COT FLOW: LEARNING OPTIMAL-TRANSPORT IMAGE SAMPLING AND EDITING BY CONTRASTIVE PAIRS

Anonymous authors

Paper under double-blind review



Figure 1: (a). Unpaired image-to-image translation by our proposed COT Flow, with one-step or multi-step sampling. (b), (c). Our proposed COT Editor enables zero-shot image editing with high flexibility. COT composition (b) allows users to composite elements and synthesize realistic images. Shape-texture coupling (c) allows users to separately draw or use shapes and textures as dual inputs, to generate fused images with high quality.

ABSTRACT

Diffusion models have demonstrated strong performance in sampling and editing multi-modal data with high generation quality, yet they suffer from the iterative generation process which is computationally expensive and slow. In addition, most methods are constrained to generate data from Gaussian noise, which limits their sampling and editing flexibility. To overcome both disadvantages, we present *Contrastive Optimal Transport Flow (COT Flow)*, a new method that achieves fast and high-quality generation with improved zero-shot editing flexibility compared to previous diffusion models. Benefiting from optimal transport (OT), our method has no limitation on the prior distribution, enabling unpaired image-to-image (I2I) translation and *doubling* the editable space (at both the start and end of the trajectory) compared to other zero-shot editing methods. In terms of quality, COT

Flow can generate competitive results in merely one step compared to previous state-of-the-art unpaired image-to-image (I2I) translation methods. To highlight the advantages of COT Flow through the introduction of OT, we introduce the *COT Editor* to perform user-guided editing with excellent flexibility and quality.

)58

060 061

054

056

1 INTRODUCTION

Diffusion models, with flexible training and sampling principles rooted in Statistical Physics, have 062 achieved unprecedented success in generating data from noise Ho et al. (2020); Song & Ermon 063 (2019); Ramesh et al. (2022); Saharia et al. (2022); Nichol et al. (2021); Rombach et al. (2021); 064 Dhariwal & Nichol (2021); Ho & Salimans (2022); Kazerouni et al. (2023); Janner et al. (2022); 065 Poole et al. (2022); Li et al. (2023); Liu et al. (2024). However, the fundamental limitations of 066 diffusion-based models, namely the sampling inefficiency and restrictive prior distribution, still 067 barricade them from wider applications, despite the recent series of improved methods Nichol & 068 Dhariwal (2021); Karras et al. (2022); Song et al. (2023). With a similar iterative sampling process, 069 flow-based methods Chen et al. (2018); Kidger et al. (2020) also suffer from the computational inefficiency problem. From a high-level perspective, the current deep generative models still cannot simultaneously satisfy three performance indicators: (1) high-quality generation, (2) mode coverage 071 and diversity, and (3) fast sampling, which is identified as the generative learning trilemma Xiao 072 et al. (2021) shown in Fig.2a. 073

074 To tackle the generative learning trilemma and eliminate the constraints on prior distribution, we 075 present a novel flow-based model called Contrastive Optimal Transport Flow (COT Flow), which 076 fundamentally addresses the computational inefficiency problem through the optimal transport (OT) formulation. We claim that OT enables the fastest sampling for diffusion/flow-based methods with 077 two key features to overcome sampling inefficiency: (1) straight lines from source to target and (2) no crossing among the trajectories. Similar principles were approached implicitly in a few latest 079 work Liu et al. (2022); Lipman et al. (2022); Tong et al. (2023); Esser et al. (2024); Karras et al. (2022). Specifically, many recent breakthroughs Song et al. (2020a); Nichol & Dhariwal (2021); 081 Karras et al. (2022) focused on the following strategies: optimizing the sampling trajectories towards straight lines, improving the time schedule of the diffusion process Song et al. (2020a); Karras et al. 083 (2022), adjusting the noise schedule or forward diffusion process Song et al. (2020b); Nichol & 084 Dhariwal (2021); Lin et al. (2023); Bartosh et al. (2024), introducing fast samplers Nichol & Dhari-085 wal (2021); Lu et al. (2022a;b); Karras et al. (2022), using distillation techniques Song et al. (2023); Liu et al. (2022); Salimans & Ho (2022); Xu et al. (2023); Meng et al. (2022), and eliminating the 087 crossing among the trajectories to improve sample stability and efficiencyLiu et al. (2022); Lipman et al. (2022); Tong et al. (2023); Esser et al. (2024). We note that these improved techniques, though 880 from different angles, approached the similar concept of OT between Gaussian and data distribution, 089 as shown in Fig.2b. Another prominent group of recent works (Korotin et al. (2022a;b); Fan et al. 090 (2021; 2022); Rout et al. (2021)) enforce direct OT by training two neural networks on saddle point 091 problems Boyd & Vandenberghe (2004). 092

⁰⁹³ The proposed COT Flow satisfies the three performance requirements in the trilemma:

Sample efficiency: The proposed COT Flow explicitly builds the bridge between diffusion/flow-based models and OT, and thus enforces straight trajectories and eliminates the crossing to improve sample efficiency. With the benefit of both diffusion/flow-based models and the OT formulation, COT Flow enables one-step or few-step sampling by design, while still producing high-quality and high-diversity results from arbitrary prior distributions. Furthermore, COT Flow allows zero-shot editing, and introduces diverse editing possibilities (Fig.1b).

100 Sample quality: COT Flow leverages the intriguing similarities between consistency models Song 101 et al. (2023); Song & Dhariwal (2023); Luo et al. (2023) and contrastive learning He et al. (2019); 102 Chen & He (2020); Chen et al. (2020); Grill et al. (2020) to produce high-quality generation using 103 indirect loss functions. In particular, the objective of consistency models consists of the similarity 104 between time-adjacent data pairs $\langle \mathbf{x}_t, \mathbf{x}_{t+1} \rangle$, which function exactly the same as the positive sample 105 pairs in contrastive learning (He et al. (2019) Eq.1). In addition, consistency models use a series of similar techniques as those in contrastive learning, such as exponential moving average (EMA) 106 weights of the teacher model and "stopgrad" operator Song et al. (2023); Song & Dhariwal (2023); 107 Chen & He (2020); Grill et al. (2020), suggesting the hidden link between the two state-of-the-art



Figure 2: (a). The generative learning trilemma. Current generative methods still cannot simulta-118 neously satisfy the three performance indicators: high quality, fast sampling, and mode coverage. 119 (b). Recent developments of the diffusion/flow-based generative models, including iDDPMNichol 120 & Dhariwal (2021), EDMKarras et al. (2022), DDIMSong et al. (2020a), DPMLu et al. (2022a), Pro-121 gressive Distillation (PD)Salimans & Ho (2022), Consistency Distillation (CD)Song et al. (2023), 122 VP ODESong et al. (2020b), Flow Matching (FM)Lipman et al. (2022), Conditional Flow Matching 123 (CFM)Tong et al. (2023), Rectified Flow (RF)Liu et al. (2022), Stable Diffusion v3 (SDv3)Esser 124 et al. (2024) All methods implicitly approach the OT formulation, either by sampling straight trajec-125 tories or avoiding crossing between the trajectories through various techniques.)

learning frameworks. Enlightened by this connection, we introduce the *Contrastive OT Pairs (COT Pairs)* for positive pair sampling during COT Flow training. By using a similar contrastive loss as in Grill et al. (2020), we consider the proposed COT Flow model as a powerful contrastive learning encoder \mathcal{E} to map all data points on the OT trajectories towards their end. We evaluate COT Flow's sample quality via the FID scores in various unpaired I2I translation tasks such as handbags—shoes, CelebA male—female, and outdoor—church (Fig.1a).

134 Mode coverage: COT Flow achieves competitive sample diversity and mode coverage compared to diffusion models, benefiting from the non-adversarial contrastive loss and the OT formulation. The 135 adversarial objectives in Generative Adversarial Nets (GAN)Goodfellow et al. (2014), Wasserstein 136 GANArjovsky et al. (2017), and StyleGANKarras et al. (2018; 2019) are susceptible to training 137 instability and mode collapse Xiao et al. (2021), which even the state-of-the-art GAN-based methods 138 still suffer from Park et al. (2020). Diffusion-based model objectives, on the other hand, are closely 139 related to the Evidence Lower Bound (ELBO) of the target data and are thus less prone to training 140 instability and mode collapse Ho et al. (2020); Kingma & Gao (2023). In addition, with the OT 141 formulation, the proposed COT Flow minimizes the transportation cost and directly maps the source 142 distribution to the target distribution, improving the faithfulness to the target data. 143

In summary, our main contributions are: (1) We tackle the generative learning trilemma by intro-144 ducing a novel framework called Contrastive Optimal Transport Flow (COT Flow), which explicitly 145 combines diffusion/flow-based model with OT to directly learn the generative flow between any two 146 unpaired data sources. (2) We present the Contrastive Optimal Transport Pair (COT Pair) formu-147 lation to train our proposed COT Flow, leveraging the intriguing connection between consistency 148 models and contrastive learning. (3) To showcase the advantages of COT Flow, we introduce the 149 COT Editor to perform controllable sampling and flexible zero-shot image editing, including COT 150 composition, shape-texture coupling, and COT augmentation, and demonstrate these functionalities 151 via diverse data and application scenarios.

152 153 154

2 BACKGROUND

155
156 COT Flow leverages the theories and concepts from (1) optimal transport Villani (2009), (2) contrastive learning He et al. (2019), and (3) consistency models Song et al. (2023), crossing these three prominent methodologies in optimization and machine learning. For a quick understanding of the proposed COT Flow, we first briefly present the three core methodologies and discuss their interconnections in Section 3.1.

161 Notations. Throughout the paper, \mathcal{X} and \mathcal{Y} denote two metric spaces of data, $\mu(\mathbf{x})$ and $\nu(\mathbf{y})$ denote the probability distributions on \mathcal{X} and \mathcal{Y} , respectively. For describing the projection between $\mu(\mathbf{x})$

and $\nu(\mathbf{y})$, we denote $T : \mathcal{X} \to \mathcal{Y}$ as a measurable map, which satisfies: for any measurable subsets $B \subset \mathcal{Y}, T^{-1}(B) \subset \mathcal{X}$. We denote $\Pi(\mu, \nu)$ as the set of joint probability distributions on $\mathcal{X} \times \mathcal{Y}$ whose marginals are μ and ν .

166 2.1 OPTIMAL TRANSPORT

The optimal Transport (OT) problem seeks the minimum overall transportation cost from one measure to another. Consider a cost function $c : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$, Kantorovitch (1958) formulates a transport coupling $\pi \in \Pi(\mu, \nu)$ and introduces the OT cost:

$$\operatorname{Cost}(\mu,\nu) := \inf_{\pi \in \Pi(\mu,\nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(\mathbf{x}, \mathbf{y}) d\pi(\mathbf{x}, \mathbf{y})$$
(1)

This is defined as the Kantorovich problem, where the infimum is taken over transport couplings $\pi \in \Pi(\mu, \nu)$. The optimal π^* is called the OT plan, which always exists under mild conditions on spaces \mathcal{X}, \mathcal{Y} and cost function c (Villani (2009)). According to the duality principle Boyd & Vandenberghe (2004), the dual problem of Kantorovich's optimization is:

$$\operatorname{Cost}(\mu,\nu) := \sup_{\varphi,\psi} \left\{ \int_{\mathcal{X}} \varphi(\mathbf{x}) d\mu(\mathbf{x}) + \int_{\mathcal{Y}} \psi(\mathbf{y}) d\nu(\mathbf{y}) \right\}$$
(2)

where $\varphi \in L^1(\mu)$ and $\psi \in L^1(\nu)$ are called Kantorovich potentials which satisfy $\varphi(\mathbf{x}) + \psi(\mathbf{y}) \leq c(\mathbf{x}, \mathbf{y})$. For $\varphi : \mathcal{X} \to \mathbb{R}$, $\psi : \mathcal{Y} \to \mathbb{R}$, and a certain cost function *c*, we replace the first potential $\varphi(\mathbf{x})$ by defining the *c*-transform of ψ : $\varphi(\mathbf{x}) = \psi^c(\mathbf{x}) = \inf_{\mathbf{y} \in \mathcal{Y}} \{c(\mathbf{x}, \mathbf{y}) - \psi(\mathbf{y})\}$, and the Kantorovich problem 2 is rewritten as:

$$\operatorname{Cost}(\mu,\nu) := \sup_{\psi} \left\{ \int_{\mathcal{X}} \inf_{\mathbf{y}} \{ c(\mathbf{x},\mathbf{y}) - \psi(\mathbf{y}) \} d\mu(\mathbf{x}) + \int_{\mathcal{Y}} \psi(\mathbf{y}) d\nu(\mathbf{y}) \right\}$$
(3)

where we denote the right side of 3 as a saddle point problem $\sup_{\psi} \inf_{\mathbf{y}} \mathcal{L}(\psi, \mathbf{y})$, whose solution (ψ^*, \mathbf{y}^*) contains the optimal choice of \mathbf{y} given a certain \mathbf{x} . In practice, \mathbf{y}^* can be estimated by optimizing a neural network $\tilde{\mathbf{y}} = T_{\theta}(\mathbf{x})$, leading to neural OT methods Korotin et al. (2022a;b); Fan et al. (2021; 2022). We further illustrate the training of $T_{\theta}(\mathbf{x})$ in Section 3.

190 191 192 193

199 200

188

189

165

171 172

173

178

179 180

181

182

183

2.2 CONTRASTIVE LEARNING

With impressive results on multiple visual tasks, contrastive learning methods learn data representations by attracting the embeddings of positive sample pairs and (optionally) repulse the embeddings of negative sample pairs in an unsupervised manner He et al. (2019); Chen et al. (2020). For the methods that only consider the positive pairs Chen & He (2020); Grill et al. (2020), the core method-ology can be described as minimizing the loss function:

$$\mathcal{L}(\theta, \theta^{-}) := d(q_{\theta}(\mathcal{E}_{\theta}(\mathbf{x})), \mathcal{E}_{\theta^{-}}(\mathbf{x}^{+}))$$
(4)

where \mathcal{E} is the target network, which we consider as an encoder. θ^- denotes the exponential moving average (EMA) of the past values of the network's weights θ . $d(\cdot, \cdot)$ is the distance function between the data embedding $\mathcal{E}(\mathbf{x})$ and its corresponding positive pairs $\mathcal{E}(\mathbf{x}^+)$, whose inputs \mathbf{x}^+ are augmented from the same sample \mathbf{x} . Combined with the EMA weights θ^- and the "stopgrad" operator, an additional prediction head q_{θ} is introduced on top of the encoder \mathcal{E}_{θ} to prevent model collapse and enable the contrastive learning methods to produce meaningful representations. In Section 3, we introduce the similarities between contrastive learning and consistency models.

208 209 2.3 CONSISTENCY MODELS

Consistency models (CMs) are an emerging family of generative models whose key idea is maintaining consistency along the ordinary differential equation (ODE) trajectory derived from the diffusion models, which we briefly introduce in Appendix E. One drawback of diffusion models is their slow sampling speed. CMs, on the other hand, learn the consistency along the trajectories $\{\hat{\mathbf{x}}_t\}_{t\in[0,T]}$ of the probability flow ODE 28 and map all the points on these trajectories to their origin $\hat{\mathbf{x}}_0$. This mapping can be described as the consistency function $\mathbf{f}^* : (\mathbf{x}_t, t) \to \mathbf{x}_0$ which satisfies the boundary condition $\mathbf{f}^*(\mathbf{x}, 0) = \mathbf{x}_0$. We then approximate $\mathbf{f}^*(\mathbf{x}, t)$ by training the consistency model $\mathbf{f}_{\theta}(\mathbf{x}_t, t)$. By discretizing the probability flow ODE 28 with a limited sequence of time steps $\epsilon < t_1 < t_2 < ... < t_N = T$, the consistency model $\mathbf{f}_{\theta}(\mathbf{x}_t, t)$ is trained by minimizing the consistency matching loss (CM loss):

231

232

233 234

235 236

237

238

239

240

241

242

243

244 245 246

247

248

249

250 251

252

 $\mathcal{L}^{N}(\theta, \theta^{-}) := \mathbb{E}[\lambda(t_{i})d(\mathbf{f}_{\theta}(\mathbf{x}_{t_{i+1}}, t_{i+1}), \mathbf{f}_{\theta^{-}}(\mathbf{x}_{t_{i}}, t_{i}))], i \sim \mathcal{U}[1, N-1]$ (5)

where $\mathbf{x}_{t_{i+1}}$ is sampled from the distribution $p_{t_{i+1}}(\mathbf{x})$ and the parameter θ^- is the EMA of θ obtained 221 with the "stopgrad" operator $\theta^- \leftarrow \text{stopgrad}(\mu\theta^- + (1-\mu)\theta)$. $0 \le \mu < 1$ denotes the EMA decay rate. $\lambda(t_i) > 0$ is a weighting function and $d(\cdot, \cdot)$ is a distance function with a typical choice 222 223 of squared l_2 . $\mathcal{U}[1, N-1]$ denotes the uniform distribution over 1, 2, ..., N-1. For \mathbf{x}_{t_i} , CMs 224 provide two approximations and correspondingly form two training algorithms called consistency 225 distillation (CD) and consistency training (CT). The approximation from CD is $\hat{\mathbf{x}}_{t_i} = \mathbf{x}_{t_{i+1}}$ – 226 $(t_i - t_{i+1})t_{i+1}\mathbf{s}_{\phi}(\mathbf{x}_{t_{i+1}}, t_{i+1})$, which relies on a pre-trained diffusion model $\mathbf{s}_{\phi}(\mathbf{x}, t)$. While the 227 approximation from CT is $\hat{\mathbf{x}}_{t_i} = \mathbf{x} + t_i \mathbf{z}$ where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is the same noise when forming $\mathbf{x}_{t_{i+1}} = \mathbf{x} + t_{i+1} \mathbf{z}$. We can directly sample the final generation by $\mathbf{x}_0 = \mathbf{f}_{\theta}(\mathbf{z}, t_N)$ or optionally 228 229 sample the intermediate results $\mathbf{x}_k = \mathbf{f}_{\theta}(\mathbf{x}_{k+1}, t_{i_{k+1}}) + \sqrt{t_N^2 - \epsilon^2} \mathbf{z}_k$ for k = K - 1, ..., 1. 230

Comparing the CM loss $\mathcal{L}^{N}(\theta, \theta^{-})$ in 5 and the contrastive learning loss in 4, we observe both structural and conceptual similarities between them, which will be discussed in Section 3.1.

3 Method



Figure 3: An overview of the training process. COT Flow minimizes the distances between the encodings of the positive pairs, which are sampled in the augmentation area between \mathbf{x} in Data 1 and its OT mapping $T_{\phi}(\mathbf{x})$ (Eq.12).

Our proposed COT Flow tackles the generative learning trilemma by fundamentally regularizing the transportation flows between two distributions. COT Flow consists of three main parts: (1) COT Pairs, (2) COT training, and (3) COT Editor. In the sections below, we first discuss the similarities between CMs and contrastive learning, which inspire our formulation of COT Pairs and COT training. We then introduce the COT Editor framework.

3.1 SIMILARITIES BETWEEN CONTRASTIVE LEARNING AND CONSISTENCY MODELS

Cone may raise a question on the mechanism of CMs: Why do they work well by simply minimizing the difference between two points on the same trajectory, especially with no guidance of the trajectory's origin \mathbf{x}_0 in the loss function? Here we put forward a hypothesis on why they learn to map to the origin by exploring the systematic similarities between CMs and contrastive learning: *The consistency function* $\mathbf{f}_{\theta}(\mathbf{x}, t)$ *is a trajectory's origin encoder* $\mathcal{E}_{\theta}(\mathbf{x})$, which has the same functionality of the encoder in contrastive learning.

Firstly, we notice the similarity between the CM loss 5 and the contrastive loss 4, which are both summarized by a distance metric $d(\cdot, \cdot)$. Specifically, the CM loss indicates the distance between the two output points $\mathbf{f}_{\theta}(\mathbf{x}_{t_{i+1}}, t_{i+1})$ and $\mathbf{f}_{\theta^-}(\mathbf{x}_{t_i}, t_i)$ from the same trajectory, while the contrastive loss indicates the distance between the embeddings of the positive pairs $\mathcal{E}_{\theta}(\mathbf{x})$ and $\mathcal{E}_{\theta^-}(\mathbf{x}^+)$ from the same image. This suggests that CMs have the capability of learning representations from complex distributions and are capable of mapping denoising trajectories $\{\mathbf{x}_t\}_{t \in [\epsilon, T]}$ to their origins \mathbf{x}_0 .

Secondly, the strategies and training recipes of the two methods are similar, especially those for preventing mode collapsing. They both utilize weight-sharing Siamese networks θ , θ^- to minimize the distance metric $d(\cdot, \cdot)$ of the entities, and they both use "stopgrad" operations to distinguish the networks and prevent collapsing:

$$\theta^{-} \leftarrow \theta^{-} - \eta \nabla_{\theta} d(\mathcal{E}_{\theta}(\cdot), \operatorname{stopgrad}(\mathcal{E}_{\theta^{-}}(\cdot)))$$
(6)

270 Furthermore, the recent work from both sides Chen & He (2020); Song & Dhariwal (2023) illus-271 trated a common improvement to optimize the results and simplify the strategies: removing the 272 EMA decay for the Siamese structure, whose weights share the same update ∇_{θ} . This improvement 273 has been proven effective from both sides Chen & He (2020); Song & Dhariwal (2023), underlining 274 the same mechanism between CMs and contrastive learning.

275 With the above observations, we explain the capability of the consistency function $f_{\theta}(\mathbf{x},t)$ to map 276 the intermediates towards the origin by considering the consistency function $f_{\theta}(\mathbf{x},t)$ as the encoder 277 $\mathcal{E}_{\theta}(\mathbf{x})$ in contrastive learning. With this foundation, we introduce COT Pairs and COT training in 278 the following sections. 279

3.2 COT PAIRS

280

281

285

286 287

297 298 299

301

303

304

306

307 308

310 311

312

318 319

282 In Section 2.1, we introduce the Kantorivich problem. The entropic regularization of the Kan-283 torovich problem, namely the entropic OT (EOT) problem Villani (2009), minimizes the transporta-284 tion cost derived from 1:

$$\operatorname{Cost}(\mu,\nu) := \inf_{\pi \in \Pi(\mu,\nu)} \left\{ \int_{\mathcal{X} \times \mathcal{Y}} c(\mathbf{x}, \mathbf{y}) d\pi(\mathbf{x}, \mathbf{y}) + \lambda H(\pi) \right\}$$
(7)

288 where the solution π_{λ}^{*} is the EOT plan. With the relative entropy $\lambda H(\pi)$, the expensive computation 289 in the exact OT problem is alleviated. For neural OT models, using EOT enables stochastic processes 290 within the OT mapping and relates OT with diffusion models Gushchin et al. (2022). In the following 291 Eq.12, we introduce noise into COT training, where Proposition 3.1 shows its relationship to the EOT plan. 292

293 We modify a neural OT model to estimate the OT map between the two data distributions. Ac-294 cording to Section 2.1, the solution (ψ^*, \mathbf{y}^*) of the Kantorovich problem 3 can be estimated by two 295 corresponding networks $(\psi_{\omega}, T_{\phi}(\mathbf{x}))$, resulting in the neural OT objective: 296

$$\operatorname{Cost}(\mu,\nu) := \sup_{\psi_{\omega}} \left\{ \inf_{T_{\phi}} \int_{\mathcal{X}} \left\{ c(\mathbf{x}, T_{\phi}(\mathbf{x})) - \psi_{\omega}(T_{\phi}(\mathbf{x})) \right\} d\mu(\mathbf{x}) + \int_{\mathcal{Y}} \psi_{\omega}(\mathbf{y}) d\nu(\mathbf{y}) \right\}$$
(8)

where ψ_{ω} denotes the Kantorovich potential in Section 2.1 and T_{ϕ} is the estimated OT map. The 300 infimum of T_{ϕ} is interchanged with the integral by Rockafellar (1976) and the OT problem 1 is derived into the optimization of the neural networks: 302

$$\sup_{\omega} \inf_{\phi} \mathcal{L}(\psi_{\omega}, T_{\phi}) \tag{9}$$

To approach 9 in implementation, we optimize the parameters ω , ϕ using stochastic gradient ascentdescent (SGAD) by sampling mini-batch data from source and target datasets $\mathbf{x} \sim \mu(\mathbf{x}), \mathbf{y} \sim \nu(\mathbf{y})$:

$$\omega \leftarrow \omega + \nabla_{\omega} \left\{ -\frac{1}{|\mathbf{x}|} \sum_{\mathbf{x} \in \mathcal{X}} \psi_{\omega} (T_{\phi}(\mathbf{x})) + \frac{1}{|\mathbf{y}|} \sum_{\mathbf{y} \in \mathcal{Y}} \psi_{\omega}(\mathbf{y}) \right\}$$
(10)

$$\phi \leftarrow \phi - \nabla_{\phi} \left\{ \frac{1}{|\mathbf{x}|} \sum_{\mathbf{x} \in \mathcal{X}} \left[c \left(\mathbf{x}, T_{\phi}(\mathbf{x}) \right) - \psi_{\omega} \left(T_{\phi}(\mathbf{x}) \right) \right] \right\}$$
(11)

313 where $|\mathbf{x}|, |\mathbf{y}|$ denote the sizes of the corresponding mini-batches $\mathbf{x} \sim \nu(\mathbf{x}), \mathbf{y} \sim \mu(\mathbf{y})$. $c(\cdot, \cdot)$ 314 denotes the cost function in 3 which is typically l_2 -norm. Based on the trained $T_{\phi}(\mathbf{x})$ in 10 and 11, 315 we interpolate an augmentation area between $\mu(\mathbf{x})$ and $\nu(\mathbf{y})$ for training COT Flow, whose concept 316 "augmentation" derives from contrastive learning: 317

$$\{\tilde{\mathbf{x}}_t\}_{t\in[0,1]} = \{tT_{\phi}(\mathbf{x}) + (1-t)\mathbf{x} + t(1-t)\sigma^2 \mathbf{z}\}_{t\in[0,1]}$$
(12)

where σ is the noise scale and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is standard Gaussian noise. We prove that the OT plan 320 π^* in 1 can be extended in $t \in [0, 1]$ by formulating this augmentation area: 321

Proposition 3.1 (Eq.12 estimates the dynamic extension of the OT plan). Let π^* be the OT plan 322 between $\mu(\mathbf{x})$ and $\nu(\mathbf{y})$. Let the OT map T^* recovers π^* . The augmentation defined by Eq.12 using 323 T^* samples the same probability as the dynamic extension of the EOT plan π_{λ}^* with $\lambda = 2\sigma^2$.

We provide the proof in Appendix B. With the guarantee of Proposition 3.1 and the observation in Section 3.1, we consider the augmentations $\tilde{\mathbf{x}}_t$ as the intermediates of the entropic OT trajectory $\{\tilde{\mathbf{x}}_t\}_{t\in[0,1]}$ and formulate a set of positive pairs as in contrastive learning, which we name as COT Pairs. In particular, COT Pairs $\langle \mathbf{x}_{t_1}, \mathbf{x}_{t_2} \rangle$ are randomly selected along the trajectory $\{\tilde{\mathbf{x}}_t\}_{t\in[0,1]}$:

$$\mathbf{x}_{t_1}, \mathbf{x}_{t_2} \in \{\tilde{\mathbf{x}}_t\}_{t \in [0,1]}, \quad 0 \le t_1 < t_2 \le 1$$
(13)

Unlike CMs choosing adjacent pairs from ODE solvers, we formulate random COT pairs in the proposed augmentation area in Eq.12.

3.3 COT TRAINING

330

331

332 333

334

338 339

350

351

352

353

354

355

356

357

362

364

According to the relationship between contrastive learning and CMs discussed in section 3.1, we consider the consistency function $\mathbf{f}_{\theta}(\cdot)$ as an encoder $\mathcal{E}(\mathbf{x}_t)$ towards the origins \mathbf{y} of the entropic OT trajectories $\{\tilde{\mathbf{x}}_t\}_{t\in[0,1]}$. The COT training loss to optimize the origin encoder $\mathcal{E}(\mathbf{x}_t, t)$ is:

$$\mathcal{L}_{\text{COT}}(\theta) = d\big(\mathcal{E}_{\theta}(\mathbf{x}_{t_1}, t_1), \mathcal{E}_{\theta}(\mathbf{x}_{t_2}, t_2)\big), \quad 0 \le t_1 < t_2 \le 1$$
(14)

where $d(\cdot, \cdot)$ denotes the dissimilarity function, which is l_2 -norm by default and $\mathbf{x}_{t_1}, \mathbf{x}_{t_2}$ is the COT Pair from $\{\tilde{\mathbf{x}}_t\}_{t\in[0,1]}$. Inspired by Esser et al. (2024), the origin estimation $\mathcal{E}(\mathbf{x}_t)$ is more difficult for t in the middle of [0, 1] since we introduce additional Gaussian noise in a quadratic manner $t(1-t)\sigma^2 \mathbf{z}$. We use the mode distribution defined in Esser et al. (2024) to sample the intermediate time step with higher frequencies.

Compared to the CM loss in 4, we emphasize the consistency along the whole OT trajectory through
COT Pairs in random time steps. In addition, we use auxiliary noise to enhance the robustness of the
OT consistency, with the theoretical guarantee in EOT and Lemma 3.1. The pseudo-code of COT
Flow training pipeline is in Algorithm 1. The detailed algorithm in implementation is in Appendix
A.

Algorithm 1 COT Training

Input: source data distribution μ , neural OT map T_{ϕ} , parameters θ , noise scale σ , learning rate η . repeat

Sample $\mathbf{x} \sim \mu(\mathbf{x})$ and $t_1, t_2 \in [0, 1]$ $\tilde{\mathbf{x}}_{t_i} \leftarrow t_i T_{\phi}(\mathbf{x}) + (1 - t_i)\mathbf{x} + t_i(1 - t_i)\sigma^2 \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad i = 1, 2$ $\mathcal{L}_{\text{COT}}(\theta) \leftarrow d(\mathcal{E}_{\theta}(\mathbf{x}_{t_1}, t_1), \mathcal{E}_{\theta}(\mathbf{x}_{t_2}, t_2))$ $\theta \leftarrow \text{stopgrad}(\theta + \eta \nabla_{\theta} \mathcal{L}_{\text{COT}}(\theta))$ **until** convergence

3.4 COT EDITOR

To further illustrate the flexibility and generalizability of COT Flow, we introduce COT Editor, a zero-shot image editor that possesses various scenarios using a series of modifications of a self-augmentation sampling strategy:

$$\tilde{\mathbf{x}}_{t_k}^{(k)} = t_k \mathbf{x} + (1 - t_k) \tilde{\mathbf{y}}^{(k)} + t_k (1 - t_k) \sigma^2 \mathbf{z}_k \tag{15}$$

$$\tilde{\mathbf{y}}^{(k+1)} = \mathcal{E}_{\theta}\big(\tilde{\mathbf{x}}_{t_k}^{(k)}, t_k\big), \qquad k = 1, 2, \dots$$
(16)

369 where $\tilde{\mathbf{y}}^{(k)}$ is the last estimation of target data. $\tilde{\mathbf{x}}_{t_k}^{(k)}$ is the corresponding self-augmented sample. t_k represents a chosen time step series $0 < t_k < 1$, which is not limited to monotonically increase over 370 371 time. With a well-trained model under Eq.9, we can sample from the source distribution through 372 one-step sampling $\tilde{\mathbf{y}} = \mathcal{E}_{\theta}(\mathbf{x}, 0)$, or optionally adopt a multi-step self-augmentation sampling strat-373 egy in Eq. 15/21, which enables zero-shot editing through the intermediate sampling steps. Both 374 sampling strategies are illustrated in the left panel of Fig.4. With the benefit of unlimited input 375 distribution of COT Flow, COT Editor extends the existing zero-shot image editing scenarios, formulating a dual-channel editing space where both source and target data space \mathcal{X}, \mathcal{Y} are included. 376 We demonstrate its capability by introducing the following scenarios: (1) COT composition, (2) 377 shape-texture coupling, and (3) COT augmentation.



Figure 4: Left: The sampling strategy of our method. Given an input x, we can generate the target data $\tilde{\mathbf{y}}$ with one-step sampling $\tilde{\mathbf{y}} = \mathcal{E}_{\theta}(\mathbf{x}, 1)$, or optionally multi-step sampling using Eq.15/21, where the intermediates $\tilde{\mathbf{x}}_{t_k}$ are the augmentations between the source input **x** and the generated target \hat{y} . **Right:** Three scenarios of the proposed COT Editor, some of which have dual-channel inputs as extensions to the current editing methods. (a). COT composition. Given a target image y with an edited component or mask m, we use the guidance $\mathbf{y}^{(g)} = \mathbf{y} \oplus \mathbf{m}$ as the single input and synthesize the output \tilde{y} by Eq.17. (b). Shape-texture coupling. With a drawn stroke image \hat{x}_1 and a texture image $\hat{\mathbf{x}}_2$, the output $\tilde{\mathbf{y}}$ consists of both features. (c). COT augmentation. Given a series of auto-detected cardiac-cycle edges $\{\hat{\mathbf{x}}^{(a)}\}\$ and a single MRI y, we can generate a cycle of cardiac MRI $\{\tilde{\mathbf{v}}\}\$ with the same movements of $\{\hat{\mathbf{x}}^{(a)}\}\$ and style of \mathbf{v} .

401 For COT composition, given a target image y with an edited component or mask m, we denote the 402 combination as the guidance $\mathbf{y}^{(g)} = \mathbf{y} \oplus \mathbf{m}$ of the COT Editor and perform the following one-step 403 editing to obtain realistic outputs: 404

$$\tilde{\mathbf{y}} = \mathcal{E}_{\theta}(\mathbf{y}^{(g)} + t_a(1 - t_a)\sigma^2 \mathbf{z}, t_a), \quad t_a \in [0, 1]$$
(17)

where t_q denotes the chosen time step of the guidance editing, enabling the trade-off between faithfulness and realism as in Meng et al. (2021). For shape-texture coupling, considering a drawn shape $\hat{\mathbf{x}}_1$ and a texture image $\hat{\mathbf{x}}_2$, we can generate a realistic image using $\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2$ as the augmentation sources:

$$\tilde{\mathbf{y}} = \mathcal{E}_{\theta}(t_c \hat{\mathbf{x}}_1 + (1 - t_c) \hat{\mathbf{x}}_2 + t_c (1 - t_c) \sigma^2 \mathbf{z}, t_c), \quad t_c \in [0, 1]$$

$$\tag{18}$$

For COT augmentation, we provide a medical image synthesis scenario. We denote $\{\hat{\mathbf{x}}^{(a)}\}$ as a series of auto-detected cardiac-cycle edges and augment a fixed input cardiac MRI (cMRI) y by fusing them:

$$\{\tilde{\mathbf{y}}\} \leftarrow \mathcal{E}_{\theta}(t_a \mathbf{y} + (1 - t_a)\{\hat{\mathbf{x}}^{(a)}\} + t_a(1 - t_a)\sigma^2 \mathbf{z}, t_a), \quad t_a \in [0, 1]$$

$$(19)$$

The dual ends of the OT trajectory in COT Flow enrich these additional zero-shot editing applications, where we demonstrate the results in Section 4.2.

4 EXPERIMENTS

423

425

391

392

393

394

395

396

397

398

399

400

405 406

407

408

409

410 411

412

413

414 415 416

417

418

We employ COT Flow in various experiments compared with other popular methods. Section 4.1 shows competitive performances of COT Flow on unpaired I2I translation benchmarks. We compare the generation quality with SDEdit Meng et al. (2021) and CycleGAN Zhu et al. (2017), which are 424 popular diffusion/GAN-based methods. Section 4.2 provides the results of our proposed extended scenarios of zero-shot editing, including COT composition, shape-texture coupling, and COT aug-426 mentation. In Section 4.3, we discuss several key techniques of COT Flow by ablation studies. The implementation details of all the experiments are shown in Appendix A.

427 428 429

- UNPAIRED IMAGE-TO-IMAGE TRANSLATION 4.1
- We perform experiments on handbag \rightarrow shoes (64×64), CelebA male \rightarrow female (64×64), 431 outdoor \rightarrow church (128×128), and edges \rightarrow cardiac MRI (cMRI) (128×128) to implement unpaired

Table 1: FID \downarrow scores of the baseline methods and our proposed COT Flow on handbag \rightarrow shoes (64×64), CelebA male \rightarrow female (64×64), and outdoor \rightarrow church (128×128). Compared to SDEdit with a larger number of function evaluations (NFE), we use one-step sampling in COT Flow as the GAN-based methods.

Method	DiscoGAN	CycleGAN	MUNIT	SDEdit	COT Flow (ours)
NFE	1	1	1	500	1
handbag→shoes male→female outdoor→church	22.42 35.64 75.36	16.00 17.74 46.39	15.76 17.07 31.42	18.91 17.26 28.84	15.01 16.30 26.34

I2I translation. The formulation of these datasets is in Appendix A. With the recommendation of Karras et al. (2022) and Song et al. (2023) to train the diffusion-based methods, we choose the hyper-parameters that are unrelated to our proposed ideas to be in line with these methods, where further details can be found in Appendix A.



Figure 5: Generation comparison between our method (bottom row) and SDEdit (middle row) on CelebA male \rightarrow female (64×64), handbag \rightarrow shoes (64×64), and outdoor \rightarrow church (128×128). We use one-step sampling in our method and set t = 500 of the reverse diffusion process in SDEdit to perform the results.

As shown in Fig.1a, our method provides high-quality generations with one-step or multi-step sampling. In Fig.5, we compare the generation results between SDEdit and the proposed COT Flow, illustrating a more faithful unpaired I2I translation by our method. In Table 1, our method outperforms the other diffusion/GAN-based methods in terms of the FID↓ scores by one-step sampling.

4.2 COT EDITOR SCENARIOS

In section 3.4, we introduce three scenarios of the proposed COT Editor. Fig.1b further present editing results with the trained COT Flow on handbag \rightarrow shoes (64×64), CelebA male \rightarrow female (64×64), and outdoor \rightarrow church (128×128).

4.3 ABLATION STUDIES

We provide reasons of COT Flow's key design by the following ablation studies. In Table 2, we choose alternated contrastive pair formulations, neural OT mapping direction, and sampling strate-gies, which represent the key design of our method. In particular, we (1) train a COT Flow model with only adjacent contrastive pairs $\langle \mathbf{x}_{t_k}, \mathbf{x}_{t_k+1} \rangle$ as is implemented in Song et al. (2023), (2) use the opposite direction of neural OT mapping from target to source $(T'(\mathbf{y}))$ to form the contrastive pairs using $\{\tilde{\mathbf{x}}_t\}_{t \in [0,1]} = \{t\mathbf{y} + (1-t)T'_{\phi}(\mathbf{y}) + t(1-t)\sigma^2 \mathbf{z}\}_{t \in [0,1]}$ instead of Eq.12, and (3) try a different sampling strategy in an ancestral manner, which is commonly adopted in diffusion-based models Ho et al. (2020). As shown in Table 2, COT Flow with the paper's choice outperforms the other alternatives in one-step and multi-step sampling.

Table 2: Ablating COT pairs and sampling strategy on various datasets (evaluated by FID \downarrow scores). "Adjacent pairs" denotes training the COT Flow with only adjacent positive pairs $\langle \mathbf{x}_{t_k}, \mathbf{x}_{t_k+1} \rangle$ as is implemented in Song et al. (2023). "Reverse OT" denotes training a neural OT model $T'(\mathbf{y})$ with opposite direction mapping from target space \mathcal{Y} to source space \mathcal{X} and form the COT pairs. "Ancestral" denotes using a sampling strategy in an ancestral manner in COT Flow.

Method	Adjacent pairs	Reverse NOT	Paper's choice			
NFE	1	1	40 (Ancestral)	40	1	
handbag→shoes male→female outdoor→church	15.24 16.67 26.95	33.49 30.28 38.11	19.97 21.12 26.92	18.33 16.93 26.05	15.01 16.30 26.34	

496 497 498

499

5

CONCLUSION

500 501

502

504

505

506

We presented *COT Flow*, a new method that provides a tangible approach to tackle the generative learning trilemma, achieving fast and high-quality generation and flexible zero-shot image editing. Benefiting from OT reformulation, we achieved competitive sample quality on a great variety of unpaired I2I translation tasks, representing flow between diverse distributions. With the proposed *COT Editor*, We demonstrated flexible zero-shot editing capacities with three scenarios, namely, COT composition, shape-texture coupling, and COT augmentation.

Our method explicitly built the bridge between diffusion/flow-based models and OT by combining consistency models and contrastive learning, opening up new directions for future work. The proposed COT Editor expanded the possibility of zero-shot image editing by the dual-channel editing spaces, enabling new directions for zero-shot editing applications.

511 512

517

527

513 REFERENCES

- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan, 2017.
- ⁵¹⁶ Grigory Bartosh, Dmitry Vetrov, and Christian A. Naesseth. Neural diffusion models, 2024.
- Olivier Bernard, Alain Lalande, Clement Zotti, Frederick Cervenansky, Xin Yang, Pheng-Ann 518 Heng, Irem Cetin, Karim Lekadir, Oscar Camara, Miguel Angel Gonzalez Ballester, Gerard San-519 roma, Sandy Napel, Steffen Petersen, Georgios Tziritas, Elias Grinias, Mahendra Khened, Vargh-520 ese Alex Kollerathu, Ganapathy Krishnamurthi, Marc-Michel Rohe, Xavier Pennec, Maxime 521 Sermesant, Fabian Isensee, Paul Jager, Klaus H. Maier-Hein, Peter M. Full, Ivo Wolf, Sandy 522 Engelhardt, Christian F. Baumgartner, Lisa M. Koch, Jelmer M. Wolterink, Ivana Isgum, Yeong-523 gul Jang, Yoonmi Hong, Jay Patravali, Shubham Jain, Olivier Humbert, and Pierre-Marc Jodoin. 524 Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: 525 Is the problem solved? *IEEE Transactions on Medical Imaging*, 37(11):2514–2525, November 526 2018. ISSN 1558-254X. doi: 10.1109/tmi.2018.2837502.
- Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, March
 2004. ISBN 9780511804441. doi: 10.1017/cbo9780511804441.
- 530 Victor M. Campello, Polyxeni Gkontra, Cristian Izquierdo, Carlos Martin-Isla, Alireza Sojoudi, Pe-531 ter M. Full, Klaus Maier-Hein, Yao Zhang, Zhiqiang He, Jun Ma, Mario Parreno, Alberto Albiol, 532 Fanwei Kong, Shawn C. Shadden, Jorge Corral Acero, Vaanathi Sundaresan, Mina Saber, Mustafa 533 Elattar, Hongwei Li, Bjoern Menze, Firas Khader, Christoph Haarburger, Cian M. Scannell, Mitko 534 Veta, Adam Carscadden, Kumaradevan Punithakumar, Xiao Liu, Sotirios A. Tsaftaris, Xiaoqiong Huang, Xin Yang, Lei Li, Xiahai Zhuang, David Vilades, Martin L. Descalzo, Andrea Guala, 536 Lucia La Mura, Matthias G. Friedrich, Ria Garg, Julie Lebel, Filipe Henriques, Mahir Karakas, Ersin Cavus, Steffen E. Petersen, Sergio Escalera, Santi Segui, Jose F. Rodriguez-Palomares, and Karim Lekadir. Multi-centre, multi-vendor and multi-disease cardiac segmentation: The m&ms 538 challenge. IEEE Transactions on Medical Imaging, 40(12):3543-3554, December 2021. ISSN 1558-254X. doi: 10.1109/tmi.2021.3090082.

- 540
 541
 542
 Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations, 2018.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
 contrastive learning of visual representations, 2020.
- Xinlei Chen and Kaiming He. Exploring simple siamese representation learning, 2020.
- 547 Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis, 2021.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis, 2024.
- Jiaojiao Fan, Shu Liu, Shaojun Ma, Haomin Zhou, and Yongxin Chen. Neural monge map estimation and its applications, 2021.
- Jiaojiao Fan, Shu Liu, Shaojun Ma, Yongxin Chen, and Hao-Min Zhou. Scalable computation of monge maps with general costs. In *ICLR Workshop on Deep Generative Models for Highly Structured Data*, 2022. URL https://openreview.net/forum?id=rEnGR3VdDW5.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
 Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014.
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena
 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Ghesh laghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own
 latent: A new approach to self-supervised learning, 2020.
- Nikita Gushchin, Alexander Kolesov, Alexander Korotin, Dmitry Vetrov, and Evgeny Burnaev. Entropic neural optimal transport via diffusion processes, 2022.
- 568 Nikita Gushchin, Sergei Kholkin, Evgeny Burnaev, and Alexander Korotin. Light and optimal
 569 schrödinger bridge matching, 2024.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning, 2019.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020.
- Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 9902–9915. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/v162/janner22a.html.
- L. Kantorovitch. On the translocation of masses. *Management Science*, 5(1):1–4, October 1958.
 ISSN 1526-5501. doi: 10.1287/mnsc.5.1.1.
- Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks, 2018.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyz ing and improving the image quality of stylegan, 2019.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models, 2022.
- Amirhossein Kazerouni, Ehsan Khodapanah Aghdam, Moein Heidari, Reza Azad, Mohsen Fayyaz, Ilker Hacihaliloglu, and Dorit Merhof. Diffusion models in medical imaging: A comprehensive survey. *Medical Image Analysis*, 88:102846, August 2023. ISSN 1361-8415. doi: 10.1016/j. media.2023.102846.

612

624

- Patrick Kidger, James Morrill, James Foster, and Terry Lyons. Neural controlled differential equations for irregular time series, 2020.
- 597 Diederik P. Kingma and Ruiqi Gao. Understanding diffusion objectives as the elbo with simple data augmentation, 2023.
- Alexander Korotin, Daniil Selikhanovych, and Evgeny Burnaev. Neural optimal transport. 2022a.
 doi: 10.48550/ARXIV.2201.12220.
- Alexander Korotin, Daniil Selikhanovych, and Evgeny Burnaev. Kernel neural optimal transport,
 2022b.
- Yanyu Li, Huan Wang, Qing Jin, Ju Hu, Pavlo Chemerys, Yun Fu, Yanzhi Wang, Sergey Tulyakov, and Jian Ren. Snapfusion: Text-to-image diffusion model on mobile devices within two seconds, 2023.
- Shanchuan Lin, Bingchen Liu, Jiashi Li, and Xiao Yang. Common diffusion noise schedules and sample steps are flawed, 2023.
- Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching
 for generative modeling, 2022.
- Kingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow, 2022.
- Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang,
 Hanchi Sun, Jianfeng Gao, Lifang He, and Lichao Sun. Sora: A review on background, technology, limitations, and opportunities of large vision models, 2024.
- ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁰
 ⁶¹⁰
 ⁶¹⁰
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁰
 ⁶¹⁰
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁶
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁰
 ⁶¹⁰
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
- 622 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast
 623 ode solver for diffusion probabilistic model sampling in around 10 steps, 2022a.
- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models, 2022b.
- Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Syn thesizing high-resolution images with few-step inference, 2023.
- 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, 2021.
- Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik P. Kingma, Stefano Ermon, Jonathan Ho, and
 Tim Salimans. On distillation of guided diffusion models, 2022.
- Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models, 2021.
- Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with
 text-guided diffusion models, 2021.
- Taesung Park, Alexei A. Efros, Richard Zhang, and Jun-Yan Zhu. Contrastive learning for unpaired image-to-image translation, 2020.
- Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion, 2022.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical textconditional image generation with clip latents, 2022.
- 647 R. Tyrrell Rockafellar. Integral functionals, normal integrands and measurable selections, pp. 157– 207. Springer Berlin Heidelberg, 1976. ISBN 9783540380757. doi: 10.1007/bfb0079944.

- 648 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-649 resolution image synthesis with latent diffusion models, 2021. 650
- Litu Rout, Alexander Korotin, and Evgeny Burnaev. Generative modeling with optimal transport 651 maps, 2021. 652
- 653 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kam-654 yar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffu-655 sion models with deep language understanding, 2022. 656
- 657 Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models, 2022. 658
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2020a. 659
- 660 Yang Song and Prafulla Dhariwal. Improved techniques for training consistency models, 2023. 661
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribu-662 tion. Curran Associates Inc., Red Hook, NY, USA, 2019. 663
- 664 Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben 665 Poole. Score-based generative modeling through stochastic differential equations, 2020b.
- 666 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models, 2023. 667
- Alexander Tong, Kilian Fatras, Nikolay Malkin, Guillaume Huguet, Yanlei Zhang, Jarrid Rector-668 Brooks, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models 669 with minibatch optimal transport, 2023. 670
- 671 Cédric Villani. Optimal Transport. Springer Berlin Heidelberg, 2009. ISBN 9783540710509. doi: 672 10.1007/978-3-540-71050-9.
- Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the generative learning trilemma with 674 denoising diffusion gans, 2021. 675
- 676 Yanwu Xu, Yang Zhao, Zhisheng Xiao, and Tingbo Hou. Ufogen: You forward once large scale text-to-image generation via diffusion gans, 2023. 677
 - Aron Yu and Kristen Grauman. Fine-grained visual comparisons with local learning. In 2014 IEEE Conference on Computer Vision and Pattern Recognition. IEEE, June 2014. doi: 10.1109/cvpr. 2014.32.
- Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: 682 Construction of a large-scale image dataset using deep learning with humans in the loop, 2015. 683
- Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable 685 effectiveness of deep features as a perceptual metric, 2018.
 - Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. Learning deep features for scene recognition using places database. In Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 1, NIPS'14, pp. 487–495, Cambridge, MA, USA, 2014. MIT Press.
 - Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In Computer Vision (ICCV), 2017 IEEE International Conference on, 2017.
- 693 694

673

678

679

680

681

684

686

687

688

689

690

691

692

IMPLEMENTATION DETAILS А

697 In this section, we provide the implementation details of our method. Section A.1 provides the detailed training algorithm of our method in the implementation. Section A.2 introduces the used datasets and the construction of the unpaired I2I translation tasks. Section A.3 discusses the details 699 of the chosen hyper-parameters of our method. Section A.4 provides the training details and the 700 computational complexity of our method. Section A.5 introduces the alternative combinations of 701 our method in Section 4.3 to perform the ablation studies.

702 A.1 DETAILED ALGORITHM

In the implementation, we uniformly discretize the sampled time steps t_1, t_2 in Eq.13 with the number of the discrete time steps N. We use LPIPS Zhang et al. (2018) distance as the distance metric $d(\cdot, \cdot)$. The detailed training algorithm of our method is as follows:

Algorithm 2 COT Training

Input: source data distribution μ , neural OT map T_{ϕ} , parameters θ , noise scale σ , learning rate η , distance metric $d(\cdot, \cdot)$, and number of discretization N. **repeat**

708 709

710

711

712

Sample $\mathbf{x} \sim \mu(\mathbf{x})$ and $n_1, n_2 \in \mathcal{U}[0, N-1]$ $n_1 < n_2$ $\tilde{\mathbf{x}}_{t_i} \leftarrow \frac{n_i}{N-1} T_{\phi}(\mathbf{x}) + (1 - \frac{n_i}{N-1}) \mathbf{x} + \frac{n_i}{N-1} (1 - \frac{n_i}{N-1}) \sigma^2 \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad i = 1, 2$ $\theta_1, \theta_2 \leftarrow \theta$ $\mathcal{L}_{\text{COT}}(\theta_1, \theta_2) \leftarrow d(\mathcal{E}_{\theta_1}(\mathbf{x}_{t_1}, t_1), \mathcal{E}_{\theta_2}(\mathbf{x}_{t_2}, t_2))$ $\theta \leftarrow \theta + \eta \nabla_{\theta_1} \mathcal{L}_{\text{COT}}(\theta_1, \theta_2)$ until convergence

717 718 719

720

732

A.2 DATASETS

We use the following publicly available datasets as the sources **x** or targets **y**: Amazon handbags and shoes Yu & Grauman (2014) to perform handbag \rightarrow shoes (64×64); CelebA faces Liu et al. (2015) to perform male \rightarrow female (64×64); outdoor images of MIT places database Zhou et al. (2014) and LSUN church dataset Yu et al. (2015) to perform outdoor \rightarrow church (128×128); auto-detected edges on the M&Ms dataset Campello et al. (2021) and ACDC dataset Bernard et al. (2018) to perform edges \rightarrow cMRI (128×128). All the coupled datasets are unpaired and randomly sampled during training.

For the proposed zero-shot image editing scenarios, we utilize the trained models on the aforementioned tasks, where no additional dataset is needed.

730731 A.3 HYPER-PARAMETERS

733 Despite the differences between our method and diffusion-based models, we use the recommenda-734 tions in Karras et al. (2022) for the common hyper-parameters such as learning rate and number of 735 discrete time steps (N = 40). We use the noise scale $\sigma = 1$ for all the tasks.

736 737 A.4 TRAINING DETAILS

For the network structure and the training details of the neural OT models, we follow the recommendations of Korotin et al. (2022a). The neural OT models converge in 1-2 days on a single NVidia A40 GPU (48GB). The batch size during training is 64 for all the tasks.

For the encoder models \mathcal{E}_{θ} , the network structure uses the recommendations in Song et al. (2023), and the models converge in 3-4 days on 4×NVidia A40 GPUs (48GB). The batch size during training is 128 for all the tasks.

745 746 A.5 Ablation Study Details

We provide three alternatives as a comparison to ablate our training and/or sampling choices.

749In particular, we first train the models using adjacent positive pairs $\langle \mathbf{x}_{t_k}, \mathbf{x}_{t_k+1} \rangle$ instead of the COT750Pairs $\langle \mathbf{x}_{t_1}, \mathbf{x}_{t_2} \rangle$ provided by Eq.13. This alternative evaluates the importance of the chosen COT Pair751formulation and emphasizes the connection between consistency models and contrastive learning.

752 Secondly, we choose an opposite direction to train the neural OT models in each task. For example, 753 in the handbag \rightarrow shoes task, instead of training a neural OT model $T(\mathbf{x})$ from the handbag dataset 754 to the shoes dataset, we train a reverse neural OT model $T'(\mathbf{y})$ from shoes data \mathbf{y} to handbag data 755 \mathbf{x} . This alternative evaluates the paper's choice of the neural OT model's direction and verifies the 757 formulation of COT Pairs. Finally, we provide an optional sampling strategy to prove the effectiveness of our self-augmentation sampling strategy in COT Editor. After training the models, we implement an ancestral-like sampling strategy to generate the results:

$$\tilde{\mathbf{x}}_{t_k}^{(k)} = \frac{t_k}{t_{k-1}} \tilde{\mathbf{x}}_{t_{k-1}}^{(k-1)} + (1 - \frac{t_k}{t_{k-1}}) \tilde{\mathbf{y}}^{(k)} + t_k (1 - t_k) \sigma^2 \mathbf{z}_k$$
(20)

764

765 766

767

768

769 770

771

772 773

778

783

784

785

790 791

794

796 797

798

799 800 801

802

803 804

805

806

760

$$\tilde{\mathbf{y}}^{(k+1)} = \mathcal{E}_{\theta} \left(\tilde{\mathbf{x}}_{t_k}^{(k)}, t_k \right), \qquad k = 1, 2, \dots, \qquad \tilde{\mathbf{x}}_{t_0}^{(0)} = \mathbf{x}$$
(21)

B PROOF OF THEOREM

Proposition 3.1. Let π^* be the OT plan between $\mu(\mathbf{x})$ and $\nu(\mathbf{y})$. Let the OT map T^* recover π^* . The augmentation defined by Eq.12 using T^* samples the same probability as the dynamic extension of the EOT plan π^*_{λ} with $\lambda = 2\sigma^2$.

Proof. According to Gushchin et al. (2024), the augmentation between x and $T^*(x)$ using Eq.12 samples a probability distribution:

$$p_t(\mathbf{x}_t | \mathbf{x}, T^*(\mathbf{x})) = \mathcal{N}(\mathbf{x}_t | tT^*(\mathbf{x}) + (1-t)\mathbf{x}, t(1-t)\sigma \mathbf{I})$$
(22)

which is the time marginal of a Brownian Bridge $\mathbf{w}_{|\mathbf{x},T^*(\mathbf{x})}^{\sigma}$ (Appendix C). Using the probability distribution in 22, the Schrödinger Bridge S^* (Appendix D) between $\mu(\mathbf{x})$ and $\nu(\mathbf{y})$ can be estimated by:

$$\tilde{S}^* = \int_{\mathbb{R}\times\mathbb{R}} \mathbf{w}^{\sigma}_{|\mathbf{x},\mathbf{y}} d\tilde{\pi}^*(\mathbf{x}, T_{\phi}(\mathbf{x}))$$
(23)

Which is the dynamic extension of the entropy-regularized OT problem with optimum $\pi^*_{2\sigma^2}$ according to Tong et al. (2023), where the joint marginal distribution π^{S^*} of S^* at times 0,1 is the EOT plan $\pi^*_{2\sigma^2}$ in 7, i.e., $\pi^{S^*} = \pi^*_{2\sigma^2}$.

C BROWNIAN BRIDGE

Suppose we have a data point **x** with time intermediates \mathbf{x}_t in the processes. Given a Wiener process \mathbf{w}_t^{σ} defined by $d\mathbf{w}_t^{\sigma} = \sqrt{\sigma} d\mathbf{w}_t$ with volatility $\sigma > 0$, $t \in [0, T]$, and standard Wiener process \mathbf{w}_t . A Brownian Bridge is the conditional probability distribution $\mathbf{w}_{|\mathbf{x}_0,\mathbf{x}_T}^{\sigma}$ subject to the condition that the start and end point of the process is $\mathbf{x}_0, \mathbf{x}_T$. The probability distribution is:

$$\mathcal{N}(\mathbf{x}_t | t\mathbf{x}_T + (T-t)\mathbf{x}_0, t(T-t)\sigma \mathbf{I})$$
(24)

Intuitively, the Brownian Bridge is pinned to the values $\mathbf{x}_0, \mathbf{x}_T$ at t = 0 and t = T, and the most uncertainty lies in the middle of the bridge.

D SCHRÖDINGER BRIDGE

Given two probability distribution $\mu(\mathbf{x})$ and $\nu(\mathbf{y})$, consider the Wiener process \mathbf{w}_t^{σ} with volatility $\sigma > 0$ starts at $\mu(\mathbf{x})$ at t = 0, the Schrödinger Bridge between $\mu(\mathbf{x}), \nu(\mathbf{y})$ is:

$$S^* = \min_{S \in \mathcal{F}(\mu,\nu)} \operatorname{KL}(S \parallel \mathbf{w}_t^{\sigma})$$
(25)

where S is a stochastic process and $\mathcal{F}(\mu, \nu)$ is a set of stochastic processes with the start of $\mu(\mathbf{x})$ at t = 0 and end of $\nu(\mathbf{y})$ at t = T.

E DIFFUSION MODELS

⁸⁰⁷ Diffusion models learn to denoise the data in different noise scales and generate samples from noise ⁸⁰⁸ via an iterative denoising process. The original data distribution $\mu(\mathbf{x})$ is diffused with a stochastic ⁸⁰⁹ differential equation (SDE):

$$d\mathbf{x}_t = \mathbf{g}(\mathbf{x}_t, t)dt + \sigma(t)d\mathbf{w}_t \tag{26}$$

where $t \in [0,T]$, T > 0 is a constant, g is the drift term and $d\mathbf{w}_t$ represents a standard Wiener process. We denote the intermediate distribution of \mathbf{x}_t as $p_t(\mathbf{x})$. Then the SDE process has a dual ODE whose solution trajectories at time t are distributed according to $p_t(\mathbf{x})$:

$$d\mathbf{x}_t = \left[\mathbf{g}(\mathbf{x}_t, t) - \frac{1}{2}\sigma(t)^2 \nabla \log p_t(\mathbf{x}_t)\right] dt$$
(27)

where $\nabla \log p_t(\mathbf{x}_t)$ denotes the score function of $p_t(\mathbf{x})$, which is estimated by a neural network $\mathbf{s}_{\phi}(\mathbf{x}_t, t) \approx \nabla \log p_t(\mathbf{x}_t)$. We then sample \mathbf{x}_0 from the estimated probability flow ODE:

$$\frac{d\mathbf{x}_t}{dt} = -t\mathbf{s}_\phi(\mathbf{x}_t, t) \tag{28}$$

where we initialize $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$ and solve 28 backward in time to obtain the generation $\hat{\mathbf{x}}_0$ via various ODE solvers such as Euler and Heun solvers.

F LIMITATIONS

COT Flow explicitly builds the bridge between optimal transport and diffusion/flow-based models. However, our method requires a two-step training pipeline, including the neural OT model $T(\mathbf{x})$ and the encoder model \mathcal{E} , which may influence the training and deploying stability. A promising future direction is to design an end-to-end method with OT formulation explicitly.

BROADER IMPACTS G

COT Flow and other generative models pose a risk of synthesizing inappropriate content such as deep-fake images, violence, or privacy-related offensiveness.

Η ADDITIONAL EXPERIMENTS

We compared the zero-shot image editing ability between our model and SDEdit on different datasets. We took 600k iteration steps with a batch size of 256 on 4×NVidia A40 GPUs to train our model, and we followed the recommended training hyper-parameters in Meng et al. (2021) to ensure the convergence of the baseline SDEdit model. The results in Fig.6 show that our model (COT Editor) outperforms SDEdit in both image editing quality and fidelity on all datasets.



Figure 6: Zero-shot image editing comparison between our method (COT Editor) and SDEdit on CelebA male \rightarrow female (64×64), handbag \rightarrow shoes (64×64), and outdoor \rightarrow church (128×128). We use one-step and multi-step sampling in our method and set t = 300, 400, 500, 600 of the reverse diffusion process in SDEdit to perform the editing results.