# Coupled Flow Matching

**Wenxi Cai** [1]  **Yuheng Wang** [1]  **Naichen Shi** [2]
[1] University of Michigan   [2]Northwestern University

## Abstract

We introduce Coupled Flow Matching (CPFM), a framework that integrates controllable dimensionality reduction and high-fidelity reconstruction. CPFM learns coupled continuous flows for both the high-dimensional data $x$ and the low-dimensional embedding $y$, which enables sampling $p(y|x)$ via a latent-space flow and $p(x|y)$ via a data-space flow. Unlike classical dimension-reduction methods, where information discarded during compression is often difficult to recover, CPFM preserves the knowledge of residual information within the weights of a flow network. This design provides a bespoke controllability: users may decide which semantic factors to retain explicitly in the latent space, while the complementary information remains recoverable through the flow network. Coupled flow matching builds on two components: (i) an extended Gromov–Wasserstein optimal transport objective that establishes a probabilistic correspondence between data and embeddings, and (ii) a dual-conditional flow-matching network that extrapolates the correspondence to the underlying space. Experiments on multiple benchmarks show that CPFM yields semantically rich embeddings and reconstructs data with higher fidelity than existing baselines.

## 1 Introduction

Dimension reduction has long been a central tool in statistical learning for extracting informative low-dimensional embeddings $y$ from high-dimensional data $x$. A classical example is Principal Component Analysis (PCA) (Shlens, 2014), which efficiently captures dominant linear structures. Nonlinear approaches such as Laplacian Eigenmaps (Belkin and Niyogi, 2003), t-SNE (van der Maaten and Hinton, 2008), and UMAP (McInnes et al., 2020) extend this principle to nonlinear mappings, with an emphasis on preserving local neighborhoods. These embeddings $y$ have proven effective for visualization, denoising, and downstream analysis, and are now indispensable across many scientific domains(Kobak and Berens, 2019; Armstrong et al., 2021; Becht et al., 2019).

Nevertheless, dimension reduction inevitably loses information, as the mapping from high-dimensional $x$ to low-dimensional $y$ is non-invertible by nature. As a result, the global structure is hard to preserve. Furthermore, reconstructing $x$ from $y$ becomes challenging.

Recent advances in generative models have sought to employ deep neural networks to establish statistical connections between $x$ and $y$. Autoencoders (Hinton and Salakhutdinov, 2006) and VAEs (Kingma and Welling, 2022) enable simultaneous data compression and reconstruction using encoder and decoder neural networks. Generative Adversarial Networks (GANs) (Goodfellow et al., 2014), such as Style-GAN (Karras and Aila, 2019), further demonstrate the power of learned latent spaces in generating high-fidelity samples from compact codes. Despite their success, these approaches still face two critical challenges: (i) the latent representation $y$ often entangles multiple explanatory factors in $x$, making it difficult to interpret or disentangle task-relevant information, and (ii) the geometry of the latent space is typically an implicit byproduct of training rather than an explicitly controlled structure.

In light of such limitations, we propose **controllability** as a design principle. We advocate a controllable dimensionality reduction paradigm that makes the choices of important information explicit. Specifically, the embedding $y$ should: (i) highlight user-specified or task-relevant statistics, (ii) enforce latent geometry through constraints or priors, and (iii) neglect nuisance variation while retaining the ability to recover it. This controllability principle yields three key benefits: targeted submanifold design guided by prior knowledge, improved interpretability, and informed generation of new data $x$.

To implement this principle, we propose **Coupled**

**Flow Matching (CPFM)**, a controllable dimensionality reduction framework that learns bidirectional flows between high-dimensional data $x$ and low-dimensional embeddings $y$. CPFM operates in two stages.

In the *first stage*, we construct a probabilistic coupling between discrete samples of $x$ and $y$ via a generalized Gromov–Wasserstein optimal transport (GWOT). Our formulation incorporates a kernel function on the $x$-space to encode relational structure, semantic labels, or other user-specified priors. This kernel induces a transport cost that quantifies the distortion between relations in the original $x$-space and those in the embedding $y$-space, thereby allowing explicit control over the coupling. The resulting optimal transport plan $\pi^{\mathrm{OT}}$ is thus guided directly by user-defined knowledge. To overcome the computational hardness of generic OT, we design an alternating minimization algorithm for an entropy-regularized version of GWOT whose per-iteration complexity is $\mathcal{O}(n^2)$, where $n$ is the dataset size.

In the *second stage*, we extrapolate $\pi^{\mathrm{OT}}$ beyond the training samples to define conditional distributions $p(y|x)$ and $p(x|y)$ via DCMF. While classical flow matching transports samples between two fixed distributions, conditional variants leverage side information to learn flow fields. We extend this paradigm by introducing a dual conditional mechanism within a shared drift network: conditioning on $x$, the network predicts flows on $y$ to generate $p(y|x)$; conditioning on $y$, it predicts flows on $x$ to generate $p(x|y)$. During training, we develop a mute-masking strategy to ensure that only the active direction contributes to the loss. This design leverages the symmetry between the two conditional flows and avoids redundant modeling.

We summarize our main contributions:

- **Generalized GWOT.** We propose a new kernelized quadratic optimal transport objective that embeds semantic priors and enables controllable alignment between data $x$ and embeddings $y$.

- **Efficient OT solver.** We develop a new alternating optimization algorithm for entropy-regularized GWOT with per-iteration complexity $\mathcal{O}(n^2)$, scalable to large datasets.

- **Dual conditional flow matching.** We design a new shared drift network with dual conditioning to jointly learn $p(y|x)$ and $p(x|y)$.

- **Comprehensive evaluation.** We implement CPFM on a range of dimensionality reduction and reconstruction benchmarks, highlighting superior performance in both embedding quality and generative fidelity.

## 2 Related Work

We review several lines of research that are most relevant to our framework.

**Flow matching.** Flow Matching learns continuous-time velocity fields and has recently scaled to large-scale image generation (Lipman et al., 2023). Several extensions improve efficiency and flexibility: Rectified Flow straightens transport paths for accelerated sampling (Liu et al., 2022), with a variational extension introducing a variational objective into this framework (Guo and Schwing, 2025). Stochastic Interpolants unify diffusion and deterministic flows through interpolation-based objectives with likelihood control (Albergo et al., 2023). On the algorithmic side, Weighted CFM reduces variance via importance weighting (Calvo-Ordonez et al., 2025), Context-Varying CFM adapts conditioning to changing contexts (Generale et al., 2025), and Optimal Flow Matching directly optimizes straight-line paths (Tong et al., 2024). Despite these advances, standard flow-matching models construct probability flows in fixed-dimensional spaces and do not provide a mechanism for dimensionality reduction.

**Cross-modal generative modeling.** Multimodal diffusion and flow-based models extend generation to joint distributions across modalities. For example, UniDiffuser fits multimodal joint distributions with a unified model (Bao et al., 2023), EasyGen (Zhao et al., 2024) combines bi-directional conditional diffusion with large language models for interactive cross-modal tasks, and Dual Diffusion Transformer (Li et al., 2025) provides a backbone for joint text–image modeling. While powerful, these models aim to capture cross-modal correspondences rather than to perform controllable dimension reduction.

**Bi-directional generative modeling.** Several bi-directional architectures have been proposed to jointly learn mappings between data and latent representations. Adversarially Learned Inference (ALI) and Bi-GAN (Dumoulin et al., 2017) couple generator and inference networks through adversarial training. Cycle-GAN and MUNIT (Huang et al., 2018) promote cycle consistency to learn unpaired inverse mappings. Multimodal VAEs (MVAE) extend VAEs with product-of-experts inference to handle missing modalities (Wu and Goodman, 2018). More recently, normalizing-flow approaches explicitly reduce reconstruction error through invertible transformations (Lee et al., 2025). However, these methods provide limited ability to regulate which information is preserved or discarded, hence lacking controllability.

Wenxi Cai [1], Yuheng Wang [1], Naichen Shi [2]

**Optimal transport.** Optimal Transport (OT) offers a principled way to compare probability measures. Entropic regularization and the Sinkhorn algorithm (Cuturi, 2013) made OT scalable (Peyré and Cuturi, 2020), with further improvements such as stabilized sparse scaling (Schmitzer, 2019). Quadratic-form OT generalizes OT to quadratic objectives over couplings (Wang and Zhang, 2025). Gromov–Wasserstein (GW) OT compares relational structures, with variants for averaging (Peyré et al., 2016) and fused forms that integrate features with geometry (Vayer et al., 2018). Recent work has studied the stability and algorithmic properties of entropic GW (Rioux et al., 2024). Nonetheless, existing OT methods remain confined to aligning discrete samples and generally lack mechanisms for generative modeling.

## 3 Preliminaries

Before presenting the details of CPFM, we briefly review the necessary background in optimal transport and flow matching generative models.

### 3.1 Linear Entropic Optimal Transport

Let $\mathcal{X} \subset \mathbb{R}^d$ and $\mathcal{Y} \subset \mathbb{R}^d$ denote the source and target domains. A transport cost $c : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$ measures the cost of moving mass between points. Given marginal distributions $\mu_{\mathcal{X}}$ and $\mu_{\mathcal{Y}}$, the Linear Entropic Optimal Transport (LEOT) problem seeks an optimal coupling $\pi$ that minimizes $\inf_{\pi \in \Pi(\mu_{\mathcal{X}}, \mu_{\mathcal{Y}})} \int_{\mathcal{X} \times \mathcal{Y}} c(x,y)\, d\pi(x,y) + \varepsilon\, \mathrm{KL}(\pi \,\|\, \mu_{\mathcal{X}} \times \mu_{\mathcal{Y}})$, where $\Pi(\mu_{\mathcal{X}}, \mu_{\mathcal{Y}})$ is the set of couplings with marginals $\mu_{\mathcal{X}}$ and $\mu_{\mathcal{Y}}$, $\varepsilon > 0$ controls entropy regularization, and KL denotes Kullback–Leibler divergence.

In the finite-sample case with $n$ support points, LEOT can be efficiently solved using the Sinkhorn algorithm, with per-iteration complexity $\mathcal{O}(n^2)$ (Cuturi, 2013).

### 3.2 Quadratic Optimal Transport

Quadratic Optimal Transport (QOT) generalizes the linear case by considering pairwise relations. It minimizes $\inf_{\pi \in \Pi(\mu_{\mathcal{X}}, \mu_{\mathcal{Y}})} \iint_{\mathcal{X}^2 \times \mathcal{Y}^2} c(x, x', y, y')\, d\pi(x,y)\, d\pi(x', y')$, where $c : (\mathcal{X} \times \mathcal{X}) \times (\mathcal{Y} \times \mathcal{Y}) \to \mathbb{R}^+$ takes inputs from two domains. Unlike LEOT, QOT naturally accommodates $\mathcal{X}$ and $\mathcal{Y}$ of different dimensions.

In practice, finite-sample QOT with entropic regularization can be optimized via mirror descent, though with higher computational cost, typically $\mathcal{O}(n^3)$ per iteration (Peyré et al., 2016). A widely used special case is Gromov–Wasserstein OT (GWOT), where the cost is defined as $c_{\mathrm{GWOT}} = \left| d_{\mathcal{X}}(x, x') - d_{\mathcal{Y}}(y, y') \right|^2$,

which represents the distortion between intra-domain distances. Entropic GWOT admits an elegant variational formulation that enables alternating minimization with $\mathcal{O}(n^2)$ complexity (Rioux et al., 2024). In CPFM, we extend GWOT by introducing kernelized costs that encode richer forms of prior knowledge.

### 3.3 Flow Matching

Flow matching (FM) provides a framework for learning generative models by transporting a base distribution $p_0$ (often $\mathcal{N}(0, I)$) to a target distribution $p_1$. The goal is to learn a time-dependent vector field $u_\theta : \mathbb{R}^d \times [0,1] \to \mathbb{R}^d$ that induces an interpolating path $p_t$ between $p_0$ and $p_1$.

Following the stochastic interpolants framework (Albergo et al., 2023), we construct reference paths by sampling $x(0) \sim p_0$, $x(1) \sim p_1$, and defining $x(t) = a_t x(0) + b_t x(1)$ for differentiable scalar functions $(a_t, b_t)$ with boundary conditions $(a_0, b_0) = (1, 0)$ and $(a_1, b_1) = (0, 1)$. The instantaneous conditional velocity is $v_t(x(t); x(0), x(1)) := \frac{d}{dt} x(t) = \dot{a}_t x(0) + \dot{b}_t x(1)$, and the optimal state-dependent drift is the conditional expectation $u^*(x, t) = \mathbb{E}_{x(0), x(1)}[v_t(x(t); x(0), x(1)) \mid x(t) = x]$.

A neural network $u_\theta$ is trained to approximate $u^*$ by minimizing the conditional flow-matching loss (Lipman et al., 2023):

$$\mathcal{L}(\theta) = \mathbb{E}\left[\|u_\theta(x(t), t) - v_t(x(t) \mid x(0), x(1))\|^2\right],$$

where the expectation is taken over $t \sim \mathrm{Uniform}[0, 1]$, $x(0) \sim p_0$, and $x(1) \sim p_1$.

During inference, sampling proceeds by drawing $x(0) \sim p_0$ and integrating the ODE $\frac{d}{dt} x(t) = u_\theta(x(t), t)$. generating final samples $x(1)$ that approximate $p_1$.

## 4 Method

The goal of CPFM is to construct semantically meaningful samplers for the conditional distributions $p(x \mid y)$ and $p(y \mid x)$, given a training dataset $\mathcal{D}_x = \{x_i\}_{i=1}^n$. As discussed, our approach proceeds in two stages.

In the *first stage*, we draw $\mathcal{D}_y = \{y_i\}_{i=1}^n$ from the marginal distribution $\mu_{\mathcal{Y}}$ and establish a probabilistic correspondence between $\mathcal{D}_x$ and $\mathcal{D}_y$ via a generalized Gromov–Wasserstein optimal transport (GWOT). This coupling incorporates user-specified priors through kernelized costs to enable controllable alignment between data and embeddings. In the *second stage*, we extend this coupling from the discrete training samples to the entire joint space $\mathcal{X} \times \mathcal{Y}$ by training a *dual conditional flow matching* model.
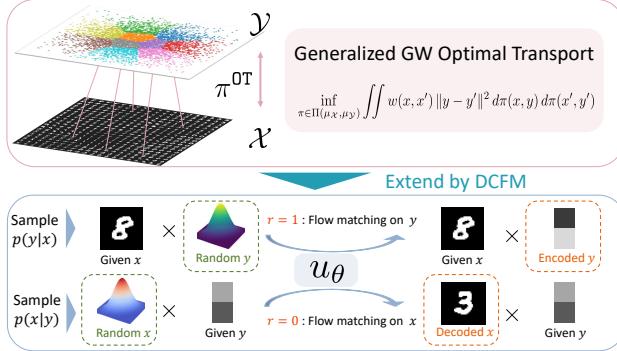
Figure 1: An overview of coupled flow matching.

The overall pipeline of CPFM is illustrated in Figure 1. The following subsections present the technical details of each component. For brevity, proofs of propositions and theorems are deferred to the appendix.

## 4.1 Generalized Gromov-Wasserstein Optimal Transport

Standard GWOT compares relational structures across domains but typically relies only on $\ell_p$ distances between raw data points. This makes it difficult to incorporate prior knowledge such as semantic labels, neighborhood graphs, or other task-relevant information. To overcome this limitation, we introduce a *generalized* formulation that leverages kernel functions to encode flexible structures in the data space $\mathcal{X}$.

Formally, given marginals $\mu_{\mathcal{X}}$ and $\mu_{\mathcal{Y}}$, we introduce the following optimization problem

$$\inf_{\pi \in \Pi(\mu_{\mathcal{X}}, \mu_{\mathcal{Y}})} \iint_{\mathcal{X}^2 \times \mathcal{Y}^2} k(x, x') \|y - y'\|^2 \, d\pi(x, y) \, d\pi(x', y'),$$
(1)

where $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a symmetric kernel that captures user-specified relational information in $\mathcal{X}$.

Formulation (1) encourages the semantic alignment between $(x, x')$ and $(y, y')$: when $k(x, x')$ is large, meaning $x$ is similar to $x'$, $\|y - y'\|$ should be small, meaning $y$ is similar to $y'$ in the latent space. The semantic information is directly encoded into the transport objective through the kernel function, thus enabling practical flexibility.

**Proposition 4.1.** *Standard GWOT is a special case of* (1) *with* $k(x, x') = -\|x - x'\|^2$.

**Finite-sample formulation.** In practice, we observe finite datasets $\mathcal{D}_x = \{x_i\}_{i=1}^n$ and $\mathcal{D}_y = \{y_j\}_{j=1}^n$, corresponding to empirical measures $\hat{\mu}_{\mathcal{X}} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$ and $\hat{\mu}_{\mathcal{Y}} = \frac{1}{n} \sum_{j=1}^n \delta_{y_j}$. With a slight abuse of notation, the coupling $\pi \in \Pi(\hat{\mu}_{\mathcal{X}}, \hat{\mu}_{\mathcal{Y}})$ can be represented as a matrix $\pi \in \mathbb{R}^{n \times n}$ with row-sum and column-sum

equal to $\frac{1}{n}$, where $\pi_{ij}$ denotes the transported mass from $x_i$ to $y_j$. The finite-sample counterpart to (1) becomes

$$\inf_{\pi \in \Pi(\hat{\mu}_{\mathcal{X}}, \hat{\mu}_{\mathcal{Y}})} \sum_{i,j,i',j'=1}^n \pi_{ij} \pi_{i'j'} \, k(x_i, x_{i'}) \, \|y_j - y_{j'}\|^2. \quad (2)$$

We denote its optimal solution by $\pi^{\text{OT}}$.

**Variational reformulation.** Although (2) seems difficult to optimize as it involves quadratic terms in entries of $\pi$, we show that it admits a tractable variational representation. Let $G \in \mathbb{R}^{n \times n}$ be the kernel Gram matrix with entries $G_{ij} = k(x_i, x_j)$. Assuming it has rank $m$, we can write its Cholesky factorization as $G = \Phi^\top \Phi$ with $\Phi \in \mathbb{R}^{m \times n}$. The following theorem proposes a reformulation of (2).

**Theorem 4.2.** *The optimal transport plan $\pi^{OT}$ is also the minimizer of*

$$\inf_{\pi \in \Pi(\hat{\mu}_{\mathcal{X}}, \hat{\mu}_{\mathcal{Y}})} \inf_{A \in \mathbb{R}^{m \times d_y}} \left\{ \|A\|^2 \right.$$
$$\left. + \sum_{i=1}^n \sum_{j=1}^n \pi_{ij} \left( \|y_j\|^2 w(x_i) - 2\langle A, \Phi_i y_j^\top \rangle \right) \right\},$$
(3)

*where $A$ is an auxiliary matrix, $\Phi_i$ is the $i$-th column of $\Phi$, $w(x_i) = \sum_{j=1}^n k(x_i, x_j)$, and the matrix inner product $\langle \cdot, \star \rangle$ is defined as the inner product of their flattened vectors.*

Theorem 4.2 transforms the quadratic coupling into a problem linear in $\pi$, at the cost of introducing an auxiliary variable $A$. This reformulation simplifies the optimization algorithm design.

**Alternating minimization.** Building on the variational form (3) of generalized GWOT, we develop an alternating optimization algorithm to obtain $\pi^{\text{OT}}$.

Step 1: Fix $\pi$, update $A$. For fixed $\pi$, the subproblem is a quadratic function of $A$, with the optimal solution

$$A^\star = \sum_{i=1}^n \sum_{j=1}^n \pi_{ij} \, \Phi_i y_j^\top.$$

Step 2: Fix $A$, update $\pi$. For fixed $A$, the cost reduces to a linear assignment problem

$$\min_{\pi \in \mathbb{R}^{n \times n}} \sum_{i=1}^n \sum_{j=1}^n \pi_{ij} \underbrace{\left( \|y_j\|^2 w(x_i) - 2\langle A, \Phi_i y_j^\top \rangle \right)}_{=: \, c(x_i, y_j)}. \quad (4)$$

We add an entropic regularization $\varepsilon \sum_{i=1}^n \sum_{j=1}^n \pi_{ij} (\log \pi_{ij} - 1)$ to ensure numerical

Wenxi Cai [1], Yuheng Wang [1], Naichen Shi [2]

**Algorithm 1:** Alternating optimization for (3)

---

**Input:** Kernel factors $\Phi \in \mathbb{R}^{n \times m}$, auxiliary
matrix $A$, embedding matrix $Y \in \mathbb{R}^{n \times d_y}$,
tolerance $\tau > 0$, regularization $\varepsilon > 0$

**Output:** transport plan $\pi \in \mathbb{R}^{n \times n}$, auxiliary
matrix $A$, OT objective $\mathcal{L}_{\mathrm{OT}}$

Initialize $\mathcal{L}_{\mathrm{OT}}^{(0)} \leftarrow +\infty$;

Compute $Y_{\mathrm{norm}} \leftarrow \left( \|y_1\|^2, \ldots, \|y_n\|^2 \right)^\top$

Compute $w \leftarrow \mathrm{Rowsum}(\Phi\Phi^\top)$

**while** *True* **do**

   $C \leftarrow w\, Y_{\mathrm{norm}}^\top - 2\,\Phi A Y^\top$ // cost matrix

   $\pi \leftarrow \mathrm{Sinkhorn}(C, \varepsilon)$

   $A \leftarrow \Phi^\top \pi Y$

   $\mathcal{L}_{\mathrm{OT}}^{(t)} \leftarrow \|A\|_F^2 + \langle C, \pi \rangle$ // current OT value

   **if** $\mathcal{L}_{\mathrm{OT}}^{(t-1)} - \mathcal{L}_{\mathrm{OT}}^{(t)} < \tau$ **then**

      ⌊ **break**

**return** $(\pi, A, \mathcal{L}_{\mathrm{OT}}^{(t)})$

---

tractability, where $\varepsilon > 0$ controls the strength of the regularization. The resulting regularized objective becomes an LEOT problem that can be solved by invoking Sinkhorn iterations:

$$\pi^\star \approx \mathrm{Sinkhorn}(c(x_i, y_j),\, \varepsilon).$$

Based on the alternating optimization scheme described above, we have the following proposition that certifies the convergence and per-iteration complexity of Algorithm 1.

**Proposition 4.3.** *Algorithm 1 produces a nonincreasing sequence of OT objectives $\{\mathcal{L}_{\mathrm{OT}}^{(t)}\}$ that converges to a stationary point, with per-iteration complexity $\mathcal{O}(n^2)$.*

**Practical considerations.** In practice, a small entropic regularization parameter $\varepsilon$ yields a more accurate approximation to the original (3), but slows down the convergence of the Sinkhorn subroutine. To address this, we develop an $\varepsilon$-scheduling scheme (Schmitzer, 2019). Specifically, we start from a larger $\varepsilon$, and then iteratively reduce it via bisection until reaching floating-point precision. This approach balances convergence speed and numerical stability. A detailed description of our $\varepsilon$-scheduling scheme is discussed in the supplementary material.

## 4.2 Dual Conditional Flow Matching

While $\pi^{\mathrm{OT}}$ provides a controllable correspondence between $\mathcal{D}_x$ and $\mathcal{D}_y$, it is defined only on finite samples. To generalize beyond the training support and construct conditional samplers on the entire space $\mathcal{X} \times \mathcal{Y}$, we introduce *Dual Conditional Flow Matching* (DCFM), the second pillar of CPFM.

**Dual conditional flows.** DCFM learns probability flows in both spaces.

To sample from $p(y|x)$ given $x$, we first draw $y(0) \sim p_0^y$ from a base distribution and integrate the ODE

$$\frac{d}{dt}y(t) = u_y(y(t); x, t)$$

to obtain $y(1)$, which serves as a sample from $p(y|x)$.

Conversely, to generate from $p(x|y)$ given an embedding $y$, we draw $x(0) \sim p_0^x$ and integrate

$$\frac{d}{dt}x(t) = u_x(x(t); y, t)$$

to obtain $x(1)$, which serves as a sample from $p(x \mid y)$.

The key observation is that $u_x$ and $u_y$ are structurally symmetric. We therefore parameterize both with a single neural network $u_\theta$, which outputs the drift for both $x$ and $y$. Specifically, we introduce a role flag $r \in \{0, 1\}$:

$$u_\theta(\cdot, r) = \begin{cases} u_x(\cdot), & r = 0, \\ u_y(\cdot), & r = 1. \end{cases}$$

In practice, $u_\theta$ is implemented as a shared backbone with two decoding heads, selected by $r$. This design leverages symmetry, avoids parameter redundancy, and enables joint parameter sharing.

**Training objective.** Training balances two objectives: learning $u_x$ and learning $u_y$. We use $a_t$ and $b_t$ to denote two differentiable functions over $t$ that satisfy the boundary condition specified in Section 3.3. We first sample $(x(1), y(1)) \sim \pi^{\mathrm{OT}}$ from the optimal coupling.

If $r = 0$, we generate an interpolant path for $x$: sample $x(0) \sim p_0^x$, set $x(t) = a_t x(0) + b_t x(1)$, and define velocity $v_x(t) = \dot{a}_t x(0) + \dot{b}_t x(1)$. The conditional flow-matching loss is

$$\ell_x(u, t, x(t), y(1), v_x(t)) = \|u(x(t), y(1), t, 0) - v_x(t)\|^2. \tag{5}$$

If $r = 1$, we instead interpolate $y$: sample $y(0) \sim p_0^y$, set $y(t) = a_t y(0) + b_t y(1)$, and define $v_y(t) = \dot{a}_t y(0) + \dot{b}_t y(1)$. The loss becomes

$$\ell_y(u, t, x(1), y(t), v_y(t)) = \|u(x(1), y(t), t, 1) - v_y(t)\|^2. \tag{6}$$

The full DCFM training objective is a weighted combination of the two roles:

$$\mathcal{L}_{\mathrm{DCFM}}(u) = (1 - \alpha)\, \mathbb{E}[\ell_x(u)] + \alpha\, \mathbb{E}[\ell_y(u)], \tag{7}$$

where $\alpha \in (0,1)$ balances the two objectives. Expectations are taken over random draws of $(t, x, y)$ and the role indicator $r \sim \text{Bernoulli}(\alpha)$.

**Algorithmic summary.** The stochastic training and sampling procedure is summarized in Algorithm 2.

---

**Algorithm 2:** Training DCFM by optimizing (7)

---

**Input:** Dataset $\mathcal{D}_x$, embedding dataset $\mathcal{D}_y$
       sampled from $\mu_{\mathcal{Y}}$, transport plan
       $\pi^{\text{OT}} \in \mathbb{R}^{|\mathcal{D}_x| \times |\mathcal{D}_x|}$, Bernoulli parameter $\alpha$

**Output:** trained drift model $u_\theta$

**while** *not converged* **do**
    Sample a pair $(x(1), y(1))$ based on $\pi^{\text{OT}}$
    Sample role $r \sim \text{Bernoulli}(\alpha)$ ;    // decide
    `train `$x$` or train `$y$`.`
    Sample time $t \sim \text{Uniform}[0, 1]$
    **if** $r = 0$ **then**
        Sample $x(0) \sim p_0^x$
        Compute interpolant $x(t) = a_t x(0) + b_t x(1)$
        Compute velocity $v_x(t) = \dot{a}_t x(0) + \dot{b}_t x(1)$
        Compute loss $\ell(u_\theta) = \ell_x(u_\theta)$ defined in (5)
    **else**
        Sample $y(0) \sim p_0^y$
        Compute interpolant $y(t) = a_t y(0) + b_t y(1)$
        Compute velocity $v_t^y = \dot{a}_t y(0) + \dot{b}_t y(1)$
        Compute loss $\ell(u_\theta) = \ell_y(u_\theta)$ defined in (6)
    Calculate $\nabla_\theta \ell(u_\theta)$ and update $\theta$ by AdamW.

---

We now argue that the two objectives in (7) are not in conflict. In fact, the joint loss admits a solution that simultaneously approximates both conditional flows.

**Theorem 4.4.** *Let $u^*$ denote the drift field corresponding to the optimal solution of (7). Then $u^*$ satisfies $u^*(z, c, t, 0) = \mathbb{E}[\dot{a}_t x(0) + \dot{b}_t x(1) \mid x(t) = z, y(1) = c]$, and $u^*(c, z, t, 1) = \mathbb{E}[\dot{a}_t y(0) + \dot{b}_t y(1) \mid y(t) = z, x(1) = c]$.*

This theorem shows that minimizing the combined loss (7) yields a drift field that is consistent with both conditional velocity fields. Thus, integrating the $\frac{d}{dt} x(t) = u^*(x(t), y(1), t, 0)$ or $\frac{d}{dt} y(t) = u^*(x(1), y(t), t, 1)$ produces samples from $p(y|x)$ and $p(x|y)$ without interference between the two directions.

The inference procedures are summarized in Algorithm 3.

## 5 Experiments

In this section, we evaluate CPFM on four image generation datasets and one molecule generation dataset to showcase its performance in both dimensionality reduction and sample reconstruction. The implementa-

---

**Algorithm 3:** Inference with DCFM

---

**Input:** trained drift model $u_\theta$, conditioning
       variable $c$ (either $x$ or $y$), direction
       $r \in \{0, 1\}$, number of integration steps $T$

**Output:** sample $\hat{z}$ in the target space

**if** $r = 0$ **then**
    Sample initial $x_0 \sim p_0^x$;
    Set $z_0 \leftarrow x_0$;
**else**
    Sample initial $y_0 \sim p_0^y$;
    Set $z_0 \leftarrow y_0$;
**for** $k = 0$ **to** $T - 1$ **do**
    Set $t = k/T$;
    Compute drift $u = u_\theta(z_k, c, t, r)$ if $r = 0$ or
    $u = u_\theta(c, z_k, t, r)$ if $r = 1$;
    Update state $z_{k+1} \leftarrow z_k + \frac{1}{T} \cdot u$ ;    // Euler
    `step`
Return $\hat{z} = z_T$;

---

tion and training code for DCFM are available in an anonymous repository: `https://anonymous.4open.science/r/AISTATS-CEF9/`. We use two NVIDIA A10 GPUs for training $u_\theta$.

For all experiments, we set the embedding space $\mathcal{Y}$ to be two-dimensional ($\mathbb{R}^2$). This choice enables direct visualization of the learned embeddings and provides a stringent test case: mapping high-dimensional data into such a severely compressed space poses a challenging scenario for both representation learning and faithful reconstruction.

### 5.1 MNIST

We begin with the MNIST dataset (Deng, 2012), a benchmark collection of handwritten digit images (0–9), to illustrate the effectiveness of CPFM.

**Kernel construction.** Since MNIST is labeled, we design a kernel function that incorporates both visual similarity and label information. Inspired by Conditional Random Field models (Krähenbühl and Koltun, 2011), we propose the composite kernel

$$k_{\text{image}}(x_i, x_j) = \underbrace{\exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)}_{\text{appearance kernel}} \cdot \underbrace{\mathbf{1}\{l_i = l_j\}}_{\text{label kernel}}, \quad (8)$$

where $l_i$ and $l_j$ are the digit labels of images $x_i$ and $x_j$, $\mathbf{1}\{\cdot\}$ is an indicator function, and $\sigma$ is the bandwidth parameter estimated as $\sigma = \frac{1}{n^2} \sum_{i,j=1}^{n} \|x_i - x_j\|$.

Here, the appearance kernel is a Gaussian (heat) kernel that captures pixel-level similarity, while the label kernel encodes semantic consistency by enforcing

Wenxi Cai [1], Yuheng Wang [1], Naichen Shi [2]

affinity only among images of the same digit. This construction ensures that the transport cost reflects both geometric similarity and label-based priors.

**Generalized GWOT.** With the kernel in (8), we compute the Gram matrix of the MNIST images, draw latent samples $y$ from $\mu_{\mathcal{Y}}$, and apply Algorithm 1 to solve the generalized GWOT problem. The resulting embeddings are plotted in Figure 2.
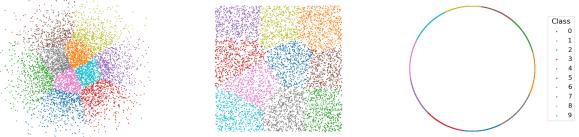


Figure 2: Based on the generalized GWOT transport plan $\pi$, each source sample $x$ is associated with a probability distribution over candidate embeddings $y$. A single embedding $y$ is randomly sampled according to its assigned weight $\pi_i$. The latent embedding distributions $\mu_{\mathcal{Y}}$ considered are: (Left) a Gaussian $\mathcal{N}(0, I_2)$; (Middle) a uniform distribution on the square $[-1, 1]^2$; and (Right) a uniform distribution on the unit circle $, y \in \mathbb{R}^2 : |y| = 1,$.

The generated embeddings align closely with the contour of the prescribed target distributions $\mu_{\mathcal{Y}}$, demonstrating that generalized GWOT can faithfully adapt to different latent marginals. Moreover, the embeddings exhibit well-separated clusters for different digit classes, reflecting the influence of the label-sensitive kernel. This separation indicates that Algorithm 1 could preserve geometric structure and semantic priors simultaneously.

**DCFM reconstruction.** Building on the transport plan provided by generalized GWOT, we train DCFM using Algorithm 2 to jointly perform dimensionality reduction into two dimensions and reconstruction back to the image space. After training, we sample 100 random images, embed them into $\mathbb{R}^2$, and then reconstruct from the embeddings using the learned drift network $u_\theta$. Figure 3 shows both the original and reconstructed images. For clarity of visualization, the figure displays 1024 embeddings in the latent space rather than 100.

From Figure 3, we observe that reconstructed digits preserve their class identity, underscoring DCFM's ability to faithfully retain label information encoded in the discrete transport plan $\pi^{\mathrm{OT}}$. At the same time, the reconstructions exhibit subtle variations in handwriting style compared to the originals. This suggests that DCFM not only captures class-relevant features
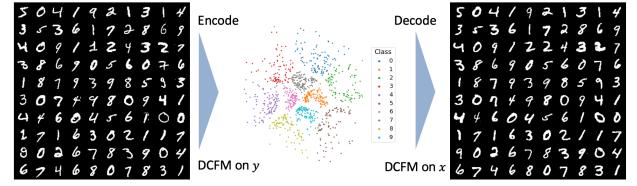


Figure 3: The MNIST dataset compressed into two dimensions by DCFM, where different classes are distinguished by labels and then reconstructed.

but also introduces diversity in nuisance attributes.

## 5.2 Image datasets

Beyond MNIST, we further evaluate CPFM on three widely used image benchmarks: **CIFAR-10** (Krizhevsky et al., 2009), **TinyImageNet** (Le and Yang, 2015), and **AFHQ** (Choi et al., 2020). Since we pretrain flow matching model $u_\theta$ from random initialization, we downsample the resolution of TinyImagenet and AFHQ to $64 \times 64$, and $128 \times 128$, respectively, to reduce the carbon footprint in the training process.

**Evaluation metrics** The primary evaluation metrics we use is the *Wasserstein distance of the latent distribution to Gaussian distribution.*

**Benchmarks.** We compare CPFM against a set of representative baselines, including KPCA (Gedon et al., 2023), VAE (Kingma and Welling, 2019), DiffAE (Preechakul et al., 2022), and Info-Diffusion (Wang et al., 2023). To ensure fairness, all diffusion and flow matching-based models are trained with the same network architecture and number of training iterations. Quantitative results are summarized in Table 1.

From Table 1, it is evident that CPFM achieves the lowest FID and OT loss, which again confirms its superior dimension reduction and sample reconstruction capability. It is worth noting that the reported FID scores are substantially higher than those in the original papers. This discrepancy is expected, as we deliberately impose an extreme dimension reduction setting by restricting the latent space to only two dimensions.

## 5.3 QM9

To further evaluate CPFM beyond image domains, we test it on the QM9 dataset (Blum and Reymond, 2009; Montavon et al., 2013), a widely used quantum chemistry benchmark. We randomly sample 4,000 molecules for training and 1,000 for evaluation. The molecule attribute of interest is the *dipole moment*

| Distance to Gaussian ↓ | MNIST | CIFAR-10 | AFHQ | TinyImagenet |
|---|---|---|---|---|
| KPCA | 245 ± 4 | 720 ± 60 | 48000 ± 5900 | 47000 ± 5000 |
| VAE | 32 ± 5 | 52 ± 7 | 48 ± 7 | 49 ± 8 |
| DiffAE | 4.96 ± 0.04 | 7.40 ± 0.02 | 28 ± 2 | 20.3 ± 0.2 |
| Info-Diffusion | 0.1608 ± 1e-06 | 0.3092 ± 2e-06 | 0.2996 ± 2e-07 | 0.253 ± 0.006 |
| **CPFM** | **0.113** ± 0.002 | **0.17** ± 0.01 | **0.138** ± 0.008 | **0.25** ± 0.03 |

Table 1: Distance to Gaussian. Mean ± std calculated over 5 independent generations.

(a continuous variable, measured in Debye). Each molecule is encoded with SELFIES (Krenn et al., 2020).A VAE is first trained to obtain 64-dimensional latent representations of molecules, which serve as $\mathcal{X}$.

**Kernel construction.** We design a composite kernel that integrates both structural similarity and property similarity:

$$k_{\mathrm{mol}}(i,j) = \frac{1}{2} \underbrace{\frac{\langle f_i, f_j \rangle}{\|f_i\|_1 + \|f_j\|_1 - \langle f_i, f_j \rangle}}_{\text{structure kernel}} + \frac{1}{2} \underbrace{|\ell_i - \ell_j|}_{\text{property kernel}} ,$$

$$(9)$$

where $f_i$ is the binary Morgan fingerprint (Rogers and Hahn, 2010) of molecule $i$, used to compute the Tanimoto similarity (Bajusz et al., 2015), and $\ell_i$ is its dipole moment (Bader et al., 1987).

**Dimension reduction.** We apply generalized GWOT with the kernel in (9), producing two-dimensional embeddings. The results are visualized in Figure 4. Compared with standard t-SNE, the embeddings obtained via Algorithm 1 exhibit smoother transitions in dipole moment values, indicating that the latent coordinates align more naturally with the underlying property. This highlights CPFM's ability to produce chemically meaningful embeddings.
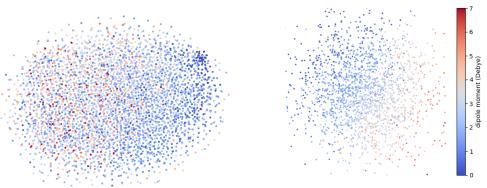


Figure 4: (Left) t-SNE visualization of the VAE latent space, where the labels do not show a clear separation. (Right) Embedding obtained from generalized GWOT, where the molecular embeddings vary continuously in label.

**Controlled generation.** Next, we train DCFM to learn a bidirectional mapping between the VAE latent space and the generalized GWOT embeddings. Figure 5 shows six examples of molecules and their reconstructions.
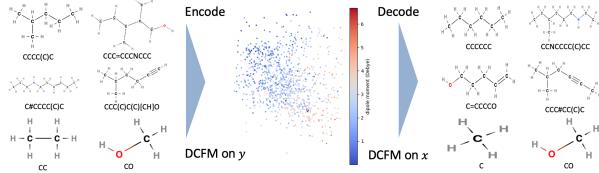


Figure 5: The QM9 dataset compressed into two dimensions by DCFM and then reconstructed.

The reconstructions largely preserve atomic composition and molecular topology, demonstrating CPFM's ability to retain chemically meaningful features even under drastic dimensionality reduction.

## 6 Discussion and Limitations

This work introduced CPFM, a framework for controllable dimensionality reduction and reconstruction. While effective, training dual conditional flows from scratch remains computationally intensive on high-resolution datasets, suggesting the need for more efficient strategies such as multi-scale architectures.

## References

Albergo, M. S., Boffi, N. M., and Vanden-Eijnden, E. (2023). Stochastic interpolants: A unifying framework for flows and diffusions.

Armstrong, G., Martino, C., Rahman, G., Gonzalez, A., Vázquez-Baeza, Y., Mishne, G., and Knight, R. (2021). Uniform manifold approximation and projection (umap) reveals composite patterns and resolves visualization artifacts in microbiome data. *mSystems*, 6(5):e00691–21.

Bader, R., Larouche, A., Gatti, C., Carroll, M., MacDougall, P., and Wiberg, K. (1987). Properties of atoms in molecules: Dipole moments and transferability of properties', j. chem. phys. 87, 1142-1152 (1987). *The Journal of chemical physics*, 87:1142–1152.

Bajusz, D., Rácz, A., and Héberger, K. (2015). Why is tanimoto index an appropriate choice for fingerprint-based similarity calculations? *Journal of Cheminformatics*, 7(20):1–13.

Bao, F., Nie, S., Xue, K., Li, C., Pu, S., Wang, Y., Yue, G., Cao, Y., Su, H., and Zhu, J. (2023). One transformer fits all distributions in multi-modal diffusion at scale.

Becht, E., McInnes, L., Healy, J., Dutertre, C.-A., Kwok, I. W., Ng, L. G., Ginhoux, F., and Newell, E. W. (2019). Dimensionality reduction for visualizing single-cell data using umap. *Nature Biotechnology*, 37(1):38–44.

Belkin, M. and Niyogi, P. (2003). Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Comput.*, 15(6):1373–1396.

Blum, L. C. and Reymond, J.-L. (2009). 970 million druglike small molecules for virtual screening in the chemical universe database GDB-13. *J. Am. Chem. Soc.*, 131:8732.

Calvo-Ordonez, S., Meunier, M., Cartea, A., Reisinger, C., Gal, Y., and Hernandez-Lobato, J. M. (2025). Weighted conditional flow matching.

Choi, Y., Uh, Y., Yoo, J., and Ha, J.-W. (2020). Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8188–8197.

Cuturi, M. (2013). Sinkhorn distances: Lightspeed computation of optimal transport. In Burges, C., Bottou, L., Welling, M., Ghahramani, Z., and Weinberger, K., editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.

Deng, L. (2012). The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE signal processing magazine*, 29(6):141–142.

Dumoulin, V., Belghazi, I., Poole, B., Mastropietro, O., Lamb, A., Arjovsky, M., and Courville, A. (2017). Adversarially learned inference.

Gedon, D., Ribeiro, A. H., Wahlström, N., and Schön, T. B. (2023). Invertible kernel pca with random fourier features. *IEEE Signal Processing Letters*, 30:563–567.

Generale, A. P., Robertson, A. E., and Kalidindi, S. R. (2025). Conditional variable flow matching: Transforming conditional densities with amortized conditional optimal transport.

Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial networks.

Guo, P. and Schwing, A. G. (2025). Variational rectified flow matching.

Hinton, G. E. and Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507.

Huang, X., Liu, M.-Y., Belongie, S., and Kautz, J. (2018). Multimodal unsupervised image-to-image translation.

Karras, Tero asnd Laine, S. and Aila, T. (2019). A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410.

Kingma, D. P. and Welling, M. (2019). An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392.

Kingma, D. P. and Welling, M. (2022). Auto-encoding variational bayes.

Kobak, D. and Berens, P. (2019). The art of using t-sne for single-cell transcriptomics. *Nature Communications*, 10(5416):1–14.

Krähenbühl, P. and Koltun, V. (2011). Efficient inference in fully connected crfs with gaussian edge potentials. *Advances in neural information processing systems*, 24.

Krenn, M., Häse, F., Nigam, A., Friederich, P., and Aspuru-Guzik, A. (2020). Self-referencing embedded strings (selfies): A 100 *Machine Learning: Science and Technology*, 1(4):045024.

Krizhevsky, A., Hinton, G., et al. (2009). Learning multiple layers of features from tiny images.

Le, Y. and Yang, X. (2015). Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3.

Lee, S., Park, J., Chu, J., Yoon, M., and Kim, H. J. (2025). Latent bayesian optimization via autoregressive normalizing flows. In *The Thirteenth International Conference on Learning Representations*.

Li, Z., Li, H., Shi, Y., Farimani, A. B., Kluger, Y., Yang, L., and Wang, P. (2025). Dual diffusion for unified image generation and understanding.

Lipman, Y., Chen, R. T. Q., Ben-Hamu, H., Nickel, M., and Le, M. (2023). Flow matching for generative modeling.

Liu, X., Gong, C., and Liu, Q. (2022). Flow straight and fast: Learning to generate and transfer data with rectified flow.

McInnes, L., Healy, J., and Melville, J. (2020). Umap: Uniform manifold approximation and projection for dimension reduction.

Montavon, G., Rupp, M., Gobre, V., Vazquez-Mayagoitia, A., Hansen, K., Tkatchenko, A., Müller, K.-R., and von Lilienfeld, O. A. (2013). Machine

learning of molecular electronic properties in chemical compound space. *New Journal of Physics*, 15(9):095003.

Peyré, G., Cuturi, M., and Solomon, J. (2016). Gromov-wasserstein averaging of kernel and distance matrices. In *Proceedings of the 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2664–2672, New York, NY, USA. PMLR.

Peyré, G. and Cuturi, M. (2020). Computational optimal transport.

Preechakul, K., Chatthee, N., Wizadwongsa, S., and Suwajanakorn, S. (2022). Diffusion autoencoders: Toward a meaningful and decodable representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10619–10629.

Rioux, G., Goldfeld, Z., and Kato, K. (2024). Entropic gromov-wasserstein distances: Stability and algorithms.

Rogers, D. and Hahn, M. (2010). Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling*, 50(5):742–754. PMID: 20426451.

Schmitzer, B. (2019). Stabilized sparse scaling algorithms for entropy regularized transport problems. *SIAM Journal on Scientific Computing*, 41(3):A1443–A1481.

Shlens, J. (2014). A tutorial on principal component analysis.

Tong, A., Fatras, K., Malkin, N., Huguet, G., Zhang, Y., Rector-Brooks, J., Wolf, G., and Bengio, Y. (2024). Improving and generalizing flow-based generative models with minibatch optimal transport.

van der Maaten, L. and Hinton, G. (2008). Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605.

Vayer, T., Chapel, L., Flamary, R., Tavenard, R., and Courty, N. (2018). Fused gromov-wasserstein distance for structured objects: theoretical foundations and mathematical properties.

Wang, R. and Zhang, Z. (2025). Quadratic-form optimal transport.

Wang, Y., Schiff, Y., Gokaslan, A., Pan, W., Wang, F., De Sa, C., and Kuleshov, V. (2023). Infodiffusion: Representation learning using information maximizing diffusion models. In *International conference on machine learning*, pages 36336–36354. PMLR.

Wu, M. and Goodman, N. (2018). Multimodal generative models for scalable weakly-supervised learning.

Zhao, X., Liu, B., Liu, Q., Shi, G., and Wu, X.-M. (2024). Easygen: Easing multimodal generation with bidiffuser and llms.

# Appendix

The supplementary materials are organized as follows. Section A elaborates on practical considerations of implementing CPFM. Section B presents the proofs to the theorems in the main paper. Section C introduces the numerical evaluation metrics in detail. Finally, Section D provides additional experiment results, including additional visualizations that illustrate the semantic coherence of the low-dimensional embeddings extracted by CPFM, visualizations of the DCFM architecture, and the choice of hyperparameters.

## A  Practical considerations

### A.1  $\varepsilon$-scheduling

We improve Algorithm 1 by adding an adaptive schedule for $\varepsilon$ to balance entropy magnitude and numerical stability. A larger $\varepsilon$ makes the Sinkhorn subroutine updates faster and more stable, but it increases bias towards blurry assignments in $\pi^{\texttt{OT}}$. A smaller $\varepsilon$ matches the unregularized quadratic objective better, but it may incur instability issues under floating-point arithmetic. Algorithm 4 keeps the inner iteration stable while improving outer iterations through warm starts. Specifically, it maintains an interval between the last stable value $\varepsilon$ and a trial value $\varepsilon_{\text{current}}$. If the inner call reaches the tolerance $\tau$ with no precision issues, we accept the iterate $(\pi, A)$ as a warm start and then halve $\varepsilon_{\text{current}}$. If a precision limit is detected, we backtrack by moving $\varepsilon_{\text{current}}$ to the midpoint between the failed value and the last stable value. The procedure stops when $|\varepsilon_{\text{current}} - \varepsilon| < \delta$ and we return the transport plan computed at the last stable $\varepsilon$.

---

**Algorithm 4:** Alternating optimization with adaptive entropic regularization

---

**Input:** kernel weights $w \in \mathbb{R}^n$, feature matrix $\Phi \in \mathbb{R}^{n \times m}$, embedding matrix $Y \in \mathbb{R}^{n \times q}$, tolerance $\tau > 0$,
        initial regularization $\varepsilon > 0$, regularization tolerance $\delta$
**Output:** transport plan $\pi \in \mathbb{R}^{n \times n}$, OT objective $\mathcal{L}_{\text{OT}}$
Initialize $\varepsilon_{\text{current}} \leftarrow \varepsilon$;
Initialize $A \leftarrow 0$;
**while** *True* **do**
    **if** *Algorithm 1 $(w, \Phi, A, Y, \tau, \varepsilon_{current})$ reaches precision limitation* **then**
        $\varepsilon_{\text{current}}, \varepsilon \leftarrow \frac{\varepsilon + \varepsilon_{\text{current}}}{2}, \varepsilon_{\text{current}}$;    // bi-section search to determine the smallest stable $\varepsilon$.
    **else**
        $\pi, A, \mathcal{L}_{\text{OT}} \leftarrow$ Algorithm 1$(w, \Phi, A, Y, \tau, \varepsilon_{\text{current}})$;
        $\varepsilon_{\text{current}}, \varepsilon \leftarrow \frac{\varepsilon_{\text{current}}}{2}, \varepsilon_{\text{current}}$;                           // halve $\varepsilon$.
    **if** $|\varepsilon_{current} - \varepsilon| < \delta$ **then**
        break;
**return** $(\pi, \mathcal{L}_{\text{OT}}^{(t)})$

---

In all experiments, we implement the generalized GWOT optimization with Algorithm 4 and set the initial $\varepsilon$ to 0.01.

### A.2  Choice of $m$

In the variational reformulation, we compute the feature map $\Phi$ using the Pivoted Cholesky feature map. Given the kernel Gram matrix $G = [k(x_i, x_j)]_{i,j=1}^n$, the algorithm compute the feature vector $\Phi_i \in \mathbb{R}^m$ for each sample such that

$$\Phi^\top \Phi \approx G,$$

where $\Phi = [\Phi_1, \ldots, \Phi_n]$.

Since $n$ can be huge in practice, Pivoted Cholesky factorization automatically selects rank $m$ to reach a target level of explained variance ratio $\eta \in [0,1]$. In all experiments, we choose the target explained variance ratio $\eta = 0.95$.

### A.3 Interpolation of $\pi^{\text{OT}}$

Despite the variational formulation introduced, Algorithm 1 still incurs $O(n^2)$ per-iteration complexity. This becomes infeasible in our computational infrastructure when $n \geq 10{,}000$. We therefore subsample 10,000 points from the entire dataset, compute $\pi^{\text{OT}}$ on this subset, and interpolate the plan to the full set to obtain $\widehat{\pi^{\text{OT}}}$. For a point $x$, let $\mathcal{N}_k(x)$ be its $k$ same-class neighbors in the subset and define weights

$$\alpha_j(x) = \frac{\exp\big(-\|x - x_j\|_2\big)}{\sum_{m \in \mathcal{N}_k(x)} \exp\big(-\|x - x_m\|_2\big)}, \quad j \in \mathcal{N}_k(x).$$

The row-interpolated transport plan is

$$\widehat{\pi^{\text{OT}}}(x, \cdot) = \sum_{j \in \mathcal{N}_k(x)} \alpha_j(x)\, \pi^{\text{OT}}(x_j, \cdot).$$

During DCFM training, we sample from a mixture of $\pi^{\text{OT}}$ and the interpolated plan $\widehat{\pi^{\text{OT}}}$.

## B  Missing proofs

We define $w_i = \sum_{i'=1}^{n} k(x_i, x_{i'})\, \hat{\mu}_{\mathcal{X}}(x_{i'})$, and $k(x_i, x_{i'})$ can be written as $k(x_i, x_j) = \langle \Phi_i, \Phi_j \rangle_{\mathbb{R}^m}$ based on $G = \Phi^{\top}\Phi$. This can be considered as the discrete version of decomposition in a Hilbert space $\mathcal{H}$:

$$k(x_i, x_{i'}) = \langle \varphi(x_i), \varphi(x_{i'}) \rangle_{\mathcal{H}}$$

for some feature map $\varphi : \mathcal{X} \to \mathcal{H}$.

**Theorem B.1.** *We have the following equivalent formulation of the optimal transport problem:*

$$\inf_{\pi \in \Pi(\hat{\mu}_{\mathcal{X}}, \hat{\mu}_{\mathcal{Y}})} \sum_{i,j,i',j'=1}^{n} \pi_{ij} \pi_{i'j'} \, k(x_i, x_{i'}) \, \|y_j - y_{j'}\|^2.$$

$$= 2 \inf_{\pi \in \Pi(\hat{\mu}_{\mathcal{X}}, \hat{\mu}_{\mathcal{Y}})} \inf_{A \in \mathbb{R}^{m \times d_y}} \left\{ \|A\|^2 \right. \tag{10}$$

$$\left. + \sum_{i=1}^{n} \sum_{j=1}^{n} \pi_{ij} \Big( \|y_j\|^2 w(x_i) - 2\langle A, \Phi_i y_j^{\top} \rangle \Big) \right\}.$$

*Proof.* We decompose the quadratic cost as

$$\sum_{i,i'=1}^{n} \sum_{j,j'=1}^{n} k(x_i, x_{i'}) \, \|y_j - y_{j'}\|^2 \, \pi_{ij}\, \pi_{i'j'} = I_1 + I_2 - 2I_3, \tag{11}$$

where

$$I_1 = \sum_{i,i'=1}^{n} \sum_{j,j'=1}^{n} k(x_i, x_{i'}) \, \|y_j\|^2 \, \pi_{ij}\, \pi_{i'j'}, \tag{12}$$

$$I_2 = \sum_{i,i'=1}^{n} \sum_{j,j'=1}^{n} k(x_i, x_{i'}) \, \|y_{j'}\|^2 \, \pi_{ij}\, \pi_{i'j'}, \tag{13}$$

$$I_3 = \sum_{i,i'=1}^{n} \sum_{j,j'=1}^{n} k(x_i, x_{i'}) \, \langle y_j, y_{j'} \rangle \, \pi_{ij}\, \pi_{i'j'}. \tag{14}$$

**Wenxi Cai** [1], **Yuheng Wang** [1], **Naichen Shi** [2]

For $I_1$, using the definition of $w_i = \sum_{i'} k(x_i, x_{i'})\, \hat{\mu}_{\mathcal{X}}(x_{i'})$,

$$I_1 = \sum_{i=1}^{n} \sum_{j=1}^{n} w_i \, \|y_j\|^2 \, \pi_{ij}. \tag{15}$$

Similarly,

$$I_2 = \sum_{i'=1}^{n} \sum_{j'=1}^{n} w_{i'} \, \|y_{j'}\|^2 \, \pi_{i'j'}. \tag{16}$$

For $I_3$, using the kernel representation $k(x_i, x_{i'}) = \langle \Phi(x_i), \Phi(x_{i'})\rangle_{\mathbb{R}^m}$, we obtain

$$I_3 = \sum_{i,i'=1}^{n} \sum_{j,j'=1}^{n} \langle \Phi(x_i), \Phi(x_{i'})\rangle_{\mathbb{R}^m} \, \langle y_j, y_{j'}\rangle \, \pi_{ij}\, \pi_{i'j'}. \tag{17}$$

Let $m_i = \sum_{j=1}^{n} \frac{y_j\, \pi_{ij}}{\hat{\mu}_{\mathcal{X}}(x_i)}$, then

$$I_3 = \left\| \sum_{i=1}^{n} \Phi(x_i) \times m_i \, \hat{\mu}_{\mathcal{X}}(x_i) \right\|^2_{\mathbb{R}^m \times \mathbb{R}^{d_y}}. \tag{18}$$

Applying the identity

$$-\|z\|^2 = \inf_{A \in \mathbb{R}^m \times \mathbb{R}^{d_y}} \left\{ \|A\|^2_{\mathbb{R}^m \times \mathbb{R}^{d_y}} - 2\langle A, z\rangle_{\mathbb{R}^m \times \mathbb{R}^{d_y}} \right\},$$

to (18), we obtain

$$-I_3 = \inf_{A \in \mathbb{R}^m \times \mathbb{R}^{d_y}} \left\{ \|A\|^2_{\mathbb{R}^m \times \mathbb{R}^{d_y}} - 2\sum_{i=1}^{n} \sum_{j=1}^{n} \langle A, \Phi(x_i) \times y_j\rangle_{\mathbb{R}^m \times \mathbb{R}^{d_y}} \, \pi_{ij} \right\}. \tag{19}$$

Combining (15), (16), and (19), the overall problem admits the equivalent formulation

$$\inf_{\pi \in \Pi(\hat{\mu}_{\mathcal{X}}, \hat{\mu}_{\mathcal{Y}})} \; \inf_{A \in \mathbb{R}^m \times \mathbb{R}^{d_y}} 2\left\{ \|A\|^2_{\mathbb{R}^m \times \mathbb{R}^{d_y}} + \sum_{i=1}^{n} \sum_{j=1}^{n} \left( \|y_j\|^2 \, w_i - 2\langle A, \Phi(x_i) \times y_j\rangle_{\mathbb{R}^m \times \mathbb{R}^{d_y}} \right) \pi_{ij} \right\},$$

. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad$ $\square$

**Proposition B.2.** *Algorithm 1 produces a non-increasing sequence of OT objectives $\{\mathcal{L}_{\mathrm{OT}}^{(t)}\}$ that converges to a stationary point, with per-iteration complexity $\mathcal{O}(n^2)$.*

*Proof.* We interpret Algorithm 1 as minimizing the entropy-regularized objective:

$$\mathcal{L}_{\mathrm{OT}}(\pi, A) = \|A\|_F^2 + \langle C(A), \pi\rangle + \varepsilon \sum_{i,j} \left( \pi_{ij} \log \pi_{ij} - \pi_{ij} \right),$$

where $C(A) = w\, Y_{\mathrm{norm}}^{\top} - 2\, \Phi A Y^{\top}$. At each iteration, the algorithm performs: (i) an exact update $A^{(t+1)} = \arg\min_A \mathcal{L}_{\mathrm{OT}}(\pi^{(t)}, A)$ (closed form), and (ii) an exact update $\pi^{(t+1)} = \arg\min_{\pi \in \Pi} \mathcal{L}_{\mathrm{OT}}(\pi, A^{(t+1)})$ (Sinkhorn run to convergence).

**Non-increasing.** Exact minimization in each block implies $\mathcal{L}_{\mathrm{OT}}(\pi^{(t)}, A^{(t+1)}) \leq \mathcal{L}_{\mathrm{OT}}(\pi^{(t)}, A^{(t)})$ and $\mathcal{L}_{\mathrm{OT}}(\pi^{(t+1)}, A^{(t+1)}) \leq \mathcal{L}_{\mathrm{OT}}(\pi^{(t)}, A^{(t+1)})$. Hence $\{\mathcal{L}_{\mathrm{OT}}^{(t)}\}$ is non-increasing.

**Lower boundedness.** Fix $\pi \in \Pi$. The $A$-subproblem is

$$\min_A \; \|A\|_F^2 - 2\langle A, \Phi^{\top} \pi Y\rangle + \underbrace{\langle w\, Y_{\mathrm{norm}}^{\top}, \pi\rangle + \varepsilon \sum_{i,j} \left( \pi_{ij} \log \pi_{ij} - \pi_{ij} \right)}_{\text{independent of } A}.$$

Its gradient and Hessian are

$$\nabla_A = 2A - 2\Phi^\top \pi Y, \qquad \nabla_A^2 = 2I \succ 0,$$

So the subproblem is strictly convex and has a unique minimizer

$$A^\star(\pi) = \Phi^\top \pi Y.$$

Plugging this back gives the reduced objective.

$$\widetilde{\mathcal{L}}_{\mathrm{OT}}(\pi) := -\left\|\Phi^\top \pi Y\right\|_F^2 + \left\langle w\, Y_{\mathrm{norm}}^\top, \pi \right\rangle + \varepsilon \sum_{i,j}\left(\pi_{ij} \log \pi_{ij} - \pi_{ij}\right),$$

where we used $\|A\|_F^2 - 2\langle A, \Phi^\top \pi Y\rangle = \|A - \Phi^\top \pi Y\|_F^2 - \|\Phi^\top \pi Y\|_F^2$.

Since the domain of $\pi$ is bounded and $\pi \mapsto \Phi^\top \pi Y$ is continuous, $-\|\Phi^\top \pi Y\|_F$ is lower-bounded. Similarly, $\langle w\, Y_{\mathrm{norm}}^\top, \pi \rangle$ is lower-bounded. For the entropy term, by Jensen's inequality,

$$\sum_{i,j} \pi_{ij} \log \pi_{ij} \;\geq\; -\log(n^2),$$

and $-\sum_{i,j} \pi_{ij} = -1$. Therefore

$$\varepsilon \sum_{i,j}\left(\pi_{ij} \log \pi_{ij} - \pi_{ij}\right) \;\geq\; -\varepsilon\big(2\log n + 1\big).$$

As a result, $\widetilde{\mathcal{L}}_{\mathrm{OT}}(\pi)$ is lower bounded and consequently, so is $\mathcal{L}_{\mathrm{OT}}(\pi, A)$.

**Stationary points.** The $A$-subproblem is a strongly convex and therefore has the unique minimizer $A^\star(\pi) = \Phi^\top \pi Y$. The $\pi$-subproblem is strictly convex because $\varepsilon > 0$ and therefore has a unique minimizer characterized by the KKT system, that is, there exist dual scalings $u, v$ such that $\log \pi^\star = -C(A^\star)/\varepsilon + u\mathbf{1}^\top + \mathbf{1}v^\top$ together with the row and column constraints (Cuturi, 2013). Standard results on two-block coordinate descent then imply that any limit point $(\pi^\star, A^\star)$ is block-optimal and satisfies the first-order conditions

$$\nabla_A \mathcal{L}_{\mathrm{OT}}(\pi^\star, A^\star) = 2A^\star - 2\Phi^\top \pi^\star Y = 0,$$

as well as the KKT conditions for $\pi^\star$. Consequently, any limit point is stationary.

**Computational complexity.** For the per-iteration complexity, forming $C(A) = wY_{\mathrm{norm}}^\top - 2\,\Phi A Y^\top$ costs $\mathcal{O}(n^2 + nmd_y)$, updating $A = \Phi^\top \pi Y$ costs $\mathcal{O}(nmd_y)$, and one Sinkhorn iteration costs $\mathcal{O}(n^2)$. With $T$ inner iterations, one outer iteration costs $\mathcal{O}\big((1 + T)n^2 + nmd_y\big)$, where $T$ is the number of iterations. When $T$ is bounded ($T$ is usually $\sim 15$ in our experiments), the computational complexity is $\mathcal{O}(n^2)$. $\qquad\square$

**Proposition B.3.** *Standard GWOT is a special case of generalized GWOT with $k(x, x') = -\|x - x'\|^2$ and can also be optimized via Algorithm 1.*

*Proof.* Consider the squared-loss Gromov–Wasserstein (GW) objective with Euclidean dissimilarities:

$$\mathrm{GW}(\pi) = \sum_{i,i'=1}^{n} \sum_{j,j'=1}^{n} \left(\, \|x_i - x_{i'}\|^2 - \|y_j - y_{j'}\|^2 \,\right)^2 \pi_{ij}\, \pi_{i'j'}. \tag{20}$$

Expanding the square yields

$$\begin{aligned}
\mathrm{GW}(\pi) = \sum_{i,i'=1}^{n} \|x_i - x_{i'}\|^4\, \hat{\mu}_{\mathcal{X}}(x_i)\, \hat{\mu}_{\mathcal{X}}(x_{i'}) &+ \sum_{j,j'=1}^{n} \|y_j - y_{j'}\|^4\, \hat{\mu}_{\mathcal{Y}}(y_j)\, \hat{\mu}_{\mathcal{Y}}(y_{j'}) \\
&- 2\sum_{i,i'=1}^{n} \sum_{j,j'=1}^{n} \|x_i - x_{i'}\|^2 \|y_j - y_{j'}\|^2\, \pi_{ij}\, \pi_{i'j'}.
\end{aligned} \tag{21}$$

Wenxi Cai [1], Yuheng Wang [1], Naichen Shi [2]

The first two terms in (21) depend only on the marginals and are therefore independent of the coupling $\pi$. They can be discarded without affecting the minimizer. Thus the minimization of $\mathrm{GW}(\pi)$ is equivalent to

$$\arg\min_{\pi} \mathrm{GW}(\pi) = \arg\min_{\pi}\left\{-2\sum_{i,i'=1}^{n}\sum_{j,j'=1}^{n}\|x_i - x_{i'}\|^2\|y_j - y_{j'}\|^2\,\pi_{ij}\,\pi_{i'j'}\right\}. \tag{22}$$

Define the kernel weight

$$k'(x_i, x_{i'}) := \|x_i - x_{i'}\|^2. \tag{23}$$

Then (22) can be written as

$$\arg\min_{\pi} -\sum_{i,i'=1}^{n}\sum_{j,j'=1}^{n} k'(x_i, x_{i'})\,\|y_j - y_{j'}\|^2\,\pi_{ij}\,\pi_{i'j'}, \tag{24}$$

which conforms with our objective (2).

This naturally leads to a variant of Theorem B.1:

$$\frac{1}{2}\inf_{\pi\in\Pi(\hat{\mu}_{\mathcal{X}},\hat{\mu}_{\mathcal{Y}})} -\sum_{i,i'=1}^{n}\sum_{j,j'=1}^{n} k'(x_i,x_{i'})\,\|y_j - y_{j'}\|^2\,\pi_{ij}\,\pi_{i'j'}$$

$$= \inf_{\pi\in\Pi(\hat{\mu}_{\mathcal{X}},\hat{\mu}_{\mathcal{Y}})}\inf_{A\in\mathbb{R}^{m\times\mathbb{R}^{d_y}}}\left\{\|A\|_{\mathbb{R}^m\times\mathbb{R}^{d_y}}^2 -\sum_{i=1}^{n}\sum_{j=1}^{n}\|y_j\|^2\,w_i\,\pi_{ij} + 2\sum_{i=1}^{n}\sum_{j=1}^{n}\langle A, \Phi(x_i)\times y_j\rangle_{\mathbb{R}^m\times\mathbb{R}^{d_y}}\,\pi_{ij}\right\}. \tag{25}$$

*where $K' \in \mathbb{R}^{n\times n}$ is the Gram matrix with $K'_{ii'} = k'(x_i, x_{i'})$ admitting a factorization $K' = \Phi\Phi^\top$ for some $\Phi \in \mathbb{R}^{n\times m}$ whose $i$-th column equals $\Phi(x_i)^\top$, and $w_i := \sum_{i'=1}^{n} k'(x_i, x_{i'})\,\hat{\mu}_{\mathcal{X}}(x_{i'})$.*

The remaining steps of the algorithm proceed exactly as before: update $\pi$ via the (entropically regularized) Sinkhorn step, and update $A$ in closed form by $A \leftarrow \Phi^\top\pi Y$. The monotone decrease argument and the $\mathcal{O}(n^2)$ per-iteration Complexity carries over verbatim.

$\square$

**Theorem B.4.** *Let $u^*$ denote the drift field corresponding to the optimal solution of (7). Then $u^*$ satisfies $u^*(z, c, t, 0) = \mathbb{E}[\dot{a}_t x(0) + \dot{b}_t x(1) \mid x(t) = z, y(1) = c]$, and $u^*(c, z, t, 1) = \mathbb{E}[\dot{a}_t y(0) + \dot{b}_t y(1) \mid y(t) = z, x(1) = c]$.*

*Proof.* We work at the population level. Let $r \in \{0, 1\}$ denote the role indicator ($r = 0$ for the $x$-direction, $r = 1$ for the $y$-direction), and define, for each training draw, the tuple.

$$(Z, C, V, t) = \begin{cases} (x(t),\ y(1),\ \dot{a}_t x(0) + \dot{b}_t x(1),\ t), & r = 0, \\ (y(t),\ x(1),\ \dot{a}_t y(0) + \dot{b}_t y(1),\ t), & r = 1, \end{cases}$$

where $(x(1), y(1)) \sim \pi^{\mathrm{OT}}$, $x(0) \sim p_0^x$, $y(0) \sim p_0^y$, and $x(t) = a_t x(0) + b_t x(1)$, $y(t) = a_t y(0) + b_t y(1)$ with differentiable $a_t, b_t$.

The DCFM population loss is

$$\mathcal{L}(u) = \mathbb{E}\big[\|u(Z, C, t, r) - V\|^2\big],$$

where the expectation is taken over the joint law of $(Z, C, r, V, t)$. Set the $\sigma$-algebra

$$\mathcal{A} := \sigma(Z, C, r, t).$$

The random vector $U := u(Z, C, t, r)$ is thus $\mathcal{A}$-measurable and square-integrable. By the orthogonal projection identity,

$$\mathbb{E}\big[\|V - U\|^2\big] = \mathbb{E}\big[\|V - \mathbb{E}(V \mid \mathcal{A})\|^2\big] + \mathbb{E}\big[\|U - \mathbb{E}(V \mid \mathcal{A})\|^2\big]. \tag{26}$$

The rightmost term is minimized when $U = \mathbb{E}(V \mid \mathcal{A})$, hence

$$u^*(Z, C, t, r) = \mathbb{E}[V \mid Z, C, r, t].$$

Thus, by specifying the conditional expectation pointwise, we obtain the role-wise form

$$u^*(z, c, t, r) = \mathbb{E}[V \mid Z = z, \ C = c, \ r, \ t]. \tag{27}$$

Unpacking the definitions of $(Z, C, V)$ under $r = 0$ and $r = 1$ gives exactly the two cases stated in the theorem:

$$u^*(z, c, t, 0) = \mathbb{E}[\dot{a}_t x(0) + \dot{b}_t x(1) \mid x(t) = z, \ y(1) = c, \ t], \qquad u^*(c, z, t, 1) = \mathbb{E}[\dot{a}_t y(0) + \dot{b}_t y(1) \mid y(t) = z, \ x(1) = c, \ t],$$

and the explicit dependence on $t$ can be suppressed in notation when no ambiguity arises.

Finally, taking total expectation in (26) and conditioning on the events $\{r = 0\}$ and $\{r = 1\}$ shows the loss decomposes additively across directions:

$$\begin{aligned} \mathcal{L}(u) - \mathcal{L}(u^*) &= \mathbb{E}\big[\|u(Z, C, t, r) - u^*(Z, C, t, r)\|^2\big] \\ &= \mathbb{E}\big[\|u(Z, C, t, 0) - u^*(Z, C, t, 0)\|^2 \mid r = 0\big] \, \mathbb{P}(r = 0) \\ &\quad + \mathbb{E}\big[\|u(Z, C, t, 1) - u^*(Z, C, t, 1)\|^2 \mid r = 1\big] \, \mathbb{P}(r = 1). \end{aligned}$$

$\square$

## C Evaluation Metrics

**Notation.** For each evaluation run, we subsample $N$ samples $\{x_i\}_{i=1}^N$ from the testing set, and generate their embedding vectors $\{y_i\}_{i=1}^N$. Then we reconstruct $\{\hat{x}_i\}_{i=1}^N$ using the trained $u_\theta$. Unless stated otherwise, metrics are reported as mean $\pm$ standard deviation aggregated over $R$ independent runs.

**Standardization.** We apply standardization procedures before computing image-based evaluation metrics: (i) if needed, single-channel images are broadcast to 3 channels; (ii) we restore dataset statistics using per-dataset $(\mu, \sigma)$ and clamp to $[0, 1]$ for FID; (iii) for LPIPS, we map images to $[-1, 1]$ after restoration. Distances in pixel space use Euclidean distance on flattened images.

### C.1 Fréchet Inception Distance (FID)

Let $\psi(\cdot) \in \mathbb{R}^{2048}$ denote Inception-V3 pool features applied to images after standardization. Let $\{\psi(x_i)\}_{i=1}^N$ and $\{\psi(\hat{x}_i)\}_{i=1}^N$ have empirical means and covariances $(\mu_r, \Sigma_r)$ and $(\mu_g, \Sigma_g)$, respectively. FID is the Fréchet distance between Gaussian approximations to these feature distributions:

$$\text{FID} = \|\mu_r - \mu_g\|_2^2 + \text{Tr}\Big(\Sigma_r + \Sigma_g - 2(\Sigma_r^{1/2} \Sigma_g \Sigma_r^{1/2})^{1/2}\Big). \tag{28}$$

The FID score is a widely adopted metric to evaluate the quality of generated samples.

### C.2 Wasserstein distance of embeddings

Since the objective of the dimension reduction algorithm is to match $\hat{\mu}_{\mathcal{X}}$ to a tractable distribution $\hat{\mu}_{\mathcal{Y}}$, it is also important to measure the distance between $\hat{\mu}_{\mathcal{Y}}$ and the target Gaussian distribution. We use the Wasserstein OT distance to characterize such a distributional distance. Let $Y = \{y_j\}_{j=1}^n$ be the embeddings generated by DCFM. We draw the same number of points $Z = \{z_j\}_{j=1}^n$ from the standard Gaussian $\mathcal{N}(0, I_2)$. Treating both sets as uniform discrete measures, we compute the 2-Wasserstein OT loss between $Y$ and $Z$ with squared Euclidean cost. In practice, we use a Sinkhorn solver with a small entropic regularization:

$$\text{WOT}_\varepsilon(Y, Z) = \min_{\gamma \in \mathbb{R}_+^{n \times n}} \sum_{i,j} \gamma_{ij} \|y_i - z_j\|^2 + \varepsilon \sum_{i,j} \gamma_{ij}(\log \gamma_{ij} - 1) \quad \text{s.t.} \quad \sum_j \gamma_{ij} = \tfrac{1}{n}, \ \sum_i \gamma_{ij} = \tfrac{1}{n}.$$

The reported value is the resulting OT objective, which quantifies how close the inferred embeddings are to the Gaussian reference. Values of FID and OT are reported in Table 1 in the main paper.

## C.3 Generalized GWOT objective

The original generalized GWOT OT loss is

$$\mathcal{L}(\pi) = \sum_{i,i',j,j'=1}^{n} \pi_{ij}\,\pi_{i'j'}\,k(x_i, x_{i'})\,\|y_j - y_{j'}\|^2. \tag{29}$$

In practice, the dimension reduction algorithms would generate a set of discrete samples $Y$ based on $X$. Therefore, we define a binary transport plan that pairs $x_i$ with its inferred $y_i$, i.e.,

$$\pi_{ij} = \tfrac{1}{n}\mathbf{1}\{i \text{ matches } j\}.$$

Substituting this into (29) gives

$$\text{GWOT} = \sum_{i,i',j,j'} \tfrac{1}{n}\mathbf{1}\{i \text{ matches } j\} \cdot \tfrac{1}{n}\mathbf{1}\{i' \text{ matches } j'\}\, k(x_i, x_{i'})\,\|y_j - y_{j'}\|^2 = \frac{1}{n^2}\sum_{i,j=1}^{n} k(x_i, x_j)\,\|y_i - y_j\|^2.$$

This objective measures how well the matching between the original samples $x$ and the embeddings $y$ align with user-defined objective.

## C.4 Aggregation Across Runs

Let $m^{(r)}$ be a scalar metric measured on run $r \in \{1, \ldots, R\}$. We report

$$\overline{m} \;=\; \frac{1}{R}\sum_{r=1}^{R} m^{(r)}, \qquad s(m) \;=\; \sqrt{\frac{1}{R-1}\sum_{r=1}^{R}\left(m^{(r)} - \overline{m}\right)^2}. \tag{30}$$

In all tables, the best results are **bolded**, while the second-best results are underlined.

## C.5 OT Loss (Generalized GWOT)

The main paper reports FID and OT loss (Wasserstein). Here we report OT loss (Generalized GWOT).

|  | MNIST | CIFAR-10 | AFHQ | TinyImagenet |
|---|---|---|---|---|
| KPCA | $23.5 \pm 0.1$ | $\underline{42 \pm 1}$ | $266 \pm 2$ | $366 \pm 3$ |
| VAE | $\underline{15.8 \pm 0.2}$ | $56 \pm 2$ | $\underline{67 \pm 1}$ | $\underline{60 \pm 9}$ |
| DiffAE | $370 \pm 1$ | $383 \pm 5$ | $398 \pm 14$ | $389 \pm 11$ |
| Info-Diffusion | $71 \pm 2$ | $298 \pm 4$ | $305 \pm 2$ | $298 \pm 3$ |
| **Coupled Flow Matching** | $\mathbf{1.393 \pm 0.001}$ | $\mathbf{1.40 \pm 0.02}$ | $\mathbf{1.34 \pm 0.03}$ | $\mathbf{1.82 \pm 0.03}$ |

Table 2: GWOT

# D Additional experiments

In this section, we provide additional experiment results and visualizations.

## D.1 Generalized GWOT on CIFAR-10, AFHQ, and TinyImageNet

We report the performance of Algorithm 4 (generalized GWOT solver) on CIFAR-10, AFHQ, and TinyImageNet, shown in Figure 6.

In Figure 6, embeddings from the same labels are clearly clustered, even for TinyImagenet with 200 classes. This highlights the effectiveness of Algorithm 4 to impose user-defined structures.
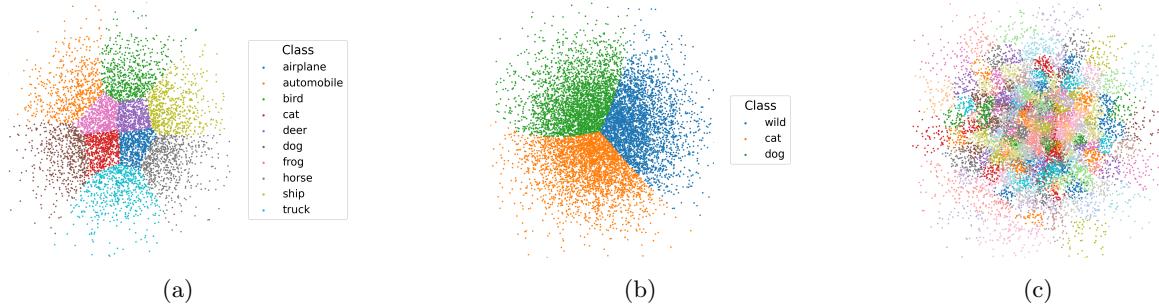
Figure 6: Based on the generalized GWOT transport plan $\pi$, each source sample $x$ is associated with a probability distribution over candidate embeddings $y$. A single embedding $y$ is randomly sampled according to its assigned weight $\pi_i$. The latent embedding distributions $\mu_{\mathcal{Y}}$ are Gaussian. The datasets used here are (a) CIFAR-10 (10 classes), (b)AFHQ (3 classes), and (c)TinyImageNet (200 classes).

## D.2 Extrapolation

To test whether CPFM internalizes the semantic structure generated by generalized GWOT from Stage 1, we sample a uniform $10 \times 10$ grid on the 2D square $[-1.5, 1.5]^2$ in $\mathcal{Y}$. Let $K = 10$ and $\Delta = 3/(K-1) = 1/3$, define $s_u = -1.5 + u\Delta$ for $u \in \{0, \ldots, 9\}$, and the grid points $y_{u,v} = (s_u, s_v)$ for $u, v \in \{0, \ldots, 9\}$. We treat each $y_{u,v}$ as the embedding and generate an image with *DCFM condition on y*, yielding samples $\{x_{u,v}\}$. We arrange $\{x_{u,v}\}$ in the same $10 \times 10$ layout (Figure 7). The generated images show clear class boundaries. Within each class, small moves on the $10 \times 10$ grid lead to smooth and local visual changes. Distant grid locations yield images that differ more. This behavior is consistent with the heat-kernel prior used in the generalized GWOT stage. Similar samples in $\mathcal{X}$ are encouraged to occupy nearby positions in the embedding, and dissimilar samples are encouraged to lie farther apart.
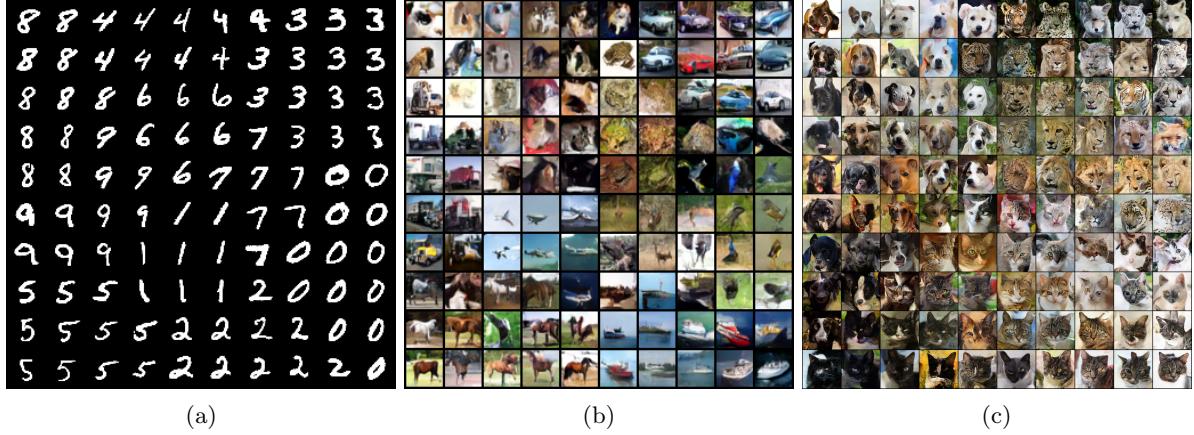


Figure 7: Image generations from a uniformly sampled 2D latent grid in $[-1.5, 1.5]^2$. Shown are (a) MNIST, (b) CIFAR-10, and (c) AFHQ.

## D.3 Architecture of the Model

Our drift network $u_\theta$ uses a U-Net backbone with time conditioning, shown in Figure 8. Based on the standard $t$ conditioned U-Net, we treat the two-dimensional embedding $y \in \mathbb{R}^2$ as a condition as well. We concatenate $[t, y]$ and inject it at every resolution. At each encoder or decoder block, we broadcast $[t, y]$ over the spatial grid, concatenate it along the channel dimension with the current feature map of $x$, and apply a linear projection before the block convolutions. We add self-attention blocks within the U-Net at matched encoder and decoder locations.

The network has two output heads selected by the role flag $r \in \{0, 1\}$. For the $y$ direction $r = 1$, the two-dimensional drift $v_y \in \mathbb{R}^2$ is produced at the bottleneck head. For the $x$ direction $r = 0$, the drift $v_x$ with the

**Wenxi Cai [1], Yuheng Wang [1], Naichen Shi [2]**

(a) U-Net backbone.



(b) Composition of each block.
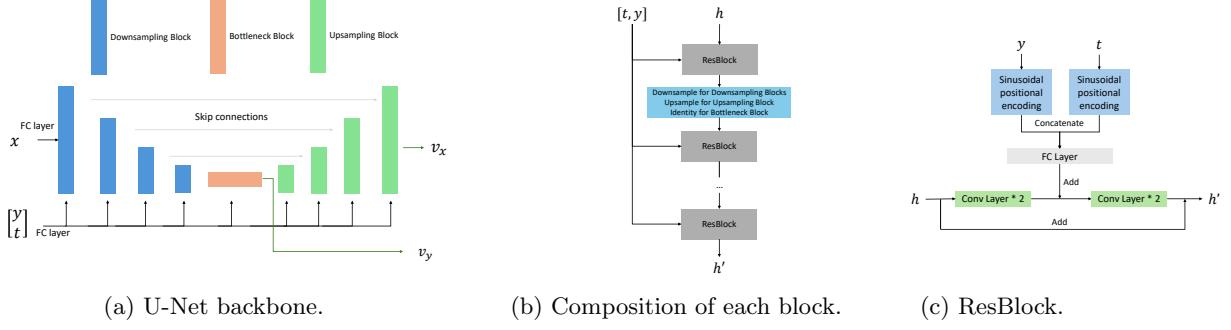


(c) ResBlock.

Figure 8: Architecture of the Model.

same shape as $x$ is produced by the final decoder head. During training, only the active head contributes to the loss while the other head is muted.

## D.4 Hyperparameters

Our model exposes several hyperparameters that control the U-Net capacity and where attention is applied. `model_channels` sets the base width of feature maps. `num_res_blocks` is the number of residual blocks per resolution level. `channel_mult` is a list of multipliers applied to `model_channels` at each level from high to low resolution. `attention_resolutions` lists the spatial sizes at which we insert self-attention blocks in both encoder and decoder. `num_heads` is the number of heads in each attention block. We use a learning rate of $1 \times 10^{-4}$ and train all models for 200 epochs.

| Dataset | model_channels | num_res_blocks | channel_mult | attention_resolutions | num_heads |
|---|---|---|---|---|---|
| MNIST | 64 | 2 | [1, 2, 2, 2] | [16] | 4 |
| CIFAR-10 | 128 | 2 | [1, 2, 2, 2] | [16] | 4 |
| TinyImageNet | 128 | 2 | [1, 2, 2, 2] | [16] | 4 |
| AFHQ | 192 | 2 | [1, 1, 2, 4] | [16, 32] | 4 |

Table 3: Default hyperparameter settings used in our runs