produce a wide range of realistic and meaningful outputs, which is essential for generalizing across diverse applications and datasets.

Several challenges persist in evaluating these metrics, particularly in aligning them with human judgments of quality and diversity. Studies have indicated discrepancies between algorithmic scores and human perceptual assessments, which highlight the need for developing more nuanced, application-specific evaluation criteria [61]. For instance, while FID and IS are robust for image-based evaluations, the domain of text-to-image generation may require additional considerations to ensure that the semantic coherence and intended meaning of the generated images meet expectations [8].

Looking forward, the intersection of diffusion models with reinforcement learning and interactive learning presents promising opportunities for evolving performance metrics that not only assess static outputs but also feedback and adaptively improve model performance over time [47]. Moreover, the integration of generative models with adversarial and autoencoder frameworks has led to emerging paradigms for evaluating hybrid model architectures, which blend the strengths of different generative approaches to achieve superior outcomes [87]. These advancements call for new metrics that capture the synthesized dimensions of performance, particularly when assessing models involved in multi-stage or conditional generation tasks [88].

In conclusion, the performance metrics for generative diffusion models must evolve to keep pace with the rapid advancements in the field. It is essential to augment traditional metrics with those tailored to specific applications,

and address emerging challenges such as aligning algorithmic evaluations with human judgments, and integrating adaptive learning mechanisms. The development of a comprehensive suite of performance metrics, grounded in both statistical rigor and practical relevance, will be pivotal for the continued success and adoption of generative diffusion models across diverse domains. With ongoing research and community input, it is possible to refine these evaluation strategies to better reflect the multidimensional nature of generative tasks and their outputs, thus steering the next wave of innovation in this exciting area.

## 6.2 Datasets and Benchmarking Protocols

Evaluating generative diffusion models involves rigorous assessments utilizing standard datasets and benchmarking protocols. This subsection delves into the various datasets and benchmarks central to evaluating these models, highlighting their crucial role in promoting consistent assessment and comparison across research studies. The discussion will navigate through the intricacies of dataset selection, benchmarking practices, and methodologies, elucidating their significance in crafting robust evaluation frameworks for generative diffusion models.

Central to the evaluation process, datasets underpin the assessment of performance and generalization capabilities of generative diffusion models. Prominent datasets often serve as standard benchmarks for experimental comparisons. For instance, CIFAR-10 and ImageNet are widely recognized standards due to their extensive and varied image collections, enabling researchers to assess a model's effectiveness in capturing complex image features [3], [4]. CIFAR-10 is appreciated for its manageable scale, facilitating the testing of computationally demanding models, while ImageNet offers a broader array of categories that test a model's scalability and capability in handling more extensive and semantically complex images.

Moreover, the landscape of datasets is expanding beyond traditional visual data to accommodate diverse domains. Datasets like CelebA provide an abundant selection of facial images, serving as a fertile ground for testing models on tasks that require nuanced understanding of human features. Similarly, LSUN, with its varied categories, provides a specific setting for deeper exploration into contextual image synthesis tasks [89]. These datasets not only push innovation in model architectures and training methodologies but also challenge existing models to keep pace with evolving data intricacies.

Benchmark suites significantly contribute to this evaluative ecosystem by offering structured methodologies for assessing the generative capabilities of diffusion models. The Generative-Discriminative Evaluation Benchmark (GDBench) is notable for facilitating a more comprehensive understanding of a model's capabilities beyond mere generative quality to encompass compositional and discriminative attributes [21]. Such benchmarks spur the development of holistic evaluation strategies that transcend traditional metrics.

Benchmarking processes are further supplemented by protocols that guide evaluation procedures to ensure consistency and dependability in performance assessments across studies. These protocols commonly define particular conditions under which models are trained and assessed, including noise levels, data pre-processing steps, and model architecture constraints. Protocols might, for example, prescribe adversarial conditions or utilize out-of-distribution data scenarios to evaluate model robustness and generalization. Ensuring adherence to these predefined protocols not only bolsters the reliability of outcomes but also enables comparative analysis across disparate research initiatives.

Through comparative analysis, it becomes evident that while the widespread adoption of CIFAR-10 and ImageNet sets a sturdy benchmark for comparability, it may inadvertently narrow innovation by steering research efforts towards specific datasets. Embracing more diverse datasets, particularly those from non-visual domains, can broaden the applicability of diffusion models to innovative fields such as medical imaging or remote sensing. Traditional benchmarks, often emphasizing static metrics, could benefit from integrating dynamic evaluations like sample diversity and temporal coherence, especially pertinent in tasks such as video synthesis.

Emerging trends within the field reflect a burgeoning interest in overcoming these limitations by leveraging domain-specific datasets that tackle the unique challenges of various applications. Temporal sequence datasets, increasingly utilized to scrutinize models in time-sensitive applications like video generation and time series forecasting, exemplify this trend. These domain-specific benchmarks furnish intricate insights into a model's adaptability to diverse data complexities and temporal nuances.

Looking forward, the evolution of datasets and benchmarking protocols for generative diffusion models will hinge on fostering inclusivity and adaptability across diverse domains. This evolution demands the curation of datasets that align with emerging applications and the development of benchmarking practices that emphasize multidimensional evaluation metrics. Collaborations between domain experts and diffusion model researchers can yield datasets that better encapsulate the nuances of specialized fields. Moreover, refining evaluation protocols to integrate ethical considerations, like bias detection and fairness evaluation, ensures that diffusion models are deployed responsibly and equitably.

In conclusion, datasets and benchmarking protocols are indispensable for the evaluation and advancement of generative diffusion models. As the field evolves, expanding datasets to span diverse domains and refining benchmarking practices to encapsulate multidimensional evaluation metrics will be crucial. Addressing challenges and capitalizing on emerging opportunities ensures that the evaluation frameworks for diffusion models remain robust, reliable, and reflective of real-world complexities, further enhancing their integration into diverse and challenging applications.

## 6.3 Methodologies for Evaluating Robustness and Bias

Generative diffusion models have emerged as a cornerstone in the domain of generative artificial intelligence, exhibiting groundbreaking performances across a range of applications, particularly in image and audio synthesis [3]. However, as these models transition from theoretical constructs

to practical applications, questions regarding their robustness and inherent biases have gained considerable attention. This subsection aims to delve into these issues, offering an analytical framework for evaluating the robustness and bias of generative diffusion models.

To commence, the notion of robustness in diffusion models pertains to their ability to maintain performance across a variety of conditions that deviate from the original training distribution. This includes scenarios such as adversarial attacks or shifts in input data distributions. The assessment of robustness is critical, as generative models are increasingly being employed in high-stakes environments where reliability is paramount. Approaches like adversarial testing, where models are exposed to systematically perturbed inputs to quantify their resilience, constitute a significant area of focus. The methodology of adversarial testing aligns with the broader landscape of score-based generative modeling, which relies on the score function's accuracy; any deviation from this can significantly affect the model's stability and robustness [80].

Diffusion models also face the challenge of bias, where the generated outputs may reflect or even amplify societal stereotypes present in the training data. Identifying and quantifying bias is non-trivial due to the complexity and high-dimensionality of diffusion processes. Emerging methodologies employ stereotypical bias assessments, which analyze generated samples to detect biases prevalent in the output. These assessments often involve comparing the generated data distributions with those of unbiased benchmarks, as highlighted in the research on unbiased diffusion modeling [90].

A pivotal approach to addressing bias involves the implementation of fairness-aware training protocols. These protocols aim to explicitly correct biases during the training phase by re-weighting the training data or adjusting the learning algorithms to minimize biased gradients. This is analogous to the reinforcement learning-based fine-tuning methods, which seek to optimize downstream reward functions, potentially steering diffusion processes towards equitable outcomes [71].

Beyond robustness and bias, the generalization capability of diffusion models underpins their potential to transfer learned properties to new, unseen tasks. Evaluating this capability involves not only measuring performance on out-of-sample data but also understanding how models can be conditioned to adapt to varying conditions [91]. Techniques such as consistency training, which enforces model outputs to remain consistent across time, have been shown to enhance model generalization by minimizing distribution drift [92].

Adding a layer of complexity is the evaluation of diffusion models in settings where the data lies on low-dimensional manifolds within the high-dimensional ambient space. Theoretical advances elucidate how models like DDPM adapt to such low-dimensional structures, crucial for high-quality generative outcomes in domains with inherent structural constraints, such as 3D shape modeling or molecular synthesis [24].

Incorporating adaptations of these models to categorical or structured data domains also highlights the adaptability of the methodology. The development of Simplex Diffusion processes for categorical data embeds the challenges of dealing with non-continuous data within the well-known framework of diffusion processes [62]. Understanding this flexibility is critical for expanding the applicability of diffusion models beyond standard continuous data formats.

In synthesizing the insights from these various methodologies, it is evident that fostering robustness and mitigating bias in diffusion models require a multifaceted approach that blends theoretical insights with practical algorithms. As research progresses, innovative methodologies such as stochastic differential equation-driven adaptations and gradient guidance [93] are likely to play pivotal roles in finetuning these models for greater fairness and stability.

Looking forward, the scope for enhancing the robustness and fairness of diffusion models is vast. Integration of more advanced adversarial learning techniques, coupled with the synthesis of insights from statistical learning theories such as optimal transport and probabilistic generative modeling, will underpin future progress. Moreover, as deployment of these models in real-world applications accelerates, ethical considerations will become increasingly central, calling for an interdisciplinary approach that bridges machine learning with fields such as social sciences and policy studies.

In conclusion, the methodologies for evaluating and enhancing the robustness and bias of generative diffusion models are at a nascent yet rapidly evolving stage. The convergence of theoretical advances, empirical insights, and algorithmic innovations is poised to expand the capability and reliability of diffusion models profoundly, promising a future where generative diffusion models not only produce high-quality outputs but do so with fairness and robustness as foundational tenets.

## 6.4 Cross-Model Comparative Analysis

In the landscape of deep generative models, such as Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs), and more recently, diffusion models, comparative analysis has become paramount to understanding their relative efficacy, complexity, robustness, and scalability. This subsection aims to provide a detailed comparative framework for evaluating generative diffusion models alongside GANs and VAEs, highlighting their strengths, limitations, emergent trends, and potential future directions.

Diffusion models have recently attracted significant attention due to their exceptional performance in generating high-fidelity images and their flexibility in model design, as reported in "Diffusion Models: A Comprehensive Survey of Methods and Applications" [3]. These models utilize a multi-step, noise-injecting denoising process that reverses a degradation procedure, contrasting with the adversarial training methodology underpinning GANs or the latent variable optimization approach in VAEs. This architectural contrast provides a robust basis for comparing their theoretical underpinnings and practical performances.

A primary concern in comparing these models is ensuring that the evaluation benchmarks reflect both qualitative and quantitative metrics pertinent to generative tasks. Metrics such as Fréchet Inception Distance (FID) and Inception Score (IS), frequently used for GANs, can extend to diffusion models, providing a baseline for evaluating image

quality and realism [81]. However, due to their iterative sampling method, diffusion models uniquely benefit from metrics that gauge process efficiency and sample diversity differently. The study "Accelerating Diffusion Models with Parallel Sampling: Inference at Sub-Linear Time Complexity" emphasizes improvements in sample efficiency, which could serve as a benchmark criterion for comparison across architectures, especially in computation-heavy applications [16].

Comparative analysis also requires a deep dive into the hyperparameter configurations of these models. In diffusion models, parameters like noise schedules, diffusion step size, and reverse process parameterization can drastically impact performance. Research in "Theoretical guarantees for sampling and inference in generative models with latent diffusions" highlights how tuning the drift coefficient in diffusion models can lead to better approximation strategies for target distributions—a feature absent in GANs' adversarial setup or the fixed prior assumption of VAEs. This examination of hyperparameters underscores the models' flexibility and adaptation to different data distributions.

Diffusion models have introduced innovative evaluation approaches, such as time-reversed Schrödinger bridges, which can fundamentally rethink how conditional synthesis tasks are measured, setting new paradigms in the evaluation space [36]. The ability to effectively integrate normalizing flows within diffusion models broadens the scope for novel metric integration, providing richer analysis dimensions compared to conventional models [73]. These developments represent a shift toward more specialized performance indicators that capture the nuanced operational characteristics of diffusion processes beyond conventional benchmarks.

Diffusion models' inherent strengths lie in their stability and ability to model complex, high-dimensional distributions without the instability often observed in GANs. GANs are prone to mode collapse, a problem diffusion models sidestep due to their progressive nature of refining samples through a reverse diffusion process. However, diffusion models are computationally expensive and can be slower due to iterative sampling needs [7]. VAEs, while straightforward in their encoding-decoding approach, often struggle with capturing detailed data distributions compared to the variance preserved in diffusion processes. Balancing these strengths and limitations is crucial for constructing a theoretical and practical framework for comparative evaluations.

As diffusion models evolve, integration with other generative paradigms is showing promise. The emerging trend of fusing diffusion models with energy-based and autoregressive models could offer new avenues for scalability and real-time applications [57]. Additionally, the convergence of diffusion models with optimal transport reflections and ODE-based structures is attracting attention for their potential to provide new solutions in domains requiring rigorous probabilistic interpretations [35].

Synthesizing these analyses identifies diffusion models as a maturing pillar in the field of generative modeling, offering distinctive advantages in robustness and stability over GANs and VAEs. While current models remain computationally intensive, innovations in fast approximation methods and hybrid model architectures will be pivotal in addressing these inefficiencies. The exploration of novel mathematical frameworks, such as those derived from stochastic control and optimal transport theory, is likely to refine the theoretical landscape of these models. Given diffusion models' promising trajectory, future research might focus on enhancing real-time applicability, integrating deep mode-learning features across latent spaces, and developing lightweight architectures for resource-constrained environments. Such directions promise to expand the boundaries and applications of this dynamic class of generative models.

# 7 CHALLENGES, LIMITATIONS, AND ETHICAL CONSIDERATIONS

## 7.1 Computational Challenges and Limitations

The computational challenges and limitations associated with generative diffusion models stem predominantly from their inherent complexity in sampling processes, scalability issues, and significant resource demands. Although these models have shown exceptional promise in various applications, such as image synthesis and molecular design, they confront substantial obstacles that restrict their broader applicability and efficiency.

One of the foremost challenges in generative diffusion models is the sampling efficiency. The process typically involves numerous iterative steps, each requiring extensive computation. The forward diffusion process progressively adds noise to the data, while the reverse process involves learning to gradually remove this noise to reconstruct the original data from a Gaussian noise input. This process necessitates multiple passes through a neural network, often resulting in high computational and time costs. This complexity is evident in the significant resource allocation necessary for these models to function, posing a barrier to both training and inference. Recent advancements have sought to address these sampling inefficiencies. For instance, techniques such as improved denoising diffusion probabilistic models have been developed to reduce the number of forward passes without compromising sample quality [2].

Scalability poses another critical limitation. As the dimensionality of data increases, so does the computational load. Diffusion models are notoriously memory-intensive, challenging their deployment in real-world scenarios where hardware resources may be constrained. Scaling these models to accommodate larger datasets or higher-resolution outputs requires substantial innovation in both algorithmic design and computational infrastructure. The complexity and sheer volume of data involved can strain memory and processing capabilities, making it crucial to develop solutions that optimize data handling and processing efficiency.

Given the high dimensionality of data these models typically engage with, the memory constraint is an overwhelming challenge that has yet to be fully resolved, despite certain breakthroughs. Memory management strategies and advanced architectural designs have been proposed to alleviate these burdens. For example, the introduction of latent space optimizations and more efficient model architectures has shown promise in improving both memory usage and processing speed [12]. However, these solutions

often involve trade-offs with other model attributes, such as inference time or flexibility.

The inference time itself is a notable limitation within diffusion models, which is closely linked to their iterative sampling nature. These prolonged inference times limit the practicality of diffusion models in applications that require real-time data processing or rapid output generation. Innovations such as parallel sampling and optimization algorithms have demonstrated potential in reducing latency, yet these approaches are still evolving and need further refinement to meet practical demands [16].

Regarding the trade-offs involved, while newer models tend to favor efficiency enhancements, they sometimes do so at the expense of accuracy or generalizability. For instance, methods that reduce sampling steps can inadvertently affect the fidelity of the generated outputs, introducing potential artifacts or reducing the model's robustness across diverse datasets [45]. Therefore, ongoing research is essential to balance these conflicting objectives effectively.

Emerging trends suggest that integrating principles from other generative paradigms, such as combining elements of generative adversarial networks (GANs) or variational autoencoders (VAEs), could offer pathways to mitigate some inherent limitations through hybrid frameworks. Similarly, leveraging insights from related fields, such as stochastic optimal control in diffusion-based generative modeling, also provides hope for novel solutions to these complexities [15].

From a technical perspective, approaches that explore alternative mathematical formulations or seek to simplify the computational processes involved, such as annealed importance sampling or score-based generative modeling, exhibit promise in enhancing both the accuracy and efficiency of diffusion models [46].

Moreover, there is active exploration into adaptive methods that dynamically adjust model complexity based on the task at hand, ensuring computational resources are allocated efficiently. Optimization techniques that involve adaptive step size control or noise scheduling, meanwhile, have gained attention for their potential to streamline processes without undermining performance [5].

In conclusion, while generative diffusion models encounter notable computational challenges and limitations, the field is witnessing substantial efforts to overcome these barriers through innovative strategies and cross-disciplinary integrations. The future directions lie in advancing computational methods that enhance sampling efficiency, scalability, and memory management while maintaining model accuracy and robustness. Effective solutions will likely emerge from a confluence of novel algorithmic designs, integration with complementary generative frameworks, and leveraging cutting-edge hardware capabilities. Such advancements will be pivotal in unlocking the full potential of diffusion models across various generative tasks and applications, ensuring they become a more viable and scalable option within the landscape of artificial intelligence.

## 7.2 Privacy Risks and Data Security

The rapid advancement and deployment of generative diffusion models have significantly reshaped artificial intelligence, introducing groundbreaking capabilities in data synthesis across diverse modalities. However, alongside these advances, they also present substantial privacy risks and data security concerns that warrant meticulous examination, particularly as these models transition to broader societal applications. This subsection delves into potential threats related to the memorization and replication of sensitive data, the implications of model sharing and deployment, and proposes strategies to mitigate these risks.

Generative diffusion models, characterized by their iterative noise injection and denoising processes, inherently process large volumes of data during training. A particularly concerning privacy risk associated with this data-intensive approach is the inadvertent memorization and subsequent reproduction of portions of training data, including sensitive information. This vulnerability is critical, especially when training data includes personally identifiable information (PII) or confidential business data. Empirical investigations have revealed instances where generative models reproduce training datasets directly, leading to potential privacy breaches [67].

The architecture and optimization processes of these models contribute to potential memorization of data. Larger models with more parameters may inadvertently store detailed patterns from the training data. Techniques like differential privacy offer a method to counteract this risk by adding calibrated noise to the data, thus obscuring individual data points without significantly impacting model performance. Regularization techniques, which promote generalization and reduce overfitting, also play a pivotal role in diminishing memorization effects. For instance, contractive diffusion probabilistic models (CDPMs) incorporate constraints that reduce the retention of specific data features during training, presenting a promising avenue for future research and application [94].

Another significant concern involves membership inference attacks (MIAs), where adversaries exploit model outputs to discern whether specific data points were part of the training set. Membership inference has been scrutinized across machine learning models, including diffusion models, with demonstrated susceptibility. One safeguard against MIAs is the adoption of regularization strategies that mask data influence while maintaining output fidelity [3]. By integrating anonymization algorithms during data preprocessing, model developers can obscure identifiers, reducing vulnerability to MIAs.

The sharing of pre-trained models across organizations introduces another vector of privacy risk due to the embedded knowledge from training datasets. Without adequate vetting and sanitization processes, there exists a risk of redistributing captured sensitive data. Techniques such as model distillation and knowledge transfer could distill essential features into smaller, less information-dense architectures, ostensibly reducing memorized data volume [27]. This underscores the need for robust model auditing and verification procedures as integral components of the model development lifecycle, ensuring compliance with ethical, legal, and organizational privacy standards.

Further scrutiny of data security within diffusion models necessitates attention to model inversion attacks. Distinct from membership inference attempts, these attacks reverse-engineer model inputs by accessing the model parameters or outputs. The latent space representations and reversible

processes in diffusion models present a unique target for inversion attacks, exploiting inherent noise reversal processes. Employing encryption methods or secure enclaves for model operations provides potential defensive strategies against unauthorized access and inversion attempts, ensuring data confidentiality [95].

Emerging trends in addressing privacy concerns within generative diffusion models highlight the vital interplay between regulatory frameworks and technical innovation. The General Data Protection Regulation (GDPR) in the European Union emphasizes the necessity for privacy-by-design principles, urging developers to incorporate data protection measures into the design phase. Concurrently, federated learning paradigms enable collaborative model training without direct data sharing between entities, mitigating data leakage risks and aligning with broader data sovereignty and compliance mandates [96].

In synthesizing the privacy challenges and opportunities presented by diffusion models, it is evident that future directions must integrate advanced cryptographic techniques, enhance federated learning frameworks, and develop proprietary model verification technologies. As the field evolves, cultivating a research and development ethos that prioritizes privacy while harnessing the transformative potential of these models is imperative. Cross-disciplinary collaborations will be essential to bridge gaps between technical feasibility, privacy legislation, and ethical implementation, ensuring diffusion models remain a sustainable and secure pillar within generative artificial intelligence.

## 7.3 Ethical Implications and Societal Impact

In recent years, generative diffusion models have emerged as a groundbreaking paradigm in artificial intelligence, hailed for their ability to produce high-quality synthetic data that can mimic human-like creativity. However, as these models transition from research environments to widespread societal applications, they bring with them a host of ethical challenges that demand careful scrutiny. This subsection explores these ethical dimensions, particularly focusing on bias amplification, misinformation, and broader societal harm.

To begin with, one of the pivotal concerns regarding generative diffusion models is their propensity to amplify existing biases inherent in training datasets. Diffusion models, like many deep learning frameworks, derive their functionality from large datasets. If these datasets are imbued with biases—whether those relate to gender, race, socioeconomic status, or any other attribute—the models trained on them are likely to learn and propagate these biases. This amplification occurs because diffusion models, through their iterative noise-reducing processes, often highlight the most prominent statistical features of the training data, which may include biased patterns [90]. As a result, there is a risk of reinforcing and perpetuating social stereotypes, thereby influencing public perception in detrimental ways [90].

Moreover, another ethical dimension of diffusion models is their potential to generate misinformation. Given their capacity to produce highly convincing synthetic media, these models can be exploited to create deepfakes or other forms of misrepresented content. Such capabilities hold the potential for misuse in malicious campaigns that sow discord or manipulate public opinion by producing misleading news footage, fake audio recordings, or fabricated images [79]. These implications are significant, as they challenge the critical mechanisms of information verification and trust in media—a cornerstone in democratic societies.

Additionally, there is the issue of societal harm that extends beyond misinformation and bias amplification. For instance, diffusion models' ability to generate realistic imagery can be used for high-stakes and sensitive applications, such as in surveillance systems. While these applications might benefit security efforts, they also pose daunting privacy concerns. The potential for surveillance tools to become pervasive, coupled with the lacking transparency in how such models operate and decide, could lead to intrusive oversight and infringement of individual freedoms. Models trained without ethical oversight may inadvertently expose vulnerable groups to additional scrutiny, exacerbating existing societal inequalities [70].

Despite these challenges, there exist emergent strategies and technologies designed to mitigate the risks associated with diffusion models. For instance, fairness-aware training techniques involve incorporating modules that detect and correct bias during the training phase. Such interventions are increasingly adopted to combat inherited biases from training datasets effectively [90]. Beyond this, advanced differential privacy techniques are being employed to limit the amount of sensitive data information that models can learn, thus protecting individual privacy [30]. Additionally, watermarking techniques that subtly tag synthetic content can serve as a means to track and verify the authenticity of media, offering a counterbalance to misinformation risks [3].

Moreover, there is a critical need for establishing ethical guidelines and regulatory frameworks that govern the deployment of diffusion models. Policymakers and researchers are actively engaging in dialogue to set these guidelines, aiming to foster responsible AI development that prioritizes public trust and safety. Regulatory bodies are urged to impose stricter mandates on content verification processes and to implement standards that demand transparency from model developers regarding data sources and model usage intentions [3].

In conclusion, while generative diffusion models hold the promise of innovation across diverse domains, they simultaneously pose profound ethical challenges that must be addressed. Future research must continue to refine these models, not only seeking to enhance their technical prowess but also ensuring their alignment with societal values and equity principles. Integrating ethical considerations into the technology's lifecycle—from design to deployment—can help strike a balance between leveraging its capabilities and safeguarding against potential harms. It is imperative that these models evolve in tandem with comprehensive ethical frameworks and technologically sound strategies to mitigate risks. As this field advances, collaborative efforts between technologists, ethicists, and policymakers will be essential to navigating this complex landscape sustainably and ethically.

## 7.4 Generalization and Robustness Challenges

Generative diffusion models have emerged as a promising computational framework in the realm of generative modeling, lauded for their ability to produce high-fidelity and diverse synthetic samples across a multitude of domains. However, these models encounter significant challenges in generalization and robustness, which are intrinsic to their deployment in diverse datasets, varying environmental conditions, and exposure to adversarial inputs. This subsection delves into these critical challenges, focusing on how diffusion models navigate the intricate balance between specificity and generalization, and examining the forefront of research committed to enhancing their robustness.

A fundamental challenge in deploying generative diffusion models is their capacity to generalize effectively. Although these models can encapsulate complex data distributions, their performance may diminish when confronted with new, unseen data. This overfitting risk is particularly acute when models are trained on limited datasets, resulting in a narrow scope within the model's learned representation. A propensity for these models to memorize specific training data points often compromises their generalization capabilities, a phenomenon noted in research concerning score-based models, whereby memorization is observed [70].

Efforts to alleviate overfitting include implementing regularization techniques and developing architectures designed to adapt nimbly to unseen data without over-reliance on available examples. Additionally, the introduction of robust evaluation metrics, as outlined in [81], equips practitioners with insights into a model's generalization during the validation phase, allowing the preemptive identification of overfitting issues prior to deployment.

Beyond generalization, robustness is paramount to model performance — specifically, a model's resilience to perturbations and ability to function under various conditions. The inherently stochastic nature of diffusion models presents both opportunities and challenges. Although this stochasticity facilitates exploration and diversity in generated outputs, it also renders models vulnerable to perturbations, potentially causing significant deviations in results. Robust training protocols and thoughtful noise schedules are essential for reinforcing diffusion model robustness. Modified training procedures that integrate adversarial noise further bolster robustness, as suggested by contrasting approaches in generative modeling, including comparisons with variational autoencoders and GANs which consider robustness directly [57].

Comparing the robustness and generalization attributes of diffusion models with other generative architectures like GANs and VAEs reveals distinct differences. GANs frequently struggle with mode collapse, leading to a lack of diversity in generated data, whereas well-regularized diffusion models typically produce more varied outputs. Conversely, while VAEs are adept at tasks requiring learned latent spaces, they may falter in generating sharp, high-fidelity samples characteristic of diffusion model outputs [57]. Each paradigm contends with unique robustness challenges, suggesting that integrating diffusion models with other architectures within hybrid frameworks may exploit the strengths of each model type [97].

Recent innovative approaches aim to address generalization and robustness in diffusion models. Transformer-based architectures are being leveraged to enhance scalability and robustness, broadening their efficacy across varied tasks [52]. Moreover, the incorporation of reinforcement learning strategies to dynamically tailor the generation process illustrates how adaptive learning enhances robustness [37].

Ensuring robustness and generalization in diffusion models necessitates a comprehensive approach that encompasses not only model architecture and training methodologies but also the implications of data selection and preprocessing. As these models progress, the focus will likely pivot towards interpretability in modeling techniques and novel architectural developments that inherently address variability and complexity in data.

In summary, generative diffusion models are at the cutting edge of artificial intelligence, offering formidable capabilities balanced by challenges in generalization and robustness. Future advancements likely reside in hybrid modeling strategies, drawing upon insights from multiple paradigms to yield more versatile and potent generative models. Moreover, the advancement of training algorithms incorporating robust learning principles derived from stochastic control and optimal transport theories may open new avenues for overcoming extant limitations [15], [82]. Engaging the research community in comparative studies that critically assess these approaches across diverse real-world scenarios will be essential to unlock the full potential of diffusion models in practical applications.

## 8 CONCLUSION

In this comprehensive survey, we have meticulously reviewed the expansive landscape of generative diffusion models, which have emerged as a pivotal component in the arsenal of machine learning and artificial intelligence. The ongoing evolution of diffusion models showcases an intersection of rigorous mathematical underpinnings with sophisticated algorithmic innovations, all while extending their reach into multifaceted applications across domains. Herein, we synthesize the critical insights drawn from preceding sections, reflecting on current advancements and proposing directions for future exploration.

Generative diffusion models, initially constrained by their theoretical complexity, have exhibited a notable trajectory of refinement and optimization. Traditional models, like the Denoising Diffusion Probabilistic Models (DDPM), have been challenged and extended to incorporate score-based generative modeling and stochastic differential equations, offering powerful frameworks to model high-dimensional data distributions [2], [5]. This has led to significant contributions in understanding the intricacies of forward and reverse diffusion processes, with various perspectives such as those elucidating connections to stochastic control theory [15].

One of the salient features of diffusion models, their capacity for producing high-quality samples, confronts the perennial trade-off between computational efficiency and accuracy. Recent innovations have emerged to tackle these challenges by optimizing sampling processes via efficient

strategies like stochastic sampling and hybrid architectures [16]. Meanwhile, architectural innovations like the integration of Vision Transformers (ViT) with diffusion models exemplify the modularity and scalability of these models, fostering adaptive frameworks that cater to diverse tasks ranging from image synthesis to DNA sequence generation [98], [99].

The adaptability of diffusion models to complex structured data is a testament to their growing sophistication. As explored through applications in fields such as molecular modeling and remote sensing, these models demonstrate remarkable prowess in circumventing the traditional barriers posed by high-dimensional, complex data spaces [12]. In particular, the extension of diffusion models to structured domains requires nuanced adjustments to account for inherent data constraints, facilitating applications like molecule design and climate modeling [61], [100].

Furthermore, empirical evidence suggests that diffusion models provide robust generative frameworks with impressive capabilities in task-specific adaptation and customization. This is particularly evident in applications that require nuanced and context-aware outputs, which are predominant in cross-modal generation tasks and require excellent coherence and fidelity [5]. Yet, the transition from generating synthetic data to assimilating it into real-world applications continues to pose formidable challenges. The multidimensional facets of generality versus specificity necessitate a deeper probing into adaptive learning mechanisms within these models.

Despite notable triumphs, the landscape of generative diffusion models is entwined with challenges that demand further scholarly attention. Ethical implications, such as potential bias dissemination and misuse for disinformation, highlight a pressing need for ethical guidelines and fairness-adaptive training protocols [44], [101]. Additionally, the expansive memory footprints and computational heavy lifting required for training and implementing large-scale models underscore the necessity for methodological improvements aimed at enhancing computational efficiency and energy usage [14].

Looking to the future, the fusion of diffusion models with other generative architectures, such as generative adversarial networks and autoencoders, presents a promising avenue for hybrid frameworks that synthesize the strengths of each approach to overcome existing limitations. Advances in reinforcement learning techniques also offer compelling prospects for enhancing alignment to human preferences, providing a bridge between model capabilities and intuitive human-centric applications [37], [47].

In conclusion, the journey of generative diffusion models is emblematic of the broader quest for intelligent systems that mirror the complexities and nuances of the world. As researchers persist in exploring theoretical foundations, computational optimizations, and ethical frameworks, diffusion models stand as a testament to the fusion of interdisciplinary insights with technological innovation. This survey serves not only as a compilation of current knowledge but also as a guidepost for future research that aims to harness the full potential of diffusion models in advancing artificial intelligence in a manner that is not only inventive but also morally and socially conscious.

## REFERENCES

[1] C. Luo, "Understanding diffusion models: A unified perspective," *ArXiv*, vol. abs/2208.11970, 2022. 1, 2, 11, 13, 16

[2] A. Nichol and P. Dhariwal, "Improved denoising diffusion probabilistic models," *ArXiv*, vol. abs/2102.09672, 2021. 1, 2, 6, 9, 16, 19, 22

[3] L. Yang, Z. Zhang, S. Hong, R. Xu, Y. Zhao, Y. Shao, W. Zhang, M.-H. Yang, and B. Cui, "Diffusion models: A comprehensive survey of methods and applications," *ACM Computing Surveys*, vol. 56, pp. 1 – 39, 2022. 1, 2, 11, 13, 14, 16, 17, 18, 20, 21

[4] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, "Diffusion models in vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, pp. 10 850–10 869, 2022. 1, 2, 9, 10, 13, 16, 17

[5] G. Batzolis, J. Stanczuk, C. Schonlieb, and C. Etmann, "Conditional image generation with score-based diffusion models," *ArXiv*, vol. abs/2111.13606, 2021. 1, 2, 6, 10, 20, 22, 23

[6] T. Karras, M. Aittala, T. Aila, and S. Laine, "Elucidating the design space of diffusion-based generative models," *ArXiv*, vol. abs/2206.00364, 2022. 1, 9, 10, 13, 16

[7] A. Ulhaq, N. Akhtar, and G. Pogrebna, "Efficient diffusion models for vision: A survey," *ArXiv*, vol. abs/2210.09292, 2022. 1, 2, 9, 10, 19

[8] C. Zhang, C. Zhang, M. Zhang, and I.-S. Kweon, "Text-to-image diffusion models in generative ai: A survey," *ArXiv*, vol. abs/2303.07909, 2023. 1, 10, 16

[9] M. Fuest, P. Ma, M. Gui, J. S. Fischer, V. T. Hu, and B. Ommer, "Diffusion models and representation learning: A survey," *ArXiv*, vol. abs/2407.00783, 2024. 1

[10] X. Yang, S.-M. Shih, Y. Fu, X. Zhao, and S. Ji, "Your vit is secretly a hybrid discriminative-generative diffusion model," *ArXiv*, vol. abs/2208.07791, 2022. 1

[11] S. Rissanen, M. Heinonen, and A. Solin, "Generative modelling with inverse heat dissipation," *ArXiv*, vol. abs/2206.13397, 2022. 1, 15

[12] W. Fan, C. Liu, Y. Liu, J. Li, H. Li, H. Liu, J. Tang, and Q. Li, "Generative diffusion models on graphs: Methods and applications," *ArXiv*, vol. abs/2302.02591, 2023. 2, 16, 19, 23

[13] L. Lin, Z. Li, R. Li, X. Li, and J. Gao, "Diffusion models for time-series applications: a survey," *Frontiers of Information Technology & Electronic Engineering*, vol. 25, pp. 19–41, 2023. 2, 13, 16

[14] Y.-H. Chen, R. Sarokin, J. Lee, J. Tang, C.-L. Chang, A. Kulik, and M. Grundmann, "Speed is all you need: On-device acceleration of large diffusion models via gpu-aware optimizations," *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 4651–4655, 2023. 2, 23

[15] J. Berner, L. Richter, and K. Ullrich, "An optimal control perspective on diffusion-based generative modeling," *Trans. Mach. Learn. Res.*, vol. 2024, 2022. 2, 4, 20, 22

[16] H. Chen, Y. Ren, L. Ying, and G. M. Rotskoff, "Accelerating diffusion models with parallel sampling: Inference at sub-linear time complexity," *ArXiv*, vol. abs/2405.15986, 2024. 2, 6, 19, 20, 23

[17] Y. Song, J. N. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," *ArXiv*, vol. abs/2011.13456, 2020. 3, 14

[18] N. F. Marshall and M. Hirn, "Time coupled diffusion maps," *ArXiv*, vol. abs/1608.03628, 2016. 3

[19] C.-W. Huang, J. H. Lim, and A. C. Courville, "A variational perspective on diffusion-based generative models and score matching," *ArXiv*, vol. abs/2106.02808, 2021. 3, 7

[20] R. S. Roman, E. Nachmani, and L. Wolf, "Noise estimation for generative diffusion models," *ArXiv*, vol. abs/2104.02600, 2021. 3, 7

[21] Y. Song, C. Durkan, I. Murray, and S. Ermon, "Maximum likelihood training of score-based diffusion models," in *Neural Information Processing Systems*, 2021, pp. 1415–1428. 3, 7, 8, 14, 15, 17

[22] E. Heitz, L. Belcour, and T. Chambon, "Iterative α -(de)blending: a minimalist deterministic diffusion model," *ACM SIGGRAPH 2023 Conference Proceedings*, 2023. 3, 10

[23] K. Zheng, C. Lu, J. Chen, and J. Zhu, "Improved techniques for maximum likelihood estimation for diffusion odes," *ArXiv*, vol. abs/2305.03935, 2023. 3

[24] G. Li and Y. Yan, "Adapting to unknown low-dimensional structures in score-based diffusion models," *ArXiv*, vol. abs/2405.14861, 2024. 3, 18

[25] B. Song, S. M. Kwon, Z. Zhang, X. Hu, Q. Qu, and L. Shen, "Solving inverse problems with latent diffusion models via hard data consistency," *ArXiv*, vol. abs/2307.08123, 2023. 3

[26] H. Chen and L. Ying, "Convergence analysis of discrete diffusion model: Exact implementation through uniformization," *ArXiv*, vol. abs/2402.08095, 2024. 4

[27] W. Tang and H. Zhao, "Score-based diffusion models via stochastic differential equations - a technical tutorial," *ArXiv*, vol. abs/2402.07487, 2024. 4, 14, 20

[28] K. Oko, S. Akiyama, and T. Suzuki, "Diffusion models are minimax optimal distribution estimators," *ArXiv*, vol. abs/2303.01861, 2023. 4

[29] M. Chen, K. Huang, T. Zhao, and M. Wang, "Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data," *ArXiv*, vol. abs/2302.07194, 2023. 4

[30] D. Kwon, Y. Fan, and K. Lee, "Score-based generative modeling secretly minimizes the wasserstein distance," *ArXiv*, vol. abs/2212.06359, 2022. 4, 21

[31] C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu, "Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps," *ArXiv*, vol. abs/2206.00927, 2022. 4, 11

[32] T. Dockhorn, A. Vahdat, and K. Kreis, "Score-based generative modeling with critically-damped langevin diffusion," *ArXiv*, vol. abs/2112.07068, 2021. 4

[33] B. Tzen and M. Raginsky, "Theoretical guarantees for sampling and inference in generative models with latent diffusions," in *Annual Conference Computational Learning Theory*, 2019, pp. 3084–3114. 4, 6

[34] S. Liu, X. Zhou, Y. Jiao, and J. Huang, "Wasserstein generative learning of conditional distribution," *ArXiv*, vol. abs/2112.10039, 2021. 4

[35] L. Zhang, E. Weinan, and L. Wang, "Monge-ampère flow for generative modeling," *ArXiv*, vol. abs/1809.10188, 2018. 4, 5, 19

[36] Y. Shi, V. D. Bortoli, G. Deligiannidis, and A. Doucet, "Conditional simulation using diffusion schrödinger bridges," in *Conference on Uncertainty in Artificial Intelligence*, 2022, pp. 1792–1802. 5, 19

[37] Y. Zhang, E. Tzeng, Y. Du, and D. Kislyuk, "Large-scale reinforcement learning for diffusion models," *ArXiv*, vol. abs/2401.12244, 2024. 5, 22, 23

[38] C.-W. Huang, M. Aghajohari, A. Bose, P. Panangaden, and A. C. Courville, "Riemannian diffusion models," *ArXiv*, vol. abs/2208.07949, 2022. 5

[39] G. Kerrigan, G. Migliorini, and P. Smyth, "Dynamic conditional optimal transport through simulation-free flows," *ArXiv*, vol. abs/2404.04240, 2024. 5

[40] S. Peluchetti, "Diffusion bridge mixture transports, schrödinger bridge problems and generative modeling," *J. Mach. Learn. Res.*, vol. 24, pp. 374:1–374:51, 2023. 5

[41] J. Choi, J. Choi, and M. joo Kang, "Generative modeling through the semi-dual formulation of unbalanced optimal transport," *ArXiv*, vol. abs/2305.14777, 2023. 5

[42] W. Deng, Y. Chen, N. T. Yang, H. Du, Q. Feng, and R. T. Q. Chen, "Reflected schrödinger bridge for constrained generative modeling," *ArXiv*, vol. abs/2401.03228, 2024. 5

[43] J.-H. Bastek, W. Sun, and D. Kochmann, "Physics-informed diffusion models," *ArXiv*, vol. abs/2403.14404, 2024. 6

[44] H. Chung, J. Kim, G. Y. Park, H. Nam, and J. C. Ye, "Cfg++: Manifold-constrained classifier free guidance for diffusion models," *ArXiv*, vol. abs/2406.08070, 2024. 6, 23

[45] A. Bansal, E. Borgnia, H.-M. Chu, J. Li, H. Kazemi, F. Huang, M. Goldblum, J. Geiping, and T. Goldstein, "Cold diffusion: Inverting arbitrary image transforms without noise," *ArXiv*, vol. abs/2208.09392, 2022. 6, 20

[46] A. Doucet, W. Grathwohl, A. G. Matthews, and H. Strathmann, "Score-based diffusion meets annealed importance sampling," *ArXiv*, vol. abs/2208.07698, 2022. 6, 20

[47] K. Black, M. Janner, Y. Du, I. Kostrikov, and S. Levine, "Training diffusion models with reinforcement learning," *ArXiv*, vol. abs/2305.13301, 2023. 6, 13, 16, 23

[48] D. Kim, B. Na, S. Kwon, D. Lee, W. Kang, and I.-C. Moon, "Maximum likelihood training of implicit nonlinear diffusion models," *ArXiv*, vol. abs/2205.13699, 2022. 7

[49] H. Chefer, O. Lang, M. Geva, V. Polosukhin, A. Shocher, M. Irani, I. Mosseri, and L. Wolf, "The hidden language of diffusion models," *ArXiv*, vol. abs/2306.00966, 2023. 8

[50] J. Jo, S. Lee, and S. J. Hwang, "Score-based generative modeling of graphs via the system of stochastic differential equations," in *International Conference on Machine Learning*, 2022, pp. 10 362–10 383. 8

[51] W. Luo, "A comprehensive survey on knowledge distillation of diffusion models," *ArXiv*, vol. abs/2304.04262, 2023. 8

[52] W. S. Peebles and S. Xie, "Scalable diffusion models with transformers," *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4172–4182, 2022. 8, 12, 22

[53] Y. Lee, J.-Y. Kim, H. Go, M. Jeong, S. Oh, and S. Choi, "Multi-architecture multi-expert diffusion models," *ArXiv*, vol. abs/2306.04990, 2023. 9

[54] Z. Guo, J. Liu, Y. Wang, M. Chen, D. Wang, D. Xu, and J. Cheng, "Diffusion models in bioinformatics: A new wave of deep learning revolution in action," *ArXiv*, vol. abs/2302.10907, 2023. 9, 12

[55] Z. Chang, G. Koulieris, and H. P. H. Shum, "On the design fundamentals of diffusion models: A survey," *ArXiv*, vol. abs/2306.04542, 2023. 9

[56] Z. Hua, Y. He, C. Ma, and A. K. Anderson-Frey, "Weather prediction with diffusion guided by realistic forecast processes," *ArXiv*, vol. abs/2402.06666, 2024. 9

[57] S. Bond-Taylor, A. Leach, Y. Long, and C. 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*, vol. 44, pp. 7327–7347, 2021. 9, 12, 15, 19, 22

[58] Z. Xing, Q. Feng, H. Chen, Q. Dai, H.-R. Hu, H. Xu, Z. Wu, and Y.-G. Jiang, "A survey on video diffusion models," *ArXiv*, vol. abs/2310.10647, 2023. 10, 13

[59] A. Graikos, N. Malkin, N. Jojic, and D. Samaras, "Diffusion models as plug-and-play priors," *ArXiv*, vol. abs/2206.09012, 2022. 10, 13

[60] H. Cao, C. Tan, Z. Gao, Y. Xu, G. Chen, P. Heng, and S. Z. Li, "A survey on generative diffusion model," *ArXiv*, vol. abs/2209.02646, 2022. 10

[61] S. Azizi, S. Kornblith, C. Saharia, M. Norouzi, and D. J. Fleet, "Synthetic data from diffusion models improves imagenet classification," *ArXiv*, vol. abs/2304.08466, 2023. 10, 16, 23

[62] P. H. Richemond, S. Dieleman, and A. Doucet, "Categorical sdes with simplex diffusion," *ArXiv*, vol. abs/2210.14784, 2022. 10, 18

[63] H. Chung, J. C. Ye, P. Milanfar, and M. Delbracio, "Prompt-tuning latent diffusion models for inverse problems," *ArXiv*, vol. abs/2310.01110, 2023. 10

[64] V. Dutordoir, A. D. Saul, Z. Ghahramani, and F. Simpson, "Neural diffusion processes," in *International Conference on Machine Learning*, 2022, pp. 8990–9012. 10

[65] J. Hu, B. Song, X. Xu, L. Shen, and J. A. Fessler, "Learning image priors through patch-based diffusion models for solving inverse problems," *ArXiv*, vol. abs/2406.02462, 2024. 11

[66] V. D. Bortoli, "Convergence of denoising diffusion models under the manifold hypothesis," *ArXiv*, vol. abs/2208.05314, 2022. 11

[67] J. Wu, T. Le, M. Hayat, and M. Harandi, "Erasediff: Erasing data influence in diffusion models," *ArXiv*, vol. abs/2401.05779, 2024. 11, 20

[68] H. Chung and J.-C. Ye, "Score-based diffusion models for accelerated mri," *Medical image analysis*, vol. 80, p. 102479, 2021. 11, 15

[69] V. Popov, I. Vovk, V. Gogoryan, T. Sadekova, and M. Kudinov, "Grad-tts: A diffusion probabilistic model for text-to-speech," in *International Conference on Machine Learning*, 2021, pp. 8599–8608. 11, 15

[70] X. Gu, C. Du, T. Pang, C. Li, M. Lin, and Y. Wang, "On memorization in diffusion models," *ArXiv*, vol. abs/2310.02664, 2023. 12, 21, 22

[71] M. Uehara, Y. Zhao, T. Biancalani, and S. Levine, "Understanding reinforcement learning-based fine-tuning of diffusion models: A tutorial and review," *ArXiv*, vol. abs/2407.13734, 2024. 12, 18

[72] Y. Zhu and Y. Zhao, "Diffusion models in nlp: A survey," *ArXiv*, vol. abs/2303.07576, 2023. 12

[73] I. Kobyzev, S. Prince, and M. A. Brubaker, "Normalizing flows: An introduction and review of current methods," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, pp. 3964–3979, 2020. 12, 15, 19

[74] Y. Yang, M. Jin, H. Wen, C. Zhang, Y. Liang, L. Ma, Y. Wang, C.-M. Liu, B. Yang, Z. Xu, J. Bian, S. Pan, and Q. Wen, "A survey

on diffusion models for time series and spatio-temporal data," *ArXiv*, vol. abs/2404.18886, 2024. 12

[75] Y. Zhou, C. Shi, L. Li, and Q. Yao, "Testing for the markov property in time series via deep conditional generative learning," *Journal of the Royal Statistical Society. Series B, Statistical Methodology*, vol. 85, pp. 1204 – 1222, 2023. 12

[76] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4396–4405, 2018. 12, 15

[77] C. Frogner and T. Poggio, "Approximate inference with wasserstein gradient flows," in *International Conference on Artificial Intelligence and Statistics*, 2018, pp. 2581–2590. 14

[78] S. Xue, M. Yi, W. Luo, S. Zhang, J. Sun, Z. Li, and Z.-M. Ma, "Sa-solver: Stochastic adams solver for fast sampling of diffusion models," *ArXiv*, vol. abs/2309.05019, 2023. 14

[79] D. Kim, Y. Kim, W. Kang, and I.-C. Moon, "Refining generative process with discriminator guidance in score-based diffusion models," in *International Conference on Machine Learning*, 2022, pp. 16 567–16 598. 15, 21

[80] Y. Wu, M. Chen, Z. Li, M. Wang, and Y. Wei, "Theoretical insights for diffusion guidance: A case study for gaussian mixture models," *ArXiv*, vol. abs/2403.01639, 2024. 15, 18

[81] L. Theis, A. van den Oord, and M. Bethge, "A note on the evaluation of generative models," *CoRR*, vol. abs/1511.01844, 2015. 15, 19, 22

[82] F. Briol, A. Barp, A. Duncan, and M. Girolami, "Statistical inference for generative models with maximum mean discrepancy," *ArXiv*, vol. abs/1906.05944, 2019. 15, 22

[83] M. Gomez-Rodriguez and B. Scholkopf, "Submodular inference of diffusion networks from multiple trees," *ArXiv*, vol. abs/1205.1671, 2012. 16

[84] H. Li, C. Xia, T. Wang, S. Wen, C. Chen, and Y. Xiang, "Capturing dynamics of information diffusion in sns: A survey of methodology and techniques," *ACM Computing Surveys (CSUR)*, vol. 55, pp. 1 – 51, 2021. 16

[85] C. Zhu, J. Tang, H. Brouwer, J. F. P'erez, M. van Dijk, and L. Y. Chen, "Quantifying and mitigating privacy risks for tabular generative models," *ArXiv*, vol. abs/2403.07842, 2024. 16

[86] A. Luccioni, C. Akiki, M. Mitchell, and Y. Jernite, "Stable bias: Analyzing societal representations in diffusion models," *ArXiv*, vol. abs/2303.11408, 2023. 16

[87] B. Liu, S. Shao, B. Li, L. Bai, Z. Xu, H. Xiong, J. Kwok, A. Helal, and Z. Xie, "Alignment of diffusion models: Fundamentals, challenges, and future," *ArXiv*, vol. abs/2409.07253, 2024. 16

[88] C. Weilbach, W. Harvey, and F. Wood, "Graphically structured diffusion models," in *International Conference on Machine Learning*, 2022, pp. 36 887–36 909. 16

[89] L. Liu, Y. Ren, Z. Lin, and Z. Zhao, "Pseudo numerical methods for diffusion models on manifolds," *ArXiv*, vol. abs/2202.09778, 2022. 17

[90] Y. Kim, B. Na, M. Park, J. Jang, D. Kim, W. Kang, and I.-C. Moon, "Training unbiased diffusion models from biased dataset," *ArXiv*, vol. abs/2403.01189, 2024. 18, 21

[91] H. Fu, Z. Yang, M. Wang, and M. Chen, "Unveil conditional diffusion models with classifier-free guidance: A sharp statistical theory," *ArXiv*, vol. abs/2403.11968, 2024. 18

[92] G. Daras, Y. Dagan, A. Dimakis, and C. Daskalakis, "Consistent diffusion models: Mitigating sampling drift by learning to be consistent," *ArXiv*, vol. abs/2302.09057, 2023. 18

[93] Y. Guo, H. Yuan, Y. Yang, M. Chen, and M. Wang, "Gradient guidance for diffusion models: An optimization perspective," *ArXiv*, vol. abs/2404.14743, 2024. 18

[94] W. Tang and H. Zhao, "Contractive diffusion probabilistic models," *ArXiv*, vol. abs/2401.13115, 2024. 20

[95] A. Lou and S. Ermon, "Reflected diffusion models," *ArXiv*, vol. abs/2304.04740, 2023. 21

[96] F. Rozet, G. Andry, F. Lanusse, and G. Louppe, "Learning diffusion priors from observations by expectation maximization," *ArXiv*, vol. abs/2405.13712, 2024. 21

[97] L. Ruthotto and E. Haber, "An introduction to deep generative modeling," *GAMM-Mitteilungen*, vol. 44, 2021. 22

[98] F. Bao, S. Nie, K. Xue, Y. Cao, C. Li, H. Su, and J. Zhu, "All are worth words: A vit backbone for diffusion models," *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 22 669–22 679, 2022. 23

[99] Z. Li, Y. Ni, W. A. V. Beardall, G. Xia, A. Das, G. Stan, and Y. Zhao, "Discdiff: Latent diffusion model for dna sequence generation," *ArXiv*, vol. abs/2402.06079, 2024. 23

[100] N. Fishman, L. Klarner, V. D. Bortoli, E. Mathieu, and M. Hutchinson, "Diffusion models for constrained domains," *ArXiv*, vol. abs/2304.05364, 2023. 23

[101] A. C. Li, M. Prabhudesai, S. Duggal, E. L. Brown, and D. Pathak, "Your diffusion model is secretly a zero-shot classifier," *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 2206–2217, 2023. 23