# **DELAY FLOW MATCHING**

Paper under double-blind review

Anonymous authors

000

001 002 003

004

006

008 009

010

011

012

013

014

015

016

017

018

019

021

024

026027028

029

031

033

034

037

040

041

042

043

044

046 047

048

051

052

# **ABSTRACT**

Flow matching (FM) based on Ordinary Differential Equations (ODEs) has achieved significant success in generative tasks. However, it faces several inherent limitations, including an inability to model trajectory intersections, capture delay dynamics, and handle transfer between heterogeneous distributions. These limitations often result in a significant mismatch between the modeled transfer process and real-world phenomena, particularly when key coupling or inherent structural information between distributions must be preserved. To address these issues, we propose Delay Flow Matching (DFM), a new FM framework based on Delay Differential Equations (DDEs). Theoretically, we show that DFM possesses universal approximation capability for continuous transfer maps. By incorporating delay terms into the vector field, DFM enables trajectory intersections and better captures delay dynamics. Moreover, by designing appropriate initial functions, DFM ensures accurate transfer between heterogeneous distributions. Consequently, our framework preserves essential coupling relationships and achieves more flexible distribution transfer strategies. We validate DFM's effectiveness across synthetic datasets, single-cell data, and image-generation tasks.

# 1 Introduction

Generative modeling is a key and rapidly advancing field in machine learning, focusing on learning transformations between different distributions. It underpins a wide range of applications across many tasks in diverse fields, including image generation (Nichol et al., 2021), molecule design (Sanchez-Lengeling & Aspuru-Guzik, 2018), and single-cell trajectory inference (Saelens et al., 2019). Traditional approaches, such as Variational Autoencoders (Kingma, 2013; Rezende et al., 2014), Generative Adversarial Networks (Goodfellow et al., 2014), and Normalizing Flows (Dinh et al., 2014; 2022; Papamakarios et al., 2021), have achieved notable success, yet they still face challenges like instability and computational inefficiency.

Recently, diffusion models (Sohl-Dickstein et al., 2015; Song & Ermon, 2019; Ho et al., 2020; Song et al., 2020a; Yang et al., 2023) have garnered considerable attention for their ability to model complex data distributions via a diffusion process. These models incrementally add noise to data in a forward process and then learn to reverse this process in order to regenerate data samples, which can be understood as learning stochastic dynamics that interpolate between the prior distribution and the target data distribution (Song et al., 2020b; 2021). Additionally, Flow Matching (FM), a simulation-free method for training continuous normalizing flows (Chen et al., 2018), is proposed to model the transformation between distributions as the flow maps of Neural ODEs (Chen et al., 2018; Grathwohl et al., 2018). FM achieves efficient training by directly regressing on an explicitly constructed conditional vector field. Concurrently, the stochastic interpolant (Albergo & Vanden-Eijnden, 2022) and rectified flow (Liu et al., 2022) are introduced, both of which employ flow maps for matching distributions, though from distinct conceptual frameworks.

As discussed above, most existing continuous generative models for distribution transfer rely on ODEs or SDEs. However, ODE-based models have limited representational capacity, restricting their ability to capture a broad range of distribution transfer strategies (Dupont et al., 2019). In contrast, Delay Differential Equations (DDEs) have been widely adopted to model various real-world systems, including neural dynamics (Campbell, 2007), electro-optical systems (Chembo Kouomou et al., 2005), population dynamics (Lotka, 1925), and many biological network motifs (Glass et al., 2021), where delayed feedback mechanisms naturally give rise to DDE-based formulations. To address the limitations of ODEs in dynamic modeling, Neural DDEs (Zhu et al., 2020) and their

056

060

061 062

063

064

065

067

069

071

073

074 075

076

077

078

079 080

081 082

083 084

087

090

091

092

094

096

098

099

100

101

102

103

104

105

106

107

variants (Ji & Orosz, 2024; Zhu et al., 2022) are proposed, which explicitly incorporate the effects of historical states and thus exhibit a significantly enhanced representational capacity. Despite these advancements, however, there is currently no approach that utilizes the probability flows of DDEs for distribution transport.

**Contributions**. We present Delay Flow Matching (DFM), a generalized framework that enables more flexible and precise distribution transport strategies using DDEs. The key contributions of this study are summarized as follows:

- 1. Development of DFM: We introduce DFM, a novel generative model framework based on Neural DDEs, which overcomes the inherent limitations of ODE-based models by incorporating delay terms in the vector field and designing appropriate initial functions. DFM can model a broader range of transport strategies, including those with trajectory intersections, and achieves more precise transport between heterogeneous distributions. It also naturally adapts to the probability flow generated by delay dynamical systems.
- 2. Theoretical insights: We rigorously prove that DFM can universally approximate any continuous transport map between source and target distributions. In contrast, ODE-based models cannot represent certain simple transport maps, such as those involving trajectory intersections, and cannot achieve exact transport between heterogeneous distributions.
- 3. **Integration with advanced techniques**: DFM can be seamlessly integrated with existing methods, such as keypoint-guided optimal transport, enabling more effective alignment with known coupling information during transport.
- 4. **Empirical validation**: We validate DFM's effectiveness on both synthetic and real-world datasets. DFM accurately recovers underlying delay dynamics from snapshot data, significantly outperforms existing methods in single-cell trajectory inference, and surpasses ODE-based FM in image generation tasks.

# 2 PRELIMINARIES

# 2.1 DDES AND PROBABILITY FLOWS

DDEs are widely used to model systems in which the time evolution depends not only on the current state but also on past states, such as those governed by delayed feedback. For time-dependent DDEs with a single delay term, the general formulation is given by:

$$\frac{\mathrm{d}\boldsymbol{x}(t)}{\mathrm{d}t} = \boldsymbol{u}[t, \boldsymbol{x}(t), \boldsymbol{x}(t-\tau)], t \in [0, T], \quad (1)$$
$$\boldsymbol{x}(h) = \boldsymbol{\psi}(h), h \in [-\tau, 0],$$

where  $\boldsymbol{u}[t,\boldsymbol{x}(t),\boldsymbol{x}(t-\tau)]:[0,T]\times\mathbb{R}^d\times\mathbb{R}^d\to\mathbb{R}^d$  is a smooth vector field which is abbreviated as  $\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_\tau)$  in the following,  $\boldsymbol{\psi}(h)$  represents the continuous initial function. We further denote  $\boldsymbol{\psi}(h;\boldsymbol{x}_0)$  as the initial function that takes the value  $\boldsymbol{x}_0$  at time t=0, i.e.  $\boldsymbol{\psi}(0;\boldsymbol{x}_0)=\boldsymbol{x}_0$ .

For a given initial distribution  $x_0 \sim p_0(x_0)$  and the corresponding initial functions  $\psi(t;x_0)$ , the above DDE (1) induces the associated probability flows  $p(x,t|\psi): \mathbb{R}^d \times [0,T] \to \mathbb{R}^+$ , satisfying the delay Fokker-Planck equation (refer to Appendix B) (Guillouzic et al., 1999; Frank, 2005):

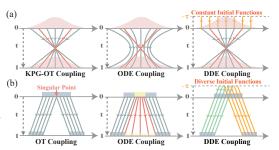


Figure 1: Comparison of ODE coupling and DDE coupling. (a) A Gaussian distribution mapped to itself via KPG-OT, with the two solid red lines and arrows representing the source-target keypoint pairs that must be satisfied, results in the mapping  $x \to -x$ . This induces trajectory intersections during the transport process, preventing ODEs from maintaining the desired mapping. In contrast, a DDE  $\dot{x} = -2x(t-1)$  with constant initial functions achieves exact coupling. (b) An uniform distribution  $\mathcal{U}(-2,2)$  is transported via OT to distribution  $\frac{1}{2}\mathcal{U}(-3,-1)+\frac{1}{2}\mathcal{U}(1,3)$ . ODEs preserve the connectivity of sets, causing unavoidable transport a small neighborhood of 0 to (-1,1). DDEs achieve precise transport by assigning different initial functions to [-2,0] and (0,2].

$$\frac{\partial p(\boldsymbol{x}, t|\boldsymbol{\psi})}{\partial t} = -\nabla \cdot \left\{ \int d\boldsymbol{x}_{\tau} [\boldsymbol{u}(t, \boldsymbol{x}, \boldsymbol{x}_{\tau}) \cdot p(\boldsymbol{x}, t; \boldsymbol{x}_{\tau}, t - \tau | \boldsymbol{\psi})] \right\}, \tag{2}$$

where  $p(x, t|\psi)$  is the probability density given the initial function,  $p(x, t; x_{\tau}, t - \tau | \psi)$  is the joint probability density of being at x at time t and at  $x_{\tau}$  at time  $t-\tau$ . The specific relationships and distinctions between the probability flows of ODEs and DDEs are elaborated in Appendix A.

#### 2.2 OPTIMAL TRANSPORT

Optimal transport (OT) provides a mathematical framework for transforming one probability distribution to another in the most cost-efficient manner (Villani et al., 2009; Santambrogio, 2015; Chen, 2016; Peyré et al., 2019). It seeks the optimal strategy for reallocating mass while minimizing the transportation cost. Formally, given two separable metric spaces  $\mathcal{X}$  and  $\mathcal{Y}$  with probability measures  $\mu$  on  $\mathcal{X}$  and  $\nu$  on  $\mathcal{Y}$ , the objective is to determine a transport plan  $\pi^*$  that minimizes the total transport cost, as defined in the Kantorovich problem (Kantorovich, 1942):

$$C(\mu, \nu) = \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(\mathbf{x}, \mathbf{y}) d\pi(\mathbf{x}, \mathbf{y}), \ \Pi(\mu, \nu) = \{\pi : P_{\#}^{\mathbf{x}}(\pi) = \mu, P_{\#}^{\mathbf{y}}(\pi) = \nu\}, \quad (3)$$

where c(x, y) denotes the transportation cost of moving a unit mass from x to y,  $P_\#^x(\pi)$  and  $P_\#^y(\pi)$  denote the marginal distributions of  $\pi$  with respect to x and y. In most cases, the transport cost between two points is defined as the squared Euclidean distance. The minimum total transport cost in this scenario corresponds to the squared 2-Wasserstein distance between the two distributions:

$$W_2(\mu,\nu)^2 = \inf_{\pi \in \Pi(\mu,\nu)} \int \|\boldsymbol{x} - \boldsymbol{y}\|^2 d\pi(\boldsymbol{x},\boldsymbol{y}). \tag{4}$$

## 2.3 KEYPOINT-GUIDED OPTIMAL TRANSPORT

Traditional OT, which focuses solely on minimizing transport costs, often neglects other crucial constraints, resulting in suboptimal matching strategies. Thus, it is important to use a small set of well-matched source-target keypoint pairs,  $\mathcal{K} = \{(\boldsymbol{x}_k, \boldsymbol{y}_k)\}_{k=1}^K$ , to semi-supervise the transport strategy, especially when inherent structures and features of data need to be preserved.

To ensure correct transport with these keypoints, Keypoint-guided Optimal Transport (KPG-OT) is proposed (Gu et al., 2022; 2023), which guarantees that the relationships between each data point and the keypoints are preserved during transport. Formally, the objective can be expressed as:

$$\inf_{\tilde{\pi}\in\tilde{\Pi}(\mu,\nu)}\int_{\mathcal{X}\times\mathcal{Y}}g(\boldsymbol{x},\boldsymbol{y})\mathrm{d}(\boldsymbol{w}\circ\tilde{\pi})(\boldsymbol{x},\boldsymbol{y}),\ \tilde{\Pi}(\mu,\nu)=\{\tilde{\pi}:P_{\#}^{\boldsymbol{x}}(\boldsymbol{w}\circ\tilde{\pi})=\mu,P_{\#}^{\boldsymbol{y}}(\boldsymbol{w}\circ\tilde{\pi})=\mu\},\ (5)$$

where  $w \circ \tilde{\pi} = w(\boldsymbol{x}, \boldsymbol{y}) \tilde{\pi}(\boldsymbol{x}, \boldsymbol{y})$  is the keypoint-masked transport plan, w denotes the mask function. For a pair of keypoints  $(\boldsymbol{x}_k, \boldsymbol{y}_k) \in \mathcal{K}$ , the mask function is defined as:  $w(\boldsymbol{x}_k, \boldsymbol{y}_k) = 1, w(\boldsymbol{x}_k, \boldsymbol{y}) = 0$  for  $\boldsymbol{y} \neq \boldsymbol{y}_k, w(\boldsymbol{x}, \boldsymbol{y}_k) = 0$  for  $\boldsymbol{x} \neq \boldsymbol{x}_k$ , and  $w(\boldsymbol{x}, \boldsymbol{y}) = 1$  if neither  $\boldsymbol{x}$  nor  $\boldsymbol{y}$  matches any keypoint. Obviously, the mask function defined above exactly preserves the matching of the keypoints in  $\mathcal{K}$ .

The cost function in Eq. (5) is defined as  $g(x, y) = d[R^s(x), R^t(y)]$ , where d denotes the Jensen-Shannon divergence.  $R^s(x) \in (0, 1)^K$  (resp.  $R^t(y) \in (0, 1)^K$ ) captures the relationship between x (resp. y) and all source points (resp. target points) in K, with its k-th dimension representing the relation score between x (resp. y) and the k-th source keypoint  $x_k$  (resp. target keypoint  $y_k$ ):

$$R_k^{\rm s}(\boldsymbol{x}) = \frac{e^{-c(\boldsymbol{x}, \boldsymbol{x}_k)/\tau}}{\sum_{i=1}^K e^{-c(\boldsymbol{x}, \boldsymbol{x}_i)/\tau}}, R_k^{\rm t}(\boldsymbol{y}) = \frac{e^{-c(\boldsymbol{y}, \boldsymbol{y}_k)/\tau}}{\sum_{i=1}^K e^{-c(\boldsymbol{y}, \boldsymbol{y}_i)/\tau}},$$
(6)

where  $\tau$  is the temperature hyperparameter. Note that g quantifies the similarity distance between the two relationship vectors outlined above. Therefore, minimizing the cost function in Eq. (5) effectively promotes the preservation of the relationship between each point and the keypoints.

# 2.4 ODE-BASED FLOW MATCHING

Flow Matching (FM) (Lipman et al., 2022) is a simulation-free training method for CNF, which trains the parameterized vector field  $v(t, x; \theta)$  by regression on the target vector field u(t, x) that generates the probability paths p(x, t) with marginals  $p_0 = q_0$  and  $p_1 = q_1$ . The training objective is

$$\mathcal{L}_{FM}(\theta) = \mathbb{E}_{t,p(\boldsymbol{x},t)} \| \boldsymbol{v}(t,\boldsymbol{x};\theta) - \boldsymbol{u}(t,\boldsymbol{x}) \|^{2},$$
(7)

where  $t \sim \mathcal{U}(0,1)$  follows the uniform distribution.

However, the explicit forms of u(t, x) and p(x, t) are intractable to compute. To address this, Conditional FM (CFM) (Tong et al., 2023a; Pooladian et al., 2023) introduces a latent variable z to construct the target probability path as a mixture of conditional probability paths  $p(x, t) = \mathbb{E}_{q(z)}[p(x, t|z)]$  and modifies the training objective to:

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t,q(\boldsymbol{z}),p(\boldsymbol{x},t|\boldsymbol{z})} \|\boldsymbol{v}(t,\boldsymbol{x};\theta) - \boldsymbol{u}(t,\boldsymbol{x}|\boldsymbol{z})\|^{2},$$
(8)

where u(t, x|z) denotes the conditional vector field generating the conditional probability path p(x, t|z). Then, it can be shown that the gradient w.r.t  $\theta$  of the CFM objective (8) is the same as that of the FM objective (7). Typically, we set the latent condition  $z := (x_0, x_1)$ , with q(z) being the coupling between distributions  $q_0(x_0)$  and  $q_1(x_1)$ , based on a coupling strategy such as OT or KPG-OT. The conditional probability path is then modeled as a Gaussian flow:

$$p(\boldsymbol{x}, t|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}|t\boldsymbol{x}_1 + (1-t)\boldsymbol{x}_0, \sigma^2), \tag{9}$$

with the corresponding conditional vector field:

$$\boldsymbol{u}(t, \boldsymbol{x}|\boldsymbol{z}) = \boldsymbol{x}_1 - \boldsymbol{x}_0. \tag{10}$$

# 3 LIMITATIONS OF ODE-BASED FLOW MATCHING

Current continuous-time generative models utilize the parameterized vector field of ODEs to facilitate distribution transport. However, the inherent constraints of ODEs limit their ability to model specific transport strategies, leading to significant discrepancies between the modeled and true transport processes. In this section, we theoretically demonstrate the limitations of the ODE-based FM frameworks.

**Proposition 3.1.** (Restriction on trajectory intersections) Suppose that the target transport strategy, guided by certain key points or specific coupling constraints, inevitably leads to trajectory intersections during the transport process, the flow map corresponding to the ODE-based FM cannot precisely preserve the transport strategy.

Remark 3.2. As demonstrated by Liu et al. (2022), when trajectory intersections occur in the transport strategy, FM tends to learn a *rectified flow*, with the targets of the trajectories being "rewired" at the intersection points. As a result, the given transport strategy cannot be accurately preserved. A simple example is shown in Fig. 1 (a).

**Proposition 3.3.** (Heterogeneity in distributions) Assume that the source (resp., target) distribution  $q_0$  (resp.,  $q_1$ ) is supported on M (resp., N) disjoint compact sets  $\{U_i^0\}_{i=1}^M$  (resp.,  $\{U_j^1\}_{j=1}^N$ ), where each  $U_i^0 \in \mathbb{R}^d$  (resp.,  $U_j^1 \in \mathbb{R}^d$ ) is a path-connected set with non-zero measure and  $q_0(\mathbf{x}_0) > 0$  (resp.,  $q_1(\mathbf{x}_1) > 0$ ) for all  $\mathbf{x}_0 \in \bigcup_{i=1}^M U_i^0$  (resp.,  $\mathbf{x}_1 \in \bigcup_{j=1}^N U_j^1$ ). If  $M \neq N$ , for any transport map  $T: \bigcup_{i=1}^M U_i^0 \to \bigcup_{j=1}^N U_j^1$  such that  $T_\# q_0 = q_1$ , we have that: (1). Neural ODEs with Lipschitz continuous vector field cannot exactly represent T; (2). If the flow of a Neural ODE is equal to the transport map T almost everywhere, then the vector field is not Lipschitz continuous.

Remark 3.4. The conclusion emphasizes that accurate transport using an ODE-based FM requires regularity assumptions on distributions, which are often violated in real-world scenarios. For example, in single-cell dynamics, cells differentiate from one type into multiple types over time, leading to heterogeneous distributions before and after differentiation. A simple example is shown in Fig. 1 (b).

Moreover, when snapshot data are generated by a delay dynamical system, ODE-based FM fails to recover the true vector field with the delay term, resulting in inaccurate distribution transfer, interpolation, and extrapolation predictions.

# 4 DELAY FLOW MATCHING

To overcome the limitations of ODE-based FM, we introduce the Delay FM (DFM) framework, a new class of generative models based on Neural DDEs for distribution transport. Since the vector field incorporates delay terms, DFM allows trajectory intersections, enabling more accurate transport for tasks requiring keypoint-guided strategies. Additionally, by designing appropriate initial functions, DFM addresses singularities caused by distributional heterogeneity. A comparison between DDE-based and ODE-based FM under various scenarios is shown in Table 1.

Table 1: Comparison of ODE-based and DDE-based FM. DFM can handle not only KPG-OT coupling and trajectory intersections but also distribution heterogeneity by employing diverse initial functions. Furthermore, DFM can accurately model and recover the dynamics of delay dynamical systems.

Methods	KPG-OT	Intersection	Heterogeneity	Delay Dynamics
FM CFM	X	X	X	X
CFM	^	^	^	^
DFM(C)	✓	$\checkmark$	X	$\checkmark$
DFM(D)	✓	✓	✓	✓

#### 4.1 FORMULATION

Consider a target probability flow  $p(\boldsymbol{x},t)$ , generated by a vector field with a single delay term  $\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau})$  and initial functions  $\boldsymbol{\psi} \sim q^{\circ}(\boldsymbol{\psi})$ . The vector field and initial functions naturally define a joint probability flow at time t and  $t-\tau$ , denoted as  $p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau)$ , which satisfies

$$\int p(\boldsymbol{x}, t; \boldsymbol{x}_{\tau}, t - \tau) d\boldsymbol{x}_{\tau} = p(\boldsymbol{x}, t).$$
(11)

Note that p(x, t) and  $p(x, t; x_{\tau}, t - \tau)$  can be modeled as a mixture of conditional distributions  $p(x, t|\psi)$  and  $p(x, t; x_{\tau}, t - \tau|\psi)$ , respectively, as follow:

$$p(\boldsymbol{x},t) = \mathbb{E}_{q^{\circ}(\boldsymbol{\psi})} p(\boldsymbol{x},t|\boldsymbol{\psi}), \quad p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau) = \mathbb{E}_{q^{\circ}(\boldsymbol{\psi})} p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi}), \tag{12}$$

where  $p(\boldsymbol{x},t|\boldsymbol{\psi})$  and  $p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi})$  represent the probability flow generated by  $\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau})$  and the initial function  $\boldsymbol{\psi}$ , and they satisfy the delay Fokker-Planck equation (2). Based on the joint probability flow, we aim to learn a parameterized vector field  $\boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta})$ , which can generate the target probability flow, by minimizing the following regression objective against the target vector field:

$$\mathcal{L}_{\text{DFM}}(\theta) = \mathbb{E}_{t,q^{\circ}(\boldsymbol{\psi}),p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi})} ||\boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\theta) - \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau})||^{2}.$$
(13)

In most cases, both  $p(x,t;x_{\tau},t-\tau|\psi)$  and  $u(t,x,x_{\tau})$  are computationally intractable. To address this, inspired by CFM, we introduce a latent variable z and further decompose the target probability path and joint probability path conditioned on the initial function into mixtures of simple probability paths and simple joint probability paths conditioned on both  $\psi$  and z, respectively, as follow:

$$p(\boldsymbol{x}, t|\boldsymbol{\psi}) = \mathbb{E}_{q(\boldsymbol{z})}p(\boldsymbol{x}, t|\boldsymbol{z}, \boldsymbol{\psi}), \quad p(\boldsymbol{x}, t; \boldsymbol{x}_{\tau}, t - \tau|\boldsymbol{\psi}) = \mathbb{E}_{q(\boldsymbol{z})}p(\boldsymbol{x}, t; \boldsymbol{x}_{\tau}, t - \tau|\boldsymbol{z}, \boldsymbol{\psi}), \tag{14}$$

where  $p(x, t|z, \psi) = \int p(x, t; x_{\tau}, t - \tau|z, \psi) dx_{\tau}$ . We can define the marginal vector field by marginalizing over the conditional vector field  $u(t, x, x_{\tau}|z)$  as follow:

$$\boldsymbol{u}(t, \boldsymbol{x}, \boldsymbol{x}_{\tau}) = \mathbb{E}_{q(\boldsymbol{z})} \frac{\boldsymbol{u}(t, \boldsymbol{x}, \boldsymbol{x}_{\tau} | \boldsymbol{z}) p(\boldsymbol{x}, t; \boldsymbol{x}_{\tau}, t - \tau | \boldsymbol{z}, \boldsymbol{\psi})}{p(\boldsymbol{x}, t; \boldsymbol{x}_{\tau}, t - \tau | \boldsymbol{\psi})},$$
(15)

where  $u(t, x, x_{\tau}|z)$  represents the conditional vector field that generates  $p(x, t|z, \psi)$  and  $p(x, t; x_{\tau}, t - \tau|z, \psi)$ . Under this setup, we can derive the following result.

**Proposition 4.1.** The marginal vector field  $\mathbf{u}(t, \mathbf{x}, \mathbf{x}_{\tau})$  given in Eq. (15), together with the selected initial function  $\boldsymbol{\psi}$ , generates the probability path  $p(\mathbf{x}, t | \boldsymbol{\psi})$  and the joint probability path  $p(\mathbf{x}, t; \mathbf{x}_{\tau}, t - \tau | \boldsymbol{\psi})$  in Eq. (14).

Then, we can train the parameterized vector field by minimizing the following Delay Conditional FM (DCFM) objective:

$$\mathcal{L}_{\text{DCFM}}(\theta) = \mathbb{E}_{t,q(\boldsymbol{z}),q^{\circ}(\boldsymbol{\psi}),p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi})}||\boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\theta) - \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z})||^{2},$$
(16)

because the following proposition holds.

**Proposition 4.2.** Given that  $p(\mathbf{x}, t; \mathbf{x}_{\tau}, t - \tau | \psi) > 0$  for all  $\mathbf{x}, \mathbf{x}_{\tau} \in \mathbb{R}^d$  and  $t \in [0, 1]$ , up to a constant independent of  $\boldsymbol{\theta}$ ,  $\mathcal{L}_{DCFM}(\boldsymbol{\theta})$  and  $\mathcal{L}_{DFM}(\boldsymbol{\theta})$  are equal, which further implies that  $\nabla_{\boldsymbol{\theta}} \mathcal{L}_{DCFM}(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \mathcal{L}_{DFM}(\boldsymbol{\theta})$ .

## SELECTION OF THE LATENT VARIABLE z

270

271 272

273

274

275

276

277

278

279

280

281

282

283 284

285

287

288

289

290

291

292 293

295

296 297

298

299 300

301

303 304

305

306

307

308

309 310

311 312

313

314

315

316

317

318

319

320

321

322

323

#### 4.2.1 TRANSPORT BETWEEN TWO MEASURES

As outlined in Section 2.4, ODE-based FM typically selects latent variables as two endpoints  $z := (x_0, x_1)$  jointly sampled from the source and target distributions. In contrast, DFM requires the construction of a conditional joint probability density path for x(t) and  $x(t-\tau)$  based on the latent variable. Hence, we define the latent variable z as a entire path  $\gamma(t; x_0, x_1)$  connecting  $x_0$  and  $x_1$ .

Formally, we define the probability distribution of the latent variable as  $q[\gamma(t; x_0, x_1)] :=$  $\pi(x_0, x_1) \mathcal{P}(\gamma; x_0, x_1)$ , where  $\pi(x_0, x_1)$  represents the OT or KPG-OT coupling between the source and target distributions, and  $\mathcal{P}(\gamma; x_0, x_1)$  denotes the path measure pinned at  $x_0$  and  $x_1$ . In practice, we can simply construct the path measure as a Dirac Delta distribution at a given path  $\gamma^*$  connecting  $\boldsymbol{x}_0$  and  $\boldsymbol{x}_1$ :

$$\mathcal{P}(\gamma; \boldsymbol{x}_0, \boldsymbol{x}_1) = \delta(\gamma - \gamma^*), \ \gamma_0^* = \boldsymbol{x}_0, \gamma_1^* = \boldsymbol{x}_1, \tag{17}$$

 $\mathcal{P}(\gamma; \boldsymbol{x}_0, \boldsymbol{x}_1) = \delta(\gamma - \gamma^*), \ \gamma_0^* = \boldsymbol{x}_0, \gamma_1^* = \boldsymbol{x}_1,$  where  $\gamma^*$  is constructed using a specific interpolation method. For instance, it can be taken as the linear interpolation  $\gamma_t^* = (1-t)x_0 + tx_1$ , or alternatively, as a geodesic interpolation on the data manifold, based on appropriate manifold learning techniques.

After sampling a path  $\gamma(t; x_0, x_1) \sim q(z)$ , the conditional joint probability density naturally degenerates into the following form:

$$p(\boldsymbol{x}, t; \boldsymbol{x}_{\tau}, t - \tau | \boldsymbol{\gamma}) = \delta[\boldsymbol{x} - \boldsymbol{\gamma}(t)] \delta[\boldsymbol{x}_{\tau} - \boldsymbol{\gamma}(t - \tau)], \tag{18}$$

which is exactly generated by the conditional vector field:

$$u(t, \boldsymbol{x}, \boldsymbol{x}_{\tau} | \boldsymbol{\gamma}(t; \boldsymbol{x}_0, \boldsymbol{x}_1)) = \frac{\partial \boldsymbol{\gamma}(t; \boldsymbol{x}_0, \boldsymbol{x}_1)}{\partial t}.$$
 (19)

#### 4.2.2 TRANSPORT BETWEEN MULTIPLE PROBABILITY MEASURES OVER TIME

For tasks with multiple target probability distributions over time  $\{q_{t_i}\}_{i=0}^J$ , where  $t_0=0, t_J=T$ , we define the latent variable  $z := \gamma(t; \{x_{t_j}\}_{j=0}^J)$  as a trajectory passing through  $x_{t_j}$  at time  $t_j$ , sampled

$$q(\gamma(t; \{\boldsymbol{x}_{t_j}\}_{j=0}^J)) := \mathcal{P}(\gamma; \{\boldsymbol{x}_{t_j}\}_{j=0}^J) \prod_{i=0}^{J-1} \pi(\boldsymbol{x}_{t_i}, \boldsymbol{x}_{t_{i+1}}),$$
(20)

where  $\{x_{t_j}\}_{j=0}^J \sim \prod_{i=0}^{J-1} \pi(x_{t_i}, x_{t_{i+1}})$ , with  $\pi(x_{t_i}, x_{t_{i+1}})$  representing the OT or KPG-OT coupling between adjacent probability distributions at time  $t_i$  and  $t_{i+1}$ , and  $\mathcal{P}(\gamma; \{x_{t_i}\}_{i=0}^J)$  denotes the path measure pinned at  $\{x_{t_j}\}_{j=0}^{J}$ . Similarly, we construct the path measure as:

$$\mathcal{P}(\gamma; \{\boldsymbol{x}_{t_j}\}_{j=0}^J) = \delta(\gamma - \gamma^*), \ \gamma_{t_j}^* = \boldsymbol{x}_{t_j},$$
(21)

where  $\gamma^*$  represents a path passing through  $x_{t_i}$  at time  $t_i$ , which can be constructed using the cubic spline (CSpline) interpolation method. With this path as the latent variable, the resulting conditional joint probability path and conditional vector field are identical to those in Eq. (18) and Eq. (19), respectively.

## SELECTION OF THE INITIAL FUNCTION

# 4.3.1 DFM WITH CONSTANT INITIAL FUNCTIONS

In general, the initial function can be simply chosen as a constant (DFM(C)), i.e.  $q^{\circ}(\psi) = \delta(\psi - \psi^*)$ , where  $\psi^*(t; x_0) \equiv x_0$  for  $t \in [-\tau, 0]$  and  $x_0 \sim q_0$ . Under this setting, we can theoretically prove that DFM can approximate any continuous transport strategy to arbitrary precision.

**Proposition 4.3.** (Universal approximating capability of DFM). For any given continuous transport map  $F: \mathbb{R}^d \to \mathbb{R}^d$ , which push-forward the source distribution  $q_0$  to the target distribution  $q_1$ , i.e.  $F_{\#}q_0 = q_1$ , if there exists a neural network that can approximate V(x) = F(x) - x, then we can construct a vector field with a single delay term, where the corresponding flow map, under the constant initial function condition, can approximately push-forward q<sub>0</sub> to q<sub>1</sub>, while preserving the target transport strategy  $x \to F(x)$ .

This indicates the exceptional representational capacity of DFM, enabling the modeling of a wider range of transport processes than ODE-based FM.

## 4.3.2 DFM WITH DIVERSE INITIAL FUNCTIONS

In Section 3, we rigorously prove that ODE-based FM cannot effectively handle tasks with heterogeneous source and target distributions. Here, we demonstrate that DFM addresses this by designing diverse initial functions (DFM(D)).

Specifically, we employ clustering methods, such as Gaussian Mixture Model, DBSCAN, to partition the source dataset  $X_0 \sim q_0$  (resp. target dataset  $X_1 \sim q_1$ ) into M (resp. N) mutually exclusive subsets, denoted as  $X_0^{(1)},...,X_0^{(M)}$  (resp.  $X_1^{(1)},...,X_1^{(N)}$ ). We can assign a normalized mass to each subset  $X_0^{(m)}$  (resp.  $X_1^{(n)}$ ) as  $\rho_0^{(m)} = |X_0^{(m)}|/|X_0|$  (resp.,  $\rho_1^{(n)} = |X_1^{(n)}|/|X_1|$ ). If the endpoints  $(x_0,x_1)$  of a sampled trajectory  $\gamma$  are drawn from  $X_0^{(m)}$  and  $X_1^{(n)}$ , we assign it an initial function  $\psi_{mn}^*$  which has a constant time derivative  $C_{mn}$ , i.e.:

$$\frac{\mathrm{d}\psi_{mn}^{*}(t; \boldsymbol{x}_{0})}{\mathrm{d}t} = C_{mn}, \ \psi_{mn}^{*}(0; \boldsymbol{x}_{0}) = \boldsymbol{x}_{0}, \tag{22}$$

where  $t \in [-\tau, 0]$ . In this case,  $q^{\circ}(\psi)$  is a discrete distribution which satisfies:

$$\sum_{n=1}^{N} q^{\circ}(\psi_{mn}^{*}) = \rho_{0}^{(m)}, \quad \sum_{m=1}^{M} q^{\circ}(\psi_{mn}^{*}) = \rho_{1}^{(n)}. \tag{23}$$

 $p(\boldsymbol{x}_0,0|\boldsymbol{\psi}=\psi_{mn}^*)=q_0(\boldsymbol{x}_0|\boldsymbol{\psi}=\psi_{mn}^*)$  represents empirical data distribution available as data points in  $\boldsymbol{X}_0^{(m)}$  whose corresponding transport target is in  $\boldsymbol{X}_1^{(n)}$ . By coupling the source and target data through OT or KPG-OT, we can obtain the transport target for any initial point. This enables the construction of  $q^\circ(\psi)$  as described above, from which we can sample to obtain the corresponding initial function.

In summary, we design distinct initial functions for different subsets of the source and target data, guiding the vector field from different initial subsets to corresponding target subsets, thereby effectively handling distributional heterogeneity.

## 4.4 GENERATION PROCESS BASED ON NEURAL DDES

After training, we sample data from  $q_0$  and initial functions from  $q^{\circ}(\psi)$ , then generate the target data by solving the forward pass of the trained vector field with a single delay term using a piecewise ODE solver, as in Neural DDEs (Zhu et al., 2020).

# 5 EXPERIMENTS

We demonstrate the advantages of DFM over ODE-based FM across various tasks, including reconstructing delay dynamics from snapshots, inferring differentiation trajectories from single-cell RNA sequencing (scRNA-seq) data, and image generation. The experimental details and additional results can be found in Appendix C. Sensitivity analysis of the time delay  $\tau$  is discussed in Appendix D.

# 5.1 SYNTHETIC DATASETS OF DELAY DYNAMICAL SYSTEMS

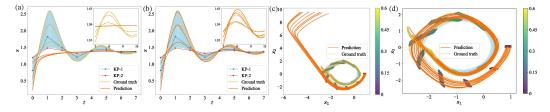


Figure 2: Comparison of predicted trajectories of CFM and DFM trained on snapshots of delay dynamical systems. (a-b) Interpolation and extrapolation (insets) results on the biological autoregulation motif dataset using KP-CFM (a) and KP-DFM(C) (b), both with two keypoints. (c-d) Interpolation results on the spiral DDE using OT-CFM (c) and OT-DFM(C) (d).

Biological autoregulation motif. DDEs provide a simplified framework for modeling biological network motifs (Glass et al., 2021). We consider recovering the delay dynamics from snapshot data generated by the autoregulation model:  $\dot{x}(t) = \frac{\eta}{1+x^n(t-\tau)}$  with  $\tau=1, n=2, \eta=5$ , which produces damped oscillations (Fig. 2 (a,b)). Using 1,000 initial values sampled from  $\mathcal{U}(0.2, 1.2)$  and constant initial functions, we generate trajectories and collect snapshots at  $t=0,1,\ldots,7$  for training. During training, KPG-OT (KP-) matches minibatches at adjacent time steps using two known keyoints. After training, forward integration from t=0 shows that KP-CFM fails to recover the dynamics, while KP-DFM(C) with  $\tau=1$  accurately captures and extrapolates the oscillatory behavior (insets in Fig. 2 (a,b)). Results with other delays and error metrics are provided in Appendix C.1 and D.

**Spiral DDE**. We next consider a 2-d DDE (Zhu et al., 2020):  $\dot{x}(t) = A \tanh(x(t) + x(t-\tau))$  with  $\tau$ =0.5 and  $A \in \mathbb{R}^{2\times 2}$ , which produces crossing spiral trajectories (Fig. 2 (c,d)). Snapshots are taken every 0.05 in  $t \in [0,0.6]$  and coupled using minibatch-OT (OT-) between adjacent steps. After training, OT-DFM(C) with  $\tau$ =0.5 successfully reproduces the dynamics, while OT-CFM fails around the crossing region. DFM also generalizes well across a range of  $\tau$  values (Appendix C.2, D).

## 5.2 Trajectory inference of single-cell

We investigate the inference of differentiation trajectories from real scRNA-seq data, where heterogeneity increases as cells transition from one type to multiple types during development, complicating modeling with ODE-based approaches. To evaluate model performance, we train on all time points except one intermediate time point, and then sample from the initial distribution, performing forward integration to predict distributions at each time point. Predictions are validated through unsupervised leave-one-out validation (L) by comparing to the true distribution at the held-out intermediate time point, and supervised final-time validation (F) by comparing to the true distribution at the final time point. Trajectory inference accuracy is assessed using the 2-Wasserstein distance ( $W_2$ ) and Maximum Mean Discrepancy with a Gaussian kernel (MMD(G)).

Table 2: Trajectory inference results on the scRNA-seq dataset in mouse hematopoiesis and the single-cell qPCR iPSC dataset with bifurcation. All results are averages of 10 runs.

Methods		Mouse hen	natopoies	is		qPCR	iPSC	
Wicthous	$ \overline{W_2(L)} $	MMD(L)	$W_2(F)$	MMD(F)	$W_2(L)$	MMD(L)	$W_2(F)$	MMD(F)
TIGON	0.519	0.563	0.264	0.155	0.733	0.791	0.695	0.405
MIOFlow	0.514	0.629	0.220	0.056	0.770	1.039	0.345	0.155
OT-CFM	0.378	0.357	0.192	0.047	0.579	0.492	0.226	0.030
OT-DFM(C)	0.379	0.384	0.136	0.021	0.553	0.447	0.234	0.041
OT-DFM(D)	0.372	0.341	0.095	0.010	0.532	0.399	0.213	0.027

Mouse hematopoiesis dataset. We evaluate DFM on scRNA-seq data from mouse hematopoiesis (Weinreb et al., 2020), focusing on cells differentiating into neutrophils (Neu) and monocytes (Mo) at Day 2, 4, and 6. Models are trained on Day 2 and 6 data; after training, samples from Day 2 are integrated forward to predict distributions at Days 4 and 6. Geodesic interpolation is used to construct trajectories for both CFM and DFM (Appendix C.3). Due to the heterogeneity between pre- and post-differentiation distributions, trajectories inferred by ODE-based models (TIGON (Sha et al., 2024), MIOFlow (Huguet et al., 2022), OT-CFM) partially deviate from the data manifold, falling into regions between Neu and Mo fates (Fig. 3 (b-e)). In contrast, OT-DFM(D), using distinct initial functions for Neu and Mo (Appendix C.3), preserves alignment with the data manifold (Fig. 3 (f)) and achieves more accurate distribution predictions (Table 2).

**qPCR iPSC dataset**. We further apply DFM to model the bifurcation process of iPSCs in cardiomyocytes (Bargaje et al., 2017). Data from Day 2, 3, 4, and 5 are selected, with bifurcation observed from Day 3, where progenitor cells differentiate into mesodermal (M) and endodermal (En) fates. Models are trained on all time points except Day 3. After training, samples from Day 2 are forward-integrated to infer trajectories and predict distributions at Day 3 and 5. Due to distributional heterogeneity at the bifurcation point, ODE-based methods and OT-DFM(C) produce trajectories that partially misalign with the data manifold (Fig. 3 (h–k)). In contrast, OT-DFM(D), by assigning separate initial functions for each fate (Appendix C.4), preserves trajectory alignment with the corresponding manifolds and yields more accurate predictions at both Day 3 and 5 compared to other methods (Table 2).

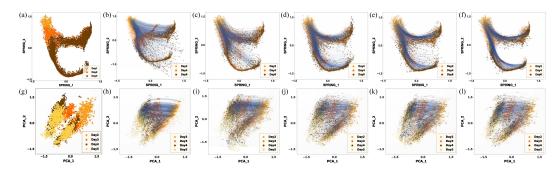


Figure 3: Comparisons between predicted trajectories of MIOFlow (b, h), TIGON (c, i), OT-CFM (d, j), OT-DFM(C) (e, k), and OT-DFM(D) (f, l) on scRNA-seq dataset in mouse hematopoiesis (a-f) and qPCR iPSC dataset (g-l). (a) and (g) illustrate the true data distribution.

# 5.3 IMAGE GENERATION

**MNIST dataset**. We design a *Semi-paired Image-to-Image Translation* task on the MNIST dataset, where the source domain consists of original images, and the target domain includes their negative images (the normalized pixel values transformed via  $X \mapsto 1-X$ ). To provide partial supervision, 10% of the training data are paired with their negative counterparts as keypoints. During training, minibatches are inde-

Table 3: Comparison of FID between KP-CFM ( $\tau=0$ ) and KP-DFM(C) with different time delay on the MNIST dataset.

$\overline{\tau}$	0 (CFM)	0.125	0.250	0.500	1.000
FID	45.020	28.497	11.747	12.653	12.031

pendently sampled from the source and target distributions, and coupled via KPG-OT (KP-). The task aims to transform source images into their negatives, where, under linear interpolation, all transport paths intersect at an image with uniform pixel values of 0.5. While KP-CFM struggles with this transformation, KP-DFM(C) effectively handles it. As shown in Table 3, across varying time delays  $\tau$ , KP-DFM(C) consistently achieves lower FID scores and better distribution alignment than KP-CFM.

CIFAR-10 dataset. We further evaluate DFM on the CIFAR-10 dataset under a more general generation setup (Zhu & Lin, 2024), where the source distribution is a two-component Gaussian mixture rather than a standard Gaussian. DFM employs trainable initial functions with distinct constant time derivatives (Appendix C.6), enabling generation from different mixture components to specific image classes. We compare CFM and DFM(D) under two coupling strate-

Table 4: Comparison of FID between CFM and DFM on the CIFAR-10 dataset.

NFE	10	20	30	40	Adap.
I-CFM OT-CFM	108.291 78.165	$94.629 \\ 27.512$	$91.404 \\ 16.409$	$90.254 \\ 12.026$	$88.306 \\ 6.162$
I-DFM(D) OT-DFM(D)	54.064 54.222	18.248 18.598	<b>11.429</b> 11.894	9.008 9.287	<b>4.980</b> 5.191

gies: independent coupling (I-) and OT coupling (OT-). As shown in Table 4, I-CFM struggles with mode heterogeneity, while I-DFM(D) generates higher-quality images. With OT coupling, OT-DFM(D) consistently outperforms OT-CFM, especially when the number of function evaluations (NFE) is small. Additional results are provided in Appendix C.6.

#### 6 Conclusion

We introduce DFM, a novel continuous-time generative framework based on Neural DDEs. Through theoretical analysis, we highlight the limitations of ODE-based generative models, particularly their inability to capture certain transport strategies and preserve critical coupling information. In contrast, DFM offers universal approximation for arbitrary continuous transport strategies, addressing these shortcomings effectively. DFM also overcomes the challenge of transport between heterogeneous distributions by incorporating task-specific initial functions. Furthermore, it is naturally suited for modeling delay dynamical systems, a feature beyond the capability of ODEs. Extensive experiments on both synthetic and real-world datasets demonstrate that DFM achieves significantly more precise and versatile distribution transport strategies compared to FM.

# REFERENCES

- Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. *Arxiv Preprint Arxiv:2209.15571*, 2022.
- Rhishikesh Bargaje, Kalliopi Trachana, Martin N Shelton, Christopher S McGinnis, Joseph X Zhou, Cora Chadick, Savannah Cook, Christopher Cavanaugh, Sui Huang, and Leroy Hood. Cell population structure prior to bifurcation predicts efficiency of directed differentiation in human induced pluripotent cells. *Proceedings of the National Academy of Sciences*, 114(9):2271–2276, 2017.
- Sue Ann Campbell. Time delays in neural systems. In *Handbook of brain connectivity*, pp. 65–90. Springer, 2007.
  - Yanne Chembo Kouomou, Pere Colet, Laurent Larger, and Nicolas Gastaud. Chaotic breathers in delayed electro-optical systems. *Physical Review Letters*, 95(20):203903, 2005.
  - Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in Neural Information Processing Systems*, 31, 2018.
  - Yongxin Chen. Modeling and control of collective dynamics: From Schrödinger bridges to optimal mass transport. PhD thesis, University of Minnesota, 2016.
  - Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.
  - Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. In *International Conference on Learning Representations*, 2022.
  - Emilien Dupont, Arnaud Doucet, and Yee Whye Teh. Augmented neural odes. *Advances in Neural Information Processing Systems*, 32, 2019.
  - TD Frank. Delay fokker-planck equations, novikov's theorem, and boltzmann distributions as small delay approximations. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, 72(1): 011112, 2005.
  - David S Glass, Xiaofan Jin, and Ingmar H Riedel-Kruse. Nonlinear delay differential equations and their application to modeling biological network motifs. *Nature Communications*, 12(1):1788, 2021.
  - Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27, 2014.
  - Will Grathwohl, Ricky TQ Chen, Jesse Bettencourt, Ilya Sutskever, and David Duvenaud. Ffjord: Free-form continuous dynamics for scalable reversible generative models. *Arxiv Preprint Arxiv:1810.01367*, 2018.
  - Xiang Gu, Yucheng Yang, Wei Zeng, Jian Sun, and Zongben Xu. Keypoint-guided optimal transport with applications in heterogeneous domain adaptation. *Advances in Neural Information Processing Systems*, 35:14972–14985, 2022.
  - Xiang Gu, Liwei Yang, Jian Sun, and Zongben Xu. Optimal transport-guided conditional score-based diffusion model. *Advances in Neural Information Processing Systems*, 36:36540–36552, 2023.
- Steve Guillouzic, Ivan L'Heureux, and André Longtin. Small delay approximation of stochastic delay differential equations. *Physical Review E*, 59(4):3970, 1999.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
  - Guillaume Huguet, Daniel Sumner Magruder, Alexander Tong, Oluwadamilola Fasina, Manik Kuchroo, Guy Wolf, and Smita Krishnaswamy. Manifold interpolating optimal-transport flows for trajectory inference. *Advances in neural information processing systems*, 35:29705–29718, 2022.

- Xunbi A. Ji and Gábor Orosz. Trainable delays in time delay neural networks for learning delayed dynamics. *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–11, 2024. doi: 10.1109/TNNLS.2024.3379020.
- Leonid V Kantorovich. On the translocation of masses. In *Dokl. Akad. Nauk. USSR (NS)*, volume 37, pp. 199–201, 1942.
  - Kacper Kapuśniak, Peter Potaptchik, Teodora Reu, Leo Zhang, Alexander Tong, Michael Bronstein, Avishek Joey Bose, and Francesco Di Giovanni. Metric flow matching for smooth interpolations on the data manifold. *arXiv preprint arXiv:2405.14780*, 2024.
  - Diederik P Kingma. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114, 2013.
  - Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
  - Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *Arxiv Preprint Arxiv*:2209.03003, 2022.
  - AJ Lotka. Elements of physical biology. Williams and Wilkins, 1925.
  - 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. *arXiv preprint arXiv:2112.10741*, 2021.
  - George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *Journal of Machine Learning Research*, 22(57):1–64, 2021.
  - Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data science. *Foundations and Trends*® *in Machine Learning*, 11(5-6):355–607, 2019.
  - Aram-Alexandre Pooladian, Heli Ben-Hamu, Carles Domingo-Enrich, Brandon Amos, Yaron Lipman, and Ricky TQ Chen. Multisample flow matching: Straightening flows with minibatch couplings. *arXiv preprint arXiv:2304.14772*, 2023.
  - Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, pp. 1278–1286. PMLR, 2014.
  - Wouter Saelens, Robrecht Cannoodt, Helena Todorov, and Yvan Saeys. A comparison of single-cell trajectory inference methods. *Nature Biotechnology*, 37(5):547–554, 2019.
  - Benjamin Sanchez-Lengeling and Alán Aspuru-Guzik. Inverse molecular design using machine learning: Generative models for matter engineering. *Science*, 361(6400):360–365, 2018.
  - Filippo Santambrogio. Optimal transport for applied mathematicians. *Birkäuser*, *NY*, 55(58-63):94, 2015.
  - Yutong Sha, Yuchi Qiu, Peijie Zhou, and Qing Nie. Reconstructing growth and dynamic trajectories from single-cell transcriptomics data. *Nature Machine Intelligence*, 6(1):25–39, 2024.
  - Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pp. 2256–2265. PMLR, 2015.
  - Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *Arxiv Preprint Arxiv:2010.02502*, 2020a.
  - Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019.
  - Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv* preprint *arXiv*:2011.13456, 2020b.

- Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. *Advances in neural information processing systems*, 34:1415–1428, 2021.
- Alexander Tong, Kilian Fatras, Nikolay Malkin, Guillaume Huguet, Yanlei Zhang, Jarrid Rector-Brooks, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models with minibatch optimal transport. *arXiv preprint arXiv:2302.00482*, 2023a.
- Alexander Tong, Nikolay Malkin, Kilian Fatras, Lazar Atanackovic, Yanlei Zhang, Guillaume Huguet, Guy Wolf, and Yoshua Bengio. Simulation-free schrödinger bridges via score and flow matching. *arXiv* preprint arXiv:2307.03672, 2023b.
- Cédric Villani et al. Optimal transport: old and new, volume 338. Springer, 2009.
- Caleb Weinreb, Alejo Rodriguez-Fraticelli, Fernando D Camargo, and Allon M Klein. Lineage tracing on transcriptional landscapes links state to fate during differentiation. *Science*, 367(6479): eaaw3381, 2020.
- Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 56(4):1–39, 2023.
- Qunxi Zhu and Wei Lin. Switched flow matching: Eliminating singularities via switching odes. In *Forty-first International Conference on Machine Learning*, 2024.
- Qunxi Zhu, Yao Guo, and Wei Lin. Neural delay differential equations. In *International Conference on Learning Representations*, 2020.
- Qunxi Zhu, Yifei Shen, Dongsheng Li, and Wei Lin. Neural piecewise-constant delay differential equations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 9242–9250, 2022.

The structure of the appendix is as follows:

- Appendix A provides the connection between the probability flows of ODEs and DDEs.
- Appendix B provides the formal proofs for the theoretical results presented in the main text.
- Appendix C includes the experimential setup details and additional experimental results.
- Appendix D provides the sensitivity analysis of time delay parameter τ on both synthetic and real-world datasets.

# A CONNECTION BETWEEN THE PROBABILITY FLOWS OF ODES AND DDES

Based on the vector field  $u[t, x(t), x(t-\tau)]$  in Eq. (1), we can define the conditional average drift (CAD) without delay terms as follows (Guillouzic et al., 1999):

$$\overline{\boldsymbol{u}}(t,\boldsymbol{x}|\boldsymbol{\psi}) = \int d\boldsymbol{x}_{\tau}[\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau})p(\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{x},t;\boldsymbol{\psi})], \tag{24}$$

where  $p(x_{\tau}, t - \tau | x, t; \psi)$  denotes the conditional probability given x(t) = x. Note that  $p(x, t; x_{\tau}, t - \tau | \psi) = p(x_{\tau}, t - \tau | x, t; \psi) p(x, t | \psi)$ , so Eq. (2) is equivalent to:

$$\frac{\partial p(\boldsymbol{x}, t|\boldsymbol{\psi})}{\partial t} = -\nabla \cdot \{ \overline{\boldsymbol{u}}(t, \boldsymbol{x}|\boldsymbol{\psi}) p(\boldsymbol{x}, t|\boldsymbol{\psi}) \}, \tag{25}$$

which is precisely the continuity equation satisfied by the probability flows of ODEs.

Remark A.1. This implies that a non-delayed vector field can be constructed to match the probability flow of the delayed vector field in Eq. (1). However, the integration in Eq. (24) eliminates the coupling information between the states at times t and  $t-\tau$ . Consequently, while the ODE preserves the marginal probability flow, it fails to maintain the coupling relationship during distribution transfer, as illustrated by a simple example in Fig. 1 (a).

# B PROOFS OF THEORETICAL RESULTS

#### B.1 FOKKER-PLANCK EQUATION OF DDES

**Theorem B.1** (Fokker-Planck equation of DDEs, (Guillouzic et al., 1999; Frank, 2005)). *Consider a time-dependent DDE with a single delay term:* 

$$\frac{\mathrm{d}\boldsymbol{x}(t)}{\mathrm{d}t} = \boldsymbol{u}[t, \boldsymbol{x}(t), \boldsymbol{x}(t-\tau)], \quad t \in [0, T],$$

$$\boldsymbol{x}(h) = \boldsymbol{\psi}(h), \quad h \in [-\tau, 0],$$
(26)

where  $\mathbf{u}[t, \mathbf{x}(t), \mathbf{x}(t-\tau)] : [0, T] \times \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}^d$  is a smooth vector field which is abbreviated as  $\mathbf{u}(t, \mathbf{x}, \mathbf{x}_{\tau})$  in the following,  $\psi(h)$  represents the continuous initial function. The associated probability flows  $p(\mathbf{x}, t|\psi) : \mathbb{R}^d \times [0, T] \to \mathbb{R}^+$  satisfy the delay Fokker-Planck equation:

$$\frac{\partial p(\boldsymbol{x}, t | \boldsymbol{\psi})}{\partial t} = -\nabla \cdot \left\{ \int d\boldsymbol{x}_{\tau} [\boldsymbol{u}(t, \boldsymbol{x}, \boldsymbol{x}_{\tau}) p(\boldsymbol{x}, t; \boldsymbol{x}_{\tau}, t - \tau | \boldsymbol{\psi})] \right\}, \tag{27}$$

where  $p(x, t|\psi)$  is the probability density at time t given the initial function, while  $p(x, t; x_{\tau}, t - \tau | \psi)$  is the joint probability density representing the likelihood of the system being at x at time t and at  $x_{\tau}$  at time  $t - \tau$ .

*Proof.* The proof follows a similar approach as presented in Guillouzic et al. (1999). Without loss of generality, we assume that d=1 and the state variable  $x\in [a,b]$  in Eq. (26). Consider an arbitrary  $C^2$  function F(x) defined on the interval [a,b], i.e.  $F\in C^2([a,b])$ , which satisfies the following conditions:

$$\lim_{x \to a} F(x) = \lim_{x \to b} F(x) = 0,$$
(28)

$$\lim_{x \to a} \frac{\mathrm{d}}{\mathrm{d}x} F(x) = \lim_{x \to b} \frac{\mathrm{d}}{\mathrm{d}x} F(x) = 0. \tag{29}$$

Then, by applying the Taylor expansion, we obtain:

$$dF[x(t)] = F[x(t) + dx(t)] - F[x(t)] = \left\{ u[t, x(t), x(t - \tau)] \frac{d}{dx} F[x(t)] \right\} dt.$$
 (30)

The ensemble average (average over realizations) of  $\mathrm{d}F[x(t)]$  can be written as

$$\left\langle \frac{\mathrm{d}}{\mathrm{d}t} F[x(t)] \right\rangle = \left\langle u[t, x(t), x(t-\tau)] \frac{\mathrm{d}}{\mathrm{d}x} F[x(t)] \right\rangle. \tag{31}$$

We denote  $p(x,t;x_{\tau},t-\tau|\psi)\mathrm{d}x\mathrm{d}x_{\tau}$  as the probability that  $x(t)\in[x,x+\mathrm{d}x]$  and  $x(t-\tau)\in[x_{\tau},x_{\tau}+\mathrm{d}x_{\tau}]$  given the initial function  $\psi$ . Then, Eq. (31) is equivalent to

$$\int_{a}^{b} dx F(x) \int_{a}^{b} dx_{\tau} \frac{\partial}{\partial t} p(x, t; x_{\tau}, t - \tau | \psi)$$

$$= \int_{a}^{b} dx \frac{d}{dx} F[x(t)] \int_{a}^{b} dx_{\tau} u(t, x, x_{\tau}) p(x, t; x_{\tau}, t - \tau | \psi).$$
(32)

By applying the integration by parts formula to the right-hand side, we obtain:

$$\int_{a}^{b} dx F(x) \int_{a}^{b} dx_{\tau} \frac{\partial}{\partial t} p(x, t; x_{\tau}, t - \tau | \psi)$$

$$= \int_{a}^{b} dx F(x) \int_{a}^{b} dx_{\tau} \left\{ -\frac{\partial}{\partial x} [u(t, x, x_{\tau}) p(x, t; x_{\tau}, t - \tau | \psi)] \right\}, \tag{33}$$

where the surface terms are neglected due to Eq. (28) and Eq. (29). Since F(x) is arbitrary, Eq. (33) leads to

$$\frac{\partial}{\partial t}p(x,t|\psi) = -\frac{\partial}{\partial x} \left\{ \int_a^b \mathrm{d}x_\tau [u(t,x,x_\tau)p(x,t;x_\tau,t-\tau|\psi)] \right\}.$$
 (34)

# B.2 Proofs of Propositions in the main text

**Proposition 3.1.** (Restriction on trajectory intersections) Suppose that the target transport strategy, guided by certain key points or specific coupling constraints, inevitably leads to trajectory intersections during the transport process, the flow map corresponding to the ODE-based FM cannot precisely preserve the transport strategy.

*Proof of Proposition 3.1.* Consider a time-dependent Neural ODE corresponding to a ODE-based FM with the following form:

$$\frac{\mathrm{d}\boldsymbol{x}(t)}{\mathrm{d}t} = \boldsymbol{v}[t, \boldsymbol{x}(t); \theta], \quad t \in [0, T],$$

$$\boldsymbol{x}(0) = \boldsymbol{x}_0,$$
(35)

where  $v[t, \boldsymbol{x}(t); \theta] : [0, T] \times \mathbb{R}^d \to \mathbb{R}^d$  is a parameterized vector field. We further assume that the vector field is (locally) Lipschitz continuous. Then, by the Picard-Lindelöf Theorem, the corresponding initial value problem (35) has a unique solution on the interval [0, T]. Suppose that there exist two distinct solutions  $\boldsymbol{x}^1(t)$  and  $\boldsymbol{x}^2(t)$  corresponding to different initial values  $\boldsymbol{x}^1_0$  and  $\boldsymbol{x}^2_0$ , which intersect at  $(t^*, \boldsymbol{x}^*)$ . By the uniqueness of the solution, these two trajectories must correspond to the same solution, leading to a contradiction. Therefore, given a transport strategy, which inevitably leads to trajectory intersections during the transport process, the flow map corresponding to the ODE-based FM cannot precisely ensure the transport strategy.

**Proposition 3.3.** (Heterogeneity in distributions) Assume that the source (resp., target) distribution  $q_0$  (resp.,  $q_1$ ) is supported on M (resp., N) disjoint compact sets  $\{U_i^0\}_{i=1}^M$  (resp.,  $\{U_j^1\}_{j=1}^N$ ), where each  $U_i^0 \in \mathbb{R}^d$  (resp.,  $U_j^1 \in \mathbb{R}^d$ ) is a path-connected set with non-zero measure and  $q_0(\mathbf{x}_0) > 0$  (resp.,  $q_1(\mathbf{x}_1) > 0$ ) for all  $\mathbf{x}_0 \in \bigcup_{i=1}^M U_i^0$  (resp.,  $\mathbf{x}_1 \in \bigcup_{j=1}^N U_j^1$ ). If  $M \neq N$ , for any transport map  $T: \bigcup_{i=1}^M U_i^0 \to \bigcup_{j=1}^N U_j^1$  such that  $T_\# q_0 = q_1$ , we have that: (1). Neural ODEs with Lipschitz continuous vector field cannot exactly represent T; (2). If the flow of a Neural ODE is equal to the transport map T almost everywhere, then the vector field is not Lipchitz continuous.

Proof of Proposition 3.3. Without loss of generality, assume M < N. Let  $T: \bigcup_{i=1}^M U_i^0 \to \bigcup_{j=1}^N U_j^1$  be an arbitrary transport map such that  $T_\# q_0 = q_1$ .

- (1). By Dirichlet's drawer principle, there must exist a non-empty path-connected compact set  $U^0_s \in \{U^0_i\}_{i=1}^M$  such that  $T(U^0_s) \not\subseteq U^1_t$  for any  $U^1_t \in \{U^1_j\}_{j=1}^N$ . Since  $T(U^0_s) \subseteq \bigcup_{j=1}^N U^1_j$ , there must exist two different points  $\boldsymbol{x}_0, \boldsymbol{y}_0 \in U^0_s$  and two disjoint sets  $U^1_k, U^1_t$  such that  $T(\boldsymbol{x}_0) \in U^1_k$  and  $T(\boldsymbol{y}_0) \in U^1_t$ . Consider a Neural ODE as Eq. (35), where the vector field  $\boldsymbol{v}[t, \boldsymbol{x}(t); \theta]$  is sufficiently smooth with a bounded Lipschitz constant L, it is well-known that its solution exists and is unique over the entire time interval, and the associated flow map  $\Phi_t(\boldsymbol{x}_0)$  is a diffeomorphism. Since  $U^0_s$  is path-connected, there exists a continuous function  $f:[0,1]\to U^0_s$  such that  $f(0)=\boldsymbol{x}_0, f(1)=\boldsymbol{y}_0$ . Due to the diffeomorphic property of the flow, we have that  $\Phi_T f:[0,1]\to\Phi_T(U^0_s)$  is also a continuous function such that  $\Phi_T f(0)=\Phi_T(\boldsymbol{x}_0), \Phi_T f(1)=\Phi_T(\boldsymbol{y}_0)$ . If  $\Phi_T=T,\Phi_T f$  is a path connecting  $T(\boldsymbol{x}_0)\in U^1_k$  and  $T(\boldsymbol{y}_0)\in U^1_k$ . Since the compact sets  $\bigcup_{j=1}^N U^1_j$  are disjoint, there exists a  $t^*\in(0,1)$  such that  $\Phi_T f(t^*)\not\subseteq\bigcup_{j=1}^N U^1_j$ . Therefore,  $\Phi_T f(t^*)\not=T(f(t^*))$ , leading to a contradiction.
- (2). Without loss of generality, we further consider the case where M=1 and N=2. Under the given condition, there exists a Neural ODE with the flow map  $\Phi_T=T$  on  $U_1^0\backslash Z$ , where Z is a set with zero measure. Since  $U_1^0$  is compact, for any  $\epsilon>0$ , there exist finitely many open balls  $\{O(\boldsymbol{x}_0^k,\epsilon)\}_{k=1}^K$  such that  $U_1^0\subseteq\bigcup_{k=1}^KO(\boldsymbol{x}_0^k,\epsilon)$ , where  $O(\boldsymbol{x}_0^k,\epsilon)$  is the open ball centered at  $\boldsymbol{x}_0^k\in U_1^0$  with radius  $\epsilon$ . From the compactness and path-connectedness of  $U_1^0$ , the following conclusion can be easily derived:

$$\exists O(\boldsymbol{x}_0^s, \epsilon) \in \{O(\boldsymbol{x}_0^k, \epsilon)\}_{k=1}^K$$
s.t. 
$$\exists \boldsymbol{x}, \boldsymbol{y} \in O(\boldsymbol{x}_0^s, \epsilon) \backslash Z, \quad \Phi_T(\boldsymbol{x}) = T(\boldsymbol{x}) \in U_1^1, \Phi_T(\boldsymbol{y}) = T(\boldsymbol{y}) \in U_2^1.$$
(36)

We define the distance between  $U_1^1$  and  $U_2^1$  as d. It is easy to see that d>0. Therefore, we have:

$$\frac{||\Phi_T(\boldsymbol{x}) - \Phi_T(\boldsymbol{y})||}{||\boldsymbol{x} - \boldsymbol{y}||} \ge \frac{d}{2\epsilon}.$$
(37)

By the arbitrariness of the choice of  $\epsilon$ , it follows that the flow map of the Neural ODE is not Lipschitz continuous.

Assume that the vector field is Lipschitz continuous, then there is a constant  $L \geq 0$  such that

$$||\boldsymbol{v}_t(\boldsymbol{x};\theta) - \boldsymbol{v}_t(\boldsymbol{y};\theta)|| < L||\boldsymbol{x} - \boldsymbol{y}||, \tag{38}$$

Then, we have

$$||\Phi_T(\boldsymbol{x}) - \Phi_T(\boldsymbol{y})|| \le e^{LT} ||\boldsymbol{x} - \boldsymbol{y}||, \tag{39}$$

which implies that the flow map is Lipschitz continuous, leading to a contradiction. Therefore, the vector field of the Neural ODE is not Lipschitz continuous.  $\Box$ 

**Proposition 4.1.** The marginal vector field  $\mathbf{u}(t, \mathbf{x}, \mathbf{x}_{\tau})$  given in Eq. (15), together with the selected initial function  $\boldsymbol{\psi}$ , generates the probability path  $p(\mathbf{x}, t | \boldsymbol{\psi})$  and the joint probability path  $p(\mathbf{x}, t; \mathbf{x}_{\tau}, t - \tau | \boldsymbol{\psi})$  in Eq. (14).

*Proof of Proposition 4.1.* Given that  $u(t, x, x_{\tau}|z)$  is the conditional vector field that generates  $p(x, t|z, \psi)$  and  $p(x, t; x_{\tau}, t - \tau|z, \psi)$ , it is evident that the following delay Fokker-Planck equation holds:

$$\frac{\partial p(\boldsymbol{x}, t | \boldsymbol{z}, \boldsymbol{\psi})}{\partial t} = -\nabla \cdot \left\{ \int d\boldsymbol{x}_{\tau} [\boldsymbol{u}(t, \boldsymbol{x}, \boldsymbol{x}_{\tau} | \boldsymbol{z}) p(\boldsymbol{x}, t; \boldsymbol{x}_{\tau}, t - \tau | \boldsymbol{z}, \boldsymbol{\psi})] \right\}. \tag{40}$$

To verify this proposition, we only need to check that the marginal vector field  $\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau})$  given by Eq. (15), the marginal probability density path  $p(\boldsymbol{x},t|\boldsymbol{\psi})$  and the joint probability density  $p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi})$  satisfy the following delay Fokker-Planck equation:

$$\frac{\partial p(\boldsymbol{x}, t | \boldsymbol{\psi})}{\partial t} = -\nabla \cdot \left\{ \int d\boldsymbol{x}_{\tau} [\boldsymbol{u}(t, \boldsymbol{x}, \boldsymbol{x}_{\tau}) p(\boldsymbol{x}, t; \boldsymbol{x}_{\tau}, t - \tau | \boldsymbol{\psi})] \right\}. \tag{41}$$

Assuming that all the functions involved satisfy the regularity conditions necessary for the interchange of integration and differentiation, we have that

$$\frac{\partial p(\boldsymbol{x},t|\boldsymbol{\psi})}{\partial t} = \frac{\partial}{\partial t} \int p(\boldsymbol{x},t|\boldsymbol{z},\boldsymbol{\psi})q(\boldsymbol{z})d\boldsymbol{z} 
= \int \left[\frac{\partial}{\partial t}p(\boldsymbol{x},t|\boldsymbol{z},\boldsymbol{\psi})\right]q(\boldsymbol{z})d\boldsymbol{z} 
= -\int \left(\nabla \cdot \left\{\int d\boldsymbol{x}_{\tau}[\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z})p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi})]\right\}\right)q(\boldsymbol{z})d\boldsymbol{z} 
= -\nabla \cdot \left\{\int d\boldsymbol{x}_{\tau}\left[\int \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z})p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi})q(\boldsymbol{z})d\boldsymbol{z}\right]\right\} 
= -\nabla \cdot \left\{\int d\boldsymbol{x}_{\tau}\left[\int \frac{\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z})p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi})}{p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi})}q(\boldsymbol{z})d\boldsymbol{z}\cdot p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi})\right]\right\} 
= -\nabla \cdot \left\{\int d\boldsymbol{x}_{\tau}\left[\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau})p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi})\right]\right\}.$$
(42)

**Proposition 4.2.** Given that  $p(\mathbf{x}, t; \mathbf{x}_{\tau}, t - \tau | \boldsymbol{\psi}) > 0$  for all  $\mathbf{x}, \mathbf{x}_{\tau} \in \mathbb{R}^d$  and  $t \in [0, 1]$ , up to a constant independent of  $\boldsymbol{\theta}$ ,  $\mathcal{L}_{DCFM}(\boldsymbol{\theta})$  and  $\mathcal{L}_{DFM}(\boldsymbol{\theta})$  are equal, which further implies that  $\nabla_{\boldsymbol{\theta}} \mathcal{L}_{DCFM}(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \mathcal{L}_{DFM}(\boldsymbol{\theta})$ .

Proof of Proposition 4.2. To guarantee the existence of all integrals and the validity of changing the order of integration (as justified by Fubini's Theorem), we assume that  $p(\boldsymbol{x},t|\boldsymbol{z},\psi)$ ,  $p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\psi)$  decay to zero sufficiently fast as  $\|\boldsymbol{x}\|\to\infty$  and  $\|\boldsymbol{x}_{\tau}\|\to\infty$ , and further assume that  $\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z}),\boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\theta),\nabla_{\boldsymbol{\theta}}\boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\theta)$  are bounded. Note that  $t\sim\mathcal{U}(0,1)$  and  $\psi\sim q^{\circ}(\psi)$  are both independent of  $\boldsymbol{\theta}$ , so we fixed t and  $\psi$  in the following analysis. Using the bilinearity of the Euclidean norm and the independence of  $\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z})$  from  $\boldsymbol{\theta}$ , we have:

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi})} \|\boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}) - \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau})\|^{2}$$

$$= \nabla_{\boldsymbol{\theta}} \mathbb{E}_{p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi})} \left( \|\boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta})\|^{2} - 2 \left\langle \boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}),\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau})\right\rangle + \|\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau})\|^{2} \right)$$

$$= \nabla_{\boldsymbol{\theta}} \mathbb{E}_{p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi})} \left( \|\boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta})\|^{2} - 2 \left\langle \boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}),\boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau})\right\rangle \right),$$

$$(43)$$

and

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{q(\boldsymbol{z}),p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi})} \| \boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}) - \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z}) \|^{2}$$

$$= \nabla_{\boldsymbol{\theta}} \mathbb{E}_{q(\boldsymbol{z}),p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi})} (\| \boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}) - 2 \langle \boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}), \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z}) \rangle + \| \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z}) \|^{2})$$

$$= \nabla_{\boldsymbol{\theta}} \mathbb{E}_{q(\boldsymbol{z}),p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi})} (\| \boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}) \|^{2} - 2 \langle \boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}), \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z}) \rangle).$$

$$(44)$$

Next,

$$\mathbb{E}_{p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi})} \|\boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta})\|^{2}$$

$$= \iint \|\boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta})\|^{2} p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi}) d\boldsymbol{x} d\boldsymbol{x}_{\tau}$$

$$= \iiint \|\boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta})\|^{2} p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi}) q(\boldsymbol{z}) d\boldsymbol{z} d\boldsymbol{x} d\boldsymbol{x}_{\tau}$$

$$= \mathbb{E}_{q(\boldsymbol{z}),p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi})} \|\boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta})\|^{2}.$$

$$(45)$$

Finally,

$$\mathbb{E}_{p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi})} \langle \boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}), \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}) \rangle 
= \iint \left\langle \boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}), \frac{\int \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z})p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi})q(\boldsymbol{z})d\boldsymbol{z}}{p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi})} \right\rangle p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{\psi}) d\boldsymbol{x}d\boldsymbol{x}_{\tau} 
= \iint \left\langle \boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}), \int \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z})p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi})q(\boldsymbol{z})d\boldsymbol{z} \right\rangle d\boldsymbol{x}d\boldsymbol{x}_{\tau} 
= \iiint \left\langle \boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}), \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau}|\boldsymbol{z}) \right\rangle p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi})q(\boldsymbol{z})d\boldsymbol{z}d\boldsymbol{x}d\boldsymbol{x}_{\tau} 
= \mathbb{E}_{q(\boldsymbol{z}),p(\boldsymbol{x},t;\boldsymbol{x}_{\tau},t-\tau|\boldsymbol{z},\boldsymbol{\psi})} \left\langle \boldsymbol{v}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{\theta}), \boldsymbol{u}(t,\boldsymbol{x},\boldsymbol{x}_{\tau};\boldsymbol{z}) \right\rangle. \tag{46}$$

From Eq. (45) and Eq. (46), it follows that Eq. (43) is equal to Eq. (44) for any latent variable z, which can be further deduced that  $\nabla_{\theta} \mathcal{L}_{DCFM}(\theta) = \nabla_{\theta} \mathcal{L}_{DFM}(\theta)$ .

**Proposition 4.3.** (Universal approximating capability of DFM). For any given continuous transport map  $F: \mathbb{R}^d \to \mathbb{R}^d$ , which push-forward the source distribution  $q_0$  to the target distribution  $q_1$ , i.e.  $F_\#q_0 = q_1$ , if there exists a neural network that can approximate V(x) = F(x) - x, then we can construct a vector field with a single delay term, where the corresponding flow map, under the constant initial function condition, can approximately push-forward  $q_0$  to  $q_1$ , while preserving the target transport strategy  $x \to F(x)$ .

*Proof of Proposition 4.3.* The proof is straightforward. We consider the following DDE with a constant initial function and a parameterized vector field with a single delay term:

$$\frac{\mathrm{d}\boldsymbol{x}(t)}{\mathrm{d}t} = \boldsymbol{v}[\boldsymbol{x}(t-\tau);\theta], \quad t \in [0,1],$$

$$\boldsymbol{x}(t) \equiv \boldsymbol{x}_0, \quad t \in [-\tau,0],$$
(47)

where the time delay  $\tau=1$ , and the vector field depends solely on the delay term, independent of the current state and time, which is a special degenerate case of DFM. In this case, for any point sampled from the initial distribution  $\boldsymbol{x}_0 \sim q_0$ , the corresponding vector field  $\boldsymbol{v}[\boldsymbol{x}(t-\tau);\theta] = \boldsymbol{v}[\boldsymbol{x}_0;\theta]$  remains constant over the time interval  $t \in [0,1]$ , which means that the flow map  $G: \mathbb{R}^d \to \mathbb{R}^d$  associated with the above DDE will map  $\boldsymbol{x}_0$  to  $\boldsymbol{x}_0 + \boldsymbol{v}[\boldsymbol{x}_0;\theta] \cdot 1$ . We further assume that the neural network  $\boldsymbol{v}[\boldsymbol{x}_0;\theta]$  can approximate  $V(\boldsymbol{x}) = F(\boldsymbol{x}) - \boldsymbol{x}$ , which implies that  $G(\boldsymbol{x}_0) = \boldsymbol{x}_0 + \boldsymbol{v}[\boldsymbol{x}_0;\theta] \cdot 1 \approx \boldsymbol{x}_0 + [F(\boldsymbol{x}_0) - \boldsymbol{x}_0] = F(\boldsymbol{x}_0)$ . This implies that the flow map of the Neural DDEs approximately preserves the target transport strategy, and naturally, it can approximately push-forward  $q_0$  to  $q_1$ .  $\square$ 

## C EXPERIMENTAL SETUP DETAILS AND ADDITIONAL RESULTS

In this section, we provide a detailed explanation of the experimental settings for different datasets described in the main text and present additional experimental results. The experimental details for the delay dynamical systems and single-cell datasets are summarized in Table 5).

Table 5: The experimental setup details for delay dynamical systems and single-cell datasets.

	Setup	Autoregulation	Spiral DDE	Mouse hematopoiesis	iPSCs
Data	Dimension	1	2	2	4
	Hidden layer	3	3	3	3
	Hidden neuron	64	64	64	64
Structure	Activation	Tanh	Tanh	SELU	SELU
	Time input	False	False	True	True
	Delay term input	True (DFM) / False (CFM)			
	Coupling	KPG-OT	OT	OT	OT
	Initical function (DFM)	Constant	Constant	Diverse	Diverse
	Latent variable $\gamma$	CSpline	CSpline	Geodesic	CSpline
Training	Batch size	256	256	128	128
	Iteration	2k	2k	10k	2k
	Optimizer	Adam	Adam	Adam	Adam
	Learning rate	$10^{-3}$	$10^{-3}$	$10^{-3}$	$10^{-3}$
Results	Figures & tables	Figs. 2 & Tab. 6, 7, 8	Figs. 2 & Tab.9, 10	Fig. 3 & Tab. 2	Fig. 3 & Tab. 2

## C.1 BIOLOGICAL AUTOREGULATION MOTIF

**Dataset generation**. The autoregulation motif, which is one of the most common biological network motifs, can be modeled as the following 1-d DDE:

$$\dot{x}(t) = \frac{\eta}{1 + x^n(t - \tau)},\tag{48}$$

where  $\tau=1, n=2, \eta=5$ . As shown in Fig. 2, the dynamics exhibit damped oscillations. We sample 1,000 initial points from the uniform distribution  $\mathcal{U}(0.2,1.2)$  and perform forward integration based on constant initial functions  $\psi(h; \boldsymbol{x}_0) = \boldsymbol{x}_0, h \in [-1,0]$ , to obtain trajectories. After that, we select snapshots at  $t=0,1,2,\cdots,7$  as the training dataset.

**Details for training**. For each training iteration, we randomly and independently sample 256 data points from each snapshot. Pairing between minibatch data points from adjacent time steps is performed using KPG-OT. For this task, we utilize the true coupling of only two keypoint pairs between adjacent time steps as guidance, although incorporating more keypoint pairs would enable more accurate pairing. Once data points across time steps are fully matched, we construct 256 transport trajectories using cubic spline (CSpline) interpolation. These trajectories provide the states, delayed states, and vector fields at various time steps, enabling the computation of the training objective Eq. (16), which is then optimized via backpropagation. Further details on the experimental setup can be found in Table 5.

**Details for testing.** After training, we evaluate the performance of KP-DFM(C) and KP-CFM on both interpolation and extrapolation tasks. Specifically, we randomly select initial points from the initial distribution  $\mathcal{U}(0.2, 1.2)$  and generate predicted trajectories by performing forward integration on the learned vector field, obtaining the predicted trajectories and snapshots at various time steps. For interpolation, we compute the mean 2-Wasserstein distance  $W_2$  between the predicted and ground-truth distributions at  $t=0.5, 1.5, 2.5, \cdots, 6.5$ . Additionally, since the true dynamics are known, we calculate the error between the predicted and ground-truth trajectories at these points, quantified by the Mean Squared Error (MSE) and Mean Relative Error (MRE). For extrapolation, we similarly compute the  $W_2$  distance between the predicted and ground-truth distributions at t=8,9,10, along with the trajectory-wise error using MSE and MRE.

**Results**. As shown in Table 7 and 8, KP-DFM(C) significantly outperforms CFM in both interpolation and extrapolation tasks, achieving much lower prediction errors for both distributions and individual trajectories, demonstrating superior representational capability. Furthermore, as illustrated in Table 6, CFM fails to reconstruct the underlying dynamics, whereas KP-DFM(C) successfully captures the damped oscillatory dynamics, even when the time delay  $\tau$  deviates from the ground truth.

## C.2 SPIRAL DDE

**Dataset generation**. We consider the following 2-d DDE:

$$\dot{\boldsymbol{x}}(t) = \boldsymbol{A} \tanh(\boldsymbol{x}(t) + \boldsymbol{x}(t - \tau)), \tag{49}$$

where  $\tau=0.5$  and  $\boldsymbol{A}=[[-1,20],[-30,-1]]\in\mathbb{R}^{2\times2}$ . As shown in Fig. 2, the dynamics exhibit trajectory crossings. We sample 1,000 initial points from the uniform distribution  $\mathcal{U}(0.8,1)\times\mathcal{U}(0,0.1)$  and perform forward integration based on constant initial functions  $\psi(h;\boldsymbol{x}_0)=\boldsymbol{x}_0,h\in[-0.5,0]$ , to obtain trajectories. After that, snapshots at intervals of 0.05 within  $t\in[0,0.6]$  are selected as training dataset.

**Details for training**. In each training iteration, 256 data points are randomly and independently sampled from each snapshot. Pairing between the minibatch data points at adjacent time steps is performed using OT. After matching data points across time steps, 256 transport trajectories are constructed using CSpline interpolation. These trajectories provide the states, delayed states, and vector fields at various time steps, enabling the computation of the training objective in Eq. (16), which is subsequently optimized via backpropagation. Additional details on the experimental setup are provided in Table 5.

**Details for testing**. After training, we assess the performance of OT-DFM(C) and OT-CFM on the interpolation task. Specifically, 100 initial points are randomly sampled from the initial distribution  $\mathcal{U}(0.8,1)\times\mathcal{U}(0,0.1)$ , and predicted trajectories are generated via forward integration of the learned

Table 6: Additional results on biological autoregulation motif based on CFM with keypoint-guided coupling (KP-CFM), as well as DFM with constant initial functions and keypoint-guided coupling (KP-DFM(C)). For DFM, the generation results are illustrated with various time delays, specifically  $\tau = 0.6, 0.8, 1.0$  (ground truth), 1.2, 1.4.

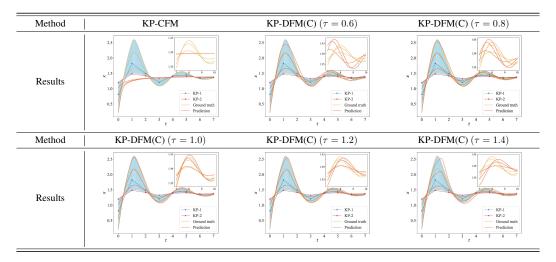


Table 7: Comparison of the interpolation prediction error between KP-CFM and KP-DFM(C) with different time delay  $\tau$  on the biological autoregulation dataset.

$\overline{\tau}$	0(CFM)	0.6	0.8	1.0	1.2	1.4
$W_2$ MSE MRE	0.1842 9.21e-2 0.117	1.51e-3	4.42e-4	3.45e-4	5.36e-4	2.68e-3

vector field, yielding trajectories and snapshots at various time steps. The mean  $W_2$  distance between the predicted and ground-truth distributions is computed at intervals of 0.01 within  $t \in [0,0.6]$ . Moreover, we calculate the error between the predicted and ground-truth trajectories at these points, quantified by MSE and MRE.

**Results**. As shown in Table 10, OT-DFM(C) significantly outperforms CFM in the interpolation prediction task, achieving substantially lower errors for prediction of both distributions and individual trajectories. Additionally, as illustrated in Table 9, CFM collapses around the trajectory crossing area, while OT-DFM(C) successfully captures the spiral dynamics, even when the time delay  $\tau$  deviates from the ground truth.

Table 8: Comparison of the extrapolation prediction error between KP-CFM and KP-DFM with different time delay  $\tau$  on the biological autoregulation dataset.

$\overline{\tau}$	0(CFM)	0.6	0.8	1.0	1.2	1.4
$\overline{W_2}$	2.13e-2	2.25e-2	8.67e-3	3.42e-3	4.55e-3	1.23e-2
<b>MSE</b>	5.51e-4	6.31e-4	1.28e-4	1.46e-5	4.38e-5	2.13e-4
MRE	1.59e-2	1.68e-2	6.75e-3	2.26e-3	3.72e-3	9.37e-3

Table 9: Additional results on Spiral DDE based on CFM with OT coupling (OT-CFM), as well as DFM with constant initial functions and OT coupling (OT-DFM(C)). For DFM, the generation results are illustrated with various time delays, specifically  $\tau = 0.35, 0.40, 0.45, 0.50$  (ground truth), 0.55.

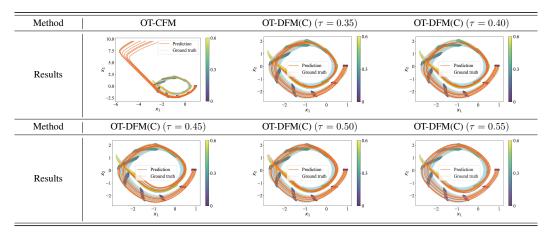


Table 10: Comparison of the interpolation prediction error between OT-CFM and OT-DFM with different time delay  $\tau$  on the spiral DDE dataset.

$\overline{\tau}$	0(CFM)	0.35	0.40	0.45	0.50	0.55
$W_2$ MSE MRE	3.95 15.81 13.14	0.07	0.43 0.13 1.66	0.12	0.03	0.45 0.14 1.96

# C.3 MOUSE HEMATOPOIESIS DATASET

**Data preprocessing.** To efficiently apply various trajectory inference methods to scRNA-seq data from mouse hematopoiesis, we first project the data into a low-dimensional space. Specifically, we use the reduced two force-directed layouts (SPRING) space for this dataset with batch correction, following the same preprocessing pipeline as TIGON (Sha et al., 2024). Besides, we use the DBSCAN clustering algorithm to categorize all cells at Day 6 into two distinct clusters, corresponding to the Neu fate and Mo fate, respectively.

**Selection of initial functions.** During training, we match the sampled minibatch cells at Day 2 and 6 using OT. For each matched cell pairs  $(x_0, x_1)$ , where  $x_0$  is sampled from data of Day 2 and  $x_1$  is sampled from data of Day 6, the initial function is defined as:

$$\frac{\mathrm{d}\boldsymbol{\psi}^*(t;\boldsymbol{x}_0)}{\mathrm{d}t} = \begin{cases} [1,0], & \text{if } \boldsymbol{x}_1 \text{ is in Neu cluster,} \\ [0,1], & \text{if } \boldsymbol{x}_1 \text{ is in Mo cluster.} \end{cases}$$
(50)

**Details of training.** During training, we align Day 2 and Day 6 with t=0 and t=1, respectively. For testing, the distribution of Day 4 is predicted at t=0.5. In this task, the time delay is set to  $\tau=1$ . In each training iteration, we randomly and independently sample 128 cells from each snapshot at Day 2 and 6. Pairing between minibatch data points at adjacent time steps is performed using OT. After matching data points across time steps, 128 transport trajectories  $\gamma$  are constructed using geodesic interpolation (refer to the next paragraph). These trajectories provide the states, delayed states, and vector fields at various time steps, enabling the computation of the training objective in Eq. (16), which is subsequently optimized via backpropagation. Additional details on the experimental setup are provided in Table 5.

Construct geodesic interpolation  $\gamma$ . For each matched cell pairs, geodesic interpolation is employed to generate the interpolation trajectory (Kapuśniak et al., 2024). Here, we first build a data-induced Riemannian metric g on the data manifold  $\mathcal{M}$ . Formally, a Riemannian metric g on  $\mathcal{M}$  is a smooth

family of inner products on the tangent spaces of  $\mathcal{M}$ . Specifically, g provides each  $x \in \mathcal{M}$  a positive defined symmetric bilinear form on  $T_x \mathcal{M}$ ,

$$g_x: T_x \mathcal{M} \times T_x \mathcal{M} \to \mathbb{R},$$
 (51)

which induces a norm  $||\cdot||_{g_x}: T_x\mathcal{M} \to \mathbb{R}$  defined by  $||v||_{g_x} = \sqrt{g_x(v,v)}$ . Based on this norm, we can calculate the length of any path connecting two points on the manifold and construct the geodesic between them. We aim to equip the ambient space  $\mathbb{R}^d$  with an appropriate metric that constrains geodesics to remain close to the data manifold. Specifically, given coordinates, the Riemannian metric g can be equivalently represented as a state-dependent  $d \times d$  symmetric positive definite matrix  $\mathbf{G}(x)$ , such that  $g_x(u,v) = u^{\top}\mathbf{G}(x)v$  for any  $u,v \in T_x\mathcal{M}$ . Intuitively, to impose manifold constraints on the transport path,  $||\mathbf{G}(x)||$  is supposed to take smaller values around data points and larger values in regions far from all points in the dataset  $\mathcal{D}$ . Therefore, we employ the LAND metric, where  $\mathbf{G}(x;\mathcal{D}) = (\operatorname{diag}(\mathbf{h}(x;\mathcal{D}) + \epsilon \mathbf{I})^{-1}$ . Here, the components of  $\mathbf{h}(x;\mathcal{D})$  are defined as

$$h_k(\mathbf{x}; \mathcal{D}) = \sum_{i=1}^{N} (x_k^i - x_k)^2 \exp(-\frac{1}{2\sigma^2} ||\mathbf{x} - \mathbf{x}^i||_2^2),$$
 (52)

where  $1 \le k \le d$  denotes the k-th component, the superscript i represents the i-th data point in the training dataset  $\mathcal{D}$  and  $\sigma$  is the kernel bandwidth. Subsequently, we utilize a neural network  $\phi_{\eta}$  to parameterize the geodesic interpolation path between  $\boldsymbol{x}_0$  and  $\boldsymbol{x}_1$  as

$$\gamma_{\eta}(t, \mathbf{x}_0, \mathbf{x}_1) = (1 - t)\mathbf{x}_0 + t\mathbf{x}_1 + t(1 - t)\phi_{\eta}(t, \mathbf{x}_0, \mathbf{x}_1). \tag{53}$$

The trajectory  $\gamma_{\eta}$  approximates the geodesic between any two points sampled at adjacent time steps on the Riemannian manifold  $(\mathcal{M},g)$  implies learning parameters  $\eta$  to minimize the following loss function:

$$\mathcal{L}_{\mathbf{G}}(\eta) := \mathbb{E}_{\pi(\boldsymbol{x}_0, \boldsymbol{x}_1)} \int_0^1 \dot{\boldsymbol{\gamma}}_{\eta}(t, \boldsymbol{x}_0, \boldsymbol{x}_1)^{\top} \mathbf{G}(\boldsymbol{\gamma}_{\eta, t}; \mathcal{D}) \dot{\boldsymbol{\gamma}}_{\eta}(t, \boldsymbol{x}_0, \boldsymbol{x}_1) dt, \tag{54}$$

where  $\pi$  represents the OT coupling. The optimal parameter is obtained as:

$$\eta^* = \arg\min_{\eta} \mathcal{L}_{\mathbf{G}}(\eta), \tag{55}$$

By substituting  $\eta^*$  into  $\gamma_{\eta}$  in Eq. (53), the resulting trajectory  $\gamma_{\eta^*}$  provides an approximation of the geodesic interpolation on the data manifold.

# C.4 iPSCs dataset

**Data Preprocessing**. To apply trajectory inference methods to the iPSC scRNA-seq data efficiently, we first reduce the data to a 4-dimensional space using Principal Component Analysis (PCA), following the preprocessing pipeline of TIGON (Sha et al., 2024). For cells after differentiation (from Day 3), we categorize them into two groups corresponding to the M and En fates using a Gaussian Mixture Model (GMM) with two components. For Day 2, we model the distribution with a single-component GMM and denote its mean as  $\mu_0$ . As described above, for Day 5, the distribution is modeled with a two-component GMM, where the means of the M and En fates are denoted as  $\mu_{11}$  and  $\mu_{12}$ , respectively.

**Selection of initial functions.** During training, we match the sampled minibatch cells at Day 2, 4 and 5 using OT. For each matched triplet  $(x_0, x_1, x_2)$ , where  $x_0$  is sampled from the Day 2 dataset,  $x_1$  from the Day 4 dataset, and  $x_2$  from the Day 5 dataset, the initial function is defined as:

$$\frac{\mathrm{d}\psi^*(t; \boldsymbol{x}_0)}{\mathrm{d}t} = \begin{cases} \boldsymbol{\mu}_{11} - \boldsymbol{\mu}_0, & \text{if } \boldsymbol{x}_2 \text{ is in M component,} \\ \boldsymbol{\mu}_{12} - \boldsymbol{\mu}_0, & \text{if } \boldsymbol{x}_1 \text{ is in En component.} \end{cases}$$
(56)

**Details of training.** During training, we align Day 2, 4 and 5 with t=0, t=2 and t=3, respectively. For testing, the distribution of Day 3 is predicted at t=1. In this task, the time delay is set to  $\tau=3$ . In each training iteration, we randomly and independently sample 128 cells from each snapshot at Day 2, 4, and 5. Pairing between minibatch data points at adjacent time steps is performed using OT. After matching data points across time steps, 128 transport trajectories  $\gamma$  are constructed using CSpline interpolation. These trajectories provide the states, delayed states, and vector fields at various time steps, enabling the computation of the training objective in Eq. (16), which is subsequently optimized via backpropagation. Additional details on the experimental setup are provided in Table 5.

## C.5 MNIST DATASET

**Dataset Generation**. We propose a *Semi-paired Image-to-Image Translation* task using the MNIST dataset. The source domain consists of the original MNIST images, while the target domain comprises their corresponding negative images. Specifically, the pixel values X of each image in the original MNIST dataset are inverted by replacing them with 1-X, which flips the brightness. The objective is to map each image to its corresponding negative counterpart, as demonstrated in Table 11.

Minibatch coupling. To provide partial supervision, 10% of the training data are paired with their negative counterparts as keypoints. During training, minibatches are independently sampled from the source and target distributions, and KPG-OT coupling (KP-) is applied. Specifically, for each batch, we sample the same number of images from the source distribution and the target distribution. Out of these, around 10% image pairs are selected from the keypoint set, while the remaining images are independently sampled from the respective distributions. It is important to note that the negative counterparts of the remaining images sampled from the initial distribution are not guaranteed to appear in the target batch. Subsequently, we perform KPG-OT pairing on the remaining images in the batch based on the selected keypoint pairs, which are then used for further training.

Table 11: Samples of generation process and results on MNIST data based on CFM with keypoint-guided coupling (KP-CFM), as well as DFM with constant initial functions and keypoint-guided coupling (KP-DFM(C)). For DFM, the generation results are illustrated with various time delays, specifically  $\tau=0.125,\,0.25,\,0.5,\,$  and 1.0. In each image, the first and last columns represent samples obtained from the source data and their negative counterparts, respectively, while the 10 intermediate columns depict the generation process.

Method	KP-CFM	KP-DFM(C) ( $\tau = 0.125$ )	KP-DFM(C) ( $\tau = 0.25$ )
Generation	00000000000000000000000000000000000000	00000000000000000000000000000000000000	00000000000000000000000000000000000000
Method	KP-CFM	$\text{KP-DFM(C)} (\tau = 0.5)$	$\text{KP-DFM(C)} \left(\tau = 1.0\right)$
Generation	00000000000000000000000000000000000000	00000000000000000000000000000000000000	00000000000000000000000000000000000000

## C.6 CIFAR-10 DATASET

**Source distribution**. The source distribution for  $x_0$  is defined as follows: with a probability of 50%,  $x_0$  is sampled as  $x_0 = torch.randn(3, 32, 32)/4 - 0.5$ , and with the remaining 50% probability, it is sampled as  $x_0 = torch.randn(3, 32, 32)/4 + 0.5$ , following the setting in Zhu & Lin (2024). A sample drawn from the source distribution is shown in Table 12.

**Trainable initial functions.** Designing appropriate initial values for high-dimensional images in image generation tasks is a particularly challenging problem. To address this, we employ torch.nn.Embedding to automatically learn the initial functions. Specifically, the source distribution consists of two modules, corresponding to two Gaussian distributions, while the target distribution consists of ten modules, corresponding to ten image categories. Therefore, for each pair of modules (m,n), where  $m \in \{0,1\}$  and  $n \in \{0,1,\cdots,9\}$ , we need to design the time derivative  $C_{mn}$  of the initial function (Eq. (22)). To achieve this, we begin by defining Emb=torch.nn.Embedding(20,  $3 \times 32 \times 32$ ), which maps 20 discrete indices to corresponding tensors of the same shape as the

image. Then, for each pair (m, n), we define the initial function as:

$$\frac{\mathrm{d}\boldsymbol{\psi}_{mn}^{*}(t;\boldsymbol{x}_{0})}{\mathrm{d}t} = \mathrm{Emb}(10m+n), \ \boldsymbol{\psi}_{mn}^{*}(0;\boldsymbol{x}_{0}) = \boldsymbol{x}_{0},$$
 (57)

where  $h \in [-\tau, 0]$ . This approach allows the derivative of the initial functions across different modules to adaptively adjust during training, thereby enabling more flexible and efficient learning.

**Experimental setup details**. For the image generation tasks, both CFM and DFM are trained following the configurations outlined in Tong et al. (2023a;b). In particular, we utilize a UNet architecture with the following structures and training parameters:

- channels = 128,
- depth = 2,

- channels multiple = [1, 2, 2, 2],
- heads = 4,
- heads channels = 64,
- attention resolution = 16,
- dropout = 0.1,
- batch size per gpu = 128, gpus = 1,
- Adam optimizer with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$ , and no weight decay,
- learning rate =  $2 \times 10^{-4}$ ,
- gradient clipping with norm = 1.0,
- exponential moving average weights with decay = 0.9999.

Table 12: Samples from the source distribution and generation results based on CFM and DFM(D) using independent coupling (I-) or OT coupling (OT-) for adaptive NFE on transporting the Gaussian mixture model with 2 components to the CIFAR-10 dataset.

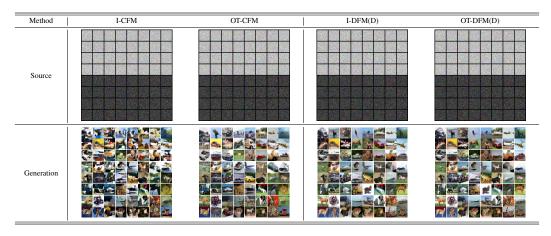


Table 13: Samples from the source distribution and generation results based on CFM and DFM(D) using independent coupling (I-) or OT coupling (OT-) for adaptive NFE on transporting the Gaussian mixture model with 2 components to the CIFAR-10 dataset.

Method	NFE=10	NFE=20	NFE=30	NFE=40
I-CFM				
OT-CFM				
I-DFM(D)				
OT-DFM(D)				

#### D SENSITIVITY ANALYSIS OF TIME DELAY au

To evaluate the effect of time delay  $\tau$  selection on DFM performance, we perform a sensitivity analysis using both synthetic and real-world datasets.

First, we examine how different values of  $\tau$  impact the reconstruction of delay dynamical systems. For the biological autoregulation motif, where the true time delay is  $\tau=1$ , we conduct extensive experiments with  $\tau$  values ranging from 0.6 to 1.4. As shown in Table 6, 7 and 8, while the prediction errors for both the distribution and trajectories increase as the delay deviates from the true value, KP-DFM(C) consistently captures the damped oscillation pattern across a wide range of delay values. In contrast, CFM fails to capture this behavior accurately. For the spiral DDE, with a true delay of  $\tau=0.5$ , we experiment with delays from 0.35 to 0.55. Although the prediction errors become larger when the delay does not match the true value, OT-DFM(C) consistently outperforms CFM, achieving significantly better results, as shown in Table 10. Additionally, as illustrated in Table 9, OT-DFM(C) successfully captures the spiral dynamics across different delay values, whereas CFM exhibits divergence at the intersections of the trajectories.

Furthermore, we explore the effect of different time delay values on image generation performance. Using the MNIST dataset, we compare the generated images for  $\tau=0.125, 0.25, 0.5,$  and 1.0. As shown in Table 11, the generated images from DFM significantly outperform those generated by CFM, with FID scores also notably lower than CFM's, as summarized in Table 3. Upon examining the impact of varying time delays, we find that performance is least favorable at  $\tau=0.125,$  which can be attributed to DFM's behavior approaching that of CFM as  $\tau$  tends towards 0. However, for values of  $\tau>0.25,$  the generation quality improves substantially.