

GENERALIZED CONSISTENCY TRAJECTORY MODELS FOR IMAGE MANIPULATION

Beomsu Kim*, Jaemin Kim*, Jeongsol Kim & Jong Chul Ye

KAIST

{beomsu.kim, kjm981995, jeongsol, jong.ye}@kaist.ac.kr

ABSTRACT

Diffusion models (DMs) excel in unconditional generation, as well as on applications such as image editing and restoration. The success of DMs lies in the iterative nature of diffusion: diffusion breaks down the complex process of mapping noise to data into a sequence of simple denoising tasks. Moreover, we are able to exert fine-grained control over the generation process by injecting guidance terms into each denoising step. However, the iterative process is also computationally intensive, often taking from tens up to thousands of function evaluations. Although consistency trajectory models (CTMs) enable traversal between any time points along the probability flow ODE (PFODE) and score inference with a single function evaluation, CTMs only allow translation from Gaussian noise to data. This work aims to unlock the full potential of CTMs by proposing generalized CTMs (GCTMs), which translate between arbitrary distributions via ODEs. We discuss the design space of GCTMs and demonstrate their efficacy in various image manipulation tasks such as image-to-image translation, restoration, and editing. Code is available at <https://github.com/1202kbs/GCTM>.

1 INTRODUCTION

Diffusion-based generative models (DMs) learn the scores of noise-perturbed data distributions, which can be used to translate samples between two distributions by numerically integrating an SDE or a probability flow ODE (PFODE) (Ho et al., 2020; Dhariwal & Nichol, 2021; Song et al., 2021). They have achieved remarkable progress over recent years, even surpassing well-known generative models such as Generative Adversarial Networks (GANs) (Goodfellow et al., 2014) or Variational Autoencoders (VAEs) (Kingma & Welling, 2014) in terms of sample quality. Moreover, diffusion models have found wide application in areas such as image-to-image translation (Saharia et al., 2022), image restoration (Chung et al., 2022; 2023), image editing (Meng et al., 2022), etc.

The success of DMs can largely be attributed to the iterative nature of diffusion, arising from its foundation on differential equations – multi-step generation grants high-quality image synthesis by breaking down the complex process of mapping noise to data into a composition of simple denoising steps. We are also able to exert fine-grained control over the generation process by injecting minute guidance terms into each step (Chung et al., 2022; Ho & Salimans, 2022). Indeed, guidance is an underlying principle behind numerous diffusion-based image editing and restoration algorithms.

However, its iterative nature is also a curse, as diffusion inference often demands from tens to thousands of number of neural function evaluations (NFEs) per sample, rendering practical usage difficult. Consequently, there is now a large body of works on improving the inference speed of DMs. Among them, distillation refers to methods which train a neural network to translate samples along PFODE trajectories generated by a pre-trained teacher DM in one or two NFEs. Representative distillation methods include progressive distillation (PD) (Salimans & Ho, 2022), consistency models (CMs) (Song et al., 2023), and consistency trajectory models (CTMs) (Kim et al., 2024b).

In contrast to PD or CMs which only allow traversal to the terminal point of the PFODE, CTMs enable traversal between any pair of time points along the PFODE as well as score inference, all in a single inference step. Thus, in theory, CTMs are more amenable to guidance, and are applicable to a wider variety of downstream image manipulation tasks. Yet, there is a lack of works exploring the effectiveness of CTMs in such context.

*Equal Contribution

Zero-shot image restoration via diffusion. Image restoration such as super-resolution, deblurring, and inpainting can be formulated as inverse problems, which obtain true signals from given observations. With the advancements in DMs serving as powerful priors, diffusion based inverse solvers have been explored actively. DDRM (Kawar et al., 2022) performs denoising steps on the spectral space of a linear corrupting matrix. DPS (Chung et al., 2022) and Π GDM (Song et al., 2022) propose posterior sampling by estimating the likelihood distribution through Jensen’s approximation and Gaussian assumption, respectively. While diffusion-based inverse solvers facilitate zero-shot image restoration, they often need prolonged sampling times. CoSIGN (Zhao et al., 2024) addresses this problem by using CMs as generative priors, but we show that GCTMs can be better priors.

Image translation via diffusion. The seminal work Pix2Pix (Isola et al., 2017) achieves image-to-image translation with conditional GANs. SDEdit (Meng et al., 2022) avoids mode collapse and learning instabilities with GANs by utilizing DMs to translate edited images along SDEs. Palette (Saharia et al., 2022) proposed conditional DMs for image-to-image translation tasks. To address the Gaussian prior constraint with DMs, Schrödinger bridge (SB) (Liu et al., 2023; Kim et al., 2024a), direct diffusion bridge (DDB) (Delbracio & Milanfar, 2023), or denoising diffusion implicit bridge (DDIB) (Su et al., 2023) methods have been proposed to learn SDEs or ODEs between arbitrary two distributions. However, such models often require large NFEs. This has inspired models which distill conditional ODE trajectories (Mei et al., 2024; Xiao et al., 2024), DDB trajectories (He et al., 2024), or DDIB trajectories (Starodubcev et al., 2024). GCTMs are more general in the sense that they enable velocity evaluation and translation between two arbitrary timesteps.

3 BACKGROUND

3.1 DIFFUSION MODELS (DMs)

DMs (Song et al., 2021; Ho et al., 2020) learn to reverse the process of corrupting data into Gaussian noise. Formally, the corruption process can be described by a forward SDE

$$d\mathbf{x}_\tau = \sqrt{2\tau} d\mathbf{w}_\tau \quad (1)$$

defined on the time interval $\tau \in (0, \infty)$. Given \mathbf{x}_0 distributed according to a data distribution $p(\mathbf{x}_0)$, (1) sends \mathbf{x}_0 to Gaussian noise as τ increases from 0 to ∞ . The reverse of the corruption process can be described by the reverse SDE

$$d\mathbf{x}_\tau = -2\tau \nabla \log p(\mathbf{x}_\tau) d\tau + \sqrt{2\tau} d\bar{\mathbf{w}}_\tau \quad (2)$$

or its deterministic counterpart, the probability flow ODE (PFODE)

$$d\mathbf{x}_\tau = -\tau \nabla \log p(\mathbf{x}_\tau) d\tau = \tau^{-1} (\mathbf{x}_\tau - \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_\tau)}[\mathbf{x}_0]) d\tau \quad (3)$$

where $p(\mathbf{x}_\tau)$ is the distribution of \mathbf{x}_τ following (1), and $\bar{\mathbf{w}}_\sigma$ is the standard Wiener process in reverse-time. Given noise $\mathbf{x}_{\hat{\tau}} \sim \mathcal{N}(\mathbf{x}_{\hat{\tau}}|\mathbf{0}, \hat{\tau}^2 \mathbf{I}) \approx p(\mathbf{x}_{\hat{\tau}})$ for some large $\hat{\tau}$, \mathbf{x}_τ following (2) or (3) is distributed $p(\mathbf{x}_\tau)$ as τ decreases from $\hat{\tau}$ to 0. Thus, DMs are able to generate data from noise by approximating the scores $\nabla \log p(\mathbf{x}_\tau)$ via score matching and numerically integrating (2) or (3).

3.2 CONSISTENCY TRAJECTORY MODELS (CTMs)

CTMs (Kim et al., 2024b) learn to translate samples between arbitrary time points of PFODE trajectories, *i.e.*, the goal of CTMs is to learn the integral of the PFODE

$$G(\mathbf{x}_\tau, \tau, \sigma) := \mathbf{x}_\tau + \int_\tau^\sigma u^{-1}(\mathbf{x}_u - \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_u)}[\mathbf{x}_0]) du \quad (4)$$

for $\sigma \leq \tau$, where the terminal distribution $p(\mathbf{x}_{\hat{\sigma}})$ is assumed to be Gaussian. The parametrization

$$\begin{cases} G(\mathbf{x}_\tau, \tau, \sigma) = \frac{\sigma}{\tau} \mathbf{x}_\tau + (1 - \frac{\sigma}{\tau}) g(\mathbf{x}_\tau, \tau, \sigma) \\ g(\mathbf{x}_\tau, \tau, \sigma) := \mathbf{x}_\tau + \frac{\tau}{\tau - \sigma} \int_\tau^\sigma u^{-1}(\mathbf{x}_u - \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_u)}[\mathbf{x}_0]) du \end{cases} \quad (5)$$

enables both traversal along the PFODE as well as score inference, since

$$\lim_{\sigma \rightarrow \tau} g(\mathbf{x}_\tau, \tau, \sigma) = \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_\tau)}[\mathbf{x}_0] \quad (6)$$

so we may define $g(\mathbf{x}_\tau, \tau, \tau) := \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_\tau)}[\mathbf{x}_0]$.

Given a pre-trained DM, CTMs approximate g with a neural net g_θ by simultaneously minimizing a distillation loss and a denoising score matching (DSM) loss. The distillation loss is

$$\mathcal{L}_{\text{CTM}}(\theta) := \mathbb{E}_{0 \leq \sigma \leq u < \tau \leq \hat{\sigma}} \mathbb{E}_{p(\mathbf{x}_\tau)} \left[d(G_\theta(\mathbf{x}_\tau, \tau, \sigma), G_{\text{sg}(\theta)}(\mathbf{x}_{\tau \rightarrow u}, u, \sigma)) \right] \quad (7)$$

where G_θ is the G -function with g_θ in place of g , $d(\cdot, \cdot)$ is a measure of similarity between inputs, sg is the stop-gradient operation, and $\mathbf{x}_{\tau \rightarrow u}$ is the integral of the PFODE from time τ to u starting from \mathbf{x}_τ using score estimates from the pre-trained diffusion model. Minimization of (7) causes G_θ to adhere to PFODE trajectories generated by the pre-trained diffusion model. The DSM loss is

$$\mathcal{L}_{\text{DSM}}(\theta) := \mathbb{E}_{0 \leq \tau \leq \hat{\tau}} \mathbb{E}_{p(\mathbf{x}_0) \mathcal{N}(\epsilon | \mathbf{0}, \mathbf{I})} \mathbb{E}_{p(\mathbf{x}_\tau | \mathbf{x}_0, \epsilon)} \left[\|\mathbf{x}_0 - g_\theta(\mathbf{x}_\tau, \tau, \tau)\|_2^2 \right] \quad (8)$$

where $p(\mathbf{x}_\tau | \mathbf{x}_0, \epsilon) = \delta_{\mathbf{x}_0 + \tau \epsilon}(\mathbf{x}_\tau)$ and $\delta_{\mathbf{y}}(\cdot)$ is a Dirac delta at \mathbf{y} . Minimization of (8) causes g_θ to satisfy (6) (Vincent, 2011). This loss acts as a regularizer which improves score accuracy, and is crucial for sampling with large NFEs (Kim et al., 2024b). Thus, the final CTM training objective is

$$\mathcal{L}_{\text{CTM}}(\theta) + \lambda_{\text{DSM}} \mathcal{L}_{\text{DSM}}(\theta), \quad (9)$$

and it is possible to further improve sample quality by adding a GAN loss.

3.3 FLOW MATCHING (FM)

FM (Lipman et al., 2023; Tong et al., 2024; Pooladian et al., 2023) is another technique for learning PFODEs between two distributions $q(\mathbf{x}_0)$ and $q(\mathbf{x}_1)$. Specifically, let $q(\mathbf{x}_0, \mathbf{x}_1)$ be a joint distribution of $q(\mathbf{x}_0)$ and $q(\mathbf{x}_1)$. Define

$$q(\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_1) := \delta_{(1-t)\mathbf{x}_0 + t\mathbf{x}_1}(\mathbf{x}_t), \quad q(\mathbf{x}_t) := \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_1)} [q(\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_1)] \quad (10)$$

where $t \in (0, 1)$. Then, by Theorem 3.1 in Tong et al. (2024), the ODE given by

$$d\mathbf{x}_t = \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_1 | \mathbf{x}_t)} [\mathbf{x}_1 - \mathbf{x}_0] dt \quad (11)$$

generates the probability path $q(\mathbf{x}_t)$, *i.e.*, with terminal condition $\mathbf{x}_1 \sim q(\mathbf{x}_1)$, \mathbf{x}_t following (11) is distributed according to $q(\mathbf{x}_t)$. Analogous to DSM, the velocity term in (11) can be approximated by a neural network \mathbf{v}_ϕ which solves a regression problem (see Theorem 3.2 in Tong et al. (2024))

$$\min_{\phi} \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_t)} \left[\|\mathbf{x}_1 - \mathbf{x}_0 - \mathbf{v}_\phi(\mathbf{x}_t, t)\|_2^2 \right]. \quad (12)$$

However, unlike diffusion whose terminal distribution $p(\mathbf{x}_{\hat{\sigma}})$ is Gaussian, $q(\mathbf{x}_1)$ can be arbitrary. We provide a complete proof of correctness of this section in Appendix C.1.

4 GENERALIZED CONSISTENCY TRAJECTORY MODELS (GCTMs)

We now present GCTMs, which generalize CTMs to enable translation between arbitrary distributions. We begin with a crucial proposition which proves we can parametrize the solution to the FM ODE (11) in a form analogous to CTMs. The proof is deferred to Appendix C.2.

Proposition 1. *The ODE (11) is equivalent to*

$$d\mathbf{x}_t = t^{-1}(\mathbf{x}_t - \mathbb{E}_{q(\mathbf{x}_0 | \mathbf{x}_t)}[\mathbf{x}_0]) dt \quad (13)$$

defined on $t \in (0, 1)$. Hence, we can express the solution to (11) as

$$\begin{cases} G(\mathbf{x}_t, t, s) = \frac{s}{t} \mathbf{x}_t + \left(1 - \frac{s}{t}\right) g(\mathbf{x}_t, t, s), \\ g(\mathbf{x}_t, t, s) := \mathbf{x}_t + \frac{t}{t-s} \int_t^s u^{-1} (\mathbf{x}_u - \mathbb{E}_{q(\mathbf{x}_0 | \mathbf{x}_u)}[\mathbf{x}_0]) du. \end{cases} \quad (14)$$

There are two differences between (5) and (14). First, the time variables t and s now lie in the unit interval $(0, 1)$ instead of $(0, \infty)$, and second, $p(\mathbf{x}_0 | \mathbf{x}_u)$ is replaced with $q(\mathbf{x}_0 | \mathbf{x}_u)$. The second difference is what enables translation between arbitrary distributions, as $q(\mathbf{x}_0 | \mathbf{x}_u)$ recovers clean images \mathbf{x}_0 given images \mathbf{x}_u perturbed by arbitrary type of vectors (e.g., Gaussian noise, images, etc.), while $p(\mathbf{x}_0 | \mathbf{x}_u)$ recovers clean images \mathbf{x}_0 only for Gaussian-perturbed samples \mathbf{x}_u . We call a neural network g_θ which approximates g in (14) a GCTM, and we can train such a network by optimizing the FM counterparts of \mathcal{L}_{CTM} and \mathcal{L}_{DSM} :

$$\mathcal{L}_{\text{GCTM}}(\theta) := \mathbb{E}_{0 \leq s \leq u < t \leq 1} \mathbb{E}_{q(\mathbf{x}_t)} \left[d(G_\theta(\mathbf{x}_t, t, s), G_{\text{sg}(\theta)}(\mathbf{x}_{t \rightarrow u}, u, s)) \right], \quad (15)$$

$$\mathcal{L}_{\text{FM}}(\theta) := \mathbb{E}_{0 \leq t \leq 1} \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_1)} \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_1)} \left[\|\mathbf{x}_0 - g_\theta(\mathbf{x}_t, t, t)\|_2^2 \right]. \quad (16)$$

The next proposition shows that the PFODE (3) learned by CTMs is a special case of the ODE (13) learned by GCTMs, so GCTMs indeed generalize CTMs. The proof is deferred to Appendix C.3.

Proposition 2. Consider the choice of $q(\mathbf{x}_0, \mathbf{x}_1) = p(\mathbf{x}_0) \cdot \mathcal{N}(\mathbf{x}_1 | \mathbf{0}, \mathbf{I})$. Let

$$t := \tau / (1 + \tau), \quad \mathbf{x}_t := \mathbf{x}_\tau / (1 + \tau) \quad (17)$$

where $\tau \in (0, \infty)$ and \mathbf{x}_τ follows the PFODE (3). Then

$$\mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_\tau)}[\mathbf{x}_0] = \mathbb{E}_{q(\mathbf{x}_0 | \mathbf{x}_t)}[\mathbf{x}_0] \quad (18)$$

and \mathbf{x}_t follows the ODE

$$d\mathbf{x}_t = t^{-1}(\mathbf{x}_t - \mathbb{E}_{q(\mathbf{x}_0 | \mathbf{x}_t)}[\mathbf{x}_0]) dt \quad (19)$$

on $t \in (0, 1)$. Furthermore, let $G_{\text{CTM}}, g_{\text{CTM}}$ denote CTM solutions and let $G_{\text{GCTM}}, g_{\text{GCTM}}$ denote GCTM solutions. Then with $s = \sigma / (1 + \sigma)$,

$$\begin{cases} G_{\text{CTM}}(\mathbf{x}_\tau, \tau, \sigma) = G_{\text{GCTM}}(\mathbf{x}_t, t, s) \cdot (1 + s), \\ g_{\text{CTM}}(\mathbf{x}_\tau, \tau, \sigma) = g_{\text{GCTM}}(\mathbf{x}_t, t, s). \end{cases} \quad (20)$$

In short, (18) shows the equivalence of scores, and (19) shows the equivalence of ODEs. Thus, given g_θ trained with \mathcal{L}_{FM} and $\mathcal{L}_{\text{GCTM}}$ with the setting of Prop. 2, we are able to evaluate diffusion scores and simulate diffusion PFODE trajectories with a simple change of variables (17), as shown in (20).

Given GCTM’s capability to replicate CTM, we will now outline the key components of GCTM that enable its significant extension for various downstream tasks. This flexibility offers a notable advantage of GCTM over CTM.

4.1 THE DESIGN SPACE OF GCTMS

Coupling $q(\mathbf{x}_0, \mathbf{x}_1)$. In contrast to diffusion which only uses the trivial coupling $q(\mathbf{x}_0, \mathbf{x}_1) = q(\mathbf{x}_0)q(\mathbf{x}_1)$ in $\mathcal{L}_{\text{DSM}}(\theta)$, FM allows us to use arbitrary joint distributions of $q(\mathbf{x}_0)$ and $q(\mathbf{x}_1)$ in $\mathcal{L}_{\text{FM}}(\theta)$. Intuitively, $q(\mathbf{x}_0, \mathbf{x}_1)$ encodes our inductive bias for what kind of pairs $(\mathbf{x}_0, \mathbf{x}_1)$ we wish the model to learn, since FM ODE is distributed $q(\mathbf{x}_t)$ at each time t , and $q(\mathbf{x}_t)$ is the distribution of $(1-t)\mathbf{x}_0 + t\mathbf{x}_1$ for $(\mathbf{x}_0, \mathbf{x}_1) \sim q(\mathbf{x}_0, \mathbf{x}_1)$. Here, we list three valid couplings of GCTM as examples (see Alg. 1 for code). In contrast, CTM only uses a special case of the independent coupling.

- *Independent coupling:*

$$q(\mathbf{x}_0, \mathbf{x}_1) = q(\mathbf{x}_0)q(\mathbf{x}_1) \quad (21)$$

This coupling reflects no prior assumption about the relation between \mathbf{x}_0 and \mathbf{x}_1 . As shown earlier, diffusion models use this type of coupling with standard normal $q(\mathbf{x}_1)$.

- *Minibatch entropic optimal transport (EOT) coupling:* in practice, FM loss (16) is approximated by an average over minibatch of pairs $(\mathbf{x}_0, \mathbf{x}_1)$. We can consider minibatch EOT coupling samples (Pooladian et al., 2023) which are generated by sampling $\{\mathbf{x}_0^i\}_{i=1}^K$ from $q(\mathbf{x}_0)$, sampling $\{\mathbf{x}_1^i\}_{i=1}^K$ from $q(\mathbf{x}_1)$, running the Sinkhorn algorithm (Cuturi, 2013) (see Alg. 3) to create a doubly-stochastic EOT matrix \mathbf{P}^{EOT} between the two batches, and sampling $(\mathbf{x}_0, \mathbf{x}_1)$ pairs from \mathbf{P}^{EOT} . As observed by Pooladian et al. (2023), minibatch EOT coupling can accelerate flow matching optimization by reducing gradient variance, and we expect similar benefits for GCTM training as well.

- *Supervised coupling:*

$$q(\mathbf{x}_0, \mathbf{x}_1) = \int q(\mathbf{x}_0)q(\mathbf{H}|\mathbf{x}_0)\delta_{\mathbf{H}\mathbf{x}_0}(\mathbf{x}_1) d\mathbf{H} \quad (22)$$

where $\mathbf{H} \sim q(\mathbf{H}|\mathbf{x}_0)$ is a random operator, possibly dependent on \mathbf{x}_0 , which maps ground-truth data \mathbf{x}_0 to observations \mathbf{x}_1 , i.e., $\mathbf{x}_1 = \mathbf{H}\mathbf{x}_0$. For instance, in the context of learning an inpainting model, \mathbf{H} is could be a random masking operator.

Gaussian perturbation. The cardinality of the support of $q(\mathbf{x}_1)$ must be larger than or equal to the cardinality of the support of $q(\mathbf{x}_0)$ for there to be a well-defined ODE from $q(\mathbf{x}_1)$ to $q(\mathbf{x}_0)$. This is because the ODE trajectory given an initial condition is unique, so a single sample $\mathbf{x}_1 \sim q(\mathbf{x}_1)$ cannot be transported to multiple points in the support of $q(\mathbf{x}_0)$. A simple way to address this problem is to add small Gaussian noise to $q(\mathbf{x}_1)$ samples such that $q(\mathbf{x}_1)$ is supported everywhere.

Algorithm 1 $q(x_0, x_1)$ Sampling

```

1: Assume  $m = 1, \dots, M$ , Batch size  $M$ 
2: if Coupling is Independent then
3:    $\{x_0^m\}_m \sim q(x_0), \{x_1^m\}_m \sim q(x_1)$ 
4:   Return  $\{(x_0^m, x_1^m)\}_m$ 
5: else if Coupling is Minibatch EOT then
6:    $\{x_0^m\}_m \sim q(x_0), \{x_1^m\}_m \sim q(x_1)$ 
7:   Return SK( $\{x_0^m\}_m, \{x_1^m\}_m, \tau$ )
8: else if Coupling is Supervised then
9:    $\{x_0^m\}_m \sim q(x_0), \mathbf{H}^m \sim q(\mathbf{H}|x_0^m)$ 
10:  Return  $\{(x_0^m, \mathbf{H}^m x_0^m)\}_m$ 
11: end if

```

Algorithm 2 GCTM Training

```

1: while training do
2:   Sample times  $\{\hat{t}^m\}_m, \{(t^m, s^m, u^m)\}_m$ 
3:   With Alg. 1,  $\{(x_0^m, x_1^m)\}_m \sim q(x_0, x_1)$ 
4:    $x_{\hat{t}^m}^m \leftarrow (1 - \hat{t}^m)x_0^m + \hat{t}^m x_1^m$ 
5:    $x_{\hat{t}^m}^m \leftarrow (1 - t^m)x_0^m + t^m x_1^m$ 
6:    $\mathcal{L}_{\text{FM}}(\theta) = \frac{1}{M} \sum_m \|x_0^m - g_\theta(x_{\hat{t}^m}^m, \hat{t}^m, \hat{t}^m)\|_2^2$ 
7:    $\tilde{x}_{s^m}^m \leftarrow G_{\text{sg}(\theta)}(x_{\hat{t}^m \rightarrow u^m}^m, u^m, s^m)$ 
8:    $\mathcal{L}_{\text{GCTM}}(\theta) = \frac{1}{M} \sum_{m=1}^M d(G_\theta(x_{\hat{t}^m}^m, t^m, s^m), \tilde{x}_{s^m}^m)$ 
9:   Minimize  $\mathcal{L}_{\text{GCTM}}(\theta) + \lambda_{\text{FM}} \mathcal{L}_{\text{FM}}(\theta)$ 
10: end while

```

We emphasize that Gaussian perturbation allows GCTMs to achieve one-to-many generation when we use the supervised coupling. Concretely, consider the scenario where there are multiple labels $x_0 \sim q(x_0|x_1)$ which correspond to an observed x_1 . Then, the perturbation ϵ added to x_1 acts as a source of randomness, allowing the GCTM network to map $x_1 + \epsilon$ to distinct labels x_0 for distinct ϵ . This stands in contrast to simply regressing the neural network output of x_1 to corresponding labels $x_0 \sim q(x_0|x_1)$ with ℓ_2 loss, as this will cause the network to map x_1 to the blurry posterior mean $\mathbb{E}_{q(x_0|x_1)}[x_0]$ instead of a sharp image x_0 . Indeed, in Section 5.2, we observe blurry outputs if we use regression instead of GCTMs.

Time discretization. In practice, to optimize $\mathcal{L}_{\text{GCTM}}(\theta)$, we sample time variables s, t, u from a discretization $t_0 < t_1 < \dots < t_N$ of the unit interval $[0, 1]$, and simulate $x_{t \rightarrow u}$ with respect to the discretization as well. Given the success of the EDM time discretization Karras et al. (2022) for fast sampling of diffusion models, we propose using the EDM time discretization converted to FM time discretization via change of variables in Proposition 2,

$$t_0 = 0, \quad t_i = \frac{\tau_i}{\tau_i + 1} \quad \text{where} \quad \tau_i = (\sigma_{\min}^{1/\rho} + \frac{i}{N}(\sigma_{\max}^{1/\rho} - \sigma_{\min}^{1/\rho}))^\rho. \quad (23)$$

We fix $\rho = 7$ and $\sigma_{\min} = 0.002$ as proposed in Karras et al. (2022), and control σ_{\max} . We note that σ_{\max} controls the amount of emphasis on time near $t = 1$, *i.e.*, larger σ_{\max} places more time discretization points near $t = 1$.

5 EXPERIMENTS

We now explore the possibilities of GCTMs on unconditional generation, image-to-image translation, image restoration, image editing, and latent manipulation. In particular, GCTM admits NFE = 1 sampling via $x_t \mapsto G_\theta(x_t, t, 0)$. Due to the similarities between CTMs and GCTMs as detailed in Thm. 1, GCTMs can be trained using CTM training methods. In fact, we run Alg. 2 with the method in Section 5.2 of (Kim et al., 2024b) to train all GCTMs without pre-trained teacher models. A complete description of training settings are deferred to Appendix A.

5.1 FAST UNCONDITIONAL GENERATION

In the scenario where we do not have access to data pairs, we must resort to either the independent coupling or the OT coupling. Here, we show that the optimal transport coupling can significantly accelerate the convergence speed of GCTMs during training, especially when we use a smaller number of timesteps N in time discretization during GCTM training (see Section 4.1). Using small N may be of interest when we wish to trade-off training speed for performance, since per-iteration training cost of GCTMs increases linearly with N . For instance, when $t = 1$ and $u = s = 0$ in the GCTM loss (15), we need to integrate along the entire time interval $(0, 1)$, which requires N steps of ODE integration.

In Figure 2, we observe up to $\times 2.5$ acceleration in terms of training iterations when we use OT coupling instead of independent coupling. Indeed, in Figure 3, OT coupling samples are visually sharper than independent coupling samples. We postulate this is because (1) OT coupling leads to straighter

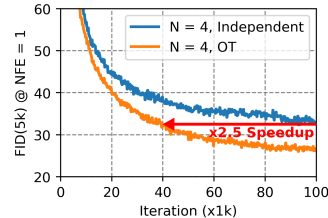


Figure 2: Training acceleration.

Method	Teacher	FID ↓
CTM	✓	5.28
	✗	9.00
CM	✓	3.55
	✗	8.70
iCM	✗	2.51
GCTM (OT)	✗	5.32

Table 1: FID at NFE = 1.



Figure 3: CIFAR10 unconditional samples with NFE = 1.

Method	NFE	Time (ms)	Edges→Shoes			Night→Day			Facades		
			FID ↓	IS ↑	LPIPS ↓	FID ↓	IS ↑	LPIPS ↓	FID ↓	IS ↑	LPIPS ↓
Regression	1	87	54.3	<u>3.41</u>	0.100	189.2	<u>1.85</u>	<u>0.373</u>	<u>121.8</u>	3.28	0.274
Pix2Pix (Isola et al., 2017)	1	33	77.0	3.17	0.208	158.0	1.68	0.418	134.1	2.74	0.288
Palette (Saharia et al., 2022)	5	166	334.1	1.90	0.861	350.2	1.16	0.707	259.3	2.47	0.394
I ² SB(Liu et al., 2023)	5	284	<u>53.9</u>	3.23	0.154	145.8	1.79	0.376	135.2	2.51	<u>0.269</u>
GCTM	1	87	40.3	3.54	0.097	<u>148.8</u>	2.00	0.317	111.3	<u>2.99</u>	0.230

Table 2: I2I translation results (64×64 resolution). Best is in **bold**, second best is underlined.

ODE trajectories, so we can accurately integrate ODEs with smaller N , and (2) lower variance from OT pairs leads to smaller variance in loss gradients, as discussed in (Pooladian et al., 2023).

In Table 1, we compare the Fréchet Inception Distance (FID) (Heusel et al., 2017) of GCTM and relevant baselines on CIFAR10 with NFE = 1. In the setting where we do not use a pre-trained teacher diffusion model, GCTM with OT coupling outperforms all methods with the exception of iCM (Song & Dhariwal, 2024), which is an improved variant of CM. Moreover, GCTM is on par with CTM trained with a teacher. Here, the numbers for CTM are our reproduced results without GAN loss for fair comparison. We speculate that further tuning of hyper-parameters or addition of a GAN loss could push the performance of GCTMs to match that of iCMs.

5.2 FAST IMAGE-TO-IMAGE TRANSLATION

Unlike previous distillation methods such as CM or CTM, GCTM can learn ODEs between arbitrary distributions, enabling image-to-image translation. To numerically validate this theoretical improvement, we train GCTMs on three translation tasks Edges→Shoes, Night→Day, and Facades (Isola et al., 2017), scaled to 64×64 , with the supervised coupling. We consider three baseline methods: ℓ_2 -regression, Pix2Pix (Isola et al., 2017), Palette (Saharia et al., 2022) and I²SB(Liu et al., 2023). To evaluate translation performance, we use FID and Inception Score (IS) (Barratt & Sharma, 2018) to rate translation quality and LPIPS (Zhang et al., 2018) to assess faithfulness to input. We control NFEs such that all methods have similar inference times, and calculate all metrics on validation set.

In Table 2, we see GCTM shows strong performance on all tasks. SDE-based methods I²SB and Palette show poor performance at low NFEs, even when trained with pairs. Qualitative results in Figure 4 are in line with the metrics. Baselines produce blurry or nonsensical samples, while GCTM produces sharp and realistic images that are faithful to the input. In Table 5, we show results on higher resolution (256×256) data, and observe similar trends.

5.3 FAST IMAGE RESTORATION

We consider two settings on the FFHQ 64×64 dataset, where we either know or do not know the corruption operator. In the former case, we train an unconditional GCTM with the independent coupling, with which we implement three zero-shot image restoration algorithms: DPS, CM-based image restoration, and the guided generation algorithm illustrated in Figure 1, where the loss is given as inconsistency between observations (see Append. B.2 for pseudo-codes and a detailed discussion of the differences). In the latter case, we train a GCTM with the supervised coupling

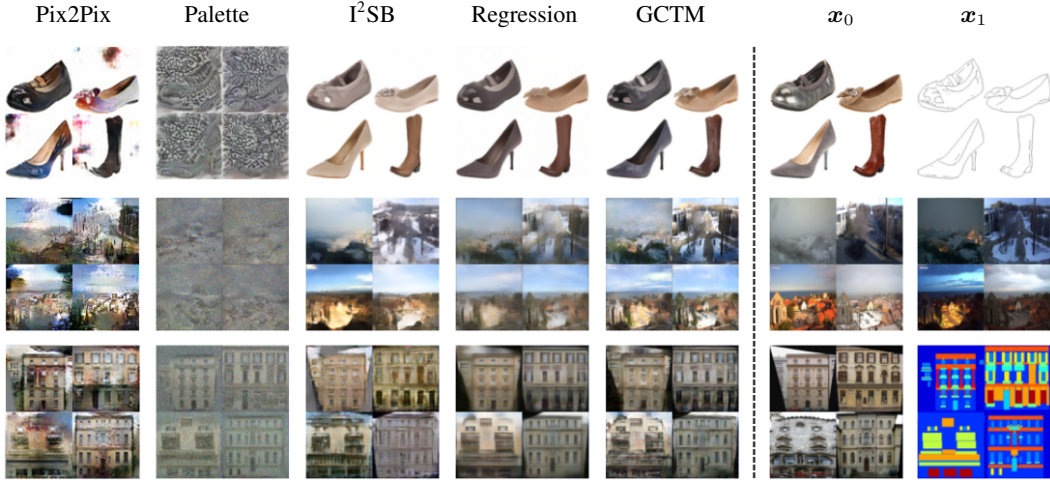


Figure 4: Qualitative evaluation of I2I translation (64×64 resolution) on Edges→Shoes (top), Night→Day (middle) and Facades (bottom). NFE = 5 for I²SB and Palette.

	Method	NFE	Time (ms)	SR2 - Bicubic			Deblur - Gaussian			Inpaint - Center		
				PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
<i>0-Shot</i>	DPS	32	1079	31.19	0.935	0.015	27.88	0.878	0.041	24.69	0.876	0.042
	CM	32	1074	30.80	0.930	0.010	27.85	0.871	0.027	23.02	0.857	0.050
	GCTM	32	1382	31.61	0.939	0.015	28.19	0.885	0.037	24.47	0.876	0.042
<i>Superv.</i>	Regression	1	87	33.46	0.964	0.015	31.19	0.942	0.015	28.76	0.922	0.028
	Palette	5	166	17.88	0.556	0.234	17.81	0.571	0.234	16.12	0.489	0.357
	I ² SB	5	284	26.74	0.869	0.033	26.20	0.853	0.038	26.01	0.874	0.038
	GCTM	1	87	32.37	0.954	0.009	30.56	0.935	0.009	27.37	0.896	0.027

Table 3: Quantitative evaluation of image restoration on FFHQ (64×64 resolution).

and ℓ_2 -regression, I²SB and Palette for comparison. Notably, GCTM is the only model applicable to both situations, thanks to the flexible choice of couplings. We again control NFEs such that all methods have similar inference speed.

Table 3 presents the numerical results in both settings. In the zero-shot setting, we see GCTM outperforming both DPS and CM. In particular, CM is slightly worse than DPS. Sample quality degradation due to error accumulation for CMs at large NFEs have already been observed in unconditional generation (e.g., see Fig. 9 in (Kim et al., 2024b)), and we speculate a similar problem occurs for CMs in image restoration as well. On the other hand, GCTMs avoid this problem, as they are able to traverse to a smaller time using the ODE velocity approximated via g_θ . In the supervised setting, we see regression attains the best PSNR and SSIM. This is a natural consequence of perception-distortion trade-off. Specifically, regression minimizes the MSE loss, so it leads to best distortion metrics (Debracio & Milanfar, 2023) while producing blurry results. GCTM, which provides best results if we exclude regression on distortion metrics (PSNR and SSIM) and best results on perception metrics (LPIPS), strikes the best balance between perception and distortion. For instance, in Fig. 5 inpainting results, regression sample lacks detail (e.g., wrinkles) while GCTM sample is sharp. We show more samples in Appendix E. In particular, in Table 6, we demonstrate image restoration task of GCTM on ImageNet with higher resolution (256×256 resolution) images to demonstrate its scalability.

5.4 FAST IMAGE EDITING

In this section, we demonstrate that GCTM can perform realistic and faithful image editing without any special purpose training. Figure 6 shows image editing with an Edges→Shoes model and an unconditional FFHQ model. On Edges→Shoes, to edit an image, a user creates an edited input,

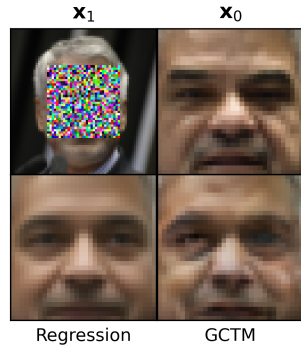


Figure 5: Reg. vs. GCTM.

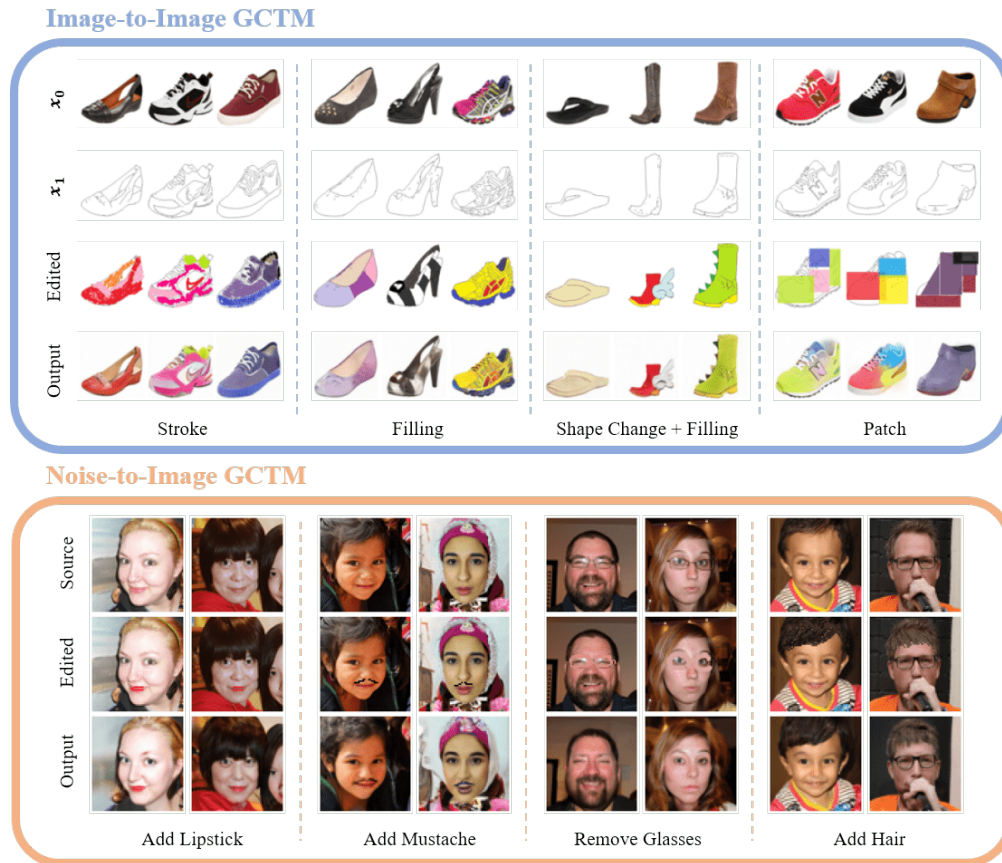


Figure 6: Image editing with GCTM, NFE = 1.

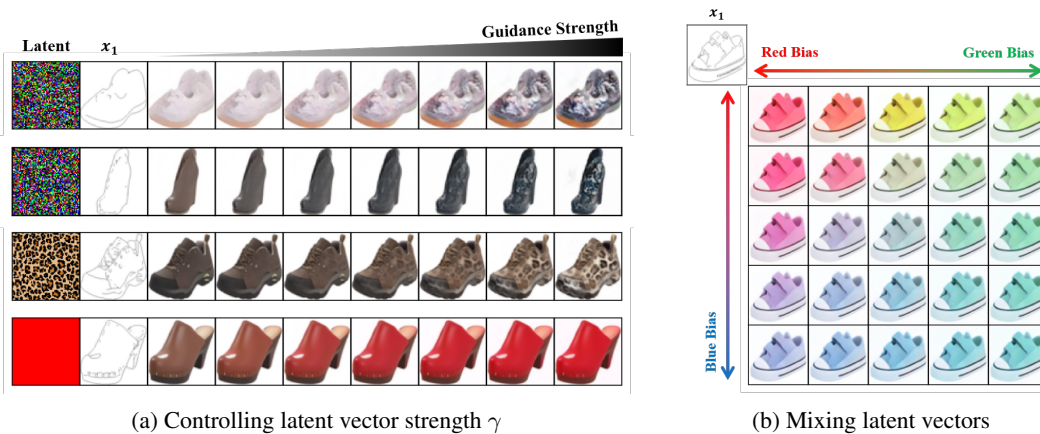


Figure 7: Latent manipulation with image-to-image GCTM, NFE = 1.

which is an edge image painted to have a desired color and / or modified to have a desired outline. We then interpolate the edited input and the original edge image to a certain time point $t = s$ and send it to time $t = 0$ with GCTM to produce the output. On FFHQ, analogous to SDEdit (Meng et al., 2022), we interpolate an edited image with Gaussian noise and send it to time $t = 0$ with GCTM to generate the output. In contrast to previous image editing models such as SDEdit, GCTM requires only a single step to edit an image. Moreover, we observe that GCTM faithfully preserves source image structure while making the desired changes to the image.

5.5 FAST LATENT MANIPULATION

In this section, we demonstrate that GCTMs have a highly controllable latent space. Since there are plenty of works on latent manipulation with unconditional diffusion models, we focus on latent manipulation with GCTMs trained for image-to-image translation. For an image-to-image translation GCTM trained with Gaussian perturbation in Section 4.1, we assert that the perturbation added to \mathbf{x}_1 can be manipulated to produce desired outputs \mathbf{x}_0 . In other words, the perturbation acts as a “latent vector” which controls the factors of variation in \mathbf{x}_0 . To test this hypothesis, in Figure 7, we display outputs $G_\theta(\mathbf{x}_1 + \gamma\epsilon, 1, 0)$ for particular choices of ϵ . In the left panel, we observe generated outputs increasingly adhere to the texture of latent ϵ as we increase guidance strength γ . Interestingly, GCTM generalizes well to latent vectors unseen during training, such as leopard spots or the color red. In the right panel, we explore the effect of linearly combining red, green, and blue latent vectors. We see that the desired color change is reflected faithfully in the outputs. These observations validate our hypothesis that image-to-image GCTMs have an interpretable latent space.

5.6 ABLATION STUDY

We now perform an ablation study on the design choices of Section 4.1. We have already illustrated the power of using appropriate couplings in previous sections, so we explore the importance of σ_{\max} . A robust choice for σ_{\max} for unconditional generation is well-known to be $\sigma_{\max} = 80$ (Karras et al., 2022; Kim et al., 2024b), and we found using this choice to perform sufficiently well for GCTMs when learning to translate noise to data with independent or OT couplings. So, we restrict our attention to image-to-image translation.

In Figure 8, we display the learning curves on Edges \rightarrow Shoes for GCTMs trained without and with Gaussian perturbation, and $\sigma_{\max} \in \{80, 500\}$. We observe that GCTM trained without perturbation and $\sigma_{\max} = 80$ exhibits unstable dynamics, and is unable to minimize the FID below 30. On other hand, GCTM trained with perturbation and $\sigma_{\max} = 80$ surpasses the model trained without perturbation. This demonstrates Gaussian perturbation is indeed crucial for one-to-many generation, as noted in the last paragraph of Section 4.1. Finally, GCTM with both perturbation and $\sigma_{\max} = 500$ minimizes FID the fastest. This shows high-curvature regions for image-to-image ODEs lie near \mathbf{x}_1 , so we need to use a large σ_{\max} which places more discretization points near $t = 1$.

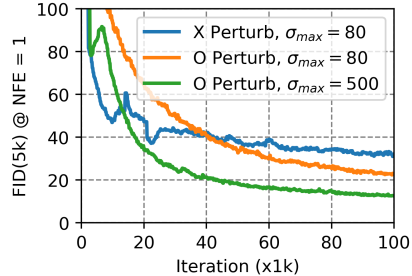


Figure 8: Ablation study of GCTM.

6 CONCLUSION

Our work marks a significant advancement in the realm of ODE-based generative models, particularly on the transformative capabilities of Consistency Trajectory Models (CTMs). While the iterative nature of diffusion has proven to be a powerful foundation for high-quality image synthesis and nuanced control, the computational demands associated with numerous neural function evaluations (NFEs) per sample have posed challenges for practical implementation. Our proposal of Generalized CTMs (GCTMs) extends the reach of CTMs by enabling one-step translation between arbitrary distributions, surpassing the limitations of traditional CTMs confined to Gaussian noise to data transformations. Through an insightful exploration of the design space, we elucidate the impact of various components on downstream task performance, providing a comprehensive understanding that contributes to a broadly applicable and stable training scheme. Empirical validation across diverse image manipulation tasks demonstrates the potency of GCTMs, showcasing their ability to accelerate and enhance diffusion-based algorithms. In summary, our work not only contributes to theoretical advancements but also delivers tangible benefits, showcasing GCTMs as a key element in unlocking the full potential of diffusion models for practical, real-world applications in image synthesis, translation, restoration, and editing.

ACKNOWLEDGMENTS

This work was supported by the National Research Foundation of Korea under Grant RS-2024-00336454 and by the Institute for Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (RS-2019-II190075, Artificial Intelligence Graduate School Program, KAIST).

REFERENCES

- Shane Barratt and Rishi Sharma. A Note on the Inception Score. *arXiv preprint arXiv:1801.01973*, 2018.
- Hyungjin Chung, Jeongsol Kim, Michael T Mccann, Marc L Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. *ICLR*, 2022.
- Hyungjin Chung, Jeongsol Kim, and Jong Chul Ye. Direct diffusion bridge using data consistency for inverse problems. *NeurIPS*, 2023.
- Hyungjin Chung, Suhyeon Lee, and Jong Chul Ye. Decomposed Diffusion Sampler for Accelerating Large-Scale Inverse Problems. *ICLR*, 2024.
- Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *NeurIPS*, 2013.
- Mauricio Delbracio and Peyman Milanfar. Inversion by direct iteration: An alternative to denoising diffusion for image restoration. *TMLR*, 2023.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. *NeurIPS*, 2021.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Networks. *NeurIPS*, 2014.
- Guande He, Kaiwen Zheng, Jianfei Chen, Fan Bao, and Jun Zhu. Consistency Diffusion Bridge Models. *NeurIPS*, 2024.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NeurIPS*, 2017.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020.
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *NeurIPS*, 2022.
- Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. *NeurIPS*, 2022.
- Beomsu Kim, Gihyun Kwon, Kwanyoung Kim, and Jong Chul Ye. Unpaired image-to-image translation via neural Schrödinger bridge. *ICLR*, 2024a.
- Dongjun Kim, Chieh-Hsin Lai, Wei-Hsiang Liao, Naoki Murata, Yuhta Takida, Toshimitsu Uesaka, Yutong He, Yuki Mitsufuji, and Stefano Ermon. Consistency Trajectory Models: Learning Probability Flow ODE Trajectory of Diffusion. *ICLR*, 2024b.
- Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *ICLR*, 2015.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *ICLR*, 2014.
- Bo Li, Kaitao Xue, Bin Liu, and Yu-Kun Lai. BBDM: Image-to-image Translation with Brownian Bridge Diffusion Models. *CVPR*, 2023.
- Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. *ICLR*, 2023.
- Guan-Horng Liu, Arash Vahdat, De-An Huang, Evangelos A Theodorou, Weili Nie, and Anima Anandkumar. I²SB: Image-to-Image Schrödinger Bridge. *ICML*, 2023.

- Kangfu Mei, Mauricio Delbracio, Hossein Talebi, Zhengzhong Tu, Vishal M. Patel, and Peyman Milanfar. CoDi: Conditional Diffusion Distillation for Higher-Fidelity and Faster Image Generation. In *CVPR*, 2024.
- Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. SDEdit: Guided image synthesis and editing with stochastic differential equations. *ICLR*, 2022.
- Aram-Alexandre Pooladian, Heli Ben-Hamu, Carles Domingo-Enrich, Brandon Amos, Yaron Lipman, and Ricky T. Q. Chen. Multisample flow matching: Straightening flows with minibatch couplings. *ICML*, 2023.
- Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *SIGGRAPH*, 2022.
- Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *ICLR*, 2022.
- Jiaming Song, Arash Vahdat, Morteza Mardani, and Jan Kautz. Pseudoinverse-guided diffusion models for inverse problems. *ICLR*, 2022.
- Yang Song and Prafulla Dhariwal. Improved techniques for training consistency models. *ICLR*, 2024.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-Based Generative Modeling through Stochastic Differential Equations. *ICLR*, 2021.
- Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency Models. *ICML*, 2023.
- Nikita Starodubcev, Mikhail Khoroshikh, Artem Babenko, and Dmitry Baranchuk. Invertible Consistency Distillation for Text-Guided Image Editing in Around 7 Steps. *arXiv preprint arXiv:2406.14539*, 2024.
- Xuan Su, Jiaming Song, Chenlin Meng, and Stefano Ermon. Dual Diffusion Implicit Bridges for Image-to-Image Translation. In *ICLR*, 2023.
- Alexander Tong, Kilian Fatras, Nikolay Malkin, Guillaume Hugué, Yanlei Zhang, Jarrid Rector-Brooks, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models with minibatch optimal transport. *TMLR*, 2024.
- Pascal Vincent. A Connection Between Score Matching and Denoising Autoencoders. *Neural Computation*, 23(7):1661–74, 2011.
- Jie Xiao, Kai Zhu, Han Zhang, Zhiheng Liu, Yujun Shen, Zhantao Yang, Ruili Feng, Yu Liu, Xueyang Fu, and Zheng-Jun Zha. CCM: Real-Time Controllable Visual Content Creation Using Text-to-Image Consistency Models. *ICML*, 2024.
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. *CVPR*, 2018.
- Jiankun Zhao, Bowen Song, and Liyue Shen. CoSIGN: Few-Step Guidance of ConSistency Model to Solve General INverse Problems. In *ECCV*, 2024.
- Min Zhao, Fan Bao, Chongxuan Li, and Jun Zhu. EGSDE: Unpaired image-to-image translation via energy-guided stochastic differential equations. *NeurIPS*, 2022.

A FULL EXPERIMENT SETTINGS

A.1 TRAINING

In this section, we introduce training choices which provided reliable performance across all experiments in our paper.

Bootstrapping scores. In all our experiments, we train GCTMs without a pre-trained score model. So, analogous to CTMs, we use velocity estimates given by an exponential moving average θ_{EMA} of θ to solve ODEs. We use exponential moving average decay rate 0.999.

Time discretization. In practice, we discretize the unit interval into a finite number of timesteps $\{t_n\}_{n=0}^N$ where

$$t_0 = 0 < t_1 < \dots < t_N = 1 \quad (24)$$

and learn ODE trajectories integrated with respect to the discretization schedule. EDM (Karras et al., 2022), which has shown robust performance on a variety of generation tasks, solves the PFODE on the time interval $(\sigma_{\min}, \sigma_{\max})$ for $0 < \sigma_{\min} < \sigma_{\max}$ according to the discretization schedule

$$\sigma_n = (\sigma_{\min}^{1/\rho} + (n/N)(\sigma_{\max}^{1/\rho} - \sigma_{\min}^{1/\rho}))^\rho \quad (25)$$

for $n = 0, \dots, N$ and $\rho = 7$. Thus, using the change of time variable (17) derived in Theorem 1, we convert PFODE EDM schedule to FM ODE discretization

$$t_0 = 0, \quad t_n = \sigma_n / (1 + \sigma_n) \quad \text{for } n = 1, \dots, N - 1, \quad t_N = 1. \quad (26)$$

In our experiments, we fix $\sigma_{\min} = 0.002$ and control σ_{\max} . We note that σ_{\max} controls the amount of emphasis on time near $t = 1$, i.e., larger σ_{\max} places more time discretization points near $t = 1$.

Number of discretization steps N . CTMs use fixed $N = 18$. In contrast, analogous to iCMs, we double N every 100k iterations, starting from $N = 4$.

Time \hat{t} distribution. For unconditional generation, we sample

$$t = \sigma / (1 + \sigma), \quad \log \sigma \sim \mathcal{N}(-1.2, 1.2^2) \quad (27)$$

in accordance with EDM. For image-to-image translation, we sample

$$t \sim \text{beta}(3, 1). \quad (28)$$

Network conditioning. We use the EDM conditioning, following CTMs.

Distance d . CTMs use d defined as

$$d(\mathbf{x}_t, \hat{\mathbf{x}}_t) = \text{LPIPS}(G_{\theta_{\text{EMA}}}(\mathbf{x}_t, t, 0), G_{\theta_{\text{EMA}}}(\hat{\mathbf{x}}_t, t, 0)) \quad (29)$$

which compares the perceptual distance of samples projected to time $t = 0$. In contrast, following iCMs, we use the pseudo-huber loss

$$d(\mathbf{x}_t, \hat{\mathbf{x}}_t) = \sqrt{\|\mathbf{x}_t - \hat{\mathbf{x}}_t\|_2^2 + c^2} - c \quad (30)$$

where $c = 0.00054\sqrt{d}$, where d is the dimension of \mathbf{x}_t .

Batch size. We use batch size 128 for 32×32 resolution images and batch size 64 for 64×64 resolution images.

Optimizer. We use the Adam optimizer (Kingma & Ba, 2015) with learning rate

$$\eta = 0.0002 / (128 / \text{batch_size}) \quad (31)$$

and default $(\beta_1, \beta_2) = (0.9, 0.999)$.

Coefficient for $\mathcal{L}_{\text{FM}}(\theta)$. We use $\lambda_{\text{FM}} = 0.1$ for all experiments.

Network. We modify SongUNet provided at <https://github.com/NVlabs/edm> to accept two time conditions t and s by using two time embedding layers.

ODE Solver. We use the second order Heun solver to calculate $\mathcal{L}_{\text{GCTM}}(\theta)$.

Gaussian perturbation. We apply a Gaussian perturbation from a normal distribution multiplied by 0.05 to sample \mathbf{x}_1 , excluding inpainting task.

A.2 EVALUATION

In this section, we describe the details of the evaluation to ensure reproducibility of our experiments.

Datasets. In unconditional generation task, we compare our GCTM generation performance using CIFAR10 training dataset. In image-to-image translation task, we evaluate the performance of models using test sets of Edges→Shoes, Night→Day, Facades from Pix2Pix. In image restoration task, we use FFHQ and apply following corruption operators H from I²SB to obtain measurement: bicubic super-resolution with a factor of 2, Gaussian deblurring with $\sigma = 0.8$, and center inpainting with Gaussian. We then assess model performance using test dataset.

Baselines. For I2I translation tasks (64×64 resolution), we compare three baselines: Pix2Pix from <https://github.com/junyanz/pytorch-CycleGAN-and-pix2pix>, Palette model from <https://github.com/Janspiry/Palette-Image-to-Image-Diffusion-Models>, and I²SB from <https://github.com/NVlabs/I2SB>. In the case of Pix2Pix, during training, we optimize the generator and discriminator with 256×256 resolution images using the recommended hyper-parameters. During inference, we apply bilinear interpolation to generator output images to resize them to 64×64 . For other baselines, we modify the model image resolution to 64×64 and use the recommended hyper-parameters. Same configuration is used in supervised image restoration task.

Metrics details. We calculate FID using <https://github.com/mseitzer/pytorch-fid> and IS from <https://github.com/pytorch/vision/blob/main/torchvision/models/inception.py>. We assess LPIPS from <https://github.com/richzhang/PerceptualSimilarity> with AlexNet version 0.1. In generation task, we employ the entire training dataset to obtain FID scores, and in the other task, we sample 5,000 test datasets. To obtain PSNR and SSIM, we convert the data type of model output to `uint8` and normalize it. We use <https://github.com/scikit-image/scikit-image> for PSNR and SSIM.

Sampling time. To compare inference speed, we measure the average time between the model taking in one batch size as input and outputting it.

B ALGORITHMS

B.1 OPTIMAL TRANSPORT

Algorithm 3 Sinkhorn-Knopp (SK)

- 1: **Input:** $\{\mathbf{x}_0^m\}_{m=1}^M, \{\mathbf{x}_1^m\}_{m=1}^M, \tau$
 - 2: Compute cost matrix $\mathbf{C} \in \mathbb{R}^{M \times M}$ subject to $C_{i,j} = \|\mathbf{x}_0^i - \mathbf{x}_1^j\|_2^2$
 - 3: Solve $\mathbf{P}^{\text{EOT}} = \arg \min_{\mathbf{P}} \langle \mathbf{P}, \mathbf{C} \rangle - \tau H(\mathbf{P})$ such that $\mathbf{P}\mathbf{1} = \mathbf{P}^\top \mathbf{1} = \frac{1}{n} \mathbf{1}$ with Algorithm 1 in Cuturi (2013)
 - 4: Treat \mathbf{P}^{EOT} as a discrete distribution over $\{1, \dots, M\} \times \{1, \dots, M\}$
 - 5: Sample $\{(i^m, j^m)\}_{m=1}^M \sim \mathbf{P}^{\text{EOT}}$
 - 6: **Return:** $\{(\mathbf{x}_0^{i^m}, \mathbf{x}_1^{j^m})\}_{m=1}^M$
-

B.2 IMAGE RESTORATION

In Alg. 4, we describe three zero-shot image restoration algorithms, DPS, CM, and GCTM. DPS uses the posterior mean $\mathbb{E}_{q(\mathbf{x}_0|\mathbf{x}'_{t_i})}[\mathbf{x}_0]$ to both traverse to a smaller time t_{i-1} and to approximate measurement inconsistency. As the posterior mean generally do not lie in the data domain, using it to calculate measurement inconsistency can be problematic. Indeed, approximation error in DPS is closely related to the discrepancy between the posterior mean and $\mathbf{x}'_{t_i \rightarrow 0}$ (see Theorem 1 in (Chung et al., 2022) for a formal statement). On the other hand, CM uses the ODE terminal point $\mathbf{x}'_{t_i \rightarrow 0}$ to traverse to a smaller time t_{i-1} and to approximate measurement inconsistency. While CM can have better guidance gradients as $\mathbf{x}'_{t_i \rightarrow 0}$ lie within the data domain, using $\mathbf{x}'_{t_i \rightarrow 0}$ to traverse to t_{i-1} can accumulate truncation error and degrade sample quality. For instance, see Figure 9 (a) in (Kim et al.,

Algorithm 4 Zero-shot Image Restoration

```

1: Input: Measurement  $\mathbf{x}_1$ , corruption  $\mathbf{H}$ , discretization  $\{t_i\}_{i=0}^M$ 
2:  $\mathbf{x}'_{t_M} \sim \mathcal{N}(0, \mathbf{I})$ 
3: for  $i = M$  to 1 do
4:    $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ 
5:   if Method is DPS then
6:      $\hat{\mathbf{x}}_0 = g_\theta(\mathbf{x}'_{t_i}, t_i, t_i)$ 
7:      $\mathbf{x}'_{t_{i-1}} = (1 - t_{i-1})\hat{\mathbf{x}}_0 + t_{i-1}\epsilon$ 
8:   else if Method is CM then
9:      $\hat{\mathbf{x}}_0 = G_\theta(\mathbf{x}'_{t_i}, t_i, 0)$ 
10:     $\mathbf{x}'_{t_{i-1}} = (1 - t_{i-1})\hat{\mathbf{x}}_0 + t_{i-1}\epsilon$ 
11:  else if Method is GCTM then
12:    Evaluate score and ODE endpoint in parallel by  $t = (t_i, t_i)$ ,  $s = (t_i, 0)$  :
13:     $\tilde{\mathbf{x}}_0, \hat{\mathbf{x}}_0 = g_\theta(\mathbf{x}'_{t_i}, t_i, t_i), G_\theta(\mathbf{x}'_{t_i}, t_i, 0)$ 
14:     $\mathbf{x}'_{t_{i-1}} = (1 - t_{i-1})\tilde{\mathbf{x}}_0 + t_{i-1}\epsilon$ 
15:  end if
16:   $\mathbf{x}'_{t_{i-1}} \leftarrow \mathbf{x}'_{t_{i-1}} - \lambda \nabla_{\mathbf{x}'_{t_i}} \|\mathbf{x}_1 - \mathbf{H}\hat{\mathbf{x}}_0\|_2^2$ 
17: end for
18: Return:  $\mathbf{x}'_0$ 

```

2024b). GCTM mitigates both problems by enabling parallel evaluation of posterior mean and ODE endpoint, as shown in Line 12-13 of Alg. 4.

B.3 IMAGE EDITING

Algorithm 5 Image Editing

```

1: Input:  $(\mathbf{x}_0, \mathbf{x}_1) \sim q(\mathbf{x}_0, \mathbf{x}_1)$ ,  $t$ 
2:  $\hat{\mathbf{x}}_t = (1 - t)\text{Edit}(\mathbf{x}_0) + t\mathbf{x}_1$ 
3: Return:  $G_\theta(\hat{\mathbf{x}}_t, t, 0)$ 

```

C PROOFS

C.1 CORRECTNESS OF SECTION 3.3

We show that our exposition in Section 3.3 adheres to Conditional Flow Matching (CFM) theory. Specifically, the notations

$$z, \quad q(z), \quad p_t(x|z), \quad u_t(x|z) \quad (32)$$

in Section 3 of Tong et al. (2024) are expressed in our paper as

$$(\mathbf{x}_0, \mathbf{x}_1), \quad q(\mathbf{x}_0, \mathbf{x}_1), \quad q(\mathbf{x}_t|\mathbf{x}_0, \mathbf{x}_1), \quad \mathbf{x}_1 - \mathbf{x}_0, \quad (33)$$

respectively. It follows that $p_t(x)$ and $u_t(x)$ in Tong et al. (2024) are expressed in our notation as

$$p_t(x) := \int p_t(x|z)q(z) dz = \int q(\mathbf{x}_t|\mathbf{x}_0, \mathbf{x}_1)q(\mathbf{x}_0, \mathbf{x}_1) d\mathbf{x}_0 d\mathbf{x}_1 = q(\mathbf{x}_t) \quad (34)$$

and

$$u_t(x) := \mathbb{E}_{q(z)} \frac{u_t(x|z)p_t(x|z)}{p_t(x)} = \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_1)} \frac{(\mathbf{x}_1 - \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0, \mathbf{x}_1)}{q(\mathbf{x}_t)} = \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_1|\mathbf{x}_t)} [\mathbf{x}_1 - \mathbf{x}_0] \quad (35)$$

respectively. By Theorem 3.1 in Tong et al. (2024), the ODE Eq. (11) indeed generates $q(\mathbf{x}_t)$, and Eq. (12) is equivalent to the CFM objective Eq. (10) in Tong et al. (2024).

C.2 PROOF OF PROPOSITION 1

Proof. We observe that the velocity term in (11) may be expressed as

$$\mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_1 | \mathbf{x}_t)}[\mathbf{x}_1 - \mathbf{x}_0] = \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_1 | \mathbf{x}_t)}[(\mathbf{x}_t - \mathbf{x}_0)/t] \quad (36)$$

$$= \mathbb{E}_{q(\mathbf{x}_0 | \mathbf{x}_t)}[(\mathbf{x}_t - \mathbf{x}_0)/t] \quad (37)$$

$$= (\mathbf{x}_t - \mathbb{E}_{q(\mathbf{x}_0 | \mathbf{x}_t)}[\mathbf{x}_0])/t \quad (38)$$

since \mathbf{x}_1 is determined given \mathbf{x}_0 and \mathbf{x}_t . This shows the equivalence between (11) and (13). Eq. (14) is then a straightforward consequence of the equivalence between ODEs. \square

C.3 PROOF OF PROPOSITION 2

Proof. We first show equivalence of scores. We note that

$$\mathbf{x}_\tau \mapsto \mathbf{x}_t \quad (39)$$

is a bijective transformation, so by change of variables,

$$q(\mathbf{x}_t | \mathbf{x}_0) = (1 + \tau) \cdot \mathcal{N}(\mathbf{x}_\tau | \mathbf{x}_0, \tau \mathbf{I}) = (1 + \tau) \cdot p(\mathbf{x}_\tau | \mathbf{x}_0) \quad (40)$$

and marginalizing out \mathbf{x}_0 , we get

$$q(\mathbf{x}_t) = (1 + \tau) \cdot p(\mathbf{x}_\tau). \quad (41)$$

It follows by Bayes' rule that

$$p(\mathbf{x}_0 | \mathbf{x}_\tau) = \frac{p(\mathbf{x}_\tau | \mathbf{x}_0)p(\mathbf{x}_0)}{p(\mathbf{x}_\tau)} \quad (42)$$

$$= \frac{(1 + \tau)^{-1} q(\mathbf{x}_t | \mathbf{x}_0) q(\mathbf{x}_0)}{(1 + \tau)^{-1} q(\mathbf{x}_t)} \quad (43)$$

$$= \frac{q(\mathbf{x}_t | \mathbf{x}_0) q(\mathbf{x}_0)}{q(\mathbf{x}_t)} \quad (44)$$

$$= q(\mathbf{x}_0 | \mathbf{x}_t) \quad (45)$$

and thus

$$\mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_\tau)}[\mathbf{x}_0] = \mathbb{E}_{q(\mathbf{x}_0 | \mathbf{x}_t)}[\mathbf{x}_0]. \quad (46)$$

for all $\tau \in (0, \infty)$ and \mathbf{x}_τ . We now show the equivalence of ODEs. Diffusion PFODE is

$$d\mathbf{x}_\tau = \frac{\mathbf{x}_\tau - \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_\tau)}[\mathbf{x}_0]}{\tau} d\tau. \quad (47)$$

With the change of variable

$$\mathbf{x}_t = \mathbf{x}_\tau / (1 + \tau), \quad (48)$$

we have

$$d\mathbf{x}_t = -\frac{\mathbf{x}_\tau}{(1 + \tau)^2} d\tau + \frac{1}{1 + \tau} d\mathbf{x}_\tau \quad (49)$$

$$= -\frac{\mathbf{x}_\tau}{(1 + \tau)^2} d\tau + \frac{\mathbf{x}_\tau - \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_\tau)}[\mathbf{x}_0]}{\tau(1 + \tau)} d\tau \quad (50)$$

$$= -\frac{\mathbf{x}_t}{1 + t} d\tau + \frac{(1 + t)\mathbf{x}_t - \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_\tau)}[\mathbf{x}_0]}{t(1 + t)} d\tau \quad (51)$$

$$= \frac{\mathbf{x}_t - \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_\tau)}[\mathbf{x}_0]}{\tau(1 + \tau)} d\tau \quad (52)$$

$$= \frac{\mathbf{x}_t - \mathbb{E}_{q(\mathbf{x}_0 | \mathbf{x}_t)}[\mathbf{x}_0]}{\tau(1 + \tau)} d\tau \quad (53)$$

where we have used equivalence of scores at the last line. We then make the change of time variable

$$t = \tau/(1 + \tau) \implies dt = \frac{1}{(1 + \tau)^2} d\tau \quad (54)$$

which gives us

$$d\mathbf{x}_t = \frac{\mathbf{x}_t - \mathbb{E}_{q(\mathbf{x}_0|\mathbf{x}_t)}[\mathbf{x}_0]}{\tau/(1 + \tau)} dt \quad (55)$$

$$= \frac{\mathbf{x}_t - \mathbb{E}_{q(\mathbf{x}_0|\mathbf{x}_t)}[\mathbf{x}_0]}{t} dt. \quad (56)$$

For the first equality in (20), transform PFODE variables (\mathbf{x}_τ, τ) into FM ODE variables (\mathbf{x}_t, t) with (17), transport \mathbf{x}_t to \mathbf{x}_s with G_{GCTM} , and then transform FM ODE variables (\mathbf{x}_s, s) into PFODE variables $(\mathbf{x}_\sigma, \sigma)$ with the inverse of (17). Second equality in (20) follows directly from (18). \square

D LIMITATION, SOCIAL IMPACTS, AND REPRODUCIBILITY

Limitations. GCTMs are yet unable to reach state-of-the-art unconditional generative performance. We speculate further tuning of hyper-parameters in the manner of iCMs could improve the performance, and leave this for future work.

Social impacts. GCTM generalizes CTM to achieve fast translation between any two distributions. Hence, GCTM may be used for beneficial purposes, such as fast medical image restoration. However, GCTM may also be used for malicious purposes, such as generation of malicious images, and this must be regulated.

Reproducibility. We open-source our code at <https://github.com/1202kbs/GCTM> including training code for unconditional generation, image-to-image translation, and supervised image restoration models.

E ADDITIONAL EXPERIMENTS

E.1 COMPARING I2I PERFORMANCE WITH OTHER BASELINE MODELS

We compare the image-to-image (I2I) performance of our model with two baseline approaches: EGSDE (Zhao et al., 2022) and BBDM (Li et al., 2023). Since BBDM, an I2I framework based on the Brownian Bridge process, leverages a latent diffusion model, we train it with a pixel-space diffusion model for a fair comparison. Both BBDM and EGSDE are trained on the Edges→Shoes dataset. As shown in Table 4, our GCTM outperforms all baselines across various metrics, even when evaluated with fewer sampling steps.

In addition, we visualize the image editing results in Fig. 9. While EGSDE generates realistic images, it fails to faithfully preserve the given conditions. BBDM, on the other hand, struggles to perform robustly on (unseen) conditional images. In contrast, GCTM produces realistic images while accurately maintaining the original conditions.

Method	NFE	Time (ms)	FID ↓	IS ↑	LPIPS ↓
BBDM (Li et al., 2023)	5	75	43.7	3.43	0.099
EGSDE (Zhao et al., 2022)	500	2590	198.1	2.87	0.476
GCTM	1	87	40.3	3.54	0.097

Table 4: Evaluation of I2I translation on Edges→Shoes with other baselines.



Figure 9: Comparison on image editing with GCTM and other baselines

E.2 CONTROLLABLE IMAGE EDITING

In this section, we demonstrate that effectiveness of image editing can be controlled. In Algorithm 5, we control the time point t to determine how much of the edited image to reflect. In Fig. 10, the results visualize how t effect the output of model output. We observe that the larger t , the more realistic the image, and the smaller t , the more faithful the edit feature. We set $t = 0.95$ and $t = 0.4$ at supervised coupling and independent coupling, respectively.

Figure 10: Controllability of image editing by t .

E.3 HIGH-RESOLUTION IMAGE-TO-IMAGE TRANSLATION

To verify the robustness of our framework to scalability of resolution, we experiment with image-to-image translation of the Facades dataset with 256×256 resolution (Fig. 11). In Table 5, we see that GCTM achieves high performance despite using fewer NFEs compared to the baselines, generating realistic, diverse, and faithful translated images.

Method	NFE	Facades-256		
		FID ↓	IS ↑	LPIPS ↓
Pix2Pix (Isola et al., 2017)	1	<u>117.2</u>	1.60	0.414
Palette (Saharia et al., 2022)	5	396.7	1.14	1.089
I ² SB(Liu et al., 2023)	5	128.6	<u>2.23</u>	0.454
GCTM	1	107.0	2.24	<u>0.426</u>

Table 5: Quantitative evaluation of I2I translation with 256×256 resolution images. Best is in **bold**, and second best is underlined.

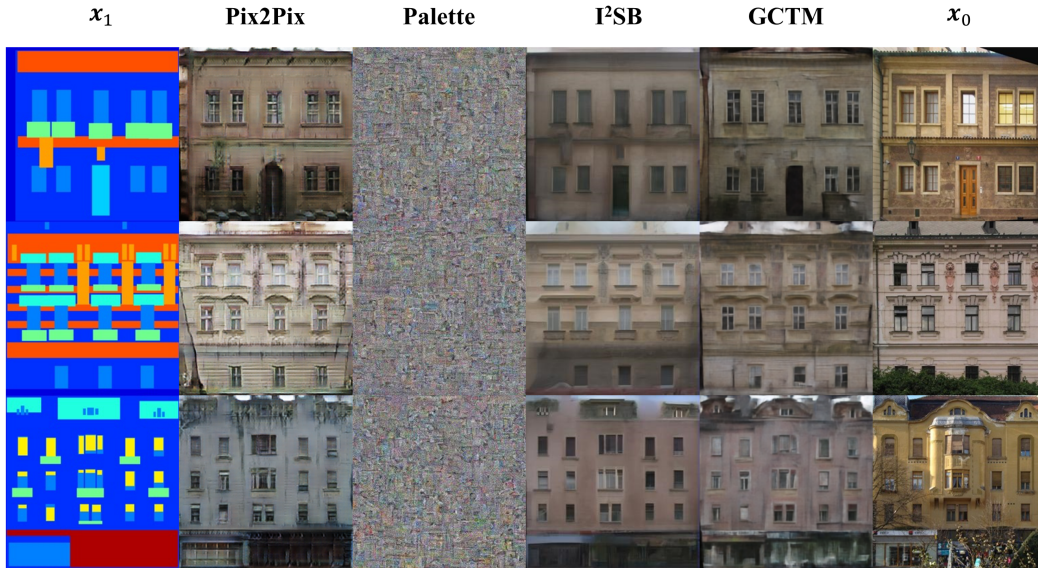


Figure 11: Qualitative comparison of I2I task on Facades 256×256 .

E.4 HIGH-RESOLUTION IMAGE RESTORATION

In Table 6, we demonstrate image restoration task of GCTM on ImageNet with higher resolution images. As diffusion-based solvers for inverse problems, both DPS (Chung et al., 2022) and DDS (Chung et al., 2024) require a sufficient number of NFEs to achieve effective reconstruction. Specifically for DPS, using significantly fewer NFEs than the 1000 NFEs suggested in the original paper, combined with the absence of measurement noise, disrupts the reconstruction process, resulting in outputs that are worse than GCTM. Although DDS shows improved performance compared to DPS, it requires more NFEs to achieve comparable performance with GCTM, which emphasizes the efficiency of GCTM in solving image restoration tasks.

Method	NFE	SR4 - Bicubic			Deblur - Gaussian			Inpaint - Center		
		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
DPS	10	10.37	0.357	0.727	10.27	0.256	0.830	9.98	0.247	0.841
	50	16.15	0.392	0.654	19.19	0.520	0.523	13.61	0.526	0.522
	1000	<u>22.36</u>	<u>0.601</u>	<u>0.327</u>	26.29	<u>0.739</u>	0.246	<u>18.53</u>	<u>0.681</u>	<u>0.288</u>
DDS	10	19.79	0.569	0.491	21.12	0.634	0.394	13.09	0.503	0.531
	50	21.25	0.571	0.409	23.33	0.704	<u>0.245</u>	13.57	0.485	0.511
GCTM	1	26.70	0.771	0.223	34.65	0.948	0.032	21.56	0.808	0.229

Table 6: GCTM evaluation of image restoration on ImageNet with 256×256 resolution. Best is in **bold**, and second best is underlined.

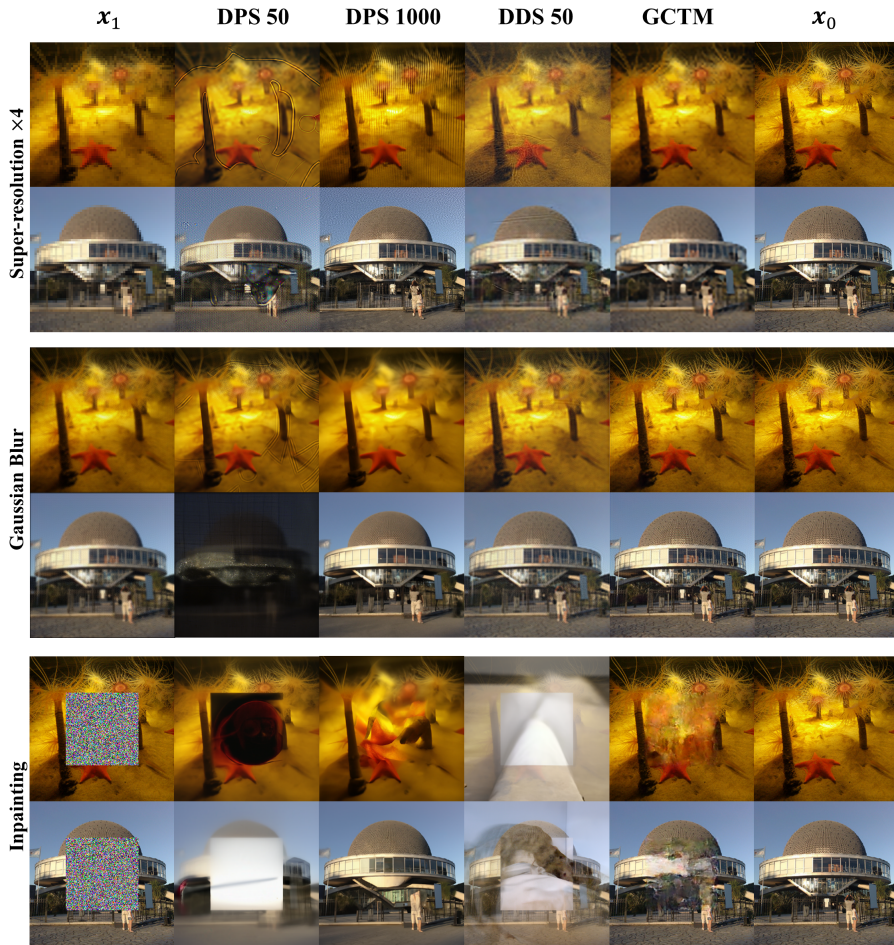


Figure 12: Qualitative comparison of image restoration task on ImageNet 256×256 .

E.5 ADDITIONAL IMAGE-TO-IMAGE TRANSLATION SAMPLES



Figure 13: Additional results on image-to-image translation task on Edges→Shoes (top), Facades (middle) and Night→Day (bottom).

E.6 ADDITIONAL IMAGE RESTORATION SAMPLES



Figure 14: Additional results of supervised image restoration task on FFHQ 64×64 .



Figure 15: Additional results of zero-shot image restoration task on FFHQ 64×64 .

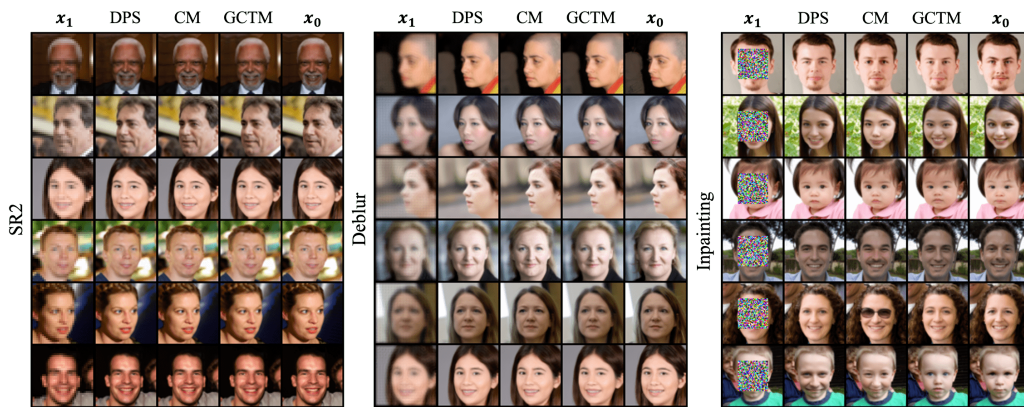


Figure 16: Qualitative comparison of zero-shot algorithms on FFHQ 64×64 .