

# NEURAL OPTIMAL TRANSPORT FOR SUBSET ALIGNMENT

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We propose approaches for static and dynamic neural optimal transport with a relaxed Monge formulation to create optimal transport maps from a source distribution to an optimized distribution constrained to have an upper-bounded density ratio to the target distribution. In machine learning applications, this allows to learn the mappings between imbalanced datasets, such that one dataset can be mapped to a reweighted subset of a target dataset, with the reweighting governed by the density ratio constraint. The density ratio is constrained to lie in  $[0, c]$  by the  $f$ -divergence associated with the indicator function for  $[0, c]$ , where  $c$  denotes the maximum allowable upweighting factor. In the static case, neural networks are employed to parameterize the Monge map between source and selected subset of the target distribution and the dual function for the constraint. In the dynamic case, two networks are also employed: first neural network parametrizes the time dependent potential whose gradient defines the velocity field and terminal value enforces the density ratio constraint, while the second parametrizes the interpolation between the samples from source and optimized terminal distribution satisfying both the density ratio bound and the continuity equation. Since the terminal distribution in subset alignment need not be equal to the target distribution, which is distinct from prior work on dynamic neural optimal transport, we explore an efficient sampling scheme guided by the terminal potential. We apply both the static and dynamic formulations on domain translations problems, and demonstrate that the relaxed problem yields a more meaningful Monge map in cases where there is natural alignment between source and target distributions, but the distributions are imbalanced.

## 1 INTRODUCTION

Gaspard Monge proposed the original idea of optimal transport as mathematical model for the problem of minimum-cost transportation of dirt from source location to a destination Monge (1781). In more modern parlance, given probability measures,  $\mu$  defined on compact set  $\mathcal{X} \subseteq \mathbb{R}^d$ ,  $\nu$  defined on compact set  $\mathcal{Y} \subseteq \mathbb{R}^d$ , and the bounded uniformly continuous cost  $c(\cdot, \cdot) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ , Monge formulation of optimal transport is stated as

$$\begin{aligned} \mathcal{D}_{\text{Monge}}(\mu, \nu) = & \inf_{T \in \mathcal{T}(\mathcal{X}, \mathcal{Y})} \int_{\mathcal{X}} c(\mathbf{x}, T(\mathbf{x})) \mu(\mathbf{x}) d\mathbf{x} \\ \text{s.t. } & T_{\#}\mu = \nu \end{aligned} \quad (1)$$

where the set  $\mathcal{T}(\mathcal{X}, \mathcal{Y})$  denotes the set of measurable maps between  $\mathcal{X}$  and  $\mathcal{Y}$ . Monge formulation of the optimal transport problem requires that the transport map of  $T$  to be a deterministic function. In order to satisfy the constraint in the Monge problem 1, the transport map  $T$  must cover  $\nu$  upto some  $\nu$ -null sets. Usually, the cost  $c$  is non-linearly dependent on the transportation map  $T$ , making the problem 1 very cumbersome and very difficult to solve (Santambrogio, 2015; Villani et al., 2009).

Recently, neural networks have been widely employed to solve optimal transport problems. Seguy et al. (2018) employed stochastic gradient-based approaches to estimate the optimal transport (Monge) map for large-scale data. In comparison, earlier work (Genevay et al., 2016) only minimized the optimal transport loss using stochastic gradient-based methods, or, as in well-known



Wasserstein-GAN (Arjovsky et al., 2017; Gulrajani et al., 2017), employed the Kantorovich-Rubinstein duality to minimize the Wasserstein-1 loss function for generative modeling; however, the resulting generator is not trained to minimize distance as in the Monge formulation. Conversely, optimal transport maps can realize generative models (Daniels et al., 2021; Rout et al., 2022; Korotin et al., 2023b; Amos, 2023). For squared Euclidean transport cost, transport plans have been either directly parameterized using non-convex neural networks, (Rout et al., 2022; Korotin et al., 2023b) or obtained by amortizing the convex conjugate as gradients of convex functions parameterized by input convex neural networks (Amos et al., 2017; Makkuva et al., 2020; Korotin et al., 2021a; Amos, 2023; Vesseron & Cuturi, 2024). With recent developments in the development of flow matching (Lipman et al., 2023; Liu et al., 2023; Albergo & Vanden-Eijnden, 2023) as a state-of-the-art method for image generation, considerable recent efforts have been made to develop an efficient neural network-based framework for dynamic optimal transport for a variety of trajectory inference and generative modeling problems (Pooladian et al., 2024; Neklyudov et al., 2023; 2024b).

While distinct from generative modeling, the Monge map is a meaningful concept for the alignment of two real distributions (neither of which is noise) from slightly different domains, as in unsupervised domain adaptation. In these cases, distributional imbalance creates challenges (Wu et al., 2019). There has been substantial theoretical work on partial optimal transport (Figalli, 2010; Caffarelli & McCann, 2010; Chizat et al., 2018b;a) where two measures are not required to be of equal mass, and Wasserstein Fisher-Rao distance (Chizat et al., 2018a;b; Bauer et al., 2016) which allows for mass growth and destruction during the transfer process. Recent work on neural optimal transport in these cases (Gazdieva et al., 2023; Choi et al., 2023; Yang & Uhler, 2019). In this work, we formulate a relaxed version of optimal transport that creates a new distribution whose density ratio to the target distribution is bounded.

We propose static and dynamic neural optimal transport formulations, under the constraint density ratio constraint. To minimize the expected ground distance<sup>1</sup>, the transported distribution can have a support that is subset of the target support. This can be interpreted as a reweighted target distribution with mass concentrated entirely on the selected subset. Our key contributions are as follows: we formulate both static and dynamic subset alignment problems by replacing the target marginal constraint with a penalty based on an  $f$ -divergence corresponding to the convex indicator function of the set  $[0, c]$ , where  $c = 1$  recovers standard optimal transport; we leverage dual formulations of our problems using neural networks, in particular, we employ Benamou-Brenier formulation (see equation 22 in the appendix) along with the Lagrange multiplier method to obtain the dual form of dynamic subset selection; we show that the dual formulations in both the static and dynamic yield a potential function defined over the target support, whose sign effectively distinguishes points within the selected subset from those outside it; and we apply our framework to unpaired domain translation problems and use the potential function for PU-learning.

## 2 METHODOLOGY

### 2.1 STATIC SUPPORT SUBSET-SELECTION

The Kantorovich formulation (Kantorovich, 1942) for the optimal transport problem is

$$\mathcal{W}(\mu, \nu) = \inf_{\pi} \int_{\mathcal{X} \times \mathcal{Y}} c(\mathbf{x}, \mathbf{y}) \pi(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}, \quad \text{s.t.} \int_{\mathcal{Y}} d\pi(\mathbf{x}, \mathbf{y}) = \mu(\mathbf{y}), \quad \int_{\mathcal{X}} d\pi(\mathbf{x}, \mathbf{y}) = \nu(\mathbf{x}), \quad (2)$$

where  $\pi$  is a density defined on  $\mathcal{X} \times \mathcal{Y}$ . Our formulation of static support subset-selection for optimal transport is derived from a relaxed problem where the constraint on the first marginal of the joint density  $\pi$  is maintained, while the second marginal  $\int_{\mathcal{X}} \pi(\mathbf{x}, \mathbf{y}) d\mathbf{x} = \tilde{\nu}(\mathbf{y})$  is allowed to vary from within a range  $[0, c]$  of the target density  $\nu$ , such that  $0 \leq \frac{\tilde{\nu}(\mathbf{y})}{\nu(\mathbf{y})} \leq c$ . The density  $\tilde{\nu}$  can be interpreted as a reweighted target density  $\tilde{\nu}(\mathbf{y}) = \omega(\mathbf{y})\nu(\mathbf{y})$ ,  $0 \leq \omega(\mathbf{y}) \leq c$ , where portions of the support can be up-weighted while others are down-weighted or removed. The relaxed constraint is equivalent to a case of the partial optimal transport relaxations using  $f$ -divergences introduced by

<sup>1</sup>While we focus on the Euclidean distance, more general distances can be considered.



(Séjourné et al., 2023)

$$\inf_{\pi} \int_{\mathcal{X} \times \mathcal{Y}} c(\mathbf{x}, \mathbf{y}) \pi(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} + \mathcal{D}_{\iota_{[a,b]}}(\tilde{\nu} \| \nu) \quad \text{s.t.} \quad \int_{\mathcal{Y}} \pi(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \mu(\mathbf{x}), \quad (3)$$

where  $\mathcal{D}_{\iota_{[a,b]}}$  is the range divergence with  $\iota_{[a,b]}$  being the convex indicator function

$$\iota_{[a,b]}(r) = \begin{cases} 0, & r \in [a, b] \\ +\infty, & \text{o.w.} \end{cases}, \quad \iota_{[a,b]}^*(t) = \sup_{u \in [a,b]} (u \cdot t) = \max(-at, bt), \quad (4)$$

and  $\iota_{[a,b]}^*$  denotes its Legendre-Fenchel conjugate. Since the function  $\iota_{[a,b]}$  is convex lower semi-continuous, therefore  $\iota_{[a,b]} = \iota_{[a,b]}^{**}$ , we can apply the variational form of the  $f$ -divergence Nguyen et al. (2010), (Polyanskiy & Wu, 2025, Theorem 7.26), exploited by  $f$ -GAN (Nowozin et al., 2016), leading to a form requiring only expected values

$$\mathcal{D}_{\varphi}(\tilde{\nu} \| \nu) = \int_{\mathcal{Y}} \iota_{[a,b]} \left( \frac{\tilde{\nu}}{\nu}(\mathbf{y}) \right) \nu(\mathbf{y}) d\mathbf{y} = \sup_{\eta} \underbrace{\int_{\mathcal{Y}} \eta(\mathbf{y}) \tilde{\nu}(\mathbf{y}) d\mathbf{y}}_{\mathbb{E}_{\tilde{\mathbf{y}} \sim \tilde{\nu}}[\eta(\tilde{\mathbf{y}})]} - \underbrace{\int_{\mathcal{Y}} \iota_{[a,b]}^*(\eta(\mathbf{y})) \nu(\mathbf{y}) d\mathbf{y}}_{\mathbb{E}_{\mathbf{y} \sim \nu}[\iota_{[a,b]}^*(\eta(\mathbf{y}))]}. \quad (5)$$

To match 3, we focus on  $a = 0$  and  $b = c \geq 1$ , such that  $\iota_{[0,c]}^*(t) = c \cdot \max(0, t)$  and for compactness denote  $\eta_+(\mathbf{y}) = \max(0, \eta(\mathbf{y}))$ . Introducing  $\psi$  as a measurable function to act as a Lagrange multiplier to enforce the constraint in 3 and combining with equation 5 yields the problem

$$\inf_{\pi} \sup_{\psi, \eta} \int_{\mathcal{X} \times \mathcal{Y}} (c(\mathbf{x}, \mathbf{y}) + \eta(\mathbf{y}) - \psi(\mathbf{x})) \pi(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} + \int \psi(\mathbf{x}) \mu(\mathbf{x}) d\mathbf{x} - c \int \eta_+(\mathbf{y}) \nu(\mathbf{y}) d\mathbf{y}. \quad (6)$$

As described in App. A.1, since  $c$  is convex and lower semi-continuous, we interchange the  $\inf_{\pi}$  and  $\sup_{\eta}$  and apply what is known as the  $c$ -transform of  $-\eta(\mathbf{y})$  (Santambrogio, 2015; Villani et al., 2009) to obtain the dual problem with measurable map  $T : \mathcal{X} \rightarrow \mathcal{Y}$

$$\sup_{\eta} \inf_T \underbrace{\int_{\mathcal{X}} (c(\mathbf{x}, T(\mathbf{x})) + \eta(T(\mathbf{x}))) \mu(\mathbf{x}) d\mathbf{x}}_{\mathbb{E}_{\mathbf{x} \sim \mu}[c(\mathbf{x}, T(\mathbf{x})) + \eta(T(\mathbf{x}))]} - c \underbrace{\int_{\mathcal{Y}} \eta_+(\mathbf{y}) \nu(\mathbf{y}) d\mathbf{y}}_{\mathbb{E}_{\mathbf{y} \sim \nu}[c \cdot \max(0, \eta(\mathbf{y}))]}, \quad (7)$$

For the computational implementation  $T$  and  $\eta$  are parameterized using neural networks with associated parameters  $\theta_T$  and  $\theta_{\eta}$  and expectations are estimated using samples from  $\mu$  and  $\nu$  as described in the Algorithm 1.

---

**Algorithm 1:** (Static-Neural-SS) Learning Algorithm for Static Subset Selection

---

**Inputs** : Source distribution  $\mu$  and target distributions  $\nu$ , cost function  $c(\cdot, \cdot)$ , reweighting bound  $c$ , neural networks  $T(\cdot, \theta_T)$  and  $\eta(\cdot, \theta_{\eta})$ , batch size  $N$ , number of updates  $n_T$  and  $n_{\eta}$ , and optimizers  $\text{optim}_T$  and  $\text{optim}_{\eta}$ .

**Outputs** : Sample based neural estimate for transport map  $T$

```

1 for all learning iterations do
2   for  $n_T$  update steps do
3     sample  $\{\mathbf{x}_i\}_{i=1}^N \sim \mu$  and  $\{\mathbf{y}_j\}_{j=1}^N \sim \nu$ 
4     compute  $\text{grad}_{\theta_T} = \nabla_{\theta_T} \frac{1}{N} \sum_{i=1}^N [c(\mathbf{x}_i, T(\mathbf{x}_i, \theta_T)) + \eta(T(\mathbf{x}_i, \theta_T), \theta_{\eta})]$ 
5     use  $\text{grad}_{\theta_T}$  to update  $\theta_T$  with  $\text{optim}_T$ 
6   end
7   for  $n_{\eta}$  update steps do
8     sample  $\{\mathbf{x}_i\}_{i=1}^N \sim \mu$  and  $\{\mathbf{y}_j\}_{j=1}^N \sim \nu$ 
9     compute  $\text{grad}_{\theta_{\eta}} = \nabla_{\theta_{\eta}} \frac{1}{N} \sum_{i=1}^N [c \cdot \max(0, \eta(\mathbf{y}_j, \theta_{\eta})) - \eta(T(\mathbf{x}_i, \theta_T), \theta_{\eta})]$ 
10    use  $\text{grad}_{\theta_{\eta}}$  to update  $\theta_{\eta}$  with  $\text{optim}_{\eta}$ 
11  end
12 end

```

---



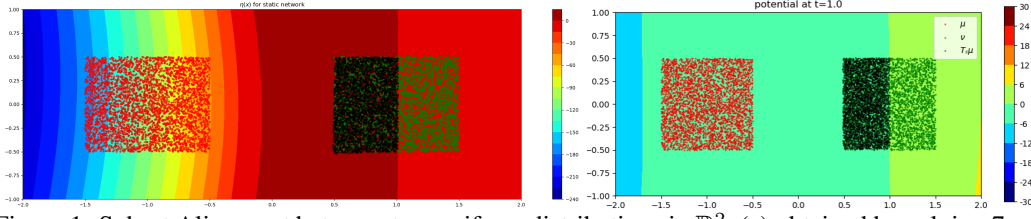


Figure 1: Subset Alignment between two uniform distributions in  $\mathbb{R}^2$ , (a) obtained by solving 7 and (b) obtained by solving 9 at  $c = 2$  by using fully connected neural networks to parametrize  $\eta$ ,  $T$ ,  $\varphi_t$  and  $\rho_t$ .

## 2.2 DYNAMIC SUBSET-SELECTION

Our formulation of dynamic support subset-selection for optimal transport is directly related to the Benamou-Brenier formulation of Wasserstein-2 distance. Similar to the static case, we replace the second marginal by a penalty based on the range divergence. The modified Benamou-Brenier problem is

$$\begin{aligned} \inf_{\rho_t, v_t} \int_0^1 \int_{\Omega} \frac{\|v_t(x)\|^2}{2} \rho_t(x) dx dt + \mathcal{D}_{[0, c]}(\rho_1 \| \nu) \\ \text{s.t.} \quad \frac{\partial}{\partial t} \rho_t(x) + \text{div}(\rho_t(x) v_t(x)) = 0, \quad \rho_0(x) = \mu(x) \end{aligned} \quad (8)$$

By introducing the Lagrange multiplier for  $\varphi_t$  for continuity equation constraint, one can write the dual form of equation 8 as (see Appendix A.2 for details)

$$\sup_{\rho_t} \inf_{\varphi_t} \mathbb{E}_{x \sim \mu} [\varphi_0(x)] + \mathbb{E}_{x \sim \nu} [c \cdot \max(0, -\varphi_1(x))] + \int_0^1 \mathbb{E}_{x_t \sim \rho_t} \left[ \frac{\partial}{\partial t} \varphi_t(x_t) + \frac{\|\nabla \varphi_t(x_t)\|^2}{2} \right] dt. \quad (9)$$

From equation 41 and equation 9, one can see that, in addition to samples from source and target distributions, one additionally needs to have a mechanism to sample from an optimized distribution that interpolates between the source distribution and the terminal distribution that satisfies the range divergence to the target. This is essentially a generative modeling problem and the subject of many recent studies (Neklyudov et al., 2024a; Atanackovic et al., 2025; Du et al., 2024).

In flow-based models, instead of explicitly modeling  $\rho_t$ , samples  $x_0 \sim \mu$  and  $x_1 \sim \nu$  are used to generate  $x_t$  using an analytically defined interpolant Lipman et al. (2023); Liu et al. (2023); Albergo & Vanden-Eijnden (2023). In this work, we adapt the computational framework for learning Wasserstein-Lagrangian flows (WLF) (Neklyudov et al., 2024b) to parameterize  $\rho_t$  in terms of  $\mu$  and  $\nu$ . For a given  $t \in [0, 1]$ , WLF creates an interpolant  $x_t \sim \rho_t$  from  $x_0 \sim \mu$  and  $x_1 \sim \nu$  (independently sampled) as

$$x_t = (1 - t)x_0 + tx_1 + t(1 - t)Q_t(x_0, x_1), \quad (10)$$

where  $Q_t$  is time-dependent neural network, which internally uses an additional Heaviside step function input  $t \geq 0.5$  (Neklyudov et al., 2024b). In the case when  $c = 1$ , subset alignment is equivalent to the optimal transport problem, therefore optimally  $\rho_1^* = \nu$ , also given the optimal velocity field  $v_t^* = \nabla \varphi_t^*$ , the optimal interpolant  $x_t^* \sim \rho_t^*$  is related to  $v_t^*$  by

$$x_t^* = \begin{cases} x_0 + \int_0^t v_{\tau}^*(x_{\tau}) d\tau & t < 0.5 \\ x_1 + \int_1^t v_{\tau}^*(x_{\tau}) d\tau & t \geq 0.5 \end{cases}, \text{ resulting in forward integration from } x_0 \text{ for } t < 0.5, \text{ and}$$

backward integration from  $x_1$  otherwise. However, for  $c > 1$   $\rho_1^* \neq \nu$ , therefore we can not directly draw samples  $x \sim \nu$  and propagate them backward for  $t \geq 0.5$ . Instead, an optimal interpolant could simply use the forward integration from  $x_0$ . This means that  $Q_t^*$  would require the capacity to be a one-step integrator, which is not different from the  $t < 0.5$  case for  $c = 1$ . However, in practice, the optimization of  $\rho_t$  lags behind  $\varphi_t$ , and it may be advantageous to map samples from  $\nu$  (or a distribution close to  $\nu$ ) in order sample from  $\rho_1$ . We propose to sample  $\tilde{x}_1 \sim \tilde{\nu}$ , where  $\tilde{\nu}$  is chosen judiciously, and replace  $x_1$  with  $\tilde{x}_1$  in equation 10. In this case, the optimal interpolant  $x_t^* \sim \rho_t^*$  is still forward integration from  $x_0$  for  $t < 0.5$ , but for  $t \geq 0.5$ ,  $Q_t^*(x_0, \tilde{x}_1)$  needs an internal map  $S^*$  such that the backward integration starts from a point sampled from the optimal



terminal distribution  $\mathbf{x}_1^* = S^*(\tilde{\mathbf{x}}_1) \sim \rho_1^*$  for  $t \geq 0.5$ , where  $S_{\#}^* \tilde{\nu} = \rho_1^*$ . If  $\tilde{\nu} = \nu$  then  $S^*$  maps the original target to  $\rho_1^*$ .

Our insight is to create  $\tilde{\nu}$  by leveraging the fact that the optimal potential  $\varphi_1^*$  satisfies  $\varphi_1^* \leq 0$  almost surely on  $\text{supp}(\tilde{\nu})$  and  $\varphi_1^* > 0$  almost surely on  $\text{supp}(\nu) \setminus \text{supp}(\tilde{\nu})$  (see Appendix B). Conditioning on the sign of  $\varphi_1^*$  allows us to sample from the selected subset of  $\text{supp}(\nu)$ . Given a current estimate  $\varphi_1$ , we create  $\nu_{\varphi_1}$ , a distribution supported on the subset of the target where  $\varphi_1 \leq 0$ , as  $\nu_{\varphi_1}(\mathbf{x}) = \nu(\mathbf{x} \mid \varphi_1(\mathbf{x}) \leq 0)$ . When  $\tilde{\nu} = \nu_{\varphi_1} = \rho_1^*$  then  $S^*(\tilde{\mathbf{x}}_1) = \tilde{\mathbf{x}}_1$ . During training, however,  $\varphi_1$  is suboptimal and may miss part of the support of the original target  $\nu$ , so we sample from the mixture  $\alpha \nu_{\varphi_1} + (1 - \alpha)\nu$ ,  $\alpha \in [0, 1]$ . Assuming  $\varphi_1$  improves with training, we create a sequence of distributions, where at the  $k$ -th learning iteration, we can sample from the mixture

$$\tilde{\nu}^{(k)} = \alpha^{(k)} \nu_{\varphi_1^{(k)}} + (1 - \alpha^{(k)}) \nu, \quad (11)$$

where  $\alpha^{(k)}$  follows a monotonically non-decreasing scheduler with  $\alpha^{(0)} = 0$  and  $\alpha^{(\infty)} = 1$ .<sup>2</sup> The complete procedure for solving the dynamic subset selection problem is outlined in Algorithm 2, wherein optimized parameters are  $\theta_\varphi$  and  $\theta_\rho$  (variables that are functions of parameters whose gradients are needed are explicitly noted).

---

**Algorithm 2:** (Dynamic-Neural-SS) Learning Algorithm for Dynamic Subset Selection

---

**Inputs** : Source distribution  $\mu$  and target distributions  $\nu$ , time-dependent neural network  $\varphi_t(\cdot, \theta_\varphi)$ , network for the interpolant  $Q_t(\cdot, \cdot, \theta_\rho)$  along with mixture schedule  $\alpha^{(k)}$ , batch size  $N$ , number of updates  $n_\varphi$  and  $n_\rho$ , and optimizers  $\text{optim}_\varphi$  and  $\text{optim}_\rho$ .

**Outputs** : Sample based neural estimate for  $\varphi_t(\cdot, \theta_\varphi)$

```

1 for learning iteration  $k = 0, 1, \dots$  do
2   for  $\varphi_t$  update steps do
3     sample  $\{\mathbf{x}_0^i\}_{i=1}^N \sim \mu$ ,  $\{\mathbf{x}_1^i\}_{i=1}^N \sim \nu$ ,  $\{\tilde{\mathbf{x}}_1^i\}_{i=1}^N \sim \tilde{\nu}^{(k)}$ , and  $\{t^i\}_{i=1}^N \sim \text{Uniform}([0, 1])$ 
4     compute  $\tilde{\mathbf{x}}_t^i = (1 - t^i)\mathbf{x}_0^i + t^i\mathbf{x}_1^i + t^i(1 - t^i)Q_{t^i}(\mathbf{x}_0^i, \tilde{\mathbf{x}}_1^i, \theta_\rho)$ ,  $\forall i \in \{1, \dots, N\}$ .
5     compute
        
$$\text{grad}_{\theta_\varphi} = \nabla_{\theta_\varphi} \frac{1}{N} \sum_{i=1}^N \left[ \frac{\partial}{\partial t} \varphi_{t^i}(\tilde{\mathbf{x}}_t^i, \theta_\varphi) + \frac{\|\nabla \varphi_{t^i}(\tilde{\mathbf{x}}_t^i, \theta_\varphi)\|^2}{2} \right. \\
        \left. + \varphi_0(\mathbf{x}_0^i, \theta_\varphi) + c \cdot \max(0, -\varphi_1(\mathbf{x}_1^i, \theta_\varphi)) \right].$$

6     use  $\text{grad}_{\theta_\varphi}$  to update  $\theta_\varphi$  with  $\text{optim}_\varphi$ .
7   end
8   for  $\rho_t$  update steps do
9     sample  $\{\mathbf{x}_0^i\}_{i=1}^N \sim \mu$ ,  $\{\tilde{\mathbf{x}}_1^i\}_{i=1}^N \sim \tilde{\nu}^{(k)}$ ,  $\{t^i\}_{i=1}^N \sim \text{Uniform}([0, 1])$ .
10    compute  $\tilde{\mathbf{x}}_t^i(\theta_\rho) = (1 - t^i)\mathbf{x}_0^i + t^i\tilde{\mathbf{x}}_1^i + t^i(1 - t^i)Q_{t^i}(\mathbf{x}_0^i, \tilde{\mathbf{x}}_1^i, \theta_\rho)$ ,  $\forall i \in \{1, \dots, N\}$ .
11    compute
        
$$\text{grad}_{\theta_\rho} = \nabla_{\theta_\rho} \frac{1}{N} \sum_{i=1}^N \left[ \frac{\partial}{\partial t} \varphi_{t^i}(\tilde{\mathbf{x}}_t^i(\theta_\rho), \theta_\varphi) + \frac{\|\nabla \varphi_{t^i}(\tilde{\mathbf{x}}_t^i(\theta_\rho), \theta_\varphi)\|^2}{2} \right].$$

12    use  $\text{grad}_{\theta_\rho}$  to update  $\theta_\rho$  with  $\text{optim}_\rho$ .
13  end
14 end
```

---

### 3 RELATED WORK

In addition to approaches mentioned in the introduction, we review advances in static neural optimal transport in the Appendix C.1. Our work on dynamic subset selection is most directly related to Lagrangian neural optimal transport (Pooladian et al., 2024), action-matching (Neklyudov et al., 2023) and Wasserstein Lagrangian flows (Neklyudov et al., 2024a). Pooladian et al. (2024). The neural optimal transport with Lagrangian costs framework (Pooladian et al., 2024) focuses on optimal

<sup>2</sup>Instead of trusting the sign directly, for small finite target datasets, we evaluate  $\varphi_1$  for all  $\mathbf{x}_1$  and retain the fraction  $\frac{1}{c}$  of points with smallest value of  $\varphi_1$  to obtain the sample from  $\tilde{\nu}_\varphi$ .



transport with different potentials in Euclidean space. Wasserstein-Lagrangian flows (Neklyudov et al., 2023) is mainly developed for the applications in cellular trajectory inference and quantum many body problems (Neklyudov et al., 2024b), and extends to more general settings on Wasserstein Fisher-Rao (Chizat et al., 2018a;b; Séjourné et al., 2023), with the ability to deal with mass growth/destruction, and different types of dynamics.

All these approaches and all flow-based models are developed for the cases when the marginals are to be preserved. (A more extensive review of recent work in dynamic neural optimal transport is included in Appendix C.2; additionally, since the optimal transport is intimately related to recent developments in generative modeling such as flow-matching and Schrödinger bridges, we also discuss the development in relation to optimal transport.) In contrast, with our proposed dynamic support subset-selection it is desirable to preserve one marginal and dynamically transfer that mass to the subset of the support of the other while minimizing the transport cost. Therefore our approach is an novel extension of prior work (Neklyudov et al., 2024a), and although we focused on the  $\ell_2^2$  cost, our method is compatible with other Lagrangian costs (Pooladian et al., 2024), which could be useful for side-information as in semi-supervised domain adaptation.

## 4 EXPERIMENTS AND RESULTS

In this section we discuss the experimental results for susbet selection on an easily interpretable image-to-image case, where MNIST (Deng, 2012) is the source and EMNIST (Cohen et al., 2017) is the target. In this case, images of digits are a subset of the characters in EMNIST. We then apply our proposed approaches to domain translation on the FFHQ dataset (Karras et al., 2019) in 512-dimensional latent space of ALAE (Pidhorskyi et al., 2020).

### 4.1 MNIST $\rightarrow$ EMNIST DOMAIN TRANSLATION

MNIST data set contains 60,000 images of digits between 0-9 in training-subset and 10,000 images in test-subset. MNIST dataset is roughly balanced in the sense that the proportions of each data class in the dataset are roughly the same. EMNIST (byclass) dataset contains a set of English alphabet and numbers. EMNIST contains 62 imbalanced classes, of which 10 classes (between 0-9) represent numbers, and the rest of 52 classes represent upper and lower English case letters of the English alphabet. Roughly, 16% of EMNIST represent numbers and remaining 84% are alphabet.

Since our goal is to transfer MNIST images to EMNIST images such a way that MNIST digits are mapped to EMNIST digits while ignoring alphabet, we trained a neural network classifier to distinguish between digits and alphabet to evaluate the learned mapping (see implementation details in Appendix D.1. After training the classifier, we used both static and dynamic subset-selection approaches for domain translation between MNIST and EMNIST. Implementation details of the underlying models and there training are in Appendix D.2.

In our experiments, we trained and evaluated both static and dynamic models using both the static and dynamic subset-selection frameworks for  $c \in \{1, 2, 4, 8\}$ . For the dynamic case, similar to any flow based generative process, dynamic subset selection also requires a numerical integration (ODE integration with Euler type numerical integrator with 100 integration steps), but one-step integration can be used (Liu et al., 2023; 2024b). Figure 2 shows that perceptually, one-step integration performs worse in comparison to both static and ODE-based generation. We evaluated the classification

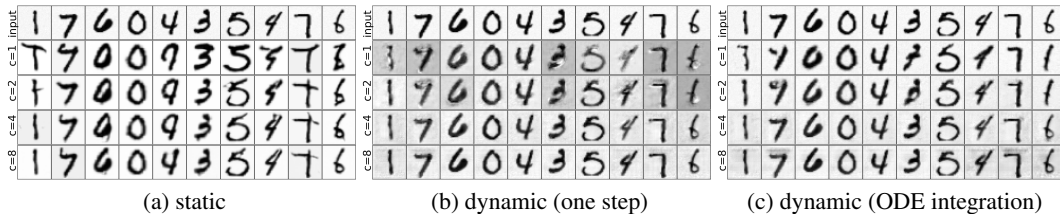


Figure 2: Image translation outputs for MNIST  $\rightarrow$  EMNIST



accuracy on translation of whole MNIST dataset, confusion matrices are given in Appendix E. A summary of accuracies of translated outputs are given in Table 1.

Method	c=1	c=2	c=4	c=8
static	46.93	75.33	82.44	87.32
dynamic (one step)	64.85	75.16	93.57	95.84
dynamic (ODE)	58.80	70.47	92.68	95.00

Table 1: Classification accuracies of translated images MNIST→EMNIST evaluated with using pretrained classifier.

## 4.2 POSITIVE-UNLABELED LEARNING

Positive Unlabeled (PU) learning is a binary classification problem in which only a subset of positive data is labeled, which is then used to train a model classifying between positive and negative data from an unlabeled (containing both positive and negative data) data set. PU Learning Bekker & Davis (2020); Kato et al. (2019); Chapel et al. (2020); Riaz et al. (2023). Since the sign of an optimal potential function in our framework differs between selected and unselected subsets, one can use it to distinguish between them positive and unlabeled datasets (see Appendix B for details). We applied applied both the static and dynamic optimal transport for PU learning on the 20 UCI-datasets (Kelly et al.) as in (Teisseyre et al., 2025), using the same settings with 75/25 train-test split on each dataset and the sampled completely randomly (SCAR) mechanism to selected and label points.

Networks were 5-layer MLPs with swish activation functions of appropriate input and output dimensions for both static and dynamic subset alignment with fixed learning rates for both static and dynamic models. Architecture and parameter details for each model are given in D.3. We trained 20 different models for each dataset using different train test splits, so in total we trained 400 models for static and 400 models for dynamic subset alignment. We adopted alternative sign and value based label assignment strategies for unlabeled dataset. Performance in terms of balanced accuracy for our approaches along with the top-performing baselines PUSB (Kato et al., 2019) and NTC-MI (Teisseyre et al., 2025) are given in Table 2.

Dataset	$\pi$	PUSB	NTC-MI	static		dynamic	
				sorted	sign	sorted	sign
Abalone	0.16	0.544 $\pm$ 0.060	<b>0.575 <math>\pm</math> 0.025</b>	<b>0.561 <math>\pm</math> 0.029</b>	0.503 $\pm$ 0.008	0.555 $\pm$ 0.033	0.532 $\pm$ 0.030
Banknote	0.44	0.829 $\pm$ 0.050	<b>0.922 <math>\pm</math> 0.019</b>	<b>0.883 <math>\pm</math> 0.037</b>	0.882 $\pm$ 0.039	0.892 $\pm$ 0.048	0.895 $\pm$ 0.044
Breast-w	0.34	0.766 $\pm$ 0.145	0.870 $\pm$ 0.028	<b>0.930 <math>\pm</math> 0.028</b>	<b>0.941 <math>\pm</math> 0.027</b>	0.839 $\pm$ 0.197	0.831 $\pm$ 0.132
Diabetes	0.35	0.546 $\pm$ 0.042	<b>0.700 <math>\pm</math> 0.039</b>	<b>0.635 <math>\pm</math> 0.044</b>	0.635 $\pm$ 0.044	0.587 $\pm$ 0.094	0.603 $\pm$ 0.066
Haberman	0.26	0.513 $\pm$ 0.023	0.532 $\pm$ 0.066	<b>0.539 <math>\pm</math> 0.066</b>	<b>0.540 <math>\pm</math> 0.067</b>	0.528 $\pm$ 0.062	0.519 $\pm$ 0.070
Heart	0.44	0.527 $\pm$ 0.033	<b>0.757 <math>\pm</math> 0.053</b>	<b>0.637 <math>\pm</math> 0.093</b>	0.623 $\pm$ 0.089	0.508 $\pm$ 0.210	0.573 $\pm$ 0.139
Ionosphere	0.64	0.440 $\pm$ 0.085	0.755 $\pm$ 0.059	<b>0.773 <math>\pm</math> 0.091</b>	<b>0.762 <math>\pm</math> 0.088</b>	0.562 $\pm$ 0.215	0.602 $\pm$ 0.149
Isolet	0.04	0.793 $\pm$ 0.072	0.725 $\pm$ 0.006	<b>0.881 <math>\pm</math> 0.028</b>	<b>0.923 <math>\pm</math> 0.030</b>	0.673 $\pm$ 0.173	0.693 $\pm$ 0.202
Jm1	0.19	<b>0.628 <math>\pm</math> 0.016</b>	<b>0.628 <math>\pm</math> 0.013</b>	0.576 $\pm$ 0.015	0.575 $\pm$ 0.010	0.573 $\pm$ 0.038	0.565 $\pm$ 0.026
Kc1	0.15	<b>0.645 <math>\pm</math> 0.075</b>	<b>0.679 <math>\pm</math> 0.030</b>	0.604 $\pm$ 0.036	0.607 $\pm$ 0.035	0.611 $\pm$ 0.063	0.606 $\pm$ 0.054
Madelon	0.5	0.496 $\pm$ 0.030	0.519 $\pm$ 0.028	<b>0.533 <math>\pm</math> 0.025</b>	<b>0.523 <math>\pm</math> 0.015</b>	0.511 $\pm$ 0.027	0.505 $\pm$ 0.017
Musk	0.15	0.712 $\pm$ 0.036	0.767 $\pm$ 0.012	<b>0.841 <math>\pm</math> 0.018</b>	<b>0.847 <math>\pm</math> 0.018</b>	0.823 $\pm$ 0.020	0.840 $\pm$ 0.019
Segment	0.14	0.848 $\pm$ 0.074	0.803 $\pm$ 0.014	0.898 $\pm$ 0.038	<b>0.927 <math>\pm</math> 0.031</b>	0.900 $\pm$ 0.042	<b>0.935 <math>\pm</math> 0.026</b>
Semeion	0.1	0.569 $\pm$ 0.055	0.755 $\pm$ 0.022	<b>0.824 <math>\pm</math> 0.044</b>	<b>0.850 <math>\pm</math> 0.067</b>	0.699 $\pm$ 0.143	0.653 $\pm$ 0.144
Sonar	0.53	0.497 $\pm$ 0.041	<b>0.573 <math>\pm</math> 0.057</b>	<b>0.561 <math>\pm</math> 0.091</b>	0.524 $\pm$ 0.074	0.515 $\pm$ 0.107	0.511 $\pm$ 0.089
Spambase	0.39	0.821 $\pm$ 0.031	<b>0.887 <math>\pm</math> 0.014</b>	<b>0.786 <math>\pm</math> 0.011</b>	0.775 $\pm$ 0.010	0.703 $\pm$ 0.057	0.664 $\pm$ 0.067
Vehicle	0.26	0.549 $\pm$ 0.067	0.804 $\pm$ 0.042	<b>0.806 <math>\pm</math> 0.037</b>	<b>0.823 <math>\pm</math> 0.032</b>	0.639 $\pm$ 0.169	0.661 $\pm$ 0.152
Waveform	0.34	<b>0.860 <math>\pm</math> 0.012</b>	<b>0.829 <math>\pm</math> 0.012</b>	0.795 $\pm$ 0.015	0.743 $\pm$ 0.013	0.676 $\pm$ 0.071	0.551 $\pm$ 0.019
Wdbc	0.37	0.798 $\pm$ 0.155	0.801 $\pm$ 0.043	<b>0.861 <math>\pm</math> 0.068</b>	<b>0.845 <math>\pm</math> 0.063</b>	0.691 $\pm$ 0.211	0.641 $\pm$ 0.149
Yeast	0.31	0.517 $\pm$ 0.051	<b>0.657 <math>\pm</math> 0.024</b>	<b>0.630 <math>\pm</math> 0.040</b>	0.612 $\pm$ 0.049	0.590 $\pm$ 0.076	0.567 $\pm$ 0.063

Table 2: Comparison of average balanced accuracies of 20 models trained using static and dynamic subset alignment methods with PUSB and NTC-MI reported Teisseyre et al. (2025). Balanced accuracies for best performing methods are colored red and second best are colored blue.

## 4.3 FFHQ IMAGE TRANSLATION

We also apply our proposed approaches to the unpaired image translation problem. We followed the experimental setup of Gazdieva et al. (2024), where the FFHQ dataset embedded in the latent space of Adversarial Latent Autoencoder (ALAE) (Pidhorskyi et al., 2020), is divided either by gender



(man or woman) or age, as two orthogonal labels. Table 3 adapted from Gazdieva et al. (2024), shows the number of images for each class, where images with age  $< 16$  are ignored, ages between 16 and 43 are labeled young, and the remainder are labeled old. Given these classes, the task is to learn to map a source distribution to a target distribution. There are four cases, young to old, old to young, man to woman, and woman to man. In order to evaluate the translation process, two classifiers pretrained in the ALAE latent space are used, one classifier is trained to classify young vs old and another to distinguish man vs woman. The target accuracy quantifies what proportion of translated images lie within the target-class boundary. The source accuracy quantifies whether the translated images retain the orthogonal label. For example, with young $\rightarrow$ old the source accuracy is whether the ‘aged’ image of a young source image retains the same gender.

Implementation details for both static and dynamic subset selection to the FFHQ dataset are given in Appendix D.4. Between young $\rightarrow$ old, old $\rightarrow$ young, man $\rightarrow$ woman and woman $\rightarrow$ man, it was observed that larger values of  $c$  tend to preserve the source accuracy, but often have lower target accuracy. This can be related to the fact that for larger values of  $c$ , it takes more training steps to achieve the optimal subset selection.

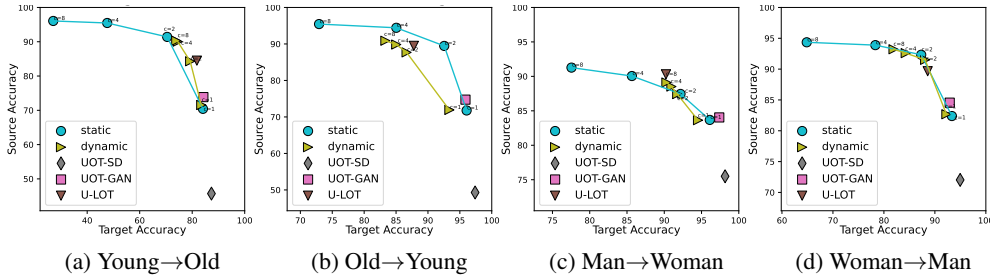


Figure 3: Accuracy curves for  $c \in \{1, 2, 4, 8\}$ , in comparison to results from LOT (Gazdieva et al., 2024), UOT-GAN (Yang & Uhler, 2019), and UOT-SD (Choi et al., 2024b).

We compared our methodology with Light Unbalanced optimal transport Gazdieva et al. (2024)(LOT), Yang & Uhler (2019)(UOT-GAN) and Choi et al. (2024b)(UOT-SD) and observed that methods which achieve better results in terms of target accuracy perform worse in terms of source class accuracy. This can be seen from Table 4 and Figure 3, using accuracy values reported by Gazdieva et al. (2024). Example translated images for static and dynamic are shown for old $\rightarrow$ young in Figure 4, with other cases provided in Appendix F.

Class	Man	Woman
Young	15K	23K
Old	7K	3.5K

Table 3: Division of FFHQ train images.

Task	Accuracy	c=1		c=2		c=4		c=8		UOT-SD	UOT-GAN	U-LOT
		static	dynamic	static	dynamic	static	dynamic	static	dynamic			
Young $\rightarrow$ Old	Target	84.09	83.45	70.47	79.23	47.63	74.51	27.07	73.93	87.33	84.25	81.78
	Class	70.43	71.55	91.41	84.31	95.47	90.03	96.06	90.30	45.71	73.85	84.49
Old $\rightarrow$ Young	Target	96.06	93.36	92.55	86.65	85.04	85.03	72.93	83.33	97.39	95.88	87.79
	Class	71.77	71.92	89.46	87.69	94.43	89.84	95.45	90.88	49.30	74.74	89.48
Man $\rightarrow$ Woman	Target	96.11	94.53	92.18	91.74	85.66	90.96	77.55	90.29	98.16	97.38	90.23
	Class	83.68	83.64	87.45	87.43	90.05	88.52	91.27	89.11	75.50	84.04	90.30
Woman $\rightarrow$ Man	Target	93.34	92.26	87.32	88.09	78.32	84.28	64.86	81.89	94.96	92.91	88.59
	Class	82.39	82.68	92.33	91.51	93.89	92.59	94.35	93.22	72.03	84.56	89.66

Table 4: Target and source accuracy (%) for different domain translations on the FFHQ dataset. Dynamic subset selection is evaluated using Euler integration with 100 steps.

## 5 DISCUSSION AND CONCLUSION

Practically, one important matter of concern for the utility of Wasserstein distances is the fact that sample estimators of Wasserstein distances are cursed by dimensionality (Weed & Bach, 2019;



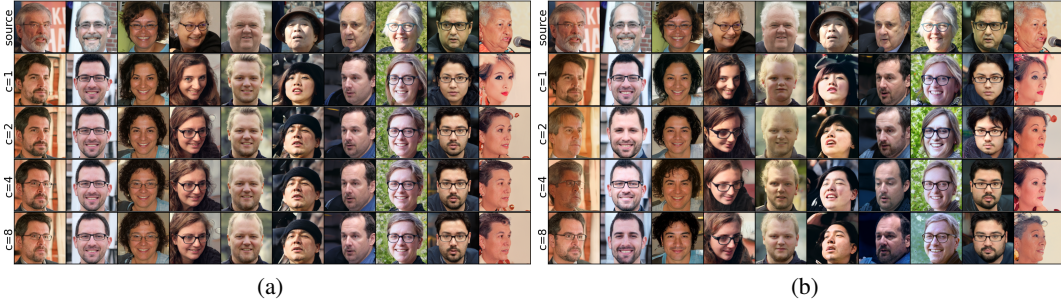


Figure 4: FFHQ old→young translation using (a) static and (b) dynamic subset selection. Dynamic subset selection is evaluated using Euler integration with 100 steps.

Fournier & Guillin, 2015), which can be alleviated to certain extent by employing the entropic regularization (Genevay et al., 2019; Feydy et al., 2019), which in the dynamic case is intimately connected with Schrödinger bridges.

Recently, unbalanced entropically-regularized optimal transport has been studied to model birth and death processes for population dynamics (Pariset et al., 2023; Neklyudov et al., 2023). Our approach can also be applied to model death processes, in cases where there is some canonical relationship between temporally ordered events, by treating  $\mu$  as the final population of survivors and  $\nu$  as the initial population.

Note the choice of  $c$  is often critical in applications. While  $c$  is interpretable, an automatic selection of  $c$  based on the resulting transport cost, which was previously conducted for partial optimal transport in the discrete case (Phatak et al., 2023), may be possible. One consolidated approach would be to sample  $c$  from a range and use multi-task learning for optimizing networks for varying  $c$ . In terms of implementation, this is possible using a scalar embedding of  $c$  as used for embeddings of the time variables in dynamic networks. We would further like to point out that one can replace range divergence with more common divergences like KL divergence but we cannot use the sign of potential in that to distinguish between selected and rejected subsets.

Finally, we note that although we focused on relaxing the target distribution; the range-divergence framework could potentially be adapted to also relax the source distribution. A fully relaxed version may be applicable to other classes of problems.

In conclusion, our approaches for neural optimal transport with subset selection are motivated by problems that require translation between two distribution with reweighting and selection of the target. The results here, limited to image translation tasks on two datasets and 20 tabular PU-learning tasks, show that both a meaningful subset can be learned simultaneously with a Monge map. Unlike previous work, our dynamic formulation of allows for variation in the terminal distribution from the original target marginal, creating flows to the nearest subset.

## REFERENCES

- Michael Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. In *The Eleventh International Conference on Learning Representations*, 2023.
- Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.
- Fabian Altekruiger, Johannes Hertrich, and Gabriele Steidl. Neural wasserstein gradient flows for discrepancies with riesz kernels. In *International Conference on Machine Learning*, pp. 664–690. PMLR, 2023.
- David Alvarez-Melis, Yair Schiff, and Youssef Mroueh. Optimizing functionals on the space of probabilities with input convex neural networks. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=dpOYN7o8Jm>.



- Brandon Amos. On amortizing convex conjugates for optimal transport. In *The Eleventh International Conference on Learning Representations*, 2023.
- Brandon Amos, Lei Xu, and J Zico Kolter. Input convex neural networks. In *International Conference on Machine Learning*, pp. 146–155. PMLR, 2017.
- Brandon Amos, Giulia Luise, Samuel Cohen, and Ievgen Redko. Meta optimal transport. In *International Conference on Machine Learning*, pp. 791–813. PMLR, 2023.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein Generative Adversarial Networks. In *International Conference on Machine Learning*, pp. 214–223. PMLR, 2017.
- Arip Asadulaev, Alexander Korotin, Vage Egiazarian, Petr Mokrov, and Evgeny Burnaev. Neural optimal transport with general cost functionals. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=gIiz7tBtYZ>.
- Lazar Atanackovic, Xi Zhang, Brandon Amos, Mathieu Blanchette, Leo J Lee, Yoshua Bengio, Alexander Tong, and Kirill Neklyudov. Meta flow matching: Integrating vector fields on the wasserstein manifold. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=9SYczU3Qgm>.
- Julio Backhoff-Veraguas, Mathias Beiglböck, and Gudmun Pammer. Existence, duality, and cyclical monotonicity for weak transport costs. *Calculus of Variations and Partial Differential Equations*, 58(6):203, 2019.
- Martin Bauer, Martins Bruveris, and Peter W Michor. Uniqueness of the fisher–rao metric on the space of smooth densities. *Bulletin of the London Mathematical Society*, 48(3):499–506, 2016.
- Jessa Bekker and Jesse Davis. Learning from Positive and Unlabeled Data: A Survey. *Machine Learning*, 109:719–760, 2020.
- Jean-David Benamou and Yann Brenier. A computational fluid mechanics solution to the monge-kantorovich mass transfer problem. *Numerische Mathematik*, 84(3):375–393, 2000.
- Yann Brenier. Polar factorization and monotone rearrangement of vector-valued functions. *Communications on pure and applied mathematics*, 44(4):375–417, 1991.
- Charlotte Bunne, Andreas Krause, and marco cuturi. Supervised training of conditional monge maps. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022a. URL <https://openreview.net/forum?id=sPNtVVUq7wi>.
- Charlotte Bunne, Laetitia Papaxanthos, Andreas Krause, and Marco Cuturi. Proximal optimal transport modeling of population dynamics. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pp. 6511–6528. PMLR, 28–30 Mar 2022b. URL <https://proceedings.mlr.press/v151/bunne22a.html>.
- Charlotte Bunne, Stefan G Stark, Gabriele Gut, Jacobo Sarabia Del Castillo, Mitch Levesque, Kjong-Van Lehmann, Lucas Pelkmans, Andreas Krause, and Gunnar Rätsch. Learning single-cell perturbation responses using neural optimal transport. *Nature methods*, 20(11):1759–1768, 2023.
- Luis A Caffarelli and Robert J McCann. Free Boundaries in Optimal Transport and Monge-Ampere Obstacle Problems. *Annals of Mathematics*, pp. 673–730, 2010.
- Laetitia Chapel, Mokhtar Z Alaya, and Gilles Gasso. Partial Optimal Transport with Applications on Positive-Unlabeled Learning. *Advances in Neural Information Processing Systems*, 33:2903–2913, 2020.
- Ricky TQ Chen and Yaron Lipman. Flow matching on general geometries. In *The Twelfth International Conference on Learning Representations*, 2024.



- Xiuyuan Cheng, Jianfeng Lu, Yixin Tan, and Yao Xie. Convergence of flow-based generative models via proximal gradient descent in wasserstein space. *IEEE Transactions on Information Theory*, 2024.
- Lenaic Chizat, Gabriel Peyré, Bernhard Schmitzer, and François-Xavier Vialard. An interpolating distance between optimal transport and fisher-rao metrics. *Foundations of Computational Mathematics*, 18:1–44, 2018a.
- Lenaic Chizat, Gabriel Peyré, Bernhard Schmitzer, and François-Xavier Vialard. Unbalanced optimal transport: Dynamic and kantorovich formulations. *Journal of Functional Analysis*, 274(11): 3090–3123, 2018b.
- Jaemoo Choi, Jaewoong Choi, and Myungjoo Kang. Generative modeling through the semi-dual formulation of unbalanced optimal transport. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=7WQt1Jl3ex>.
- Jaemoo Choi, Jaewoong Choi, and Myungjoo Kang. Analyzing and improving optimal-transport-based adversarial networks. In *The Twelfth International Conference on Learning Representations*, 2024a. URL <https://openreview.net/forum?id=jODEhvtTDx>.
- Jaemoo Choi, Jaewoong Choi, and Myungjoo Kang. Generative modeling through the semi-dual formulation of unbalanced optimal transport. *Advances in Neural Information Processing Systems*, 36, 2024b.
- Jaemoo Choi, Jaewoong Choi, and Myungjoo Kang. Scalable wasserstein gradient flow for generative modeling through unbalanced optimal transport. In *Forty-first International Conference on Machine Learning*, 2024c.
- Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*, pp. 2921–2926. IEEE, 2017.
- Max Daniels, Tyler Maunu, and PAul HAnd. Score-based generative neural networks for large-scale optimal transport. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=PPzV1H4atM4>.
- Li Deng. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012. doi: 10.1109/MSP.2012.2211477.
- Yuanqi Du, Michael Plainer, Rob Brekelmans, Chenru Duan, Frank Noe, Carla P Gomes, Alan Aspuru-Guzik, and Kirill Neklyudov. Doob’s lagrangian: A sample-efficient variational approach to transition path sampling. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=ShJWT0n7kX>.
- Luca Eyring, Dominik Klein, Théo Uscidda, Giovanni Palla, Niki Kilbertus, Zeynep Akata, and Fabian J Theis. Unbalancedness in neural monge maps improves unpaired domain translation. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=2UnCj3jeao>.
- Jiaojiao Fan, Shu Liu, Shaojun Ma, Yongxin Chen, and Hao-Min Zhou. Scalable computation of monge maps with general costs. In *ICLR Workshop on Deep Generative Models for Highly Structured Data*, 2022a. URL <https://openreview.net/forum?id=rEnGR3VdDW5>.
- Jiaojiao Fan, Qinsheng Zhang, Amirhossein Taghvaei, and Yongxin Chen. Variational wasserstein gradient flow. In *proceedings of international conference on machine learning*, 2022b.
- Jiaojiao Fan, Shu Liu, Shaojun Ma, Hao-Min Zhou, and Yongxin Chen. Neural monge map estimation and its applications. *Transactions on Machine Learning Research*, 2023.



- Jean Feydy, Thibault Séjourné, François-Xavier Vialard, Shun-ichi Amari, Alain Trounev, and Gabriel Peyré. Interpolating Between Optimal Transport and MMD using Sinkhorn Divergences. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2681–2690. PMLR, 2019.
- Alessio Figalli. The Optimal Partial Transport Problem. *Archive for Rational Mechanics and Analysis*, 195(2):533–560, 2010.
- Alessio Figalli and Federico Glaudo. *An invitation to optimal transport, Wasserstein distances, and gradient flows: Second Edition*. European Mathematical Society, 2023.
- Nicolas Fournier and Arnaud Guillin. On the rate of convergence in wasserstein distance of the empirical measure. *Probability theory and related fields*, 162(3):707–738, 2015.
- Wilfrid Gangbo and Robert J McCann. The geometry of optimal transportation. *Acta Mathematica*, 177:113–161, 1996.
- Milena Gazdieva, Alexander Korotin, Daniil Selikhanovych, and Evgeny Burnaev. Extremal domain translation with neural optimal transport. *Advances in Neural Information Processing Systems*, 36:40381–40413, 2023.
- Milena Gazdieva, Arip Asadulaev, Evgeny Burnaev, and Alexander Korotin. Light unbalanced optimal transport. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=co8KZws1YK>.
- Aude Genevay, Marco Cuturi, Gabriel Peyré, and Francis Bach. Stochastic optimization for large-scale optimal transport. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL [https://proceedings.neurips.cc/paper\\_files/paper/2016/file/2a27b8144ac02f67687f76782a3b5d8f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/2a27b8144ac02f67687f76782a3b5d8f-Paper.pdf).
- Aude Genevay, Gabriel Peyré, and Marco Cuturi. Learning Generative Models with Sinkhorn Divergences. In *International Conference on Artificial Intelligence and Statistics*, pp. 1608–1617. PMLR, 2018.
- Aude Genevay, Lénaïc Chizat, Francis Bach, Marco Cuturi, and Gabriel Peyré. Sample complexity of sinkhorn divergences. In *The 22nd international conference on artificial intelligence and statistics*, pp. 1574–1583. PMLR, 2019.
- Jonathan Geuter, Gregor Kornhardt, Ingimar Tomasson, and Vaios Laschos. Universal neural optimal transport. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=t10fde8tQ7>.
- Nathael Gozlan, Cyril Roberto, Paul-Marie Samson, and Prasad Tetali. Kantorovich duality for general transport costs and applications. *Journal of Functional Analysis*, 273(11):3327–3405, 2017.
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved Training of Wasserstein GANs. *Advances in Neural Information Processing Systems*, 30, 2017.
- Nikita Gushchin, Alexander Kolesov, Alexander Korotin, Dmitry P. Vetrov, and Evgeny Burnaev. Entropic neural optimal transport via diffusion processes. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a. URL <https://openreview.net/forum?id=fHyLsfMDIs>.
- Nikita Gushchin, Alexander Kolesov, Petr Mokrov, Polina Karpikova, Andrei Spiridonov, Evgeny Burnaev, and Alexander Korotin. Building the bridge of schrödinger: A continuous entropic optimal transport benchmark. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023b. URL <https://openreview.net/forum?id=OHimIaixXk>.



- Doron Haviv, Aram-Alexandre Pooladian, Dana Pe’er, and Brandon Amos. Wasserstein flow matching: Generative modeling over families of distributions. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=MRmI68k3gd>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Bamdad Hosseini, Alexander W Hsu, and Amirhossein Taghvaei. Conditional optimal transport on function spaces. *arXiv preprint arXiv:2311.05672*, 2023.
- Samuel Howard, George Deligiannidis, Patrick Rebeschini, and James Thornton. Differentiable cost-parameterized monge map estimators. In *ICML 2024 Workshop on Differentiable Almost Everything: Differentiable Relaxations, Algorithms, Operators, and Simulators*, 2024. URL <https://openreview.net/forum?id=UZ7lnFrwBt>.
- Chin-Wei Huang, Ricky T. Q. Chen, Christos Tsirigotis, and Aaron Courville. Convex potential flows: Universal probability distributions with optimal transport and convex optimization. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=te7PVHlsPxJ>.
- Guillaume Hugué, Daniel Sumner Magruder, Alexander Tong, Oluwadamilola Fasina, Manik Kuchroo, Guy Wolf, and Smita Krishnaswamy. Manifold interpolating optimal-transport flows for trajectory inference. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=ahAEhOtVif>.
- Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the fokker-planck equation. *SIAM Journal on Mathematical Analysis*, 29(1):1–17, 1998. doi: 10.1137/S0036141096303359. URL <https://doi.org/10.1137/S0036141096303359>.
- Leonid V Kantorovich. On the translocation of masses. In *Dokl. Akad. Nauk. USSR (NS)*, volume 37, pp. 199–201, 1942.
- Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4401–4410, 2019.
- Masahiro Kato, Takeshi Teshima, and Junya Honda. Learning from positive and unlabeled data with a selection bias. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rJzLciCqKm>.
- Markelle Kelly, Rachel Longjohn, and Kolby Nottingham. The UCI Machine Learning Repository. URL <https://archive.ics.uci.edu>.
- Gavin Kerrigan, Giosue Migliorini, and Padhraic Smyth. Dynamic conditional optimal transport through simulation-free flows. *arXiv preprint arXiv:2404.04240*, 2024.
- Valentin Khruikov, Gleb Ryzhakov, Andrei Chertkov, and Ivan Oseledets. Understanding DDPM latent codes through optimal transport. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=6PIrhAx1j4i>.
- Boah Kim, Yan Zhuang, Tejas Sudharshan Mathai, and Ronald M Summers. Otmorph: Unsupervised multi-domain abdominal medical image registration using neural optimal transport. *IEEE Transactions on Medical Imaging*, 2024.
- Diederik P Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Nikita Kornilov, Alexander Gasnikov, and Alexander Korotin. Optimal flow matching: Learning straight trajectories in just one step. *arXiv preprint arXiv:2403.13117*, 2024.



- Alexander Korotin, Vage Egiazarian, Arip Asadulaev, Alexander Safin, and Evgeny Burnaev. Wasserstein-2 generative networks. In *International Conference on Learning Representations*, 2021a. URL [https://openreview.net/forum?id=bEoxzW\\_EXsa](https://openreview.net/forum?id=bEoxzW_EXsa).
- Alexander Korotin, Lingxiao Li, Aude Genevay, Justin M Solomon, Alexander Filippov, and Evgeny Burnaev. Do neural optimal transport solvers work? a continuous wasserstein-2 benchmark. *Advances in Neural Information Processing Systems*, 34:14593–14605, 2021b.
- Alexander Korotin, Daniil Selikhanovych, and Evgeny Burnaev. Kernel neural optimal transport. In *The Eleventh International Conference on Learning Representations*, 2023a. URL [https://openreview.net/forum?id=Zuc\\_MHtUma4](https://openreview.net/forum?id=Zuc_MHtUma4).
- Alexander Korotin, Daniil Selikhanovych, and Evgeny Burnaev. Neural Optimal Transport. In *The Eleventh International Conference on Learning Representations*, 2023b.
- Dohyun Kwon, Ying Fan, and Kangwook Lee. Score-based generative modeling secretly minimizes the wasserstein distance. In *Advances in Neural Information Processing Systems*, volume 35, pp. 20205–20217. Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/7f52f6b8f107931127eefe15429ee278-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/7f52f6b8f107931127eefe15429ee278-Paper-Conference.pdf).
- Marc Lambert, Sinho Chewi, Francis Bach, Silvère Bonnabel, and Philippe Rigollet. Variational inference via wasserstein gradient flows. In *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=K2PTuvVTF1L>.
- Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023.
- Guan-Horng Liu, Tianrong Chen, Oswin So, and Evangelos Theodorou. Deep generalized schrödinger bridge. *Advances in Neural Information Processing Systems*, 35:9374–9388, 2022.
- Guan-Horng Liu, Yaron Lipman, Maximilian Nickel, Brian Karrer, Evangelos Theodorou, and Ricky T. Q. Chen. Generalized schrödinger bridge matching. In *The Twelfth International Conference on Learning Representations*, 2024a. URL <https://openreview.net/forum?id=SoismgeX7z>.
- Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=XVjTT1nw5z>.
- Xingchao Liu, Xiwen Zhang, Jianzhu Ma, Jian Peng, and qiang liu. InstafLOW: One step is enough for high-quality diffusion-based text-to-image generation. In *The Twelfth International Conference on Learning Representations*, 2024b. URL <https://openreview.net/forum?id=1k4yZbbDqX>.
- Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Skq89Scxx>.
- Frederike Lübeck, Charlotte Bunne, Gabriele Gut, Jacobo Sarabia del Castillo, Lucas Pelkmans, and David Alvarez-Melis. Neural unbalanced optimal transport via cycle-consistent semi-couplings. In *NeurIPS 2022 AI for Science: Progress and Promises*, 2022. URL <https://openreview.net/forum?id=5lflxpNymZr>.
- Shaojun Ma, Shu Liu, Hongyuan Zha, and Haomin Zhou. Learning stochastic behaviour from aggregate data. In *International Conference on Machine Learning*, pp. 7258–7267. PMLR, 2021.
- Ashok Makkuva, Amirhossein Taghvaei, Sewoong Oh, and Jason Lee. Optimal transport mapping via input convex neural networks. In *International Conference on Machine Learning*, pp. 6672–6681. PMLR, 2020.
- Robert J McCann. A convexity principle for interacting gases. *Advances in mathematics*, 128(1): 153–179, 1997.



- Petr Mokrov, Alexander Korotin, Lingxiao Li, Aude Genevay, Justin Solomon, and Evgeny Burnaev. Large-scale wasserstein gradient flows. In *Advances in Neural Information Processing Systems*, 2021.
- Petr Mokrov, Alexander Korotin, Alexander Kolesov, Nikita Gushchin, and Evgeny Burnaev. Energy-guided entropic neural optimal transport. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=d6tUsZeVs7>.
- Gaspard Monge. Mémoire sur la théorie des déblais et des remblais. *Mem. Math. Phys. Acad. Royale Sci.*, pp. 666–704, 1781.
- Kirill Neklyudov, Rob Brekelmans, Daniel Severo, and Alireza Makhzani. Action matching: Learning stochastic dynamics from samples. In *International conference on machine learning*, pp. 25858–25889. PMLR, 2023.
- Kirill Neklyudov, Rob Brekelmans, Alexander Tong, Lazar Atanackovic, Qiang Liu, and Alireza Makhzani. A computational framework for solving Wasserstein lagrangian flows. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 37461–37485. PMLR, 21–27 Jul 2024a. URL <https://proceedings.mlr.press/v235/neklyudov24a.html>.
- Kirill Neklyudov, Jannes Nys, Luca Thiede, Juan Carrasquilla, Qiang Liu, Max Welling, and Alireza Makhzani. Wasserstein quantum monte carlo: a novel approach for solving the quantum many-body schrödinger equation. *Advances in Neural Information Processing Systems*, 36, 2024b.
- XuanLong Nguyen, Martin J Wainwright, and Michael I Jordan. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847–5861, 2010.
- Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-gan: Training generative neural samplers using variational divergence minimization. *Advances in neural information processing systems*, 29, 2016.
- Matteo Pariset, Ya-Ping Hsieh, Charlotte Bunne, Andreas Krause, and Valentin De Bortoli. Unbalanced diffusion schrödinger bridge. *arXiv preprint arXiv:2306.09099*, 2023.
- Abhijeet Phatak, Sharath Raghvendra, Chittaranjan Tripathy, and Kaiyi Zhang. Computing all optimal partial transports. In *The Eleventh International Conference on Learning Representations*, 2023.
- Stanislav Pidhorskyi, Donald A Adjeroh, and Gianfranco Doretto. Adversarial latent autoencoders. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 14104–14113, 2020.
- Yury Polyanskiy and Yihong Wu. *Information theory: From coding to learning*. Cambridge university press, 2025.
- Aram-Alexandre Pooladian, Heli Ben-Hamu, Carles Domingo-Enrich, Brandon Amos, Yaron Lipman, and Ricky T Chen. Multisample flow matching: Straightening flows with minibatch couplings. *ICML 2023*, 2023.
- Aram-Alexandre Pooladian, Carles Domingo-Enrich, Ricky T. Q. Chen, and Brandon Amos. Neural optimal transport with lagrangian costs. In *The 40th Conference on Uncertainty in Artificial Intelligence*, 2024. URL <https://openreview.net/forum?id=x4paJ2sJyZ>.
- Bilal Riaz, Yuksel Karahan, and Austin J. Brockmeier. Partial optimal transport for support subset selection. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=75CcopPxIr>.
- R Tyrrell Rockafellar. Integral functionals, normal integrands and measurable selections. *Lecture Notes in Mathematics*, pp. 157–207, 1976.



- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III* 18, pp. 234–241. Springer, 2015.
- Litu Rout, Alexander Korotin, and Evgeny Burnaev. Generative Modeling with Optimal Transport Maps. In *International Conference on Learning Representations*, 2022.
- Filippo Santambrogio. Optimal transport for applied mathematicians. *Birkhäuser, NY*, 55(58–63):94, 2015.
- Filippo Santambrogio. {Euclidean, metric, and Wasserstein} gradient flows: an overview. *Bulletin of Mathematical Sciences*, 7:87–154, 2017.
- Christopher Scovel and Justin Solomon. Riemannian metric learning via optimal transport. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=v3y68gz-WEz>.
- Vivien Seguy, Bharath Bhushan Damodaran, Remi Flamary, Nicolas Courty, Antoine Rolet, and Mathieu Blondel. Large scale optimal transport and mapping estimation. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=Blzlp1bRW>.
- Thibault Séjourné, Gabriel Peyré, and François-Xavier Vialard. Unbalanced optimal transport, from theory to numerics. *Numerical Control: Part B*, pp. 407, 2023.
- Yuyang Shi, Valentin De Bortoli, Andrew Campbell, and Arnaud Doucet. Diffusion schrödinger bridge matching. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=qy07OHsJT5>.
- Yuyang Shi, Valentin De Bortoli, Andrew Campbell, and Arnaud Doucet. Diffusion schrödinger bridge matching. *Advances in Neural Information Processing Systems*, 36, 2024.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2256–2265, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/sohl-dickstein15.html>.
- Vignesh Ram Somnath, Matteo Pariset, Ya-Ping Hsieh, Maria Rodriguez Martinez, Andreas Krause, and Charlotte Bunne. Aligned diffusion schrödinger bridges. In *The 39th Conference on Uncertainty in Artificial Intelligence*, 2023. URL [https://openreview.net/forum?id=BkWFJN7\\_bQ](https://openreview.net/forum?id=BkWFJN7_bQ).
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=StlgiaRCHLP>.
- Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020.
- Paweł Teisseyre, Timo Martens, Jessa Bekker, and Jesse Davis. Learning from biased positive-unlabeled data via threshold calibration. In *The 28th International Conference on Artificial Intelligence and Statistics*, 2025. URL <https://openreview.net/forum?id=dT0ldWDBto>.
- Alexander Tong, Jessie Huang, Guy Wolf, David Van Dijk, and Smita Krishnaswamy. Trajectorynet: A dynamic optimal transport network for modeling cellular dynamics. In *International conference on machine learning*, pp. 9526–9536. PMLR, 2020.
- 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. In *ICML Workshop on New Frontiers in Learning, Control, and Dynamical Systems*, 2023.



- Alexander Tong, Kilian Fatras, Nikolay Malkin, Guillaume Huguet, Yanlei Zhang, Jarrod Rector-Brooks, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models with minibatch optimal transport. *Transactions on Machine Learning Research*, 2024a. ISSN 2835-8856. URL <https://openreview.net/forum?id=CD9Snc73AW>. Expert Certification.
- Alexander Y. 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. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pp. 1279–1287. PMLR, 02–04 May 2024b. URL <https://proceedings.mlr.press/v238/tong24a.html>.
- Théo Uscidda and Marco Cuturi. The monge gap: A regularizer to learn all transport maps. In *International Conference on Machine Learning*, pp. 34709–34733. PMLR, 2023.
- Nina Vesseron and Marco Cuturi. On a neural implementation of brenier’s polar factorization. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=zDCwJQY3eI>.
- Cédric Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2009.
- Wei Wan, Yuejin Zhang, Chenglong Bao, Bin Dong, and Zuoqiang Shi. A scalable deep learning approach for solving high-dimensional dynamic optimal transport. *SIAM Journal on Scientific Computing*, 45(4):B544–B563, 2023.
- Gefei Wang, Yuling Jiao, Qian Xu, Yang Wang, and Can Yang. Deep generative learning via schrödinger bridge. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 10794–10804. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/wang21l.html>.
- Jonathan Weed and Francis Bach. Sharp asymptotic and finite-sample rates of convergence of empirical measures in wasserstein distance. *Bernoulli*, 25(4A):pp. 2620–2648, 2019. ISSN 13507265, 15739759. URL <https://www.jstor.org/stable/48586009>.
- Yifan Wu, Ezra Winston, Divyansh Kaushik, and Zachary Lipton. Domain adaptation with asymmetrically-relaxed distribution alignment. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6872–6881. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/wu19f.html>.
- Chen Xu, Xiuyuan Cheng, and Yao Xie. Normalizing flow neural networks by jko scheme. In *Advances in Neural Information Processing Systems*, volume 36, pp. 47379–47405. Curran Associates, Inc., 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/93fce71def4e3cf418918805455d436f-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/93fce71def4e3cf418918805455d436f-Paper-Conference.pdf).
- Karren D. Yang and Caroline Uhler. Scalable unbalanced optimal transport using generative adversarial networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=HyexAiA5Fm>.

## A PRELIMINARIES AND PROBLEM FORMULATION

Kantorovich (1942) reformulated the Monge problem by relaxing the constraint that supports of  $\mu$  and  $\nu$  should be related to each other by a functional relation  $T$ . Instead, he allowed  $\mu$  and  $\nu$  to be related to each other by a joint measure. Kantorovich’s reformulation of the problem is a linear program and its solution exists for all convex lower-semi-continuous costs. Santambrogio (2015); Figalli & Glaudo (2023). The Kantorovich problem is,

$$\begin{aligned} \mathcal{W}(\mu, \nu) = \inf_{\pi} \int_{\mathcal{X} \times \mathcal{Y}} c(\mathbf{x}, \mathbf{y}) \pi(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} \\ \text{s.t. } \int_{\mathcal{Y}} d\pi(\mathbf{x}, \mathbf{y}) = \mu(\mathbf{y}), \quad \int_{\mathcal{X}} d\pi(\mathbf{x}, \mathbf{y}) = \nu(\mathbf{x}), \end{aligned} \quad (12)$$



where it can be observed that integrals  $\int_{\mathcal{X}}$  and  $\int_{\mathcal{Y}}$  marginalize with respect to spaces  $\mathcal{Y}$  and  $\mathcal{X}$ , respectively. Therefore,  $\pi$  must be the joint measure between  $\mu$  and  $\nu$  defined on the product space  $\mathcal{X} \times \mathcal{Y}$ , i.e.  $\pi \in \mathcal{P}(\mathcal{X} \times \mathcal{Y})$ . In other words, constraints in the Kantorovich problem ensure that every feasible  $\pi$  must be joint distribution of  $\mu$  and  $\nu$ . While in general, Kantorovich problem is much easier to solve in comparison to Monge problem, there are the conditions of practical importance where one can employ the solutions of Kantorovich problem to obtain the solution of Monge problem. Those conditions are more clearly discussed in terms of dual form of Kantorovich problem (Santambrogio, 2015), given as,

$$\begin{aligned} \mathcal{W}(\mu, \nu) = \sup_{f, g} \int_{\mathcal{X}} f(\mathbf{x})\mu(\mathbf{x})d\mathbf{x} + \int_{\mathcal{Y}} g(\mathbf{y})\nu(\mathbf{y})d\mathbf{y} \\ \text{s.t. } f(\mathbf{x}) + g(\mathbf{y}) \leq c(\mathbf{x}, \mathbf{y}) \end{aligned} \quad (13)$$

The functions  $f(\mathbf{x})$  and  $g(\mathbf{y})$  are called Kantorovich potentials. By defining  $c$ -conjugate (also called  $c$ -transform) of  $f(\mathbf{x})$ , and  $\bar{c}$ -conjugate (also called  $\bar{c}$ -transform) of  $g(\mathbf{y})$  as

$$f^c(\mathbf{y}) = \inf_{\mathbf{x} \in \mathcal{X}} c(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}), \quad (14)$$

$$g^{\bar{c}}(\mathbf{x}) = \inf_{\mathbf{y} \in \mathcal{Y}} c(\mathbf{x}, \mathbf{y}) - g(\mathbf{y}), \quad (15)$$

Using  $c$  and  $\bar{c}$  conjugates, Kantorovich problem is expressed as

$$\mathcal{W}(\mu, \nu) = \sup_f \int_{\mathcal{X}} f(\mathbf{x})\mu(\mathbf{x})d\mathbf{x} + \int_{\mathcal{Y}} f^c(\mathbf{y})\nu(\mathbf{y})d\mathbf{y} \quad (16)$$

$$= \sup_g \int_{\mathcal{X}} g^{\bar{c}}(\mathbf{x})\mu(\mathbf{x})d\mathbf{x} + \int_{\mathcal{Y}} g(\mathbf{y})\nu(\mathbf{y})d\mathbf{y} \quad (17)$$

Under very general conditions, one can relate the cost  $c$  with the support of optimal coupling solution  $\pi^*$  and optimal Kantorovich potentials  $f^*(\mathbf{x})$  and  $f^{c*}(\mathbf{y})$  (Santambrogio, 2015, Theorem 1.37) by

$$\text{supp}(\pi^*) \subset \{(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y} : f^*(\mathbf{x}) + f^{c*}(\mathbf{y}) = c(\mathbf{x}, \mathbf{y})\} \quad (18)$$

In discrete domains, above result is equivalent to Karush-Kuhn-Tucker (KKT) conditions for optimality. One can further relate the optimal solutions of Monge and Kantorovich problems using a landmark result by Gangbo & McCann (1996), (Figalli & Glaudo, 2023, Theorem 2.7.1), which states that there exists an optimal Kantorovich coupling of the form  $\pi^* = (\text{Id} \times T^*)_{\#}\mu$ , where  $T^*$  is Monge map satisfying

$$\nabla_{\mathbf{x}}c(\mathbf{x}, T^*(\mathbf{x})) + \nabla f^*(\mathbf{x}) = 0, \quad (19)$$

if the following conditions are satisfied

- $\mu$  is absolutely continuous,
- $\forall \mathbf{y} \in \mathcal{Y}$  the map  $\mathbf{x} \mapsto c(\mathbf{x}, \mathbf{y})$  is differentiable,  $\forall \mathbf{x} \in \mathcal{X}$ ,
- $\forall \mathbf{x} \in \mathcal{X}$  the gradient map  $\mathbf{y} \mapsto \nabla_{\mathbf{x}}c(\mathbf{x}, \mathbf{y})$  is injective  $\forall \mathbf{y} \in \mathcal{Y}$ ,
- and the gradient  $\nabla_{\mathbf{x}}c(\mathbf{x}, \mathbf{y})$  satisfies the local Lipschitz condition  $\|\nabla_{\mathbf{x}}c(\mathbf{x}, \mathbf{y})\| \leq C_r$  for all  $\mathbf{x} \in \mathcal{B}_r$ , where  $\mathcal{B}_r$  is ball of radius  $r$  around  $\mathbf{x}$ .

When the cost can be written as  $c(\mathbf{x}, \mathbf{y}) = h(\mathbf{x} - \mathbf{y})$ , where  $h$  is strictly convex and translation invariant function, one can further relate the Monge mapping with optimal dual potential by (Santambrogio, 2015, Theorem 1.17)

$$T^*(\mathbf{x}) = \mathbf{x} - \nabla h^* \circ \nabla f^*(\mathbf{x}), \quad (20)$$

where  $h^*$  is Legendre-Fenchel conjugate of  $h$  given by  $h^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathcal{X}} \{\langle \mathbf{y}, \mathbf{x} \rangle - h(\mathbf{x})\}$ . In the result

above, when  $h(\mathbf{x} - \mathbf{y}) = \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|_2^2$ , one obtains the result of celebrated Brenier theorem of optimal transport for squared-Euclidean costs with  $T^*(\mathbf{x}) = \nabla f^*(\mathbf{x})$ , where  $f^*(\mathbf{x})$  is convex (Brenier, 1991), (Figalli & Glaudo, 2023, Theorem 2.5.10). The Brenier theorem on optimal transport differs from another important theorem on polar factorization (Brenier, 1991) stating that under very general conditions a square integrable vector field  $v$  can be decomposed into the composition of gradient of a unique convex function  $\xi$  and a unique measure-preserving map  $u$ , i.e.  $v(\mathbf{x}) = \nabla \xi \circ u(\mathbf{x})$ . Before



the discussion on dynamic formulation of the problem, we would like to point out that much of the recent work on static neural optimal transport rely on above results.

Benamou and Brenier formulated the Wasserstein distance with squared Euclidean cost as the kinetic energy minimization problem under the assumption that both the source  $\mu$  and target  $\nu$  distributions have finite second moments (Benamou & Brenier, 2000; Santambrogio, 2015; Figalli & Glaudo, 2023). Assuming that supports of both source and target distributions lie in a convex set  $\Omega \subseteq \mathbb{R}^d$ , whose normal at the boundary is given by  $\mathbf{n} : \partial\Omega \rightarrow \mathbb{R}^d$ , for a bounded and smooth velocity-field  $\mathbf{v}_t(\mathbf{x}) : [0, 1] \times \Omega \rightarrow \mathbb{R}^d$ , such that  $\langle \mathbf{v}_t(\mathbf{x}), \mathbf{n} \rangle|_{\partial\Omega} = 0$ , the flow corresponding to  $\mathbf{v}_t$  is given by

$$\frac{d}{dt}\Phi_t(\mathbf{x}_t) = \mathbf{v}_t(\Phi_t(\mathbf{x}_t)), \quad \Phi_0(\mathbf{x}_0) = \mathbf{x}_0. \quad (21)$$

Considering that there also exists a probability path  $\rho_t(\mathbf{x}) : [0, 1] \times \Omega \rightarrow \mathbb{R}_+$ , corresponding to the flow  $\Phi_t(\mathbf{x})$  such that  $\rho_t(\mathbf{x}) = \Phi_{t\#}\rho_0(\mathbf{x})$ , Benamou-Brenier formulation of optimal transport is

$$\begin{aligned} \inf_{\rho_t, \mathbf{v}_t} \int_0^1 \int_{\Omega} \frac{\|\mathbf{v}_t(\mathbf{x})\|^2}{2} \rho_t(\mathbf{x}) d\mathbf{x} dt \\ \text{s.t. } \frac{\partial}{\partial t} \rho_t(\mathbf{x}) + \text{div}(\rho_t(\mathbf{x}) \mathbf{v}_t(\mathbf{x})) = 0, \quad \rho_0(\mathbf{x}) = \mu(\mathbf{x}), \quad \rho_1(\mathbf{x}) = \nu(\mathbf{x}), \end{aligned} \quad (22)$$

where  $\text{div}(\cdot)$  denotes divergence operator mapping scalar or vector fields to scalar, for the field  $z_t(\mathbf{x})$  by  $\text{div}(z_t(\mathbf{x})) = \sum_i \frac{\partial}{\partial x_i} z_t(\mathbf{x})$ . The optimal flow  $\Phi_t^*$  is related to Monge mapping  $T^*$  by displacement interpolation (McCann, 1997).

$$\Phi_t^* = (1 - t)\text{Id} + tT^* \quad (23)$$

It is important to mention that the Benamou-Brenier formulation can be extended to Wasserstein- $p$  distances, for  $p > 1$ , under the assumption that both source and target distributions have finite  $p$ -th moments (Santambrogio, 2015, chapters 5 & 6).

#### A.1 DERIVATION OF STATIC NEURAL SUBSET SELECTION

We denote the problem expressed in 6 as  $\inf_{\pi} \sup_{\psi} \sup_{\eta} \mathcal{L}(\pi, \psi, \eta)$ , where the Lagrangian is

$$\mathcal{L}(\pi, \psi, \eta) = \int_{\mathcal{X} \times \mathcal{Y}} (c(\mathbf{x}, \mathbf{y}) + \eta(\mathbf{y}) - \psi(\mathbf{x})) \pi(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} + \int_{\mathcal{X}} \psi(\mathbf{x}) \mu(\mathbf{x}) d\mathbf{x} \quad (24)$$

$$- c \int_{\mathcal{Y}} \eta_+(\mathbf{y}) \nu(\mathbf{y}) d\mathbf{y}. \quad (25)$$

We proceed to interchange the sup and inf,<sup>3</sup> which is allowed due to the strong duality property associated with optimal transport when the cost  $c$  is convex and lower semi-continuous,

$$\inf_{\pi} \sup_{\psi, \eta} \mathcal{L}(\pi, \psi, \eta) = \sup_{\eta, \psi} \inf_{\pi} \mathcal{L}(\pi, \psi, \eta). \quad (26)$$

Optimizing with respect to  $\pi$  for given  $\eta$  and  $\psi$ , the integrand in the first term of 24 is unbounded from below at any point  $c(\mathbf{x}, \mathbf{y}) + \eta(\mathbf{y}) - \psi(\mathbf{x}) < 0$ . Thus,  $\eta$  and  $\psi$  need to ensure that  $c(\mathbf{x}, \mathbf{y}) + \eta(\mathbf{y}) - \psi(\mathbf{x}) \geq 0$ . This constraint requires for any  $\mathbf{x} \in \text{supp}(\mu) \subseteq \mathcal{X}$ ,  $\psi(\mathbf{x}) = \inf_{\mathbf{y} \in \mathcal{Y}} (c(\mathbf{x}, \mathbf{y}) + \eta(\mathbf{y}))$  (Villani et al., 2009) (Theorem 5.10 and Remark 5.13). This definition of  $\psi$  corresponds to the  $c$ -transform of  $-\eta(\mathbf{y})$  in the optimal transport literature (Santambrogio, 2015; Villani et al., 2009). Then, the inner infimum with respect to  $\pi$  is attained with zero value if  $\forall (\mathbf{x}, \mathbf{y}) \in (\mathcal{X} \times \mathcal{Y})$ :

<sup>3</sup>This requires us to verify Slater's constraint qualifications, which are: (i) Primal is convex wrt  $\pi$ , (which is obvious), (ii) Dual is concave wrt  $\eta$ , which is also obvious (iii) relative interior for inequality constraints set is non-empty, which can be verified by looking at the fact that for any  $\tilde{\nu}$  the distribution  $\pi(\mathbf{x}, \mathbf{y}) = \mu(\mathbf{x})\tilde{\nu}(\mathbf{y})$  is feasible and one can see that if one defines the feasible set of coupling  $\Pi_c(\mu, \nu) = \{\pi \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}) : \pi_{\mathcal{X}} = \mu, \pi_{\mathcal{Y}} \leq c\nu\}$ , then for  $\forall 1 \leq c_0 \leq c_1$ ,  $\Pi_{c_0}(\mu, \nu) \subseteq \Pi_{c_1}(\mu, \nu)$ , which in other words mean that  $\Pi_{c=1}(\mu, \nu)$  is a subset of feasible solutions for all values of  $c > 1$ , therefore relative-interior is non-empty.



$\pi(\mathbf{x}, \mathbf{y}) > 0 \implies c(\mathbf{x}, \mathbf{y}) + \eta(\mathbf{y}) - \psi(\mathbf{x}) = 0$ . Therefore, the dual problem becomes

$$\sup_{\eta} \int_{\mathcal{X}} \left( \inf_{\mathbf{y} \in \mathcal{Y}} (c(\mathbf{x}, \mathbf{y}) + \eta(\mathbf{y})) \right) \mu(\mathbf{x}) d\mathbf{x} - c \int_{\mathcal{Y}} \eta_+(\mathbf{y}) \nu(\mathbf{y}) d\mathbf{y}. \quad (27)$$

$$= \sup_{\eta} \inf_T \int_{\mathcal{X}} (c(\mathbf{x}, T(\mathbf{x})) + \eta(T(\mathbf{x}))) \mu(\mathbf{x}) d\mathbf{x} - c \int_{\mathcal{Y}} \eta_+(\mathbf{y}) \nu(\mathbf{y}) d\mathbf{y} \quad (28)$$

$$= \sup_{\eta} \inf_T \mathbb{E}_{\mathbf{x} \sim \mu} [c(\mathbf{x}, T(\mathbf{x})) + \eta(T(\mathbf{x}))] - c \mathbb{E}_{\mathbf{y} \sim \nu} [\eta_+(\mathbf{y})]. \quad (29)$$

$$= \inf_T \sup_{\eta} \mathbb{E}_{\mathbf{x} \sim \mu} [c(\mathbf{x}, T(\mathbf{x}))] + \mathbb{E}_{\mathbf{x} \sim \mu} [\eta(T(\mathbf{x}))] - c \mathbb{E}_{\mathbf{y} \sim \nu} [\eta_+(\mathbf{y})]. \quad (30)$$

$$= \inf_T \mathbb{E}_{\mathbf{x} \sim \mu} [c(\mathbf{x}, T(\mathbf{x}))] + \mathcal{D}_{i_{[0, c]}}(T_{\#}\mu \| \nu) = \inf_{T_{\#}\mu \leq c\nu} \mathbb{E}_{\mathbf{x} \sim \mu} [c(\mathbf{x}, T(\mathbf{x}))]. \quad (31)$$

Equation 27 and equation 28 are equal due to a theorem by (Rockafellar, 1976, Theorem 3A); equation 28 and equation 29 are equivalent by definition; equation 29 and equation 30 are equivalent since the function is convex with respect to  $T$  and concave with respect to  $\eta$ ; and equation 30 and equation 31 are equivalent by the variational formula for the  $f$ -divergence. Thus, we obtain a relaxed Monge formulation.

## A.2 DERIVATION OF DYNAMIC NEURAL SUBSET SELECTION

Combining the objective and constraints in 8 to obtain the Lagrangian

$$\begin{aligned} \mathcal{L}(\mathbf{v}_t, \rho_t, \psi_0, \varphi_t, \eta) = & \overbrace{\int_0^1 \int_{\Omega} \frac{\|\mathbf{v}_t(\mathbf{x})\|^2}{2} \rho_t(\mathbf{x}) d\mathbf{x} dt}^{(I)} + \overbrace{\int_{\Omega} \psi_0(\mathbf{x}) (\rho_0(\mathbf{x}) - \mu(\mathbf{x})) d\mathbf{x}}^{(II)} \\ & + \overbrace{\int_0^1 \int_{\Omega} \varphi_t(\mathbf{x}) \frac{\partial}{\partial t} \rho_t(\mathbf{x}) d\mathbf{x} dt}^{(III)} + \overbrace{\int_0^1 \int_{\Omega} \varphi_t(\mathbf{x}) \operatorname{div}(\rho_t(\mathbf{x}) \mathbf{v}_t(\mathbf{x})) d\mathbf{x} dt}^{(IV)} \\ & + \overbrace{\sup_{\eta} \left( \int_{\Omega} \eta(\mathbf{x}) \rho_1(\mathbf{x}) d\mathbf{x} - c \int_{\Omega} \max(0, \eta(\mathbf{x})) \nu(\mathbf{x}) d\mathbf{x} \right)}^{(V)}. \end{aligned} \quad (32)$$

Since  $\rho_t(\mathbf{x})$  is supported on bounded subset  $\Omega \subset \mathbb{R}^d$ , one can change the order of integration. Therefore, for term (III), by changing the order of integration and then computing integration by parts one obtains,

$$\begin{aligned} (III) = & \int_{\Omega} \int_0^1 \varphi_t(\mathbf{x}) \frac{\partial}{\partial t} \rho_t(\mathbf{x}) d\mathbf{x} dt = \int_{\Omega} \varphi_1(\mathbf{x}) \rho_1(\mathbf{x}) d\mathbf{x} - \int_{\Omega} \varphi_0(\mathbf{x}) \rho_0(\mathbf{x}) d\mathbf{x} \\ & - \int_0^1 \int_{\Omega} \frac{\partial}{\partial t} \varphi_t(\mathbf{x}) \rho_t(\mathbf{x}) d\mathbf{x} dt. \end{aligned} \quad (33)$$

In order to simplify (IV), we can use product rule of derivatives to write

$$\varphi_t(\mathbf{x}) \operatorname{div}(\rho_t(\mathbf{x}) \mathbf{v}_t(\mathbf{x})) = \operatorname{div}(\varphi_t(\mathbf{x}) \rho_t(\mathbf{x}) \mathbf{v}_t(\mathbf{x})) - \rho_t(\mathbf{x}) \langle \nabla \varphi_t(\mathbf{x}), \mathbf{v}_t(\mathbf{x}) \rangle.$$

Therefore by combining above identity with Gauss's theorem one obtains

$$\begin{aligned} (IV) = & \int_0^1 \int_{\Omega} \operatorname{div}(\varphi_t(\mathbf{x}) \rho_t(\mathbf{x}) \mathbf{v}_t(\mathbf{x})) d\mathbf{x} dt - \int_0^1 \int_{\Omega} \rho_t(\mathbf{x}) \langle \nabla \varphi_t(\mathbf{x}), \mathbf{v}_t(\mathbf{x}) \rangle d\mathbf{x} dt \\ = & \underbrace{\int_0^1 \oint_{\partial\Omega} \varphi_t(\mathbf{x}_t) \rho_t(\mathbf{x}_t) \langle \mathbf{v}_t(\mathbf{x}), d\mathbf{n} \rangle dt}_{=0} - \int_0^1 \int_{\Omega} \rho_t(\mathbf{x}) \langle \nabla \varphi_t(\mathbf{x}), \mathbf{v}_t(\mathbf{x}) \rangle d\mathbf{x} dt. \end{aligned} \quad (34)$$

From the boundary condition on optimal transport (see the discussion above Equation 21, also Figalli & Glaudo (2023)-section 4.1), the first part of the right-hand side of 35 is zero; therefore,

$$(IV) = \int_0^1 \int_{\Omega} \varphi_t(\mathbf{x}) \cdot \operatorname{div}(\rho_t(\mathbf{x}) \mathbf{v}_t(\mathbf{x})) d\mathbf{x} dt = - \int_0^1 \int_{\Omega} \langle \nabla \varphi_t(\mathbf{x}), \mathbf{v}_t(\mathbf{x}) \rangle \rho_t(\mathbf{x}) d\mathbf{x} dt. \quad (35)$$



In order to eliminate primal variable  $\mathbf{v}_t(\mathbf{x})$ , substitute Equation 33 and Equation 35 into 32 and compute the variational-derivative to obtain the stationary condition. For that one can write the terms of Lagrangian depending on  $\mathbf{v}_t(\mathbf{x})$  as

$$\tilde{\mathcal{L}}(\mathbf{v}_t) = \int_0^1 \int_{\Omega} \left( \frac{\|\mathbf{v}_t(\mathbf{x})\|^2}{2} - \langle \nabla \varphi_t(\mathbf{x}), \mathbf{v}_t(\mathbf{x}) \rangle \right) \rho_t(\mathbf{x}) d\mathbf{x} dt. \quad (36)$$

With the additive perturbation function  $\boldsymbol{\tau}_t$  vanishing at  $t = 0$  and  $t = 1$  and a scalar  $\varepsilon$ , the Lagrangian  $\tilde{\mathcal{L}}(\mathbf{v}_t + \varepsilon \boldsymbol{\tau}_t)$  is

$$\begin{aligned} \tilde{\mathcal{L}}(\mathbf{v}_t + \varepsilon \boldsymbol{\tau}_t) &= \int_0^1 \int_{\Omega} \left( \frac{\|\mathbf{v}_t(\mathbf{x}) + \varepsilon \boldsymbol{\tau}_t(\mathbf{x})\|^2}{2} - \langle \nabla \varphi_t(\mathbf{x}), \mathbf{v}_t(\mathbf{x}) + \varepsilon \boldsymbol{\tau}_t(\mathbf{x}) \rangle \right) \rho_t(\mathbf{x}) d\mathbf{x} dt \\ &= \int_0^1 \int_{\Omega} \left( \frac{\|\mathbf{v}_t(\mathbf{x})\|^2}{2} - \langle \nabla \varphi_t(\mathbf{x}), \mathbf{v}_t(\mathbf{x}) \rangle \right) \rho_t(\mathbf{x}) d\mathbf{x} dt \\ &\quad + \int_0^1 \int_{\Omega} \left( \varepsilon^2 \frac{\|\boldsymbol{\tau}_t(\mathbf{x})\|^2}{2} + \varepsilon \langle \mathbf{v}_t(\mathbf{x}) - \nabla \varphi_t(\mathbf{x}), \boldsymbol{\tau}_t(\mathbf{x}) \rangle \right) \rho_t(\mathbf{x}) d\mathbf{x} dt, \end{aligned}$$

and variational derivative is

$$\delta \tilde{\mathcal{L}}(\mathbf{v}_t(\mathbf{x})) \Big|_{\mathbf{v}_t} = \frac{d}{d\varepsilon} \Big|_{\varepsilon=0} \tilde{\mathcal{L}}(\mathbf{v}_t + \varepsilon \boldsymbol{\tau}_t) = \int_0^1 \int_{\Omega} \langle \mathbf{v}_t(\mathbf{x}) - \nabla \varphi_t(\mathbf{x}), \boldsymbol{\tau}_t(\mathbf{x}) \rangle \rho_t(\mathbf{x}) d\mathbf{x} dt. \quad (37)$$

The stationarity condition requires  $\delta_{\mathbf{v}_t} \tilde{\mathcal{L}}(\mathbf{v}_t(\mathbf{x})) = 0$ . For arbitrary perturbation  $\boldsymbol{\tau}_t(\mathbf{x})$ , the variation  $\delta_{\mathbf{v}_t} \tilde{\mathcal{L}}(\mathbf{v}_t(\mathbf{x})) = 0$  if and only if

$$\mathbf{v}_t(\mathbf{x}) = \nabla \varphi_t(\mathbf{x}). \quad (38)$$

Therefore one can write the Lagrangian as

$$\begin{aligned} \mathcal{L}(\rho_t, \psi_0, \varphi_t, \eta) &= \int_{\Omega} \psi_0(\mathbf{x}) \cdot (\rho_0(\mathbf{x}) - \mu(\mathbf{x})) d\mathbf{x} + \int_{\Omega} \varphi_1(\mathbf{x}) \rho_1(\mathbf{x}) d\mathbf{x} - \int_{\Omega} \varphi_0(\mathbf{x}) \rho_0(\mathbf{x}) d\mathbf{x} \\ &\quad + \int_{\Omega} \eta(\mathbf{x}) \rho_1(\mathbf{x}) d\mathbf{x} - c \cdot \int_{\Omega} \max(0, \eta(\mathbf{x})) \nu(\mathbf{x}) d\mathbf{x} \\ &\quad - \int_0^1 \int_{\Omega} \left( \frac{\partial}{\partial t} \varphi_t(\mathbf{x}) + \frac{\|\nabla \varphi_t(\mathbf{x})\|^2}{2} \right) \rho_t(\mathbf{x}) d\mathbf{x} dt. \end{aligned}$$

Similarly, by computing  $\delta_{\psi_0} \mathcal{L}$  and  $\delta_{\rho_1} \mathcal{L}$  using stationary conditions, one obtains the condition,

$$\psi_0(\mathbf{x}) = \varphi_0(\mathbf{x}), \quad (39)$$

$$\eta(\mathbf{x}) = -\varphi_1(\mathbf{x}). \quad (40)$$

Therefore the Lagrangian is simplified to

$$\begin{aligned} \mathcal{L}(\rho_t, \varphi_t, \eta) &= - \int_{\Omega} \varphi_0(\mathbf{x}) \mu(\mathbf{x}) d\mathbf{x} - c \int_{\Omega} \max(0, -\varphi_1(\mathbf{x})) \nu(\mathbf{x}) d\mathbf{x} \\ &\quad - \int_0^1 \int_{\Omega} \left( \frac{\partial}{\partial t} \varphi_t(\mathbf{x}) + \frac{\|\nabla \varphi_t(\mathbf{x})\|^2}{2} \right) \rho_t(\mathbf{x}) d\mathbf{x} dt. \end{aligned} \quad (41)$$

The simplified problem (equation 9 in the main body) is

$$\sup_{\rho_t} \inf_{\varphi_t, \mathbf{x} \sim \mu} \mathbb{E} [\varphi_0(\mathbf{x})] + c \cdot \mathbb{E}_{\mathbf{x} \sim \nu} [\max(0, -\varphi_1(\mathbf{x}))] + \int_0^1 \mathbb{E}_{\mathbf{x}_t \sim \rho_t} \left[ \frac{\partial}{\partial t} \varphi_t(\mathbf{x}_t) + \frac{\|\nabla \varphi_t(\mathbf{x}_t)\|^2}{2} \right] dt.$$

## B THRESHOLDING FOR PU-LEARNING AND REJECTION SAMPLING

Our idea of rejection sampling and thresholding for PU-Learning is based on the fact that the dual form of range divergence is zero when the supremum in the dual is attained by the function  $\eta^*(\mathbf{x})$  with  $\tilde{\nu}(\mathbf{x}) = \rho_1^*(\mathbf{x})$  i.e.

$$\mathbb{E}_{\mathbf{x} \sim \tilde{\nu}} [\eta^*(\mathbf{x})] - c \mathbb{E}_{\mathbf{x} \sim \nu} [\text{ReLU}(\eta^*(\mathbf{x}))] = 0 \quad (42)$$



By defining  $\mathcal{A} = \text{supp}(\nu)$ ,  $\tilde{\mathcal{A}} = \text{supp}(\tilde{\nu})$ , and  $\bar{\mathcal{A}} = \mathcal{A} / \tilde{\mathcal{A}}$ , one can also see that  $\tilde{\mathcal{A}} \cap \bar{\mathcal{A}} = \emptyset$ , therefore one can write Equation 42 as

$$\int_{\tilde{\mathcal{A}}} \eta^*(x) \tilde{\nu}(x) dx - c \int_{\tilde{\mathcal{A}}} \text{ReLU}(\eta^*(x)) \nu(x) dx - c \int_{\bar{\mathcal{A}}} \text{ReLU}(\eta^*(x)) \nu(x) dx = 0. \quad (43)$$

One can further write

$$\eta^*(x) = \text{ReLU}(\eta^*(x)) - \text{ReLU}(-\eta^*(x)). \quad (44)$$

After substituting Equation 44 into Equation 43 one obtains

$$\overbrace{\int_{\tilde{\mathcal{A}}} \text{ReLU}(\eta^*(x)) (\tilde{\nu}(x) - c\nu(x)) dx}^{\text{LHS}} = \overbrace{\int_{\tilde{\mathcal{A}}} \text{ReLU}(-\eta^*(x)) \tilde{\nu}(x) dx + c \int_{\bar{\mathcal{A}}} \text{ReLU}(\eta^*(x)) \nu(x) dx}^{\text{RHS}} \quad (45)$$

The dual form Equation 43 is optimal with zero duality gap, if the primal form satisfies  $\forall x \in \mathcal{A}$ ,  $\iota_{[0,c]}(\frac{\tilde{\nu}}{\nu}(x)) = 0$ , which can also be restricted to  $\forall x \in \tilde{\mathcal{A}} \iota_{[0,c]}(\frac{\tilde{\nu}}{\nu}(x)) = 0$ . This is equivalent to  $\tilde{\nu}(x) \leq c\nu(x)$  almost-everywhere in  $\tilde{\mathcal{A}}$ . Therefore, one can say that  $0 \geq \text{LHS}$  and also

$$0 \geq \overbrace{\int_{\tilde{\mathcal{A}}} \text{ReLU}(-\eta^*(x)) \nu(x) dx + c \int_{\bar{\mathcal{A}}} \text{ReLU}(\eta^*(x)) \nu(x) dx}^{\text{RHS}} \quad (46)$$

We can now see that both integrands in Equation 46 are nonnegative and sum to a value less than or equal to zero, which is only possible if both are equal to zero. Therefore, one can write

$$0 = \int_{\tilde{\mathcal{A}}} \text{ReLU}(-\eta^*(x)) \nu(x) dx + c \int_{\bar{\mathcal{A}}} \text{ReLU}(\eta^*(x)) \nu(x) dx \quad (47)$$

Further, two non-negative integrals are evaluated on two mutually exclusive sets, therefore to have sum equal to zero value we can conclude that each integral is zero individually. Therefore, we can write

$$0 = \int_{\tilde{\mathcal{A}}} \text{ReLU}(-\eta^*(x)) \nu(x) dx = c \int_{\bar{\mathcal{A}}} \text{ReLU}(\eta^*(x)) \nu(x) dx \quad (48)$$

The Equation 48 is therefore equivalent to following element-wise test

$$\begin{aligned} \eta^*(x) &\geq 0, \text{ almost surely in } \tilde{\mathcal{A}} \\ \eta^*(x) &< 0, \text{ almost surely in } \bar{\mathcal{A}} \end{aligned} \quad (49)$$

Additionally, from the Equation 47, one can also conclude that

$$\overbrace{\int_{\tilde{\mathcal{A}}} \text{ReLU}(\eta^*(x)) (\tilde{\nu}(x) - c\nu(x)) dx}^{\text{LHS}} = 0, \quad (50)$$

which is a complementary slackness condition in the sense that  $\tilde{\nu}(x) < c\nu(x) \implies \eta^*(x) = 0$  almost every-where in  $\mathcal{A}$ . During the neural network training with finite data-points, potential function  $\eta$  is usually suboptimal and its sign cannot be relied, therefore instead of directly using the sign, one can sort values of potential at data points and select predetermined proportion (prior) of data-points. Therefore for training PU-learning models, we applied both sign and sorting based filtration of data. From the figures 1a and 1b, one can observed that for the optimal potential for static problem exactly follows equation 49, whereas in the dynamic case the sign of  $\varphi_1$  is inverted, which is due to the relation obtained in equation 40, which ensures that for the at optimal  $\varphi_1^*$  following relation holds

$$\begin{aligned} \varphi_1^*(x) &\leq 0, \text{ almost surely in } \tilde{\mathcal{A}} \\ \varphi_1^*(x) &> 0, \text{ almost surely in } \bar{\mathcal{A}}. \end{aligned} \quad (51)$$

The Figure 5 gives snapshots of the transition of  $\varphi_t(x)$  between  $t = 0$  and  $t = 1$  for the dynamic subset alignment results shown 1b.



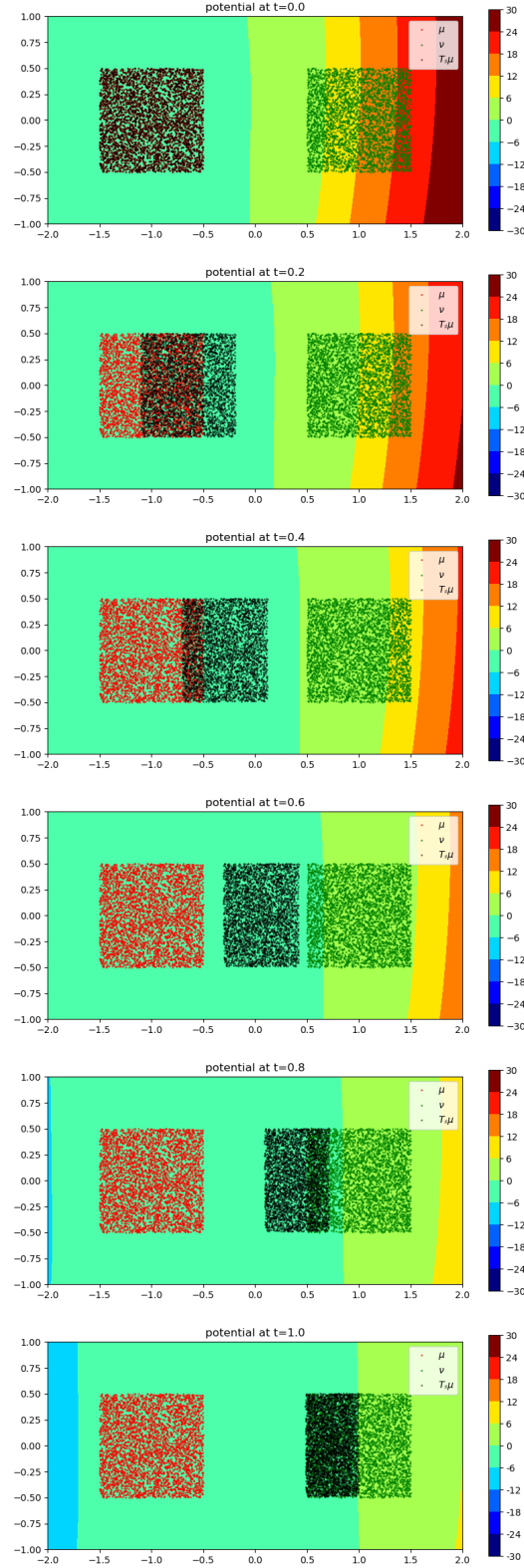


Figure 5:  $\varphi_t$  between  $t = 0$  and  $t = 1$  for subset alignment between 2D uniform distributions for which  $\varphi_1$  is also shown in Figure 1b, It can be seen that unlike  $\eta$  in static problem  $\varphi_t$  is function of time and varies with  $t$ .



## C SURVEY OF RECENT WORK ON NEURAL OPTIMAL TRANSPORT

In this section, we discuss the recent related work on computational optimal transport and its applications. More specifically, we consider the works which are related to neural estimation of optimal transport maps with occasional reference to theoretical developments.

### C.1 STATIC NEURAL OPTIMAL TRANSPORT

Seguy et al. (2018) employed stochastic gradient based approaches in one of the earliest works to estimate the optimal transport map using neural networks. Notably, the work by Seguy et al. (2018) differed from early work Genevay et al. (2016) in the sense that the later work employed stochastic gradient based methods to estimate the transport plan for large scale data, whereas earlier work Genevay et al. (2016) only minimized the optimal transport loss using stochastic-gradient based methods. This is also in contrast to the well-known Wasserstein-GAN Arjovsky et al. (2017); Gulrajani et al. (2017) that employs the Kantorovich-Rubinstein duality to minimize the Wasserstein-1 loss function for generative modeling, where neural networks are employed as parameterizations for both dual-potential and data generator, but do not provide transport plans. Finally, the Sinkhorn-GAN employs an approximation of the discrete Wasserstein distance between latent representations of data and that of samples from non-informative prior Genevay et al. (2018) for generative modeling. Now, we can see clear distinction between two different classes of approaches employing Wasserstein distances in generative modeling, the first class of works concerns with employing Wasserstein distance as a loss for generative modeling, without any explicit concern for obtaining the underlying transport plan across the distributions Arjovsky et al. (2017); Gulrajani et al. (2017). The second class seeks to learn a transport plan to realize the generative model.

Efforts to learn Monge maps were motivated by a theorem by Brenier (1991), which essentially states that, for continuous distributions with squared-Euclidean transportation cost, the optimal solution of the Monge problem is the gradient of a convex function (Figalli & Glaudo, 2023, (Theorem 2.5.10)). Therefore initially, gradient of input-convex neural networks (ICNN) Amos et al. (2017) we employed to estimate the transport plan for the Wasserstein-2 distance Makuva et al. (2020); Korotin et al. (2021a;b). This approach has also been employed to supervised conditional neural Monge maps (Bunne et al., 2022a) and unbalanced optimal transport (Lübeck et al., 2022). The study by Amos et al. (2023) focuses on the development of an efficient neural optimal solution that could be implemented quickly in more practical scenarios. This approach to solve Wasserstein-2 distances employing convex potentials involves computationally challenging evaluation of the Fenchel conjugate of a ICNN parameterized convex function. More recent work in this direction focuses on improved optimization strategies and better ICNN architectures to bypass problems related to Fenchel conjugate evaluations and ICNN training Amos (2023); Vesseron & Cuturi (2024). Recent work also focuses on some batch-based schemes have also been devised to improve the regularity of learned neural Monge maps Uscidda & Cuturi (2023); Eyring et al. (2024).

Another recent direction of work is based on the idea that ICNNs can be overly restrictive, therefore more general neural network architectures should be employed to directly parameterize the transport maps Rout et al. (2022); Korotin et al. (2023b). The work by Fan et al. (2022a; 2023) focuses on employing neural networks to approximate the solution for Monge’s transport problem also draws inspiration from the recent developments in neural-network-based parametric realizations for approximating Kantorovich plans. Recently neural optimal transport has also been extended to unbalanced transportation setting (Yang & Uhler, 2019; Choi et al., 2023). Another work directly related to static subset selection problem is (Gazdieva et al., 2023).

Unless there is a corresponding Monge mapping (Choi et al., 2024a; Mokrov et al., 2024; Geuter et al., 2025), optimal transport requires a stochastic transport plans. A recent body of work (Korotin et al., 2023b;a; Asadulaev et al., 2024) deals with learning transportation plans using a weaker formulation of optimal transport (Gozlan et al., 2017; Backhoff-Veraguas et al., 2019) along with noise outsourcing techniques, which is also extended to more general costs. Apart from the applications in image translation (Korotin et al., 2023b), neural optimal transport has been applied for bio-medical image registration (Kim et al., 2024) and to study single cell perturbations (Bunne et al., 2023). Neural optimal transport have also been employed for metric learning (Howard et al., 2024; Scarvelis & Solomon, 2023).



## C.2 DYNAMIC NEURAL OPTIMAL TRANSPORT

The potential applications of dynamic optimal transport in the cellular trajectory inference (Tong et al., 2020) and its connections with flow based models for generative modeling (Huang et al., 2021; Huguet et al., 2022) has been instrumental in the recent research developments in this direction. Jordan-Kinderlehrer-Otto flow (JKO) is time discretization scheme to solve Wasserstein gradient flows for different energy functionals (Jordan et al., 1998; Santambrogio, 2017). Therefore, a lot of effort done in that regard is focused on neural network parameterized schemes to solve JKO-flow problem for both cellular trajectory inference and generative modeling (Ma et al., 2021; Fan et al., 2022b; Lambert et al., 2022; Bunne et al., 2022b; Xu et al., 2023; Choi et al., 2023; 2024c; Altekruiger et al., 2023; Mokrov et al., 2021; Alvarez-Melis et al., 2022). JKO-scheme has also been studied for the applications related to molecular discovery (Alvarez-Melis et al., 2022). A recent study deals with convergence properties of JKO-based generative models (Cheng et al., 2024).

Recent developments in flow-matching models based on flow matching (Lipman et al., 2023; Albergo & Vanden-Eijnden, 2023; Liu et al., 2023) for generative modeling lead to even more interest in the development of algorithms to solve dynamic optimal transportation problems. Action-Matching based framework lead to the development of a more general framework to solve both trajectory inference and generative modeling problems (Neklyudov et al., 2023) for the cases where one could also sample from the trajectory between two terminal marginals. Rectified flow-matching (Liu et al., 2023; 2024b) uses the neural-optimal transport in additional rectification step to improve the linearity of flows, so that after training the model, images could be generated efficiently with only a single-step integration along straight lines paths. For generative modeling, in contrast to target-conditional flow matching (Lipman et al., 2023), where during training, flows are conditioned on target samples, discrete optimal transport conditioned flow-matching employs the mini-batch optimal transport to create the conditionals (Pooladian et al., 2023; 2024; Tong et al., 2024b). Another recent work (Kornilov et al., 2024) attempts to alleviate the error accumulation problems associated with mini-batch optimal transport by learning straight paths between source and target distributions in single step. Flow-matching (Albergo & Vanden-Eijnden, 2023; Albergo et al., 2023), diffusion models (Sohl-Dickstein et al., 2015; Song & Ermon, 2020; Song et al., 2021), and Schrödinger bridges (Wang et al., 2021; Liu et al., 2022; 2024a; Shi et al., 2023; Gushchin et al., 2023b;a), and (Somnath et al., 2023) are deeply interconnected under the framework of generalized bridge matching (Tong et al., 2023; Albergo et al., 2023; Tong et al., 2024a; Shi et al., 2024). Recently, there has also been attempts to understand diffusion models as approaches to minimize the dynamic Wasserstein distances (Kwon et al., 2022; Khrulkov et al., 2023). Another recent work extends the flow matching to the flows on Riemannian manifolds (Chen & Lipman, 2024; Atanackovic et al., 2025). Recent works generalize flow-matching from different perspectives, Chen & Lipman (2024) generalize the flow-matching to the flows on Riemannian manifolds, Atanackovic et al. (2025) attempt to extend the flow-models to return meaningful flows for the data beyond training distributions, and Haviv et al. (2025) generalize the flow matching to the cases where data can be treated as distributions of distributions.

Additionally, there has been recent dynamic extension to the conditional neural optimal transport (Hosseini et al., 2023; Kerrigan et al., 2024). There has also been efforts to study neural network based scalable approaches to solve high-dimensional partial differential equations (Wan et al., 2023).

## D IMPLEMENTATION DETAILS

### D.1 EMNIST CLASSIFIER

We merged the whole alphabet into one class and each number is treated as a separate class (digits between 0 and 9 are given same label as their value and any letter is labeled 10). In order to circumvent the effects of data imbalancedness on classifier training, we employed the class-reweighted softmax loss function. For  $k$ -class classification, consider the vector  $\mathbf{z} \in \mathbb{R}^k$  containing the counts for class in the training data, we define the reweighting vector  $\boldsymbol{\omega} \in \mathbb{R}^k$  with

$$\omega_i = \left( \sum_{j=1}^N \frac{z_i}{z_j} \right)^{-1}, \quad \forall i \in [k]. \quad (52)$$



For one hot encoded label vector  $\mathbf{y}$  and softmax activation output at neural network output  $\hat{\mathbf{y}}$ , the reweighted loss (risk) is given by

$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = \mathbf{1}_k^\top (\boldsymbol{\omega} \odot \mathbf{y} \odot \hat{\mathbf{y}}) \quad (53)$$

The classifier for EMNIST is trained with the same train/validation split as provided in EMNIST dataset (Cohen et al., 2017). We trained the classifier with ResNet-18 (He et al., 2016) architecture and class-reweighted softmax loss function in equation 53. Adam optimizer (Kingma & Ba, 2014) along with warmup-cosine learning rate scheduler (Loshchilov & Hutter, 2017) is used to train the classifier with peak learning rate of  $1 \times 10^{-3}$  with 500 warm-up steps. Total decay steps for cosine scheduler are set to 20,000 with end-value of learning rate set to be equal to  $1 \times 10^{-5}$ . The classifier training is stopped after 20,000 training steps, when classifier achieves more than 90% overall validation accuracy and 99% accuracy on digits. Confusion matrix of classifier are given in Appendix E.

## D.2 MNIST-EMNIST TRANSLATION MODELS

For the static domain translation, the transport network  $T$  is a U-Net Ronneberger et al. (2015) with base-factor of 48 and the critic network  $\eta$  is ResNet-51 He et al. (2016). In order to train both transport and critic networks, Adam optimizer Kingma & Ba (2014) is used with initial learning rate of  $1 \times 10^{-4}$ , which is scheduled to be halved after 10,000 + 5000c, 20,000 + 5000c, 30,000 + 5000c, 40,000 + 5000c and 70,000 + 5000c training steps. Algorithm 1 is used for training with 50,000 *learning iterations* with 10 *T update steps* for each  $\eta$  *update step*, our training settings for static case are very similar to those of Gazdieva et al. (2023). For dynamic subset selection, following the settings from Neklyudov et al. (2023), the vector field  $\varphi_t$  is parametrized using a U-Net with time embeddings from DDPM (Song & Ermon, 2020). Similar to action matching (Neklyudov et al., 2023),  $\varphi_t$  is parametrized to return scalar by  $\varphi_t(\mathbf{x}) = \langle \text{U-Net}(\mathbf{x}), \mathbf{x} \rangle$ . Likewise,  $Q_t$ , which parametrizes  $\rho_t$ , is also a U-Net with time embeddings. We used AdamW optimizer with learning rate scheduling for 50,000 iterations. The optimizer parameters are  $\beta = (0, 0.999)$ , weight decay = 0.1 and drop out = 0.1. Additionally, we also employed exponential moving averages (EMA) in the training with the ema-rate 0.999. These settings are very similar to rectified flow matching and action matching (Liu et al., 2023; Neklyudov et al., 2023). Learning rate linearly increases from 0 to maximum value during first 5,000 iterations and then stays constant at maximum value with maximum learning rates of  $2 \times 10^{-4}$  and  $1 \times 10^{-4}$  for  $\varphi_t$  and  $Q_t$ , respectively. Additionally, we clipped gradients to lie within [-1, 1]. Algorithm 2 is employed with 50,000 training iterations and 2  $\varphi_t$  for each  $\rho_t$  update.

## D.3 MODELS FOR PU-LEARNING USING SUBSET ALIGNMENT

For PU learning with both static and dynamic subset alignment based approaches respectively, model architectures are given in code listings D.3 and D.3, respectively. For all models num\_hid is set to be 1024, for Smodel and etamodel, the parameter num\_out is by definition 1, whereas for Qmodel and Tmodel, outputs are set to be equal to data dimension. For both static and dynamic models, we used Adam optimizer Kingma & Ba (2014), with default settings, and learning rates  $1 \times 10^{-4}$  and  $2 \times 10^{-5}$  respectively. Additionally, we used EMA with ema-rate of 0.999 to evaluate models on both the test dataset and the validation datasets. We trained the model for the total of 20,000 *learning iterations*, with 10 *T update steps* for single  $\eta$  *update step* using the Algorithm 1. Similarly, Algorithm 2 is employed to train neural networks for dynamic subset alignment. *learning iterations*, with 2  $\varphi_t$  *update steps* for single  $\rho_t$  *update step*. The dynamic models contain time embeddings with trainable parameters. We employed the Adam algorithm for gradient based updates of neural network parameters. For all the tests for PU learning we fix  $c = \frac{1}{\pi_+}$ . For each data set same batch sizes are used to train both static and dynamic models and table 5, gives the values.



Dataset	n	dim	$\pi$	batch size
Abalone	4177	8	0.16	20
Banknote	1372	4	0.44	10
Breast-w	699	9	0.34	10
Diabetes	768	8	0.35	6
Haberman	306	3	0.26	6
Heart	270	13	0.44	6
Ionosphere	351	34	0.64	6
Isolet	7797	617	0.04	4
Jml	10885	21	0.19	20
Kc1	2109	21	0.15	20
Madelon	2600	500	0.5	20
Musk	6598	166	0.15	20
Segment	2310	19	0.14	20
Semeion	1593	256	0.1	4
Sonar	208	60	0.53	4
Spambase	4601	57	0.39	20
Vehicle	846	18	0.26	6
Waveform	5000	40	0.34	20
Wdbc	569	30	0.37	6
Yeast	1484	8	0.31	10

Table 5: UCI datasetets for PU Learning, along with total number of data points (n), dimension (dim), positive prior ( $\pi$ ) and batch sizes employed in training the correspondng models.

```

1427 1 import jax
1428 2 from jax import numpy as jnp
1429 3 from flax import linen as nn
1430 4 import math
1431 5 '''
1432 6 etamodel: neural network parameterization for eta function
1433 7 Tmodel: neural network parameterization for T function
1434 8 '''
1435 9 class etamodel(nn.Module):
1436 10     num_hid : int
1437 11     num_out : int
1438 12     @nn.compact
1439 13     def __call__(self, x):
1440 14         h = nn.Dense(self.num_hid)(x)
1441 15         h = nn.swish(h)
1442 16         h = nn.Dense(self.num_hid)(h)
1443 17         h = nn.swish(h)
1444 18         h = nn.Dense(self.num_hid)(h)
1445 19         h = nn.swish(h)
1446 20         h = nn.Dense(self.num_out)(h)
1447 21         return h
1448 22
1449 23 class Tmodel(nn.Module):
1450 24     num_hid : int
1451 25     num_out : int
1452 26     @nn.compact
1453 27     def __call__(self, x):
1454 28         def transport_net(x):
1455 29             MLP_out = nn.Sequential([
1456 30                 nn.Dense(self.num_hid),
1457 31                 nn.swish,
1458 32                 nn.Dense(self.num_hid),
1459 33                 nn.swish,
1460 34                 nn.Dense(self.num_hid),
1461 35                 nn.swish,
1462 36                 nn.Dense(self.num_hid),
1463 37                 nn.swish,
1464 38                 nn.Dense(self.num_out),]) (x)

```



```
1458 39         ResConnect = nn.Dense(self.num_out)(x)
1459 40         return MLP_out + ResConnect
1460 41     output = transport_net(x)
1461 42     return output
```

Listing 1: Model architectures for PU-Learning with static subset alignment

1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511



```

1512 1 import jax
1513 2 from jax import numpy as jnp
1514 3 from flax import linen as nn
1515 4 import math
1516 5
1517 6 class Smodel(nn.Module):
1518 7     num_hid : int
1519 8     num_out : int
1520 9
1521 10 @nn.compact
1522 11 def __call__(self, t, x):
1523 12     if jnp.ndim(t) == 0:
1524 13         t = jnp.broadcast_to(t, x.shape[0:-1]+(1,))
1525 14     h = jnp.concatenate([t,x], axis=-1)
1526 15     h = nn.Dense(self.num_hid)(h)
1527 16     h = nn.swish(h)
1528 17     h = nn.Dense(self.num_hid)(h)
1529 18     h = nn.swish(h)
1530 19     h = nn.Dense(self.num_hid)(h)
1531 20     h = nn.swish(h)
1532 21     h = nn.Dense(self.num_hid)(h)
1533 22     h = nn.swish(h)
1534 23     h = nn.Dense(self.num_out)(h)
1535 24     return h
1536 25
1537 26
1538 27 class Qmodel(nn.Module):
1539 28     num_hid : int
1540 29     num_out : int
1541 30
1542 31 @nn.compact
1543 32 def __call__(self, t, x_0, x_1):
1544 33
1545 34     h = jnp.concatenate([t, x_0, x_1, t<0.5], axis=-1)
1546 35     h = nn.Dense(self.num_hid)(h)
1547 36     h = nn.swish(h)
1548 37     h = nn.Dense(self.num_hid)(h)
1549 38     h = nn.swish(h)
1550 39     h = nn.Dense(self.num_hid)(h)
1551 40     h = nn.swish(h)
1552 41     h = nn.Dense(self.num_hid)(h)
1553 42     h = nn.swish(h)
1554 43     h = nn.Dense(self.num_out)(h)
1555 44
1556 45     x_t = (1-t)*x_0 + t*(x_1) + t*(1-t)*h
1557 46
1558 47     return x_t
1559 48

```

Listing 2: Model architectures for PU learning with dynamic subset alignment

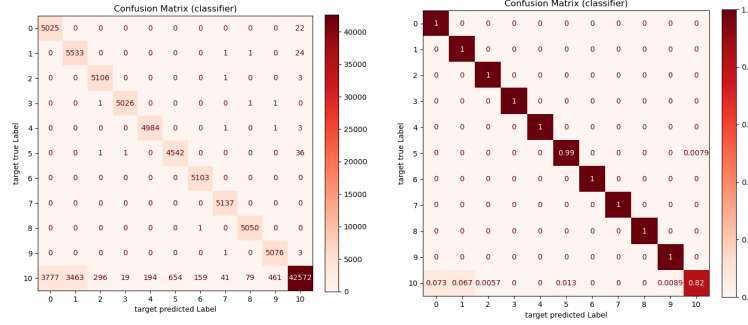
#### D.4 IMAGE-TO-IMAGE TRANSLATION ON FFHQ

In our experiments for static subset alignment, we used a three layered MLP architecture with swish activation functions in hidden layers to parameterize both the transportation map  $T$  and the potential  $\eta$ . For the network parameterizing  $T$ , an additional skip connection connecting input and output is also used, which also contains a linear mapping, without any non-linear activation. Dimension of hidden layers are set to 1,024 for both Networks. Output dimension of the transport network is same as its input dimension (512), whereas potential network returns a scalar output. The Adam optimization algorithm is used to train both networks with a fixed learning rate of  $1 \times 10^{-5}$ . We employ EMA with ema-rate 0.999 in the training process. Algorithm 1 is used in the training with 50,000 learning iterations with 5  $T$  updates for each  $\eta$  update.

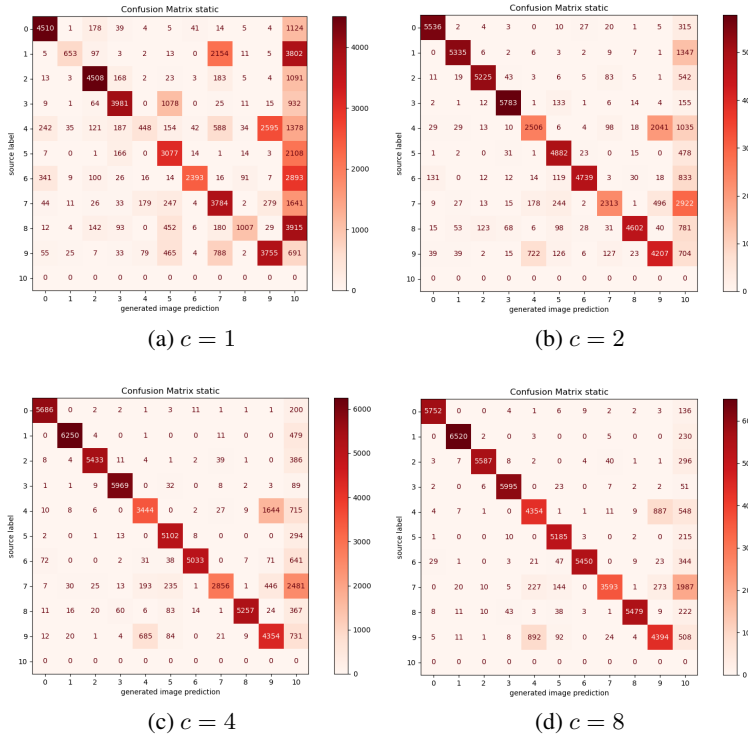


In order to train the dynamic models, the model architectures employed are also three layered MLPs but with time embeddings. The neural network parameterizing  $\varphi_t$  is a three layers MLP with 64 dimensional time embeddings, 1,024 dimensional hidden layers, and a scalar output. The neural Network parameterizing  $\rho_t$  contains two branches for static and dynamic components respectively. The dynamic part of network parameterizing  $\rho_t$  also contains 64 dimensional time embeddings. We also use EMA with ema-rate 0.999 to train both networks, and a fixed learning rate of  $1 \times 10^{-5}$ . Dynamic models are trained using algorithm 2 for 50,000 *learning iterations* with 1  $\varphi_t$  *update* for 5  $\rho_t$  *updates*.

## E CONFUSION MATRICES FOR MNIST $\rightarrow$ EMNIST DOMAIN TRANSLATION



(a) Unnormalized confusion matrix (b) Normalized confusion matrix  
Figure 6: Confusion matrices for EMNIST classifier discussed in section 4.1



(a)  $c=1$  (b)  $c=2$  (c)  $c=4$  (d)  $c=8$   
Figure 7: Confusion matrices for MNIST  $\rightarrow$  EMNIST domain translation using static subset selection. Accuracy is computed by computing ratio between trace and some of all entries of confusion matrices.



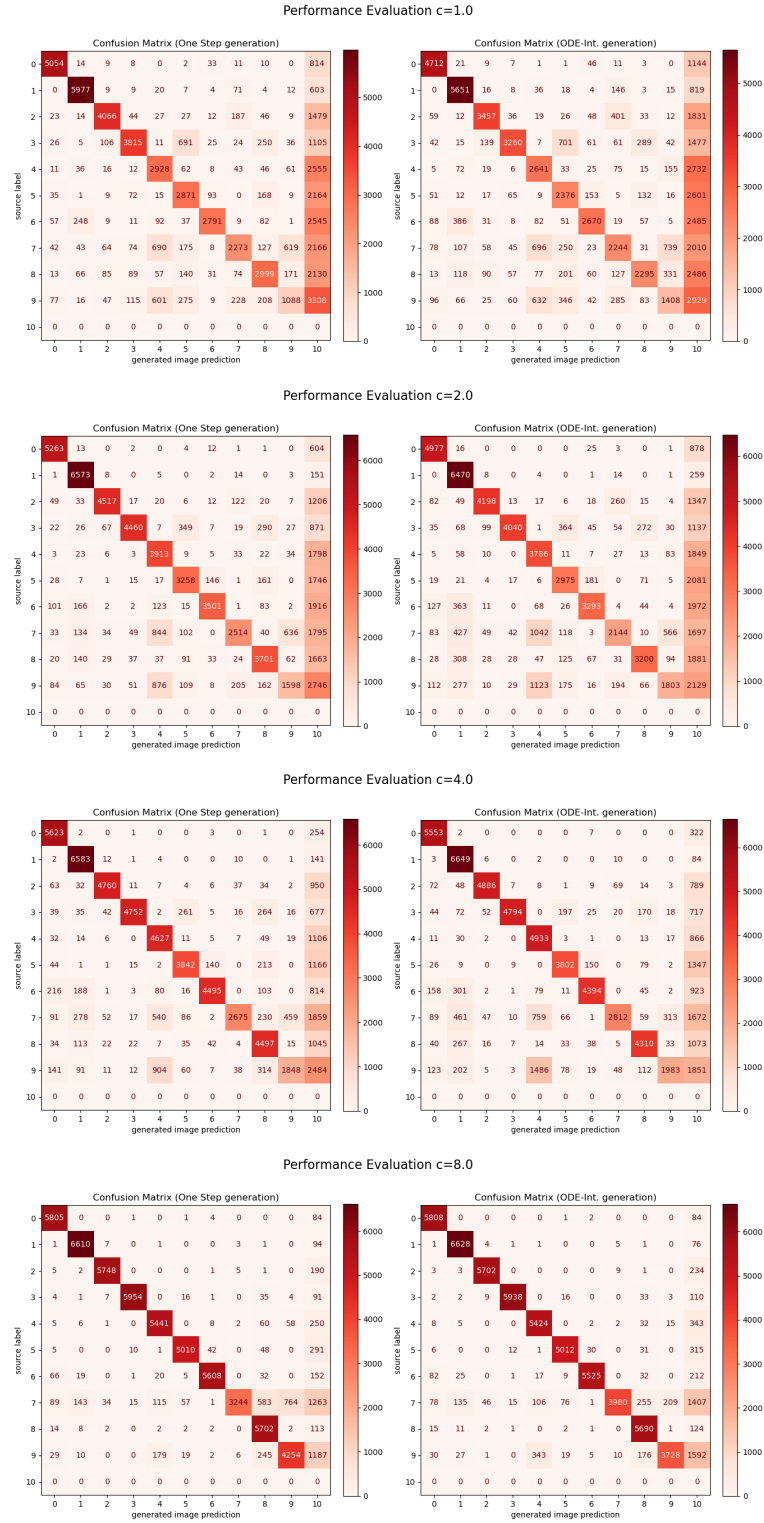


Figure 8: Confusion matrices for MNIST→EMNIST domain translation using dynamic subset selection.



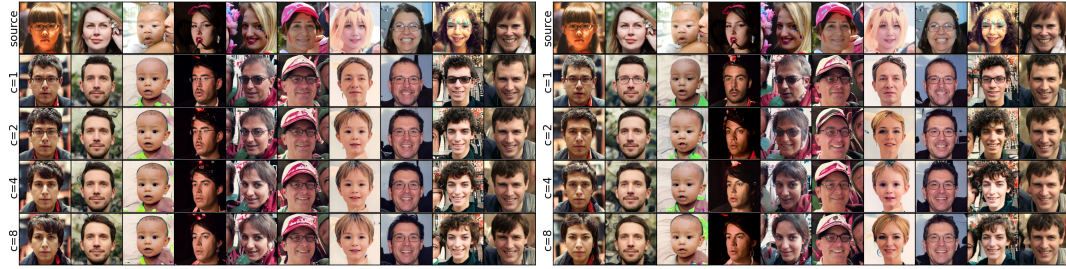
## F RESULTS FROM FFHQ



(a)

(b)

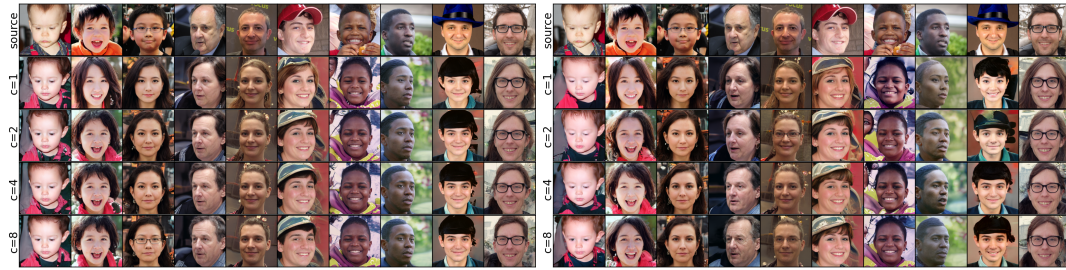
Figure 9: FFHQ young→old translation using (a) static and (b) dynamic subset selection. Dynamic subset selection is evaluated using Euler integration with 100 steps.



(a)

(b)

Figure 10: FFHQ woman→man translation using (a) static and (b) dynamic subset selection. Dynamic subset selection is evaluated using Euler integration with 100 steps.



(a)

(b)

Figure 11: FFHQ man→woman translation using (a) static and (b) dynamic subset selection. Dynamic subset selection is evaluated using Euler integration with 100 steps.