

DIFFUSION BRIDGE OR FLOW MATCHING? A UNIFYING FRAMEWORK AND COMPARATIVE ANALYSIS

Anonymous authors

Paper under double-blind review

ABSTRACT

Diffusion Bridge and Flow Matching have both demonstrated compelling empirical performance in transformation between arbitrary distributions. However, there remains confusion about which approach is generally preferable, and the substantial discrepancies in their modeling assumptions and practical implementations have hindered a unified theoretical account of their relative merits. We have, for the first time, provided a unified theoretical and experimental validation of these two models. We recast their frameworks through the lens of Stochastic Optimal Control and prove that the cost function of the Diffusion Bridge is lower, guiding the system toward more stable and natural trajectories. Simultaneously, from the perspective of Optimal Transport, interpolation coefficients t and $1 - t$ of Flow Matching become increasingly ineffective when the training data size is reduced. To corroborate these theoretical claims, we propose a novel, powerful architecture for Diffusion Bridge built on a latent Transformer, and implement a Flow Matching model with the same structure to enable a fair performance comparison in various experiments. Comprehensive experiments are conducted across Image Inpainting, Super-Resolution, Deblurring, Denoising, Translation, and Style Transfer tasks, systematically varying both the distributional discrepancy (different difficulty) and the training data size. Extensive empirical results align perfectly with our theoretical predictions and allow us to delineate the respective advantages and disadvantages of these two models. Our code is available at <https://anonymous.4open.science/r/DBFM-3E8E/>.

1 INTRODUCTION

Diffusion models have been widely used in a variety of applications, demonstrating remarkable capabilities and promising results in numerous tasks such as image generation (Ho et al., 2020; Song et al., 2020; Xia et al., 2023), video generation (Ho et al., 2022; Luo et al., 2023a), imitation learning (Wu et al., 2024; Chi et al., 2023; Ze et al., 2024; Reuss et al., 2023), and reinforcement learning (Ding et al., 2024; 2025; Wang et al., 2022; Ada et al., 2024), etc. However, standard diffusion models exhibit inherent limitations that they are difficult to achieve the conversion between any two distributions since its prior distribution is assumed to be Gaussian noise. They can also rely on meticulously designed conditioning mechanisms and classifier/loss guidance (Chung et al., 2023; Yang et al., 2024; Kawar et al., 2022) to facilitate conditional sampling and ensure output alignment with a target distribution. However, these methods can be cumbersome and may introduce manifold deviations during the sampling process.

To address this challenge, Diffusion Bridge and Flow Matching are two prevalent approaches for achieving distribution-to-distribution transformation. As for diffusion bridges, on one hand, Schrödinger Bridge (Liu et al., 2023; Shi et al., 2024; De Bortoli et al., 2021) deterministically steers one prescribed probability measure to another via the minimum-entropy re-weighting of an underlying reference stochastic process, producing a path-wise coupling whose time-marginals coincide with the specified endpoint distributions. On the other hand, DDBMs (Zhou et al., 2023), GOUB (Yue et al., 2024), and related methods enforce end-to-end exact matching by incorporating Doob’s h -transform into the forward Stochastic Differential Equations (SDE) of the diffusion process. UniDB (Zhu et al., 2025) further reformulates and unifies the diffusion bridge paradigm through a stochastic optimal control framework.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

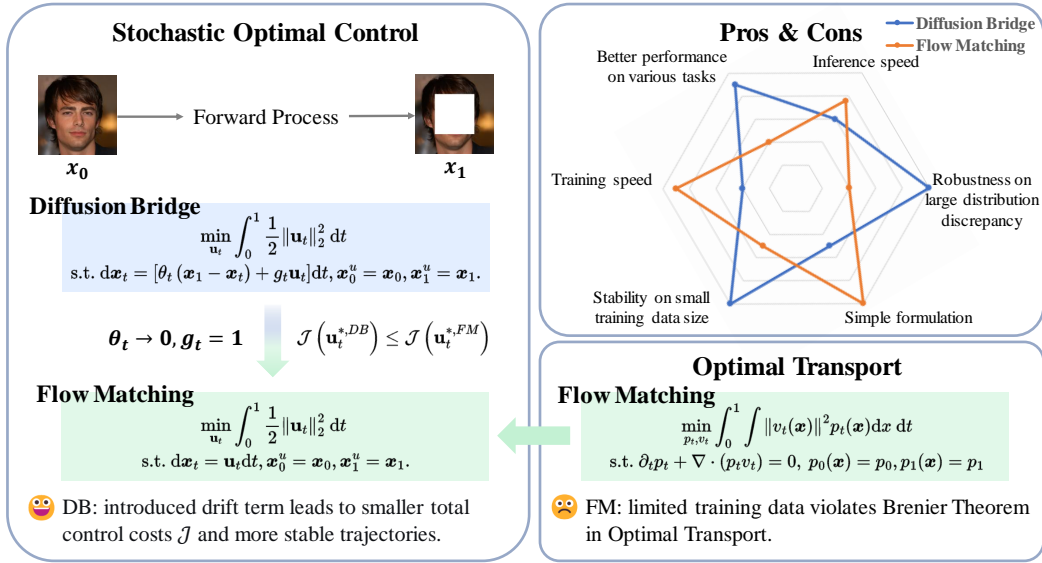


Figure 1: Overview of Diffusion Bridge (DB) versus Flow Matching (FM) on mapping one distribution to another. From the perspective of Stochastic Optimal Control and Optimal Transport, we demonstrate that 1) the cost function of DB is lower than that of FM, implying a more stable and natural trajectory of DB; 2) the linear interpolation scheme in FM (with coefficients t and $1 - t$) becomes suboptimal under limited training data regimes. The respective strengths and weaknesses of DB and FM are summarized in the radar chart.

Flow-based generative models have undergone rapid evolution, progressing from the early invertible transformations of Normalizing Flows (Papamakarios et al., 2021; Albergo & Vanden-Eijnden, 2022a; Mathieu & Nickel, 2020) to the recently emerged paradigm of Flow Matching. Continuous Normalizing Flows rely on repeated calls to high-order ODE solvers during forward-backward propagation, resulting in a training procedure dominated by costly simulation-and-differentiation loops that severely limit scalability. Flow Matching (Lipman et al., 2023; 2024; Liu et al., 2022a; Liu, 2022) circumvents this bottleneck by first designing a probability path that interpolates between a prior and the data distribution, and then directly learning a time-dependent vector field whose integral trajectories realize this path. Its simple modeling principle and high generation quality have attracted significant attention. OT-CFM (Tong et al., 2023) further stabilizes training and improves sample fidelity by performing optimal-transport pairing of noise and data samples before learning the vector field. Non-Euclidean Flow Matching (Chen & Lipman, 2023) extends this framework beyond flat spaces by replacing Euclidean straight-line interpolation with geodesic interpolation on arbitrary Riemannian manifolds, enabling high-quality generation on intrinsically curved domains.

Diffusion Bridges and Flow Matching successfully address a wide spectrum of tasks—including image restoration (Luo et al., 2023b; Yue et al., 2024; Zhu et al., 2025; Martin et al., 2024; Cohen et al., 2025), translation (Zhou et al., 2023; Zheng et al., 2024), text-to-image generation (Liu et al., 2025; Xu et al., 2025), and robotic policy synthesis (Ren et al., 2025; Zhang et al., 2025)—yet a comprehensive theoretical framework that rigorously delineates their mutual relationship, comparative advantages, and inherent limitations remains open. This naturally leads to a fundamental question:

“Mapping one distribution to another, which is better—Diffusion Bridge or Flow Matching?”

In this paper, we firstly theoretically analyze the relationship between Diffusion Bridge and Flow Matching, and experimentally compare their respective advantages and disadvantages in implementing distribution-to-distribution transformation.

Comparative Theoretical Framework for Diffusion Bridge and Flow Matching. From the perspective of Stochastic Optimal Control, we provide a unified framework for Diffusion Bridge and Flow Matching, demonstrating that the cost function of Diffusion Bridge is lower than that of Flow Matching, implying a more stable and natural trajectory of Diffusion Bridge. Furthermore, we analyze the linear interpolation scheme of Flow Matching (governed by coefficients t and $1 - t$) from

an Optimal Transport perspective. Our analysis reveals that this scheme can lead to performance degradation, particularly when the training dataset size is limited.

Comparative Experimental Evaluation of Diffusion Bridge versus Flow Matching. To enable a fair and controlled empirical evaluation, we design a new, powerful neural architecture for Diffusion Bridge models based on a latent Transformer. This architectural innovation significantly enhances the capability of DB models. Using this and an equivalent-structure Flow Matching model, we conduct extensive experiments across a diverse suite of tasks, including Image Inpainting (under varying training data sizes and distributional discrepancies), Super-Resolution, Deblurring, Denoising, and Image-to-Image Translation. This systematic evaluation clearly shows the respective strengths and weaknesses of both methodologies in Figure 1.

2 RELATED WORK

Diffusion Bridge. Diffusion Bridge models are relatively advanced methods for achieving transformation between distributions. Diffusion Schrödinger Bridges (Liu et al., 2023; Shi et al., 2024; De Bortoli et al., 2021) construct a transport mapping between p_{data} and p_{prior} by minimizing the KL divergence π^* , using IPF to alternately optimize boundary conditions. However, it suffers from high computational complexity, especially in high dimensions or with large N , making direct optimization challenging. DDBMs (Zhou et al., 2023) and GOUB (Yue et al., 2024) achieve promising results in tasks such as image restoration and translation by incorporating Doob’s h -transform into the original forward SDE of diffusion models to pin the terminal distribution to a specific target. **While DDBMs establish that Flow Matching is a special case of DDBMs-VE in its zero-noise limit, a systematic comparison of how different diffusion bridge and flow matching formulations affect generative performance remains unexplored.** UniDB (Zhu et al., 2025) reformulates the Diffusion Bridge as a Stochastic Optimal Control, proving that Doob’s h -transform is a special case when the terminal penalty coefficient $\gamma \rightarrow \infty$, thereby unifying and generalizing existing diffusion bridges.

Flow Matching. Flow-based generative models have evolved from Continuous Normalizing Flows (CNFs) (Papamakarios et al., 2021; Albergo & Vanden-Eijnden, 2022a; Mathieu & Nickel, 2020), which treat data generation as an invertible, continuous-time transformation governed by an ordinary differential equation (ODE). CNFs admit exact likelihood evaluation and an invertible mapping, but their integration cost scales with dimension and integration steps, making training and sampling expensive. To overcome this computational bottleneck, Flow Matching (Lipman et al., 2023; 2024; Liu et al., 2022a) trains the model by regressing onto the conditional optimal vector field that maps distributions to distributions paired with their conditions; the resulting conditional vector field can be computed in closed form without simulating the ODE. While the conditional path is a straight-line interpolation in Euclidean space, such paths are not optimal in terms of transport cost. Optimal Transport Conditional Flow Matching (OT-CFM) (Tong et al., 2023) replaces the independent coupling with a static optimal-transport plan, yielding conditional trajectories that are globally straighter and reduce curvature. This mitigates the back-tracking phenomenon observed in diffusion-like schedules and accelerates both training and sampling.

3 PRELIMINARIES

3.1 STOCHASTIC OPTIMAL CONTROL

Stochastic Optimal Control (SOC) offers a principled methodology for designing optimal policies in dynamical systems operating under uncertainty. Its core objective is to derive the optimal control laws that remain effective in stochastic environments, with applications spanning multiple disciplines (Geering et al., 2010; Rout et al., 2024; Zhu et al., 2025; Pandey et al., 2025b). We consider the following SOC formulation with quadratic costs (Kappen, 2008; Chen et al., 2023a):

$$\begin{aligned} \min_{\mathbf{u}_{t,\gamma} \in \mathcal{U}} \mathbb{E} \left[\int_0^T \frac{1}{2} \|\mathbf{u}_{t,\gamma}\|_2^2 dt + \frac{\gamma}{2} \|\mathbf{x}_T^u - \mathbf{x}_T\|_2^2 \right] \\ \text{s.t. } d\mathbf{x}_t^u = [\mathbf{f}(\mathbf{x}_t^u, t) + g_t \mathbf{u}_{t,\gamma}] dt + g_t d\mathbf{w}_t, \mathbf{x}_0^u = \mathbf{x}_0, \end{aligned} \quad (1)$$

where \mathbf{x}_t^u is the diffusion process under control, \mathbf{x}_0 and \mathbf{x}_T denote the predetermined initial and terminal states, respectively, $\|\mathbf{u}_{t,\gamma}\|_2^2$ is the path controlled cost, and $\frac{\gamma}{2}\|\mathbf{x}_T^u - \mathbf{x}_T\|_2^2$ corresponds to the terminal cost with its penalty coefficient γ . The expectation value is over all stochastic trajectories originating from \mathbf{x}_0 (Kappen, 2008). The whole SOC problem (1) aims to design the controller $\mathbf{u}_{t,\gamma}$ to drive the dynamic from \mathbf{x}_0 to \mathbf{x}_T while minimizing the overall cost.

3.2 FLOW MATCHING

Notably, most Flow Matching (FM) work (Lipman et al., 2022b; 2024; Liu et al., 2022a; Albergo & Vanden-Eijnden, 2022b) directly models the generative process, which is different from standard diffusion and diffusion bridges (Ho et al., 2020; Song et al., 2020; Zhou et al., 2023). To avoid discrepancy, we decouple FM with both forward and sampling processes (Liu et al., 2022a) in the context of diffusion models. The FM forward process is to transform samples $\mathbf{x}_0 \in \mathbb{R}^d$ from a source (data) distribution p_0 into $\mathbf{x}_1 \in \mathbb{R}^d$ from a target (prior) distribution p_1 by defining a flow $\psi_t : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ parameterized by a learnable vector field $v_t(\mathbf{x})$. Due to the intractability of the true vector field in practice, the conditional flow $\psi_t(\mathbf{x} | \mathbf{x}_1)$ is constructed under the probability paths $p_t(\mathbf{x} | \mathbf{x}_1)$ with the vector field $v_t(\mathbf{x} | \mathbf{x}_1)$ (Lipman et al., 2022b; Liu et al., 2022a). The conditional flow matching (Lipman et al., 2022a) training objective is formed as

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_t} \left[\|\hat{v}_\theta(\mathbf{x}_t, t) - (\mathbf{x}_1 - \mathbf{x}_0)\|^2 \right]. \quad (2)$$

To find a useful conditional flow $\psi_t(\mathbf{x} | \mathbf{x}_1)$, one popular example is the minimizer of the dynamic Optimal Transport (OT) problem with quadratic cost (Villani et al., 2008; Villani, 2021b; Peyré et al., 2019), which is formalized as

$$\begin{aligned} \min_{p_t, v_t} \int_0^1 \int \|v_t(\mathbf{x})\|^2 p_t(\mathbf{x}) d\mathbf{x} dt \\ \text{s.t. } \partial_t p_t + \nabla \cdot (p_t v_t) = 0, p_0(\mathbf{x}) = p_0, p_1(\mathbf{x}) = p_1 \end{aligned} \quad (3)$$

where the objective is the Wasserstein-2 transport cost between two distributions and the conditions are that the vector field v_t satisfies the continuity equation and the initial and terminal conditions. The solution (p_t^*, v_t^*) defines a flow called an OT displacement interpolant as $\psi_t^*(x) = tT(x) + (1-t)x$ where $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the OT map such that $T\#p_0 = p_1$ (i.e., p_1 is the push-forward of p_0 under the transport map T). By defining the random variable $\mathbf{x}_t = \psi_t^*(\mathbf{x}_0)$ with $\mathbf{x}_0 \sim p_0$ and finding a bound for the Wasserstein-2 transport cost (Liu et al., 2022a; Lipman et al., 2024), the dynamic OT problem leads to the following variational problem for $\eta_t = \psi_t(\mathbf{x}_0 | \mathbf{x}_1)$:

$$\min_{\eta_t} \int_0^1 \|\dot{\eta}_t\|^2 dt \quad \text{s.t. } \eta_0 = \mathbf{x}_0, \eta_1 = \mathbf{x}_1. \quad (4)$$

According to the Euler-Lagrange equations (Gelfand et al., 2000), the minimizer can be obtained as

$$\mathbf{x}_t = \psi_t^*(\mathbf{x}_0 | \mathbf{x}_1) = \eta_t^* = t\mathbf{x}_1 + (1-t)\mathbf{x}_0. \quad (5)$$

For backward sampling, after we get \hat{v}_θ , we solve the ODE $d\tilde{\mathbf{x}}_t = -\hat{v}_\theta dt$ starting from $\tilde{\mathbf{x}}_0 \sim p_1$ to transfer p_1 to p_0 and set $\tilde{\mathbf{x}}_t = \mathbf{x}_{1-t}$ according to the time-symmetric property (Liu et al., 2022a). For more details, please refer to Lipman et al. (2024).

3.3 DIFFUSION BRIDGES

Doob’s h -transform-based diffusion bridges, such as DDBMs (Zhou et al., 2023) and GOUB (Yue et al., 2024), achieve promising results by modifying the original forward SDE to pass through the predetermined terminal. UniDB further constructs diffusion bridges through the Stochastic Optimal Control (SOC) problem and generalizes these h -transform-based methods. Based on the Generalized Ornstein-Uhlenbeck (GOU) process (Ahmad, 1988; Yue et al., 2024),

$$d\mathbf{x}_t = \theta_t (\boldsymbol{\mu} - \mathbf{x}_t) dt + g_t d\mathbf{w}_t, \quad (6)$$

where $\boldsymbol{\mu}$ is a given state vector, \mathbf{w}_t is the Wiener process, θ_t and g_t denote the scalar drift and diffusion coefficient with the relationship $g_t^2 = 2\lambda^2\theta_t$ where λ^2 is a given positive constant scalar, a

specific example of the UniDB framework is introduced as UniDB-GOU (Zhu et al., 2025) and the related SOC problem constructs the forward bridge process as:

$$\begin{aligned} & \min_{\mathbf{u}_{t,\gamma}} \mathbb{E} \left[\int_0^1 \frac{1}{2} \|\mathbf{u}_{t,\gamma}\|_2^2 dt + \frac{\gamma}{2} \|\mathbf{x}_1^u - \mathbf{x}_1\|_2^2 \right] \\ & \text{s.t. } d\mathbf{x}_t^u = [\theta_t(\mathbf{x}_1 - \mathbf{x}_t^u) + g_t \mathbf{u}_{t,\gamma}] dt + g_t d\mathbf{w}_t, \mathbf{x}_0^u = \mathbf{x}_0, \end{aligned} \quad (7)$$

where $\boldsymbol{\mu} = \mathbf{x}_1$ in the SDE is set as the final condition. To be consistent with flow matching above, here we change the original time schedule $t \in [0, T]$ to $t \in [0, 1]$. According to the certainty equivalence principle (Chen et al., 2023b; Rout et al., 2024), UniDB derives the same optimal controller $\mathbf{u}_{t,\gamma}^*$ by modifying the SOC problem (7) to one with the deterministic ODE condition as follows, specifically,

$$\begin{aligned} & \min_{\mathbf{u}_{t,\gamma}} \int_0^1 \frac{1}{2} \|\mathbf{u}_{t,\gamma}\|_2^2 dt + \frac{\gamma}{2} \|\mathbf{x}_1^u - \mathbf{x}_1\|_2^2 \\ & \text{s.t. } d\mathbf{x}_t^u = [\theta_t(\mathbf{x}_1 - \mathbf{x}_t^u) + g_t \mathbf{u}_{t,\gamma}] dt, \mathbf{x}_0^u = \mathbf{x}_0. \end{aligned} \quad (8)$$

Previous works like GOUB (Yue et al., 2024) can be considered as taking $\gamma \rightarrow \infty$ in (8), which means the controlled dynamics would pass precisely through the preset terminal \mathbf{x}_1 (Chen et al., 2023b), therefore, we can transform the SOC problem (8) with $\gamma \rightarrow \infty$ into the following form:

$$\begin{aligned} & \min_{\mathbf{u}_t} \int_0^1 \frac{1}{2} \|\mathbf{u}_t\|_2^2 dt \\ & \text{s.t. } d\mathbf{x}_t^u = [\theta_t(\mathbf{x}_1 - \mathbf{x}_t^u) + g_t \mathbf{u}_t] dt, \mathbf{x}_0^u = \mathbf{x}_0, \mathbf{x}_1^u = \mathbf{x}_1. \end{aligned} \quad (9)$$

UniDB underscores the equivalence between the SOC formulation under this limiting behavior and Doob's h -transform (Särkkä & Solin, 2019). By solving the problem (8), the closed-form optimal controller $\mathbf{u}_{t,\gamma}^*$ can be obtained. For more details, please refer to Zhu et al. (2025).

4 COMPARATIVE THEORETICAL FRAMEWORK FOR DIFFUSION BRIDGE AND FLOW MATCHING

Recent research (Gao et al., 2025) has mainly clarified the exact equivalence between diffusion models and flow matching (for the special case that the source distribution corresponds to a Gaussian). However, there remains confusion about the connections between diffusion bridges and flow matching, since both two can achieve the transition between two arbitrary paired distributions. In the following analysis, we adopt the forward process of the UniDB-GOU model (denoted DB hereafter) and the well-known Flow Matching model (denoted FM hereafter) introduced in Section 3 above as the main diffusion bridge and flow matching model, respectively.

4.1 CONNECTIONS BETWEEN DIFFUSION BRIDGES AND FLOW MATCHING

To compare the two models within the same theoretical framework, we first consider the construction of the FM-related SOC problem to be consistent with the formulation of DB. Taking the substitution $\mathbf{u}_t = \dot{\eta}_t$ and $\mathbf{x}_t^u = \eta_t$ in the variational problem (4) leading to

$$\begin{aligned} & \min_{\mathbf{u}_t} \int_0^1 \frac{1}{2} \|\mathbf{u}_t\|_2^2 dt \\ & \text{s.t. } d\mathbf{x}_t^u = \mathbf{u}_t dt, \mathbf{x}_0^u = \mathbf{x}_0, \mathbf{x}_1^u = \mathbf{x}_1, \end{aligned} \quad (10)$$

with the optimal controller

$$\mathbf{u}_t^{*,\text{FM}} = \frac{d\eta_t^*}{dt} = \mathbf{x}_1 - \mathbf{x}_0 = \frac{\mathbf{x}_1 - \mathbf{x}_t^u}{1-t}, \quad (11)$$

where the last equation comes from the interpolation (5) (Lipman et al., 2024). Despite the derivation of this FM's SOC problem, a direct comparison between the two SOC frameworks requires addressing a fundamental discrepancy in the formulation of their optimization objectives: the objective of the SDE-based SOC problem involves the expectation over all stochastic trajectories, whereas the ODE-based one does not, due to its determinism (Kappen, 2008). Considering that 1) the derivation of

DB’s optimal controller often relies on the ODE-based SOC formulation by the certainty equivalence principle (Chen et al., 2023b; Rout et al., 2024; Zhu et al., 2025) in practice and 2) FM (10) forces exact target matching with $x_1^u = x_1$ which is the same as in (9), we adopt the ODE-based SOC problem (9) of DB to facilitate a fair comparison with the same ODE-based formulation of FM (10).

We can easily find the relation between (9) and (10) that the FM’s SOC problem is a special case of DB’s one, which leads to the following proposition:

Proposition 4.1. *Under the conditions that $\theta_t \rightarrow 0$ and $g_t = 1$ in (9), Diffusion Bridge degrades to Flow Matching.*

Details are provided in the Appendix A.1. This proposition indicates that under specific parameter constraints—namely, zero drift coefficient θ_t and unit diffusion coefficient g_t —the SOC formulation of DB (9) reduces to that of FM (10). The key distinction lies in the structure of the drift term: DB incorporates the drift term of the form $\theta_t(x_1 - x_t^u)$, which is absent in FM. We further analyze the role of this drift term and present the following theorem to demonstrate that this drift term contributes to reducing the total cost (objective function) in the SOC problem.

Theorem 4.2. *Denote $\mathcal{J}(\mathbf{u}_t) \triangleq \int_0^1 \frac{1}{2} \|\mathbf{u}_t\|_2^2 dt$ as the overall cost of the SOC problems (9) and (10) with the related optimal controller $\mathbf{u}_t^{*,DB}$ and $\mathbf{u}_t^{*,FM}$, respectively. Under the condition of diffusion coefficient $g_t = 1$ in (9) to be consistent with FM (10), then*

$$\mathcal{J}(\mathbf{u}_t^{*,DB}) \leq \mathcal{J}(\mathbf{u}_t^{*,FM}). \quad (12)$$

Please refer to Appendix A.2 for detailed proof. Theorem 4.2 mainly emphasizes that the drift term $\theta_t(x_1 - x_t^u)$ introduced in DB actively guides the system toward more stable trajectories, thereby lowering the total cost in the SOC problem.

Lower total costs typically lead to smoother and more natural SDE/ODE trajectories. On one hand, a larger total cost implies the larger controller $\|\mathbf{u}_t^*\|_2^2$, which in turn excites oscillations along the forward trajectory and may disrupt the inherent continuity and smoothness of images. Consequently, the state x_t undergoes violent fluctuations at every time step, destabilizing the evolution of each individual pixel and inevitably degrading the visual quality of the generated image. On the other hand, from the SOC viewpoint, FM employs the trivial forward ODE, which may be too simple to characterise the transition between arbitrary distributions. As the discrepancy between the terminal distributions grows, the neural network becomes harder to fit the trajectories. However, DB incorporates the drift term $\theta_t(x_1 - x_t^u)$, which helps to construct a more stable transition between two arbitrary distributions.

4.2 OVERSIMPLIFICATION AND INEFFECTIVENESS OF FM’S INTERPOLATION COEFFICIENT

We have shown that the absence of the drift term in the forward process significantly amplifies the total control cost within the SOC framework, concurrently appearing as less stable trajectories. In this section, we examine from an Optimal Transport (OT) perspective that the interpolation coefficients t and $1 - t$ in (5) of FM represent an oversimplification of McCann’s approach, and it would become ineffective when the training data size gradually decreases.

Oversimplification. Under the conditions where the Brenier Theorem holds (Santambrogio, 2015; Lei & Gu, 2021) (the source measure is absolutely continuous), the McCann interpolation in OT (Villani, 2021a) defines a deterministic dynamical system via the OT map $T : \Omega \rightarrow \Omega$, such that $T_{\#}x_0 = x_1$. Its Lagrangian trajectory and the corresponding Eulerian velocity field are given by:

$$v_{OT}(x, t) = T(\tilde{x}_0(x, t)) - \tilde{x}_0(x, t), \quad (13)$$

where $\tilde{x}_0(x, t)$ is implicitly defined by the equation $x = (1 - t)x_0 + tT(x_0)$, representing a complex inverse problem. In contrast, Flow Matching (FM) posits oversimplification by replacing the globally optimal coupling T with a probabilistic coupling $q(x_0, x_1)$. For a sample $(x_0, x_1) \sim q$, FM defines a conditional flow and its velocity field as:

$$v_{CFM}(x(t) | x_0, x_1) = x_1 - x_0. \quad (14)$$

The genuine McCann interpolant of OT allows the velocity field to vary with t while still satisfying the continuity equation, thereby minimizing the total kinetic energy. In contrast, FM’s fixed interpolation

coefficients, t and $1 - t$, freeze the velocity field, which remains viable under small distributional discrepancies but leads to rapid performance degradation as the discrepancy grows. This failure stems from its inability to capture the complex inter-manifold geometry, triggering a large number of streamline intersections and vector field conflicts, ultimately leading to blurry learning objectives and a significant decrease in generation quality of the model. (Liu et al., 2022b)

Ineffective with scarce data. In practical learning scenarios, the empirical measures are constructed from paired training data, which are discrete measures consisting of finitely many $n \in \mathbb{Z}^+$ sample points $\{(x_i, y_i)\}$ and take the form

$$\hat{\mu}_0^{(n)} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}, \quad \hat{\mu}_1^{(n)} = \frac{1}{n} \sum_{i=1}^n \delta_{y_i}, \quad (15)$$

where $\hat{\mu}_0$ and $\hat{\mu}_1$ represent the source and target measures, respectively. Crucially, this discrete setting diverges fundamentally from the continuous assumptions (the absolutely continuous source measure) discussed above in the Brenier potential theorem, which leads to the following remark:

Remark 4.3. For the finite empirical measures, they violate the absolute-continuity assumption required by Brenier’s theorem; consequently, the existence of a convex potential function pushing $\hat{\mu}_0^{(n)}$ to $\hat{\mu}_1^{(n)}$ cannot be guaranteed, and McCann’s interpolation is no longer well-defined.

Please refer to Appendix A.3 for detailed derivation. Remark 4.3 shows that

- When training data are limited, the interpolation coefficients t and $1 - t$ lose their validity and the resulting path no longer corresponds to an OT interpolation.
- As the training data size grows ($n \rightarrow \infty$), then $\hat{\mu}_0^{(n)} \rightharpoonup \mu_0, \hat{\mu}_1^{(n)} \rightharpoonup \mu_1$ weakly and the empirical interpolation regains absolute continuity; McCann’s interpolation is asymptotically restored, and the performance of FM improves. However, the coefficients remain rigidly fixed at t and $1 - t$, which cannot truly express the local velocity adjustment of optimal transmission. As the discrepancy between distributions widens and the task complexity escalates, the error rate climbs rapidly.

Simply adding more data reduces but cannot eliminate the systematic error of the frozen vector field. A learnable drift term $\theta_t(\mathbf{x}_1 - \mathbf{x}_t^u)$, as incorporated by DB, is required to approximate the true Wasserstein geodesic in regimes of scarce data or complex deformation.

5 COMPARATIVE EXPERIMENTAL EVALUATION OF DIFFUSION BRIDGE VERSUS FLOW MATCHING

In this section, we evaluate the performance of Diffusion Bridge (DB) versus Flow Matching (FM) through different image-to-image tasks including different Image Restoration tasks (Image Inpainting, Image 4×Super-Resolution, Image Deblurring and Image Denoising). For related implementation details, additional experimental and additional visual results, please refer to Appendix E, F, and H, respectively.

For image restoration tasks, we evaluated the two models on the CelebA-HQ 256×256 (Karras, 2017) dataset. We take four common image evaluation metrics: Peak Signal-to-Noise Ratio (PSNR, higher is better) (Fardo et al., 2016), Structural Similarity Index (SSIM, higher is better) (Wang et al., 2004), Learned Perceptual Image Patch Similarity (LPIPS, lower is better) (Zhang et al., 2018), and Fréchet Inception Distance (FID, lower is better) (Heusel et al., 2017).

5.1 EXPERIMENT SETUP: SAME TRANSFORMER ARCHITECTURE

FM has now widely adopted both U-Net (Lipman et al., 2024) and Transformer (Dao et al., 2023; Hu et al., 2024) as its training network architecture, while the previous DB relied only on the U-Net network (Zhou et al., 2023; Yue et al., 2024; Zhu et al., 2025). In practice, substantial variations in U-Net implementations—particularly in parameter count and architectural details—make it difficult to conduct a direct comparative evaluation. To ensure a fair empirical comparison between DB and FM, we present a Latent Transformer-based Network for DB built on the DiT architecture (Peebles

& Xie, 2023; Ma et al., 2024), unifying the network backbone of the two models. Please refer to Appendix E for detailed parameters of the network.

Taking an image restoration task as an example, the overall training procedure can be summarized as follows. For paired Low-Quality (LQ) and High-Quality (HQ) images, denoted x^{LQ} and x^{HQ} with dimension $(h \times w \times 3)$, respectively, we map them into the latent space with dimension $(\frac{h}{8} \times \frac{w}{8} \times 4)$ using a pre-trained VAE encoder with a 8 down-sampling factor adopted from Stable Diffusion (Rombach et al., 2022). Although no hard constraint forces the latent codes produced by the VAE to lie in a compact set, we prove that, even in non-compact spaces, McCann interpolation remains universally valid for arbitrary distributional transport. Please refer to Appendix A.4 for details (Lei & Gu, 2021). The latent codes are then input into the DiT block, which predicts either the score/noise (for DB) or the velocity (for FM) conditioned on the current time step, and the models are trained by minimizing the respective objective functions.

5.2 DB VS FM ON DIFFERENT TASKS

Here we demonstrate some experiments between the two models: Inpainting with a centered 64×64 and 128×128 box mask, $4 \times$ Super-Resolution using bicubic down-sampling, Deblurring using a 15×15 and 61×61 Gaussian kernel, and Denoising with Gaussian noise ($\sigma = 35$ over 255). The quantitative and qualitative results are illustrated in Table 1 and Figure 2. It can be concluded that DB consistently demonstrates significantly superior perceptual scores (LPIPS and FID) across all tasks, whereas FM exhibits relatively stronger performance on pixel-level metrics (SSIM). Visually, DB achieves more realistic images than FM across the majority of tasks and avoids some unnatural patterns in smooth regions.

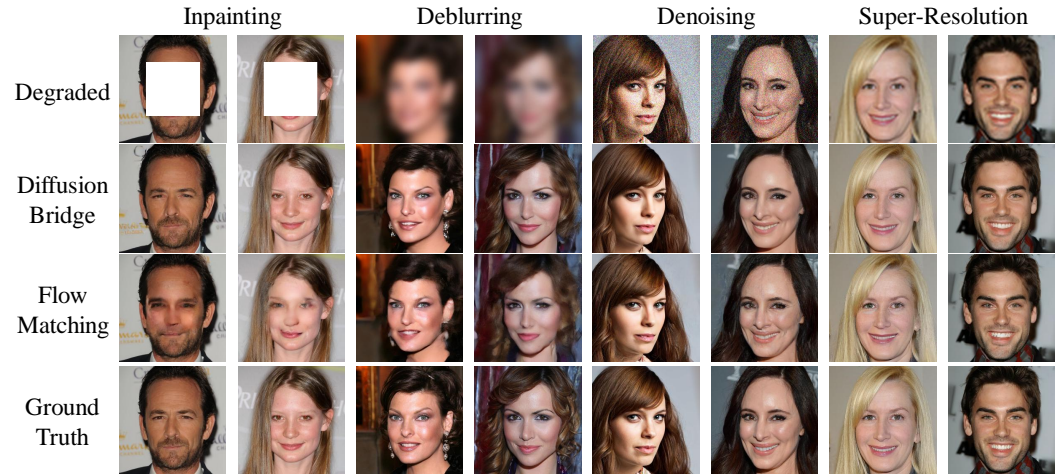


Figure 2: Qualitative comparison of visual results between Diffusion Bridge and Flow Matching on different Image Restoration tasks including Inpainting with a centered 128×128 box mask, Deblurring with a 61×61 Gaussian kernel, Denoising, and $4 \times$ Super-Resolution.

Table 1: Quantitative results for Flow Matching and Diffusion Bridge (denoted FM and DB in table, respectively) under different image restoration tasks on the CelebA-HQ dataset.

Method	Inpainting-Box64				Inpainting-Box128				4xSuper-Resolution			
	PSNR↑	SSIM↑	LPIPS↓	FID↓	PSNR↑	SSIM↑	LPIPS↓	FID↓	PSNR↑	SSIM↑	LPIPS↓	FID↓
FM	28.03	0.840	0.039	5.13	23.54	0.760	0.106	17.84	27.11	0.789	0.088	11.61
DB	27.90	0.813	0.038	5.11	23.57	0.741	0.078	7.71	27.47	0.762	0.077	8.50
Method	Deblurring (15 × 15)				Deblurring (61 × 61)				Denoising			
	PSNR↑	SSIM↑	LPIPS↓	FID↓	PSNR↑	SSIM↑	LPIPS↓	FID↓	PSNR↑	SSIM↑	LPIPS↓	FID↓
FM	27.24	0.793	0.088	10.49	24.67	0.683	0.228	38.18	27.02	0.760	0.093	16.04
DB	27.26	0.757	0.087	8.77	24.13	0.661	0.172	19.03	27.30	0.757	0.086	10.16

5.3 DB VS FM UNDER VARYING LEVELS OF TASK DIFFICULTY

To assess the robustness of the two models under varying levels of task difficulty, we perform a series of Image Inpainting tasks with center box masks ranging from 50×50 to 128×128 . Specifically, in image inpainting tasks, enlarging the center box mask amplifies the distance between the reference and target distributions, which represents the increasing levels of difficulty. We report the results of the LPIPS and FID scores of two models in Table 2. When the task is relatively straightforward (Box50 and Box64), FM achieves a performance comparable to that of DB. As the distributional discrepancy increases, FM exhibits markedly degraded stability compared to DB. Under this progressive shift, the two perceptual scores (LPIPS and FID) for FM, particularly in FID, increase significantly faster than those for DB, indicating a pronounced sensitivity to growing distributional shifts for FM and robustness on distribution discrepancy from the introduced drift term in DB, which validates the conclusion in Section 4.1.

Table 2: Quantitative results for Flow Matching and Diffusion Bridge (denoted FM and DB in table, respectively) under Image Inpainting tasks with different center box masks ranging from 50×50 to 128×128 on the CelebA-HQ dataset.

Method	Box50		Box64		Box72		Box80		Box96		Box128	
	LPIPS↓	FID↓	LPIPS↓	FID↓	LPIPS↓	FID↓	LPIPS↓	FID↓	LPIPS↓	FID↓	LPIPS↓	FID↓
FM	0.035	4.93	0.039	5.13	0.042	5.43	0.047	5.86	0.060	8.18	0.106	17.84
DB	0.035	4.93	0.038	5.11	0.041	5.25	0.044	5.34	0.052	6.25	0.078	7.71

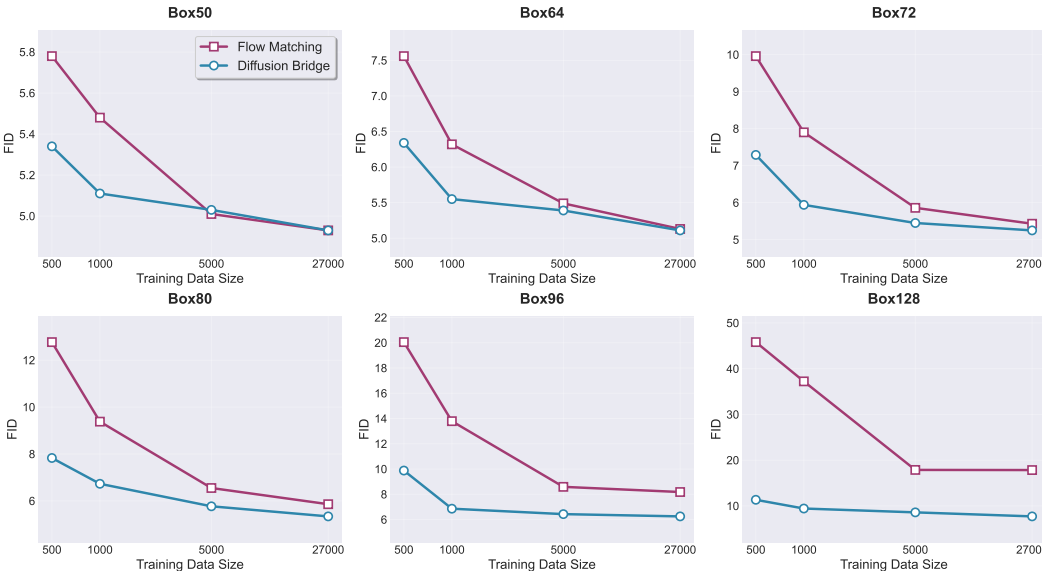


Figure 3: The performance (FID) of Flow Matching and Diffusion Bridge under Image Inpainting tasks with different center box masks and training data sizes ranging from 500 to 27000.

5.4 DB VS FM UNDER DIFFERENT TRAINING DATA SIZE

Further, we systematically conducted several Image Inpainting tasks on both DB and FM models under four different increasing training data sizes from 500, 1000, 5000, to 27000 to justify that the interpolation coefficients t and $1 - t$ in (5) of FM would become ineffective when the training data size gradually decreases, as mentioned in Section 4.2. The quantitative results are illustrated in Figure 3, which reveals that although the performance of both DB and FM deteriorates as training data become scarce, the degradation of FM is markedly steeper, while DB remains stable and maintains a consistently high performance level.

5.5 DB VS FM ON TRAINING AND INFERENCE TIME

In addition, we conducted a systematic profiling of training and wall-clock inference latency for the tasks above. As depicted in Figure 4, under identical network architectures and experimental hyperparameter settings, after convergence, the training time of FM is substantially shorter than that of DB. FM admits a short, direct gradient path that bypasses intermediate numerical complications, yielding superior numerical stability. In contrast, DB must predict a time-varying field, which is intrinsically difficult to fit, and gradients must back-propagate through a chain of compositional functions, leading to slower empirical convergence rate. Under the same generation conditions with the same Number of Function Evaluations (NFEs), the inference time of FM is basically similar to DB when the CPU frequency is relatively high. A detailed analysis of how CPU frequency affects inference latency is provided in Appendix G.

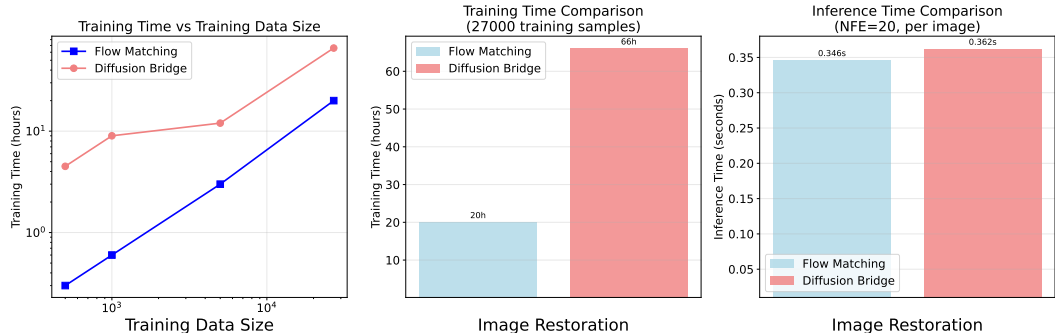


Figure 4: Training (with training data size 27000) and inference time (NFEs = 20) comparison on different Image Restoration tasks (Inpainting, Deblurring, Denoising, and Super-Resolution) between Diffusion Bridge and Flow Matching.

6 CONCLUSION

In this paper, we conduct a comprehensive analysis of Diffusion Bridge (DB) and Flow Matching (FM), the two most advanced methodologies currently available for distribution-to-distribution transformation. Under the unifying perspective of Stochastic Optimal Control (SOC), we establish that FM constitutes a special case of the DB framework, and we rigorously prove that the cost functional of FM is strictly larger than that minimized by DB because of the absence of the drift term in the forward process, thereby suggesting a less stable controlled trajectory. Concurrently, from the perspective of Optimal Transport, we demonstrate that the linear interpolation coefficients t and $1 - t$ employed by FM violate the Brenier-potential theorem when the training data size is reduced, causing the algorithm to collapse and its performance to degrade precipitously. For the first time, we implement a Latent Transformer-based Diffusion Bridge and conduct experiments under the same architecture as Flow Matching, evaluating their performance through different image-to-image tasks including different Image Restoration (Image Inpainting, Image 4x Super-Resolution, Image Deblurring, and Image Denoising), further confirming the respective advantages and disadvantages of the two models. Our theoretical research also has important reference significance for other fields, such as embodied AI, medical imaging, etc., on how to choose models that can achieve both performance and efficiency when training data size is small.

REFERENCES

- 540
541
542 Suzan Ece Ada, Erhan Oztop, and Emre Ugur. Diffusion Policies for Out-of-Distribution Generaliza-
543 tion in Offline Reinforcement Learning. *IEEE Robotics and Automation Letters*, 9(4):3116–3123,
544 2024.
- 545 R Ahmad. Introduction to Stochastic Differential Equations, 1988.
- 546 Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants.
547 *arXiv preprint arXiv:2209.15571*, 2022a.
- 548
549 Michael Samuel Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic
550 interpolants. In *The Eleventh International Conference on Learning Representations*, 2022b.
- 551
552 Ricky TQ Chen and Yaron Lipman. Flow matching on general geometries. *arXiv preprint*
553 *arXiv:2302.03660*, 2023.
- 554 Tianrong Chen, Jiatao Gu, Laurent Dinh, Evangelos A Theodorou, Joshua Susskind, and Shuangfei
555 Zhai. Generative modeling with phase stochastic bridges. *arXiv preprint arXiv:2310.07805*, 2023a.
- 556
557 Tianrong Chen, Jiatao Gu, Laurent Dinh, Evangelos A Theodorou, Joshua Susskind, and Shuangfei
558 Zhai. Generative Modeling with Phase Stochastic Bridges. *arXiv preprint arXiv:2310.07805*,
559 2023b.
- 560 Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake,
561 and Shuran Song. Diffusion Policy: Visuomotor Policy Learning via Action Diffusion. *The*
562 *International Journal of Robotics Research*, pp. 02783649241273668, 2023.
- 563
564 Hyungjin Chung, Jeongsol Kim, Michael T McCann, Marc L Klasky, and Jong Chul Ye. Diffusion
565 Posterior Sampling for General Noisy Inverse Problems. In *International Conference on Learning*
566 *Representations*, 2023.
- 567
568 Elad Cohen, Idan Achituve, Idit Diamant, Arnon Netzer, and Hai Victor Habi. Efficient Image
569 Restoration via Latent Consistency Flow Matching. *arXiv preprint arXiv:2502.03500*, 2025.
- 570
571 Quan Dao, Hao Phung, Binh Nguyen, and Anh Tran. Flow matching in latent space. *arXiv preprint*
572 *arXiv:2307.08698*, 2023.
- 573
574 Valentin De Bortoli, James Thornton, Jeremy Heng, and Arnaud Doucet. Diffusion schrödinger
575 bridge with applications to score-based generative modeling. *Advances in Neural Information*
576 *Processing Systems*, 34:17695–17709, 2021.
- 577
578 Shutong Ding, Ke Hu, Zhenhao Zhang, Kan Ren, Weinan Zhang, Jingyi Yu, Jingya Wang, and Ye Shi.
579 Diffusion-based Reinforcement Learning via Q-weighted Variational Policy Optimization. In
580 *Advances in Neural Information Processing Systems*, 2024.
- 581
582 Shutong Ding, Ke Hu, Shan Zhong, Haoyang Luo, Weinan Zhang, Jingya Wang, Jun Wang, and
583 Ye Shi. Genpo: Generative Diffusion Models Meet On-Policy Reinforcement Learning. *arXiv*
584 *preprint arXiv:2505.18763*, 2025.
- 585
586 Carles Domingo-Enrich, Michal Drozdal, Brian Karrer, and Ricky TQ Chen. Adjoint matching:
587 Fine-tuning flow and diffusion generative models with memoryless stochastic optimal control.
588 *arXiv preprint arXiv:2409.08861*, 2024.
- 589
590 Fernando A Fardo, Victor H Conforto, Francisco C de Oliveira, and Paulo S Rodrigues. A Formal
591 Evaluation of PSNR as Quality Measurement Parameter for Image Segmentation Algorithms.
592 *arXiv preprint arXiv:1605.07116*, 2016.
- 593
594 Ruiqi Gao, Emiel Hoogeboom, Jonathan Heek, Valentin De Bortoli, Kevin Patrick Murphy, and Tim
595 Salimans. Diffusion models and gaussian flow matching: Two sides of the same coin. In *The*
596 *Fourth Blogpost Track at ICLR 2025*, 2025.
- 597
598 Hans P Geering, Florian Herzog, and Gabriel Dondi. Stochastic Optimal Control with Applications in
599 Financial Engineering. *Optimization and Optimal Control: Theory and Applications*, pp. 375–408,
600 2010.

- 594 Izrail Moiseevitch Gelfand, Richard A Silverman, et al. *Calculus of variations*. Courier Corporation,
595 2000.
- 596
- 597 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs
598 Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. *Advances in*
599 *Neural Information Processing Systems*, 30, 2017.
- 600
- 601 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. *Advances in*
602 *Neural Information Processing Systems*, 33:6840–6851, 2020.
- 603
- 604 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet.
605 Video Diffusion Models. *Advances in Neural Information Processing Systems*, 35:8633–8646,
606 2022.
- 607
- 608 Vincent Tao Hu, Wei Zhang, Meng Tang, Pascal Mettes, Deli Zhao, and Cees Snoek. Latent space
609 editing in transformer-based flow matching. In *Proceedings of the AAAI conference on artificial*
610 *intelligence*, volume 38, pp. 2247–2255, 2024.
- 611
- 612 HJ Kappen. Stochastic Optimal Control Theory. *International Conference on Machine Learning,*
613 *Helsinki, Radbound University, Nijmegen, Netherlands*, 2008.
- 614
- 615 Tero Karras. Progressive Growing of GANs for Improved Quality, Stability, and Variation. *arXiv*
616 *preprint arXiv:1710.10196*, 2017.
- 617
- 618 Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising Diffusion Restoration
619 Models. *Advances in Neural Information Processing Systems*, 35:23593–23606, 2022.
- 620
- 621 Diederik P Kingma. ADAM: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
622 2014.
- 623
- 624 N Lei and X Gu. Optimal Transportation Theory and Computation, 2021.
- 625
- 626 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching
627 for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022a.
- 628
- 629 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow Matching
630 for Generative Modeling. In *The Eleventh International Conference on Learning Representations*,
631 2022b.
- 632
- 633 Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky TQ Chen,
634 David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code. *arXiv preprint*
635 *arXiv:2412.06264*, 2024.
- 636
- 637 Guan-Horng Liu, Arash Vahdat, De-An Huang, Evangelos A Theodorou, Weili Nie, and Anima
638 Anandkumar. I2SB: image-to-image Schrödinger bridge. In *Proceedings of the 40th International*
639 *Conference on Machine Learning*, pp. 22042–22062, 2023.
- 640
- 641 Qiang Liu. Rectified flow: A marginal preserving approach to optimal transport. *arXiv preprint*
642 *arXiv:2209.14577*, 2022.
- 643
- 644 Qihao Liu, Xi Yin, Alan Yuille, Andrew Brown, and Mannat Singh. Flowing from Words to Pixels:
645 A Noise-Free Framework for cross-modality evolution. In *Proceedings of the Computer Vision*
646 *and Pattern Recognition Conference*, pp. 2755–2765, 2025.
- 647
- 648 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and
649 transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022a.
- 650
- 651 Xingchao Liu, Chengyue Gong, et al. Flow straight and fast: Learning to generate and transfer data
652 with rectified flow. In *The Eleventh International Conference on Learning Representations*, 2022b.

- 648 Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool.
649 Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the*
650 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 11461–11471, 2022.
651
- 652 Zhengxiong Luo, Dayou Chen, Yingya Zhang, Yan Huang, Liang Wang, Yujun Shen, Deli Zhao,
653 Jingren Zhou, and Tieniu Tan. VideoFusion: Decomposed Diffusion Models for High-Quality
654 Video Generation. *arXiv preprint arXiv:2303.08320*, 2023a.
- 655 Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön. Image Restoration
656 with Mean-Reverting Stochastic Differential Equations. *International Conference on Machine*
657 *Learning*, 2023b.
658
- 659 Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and
660 Saining Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant
661 transformers. In *European Conference on Computer Vision*, pp. 23–40. Springer, 2024.
- 662 Ségolène Martin, Anne Gagneux, Paul Hagemann, and Gabriele Steidl. Pnp-flow: Plug-and-play
663 image restoration with flow matching. *arXiv preprint arXiv:2410.02423*, 2024.
664
- 665 Emile Mathieu and Maximilian Nickel. Riemannian continuous normalizing flows. *Advances in*
666 *neural information processing systems*, 33:2503–2515, 2020.
- 667 Yuhang Mei, Mohammad Al-Jarrah, Amirhossein Taghvaei, and Yongxin Chen. Flow matching for
668 stochastic linear control systems. *arXiv preprint arXiv:2412.00617*, 2024.
669
- 670 Mokai Pan, Kaizhen Zhu, Yuexin Ma, Yanwei Fu, Jingyi Yu, Jingya Wang, and Ye Shi. Unidb++:
671 Fast Sampling of Unified Diffusion Bridge. *arXiv preprint arXiv:2505.21528*, 2025.
672
- 673 Kushagra Pandey, Farrin Marouf Sofian, Felix Draxler, Theofanis Karaletsos, and Stephan Mandt.
674 Variational control for guidance in diffusion models. *arXiv preprint arXiv:2502.03686*, 2025a.
- 675 Kushagra Pandey, Farrin Marouf Sofian, Felix Draxler, Theofanis Karaletsos, and Stephan Mandt.
676 Variational control for guidance in diffusion models. *arXiv preprint arXiv:2502.03686*, 2025b.
677
- 678 George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji
679 Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *Journal of*
680 *Machine Learning Research*, 22(57):1–64, 2021.
- 681 Byoungwoo Park, Jungwon Choi, Sungbin Lim, and Juho Lee. Stochastic Optimal Control for
682 Diffusion Bridges in Function Spaces. *arXiv preprint arXiv:2405.20630*, 2024.
683
- 684 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
685 *the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.
- 686 Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data
687 science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
688
- 689 Hao Ren, Yiming Zeng, Zetong Bi, Zhaoliang Wan, Junlong Huang, and Hui Cheng. Prior does
690 matter: Visual navigation via denoising diffusion bridge models. In *Proceedings of the Computer*
691 *Vision and Pattern Recognition Conference*, pp. 12100–12110, 2025.
- 692 Moritz Reuss, Maximilian Li, Xiaogang Jia, and Rudolf Lioutikov. Goal-conditioned Imitation
693 Learning using Score-based Diffusion Policies. *arXiv preprint arXiv:2304.02532*, 2023.
694
- 695 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
696 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
697 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 698 Litu Rout, Yujia Chen, Nataniel Ruiz, Abhishek Kumar, Constantine Caramanis, Sanjay Shakkottai,
699 and Wen-Sheng Chu. RB-Modulation: Training-Free Personalization of Diffusion Models using
700 Stochastic Optimal Control. *arXiv preprint arXiv:2405.17401*, 2024.
701
- Filippo Santambrogio. Optimal transport for applied mathematicians. 2015.

- 702 Simo Särkkä and Arno Solin. *Applied Stochastic Differential Equations*, volume 10. Cambridge
703 University Press, 2019.
- 704
- 705 Yuyang Shi, Valentin De Bortoli, Andrew Campbell, and Arnaud Doucet. Diffusion Schrödinger
706 bridge matching. *Advances in Neural Information Processing Systems*, 36, 2024.
- 707
- 708 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and
709 Ben Poole. Score-based Generative Modeling Through Stochastic Differential Equations. *arXiv
710 preprint arXiv:2011.13456*, 2020.
- 711 Alexander Tong, Kilian Fatras, Nikolay Malkin, Guillaume Hugué, Yanlei Zhang, Jarrid Rector-
712 Brooks, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models
713 with minibatch optimal transport. *arXiv preprint arXiv:2302.00482*, 2023.
- 714 Cédric Villani. *Topics in optimal transportation*, volume 58. American Mathematical Soc., 2021a.
- 715 Cédric Villani. *Topics in optimal transportation*, volume 58. American Mathematical Soc., 2021b.
- 716 Cédric Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2008.
- 717
- 718
- 719 Zhendong Wang, Jonathan J Hunt, and Mingyuan Zhou. Diffusion Policies as an Expressive Policy
720 Class for Offline Reinforcement Learning. *arXiv preprint arXiv:2208.06193*, 2022.
- 721
- 722 Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image Quality Assessment:
723 From Error Visibility to Structural Similarity. *IEEE transactions on image processing*, 13(4):
724 600–612, 2004.
- 725 Shijie Wu, Yihang Zhu, Yunao Huang, Kaizhen Zhu, Jiayuan Gu, Jingyi Yu, Ye Shi, and Jingya
726 Wang. AffordDP: Generalizable Diffusion Policy with Transferable Affordance. *arXiv preprint
727 arXiv:2412.03142*, 2024.
- 728
- 729 Bin Xia, Yulun Zhang, Shiyin Wang, Yitong Wang, Xinglong Wu, Yapeng Tian, Wenming Yang, and
730 Luc Van Gool. DiffIR: Efficient Diffusion Model for Image Restoration. In *Proceedings of the
731 IEEE/CVF International Conference on Computer Vision*, pp. 13095–13105, 2023.
- 732 Pengcheng Xu, Qingnan Fan, Fei Kou, Shuai Qin, Hong Gu, Ruoyu Zhao, Charles Ling, and Boyu
733 Wang. Textualize Visual Prompt for Image Editing via Diffusion Bridge. In *Proceedings of the
734 AAAI Conference on Artificial Intelligence*, volume 39, pp. 21779–21787, 2025.
- 735 Lingxiao Yang, Shutong Ding, Yifan Cai, Jingyi Yu, Jingya Wang, and Ye Shi. Guidance with
736 Spherical Gaussian Constraint for Conditional Diffusion. In *Forty-first International Conference
737 on Machine Learning*, 2024.
- 738
- 739 Conghan Yue, Zhengwei Peng, Junlong Ma, Shiyan Du, Pengxu Wei, and Dongyu Zhang. Image
740 Restoration Through Generalized Ornstein-Uhlenbeck Bridge. In *International Conference on
741 Machine Learning*, 2024.
- 742 Yanjie Ze, Gu Zhang, Kangning Zhang, Chenyuan Hu, Muhan Wang, and Huazhe Xu. 3D Diffusion
743 Policy. *arXiv e-prints*, pp. arXiv–2403, 2024.
- 744
- 745 Qinglun Zhang, Zhen Liu, Haoqiang Fan, Guanghui Liu, Bing Zeng, and Shuaicheng Liu. Flowpolicy:
746 Enabling fast and robust 3d flow-based policy via consistency flow matching for robot manipulation.
747 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 14754–14762,
748 2025.
- 749 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The Unreasonable
750 Effectiveness of Deep Features as a Perceptual Metric. In *Proceedings of the IEEE conference on
751 computer vision and pattern recognition*, pp. 586–595, 2018.
- 752 Kaiwen Zheng, Guande He, Jianfei Chen, Fan Bao, and Jun Zhu. Diffusion Bridge Implicit Models.
753 *International Conference on Learning Representations*, 2024.
- 754
- 755 Linqi Zhou, Aaron Lou, Samar Khanna, and Stefano Ermon. Denoising Diffusion Bridge Models.
arXiv preprint arXiv:2309.16948, 2023.

756 Kaizhen Zhu, Mokai Pan, Yuexin Ma, Yanwei Fu, Jingyi Yu, Jingya Wang, and Ye Shi. UniDB: A Uni-
757 fied Diffusion Bridge Framework via Stochastic Optimal Control. *arXiv preprint arXiv:2502.05749*,
758 2025.
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

810 A PROOF

811 A.1 PROOF OF PROPOSITION 4.1

812 **Proposition 4.1.** *Under the conditions that $\theta_t \rightarrow 0$ and $g_t = 1$ in (9), Diffusion Bridge degrades to*
 813 *Flow Matching.*

814 *Proof.* Recall the DB’s SOC problem (9):

$$815 \min_{\mathbf{u}_t} \int_0^1 \frac{1}{2} \|\mathbf{u}_t\|_2^2 dt \quad (16)$$

$$816 \text{ s.t. } d\mathbf{x}_t^u = [\theta_t(\mathbf{x}_1 - \mathbf{x}_t^u) + g_t \mathbf{u}_t] dt, \mathbf{x}_0^u = \mathbf{x}_0, \mathbf{x}_1^u = \mathbf{x}_1.$$

817 Taking conditions $\theta \rightarrow 0$ and $g_t = 1$, then DB’s SOC problem (9) becomes

$$818 \min_{\mathbf{u}_t} \int_0^1 \frac{1}{2} \|\mathbf{u}_t\|_2^2 dt \quad (17)$$

$$819 \text{ s.t. } d\mathbf{x}_t^u = \mathbf{u}_t dt, \mathbf{x}_0^u = \mathbf{x}_0, \mathbf{x}_1^u = \mathbf{x}_1,$$

820 which is exactly equivalent to (10). \square

821 A.2 PROOF OF THEOREM 4.2

822 **Theorem 4.2.** *Denote $\mathcal{J}(\mathbf{u}_t) \triangleq \int_0^1 \frac{1}{2} \|\mathbf{u}_t\|_2^2 dt$ as the overall cost of the SOC problems (9) and (10)*
 823 *with the related optimal controller $\mathbf{u}_t^{*,DB}$ and $\mathbf{u}_t^{*,FM}$, respectively. Under the condition of diffusion*
 824 *coefficient $g_t = 1$ in (9) to be consistent with FM (10), then*

$$825 \mathcal{J}(\mathbf{u}_t^{*,DB}) \leq \mathcal{J}(\mathbf{u}_t^{*,FM}). \quad (12)$$

826 *Proof.* Recall the SOC problem (9) with $g_t = 1$ as

$$827 \min_{\mathbf{u}_t} \int_0^1 \frac{1}{2} \|\mathbf{u}_t\|_2^2 dt \quad (18)$$

$$828 \text{ s.t. } d\mathbf{x}_t^u = [\theta_t(\mathbf{x}_1 - \mathbf{x}_t^u) + \mathbf{u}_t] dt, \mathbf{x}_0^u = \mathbf{x}_0, \mathbf{x}_1^u = \mathbf{x}_1.$$

829 Denote $\bar{\theta}_{s:t} = \int_s^t \theta_z dz$, $\bar{\theta}_t = \int_0^t \theta_z dz$ for simplification when $s = 0$ and $\bar{\sigma}_{s:t}^2 = \lambda^2(1 - e^{-2\bar{\theta}_{s:t}})$.

830 Then the related optimal controller $\mathbf{u}_t^{*,DB}$ and $\mathbf{u}_t^{*,FM}$ are

$$831 \mathbf{u}_t^{*,DB} = \frac{e^{-2\bar{\theta}_{t:1}}(\mathbf{x}_1 - \mathbf{x}_t^u)}{\bar{\sigma}_{t:1}^2}, \quad (19)$$

$$832 \mathbf{u}_t^{*,FM} = \frac{\mathbf{x}_1 - \mathbf{x}_t^u}{1 - t},$$

833 where $\mathbf{u}_t^{*,DB}$ is derived from UniDB (Zhu et al., 2025) with conditions $\gamma \rightarrow \infty$ and $g_t = 1$, and
 834 $\mathbf{u}_t^{*,FM}$ is directly obtained from (11).

835 Therefore,

$$836 \frac{1}{2} \int_0^1 \|\mathbf{u}_t^{*,DB}\|_2^2 dt = \frac{1}{2} \int_0^1 \left\| \frac{e^{-2\bar{\theta}_{t:1}}(\mathbf{x}_1 - \mathbf{x}_t^u)}{\bar{\sigma}_{t:1}^2} \right\|_2^2 dt$$

$$837 = \frac{1}{2} \int_0^1 \frac{e^{-4\bar{\theta}_{t:1}}}{\lambda^4(1 - e^{-2\bar{\theta}_{t:1}})^2} \|\mathbf{x}_1 - \mathbf{x}_t^u\|_2^2 dt \quad (20)$$

$$838 = \frac{1}{2} \int_0^1 \frac{\|\mathbf{x}_1 - \mathbf{x}_t^u\|_2^2}{\lambda^4(e^{2\bar{\theta}_{t:1}} - 1)^2} dt,$$

839 and

$$840 \frac{1}{2} \int_0^1 \|\mathbf{u}_t^{*,FM}\|_2^2 dt = \frac{1}{2} \int_0^1 \left\| \frac{\mathbf{x}_1 - \mathbf{x}_t^u}{1 - t} \right\|_2^2 dt \quad (21)$$

$$841 = \frac{1}{2} \int_0^1 \frac{\|\mathbf{x}_1 - \mathbf{x}_t^u\|_2^2}{(1 - t)^2} dt.$$

Recall that $2\lambda^2\theta_t = g_t^2$ in UniDB and we have $g_t = 1$ which implies $\theta_t = \frac{1}{2\lambda^2}$ and we have $e^x - 1 \geq x$ for $\forall x \geq 0$, then

$$\begin{aligned}
\mathcal{J}(\mathbf{u}_t^{*,\text{DB}}) &= \frac{1}{2} \int_0^1 \|\mathbf{u}_t^{*,\text{DB}}\|_2^2 dt \\
&= \frac{1}{2} \int_0^1 \frac{\|\mathbf{x}_1 - \mathbf{x}_t^u\|_2^2}{\lambda^4 (e^{2\bar{\theta}_{t:1}} - 1)^2} dt \\
&= \frac{1}{2} \int_0^1 \frac{\|\mathbf{x}_1 - \mathbf{x}_t^u\|_2^2}{\lambda^4 (e^{\int_t^1 \frac{1}{\lambda^2} dz} - 1)^2} dt \\
&= \frac{1}{2} \int_0^1 \frac{\|\mathbf{x}_1 - \mathbf{x}_t^u\|_2^2}{\lambda^4 (e^{\frac{1-t}{\lambda^2}} - 1)^2} dt \\
&\leq \frac{1}{2} \int_0^1 \frac{\|\mathbf{x}_1 - \mathbf{x}_t^u\|_2^2}{\lambda^4 \left(\frac{1-t}{\lambda^2}\right)^2} dt \\
&= \frac{1}{2} \int_0^1 \frac{\|\mathbf{x}_1 - \mathbf{x}_t^u\|_2^2}{(1-t)^2} dt \\
&= \frac{1}{2} \int_0^1 \|\mathbf{u}_t^{*,\text{FM}}\|_2^2 dt = \mathcal{J}(\mathbf{u}_t^{*,\text{FM}}),
\end{aligned} \tag{22}$$

which concludes the proof. \square

A.3 RECOVERY OF ABSOLUTE CONTINUITY IN THE EMPIRICAL WASSERSTEIN GEODESIC LIMIT

For the empirical measures

$$\hat{\mu}_0^{(n)} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} \rightarrow \mu_0 \quad (n \rightarrow \infty), \quad \hat{\mu}_1^{(n)} = \frac{1}{n} \sum_{j=1}^n \delta_{y_j} \rightarrow \mu_1 \quad (n \rightarrow \infty), \tag{23}$$

As n is limited, for every $t \in (0, 1)$, $\dim_{\mathcal{H}}(\text{spt } \hat{\mu}_t^{(n)}) = 0 < d$

As $n \rightarrow \infty$ with $\hat{\mu}_0^{(n)} \rightarrow \mu_0, \hat{\mu}_1^{(n)} \rightarrow \mu_1$ weakly, one has

$$\hat{\mu}_t^{(n)} \rightarrow \mu_t \quad \text{and} \quad \limsup_{n, m \rightarrow \infty} \dim_{\mathcal{H}}(\mathcal{S}_t^{(n)}) = d \tag{24}$$

where $\mathcal{S}_t^{(n)} = \text{spt } \hat{\mu}_t^{(n)} = \{(1-t)x_i + tT_n(x_i)\}_{i=1}^n$, $\text{spt } \mu$ is the minimal closed support of the measure μ , $\dim_{\mathcal{H}} E$ is the Hausdorff dimension of the set E .

Proof. Let

$$\mu_0, \mu_1 \in \mathcal{P}_2^{\text{ac}}(\mathbb{R}^d) \quad \text{with} \quad \mu_0, \mu_1 \ll \mathcal{L}^d, \tag{25}$$

and let

$$\mu_t := [(1-t)\text{id} + t\nabla\varphi]_{\#}\mu_0, \quad 0 \leq t \leq 1, \tag{26}$$

Suppose only finite samples are available:

$$\hat{\mu}_0^{(n)} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} \rightarrow \mu_0 \quad (n \rightarrow \infty), \quad \hat{\mu}_1^{(n)} = \frac{1}{n} \sum_{j=1}^n \delta_{y_j} \rightarrow \mu_1 \quad (n \rightarrow \infty), \tag{27}$$

Assume, for contradiction, that $\hat{\mu}_0^{(n)} \ll \mathcal{L}^d$. Then there would exist $f \in L^1(\mathcal{L}^d)$ such that

$$\hat{\mu}_0^{(n)}(A) = \int_A f d\mathcal{L}^d \quad \forall A \in \mathcal{B}(\mathbb{R}^d). \tag{28}$$

918 Choosing $A = \{x_i\}$ yields $\mathcal{L}^d(A) = 0$ but $\hat{\mu}_0^{(n)}(A) = 1/n > 0$, contradicting absolute continuity.

919 Meanwhile, the support of any empirical measure is a finite set

$$920 \text{spt} \left(\hat{\mu}_0^{(n)} \right) = \{x_1, \dots, x_n\}, \quad \dim_{\mathcal{H}} \left(\text{spt} \left(\hat{\mu}_0^{(n)} \right) \right) = 0 < d. \quad (29)$$

921 If $\mu_t \ll \mathcal{L}^d$, Caffarelli's interior regularity of optimal transport implies that $\nabla\varphi$ is a locally Hölder-
922 homeomorphism; hence

$$923 \text{spt} \mu_t = \overline{\{(1-t)x + t\nabla\varphi(x) : x \in \text{spt} \mu_0\}}, \quad (30)$$

924 where the symbol \bar{A} denotes the topological closure of the set A , that is, the union of A with the set
925 of all its limit points in the ambient space. Eq. (30) has full dimension:

$$926 \dim_{\mathcal{H}} (\text{spt} \mu_t) = d. \quad (31)$$

927 As the empirical supports $\mathcal{S}_t^{(n)}$ converge weakly to $\text{spt} \mu_t$, Caffarelli's regularity ensures the Hausdorff
928 dimension of the limit superior set recovers the full dimension d , thereby restoring absolute continuity
929 in the large-sample limit, the empirical supports

$$930 \mathcal{S}_t^{(n)} := \text{spt} \hat{\mu}_t^{(n)} = \{(1-t)x_i + tT_n(x_i)\}_{i=1}^n, \quad (32)$$

931 become dense in $\text{spt} \mu_t$ as $n \rightarrow \infty$. More precisely, for every $\epsilon > 0$ and every compact $K \subset \text{spt} \mu_t$,

$$932 \mathcal{L}^d \left(K \setminus B_\epsilon \left(\mathcal{S}_t^{(n)} \right) \right) \rightarrow 0, \quad (33)$$

933 where $B_\epsilon \left(\mathcal{S}_t^{(n)} \right)$ is the ϵ -neighbourhood of the set $\mathcal{S}_t^{(n)}$. Consequently, the upper Minkowski
934 dimension (and hence the Hausdorff dimension) of the limit superior set satisfies

$$935 \limsup_{n, m \rightarrow \infty} \dim_{\mathcal{H}} \left(\mathcal{S}_t^{(n)} \right) = d. \quad (34)$$

936 Thus absolute continuity is recovered in the limit, and the empirical interpolation becomes an honest
937 Wasserstein geodesic. \square

938 A.4 OPTIMALITY IN NON-COMPACT SPACES

939 Let $\mu \in \mathcal{P}(\mathbb{R}^d)$ satisfy $\int_{\mathbb{R}^d} \|x\|^2 d\mu(x) < \infty$ and let $u : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be convex and μ a.e.
940 differentiable. For the quadratic cost

$$941 c(x, y) := \frac{1}{2} \|x - y\|^2 \quad (35)$$

942 define $T = \nabla u$. The McCann interpolation

$$943 x_t = (1-t)x_0 + tT, \quad t \in [0, 1], \quad (36)$$

944 constitutes an optimal transport map in the Wasserstein-2 sense, and this conclusion remains valid
945 even when the support of μ is non-compact (Lei & Gu, 2021).

946 This demonstrates that, as long as the source distribution possesses finite second moments, the linear
947 trajectory of FM (an oversimplification of McCann's approach) coincides with the optimal transport
948 map, without requiring any additional regularity conditions irrespective of compactness.

949 The Monge Problem (MP) with quadratic cost is equivalent to:

$$950 \sup \left\{ M(T) := \int_X \langle x, T(x) \rangle d\mu, \quad T_{\#}\mu = \nu \right\} \quad (37)$$

951 This can be transformed into the dual problem:

$$952 \inf \left\{ \int_X u d\mu + \int_Y u^* d\nu : u(x) + u^*(y) \geq \langle x, y \rangle \right\} \quad (38)$$

A convex function u satisfies

$$u(x) + u^*(y) \geq \langle x, y \rangle \quad \forall x, y \in \mathbb{R}^d, \quad u(x) + u^*(y) = \langle x, y \rangle \quad \text{if} \quad y = \nabla u(x) \quad (39)$$

For any coupling $\gamma \in \Pi(\mu, \nu)$, we have

$$\begin{aligned} \int_{\mathbb{R}^d \times \mathbb{R}^d} \langle x, y \rangle d\gamma(x, y) &\leq \int_{\mathbb{R}^d \times \mathbb{R}^d} (u(x) + u^*(y)) d\gamma(x, y) \\ &= \int_{\mathbb{R}^d} u(x) d\mu(x) + \int_{\mathbb{R}^d} u^*(T(x)) d\mu(x) \\ &= \int_{\mathbb{R}^d} \langle x, T(x) \rangle d\mu(x) \end{aligned} \quad (40)$$

Hence

$$\int_{\mathbb{R}^d \times \mathbb{R}^d} \langle x, y \rangle d\gamma \leq \int_{\mathbb{R}^d \times \mathbb{R}^d} \langle x, y \rangle d\gamma_T. \quad (41)$$

Moreover,

$$\int_{\mathbb{R}^d \times \mathbb{R}^d} \frac{1}{2} (\|x\|^2 + \|y\|^2) d\gamma = \int_{\mathbb{R}^d \times \mathbb{R}^d} \frac{1}{2} (\|x\|^2 + \|y\|^2) d\gamma_T. \quad (42)$$

Subtracting the two identities yields

$$\int_{\mathbb{R}^d \times \mathbb{R}^d} \frac{1}{2} \|x - y\|^2 d\gamma \geq \int_{\mathbb{R}^d \times \mathbb{R}^d} \frac{1}{2} \|x - y\|^2 d\gamma_T, \quad (43)$$

which establishes the optimality of T .

B MORE RELATED WORK

Stochastic Optimal Control. The incorporation of SOC principles into diffusion bridges and flow matching models has emerged as a promising paradigm for effectively guiding distribution transitions. I2SB (Liu et al., 2023) firstly introduced the SOC formulation for modeling diffusion bridges. RB-Modulation (Rout et al., 2024) and NDTM (Pandey et al., 2025a) introduced SOC via a simplified SDE structure for training-free style transfer and solving inverse problems based on pre-trained diffusion models. DBFS (Park et al., 2024) leveraged SOC to construct diffusion bridges in infinite-dimensional function spaces and also established equivalence between SOC and Doob’s h -transform. Adjoint Matching (Domingo-Enrich et al., 2024) for the first time, systematically unifies Flow Matching with SOC, and on this basis proposes a novel, theoretically unbiased reward fine-tuning framework. Mei et al. (2024) proposes a Flow-Matching-based framework for SOC, which reformulates the classical SOC problem into a data-driven optimization task, thereby circumventing the intractability of solving high-dimensional nonlinear Hamilton-Jacobi-Bellman (HJB) equation.

C USE OF LLMs

In this paper, LLMs are used only for polishing and formatting purposes. The core methodology development of this research does not involve LLMs as any important, original, or non-standard components.

D TRAINING AND INFERENCE ALGORITHMS FOR DB AND FM.

Here we demonstrate the training and inference algorithms we used in experiments for DB and FM. We take the training algorithm of UniDB with $\gamma \rightarrow \infty$ for DB and the training objective (2) as the training algorithm for FM. As for the inference method of FM, we directly take the Euler sampling algorithm. We’ve learned that the Euler sampling method for FM is equivalent to DDPMs with DDIMs sampler Gao et al. (2025). Hence, to ensure a fair comparison, we take the first-order UniDB++ with $\gamma \rightarrow \infty$ (accelerating algorithms for GOUB, similar to DDIMs accelerating for DDPMs) as the inference method of DB.

Algorithm 1 DB Training Algorithm**repeat**Take a pair of images \mathbf{x}_0 and \mathbf{x}_T Encode the images as $\mathbf{z}_0 = \text{Encoder}(\mathbf{x}_0)$ and $\mathbf{z}_T = \text{Encoder}(\mathbf{x}_T)$ $t \sim \text{Uniform}(\{1, \dots, T\})$ and $\epsilon \sim \mathcal{N}(0, I)$

$$a = \frac{e^{-\bar{\theta}_{t-1:t}} \bar{\sigma}_{t:T}^2}{\bar{\sigma}_{t-1:T}^2}$$

$$b = \frac{1}{\bar{\sigma}_T^2} \left((1 - e^{-\bar{\theta}_t}) \bar{\sigma}_{t:T}^2 + e^{-2\bar{\theta}_{t:T}} \bar{\sigma}_t^2 - ((1 - e^{-\bar{\theta}_{t-1}}) \bar{\sigma}_{t-1:T}^2 + e^{-\bar{\theta}_{t-1:T}} \bar{\sigma}_{t-1}^2) a \right)$$

$$\mathbf{z}_t = e^{-\bar{\theta}_t \frac{\bar{\sigma}_{t:T}^2}{\bar{\sigma}_T^2}} \mathbf{z}_0 + \left(1 - e^{-\bar{\theta}_t \frac{\bar{\sigma}_{t:T}^2}{\bar{\sigma}_T^2}} \right) \mathbf{z}_T + \bar{\sigma}'_t \epsilon$$

$$\bar{\boldsymbol{\mu}}_t = e^{-\bar{\theta}_t \frac{\bar{\sigma}_{t:T}^2}{\bar{\sigma}_T^2}} \mathbf{z}_0 + \left(1 - e^{-\bar{\theta}_t \frac{\bar{\sigma}_{t:T}^2}{\bar{\sigma}_T^2}} \right) \mathbf{z}_T$$

$$\boldsymbol{\mu}_{t-1,\theta} = \mathbf{z}_t - \left(\theta_t + g_t^2 \frac{e^{-2\bar{\theta}_{t:T}}}{\bar{\sigma}_{t:T}^2} \right) (\mathbf{z}_T - \mathbf{z}_t) + \frac{g_t^2}{\bar{\sigma}_t^2} \boldsymbol{\epsilon}_\theta(\mathbf{z}_t, \mathbf{z}_T, t)$$

$$\boldsymbol{\mu}_{t-1} = \frac{1}{\bar{\sigma}'_t} [\bar{\sigma}'_t{}^2 (\mathbf{x}_t - b\mathbf{x}_T) a + (\bar{\sigma}'_t{}^2 - \bar{\sigma}'_{t-1}{}^2 a^2) \bar{\boldsymbol{\mu}}_t]$$

Take gradient descent step on $\mathcal{L}_\theta = \mathbb{E} \left[\frac{1}{2g_t^2} \|\boldsymbol{\mu}_{t-1,\theta} - \boldsymbol{\mu}_{t-1}\|_1 \right]$ **until** converged**Algorithm 2** FM Training Algorithm**repeat**Take a pair of images \mathbf{x}_0 and \mathbf{x}_1 Encode the images as $\mathbf{z}_0 = \text{Encoder}(\mathbf{x}_0)$ and $\mathbf{z}_1 = \text{Encoder}(\mathbf{x}_1)$ $t \sim \text{Uniform}(\{\frac{1}{T}, \dots, \frac{T}{T} = 1\})$

$$\mathbf{z}_t = (1 - t)\mathbf{z}_0 + t\mathbf{z}_1$$

Take gradient descent step on $\mathcal{L}_\theta = \mathbb{E} [\|\mathbf{v}_\theta(\mathbf{z}_t, t) - (\mathbf{z}_1 - \mathbf{z}_0)\|_2]$ **until** converged**Algorithm 3** DB Inference Algorithm**Input:** LQ images \mathbf{x}_T , noise predicted model $\epsilon_\theta(\mathbf{x}_t, \mathbf{x}_T, t)$, and $M + 1$ time steps $\{t_i\}_{i=0}^M$ decreasing from $t_0 = T$ to $t_M = 0$. Initialize $\mathbf{z}_{t_0} = \text{Encoder}(\mathbf{x}_T)$.**for** $i = 1$ **to** M **do**Sample $\epsilon \sim \mathcal{N}(0, I)$ if $i < M$, else $\epsilon = 0$.

$$\kappa_{t_{i-1}} = e^{\bar{\theta}_{t_{i-1}:T}} (1 - e^{-2\bar{\theta}_{t_{i-1}:T}}), \kappa_{t_i} = e^{\bar{\theta}_{t_i:T}} (1 - e^{-2\bar{\theta}_{t_i:T}})$$

$$\rho_{t_{i-1}} = e^{\bar{\theta}_{t_{i-1}}} (1 - e^{-2\bar{\theta}_{t_{i-1}}}), \rho_{t_i} = e^{\bar{\theta}_{t_i}} (1 - e^{-2\bar{\theta}_{t_i}})$$

$$\delta_{t_{i-1}:t_i}^n = \lambda \kappa_{t_i, \gamma} \sqrt{\frac{1}{1 - e^{-2\bar{\theta}_{t_{i-1}:T}}} - \frac{1}{1 - e^{-2\bar{\theta}_{t_i:T}}}}$$

$$\mathbf{z}_{t_i} = \frac{\kappa_{t_i}}{\kappa_{t_{i-1}}} \mathbf{z}_{t_{i-1}} + \left(1 - \frac{\kappa_{t_i}}{\kappa_{t_{i-1}}} \right) \mathbf{z}_T - 2 \frac{\lambda \kappa_{t_i}}{\sqrt{\rho T}} \left(\sqrt{\frac{\rho_{t_{i-1}}}{\kappa_{t_{i-1}}}} - \sqrt{\frac{\rho_{t_i}}{\kappa_{t_i}}} \right) \boldsymbol{\epsilon}_\theta(\mathbf{z}_{t_{i-1}}, \mathbf{z}_T, t_{i-1}) + \delta_{t_{i-1}:t_i}^n \boldsymbol{\epsilon}$$

end for $\tilde{\mathbf{x}}_0 = \text{Decoder}(\mathbf{z}_0)$ **Return** HQ images $\tilde{\mathbf{x}}_0$ **Algorithm 4** FM Inference Algorithm**Input:** LQ images \mathbf{x}_1 , velocity predicted model $\mathbf{v}_\theta(\mathbf{x}_t, t)$, and $M + 1$ time steps $\{t_i\}_{i=0}^M$ decreasing from $t_0 = 1$ to $t_M = 0$. $\mathbf{z}_{t_0} = \text{Encoder}(\mathbf{x}_1)$.**for** $i = 1$ **to** M **do**

$$\mathbf{z}_{t_i} = \mathbf{z}_{t_{i-1}} + (t_{i-1} - t_i) \mathbf{v}_\theta(\mathbf{z}_{t_{i-1}}, t_{i-1})$$

end for $\tilde{\mathbf{x}}_0 = \text{Decoder}(\mathbf{z}_0)$ **Return** HQ images $\tilde{\mathbf{x}}_0$

E EXPERIMENTAL AND IMPLEMENTATION DETAILS

Details about Datasets, Training time and Inference time (NFE=20, per image). We list details about the used datasets and computational times in all image restoration tasks in Table 3.

Task	Dataset	Training Data Size	Method	Training Time	Inference Time
Inpainting	CelebA-HQ	500	Flow Matching	18 minutes	0.287 s
		1000		36 minutes	0.287 s
		5000		3 hours	0.287 s
		27000		16 hours	0.287 s
		500	Diffusion Bridge	4.5 hours	0.668 s
		1000		9 hours	0.668 s
		5000		12 hours	0.668 s
		27000		66 hours	0.668 s
Super-Resolution	CelebA-HQ	27000	Flow Matching	16 hours	0.287 s
		27000	Diffusion Bridge	66 hours	0.668 s
Deblurring	CelebA-HQ	27000	Flow Matching	16 hours	0.287 s
		27000	Diffusion Bridge	66 hours	0.668 s
Denoising	CelebA-HQ	27000	Flow Matching	16 hours	0.287 s
		27000	Diffusion Bridge	66 hours	0.668 s
Style Transfer	CelebA-HQ	27000	Flow Matching	16 hours	0.287 s
		27000	Diffusion Bridge	66 hours	0.668 s
Image Translation	CelebAMask-HQ	26000	Flow Matching	48 hours	1.164 s
		26000	Diffusion Bridge	212 hours	2.894 s

Table 3: Details about the used datasets and computational times in all image restoration tasks. Both training time and inference time are evaluated on a single NVIDIA H20 GPU.

Unified Transformer Architecture for Diffusion Bridge and Flow Matching. Here we list some shared Transformer Hyper-parameters.

Table 4: Shared Transformer Hyper-parameters.

Hyper-parameter	Value
Patch size	2
Hidden size	1024
Depth	24
Attention heads	16
MLP ratio	4.0

All Hyper-parameters. We set steady variance level $\lambda^2 = 30^2/255^2$, coefficient $e^{-\bar{\theta}_T} = 0.005$ instead of zero, 8 batch size when training, ADAM optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.99$ (Kingma, 2014), 600 thousand total training steps with 10^{-4} initial learning rate and decaying by half at 300 and 500 thousand iterations. We choose a flipped version of cosine noise schedule for θ_t (Luo et al., 2023b),

$$\theta_t = 1 - \frac{\cos^2\left(\frac{t/T+s}{1+s} \frac{\pi}{2}\right)}{\cos^2\left(\frac{s}{1+s} \frac{\pi}{2}\right)}, \quad (44)$$

where $s = 0.008$ is followed from Yue et al. (2024); Zhu et al. (2025); Pan et al. (2025) to achieve a smooth noise schedule. As for time schedule, we directly take the naive uniform time schedule.

Implementation Details. All experiments are trained and tested on a single NVIDIA H20 GPU with 141GB memory.

F ADDITIONAL EXPERIMENTAL RESULTS

We adopted other inpainting masks like thin, thick, and every-second-line (ev2li) masks from RePaint (Lugmayr et al., 2022).

Table 5: Quantitative results for Flow Matching and Diffusion Bridge (denoted FM and DB in table, respectively) under different Image Inpainting tasks with different training data sizes.

Training data size 500												
Method	Box50				Box64				Box72			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
FM	27.33	0.833	0.042	5.78	26.16	0.819	0.050	7.56	25.79	0.812	0.055	9.96
DB	28.04	0.813	0.038	5.34	26.43	0.790	0.047	6.34	26.54	0.779	0.056	7.29
Method	Box80				Box96				Box128			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
FM	25.26	0.803	0.062	12.77	23.93	0.779	0.078	20.05	21.43	0.719	0.133	45.81
DB	25.45	0.777	0.058	7.83	23.95	0.753	0.070	9.88	21.70	0.705	0.103	11.34
Method	Thin mask				Thick mask				ev2li			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
FM	18.73	0.588	0.239	69.47	17.95	0.678	0.200	66.40	24.27	0.681	0.178	39.56
DB	23.41	0.690	0.120	14.83	21.10	0.700	0.131	12.32	25.42	0.702	0.128	16.03
Training data size 1000												
Method	Box50				Box64				Box72			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
FM	27.62	0.836	0.047	5.48	26.68	0.824	0.047	6.32	26.18	0.816	0.052	7.90
DB	27.80	0.810	0.039	5.11	27.18	0.802	0.044	5.55	26.68	0.797	0.048	5.94
Method	Box80				Box96				Box128			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
FM	25.62	0.807	0.058	9.38	24.34	0.785	0.073	13.79	21.76	0.724	0.124	37.23
DB	26.03	0.786	0.053	6.73	24.97	0.772	0.061	6.86	22.23	0.717	0.096	9.43
Method	Thin mask				Thick mask				ev2li			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
FM	20.00	0.625	0.203	52.81	19.22	0.689	0.179	48.29	24.56	0.689	0.162	34.70
DB	24.13	0.711	0.103	11.45	21.96	0.712	0.120	9.98	26.01	0.729	0.113	13.07
Training data size 5000												
Method	Box50				Box64				Box72			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
FM	28.32	0.843	0.037	5.01	27.54	0.834	0.042	5.49	27.07	0.828	0.046	5.86
DB	28.04	0.881	0.037	5.03	27.62	0.809	0.041	5.39	27.04	0.802	0.045	5.45
Method	Box80				Box96				Box128			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
FM	26.50	0.819	0.051	6.55	25.24	0.799	0.064	8.59	22.65	0.744	0.107	17.87
DB	26.19	0.794	0.048	5.77	25.39	0.778	0.058	6.43	22.96	0.733	0.085	8.59
Method	Thin mask				Thick mask				ev2li			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
FM	23.38	0.713	0.129	23.69	21.35	0.725	0.146	24.52	25.66	0.724	0.126	19.99
DB	24.66	0.724	0.095	9.68	22.16	0.719	0.110	8.64	26.23	0.731	0.102	11.28

Table 6: Quantitative results for Flow Matching and Diffusion Bridge (denoted FM and DB in table, respectively) under different Image Inpainting tasks on the CelebA-HQ dataset.

Method	Box50				Box64				Box72			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
FM	28.66	0.847	0.035	4.93	28.03	0.840	0.039	5.13	27.68	0.835	0.042	5.43
DB	28.65	0.820	0.035	4.93	27.90	0.813	0.038	5.11	27.45	0.807	0.041	5.25
Method	Box80				Box96				Box128			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
FM	27.17	0.828	0.047	5.86	26.02	0.810	0.060	8.18	23.54	0.760	0.106	17.84
DB	27.08	0.802	0.044	5.34	25.85	0.786	0.052	6.25	23.57	0.741	0.078	7.71
Method	Thin mask				Thick mask				ev2li			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
FM	24.68	0.746	0.105	14.90	22.51	0.748	0.129	15.82	25.69	0.731	0.122	19.53
DB	25.28	0.739	0.083	8.35	22.86	0.736	0.096	7.34	26.99	0.749	0.084	9.59

Table 7: Quantitative results for Flow Matching and Diffusion Bridge (denoted FM and DB in table, respectively) under different Image Inpainting tasks on the FFHQ dataset.

Method	Box64		Box96		Box128	
	LPIPS↓	FID↓	LPIPS↓	FID↓	LPIPS↓	FID↓
FM	0.047	11.16	0.069	14.89	0.118	25.37
DB	0.044	11.11	0.059	12.02	0.089	16.72

Table 8: Quantitative results for Flow Matching and Diffusion Bridge (denoted FM and DB in table, respectively) under different NFEs of Image Inpainting tasks (Box 128) on the CelebA-HQ dataset.

Method	10 NFEs		20 NFEs		50 NFEs		100 NFEs	
	LPIPS↓	FID↓	LPIPS↓	FID↓	LPIPS↓	FID↓	LPIPS↓	FID↓
FM	0.111	23.01	0.106	17.84	0.099	17.40	0.101	17.51
DB	0.078	9.32	0.078	7.71	0.079	7.20	0.080	7.14

G AN ANALYSIS OF THE ROLE OF CPU IN INFERENCE TIME:

All experiments reported in this paper were conducted on the Server (GPU: 8xH20-3e) equipped with CPUs (Intel Xeon Platinum 8558: base frequency 800 MHz, max turbo 2.1 GHz). We further benchmarked the two models on the Server (GPU: 8x4090) equipped with more powerful CPUs (Intel Xeon Gold 6430: base 1.7 GHz, max turbo 3.4 GHz) and Desktop (GPU: 3090) equipped with Consumer grade CPU (Intel i7-14700K: base 3.4 GHz, max turbo 5.1 GHz). The wall-clock decompositions for neural-network evaluation and iterative updates are reported below.

Table 9: Wall-clock breakdown comparison across different hardware configurations. CPU frequencies (base/turbo) are shown for each platform: H20 Server (800 MHz/2.1 GHz), 4090 Server (1.7 GHz/3.4 GHz), and 3090 Desktop (3.4 GHz/5.1 GHz).

Method	H20 Server		4090 Server		3090 Desktop	
	800 MHz, 2.1 GHz		1.7 GHz, 3.4 GHz		3.4 GHz, 5.1 GHz	
	Network prediction (ms)	Per-NFE prediction (ms)	Network prediction (ms)	Per-NFE prediction (ms)	Network prediction (ms)	Per-NFE prediction (ms)
FM	157.73	7.88	182.37	9.11	346.08	17.30
DB	279.16	13.95	217.53	10.87	362.29	18.11

It can be observed that CPUs with different base frequencies have a significant impact on network prediction time. Specifically, the CPUs with larger base and max turbo frequencies would significantly reduce the gap of time overhead for FM and DB in network prediction. The reason that FM still achieved about a bit fewer network prediction costs than DB lies in the doubled channel, which forces the patch-embedding layer to process double data, slightly raising the inference time.

H ADDITIONAL VISUAL EXPERIMENTAL RESULTS

Here we will illustrate more experimental results.

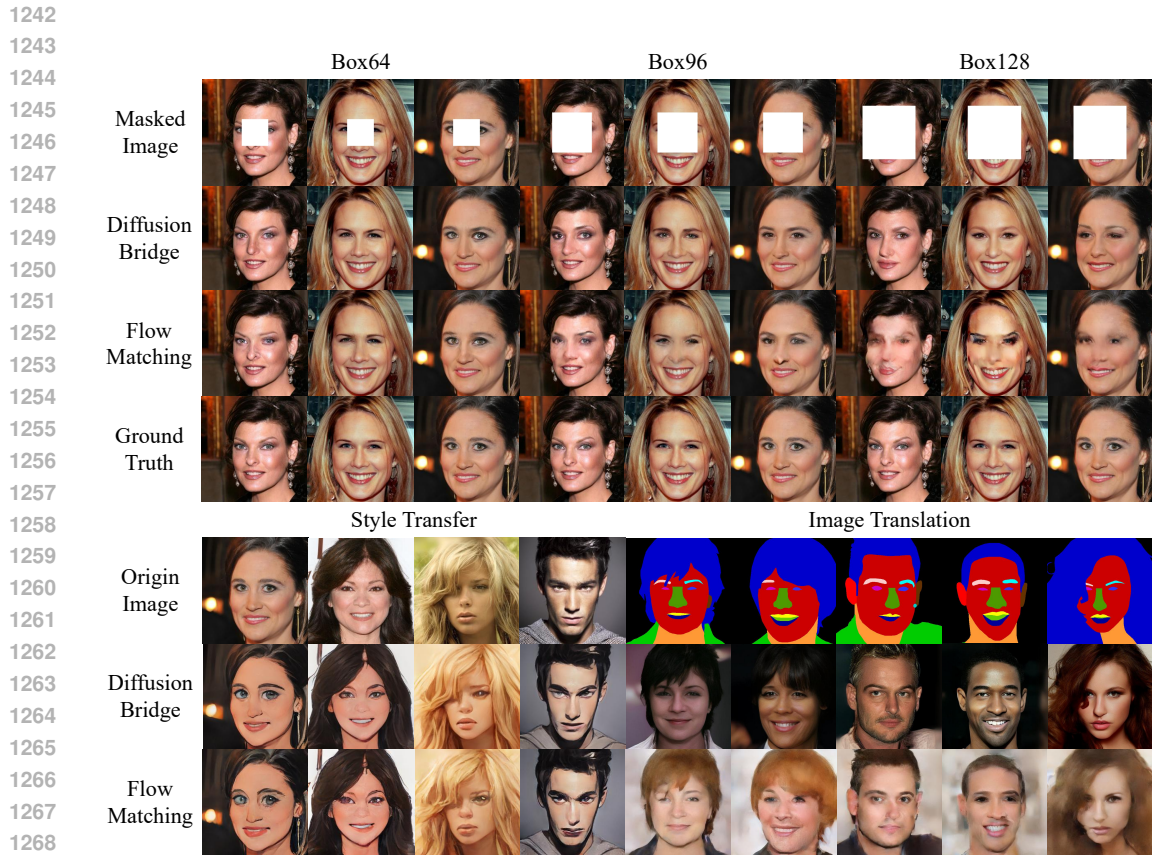


Figure 5: Qualitative comparison of visual results between Diffusion Bridge and Flow Matching in Image Inpainting, Style Transfer and Image Translation tasks.

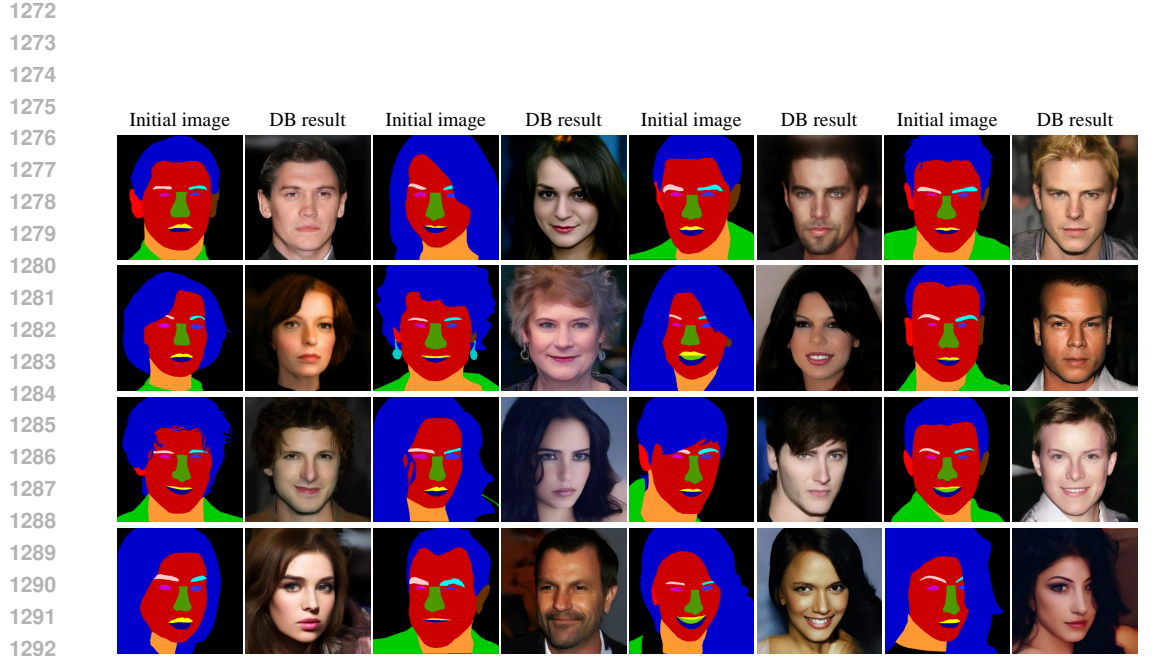


Figure 6: Additional Diffusion Bridge visual results on Image Translation with CelebFaces dataset.

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

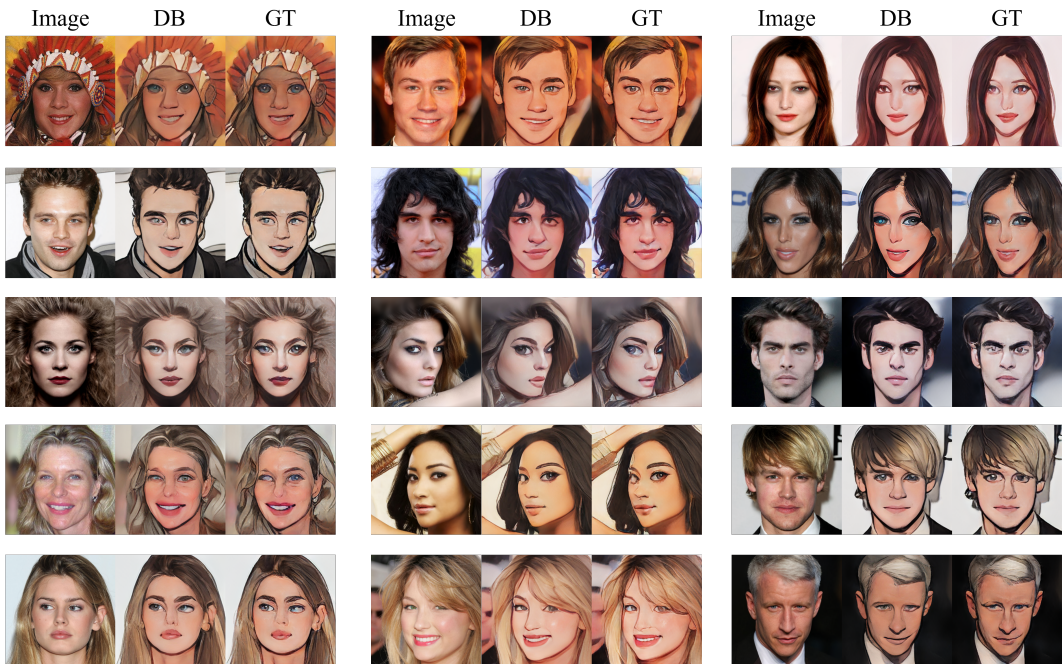


Figure 7: Additional Diffusion Bridge visual results on Style-Transfer with CelebA-HQ dataset.



Figure 8: Additional visual comparison between Diffusion Bridge and Flow Matching on 4×Super-Resolution on the CelebFaces dataset.



Figure 9: Additional Diffusion Bridge and Flow Matching visual results on Image Inpainting (different masks) with CelebA-HQ dataset.



1456 Figure 10: Additional Flow Matching visual results on Image Inpainting (different masks and different
1457 data size (500, 1000, 5000)) with CelebA-HQ dataset.



1510 Figure 11: Additional Flow Matching visual results on Image Inpainting (different masks and different
 1511 data size) with CelebA-HQ dataset.