

# Distribution-Conditioned Transport

Anonymous Authors<sup>1</sup>

## Abstract

Learning a transport model that maps a source distribution to a target distribution is a canonical problem in machine learning, but scientific applications increasingly require models that can generalize to source and target distributions unseen during training. We introduce distribution-conditioned transport (DCT), a framework that conditions transport maps on learned embeddings of source and target distributions, enabling generalization to unseen distribution pairs. DCT also enables semi-supervised learning for distributional forecasting problems: because it learns from arbitrary distribution pairs, it can leverage distributions observed at only one timepoint to improve transport prediction. DCT is agnostic to the underlying transport mechanism, supporting models ranging from flow matching to distributional divergence-based models (e.g. Wasserstein, MMD). We demonstrate the practical performance benefits of DCT on synthetic benchmarks and four applications in biology: batch effect transfer in single-cell genomics, perturbation prediction from mass cytometry data, learning clonal transcriptional dynamics in hematopoiesis, and modeling T-cell receptor sequence evolution.

## 1. Introduction

Learning transport maps between probability distributions is a central challenge across machine learning and the sciences. A wide variety of approaches have addressed the one-to-one transport problem: learning a map that pushes a source density  $P_0$  to a target density  $P_1$  (Goodfellow et al., 2014; Rezende & Mohamed, 2015; Genevay et al., 2018; Liu et al., 2022; Lipman et al., 2023; Albergo et al., 2023b). However, a new class of problems is emerging where the goal is not only to model the evolution between a single pair

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

of distributions, but rather to model dynamics in a way that can generalize across source and target distributions.

A concrete motivating example of this shift can be seen in the study of cellular dynamics. Experimental techniques such as clonal lineage-traced genomics have enabled the generation of datasets containing snapshots not just of a single cell population, but of thousands of distinct populations (clones) evolving in parallel (Biddy et al., 2018; Weinreb et al., 2020; Wagner & Klein, 2020). These datasets can be sparse, in the sense that not all time-marginals are observed for all populations. For example, we may observe some populations at both an initial and final times  $t_0$  and  $t_1$ , while others are “orphans” observed only at either  $t_0$  or  $t_1$ .

Two lines of literature have begun to develop tools for these settings. Multimarginal stochastic interpolants offer an approach to learn dynamics between any pair of a fixed set of  $K$  distributions, solving a  $K$ -to- $K$  transport problem (Albergo et al., 2023a), but cannot condition on a continuous space of distributions nor generalize to distributions unseen during training. In another direction, Fishman et al. (2025) develop approaches to learn “autoencoders” on the space of distributions by simultaneously learning to embed distributions and sample from the distribution conditional on that embedding. Meta flow matching (Atanackovic et al., 2024) is an early approach that couples a population encoder with a flow-matching transport map for generalization across source populations. While effective in settings where many source-target pairs are available (e.g., if we observed many pairs of clonal lineages at both  $t_0$  and  $t_1$ ), the MFM framework cannot ingest unpaired marginal distributions (e.g., cell populations observed at a single timepoint).

In this work, we unify these perspectives. We first generalize and formalize the approach in Atanackovic et al. (2024), demonstrating how distribution encoders developed in Fishman et al. (2025) can be coupled with a broad class of transport models for source-conditioned transport. This immediately enables us to handle any-to-any transport by conditioning on both source and target distribution embeddings. This formulation enables us to generalize across distributions, predicting transport maps between population pairs unseen during training, as well as allowing us to make use of unstructured, partial observations (such as orphan marginals from incomplete time-series). We demonstrate

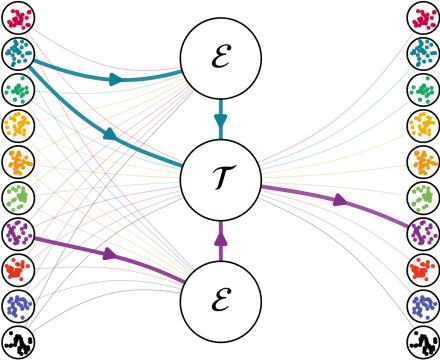


Figure 1. The distribution-conditioned transport framework. A source distribution (teal) is pushed to a target distribution (purple) by a transport model  $\mathcal{T}$  which is conditioned on distribution embeddings learned by an encoder  $\mathcal{E}$ . The learned transport map is universal in the sense that any distribution can in principle be pushed to any distribution by conditioning on the corresponding source and target embeddings.

the effectiveness of our approach first using a set of synthetic benchmark datasets, then on real-world applications in biology. The codebase is available [here](#).

## 2. Problem setting

Modern datasets in the sciences increasingly have multiscale structure. An example of this can be seen in single-cell RNA-sequencing (scRNA-seq) data. Rather than making measurements of cells from a single experimental condition, modern datasets increasingly consist of cells from many distinct donors (Yazar et al., 2022), timepoints (Qiu et al., 2024), perturbations (Zapatero et al., 2023), or clones (Weinreb et al., 2020), each of which can be thought of as inducing its own distribution over cell states.

This multiscale structure can be formalized as a hierarchical model. Suppose we have sets of samples (e.g., cell gene expression profiles) from  $n$  different conditions (e.g., donors). For each condition  $i$ , we observe a set of  $m$  samples  $S_i = \{x_{ij}\}_{j=1}^m$ , where each sample is a vector  $x_{ij} \in \mathcal{X}$ . Each condition induces a distribution  $P_i$ , and we can think of these condition-specific distributions as being drawn from a shared metadistribution  $Q$  over the space of probability measures  $\mathcal{P}(\mathcal{X})$ :  $x_{ij} \sim P_i$  and  $P_i \sim Q$ . The hierarchical dataset consists of these  $n$  sample sets:  $\mathcal{D} = \{S_i\}_{i=1}^n$ .

Many practical challenges in analyzing these hierarchical datasets can be framed as transport problems. For example, in scRNA-seq data, one often is interested in how a distribution would appear under different technical effects (batch integration), how distributions change in response to unseen treatment conditions (perturbation prediction), or forecasting how distributions evolve (dynamic inference).

While a wide variety of tools exist to model transport between a single pair of distributions, these hierarchical datasets and distributional tasks require transport models which can generalize across a broad range of distributions. In this work, we introduce *distribution-conditioned transport*, a framework for solving transport problems on the space of distributions. Concretely, this framework will allow us to solve three classes of practically relevant problems:

**Supervised (one-to-one) transport** between prescribed pairs of source and target distributions. For example, in clonal lineage-traced scRNA-seq data we observe many clonal populations at multiple timepoints. Each clone has a pair of population-level observations  $(P_i^{t=1}, P_i^{t=2})$ , and the task of interest is to forecast  $P_i^{t=2}$  given  $P_i^{t=1}$ .

**Unsupervised (any-to-any) transport** between any pair of distributions. In some cases, we may want to transport between any pair of distributions at test time, with no particular source-target coupling. For example, in scRNA-seq batch integration, we wish to transport cells between any pair of experimental batches, enabling any-to-any batch correction even for experimental batches unseen during training.

**Semi-supervised transport** with access to orphan marginals. For example, real lineage tracing data often contains incomplete time-series with many populations observed at only a single timepoint (Weinreb et al., 2020). Our goal is still to forecast  $P_i^{t=2}$  for any  $P_i^{t=1}$  while making use of these orphan marginals.

In the following section, we will show how distribution-conditioned transport provides a framework for addressing each of these problem settings. First, we will show that supervised transport problems can be solved by source-conditioned transport models, formalizing and generalizing the approach of Meta Flow Matching (Atanackovic et al., 2024). Then, we will introduce source-target-conditioned transport models for the unsupervised transport setting, which subsumes the setting addressed by multimarginal stochastic interpolants (Albergo et al., 2023a). Finally, we will show how these two perspectives can be combined to solve the semi-supervised transport problem when one has access to sparse observations of population pairs.

## 3. Methods

### 3.1. Distribution encoders

We follow Fishman et al. (2025), defining a *distribution encoder*  $\mathcal{E} : S_i \mapsto z_i \in \mathbb{R}^d$ , which produces a fixed-dimensional embedding  $z_i$  intended to summarize the entire distribution  $P_i$ , rather than any particular cell. The key aim of distribution encoders is that  $z_i$  reflects only the underlying distributional signal and not sampling noise in  $S_i$ . To enforce this the encoder must be *distributionally invariant*:

(i) *permutation invariant*, so that reordering samples does not change  $\mathcal{E}(S_i)$ , and (ii) *proportionally invariant*, so that uniformly duplicating samples does not change  $\mathcal{E}(S_i)$ .

When these hold, the encoder depends only on the empirical measure

$$\widehat{P}_i = \frac{1}{m_i} \sum_{j=1}^{m_i} \delta_{x_{ij}}.$$

In particular, there exists a measurable functional  $\phi$  such that  $\mathcal{E}(S_i) = \phi(\widehat{P}_i)$ . We refer to  $z_i = \phi(\widehat{P}_{i,m})$  as a *distribution embedding*. A key property of such encoders is that they admit a central limit theorem (CLT). Let

$$z_i^* = \lim_{m_i \rightarrow \infty} \phi(\widehat{P}_i) = \phi(P_i)$$

denote the population-level embedding obtained in the infinite-sample limit. Under the invariances above and Hadamard differentiability of  $\phi$ , we have

$$\sqrt{m_i} (\mathcal{E}(S_i) - z_i^*) \xrightarrow{d} \mathcal{N}(0, \Sigma_i) \quad (1)$$

for a covariance  $\Sigma_i$  depending on  $P_i$  and the encoder  $\mathcal{E}$ .

For our purposes, this CLT is the most important feature of distribution encoders because, as proved in (Fishman et al., 2025), it implies that for any sufficiently smooth downstream objective we can train using moderate-sized subsamples and recover population-level behavior up to  $O(m^{-1/2})$  error. We formalize this plug-in limit theory for DCT in Prop. 3.1 and App. A.

### 3.2. Supervised (one-to-one) transport

A key setting for distributional transport is when we observe “pairs” of distributions. For example, in perturbation prediction, we observe cell populations of various types (e.g., cell lines, donors) in their untreated and perturbed states. Here, the goal is to transport unperturbed cell populations into their perturbed condition. This transport setting is analogous to supervised learning: we observe many source sets of cells and their corresponding target distribution, and we want to learn a one-to-one map from the source to the target. We show that these supervised problems can be addressed by conditioning transport maps on embeddings of the source distribution. We refer to these as source-conditioned transport models.

A number of recent works have developed particular instances source-conditioned models such as (Atanackovic et al., 2024; Klein et al., 2025) in the context of flow matching or (He et al., 2025) using a direct MMD objective. Our notion of source-conditioning generalizes these under a common framework and clarifies the underlying assumptions under which they enable distributional transport.

Formally, we will assume throughout this section that we have a distribution of pairs sampled from some meta distribution (e.g., control cell populations and their corresponding

perturbed populations)

$$(P_1^{\text{src}}, P_1^{\text{tgt}}), \dots (P_n^{\text{src}}, P_n^{\text{tgt}}) \sim Q_{\text{joint}},$$

and we observe empirical sets

$$S_i^{\text{src}} = \{x_{ij}^{\text{src}}\}_{j=1}^{m_i^{\text{src}}} \sim P_i^{\text{src}} \text{ and } S_i^{\text{tgt}} = \{x_{ij}^{\text{tgt}}\}_{j=1}^{m_i^{\text{tgt}}} \sim P_i^{\text{tgt}}.$$

We encode the *source only*, so  $z_i^{\text{src}} = \mathcal{E}(S_i^{\text{src}})$ , where  $\mathcal{E}$  is a distribution encoder as laid out in Sec. 3.1. We can then learn a source-conditioned transport map acting on individual samples,

$$\mathcal{T} : \mathcal{X} \times \mathbb{R}^d \rightarrow \mathcal{X}, \quad x_i^{\text{src}} \mapsto \hat{x}_i^{\text{tgt}} = \mathcal{T}(x_i^{\text{src}} | z_{\text{src}}),$$

so that transported samples asymptotically follow the target distribution:

$$\mathcal{T}(S_i^{\text{src}} | \mathcal{E}(S_i^{\text{src}})) \xrightarrow[m_i^{\text{src}} \rightarrow \infty]{d} P_i^{\text{tgt}} \quad (2)$$

Here  $\mathcal{T}$  can be any conditional distributional transport model, and under a supervised pairing policy, the correct destination  $P_i^{\text{tgt}}$  is implied once the source is specified.

We jointly fit  $\mathcal{E}$  and  $\mathcal{T}$  using the native loss of the chosen transport mechanism, written abstractly as

$$\mathcal{L}_{\text{sc}} = \ell(S_i^{\text{tgt}}, \mathcal{T}(S_i^{\text{src}} | \mathcal{E}(S_i^{\text{src}}))), \quad (3)$$

where  $\ell$  may be a flow-matching objective (Atanackovic et al., 2024), a distributional divergence (e.g. Sinkhorn/MMD) (He et al., 2025), or any other generative transport model (see Sec 4.2). Because plug-in losses are mean consistent and inherit CLTs (Prop. 3.1 and App. A; see also Cor. A.6), we train with minibatches from  $S_{\text{src}}$  and  $S_{\text{tgt}}$  without biasing the population-level objective in the large- $m$  limit (Alg. 1). After training,  $\mathcal{T}(\cdot | \mathcal{E}(\cdot))$  can be applied to unseen source samples to generate counterfactual realizations under their implied targets.

Our notion of a *transport map* (see Sec. 4.2) is intentionally broad and includes models that, in some regimes, behave more like conditional generation than a coupling-selecting pushforward. In particular, conditional generative models with access to additional noise can achieve Eqn. 2 and match the target distributions while ignoring the source sample  $x_i^{\text{src}}$  they are supposed to condition on. We prove this in App. A.2 and provide a practical alignment diagnostic to detect this failure mode (Tab. 5). Importantly, this is a property of the underlying mechanism/objective rather than the DCT conditioning framework: DCT enables models to condition on distributions, but does not impose a within-distribution coupling structure beyond what the transport map would learn between a single pair of distributions.

### 3.3. Unsupervised (any-to-any) transport

In some settings, we wish to transport not between specific pairs of source-target distributions, but rather any pair of distributions. For example, in batch integration for scRNA-seq data, we are interested in the counterfactual of what one

165 dataset would look like under a different experimental batch  
 166 condition – and this counterfactual is of interest for any pair of  
 167 experimental batches.

168 In other settings, we are interested in particular source-  
 169 target pairings, but we do not have access to many true  
 170 pairs. Returning to our running lineage tracing case, in  
 171 many datasets, we may only observe the same clone in both  
 172 timepoints a small fraction of the time. In Weinreb et al.  
 173 (2020), only about 10% of clones are observed at multiple  
 174 timepoints. Here, we would ideally like a model to also  
 175 make use of the abundance of unpaired data.

176 To address both of these settings, we develop an “unsupervised” analogue to the “supervised” transport setting  
 177 considered above. In this unsupervised setting, we simply  
 178 have a set of unstructured distributions:  $P_1, \dots, P_n \sim Q$  and  
 179 associated samples:  $S_i = \{x_{ij}\}_{j=1}^{m_i} \sim P_i$ .

180 Our goal is to develop a model capable of learning to trans-  
 181 port between any two distributions  $i$  and  $i'$ . Multimarginal  
 182 stochastic interpolants (Albergo et al., 2023a) developed a  
 183 class of models that can flow between any pair of distribu-  
 184 tions from a fixed initial set of  $K$  distributions by effectively  
 185 embedding each distribution as a corner of the  $K$ -simplex  
 186 and conditioning on the source and target corners. Our no-  
 187 tion of source-target conditioned transport models leverages  
 188 distribution encoders to learn distribution embeddings, en-  
 189 abling us to generalize to a new  $(K + 1)^{\text{th}}$  distribution, and  
 190 to embed a continuous number of distributions.

191  $\mathcal{T} : \mathcal{X} \times \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathcal{X}$ ,  $x_i \mapsto \hat{x}_{i \rightarrow i'} = \mathcal{T}(x_i | z_i, z_{i'})$ ,  
 192 to satisfy

$$\mathcal{T}(S_i | \mathcal{E}(S_i), \mathcal{E}(S_{i'})) \xrightarrow[m_i, m_{i'} \rightarrow \infty]{d} P_{i'} \quad (4)$$

193 Then, using the same plug-in loss logic (Prop. 3.1), we can  
 194 train on samples using the objective conditioning on the  
 195 encoder for the transport map of choice

$$\mathcal{L}_{\text{stc}} = \ell(S_{i'}, \mathcal{T}(S_i | \mathcal{E}(S_i), \mathcal{E}(S_{i'}))) \quad (5)$$

196 An important point is that our Alg. 1 does not incur any  
 197 quadratic cost *per gradient step*, even when training on all  
 198 pairs of distributions. The computational efficiency is the  
 199 same as the base transport model, up to a forward pass of  
 200 the distribution encoder.

### 3.4. Semi-supervised transport

201 An important use of the unsupervised transport modelling  
 202 objective is in partial supervision: we want to use all avail-  
 203 able data to learn a high-quality transport model, but at test  
 204 time, we are concerned with prediction of a specific target  
 205 distribution for each source. Source-target-conditioned  
 206 models can be applied in this semi-supervised setting through  
 207 two practical adaptations: (1) latent predictions of target  
 208 distribution embeddings and (2) sampling distribution pairs

---

**Algorithm 1** Training distribution-conditioned models

---

```

1: Input: Dataset  $\mathcal{D}$ , loss  $\ell$ , metadistribution  $Q_{\text{joint}}$ 
2: for each training step do
3:   Sample indices  $(u, v) \sim Q_{\text{joint}}$ 
4:   Subsample  $\hat{S}_u \subset S_u, \hat{S}_v \subset S_v$  // justified by Prop. 3.1
5:   // embed distributions and compute loss
6:   source-conditioned      source-target-conditioned
     $z_u \leftarrow \mathcal{E}(\hat{S}_u)$        $z_u, z_v \leftarrow \mathcal{E}(\hat{S}_u), \mathcal{E}(\hat{S}_v)$ 
     $\mathcal{L} \leftarrow \ell(\hat{S}_v, \mathcal{T}(\hat{S}_u | z_u))$        $\mathcal{L} \leftarrow \ell(\hat{S}_v, \mathcal{T}(\hat{S}_u | z_u, z_v))$ 
7:   Backpropagate  $\mathcal{L}$ 
8: end for

```

---

during training (i.e., choosing  $Q_{\text{joint}}$ ) in a way that aligns  
 with the supervised task at hand.

To make this concrete, we return to the lineage tracing case  
 where we want to leverage all the clones, but the test-time  
 objective is still clonal fate prediction, as in the supervised  
 setting. To convert the any-to-any model into a “supervised”  
 model we can fit a lightweight (e.g., linear) model to predict  
 $z_{\text{src}} \mapsto z_{\text{tgt}}$  using the subset of available (src, tgt) pairs  
 for the task of interest; we then evaluate  $\mathcal{T}(\cdot | z_{\text{src}}, \hat{z}_{\text{tgt}})$  to  
 generate counterfactual samples “as if” drawn from  $P_{\text{tgt}}$  in  
 a semi-supervised manner.

While the source-target conditioned model provides an interface  
 for any-to-any transport, we need not train on all possible pairs. In some scenarios, only certain pairs are of interest. For example, in lineage tracing, it makes more sense to learn forward-time transport between clones ( $t_1 \rightarrow t_2$ ) rather than all pairs including potentially uninformative  
 within-timepoint and reverse-time transport. In general, we can define a meta-distribution  $Q_{\text{joint}}$  over source-target pairs, allowing practitioners to encode domain-specific structure into the model.

**Proposition 3.1: Loss CLT (informal)**

Fix indices  $(u, v)$ . Let  $z_u = \mathcal{E}(\hat{S}_u)$  and  $z_v = \mathcal{E}(\hat{S}_v)$  be  
 embeddings from minibatches of sizes  $m$  and  $m'$ , with limits  
 $z_u^* = \phi(P_u)$  and  $z_v^* = \phi(P_v)$ . Under Fishman et al. (2025)  
 and the sampling schemes in App. A,

$$(\sqrt{m}(z_u - z_u^*), \sqrt{m'}(z_v - z_v^*)) \xrightarrow{d} \mathcal{N}(0, \Sigma_{uv}),$$

where  $\Sigma_{uv} = \text{diag}(\Sigma_u, \Sigma_v)$  for independent minibatches  
 (product coupling) and may have off-diagonal blocks for  
 paired/coupled minibatches. Consequently, losses are mean  
 consistent and admit plug-in CLTs; see App. A.

## 4. Related work

### 4.1. Distribution embeddings

Distribution conditioned transport relies on a vector rep-  
 resentation of a distribution, or a distribution embedding.

Kernel methods, including kernel mean embeddings, provide a toolkit for representing probability measures as points in a reproducing kernel Hilbert space (Smola et al., 2007; Muandet et al., 2012; Oliva et al., 2013; Szabo et al., 2015; Muandet et al., 2017). In parallel, others have explored learning distribution embeddings which preserve Wasserstein distances between distributions (Haviv et al., 2024), in analogy to multidimensional scaling. Our distribution encoders follows the formal framework in Fishman et al. (2025) to learn distributionally-invariant embeddings by leveraging permutation-invariant architectures with a particular class of pooling operators (Zaheer et al., 2017; Wagstaff et al., 2021; Zhang et al., 2022).

## 4.2. Transport models

A range of generative transport mechanisms can be used as the conditional map  $\mathcal{T}(\cdot | z)$  or  $\mathcal{T}(\cdot | z, z')$  in our framework. Classical adversarial approaches learn source–target mappings via discriminators, including GANs (Goodfellow et al., 2014) and Wasserstein GANs (Arjovsky et al., 2017). We can also transport distributions by minimizing distributional divergences (Gneiting & Raftery, 2007; Li et al., 2015; Sutherland et al., 2016; Genevay et al., 2018; Kolouri et al., 2019). Normalizing flows provide invertible transport parameterizations through continuous change-of-variable models (Rezende & Mohamed, 2015). Continuous-time flow-based models, such as flow matching, rectified flows, and stochastic interpolants (Lipman et al., 2023; Liu et al., 2022; Albergo et al., 2023b), define transport by learning a velocity field along a path interpolating between source and target.

On discrete sequences, unsupervised machine translation methods (Lample et al., 2018a; Artetxe et al., 2018; Lample et al., 2018b) learn bidirectional maps between languages using only monolingual corpora, providing a discrete analogue of our unpaired any-to-any transport between distributions. More recently, flow-matching has been extended to the simplex and general discrete domains (Gat & Lipman, 2024; Cheng et al., 2024; Davis et al., 2024).

All of these models define mappings that take samples from a source distribution and produce samples matching a target distribution. Our framework is orthogonal to the specific transport mechanism: any such model can be conditioned on distribution embeddings to yield a source-conditioned or source–target–conditioned transport map that satisfies the population-level limits in Eqn. (2) and Eqn. (4).

## 4.3. Multi-distribution transport

Meta Flow Matching (Atanackovic et al., 2024) develops the idea of source-conditioning for learning transport maps across contexts. Our notion of source-conditioned models generalizes their approach to a broader class of transport

models and demonstrates the conditions under which we can expect it to achieve Eqn. 2. A number of recent cellular perturbation response prediction methods are also effectively source-conditioned models (Klein et al., 2025; He et al., 2025; Adduri et al., 2025).

Multimarginal stochastic interpolants (MMSI) studies a form of any-to-any transport between a fixed set of  $K$  distributions (Albergo et al., 2023a). This is a particular instance of the unsupervised modelling problem we consider, where  $Q$  is chosen to cover a discrete set of distributions that are all seen at training time. Additionally, MMSI enforces that transport paths go through all  $K$  distributions to learn a barycenter, but in DCT we do not impose any constraint.

Style-transfer and unpaired image-translation methods provide a complementary line of work. Approaches such as CycleGAN, StarGAN, and MUNIT (Zhu et al., 2017; Choi et al., 2018; Huang et al., 2018) condition the generator on a domain identifier, typically a one-hot categorical label or an instance-level style code. These approaches are restricted to a fixed, finite set of domains and cannot generalize transport to a new domain without retraining.

## 4.4. Transport models for biological data

The family of transport models described above are widely used in biology. They have attracted significant recent interest for at least two problems: predicting cellular perturbation responses from scRNA-seq, imaging, and mass cytometry data (Atanackovic et al., 2024; Klein et al., 2025; Adduri et al., 2025; He et al., 2025; Bunne et al., 2023), and inferring cellular dynamics from snapshot observations (Schiebinger et al., 2019; Lavenant et al., 2021; Tong et al., 2023a;c; Vinyard et al., 2025).

## 5. Unsupervised Transport

Using both synthetic data and real scRNA-seq data, we will show that source-target-conditioned DCTs can address the any-to-any unsupervised transport problem in practice. We will compare source-target-conditioned DCTs using a mean pooled deep-sets encoder (see App. C.2) to  $K$ -to- $K$  models implemented using a “one-hot”  $\mathcal{E}$  rather than a distributionally invariant architecture. Using the FM transport map, this is effectively a MMSI model (Albergo et al., 2023a), with minor differences we discuss in App. B. To enable test-time generalization,  $K$ -to- $K$  models are conditioned on the most similar training set distributions, determined by nearest-neighbor over distribution centroids. We implement both approaches across sliced Wasserstein distance (SWD) models, Energy score models, and flow matching (FM) models (see App. C), and evaluate using MMD in the main text with SWD and Energy distances available for all experiments in the appendix (D).

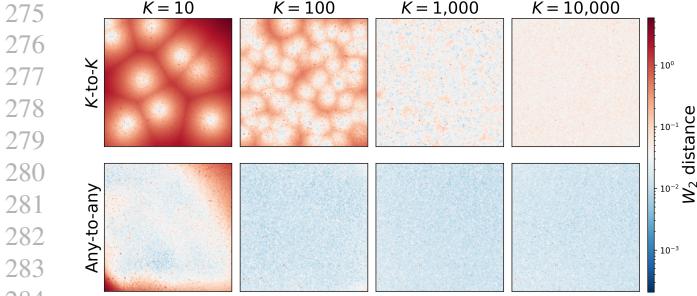


Figure 2. Bivariate normal transport error landscape.  $K$ -to- $K$  (top) vs. any-to-any (bottom), showing  $W_2$  distance for targets  $\mu \in [0, 5]^2$  across  $K$ . The  $K$ -to- $K$  model predicts via nearest training distribution, yielding Voronoi-like error patterns; the any-to-any model embeds targets directly and achieves uniformly lower error across all  $K$ .

### 5.1. Gaussian transport

We first study transport between simple distributions in two dimensions, where we can control the data-generating process and directly visualize performance.

We construct datasets of bivariate normal (MVN) distributions and Gaussian mixture models (GMM). For MVN, we sample parameters  $(\mu_i, \Sigma_i)$ ,  $i = 1, \dots, n$ , with means  $\mu_i \in [0, 5]^2$  drawn from a uniform prior and covariances  $\Sigma_i \in \mathbb{R}^{2 \times 2}$  drawn from an inverse-Wishart prior. For GMM, we additionally sample mixture weights from a Dirichlet prior and per-component parameters (see App. E.1 for details). We choose  $K \in \{10, 100, 1,000, 10,000\}$  distinct parameter sets and draw  $n = 50,000$  sample sets.

For each  $K$  we train two types of model. The first is a conventional  $K$ -to- $K$  architecture that treats each of the  $K$  distinct distributions as a discrete label (a “corner”) and learns transport conditional on distribution-specific labels. At test time, this model cannot condition on a new target distribution, so we assign each target to its nearest training distribution. The second is our *source–target conditioned* model from Sec. 3.3, which encodes both source and target sets via the distribution encoder  $z_i = \mathcal{E}(S_i)$  and learns a transport map  $\mathcal{T}(x \mid z_{\text{src}}, z_{\text{tgt}})$  that can operate between arbitrary pairs of distributions, including unseen targets.

To visualize zero-shot performance, we fix a single source distribution from the training set and evaluate transport to a dense grid of target distributions, with  $\mu \in [0, 5]^2$  and corresponding random covariances. Results (Figs. 2 and 3) show that the  $K$ -to- $K$  model outperforms DCT on in-distribution targets when  $K$  is small, but loses this advantage as  $K$  grows. And crucially, on out-of-distribution targets, DCT achieves significantly lower error: the baseline cannot extrapolate beyond its training corners, while the any-to-any model generalizes smoothly. This pattern holds not only for MVN but also GMM, suggesting that the benefits of DCT extend to complex distributions.

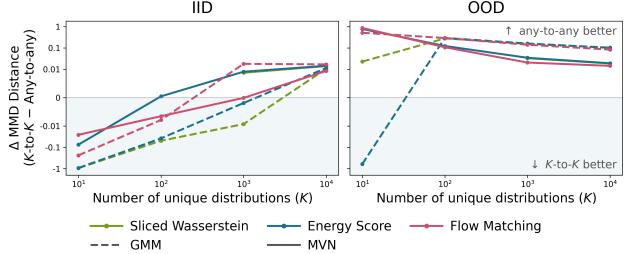


Figure 3. Unsupervised transport model generalization. The gap between  $K$ -to- $K$  based embedding and any-to-any distribution encoders across generator families, evaluated on in-distribution (IID) and out-of-distribution (OOD) test sets. Positive values indicate the distribution encoder achieves lower transport cost. At low  $K$ , the embedding encoder outperforms on IID data while the distribution encoder shows stronger OOD generalization.

Generator	$K$ -to- $K$	Any-to-any
SWD	$0.3411 \pm 0.0351$	$0.1167 \pm 0.0119$
Energy	$0.3623 \pm 0.0327$	<b><math>0.0729 \pm 0.0090</math></b>
FM	$0.4235 \pm 0.0540$	$0.1480 \pm 0.0157$
scVI		$0.8056 \pm 0.0643$

Table 1. Batch effect transfer quality on held-out samples. MMD distance ( $\downarrow$ ) averaged across all pairs of 3 held-out donors.

### 5.2. Transferring scRNA-seq batch effects

A common problem in scRNA-seq analysis is technical variation between experimental batches (Luecken et al., 2022). We show that source-target-conditioned models can be used to push cells from a source batch to a target batch condition, making a prediction for how cells appear under different technical effects. This problem, which we refer to as batch effect transfer, is closely related to batch integration.

Source-target-conditioning, in principle, allows us to fit a single batch transfer model which can be applied to unseen experimental batches at test time. We demonstrate this using a real-world dataset of murine pancreas cells from 56 mice across ages and growth conditions (Hrovatin et al., 2023). We hold out all mice from a single condition (old age), inducing a distribution shift.

We might expect that the  $K$ -to- $K$  models are not be able to adapt to distribution shift, given that they can only condition on distribution identities seen at test time, while DCTs in principle may be able to encode information about the correct target distribution.

DCTs outperform  $K$ -to- $K$  alternatives and scVI (Lopez et al., 2018) on transport of cells from held-out experimental batches across all tested generator families (Table 1, implementation details in App. G). In Fig. 4, we visualize model predictions for batch effect transfer, which similarly show that DCT predictions more closely match the target cell distribution than those of a  $K$ -to- $K$  model.

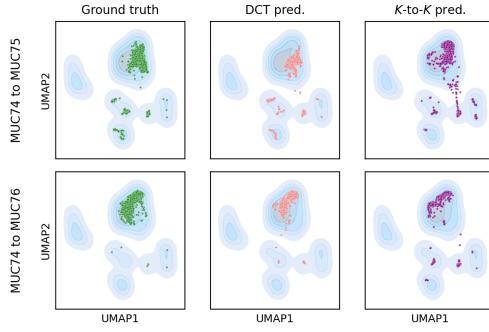


Figure 4. Batch effect transfer for two held-out donor pairs (rows). Blue contours represent the reference density ( $\sim 3 \cdot 10^5$  cells) computed from all 56 donors. (left) Ground truth cells from the target donor. (center) predictions from DCT and (right)  $K$ -to- $K$  model. DCT predictions more closely match ground truth than  $K$ -to- $K$  predictions.

## 6. Semi-supervised Transport

Across synthetic data, mass cytometry data, and protein sequence data, we will show that source-target-conditioned DCTs can improve performance in the semi-supervised setting. We compare source-target-conditioned DCTs with their source-conditioned counterparts, across SWD-based, Energy-based, and FM transport maps (see App. C). We evaluate using three distributional metrics (MMD, SWD, and Energy) with MMD presented in the main text and all three available in the App. D.

The source-conditioned models we evaluate against include the source-conditioned FM, which is identical to MFM (Atanackovic et al., 2024) aside from the choice of a CLT-satisfying population encoder (see App. B for discussion). All semi-supervised approaches use ridge regression fit on the same data as the supervised (SC) model, with the regularization selected through cross-validation. We also compare with ‘‘oracle’’ models which are STC models allowed to condition on the embedding of the true distribution. A central finding here is that in real data applications we can often improve OOD generalization of supervised models by incorporating a wider set of data without explicit paired distributions (see the orphan marginals discussion in Sec. 3.4).

### 6.1. Gaussian Transport

Returning to our Gaussian setting, we build a benchmark for semi-supervised learning. We define paired distributions via fixed transformations: for MVN, a simple shift  $Y = X + b$ ; for GMM, a bimodal target with off-axis displacements. We compare three approaches: (i) a *supervised* source-conditioned model trained on paired examples with means  $\|\mu\|_\infty \leq 2.5$ , (ii) the *semi-supervised* any-to-any model from Sec. 5.1, and (iii) an *oracle* using the true target

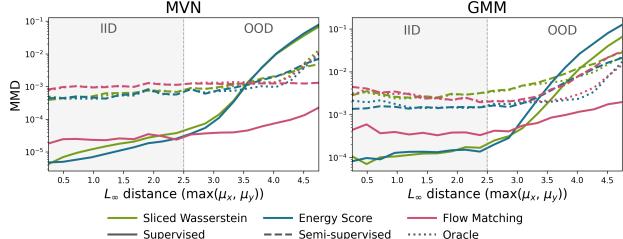


Figure 5. Semi-supervised transport model generalization. Supervised models are trained on distributions with  $\|\mu\|_\infty \leq 2.5$  (shaded) and evaluated across  $\|\mu\|_\infty \leq 5$ . Supervised models (red) degrade sharply outside the training support (except flow matching models), while semi-supervised models (blue) maintain stable performance approaching the oracle (green). Results shown for multivariate normal (left) and Gaussian mixture (right).

embedding. All models are evaluated across  $\|\mu\|_\infty \leq 5$ , testing generalization beyond the supervised training support.

Figure 5 shows energy distance as a function of  $L^\infty$  distance for both settings. Within the training support, supervised and semi-supervised methods achieve comparable performance. Beyond it, the semi-supervised approach maintains substantially lower error: the any-to-any objective learns distributional structure that enables extrapolation even when the latent predictor must generalize. The flow matching model is an exception here: it generalizes very well outside its support. The extremely strong performance of FM models here may be connected to the tightly coupled transport maps learned by FM (see App. A.2), but this performance is not replicated in any of our real-world experiments.

### 6.2. Drug perturbation prediction on organoids

We evaluate DCTs on a single-cell mass-cytometry drug screen from Zapatero et al. (2023), which profiles patient-derived organoids (PDOs) from 10 colorectal cancer patients treated with 11 chemotherapy compounds. Each experimental replicate consists of a matched control-treatment pair, yielding 927 population pairs. Cells are characterized by an abundance profile over 43 proteins.

A key challenge is predicting patient-specific treatment response: tumors from different patients exhibit distinct drug sensitivities. The task is to predict the post-treatment distribution given the control population and treatment identity. We evaluate on two splits: (i) held-out replicates within patients seen during training, and (ii) held-out patients, testing generalization.

We compare a source-conditioned baseline (including MFM) against our semi-supervised source-target conditioned model. Both models condition on treatment identity. As shown in Table 2, the source-conditioned model achieves lower error on held-out replicates within pa-

		Supervised	Semi-sup.	Oracle
IID	SWD	0.64 $\pm$ .22	0.94 $\pm$ .18	0.10 $\pm$ .00
	Energy	<b>0.54<math>\pm</math>.19</b>	1.03 $\pm$ .32	0.11 $\pm$ .02
	FM	1.01 $\pm$ .23	1.06 $\pm$ .25	0.30 $\pm$ .03
OOD	SWD	2.59 $\pm$ .08	1.97 $\pm$ .22	0.26 $\pm$ .09
	Energy	2.98 $\pm$ .19	<b>1.91<math>\pm</math>.49</b>	0.22 $\pm$ .02
	FM	2.27 $\pm$ .06	2.08 $\pm$ .18	0.53 $\pm$ .13

Table 2. Perturbation prediction performance on PDO task. MMD distance  $\times 10^2$  ( $\downarrow$ ) averaged over held-out patients.

	Supervised	Semi-supervised	Oracle
SWD	9.87 $\pm$ .27	8.67 $\pm$ .22	2.99 $\pm$ .05
Energy	9.42 $\pm$ .26	<b>8.54<math>\pm</math>.21</b>	2.88 $\pm$ .05
FM	14.10 $\pm$ .28	9.67 $\pm$ .25	5.08 $\pm$ .10

Table 3. Clonal distribution forecasting on Weinreb et al. (2020). MMD distance ( $\downarrow$ ) across 628 held-out clones.

tients, but the any-to-any model generalizes better to unseen patients—echoing the IID/OOD tradeoff observed in our Gaussian benchmarks (Fig. 3).

### 6.3. Modeling clonal population dynamics

We next apply DCTs to learn clonal population dynamics using lineage-traced single-cell RNA-sequencing (scRNA-seq) data from Weinreb et al. (2020). A critical challenge in this dataset is sparsity: there are about  $6 \cdot 10^3$  clones measured in total, however only about  $2 \cdot 10^3$  of these clones are profiled at multiple timepoints. While source-conditioned approaches can only be trained on clones observed at multiple timepoints, training DCTs with any-to-any pairings allow us to make use of all available marginals. Here, when training without coupling information we restrict source-target pairing to only allow an early timepoint clone as a source distribution and a late timepoint clone as a target distribution.

We compare source-conditioned (supervised) and source-target-conditioned (semi-supervised) approaches using an identical architectures operating on the first 50 principal components of gene expression profiles (implementation details in App. I). As shown in Table 6.3, the source-target-conditioned models significantly outperform the source-conditioned baselines, demonstrating the benefit of incorporating orphan marginals.

### 6.4. Longitudinal TCR repertoire forecasting

We evaluate DCT on longitudinal T-cell receptor (TCR) repertoire sequencing from COVID-19 patients (Schultheiß et al., 2020). The dataset contains repertoires from 37 patients, of whom only 10 are profiled at multiple timepoints –

a sparsity pattern mirroring the orphan marginal regime in lineage tracing. We hold out three longitudinal patients for evaluation and train on the remainder. The task is to forecast a patient’s repertoire at time  $t+1$  given their repertoire at  $t$ .

Each repertoire is an empirical distribution over TRB CDR3 amino acid sequences. We embed sequences using a pre-trained ESM backbone (Lin et al., 2023) with mean-pooled DeepSets aggregation. We compare two discrete transport mechanisms: a ProGen-based bridge (Nijkamp et al., 2023) and a discrete flow matching (DFM) bridge (Gat & Lipman, 2024), both trained in the same distribution-conditioned framework. The supervised baseline trains only on within-patient adjacent pairs ( $t, t+1$ ); the semi-supervised model trains on cross-patient pairs (earlier-to-later timepoints), incorporating single-timepoint patients as unpaired marginals.

	Supervised	Semi-supervised	Oracle
Progen	0.0588 $\pm$ .009	0.0587 $\pm$ .009	0.0567 $\pm$ .009
DFM	0.0579 $\pm$ .010	<b>0.0215<math>\pm</math>.004</b>	0.0078 $\pm$ .001

Table 4. TCR repertoire forecasting performance. MMD distance ( $\downarrow$ ) averaged over held-out patients.

The ProGen model appears degenerate; it learns identical embeddings for all distributions with average cosine similarity 0.998 (see App. J.5). The semi-supervised DFM model reduces energy distance by more than half, demonstrating that cross-patient distributional structure improves forecasting even with limited longitudinal supervision.

## 7. Conclusion

We introduced *distribution-conditioned transport*, a framework that couples distribution encoders with conditional transport models to learn transport maps which can generalize to unseen source and target distributions. We showed that DCT formalizes and generalizes existing approaches under a unified framework, leading to new models with empirical performance gains on synthetic benchmarks and four real world problems in biology.

**Limitations and extensions** Although source-target-conditioned models outperform source-conditioned and K-to-K baselines out-of-distribution, they sometimes underperform in-distribution. This may be due to underfitting: while all compared models are trained with equal compute budgets, source-target-conditioned models may have different scaling properties and require more training. In semi-supervised problems, we fit simple linear latent predictors for source-target-conditioned models. It may be fruitful to explore more complex latent prediction schemes.

## 440 Impact Statement

441 This paper presents work whose goal is to advance the field  
 442 of Machine Learning. There are many potential societal  
 443 consequences of our work, none which we feel must be  
 444 specifically highlighted here.

## 445 References

- 446 Adduri, A. K., Gautam, D., Bevilacqua, B., Imran, A., Shah,  
 447 R., Naghipourfar, M., Teyssier, N., Ilango, R., Nagaraj,  
 448 S., Dong, M., et al. Predicting cellular responses to  
 449 perturbation across diverse contexts with state. *bioRxiv*,  
 450 pp. 2025–06, 2025.
- 451 Albergo, M. S., Boffi, N. M., Lindsey, M., and Vanden-  
 452 Eijnden, E. Multimarginal generative modeling with  
 453 stochastic interpolants. *arXiv preprint arXiv:2310.03695*,  
 454 2023a.
- 455 Albergo, M. S., Boffi, N. M., and Vanden-Eijnden, E.  
 456 Stochastic interpolants: A unifying framework for flows  
 457 and diffusions. *Journal of Machine Learning Research*,  
 458 24:1–63, 2023b.
- 459 Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein  
 460 GAN. In *Proceedings of the 34th International Conference  
 461 on Machine Learning*, volume 70, pp. 214–223. PMLR,  
 462 2017.
- 463 Artetxe, M., Labaka, G., Agirre, E., and Cho, K. Unsu-  
 464 pervised neural machine translation. In *International  
 465 Conference on Learning Representations*, 2018.
- 466 Atanackovic, L., Zhang, X., Amos, B., Blanchette, M., Lee,  
 467 L. J., Bengio, Y., Tong, A., and Neklyudov, K. Meta Flow  
 468 Matching: Integrating Vector Fields on the Wasserstein  
 469 Manifold. In *The Thirteenth International Conference on  
 470 Learning Representations*, 4 October 2024.
- 471 Biddy, B. A., Kong, W., Kamimoto, K., Guo, C., Waye,  
 472 S. E., Sun, T., and Morris, S. A. Single-cell mapping  
 473 of lineage and identity in direct reprogramming. *Nature*,  
 474 564(7735):219–224, 2018.
- 475 Bunne, C., Stark, S. G., Gut, G., Del Castillo, J. S.,  
 476 Levesque, M., Lehmann, K.-V., Pelkmans, L., Krause,  
 477 A., and Rätsch, G. Learning single-cell perturbation  
 478 responses using neural optimal transport. *Nature Methods*,  
 479 20(11):1759–1768, 28 November 2023. ISSN 1548-  
 480 7091,1548-7105. doi: 10.1038/s41592-023-01969-x.
- 481 Cheng, X., Chen, T., Li, K., Gao, R., Zhu, J., and Liu, Q.  
 482 Categorical flow matching. In *International Conference  
 483 on Learning Representations*, 2024.
- 484 Choi, Y., Choi, M., Kim, M., Ha, J.-W., Kim, S., and Choo,  
 485 J. StarGAN: Unified generative adversarial networks for  
 486 multi-domain image-to-image translation. In *Proceedings  
 487 of the IEEE Conference on Computer Vision and Pattern  
 488 Recognition (CVPR)*, pp. 8789–8797, 2018.
- 489 Davis, N., Chen, R. T. Q., Hu, Y., and Rezende, D. J. Fisher  
 490 flow matching: Transport in the probability simplex. In  
 491 *International Conference on Machine Learning*, 2024.
- 492 Fishman, N., Gowri, G., Yin, P., Gootenberg, J., and Abu-  
 493 dayyeh, O. Generative distribution embeddings. *arXiv  
 494 preprint arXiv:2505.18150*, 2025.
- 495 Gat, I. and Lipman, Y. Discrete flow matching. In *Advances  
 496 in Neural Information Processing Systems*, 2024.
- 497 Genevay, A., Peyré, G., and Cuturi, M. Learning generative  
 498 models with sinkhorn divergences. In *Proceedings of  
 499 the Twenty-First International Conference on Artificial  
 500 Intelligence and Statistics*, volume 84, pp. 1608–1617.  
 501 PMLR, 2018.
- 502 Gneiting, T. and Raftery, A. E. Strictly proper scoring  
 503 rules, prediction, and estimation. *Journal of the American  
 504 statistical Association*, 102(477):359–378, 2007.
- 505 Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B.,  
 506 Warde-Farley, D., Ozair, S., Courville, A., and Bengio,  
 507 Y. Generative adversarial nets. In *Advances in Neural  
 508 Information Processing Systems*, volume 27, 2014.
- 509 Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B.,  
 510 and Smola, A. A kernel two-sample test. *The journal of  
 511 machine learning research*, 13(1):723–773, 2012.
- 512 Haviv, D., Kunes, R. Z., Dougherty, T., Burdziak, C., Nawy,  
 513 T., Gilbert, A., and Pe'er, D. Wasserstein Wormhole:  
 514 Scalable optimal transport distance with transformers. *In-  
 515 ternational Conference on Machine Learning*, 235:17697–  
 516 17718, 15 April 2024. doi: 10.48550/arXiv.2404.09411.
- 517 He, C., Zhang, J., Dahleh, M., and Uhler, C. Morph predicts  
 518 the single-cell outcome of genetic perturbations across  
 519 conditions and data modalities. *bioRxiv*, 2025.
- 520 Heumos, L., Schaar, A. C., Lance, C., Litinetskaya, A.,  
 521 Drost, F., Zappia, L., Lücke, M. D., Strobl, D. C., Henao,  
 522 J., Curion, F., Single-cell Best Practices Consortium,  
 523 Schiller, H. B., and Theis, F. J. Best practices for single-  
 524 cell analysis across modalities. *Nature reviews. Genetics*,  
 525 pp. 1–23, 31 March 2023. ISSN 1471-0056,1471-0064.  
 526 doi: 10.1038/s41576-023-00586-w.
- 527 Hrovatin, K., Bastidas-Ponce, A., Bakhti, M., Zappia, L.,  
 528 Büttner, M., Salinno, C., Sterr, M., Böttcher, A., Migliorini,  
 529 A., Lickert, H., and Theis, F. J. Delineating mouse  
 530  $\beta$ -cell identity during lifetime and in diabetes with a  
 531 single cell atlas. *Nature Metabolism*, 5(9):1615–1637,  
 532 11 September 2023. ISSN 2522-5812,2522-5812. doi:  
 533 10.1038/s42255-023-00876-x.

- 495 Huang, X., Liu, M.-Y., Belongie, S., and Kautz, J. Multi-  
 496 modal unsupervised image-to-image translation. In *Proceedings of the European Conference on Computer Vision*  
 497 (*ECCV*), pp. 179–196, 2018.
- 498
- 499 Klein, D., Fleck, J. S., Bobrovskiy, D., Zimmermann, L.,  
 500 Becker, S., Palma, A., Dony, L., Tejada-Lapuerta, A.,  
 501 Huguet, G., Lin, H.-C., et al. Cellflow enables generative  
 502 single-cell phenotype modeling with flow matching.  
*bioRxiv*, pp. 2025–04, 2025.
- 503
- 504 Kolouri, S., Nadjahi, K., Simsekli, U., Badeau, R., and  
 505 Rohde, G. Generalized sliced wasserstein distances. *Advances in neural information processing systems*, 32,  
 506 2019.
- 507
- 508 Lample, G., Conneau, A., Denoyer, L., and Ranzato, M.  
 509 Unsupervised machine translation using monolingual corpora only. In *International Conference on Learning Representations*, 2018a.
- 510
- 511 Lample, G., Ott, M., Conneau, A., Denoyer, L., and Ranzato, M.  
 512 Phrase-based & neural unsupervised machine  
 513 translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp.  
 514 5039–5049. Association for Computational Linguistics,  
 515 2018b.
- 516
- 517 Lavenant, H., Zhang, S., Kim, Y.-H., and Schiebinger, G.  
 518 Towards a mathematical theory of trajectory inference.  
*arXiv [stat.ML]*, (1A):428–500, 18 February 2021.
- 519
- 520 LeCun, Y., Cortes, C., and Burges, C. Mnist hand-  
 521 written digit database. *ATT Labs [Online]*. Available:  
<http://yann.lecun.com/exdb/mnist>, 2, 2010.
- 522
- 523 Li, Y., Swersky, K., and Zemel, R. S. Generative moment  
 524 matching networks. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37, pp.  
 525 1718–1727. PMLR, 2015.
- 526
- 527 Lin, Z., Akin, H., Rao, R., Hie, B., Zhu, Z., Lu, W.,  
 528 Smetanin, N., Verkuil, R., Kabeli, O., Shmueli, Y., et al.  
 529 Evolutionary-scale prediction of atomic-level protein  
 530 structure with a language model. *Science*, 379(6637):  
 531 1123–1130, 2023.
- 532
- 533 Lipman, Y., Chen, R. T. Q., Ben-Hamu, H., Nickel, M.,  
 534 and Le, M. Flow matching for generative modeling. In  
 535 *International Conference on Learning Representations*,  
 536 2023.
- 537
- 538 Liu, X., Gong, C., and Liu, Q. Flow straight and fast:  
 539 Learning to generate and transfer data with rectified flow.  
 540 In *Advances in Neural Information Processing Systems*,  
 541 volume 35, 2022.
- 542
- 543 Lopez, R., Regier, J., Cole, M. B., Jordan, M. I., and  
 544 Yosef, N. Deep generative modeling for single-cell  
 545 transcriptomics. *Nature methods*, 15(12):1053–1058,  
 546 December 2018. ISSN 1548-7091,1548-7105. doi:  
 547 10.1038/s41592-018-0229-2.
- 548
- 549 Luecken, M. D., Büttner, M., Chaichoompu, K., Danese, A., Interlandi, M., Mueller, M. F., Strobl, D. C., Zap-  
 550 pia, L., Dugas, M., Colomé-Tatché, M., and Theis, F. J. Benchmarking atlas-level data integration in single-  
 551 cell genomics. *Nature Methods*, 19(1):41–50, January 2022. ISSN 1548-7091,1548-7105. doi: 10.1038/  
 552 s41592-021-01336-8.
- 553
- 554 Muandet, K., Fukumizu, K., Dinuzzo, F., and Scholkopf, B.  
 555 Learning from distributions via support measure machines. *Neural Information Processing Systems*, 25:10–  
 556 18, 29 February 2012.
- 557
- 558 Muandet, K., Fukumizu, K., Sriperumbudur, B., and Schölkopf, B. Kernel mean embedding of distributions: A  
 559 review and beyond. *Foundations and Trends® in Machine Learning*, 10(1-2):1–141, 2017. ISSN 1935-8237,1935-  
 560 8245. doi: 10.1561/2200000060.
- 561
- 562 Nijkamp, E., Ruffolo, J. A., Weinstein, E. N., Naik, N., and Madani, A. ProGen2: Exploring the boundaries  
 563 of protein language models. *Cell systems*, 14(11):968–  
 564 978.e3, 15 November 2023. ISSN 2405-4712,2405-4720.  
 565 doi: 10.1016/j.cels.2023.10.002.
- 566
- 567 Oliva, J. B., Póczos, B., and Schneider, J. Distribution to  
 568 Distribution Regression. *International Conference on Machine Learning*, 28(3):1049–1057, 16 June 2013.
- 569
- 570 Qiu, C., Martin, B. K., Welsh, I. C., Daza, R. M., Le, T.-  
 571 M., Huang, X., Nichols, E. K., Taylor, M. L., Fulton,  
 572 O., O’Day, D. R., Gomes, A. R., Ilcisin, S., Srivatsan,  
 573 S., Deng, X., Disteche, C. M., Noble, W. S., Hamazaki,  
 574 N., Moens, C. B., Kimelman, D., Cao, J., Schier, A. F.,  
 575 Spielmann, M., Murray, S. A., Trapnell, C., and Shendure,  
 576 J. A single-cell time-lapse of mouse prenatal development  
 577 from gastrula to birth. *Nature*, 626(8001):1084–1093,  
 578 14 February 2024. ISSN 0028-0836,1476-4687. doi:  
 579 10.1038/s41586-024-07069-w.
- 580
- 581 Rezende, D. J. and Mohamed, S. Variational inference with  
 582 normalizing flows. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37, pp.  
 583 1530–1538. PMLR, 2015.
- 584
- 585 Schiebinger, G., Shu, J., Tabaka, M., Cleary, B., Subra-  
 586 manian, V., Solomon, A., Gould, J., Liu, S., Lin, S.,  
 587 Berube, P., Lee, L., Chen, J., Brumbaugh, J., Rigol-  
 588 let, P., Hochedlinger, K., Jaenisch, R., Regev, A., and  
 589 Lander, E. S. Optimal-Transport Analysis of Single-  
 590 Cell Gene Expression Identifies Developmental Trajec-  
 591 tories in Reprogramming. *Cell*, 176(4):928–943.e22,

- 550 7 February 2019. ISSN 0092-8674,1097-4172. doi:  
551 10.1016/j.cell.2019.01.006.
- 552 Schultheiß, C., Paschold, L., Simnica, D., Mohme, M.,  
553 Willscher, E., von Wenserski, L., Scholz, R., Wieters, I.,  
554 Dahlke, C., Tolosa, E., et al. Next-generation sequencing  
555 of t and b cell receptor repertoires from covid-19 patients  
556 showed signatures associated with severity of disease.  
557 *Immunity*, 53(2):442–455, 2020.
- 558
- 559 Smola, A., Gretton, A., Song, L., and Schölkopf, B. A  
560 Hilbert space embedding for distributions. In *Lecture*  
561 *Notes in Computer Science*, Lecture notes in computer sci-  
562 ence, pp. 13–31. Springer Berlin Heidelberg, Berlin, Hei-  
563 delberg, 2007. ISBN 9783540752240,9783540752257.  
564 doi: 10.1007/978-3-540-75225-7\5.
- 565 Sutherland, D. J., Tung, H.-Y., Strathmann, H., De, S., Ram-  
566 das, A., Smola, A., and Gretton, A. Generative models  
567 and model criticism via optimized maximum mean dis-  
568 crepancy. *arXiv preprint arXiv:1611.04488*, 2016.
- 569
- 570 Szabo, Z., Gretton, A., Poczos, B., and Sriperumbudur, B.  
571 Two-stage sampled learning theory on distributions. In  
572 *Artificial Intelligence and Statistics*, pp. 948–957. PMLR,  
573 21 February 2015.
- 574
- 575 Tong, A., Fatras, K., Malkin, N., Huguet, G., Zhang, Y.,  
576 Rector-Brooks, J., Wolf, G., and Bengio, Y. Improv-  
577 ing and generalizing flow-based generative models with  
578 minibatch optimal transport. *arXiv [cs.LG]*, 1 February  
579 2023a.
- 580
- 581 Tong, A., Fatras, K., Malkin, N., Huguet, G., Zhang, Y.,  
582 Rector-Brooks, J., Wolf, G., and Bengio, Y. Improv-  
583 ing and generalizing flow-based generative models with mini-  
584 batch optimal transport. *arXiv preprint arXiv:2302.00482*,  
585 2023b.
- 586
- 587 Tong, A., Malkin, N., Fatras, K., Atanackovic, L., Zhang,  
588 Y., Huguet, G., Wolf, G., and Bengio, Y. Simulation-free  
589 Schrödinger bridges via score and flow matching. *arXiv*  
590 [*cs.LG*], 7 July 2023c.
- 591
- 592 Van Der Vaart, A. W. and Wellner, J. A. *Weak convergence*.  
593 Springer, 1996.
- 594
- 595 Vinyard, M. E., Rasmussen, A. W., Li, R., Klein, A. M.,  
596 Getz, G., and Pinello, L. Learning cell dynamics with  
597 neural differential equations. *Nature Machine Intelli-  
598 gence*, 7(12):1969–1984, 18 December 2025. ISSN 2522-  
599 5839,2522-5839. doi: 10.1038/s42256-025-01150-3.
- 600
- 601 Wagner, D. E. and Klein, A. M. Lineage tracing meets  
602 single-cell omics: opportunities and challenges. *Nature*  
603 *Reviews Genetics*, 21(7):410–427, 2020.
- 604
- Wagstaff, E., Fuchs, F., Engelcke, M., Osborne, M. A., and  
Posner, I. Universal approximation of functions on sets.  
*Journal of machine learning research: JMLR*, 23(151):  
151:1–151:56, 5 July 2021. ISSN 1532-4435,1533-7928.
- Weinreb, C., Rodriguez-Fraticelli, A., Camargo, F. D., and  
Klein, A. M. Lineage tracing on transcriptional land-  
scapes links state to fate during differentiation. *Science*,  
367(6479), 14 February 2020. ISSN 0036-8075,1095-  
9203. doi: 10.1126/science.aaw3381.
- Yazar, S., Alquicira-Hernandez, J., Wing, K., Senabouth, A.,  
Gordon, M. G., Andersen, S., Lu, Q., Rowson, A., Taylor,  
T. R. P., Clarke, L., Maccora, K., Chen, C., Cook, A. L.,  
Ye, C. J., Fairfax, K. A., Hewitt, A. W., and Powell, J. E.  
Single-cell eQTL mapping identifies cell type-specific  
genetic control of autoimmune disease. *Science (New  
York, N.Y.)*, 376(6589):eabf3041, 8 April 2022. ISSN  
0036-8075,1095-9203. doi: 10.1126/science.abf3041.
- Zaheer, M., Kottur, S., Ravanbakhsh, S., Póczos, B.,  
Salakhutdinov, R., and Smola, A. Deep Sets. *Advances  
in neural information processing systems*, 30, 10 March  
2017.
- Zapatero, M. R., Tong, A., Opzoomer, J. W., O’Sullivan,  
R., Rodriguez, F. C., Sufi, J., Vlckova, P., Nattress, C.,  
Qin, X., Claus, J., et al. Trellis tree-based analysis re-  
veals stromal regulation of patient-derived organoid drug  
responses. *Cell*, 186(25):5606–5619, 2023.
- Zhang, L. H., Tozzo, V., Higgins, J., and Ranganath, R.  
Set norm and equivariant skip connections: Putting the  
deep in Deep Sets. *International Conference on Machine  
Learning*, 162:26559–26574, 23 June 2022. doi: 10.  
48550/arXiv.2206.11925.
- Zhu, J.-Y., Park, T., Isola, P., and Efros, A. A. Unpaired  
image-to-image translation using cycle-consistent adver-  
sarial networks. In *Proceedings of the IEEE International  
Conference on Computer Vision (ICCV)*, pp. 2223–2232,  
2017.

## A. Theory

### A.1. Why can we train on minibatches?

This appendix formalizes two ingredients used throughout the main text: (i) distribution encoders factor through the empirical measure, and (ii) encoder CLTs yield plug-in CLTs for objectives that condition on the encoder state. The plug-in loss theory is adapted from Fishman et al. (2025, Assump. (C.1)–(C.3)); we include short proofs to make the manuscript self-contained and to handle two DCT-specific points: (a) joint CLTs for source–target conditioned models when the source and target minibatches are *coupled* (paired data), and (b) the small dependence introduced when a loss conditions on a point  $x$  that is reused inside the same minibatch used to compute  $\mathcal{E}(\hat{S})$ . We refer to Fishman et al. (2025); Van Der Vaart & Wellner (1996) for broader context.

#### A.1.1. EMPIRICAL-MEASURE FACTORIZATION

Let  $(\mathcal{X}, \mathcal{B})$  be a measurable space and let  $\mathcal{P}(\mathcal{X})$  denote probability measures on it. For  $m \in \mathbb{N}$  and  $S_m = (x_1, \dots, x_m) \in \mathcal{X}^m$ , write the empirical measure

$$\hat{P}_m := \frac{1}{m} \sum_{j=1}^m \delta_{x_j}.$$

These factorization statements mirror the empirical-measure sufficiency/maximal-invariance viewpoint in Fishman et al. (2025, Thm. “Empirical measure is sufficient, minimal, and a maximal invariant”); we include short proofs for completeness and to match DCT notation.

**Definition A.1** (Permutation and proportional invariance). Fix  $m \in \mathbb{N}$ . A map  $\mathcal{E}_m : \mathcal{X}^m \rightarrow \mathbb{R}^d$  is *permutation invariant* if

$$\mathcal{E}_m(x_1, \dots, x_m) = \mathcal{E}_m(x_{\pi(1)}, \dots, x_{\pi(m)}) \quad \text{for all permutations } \pi \in \mathfrak{S}_m.$$

A family  $(\mathcal{E}_m)_{m \geq 1}$  is *proportionally invariant* if for every integer  $k \geq 1$  and every  $(x_1, \dots, x_m) \in \mathcal{X}^m$ ,

$$\mathcal{E}_m(x_1, \dots, x_m) = \mathcal{E}_{km}(\underbrace{x_1, \dots, x_m, \dots, x_1, \dots, x_m}_{k \text{ copies}}).$$

**Proposition A.2** (Permutation invariance  $\Rightarrow$  factorization). Fix  $m \in \mathbb{N}$  and let  $\mathcal{E}_m : \mathcal{X}^m \rightarrow \mathbb{R}^d$  be permutation invariant. Then there exists a measurable map  $\phi_m$  defined on the set of empirical measures  $\{\hat{P}_m : S_m \in \mathcal{X}^m\}$  such that

$$\mathcal{E}_m(S_m) = \phi_m(\hat{P}_m) \quad \text{for all } S_m \in \mathcal{X}^m.$$

*Proof.* Let  $\hat{P}_m : \mathcal{X}^m \rightarrow \mathcal{P}(\mathcal{X})$  be the empirical-measure map and equip its image  $\{\hat{P}_m(S_m) : S_m \in \mathcal{X}^m\}$  with the quotient  $\sigma$ -algebra induced by  $\hat{P}_m$ . If  $\hat{P}_m(S_m) = \hat{P}_m(S'_m)$ , then  $S'_m$  is a permutation of  $S_m$  (as a multiset with multiplicity), so permutation invariance implies  $\mathcal{E}_m(S_m) = \mathcal{E}_m(S'_m)$ ; hence  $\mathcal{E}_m$  is constant on fibers of  $\hat{P}_m$ . Therefore there exists a well-defined measurable  $\phi_m$  on the image of  $\hat{P}_m$  such that  $\mathcal{E}_m = \phi_m \circ \hat{P}_m$ , i.e.  $\mathcal{E}_m(S_m) = \phi_m(\hat{P}_m)$ .  $\square$

**Proposition A.3** (Permutation + proportional invariance  $\Rightarrow$  a single functional). Let  $(\mathcal{E}_m)_{m \geq 1}$  be permutation and proportionally invariant. Then there exists a single measurable map  $\phi$  defined on the set of all empirical measures  $\bigcup_{m \geq 1} \{\hat{P}_m : S_m \in \mathcal{X}^m\}$  such that

$$\mathcal{E}_m(S_m) = \phi(\hat{P}_m) \quad \text{for all } m \geq 1, S_m \in \mathcal{X}^m.$$

*Proof.* By Prop. A.2, for each  $m$  there exists  $\phi_m$  with  $\mathcal{E}_m(S_m) = \phi_m(\hat{P}_m)$ . Define  $\phi$  on empirical measures by  $\phi(\hat{P}_m) := \phi_m(\hat{P}_m)$ . If the same empirical measure  $\nu$  can be represented at two sample sizes  $m$  and  $m'$ , then necessarily  $m' = km$  for some integer  $k$  and the  $m'$ -sample is obtained by duplicating each of the  $m$  atoms  $k$  times. Proportional invariance then implies  $\phi_m(\nu) = \phi_{m'}(\nu)$ , so  $\phi$  is well defined.  $\square$

#### A.1.2. ENCODER CLT VIA THE FUNCTIONAL DELTA METHOD

Let  $P \in \mathcal{P}(\mathcal{X})$  and let  $X_1, X_2, \dots \stackrel{\text{iid}}{\sim} P$  with empirical measure  $P_m = \frac{1}{m} \sum_{i=1}^m \delta_{X_i}$ . Assume the encoder is distributionally invariant, so that  $\mathcal{E}(S_m) = \phi(P_m)$  for a functional  $\phi : \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R}^d$ .

660 **Theorem A.4** (Encoder CLT). Suppose  $\phi : \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R}^d$  is Hadamard differentiable at  $P$  tangentially to a subspace  
 661 containing the empirical process limit, and denote its derivative by  $D\phi_P$ . Then

$$662 \quad 663 \quad 664 \quad 665 \quad 666 \quad 667 \quad 668 \quad 669 \quad 670 \quad 671 \quad 672 \quad 673 \quad 674 \quad 675 \quad 676 \quad 677 \quad 678 \quad 679 \quad 680 \quad 681 \quad 682 \quad 683 \quad 684 \quad 685 \quad 686 \quad 687 \quad 688 \quad 689 \quad 690 \quad 691 \quad 692 \quad 693 \quad 694 \quad 695 \quad 696 \quad 697 \quad 698 \quad 699 \quad 700 \quad 701 \quad 702 \quad 703 \quad 704 \quad 705 \quad 706 \quad 707 \quad 708 \quad 709 \quad 710 \quad 711 \quad 712 \quad 713 \quad 714$$

$$\sqrt{m} (\phi(P_m) - \phi(P)) \xrightarrow{d} D\phi_P(\mathbb{G}_P),$$

where  $\mathbb{G}_P$  is the  $P$ -Brownian bridge. In particular, if  $D\phi_P(\mathbb{G}_P)$  is a mean-zero Gaussian in  $\mathbb{R}^d$  with covariance  $\Sigma_P$ , then

$$\sqrt{m} (\mathcal{E}(S_m) - \phi(P)) \xrightarrow{d} \mathcal{N}(0, \Sigma_P),$$

which is Eqn. I.

*Proof.* By the classical empirical process CLT,  $\sqrt{m} (P_m - P) \Rightarrow \mathbb{G}_P$  in an appropriate function space (e.g.  $\ell^\infty(\mathcal{F})$  for a Donsker class  $\mathcal{F}$ ), see Van Der Vaart & Wellner (1996, Ch. 2). Hadamard differentiability of  $\phi$  at  $P$  allows application of the functional delta method Van Der Vaart & Wellner (1996, Thm. 3.9.4), yielding the stated limit. See Fishman et al. (2025) for this argument in the GDE setting.  $\square$

### A.1.3. PLUG-IN LOSSES FOR CONDITIONED MODELS

We now formalize the plug-in CLTs used to justify minibatch training in Alg. 1 and Prop. 3.1. The core plug-in loss CLT and bias expansion are proved in Fishman et al. (2025, Assump. (C.1)–(C.3) and Thm. “Large- $m$  behaviour of the plug-in loss”). Here we record a DCT-specialized statement and explicitly handle two technical points needed for conditional transport: (i) joint CLTs for pairs of embeddings when the source and target minibatches are either independent (product coupling) or paired/coupled (non-product coupling), and (ii) the fact that reusing a finite number of points from the same minibatch as the embedding does not affect the limiting distribution.

**Theorem A.5** (Plug-in loss CLTs for DCT). Fix indices  $(u, v)$ . Let  $z_u^* := \phi(P_u)$  and  $z_v^* := \phi(P_v)$ . Consider either of the following sampling schemes:

- (A) (**Independent minibatches**)  $\hat{S}_u = (X_{u1}, \dots, X_{um}) \sim P_u^{\otimes m}$  and  $\hat{S}_v = (X_{v1}, \dots, X_{vm'}) \sim P_v^{\otimes m'}$  are independent.
- (B) (**Coupled/paired minibatches**)  $m' = m$  and  $(X_{ui}, X_{vi})_{i=1}^m \stackrel{\text{iid}}{\sim} \Pi_{uv}$  for some coupling  $\Pi_{uv}$  with marginals  $P_u$  and  $P_v$ , with  $\hat{S}_u = (X_{u1}, \dots, X_{um})$  and  $\hat{S}_v = (X_{v1}, \dots, X_{vm})$ .

Define  $z_u := \mathcal{E}(\hat{S}_u)$  and  $z_v := \mathcal{E}(\hat{S}_v)$ . When  $m' = m$ , case (A) is obtained from case (B) by taking the product coupling  $\Pi_{uv} = P_u \otimes P_v$ .

Assume the encoder is asymptotically linear at  $P_u$  and  $P_v$  as in Fishman et al. (2025, Assump. (C.1)): there exist  $\psi_u \in L^2(P_u)$  and  $\psi_v \in L^2(P_v)$  with mean zero such that

$$\sqrt{m} (z_u - z_u^*) = \frac{1}{\sqrt{m}} \sum_{i=1}^m \psi_u(X_{ui}) + o_{\mathbb{P}}(1), \quad \sqrt{m'} (z_v - z_v^*) = \frac{1}{\sqrt{m'}} \sum_{i=1}^{m'} \psi_v(X_{vi}) + o_{\mathbb{P}}(1).$$

Let  $\Sigma_u := \text{Var}_{X \sim P_u} [\psi_u(X)]$  and  $\Sigma_v := \text{Var}_{X \sim P_v} [\psi_v(X)]$ . In case (B), define the cross-covariance  $\Sigma_{uv}^{12} := \text{Cov}_{(X, Y) \sim \Pi_{uv}} [\psi_u(X), \psi_v(Y)]$ .

Fix  $x \in \mathcal{X}$  and let  $\ell^{\text{sc}}(x, \cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $\ell^{\text{stc}}(x, \cdot, \cdot) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  be differentiable at  $z_u^*$  and  $(z_u^*, z_v^*)$ , respectively. Define the plug-in losses

$$\hat{\ell}_m^{\text{sc}}(x) := \ell^{\text{sc}}(x, z_u), \quad \hat{\ell}_{m,m'}^{\text{stc}}(x) := \ell^{\text{stc}}(x, z_u, z_v).$$

Then:

- (i) **Joint embedding CLT.**  $(\sqrt{m}(z_u - z_u^*), \sqrt{m'}(z_v - z_v^*)) \Rightarrow \mathcal{N}(0, \Sigma_{uv})$ . In case (A),  $\Sigma_{uv} = \text{diag}(\Sigma_u, \Sigma_v)$ , while in case (B) (with  $m' = m$ ) we have

$$\Sigma_{uv} = \begin{pmatrix} \Sigma_u & \Sigma_{uv}^{12} \\ (\Sigma_{uv}^{12})^\top & \Sigma_v \end{pmatrix}.$$

- (ii) **Source-conditioned plug-in loss.**  $\sqrt{m} (\hat{\ell}_m^{\text{sc}}(x) - \ell^{\text{sc}}(x, z_u^*)) \Rightarrow \mathcal{N}(0, \sigma_{\text{sc}}^2(x))$ .

715 (iii) **Source–target conditioned plug-in loss.** If  $m/m' \rightarrow \rho \in (0, \infty)$  then  $\sqrt{m}(\hat{\ell}_{m,m'}^{\text{stc}}(x) - \ell^{\text{stc}}(x, z_u^*, z_v^*)) \Rightarrow$   
 716  $\mathcal{N}(0, \sigma_{\text{stc},\rho}^2(x))$ .  
 717

718 Moreover, the same limits in (ii)–(iii) continue to hold if  $x$  is replaced by any fixed finite subset of points from  $\hat{S}_u$  (e.g.  
 719  $x = X_{u1}$ ): conditioning on those points changes  $\sqrt{m}(z_u - z_u^*)$  only by  $o_{\mathbb{P}}(1)$ .  
 720

721 *Proof.* (i) The asymptotic linearity expansions reduce the joint fluctuations to sums of i.i.d. terms. In case (A) the two  
 722 sums are independent, giving the block-diagonal covariance. In case (B), apply the multivariate CLT to the i.i.d. sequence  
 723  $(\psi_u(X_{ui}), \psi_v(X_{vi}))_{i=1}^m$  to obtain a joint Gaussian limit with covariance as stated.  
 724

725 (ii)–(iii) Apply the (multivariate) delta method to the map  $z \mapsto \ell^{\text{sc}}(x, z)$  in the source-conditioned case and to  $(z, z') \mapsto$   
 726  $\ell^{\text{stc}}(x, z, z')$  in the source–target-conditioned case, using the joint CLT from (i). For (iii), note that  $\sqrt{m}(z_v - z_v^*) =$   
 727  $\sqrt{m/m'} \cdot \sqrt{m'}(z_v - z_v^*) \Rightarrow \sqrt{\rho} Z_v$ .  
 728

729 (reuse-of-points dependence) The expansion in Assump. (C.1) already gives  $\sqrt{m}(z_u - z_u^*) = m^{-1/2} \sum_{i=1}^m \psi_u(X_{ui}) + o_{\mathbb{P}}(1)$ .  
 730 If we condition on  $X_{u1}$  (or any fixed finite subset), only finitely many summands are affected and they vanish at rate  $m^{-1/2}$ ,  
 731 so the conditional limit of  $\sqrt{m}(z_u - z_u^*)$  is unchanged. Applying the same delta-method argument conditional on  $X_{u1}$   
 732 yields the stated conclusion.  $\square$

733 **Corollary A.6** (Mean consistency and (asymptotic) unbiasedness of plug-in training). *Under the assumptions of Thm. A.5,  
 734 suppose additionally that for each fixed  $x$  the maps  $z \mapsto \ell^{\text{sc}}(x, z)$  and  $(z, z') \mapsto \ell^{\text{stc}}(x, z, z')$  are locally Lipschitz  
 735 near  $z_u^*$  and  $(z_u^*, z_v^*)$ , respectively, and that  $\sup_m \mathbb{E}\|\sqrt{m}(z_u - z_u^*)\|^2 < \infty$  and  $\sup_{m'} \mathbb{E}\|\sqrt{m'}(z_v - z_v^*)\|^2 < \infty$ . Then  
 736  $\mathbb{E}|\hat{\ell}_m^{\text{sc}}(x) - \ell^{\text{sc}}(x, z_u^*)| \rightarrow 0$  and  $\mathbb{E}|\hat{\ell}_{m,m'}^{\text{stc}}(x) - \ell^{\text{stc}}(x, z_u^*, z_v^*)| \rightarrow 0$ , hence the plug-in objectives are mean consistent. If,  
 737 moreover, these maps are twice continuously differentiable with bounded Hessian near the same points and  $\mathbb{E}[z_u - z_u^*] = O(m^{-1})$  and  $\mathbb{E}[z_v - z_v^*] = O(m'^{-1})$ , then the plug-in biases are  $O(m^{-1})$  (and  $O(m'^{-1})$ ).  
 738*

739 740 *Proof.* This is the same local Lipschitz / second-order Taylor expansion argument as in Fishman et al. (2025, Thm. “Large- $m$   
 741 behaviour of the plug-in loss”), applied to the scalar maps  $z \mapsto \ell^{\text{sc}}(x, z)$  and  $(z, z') \mapsto \ell^{\text{stc}}(x, z, z')$  together with the  
 742 moment bounds implied by the encoder CLTs.  $\square$

## 743 A.2. When do source samples matter?

744 The plug-in CLTs above justify minibatch training for objectives that compare *distributions*. Separately, when training  
 745 any-to-any models with an *independent* pairing policy (product coupling) and a stochastic generator that has access to  
 746 additional noise, the population objective can admit solutions that match the target distribution while effectively ignoring the  
 747 source sample. This subsection records a simple identifiability observation and a practical diagnostic.  
 748

749 **Proposition A.7** (Degenerate conditional sampler under product coupling). *Fix indices  $(u, v)$ . Let  $\xi \sim \nu$  be auxiliary noise  
 750 independent of  $X \sim P_u$ . Consider a source–target conditioned generator  $\mathcal{T}(\cdot | z_u, z_v; \xi)$  and a distribution-level loss of the  
 751 form*

$$752 \quad \mathcal{L}_{\text{stc}}(u, v) := \mathcal{L}\left(P_v, \text{Law}(\mathcal{T}(X | z_u, z_v; \xi))\right),$$

753 where  $\mathcal{L}(P, Q) = 0$  iff  $P = Q$ . If the model class contains  $x$ -ignoring generators of the form

$$754 \quad \mathcal{T}(x | z_u, z_v; \xi) = G(z_v; \xi)$$

755 and for each  $v$  there exists such a  $G$  with  $G(z_v; \xi) \sim P_v$ , then the minimum of  $\mathcal{L}_{\text{stc}}(u, v)$  equals 0 and is achieved by an  
 756  $x$ -ignoring solution.  
 757

758 *Proof.* If  $\mathcal{T}(x | z_u, z_v; \xi) = G(z_v; \xi)$  and  $G(z_v; \xi) \sim P_v$ , then the output marginal distribution equals  $P_v$  regardless of  $P_u$ .  
 759 Thus  $\mathcal{L}(P_v, \text{Law}(\mathcal{T}(X | z_u, z_v; \xi))) = 0$ .  $\square$

760 Prop. A.7 does not say a model must ignore  $x$ ; it says a product-coupling, distribution-level objective cannot generally rule  
 761 out this failure mode when the generator has sufficient stochasticity. In practice, one must break this symmetry—e.g. by  
 762 restricting to deterministic generators; using structured pairings/couplings (e.g. optimal transport) (Tong et al., 2023b); or  
 763 adding regularization that selects a coupling or geometry (cf. (Haviv et al., 2024)).  
 764

Generator	$d_{\text{pair}}$	$d_{\text{rand}}$	$d_{\text{pair}}/d_{\text{rand}}$	Rank $\rho$
Flow matching (FM)	0.011	0.629	0.017	1.000
Energy regression (Energy)	0.543	0.615	0.885	0.106
Sliced Wasserstein (SWD)	0.555	0.623	0.893	0.108
Stochastic energy sampler	0.647	0.649	0.998	0.009

Table 5. Source-sample alignment diagnostic on the MVN any-to-any benchmark (App. E,  $K=1000$  training distributions; averages over 20 random  $(u, v)$  pairs with  $N=200$  samples each;  $d_{\text{rand}}$  averaged over 50 random permutations). FM (which uses a sample-level paired objective) yields strong alignment ( $d_{\text{pair}}/d_{\text{rand}} \approx 0.02$ ,  $\rho \approx 1$ ). Energy and SWD are deterministic maps at inference but are trained with set-level distributional losses; in this experiment they show moderate alignment ( $d_{\text{pair}}/d_{\text{rand}} \approx 0.9$ ,  $\rho \approx 0.1$ ). The stochastic energy sampler injects noise at inference and achieves  $d_{\text{pair}} \approx d_{\text{rand}}$ , consistent with the  $x$ -ignoring failure mode in Prop. A.7 and Lem. A.8.

**An alignment diagnostic.** Fix a cost  $c : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$  (we use  $c(x, y) = \|x - y\|_2$ ). For a fixed pair  $(u, v)$ , let  $X \sim P_u$  and  $Y \sim P_v$  be independent and let  $\hat{Y} := \mathcal{T}(X \mid z_u, z_v; \xi)$  for independent noise  $\xi \sim \nu$ . Let  $\mu_u := \mathbb{E}[X]$  and  $\mu_v := \mathbb{E}[Y]$  and define the centered cost

$$\bar{c}_{uv}(x, y) := c(x - \mu_u, y - \mu_v).$$

We report

$$d_{\text{pair}} := \mathbb{E}[\bar{c}_{uv}(X, \hat{Y})], \quad d_{\text{rand}} := \mathbb{E}[\bar{c}_{uv}(X, Y)], \quad \text{and the ratio } d_{\text{pair}}/d_{\text{rand}}.$$

Centering removes the bulk mean-shift component (which can otherwise dominate  $c$ ) and makes the diagnostic more sensitive to the coupling.

**Lemma A.8** (Independence baseline). *If  $\hat{Y}$  is independent of  $X$  and  $\text{Law}(\hat{Y}) = P_v$ , then  $d_{\text{pair}} = d_{\text{rand}}$ .*

*Proof.* Under the assumptions,  $(X, \hat{Y}) \stackrel{d}{=} (X, Y)$  with  $Y \sim P_v$  independent of  $X$ , so the expectations agree.  $\square$

As a complementary diagnostic, let  $\Delta_{uv} := \mu_v - \mu_u$  and compute the Spearman rank correlation  $\rho$  between the one-dimensional projections  $\langle X, \Delta_{uv} \rangle$  and  $\langle \hat{Y}, \Delta_{uv} \rangle$ . High  $|\rho|$  indicates that the map preserves rank structure along the mean-shift direction, whereas  $\rho \approx 0$  is consistent with weak pointwise dependence. In practice we estimate the expectations and means with sample averages. For the *stochastic energy sampler* baseline in Table 5, we augment the deterministic Energy regression map with per-sample noise (so  $\hat{Y} = \mathcal{T}(X \mid z_u, z_v; \xi)$  with  $\xi \sim \nu$ ), making the model class expressive enough to realize  $x$ -ignoring solutions.

Taken together, these results illustrate when we can expect DCT to learn more than an independent conditional sampler: objectives with an explicit sample-level pairing signal (or additional structure selecting a coupling) can induce meaningful pointwise dependence, while purely distributional objectives need not.

### A.3. How is the latent space structured?

Following the manifold/trajectory viewpoint of Fishman et al. (2025), we use the learned embedding space as a coordinate system for a family of related distributions and probe it via *latent interpolants*. Fix a pair  $(u, v)$  with embeddings  $(z_u, z_v)$ . We linearly interpolate between them,  $z(t) = (1-t)z_u + tz_v$ , choose a grid  $0 = t_0 < \dots < t_K = 1$ , and set  $z_k = z(t_k)$ . Starting from samples  $x_0 \sim P_u$ , we form a discrete trajectory by repeatedly applying the source–target conditioned map between successive latent codes, i.e.  $x_{k+1} = \mathcal{T}(x_k \mid z_k, z_{k+1})$ . The empirical laws of  $(x_k)_{k=0}^K$  define a path in distribution space.

Fig. 6 shows the resulting STC trajectory on the MVN any-to-any benchmark (App. E) alongside the closed-form optimal transport (OT) displacement interpolation between the same endpoint Gaussians. In this example, the learned trajectory closely tracks the OT path, suggesting that linear structure in the latent space induces a meaningful geometry in distribution space (at least within this controlled family). This is not a guarantee, it merely shows that the latent spaces learned by the distribution encoders for the any-to-any transport task inherit many of the interesting and useful properties of GDEs.

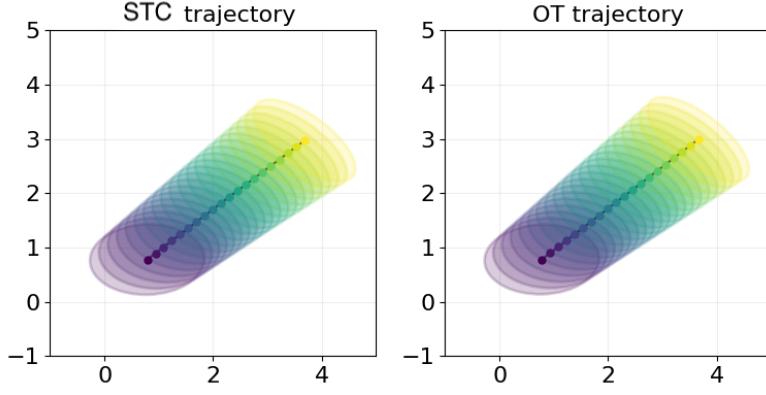


Figure 6. Latent-space interpolants on the MVN any-to-any benchmark. **Left:** STC trajectory obtained by sequentially pushing samples using intermediate latent codes along a linear interpolant between endpoint embeddings. **Right:** OT displacement interpolation (computed in closed form for Gaussians).

## B. Basline Comparisons

### B.1. Unsupervised Transport

In our unsupervised transport experiments, we compare source-target-condition models in the DCT framework to  $K$ -to- $K$  baselines, which condition on one-hot embeddings of distributions. This approach is closely related to that of MMSI (Albergo et al., 2023a), but distinct in a subtle way.

MMSI models interpolant time on the simplex, where each point in the simplex corresponds to an interpolant taking a convex combination of a sample from each of the  $K$  distributions proportional to their simplicial weight. Our  $K$ -to- $K$  baselines (when using flow matching generators) can be viewed as MMSI models which take “edge paths” on the simplex, where we flow directly between two corners without any weight on other distributions.

We choose our one-hot  $K$ -to- $K$  as a baseline, because it allows for direct comparison across a broader family of transport mechanisms rather than flow matching alone, and anecdotally we have observed more stable performance with MMSI models restricted to simplicial edge paths.

### B.2. Supervised Transport

In our supervised transport experiments, we evaluate source-conditioned models. These models are closely related to Meta Flow Matching models (Atanackovic et al., 2024). When using flow matching, source-conditioned DCTs differ from MFM only in the choice of distribution encoder. Rather than a distributionally invariant encoder, MFM instead uses a heuristic “population encoder” which involves message passing over a k-nearest-neighbor graph on the input population.

We choose to baseline against source-conditioned models with distributionally invariant encoders as these provide a CLT and robustness to variation in sample set size. Furthermore, using source-conditioned models rather than MFM alone allows us to again directly compare between source-conditioned and source-target-conditioned models across a broad range of transport mechanisms.

## C. Standard distribution-conditioned transport maps

Several experiments in this paper reuse the same generator families and only vary the dataset, encoder, or model capacity. This appendix defines the three *continuous-data* generators used throughout: sliced Wasserstein regression (**SWD**) (Kolouri et al., 2019), energy/MMD regression (**Energy**) (Gneiting & Raftery, 2007; Sutherland et al., 2016), and flow matching (**FM**) (Liu et al., 2022; Lipman et al., 2023; Albergo et al., 2023b). Discrete sequence generators for the TCR experiments are described separately in Appendix J.3. Where an experiment uses different hyperparameters (e.g. number of projections or solver tolerances), it is stated explicitly in that experiment’s appendix.

Throughout, **DCT** refers to the overall *distribution-conditioned transport* framework of coupling a distribution encoder with a conditional generator. Within DCT we use **source-conditioned (SC)** to mean generators conditioned only on the source embedding  $z_{\text{src}}$ , and **source–target-conditioned (STC)** to mean generators conditioned on both  $(z_{\text{src}}, z_{\text{tgt}})$ ; STC models are also referred to as *any-to-any* models in the main text. Unless otherwise stated, the generator families below apply to both SC and STC, with SC obtained by dropping  $z_{\text{tgt}}$ .

## C.1. Notation

Let  $S^{\text{src}} = \{x_j\}_{j=1}^n \sim P^{\text{src}}$  and  $S^{\text{tgt}} = \{y_j\}_{j=1}^n \sim P^{\text{tgt}}$  denote source and target sample sets (not necessarily paired). Let  $z_{\text{src}} = \mathcal{E}(S^{\text{src}})$  and  $z_{\text{tgt}} = \mathcal{E}(S^{\text{tgt}})$  be the corresponding distribution embeddings. All generators produce a transported sample set  $\hat{S}^{\text{tgt}} = \{\hat{y}_j\}_{j=1}^n$ ; in the STC setting they are conditioned on  $(z_{\text{src}}, z_{\text{tgt}})$ , while in the SC setting they are conditioned on  $z_{\text{src}}$  alone.

## C.2. Standard distribution encoder

Across continuous-data experiments, we use a single *standard* distribution encoder  $\mathcal{E}$  to map a finite sample set to a fixed-dimensional embedding. Concretely, given a set  $S = \{x_j\}_{j=1}^n \subset \mathbb{R}^d$ , the encoder outputs  $z = \mathcal{E}(S) \in \mathbb{R}^{d_z}$ . We adopt the mean-pooled DeepSets/GNN-style architecture used in Generative Distribution Embeddings (GDE) (Fishman et al., 2025), which is consistently strong across their benchmark suite.

**Architecture.** We first compute per-sample features with an MLP,

$$h_j^{(0)} = \text{MLP}_{\text{in}}(x_j) \in \mathbb{R}^{d_h}. \quad (6)$$

We then apply  $L$  mean-pooled update blocks. For  $\ell = 1, \dots, L$ , define the pooled representation  $\bar{h}^{(\ell-1)} = \frac{1}{n} \sum_{j=1}^n h_j^{(\ell-1)}$  and update each sample feature by concatenating with the pooled representation:

$$h_j^{(\ell)} = \text{MLP}_{\text{pool}}^{(\ell)} \left( [h_j^{(\ell-1)}; \bar{h}^{(\ell-1)}] \right). \quad (7)$$

Finally, we average over samples and project to the latent space:

$$z = \sigma \left( W \cdot \frac{1}{n} \sum_{j=1}^n h_j^{(L)} + b \right), \quad (8)$$

with  $\sigma = \text{SELU}$ . Optionally, we  $\ell_2$ -normalize  $z$ ; this is used only when stated explicitly.

**Hyperparameters.** Unless stated otherwise in an experiment appendix, we use  $L = 2$  mean-pooled blocks and 2-layer MLPs with SELU activations. The hidden dimension  $d_h$  and embedding dimension  $d_z$  are reported per experiment.

**Non-continuous data.** For discrete sequence experiments (TCR forecasting), we use an ESM2-based distribution encoder described in Appendix J.2.

## C.3. Deterministic regression generators (SWD and Energy)

Both SWD and Energy use a deterministic transport map applied pointwise:

$$\hat{y}_j = \mathcal{T}_{\theta}(x_j \mid z_{\text{src}}, z_{\text{tgt}}). \quad (9)$$

They differ only in the training objective, which compares  $\hat{S}^{\text{tgt}}$  to the observed target set  $S^{\text{tgt}}$  using a distributional distance (Appendix D).

**SWD (sliced Wasserstein regression).** We minimize the empirical sliced 2-Wasserstein distance using  $L = 100$  random projections  $\{\omega_{\ell}\}_{\ell=1}^L$  drawn uniformly from the unit sphere. For each  $\omega_{\ell}$ , we compute scalar projections and sort them,  $\text{sort}(\{\langle \omega_{\ell}, \hat{y}_j \rangle\}_{j=1}^n)$  and  $\text{sort}(\{\langle \omega_{\ell}, y_j \rangle\}_{j=1}^n)$ , then take the mean squared difference across matched order statistics and projections.

**Energy (energy/MMD regression).** We minimize the energy distance, equivalently MMD with the energy kernel  $k(x, y) = -\|x - y\|$  (Appendix D). The energy kernel is parameter-free, avoiding bandwidth selection.

935 **C.4. Flow matching generator (FM)**

936 Flow matching parameterizes transport with a time-conditional velocity field  $v_\theta(x, t, z_{\text{src}}, z_{\text{tgt}})$ . For training we sample  
 937  $t \sim \text{Uniform}[0, 1]$ ,  $x_0 \sim P^{\text{src}}$ ,  $x_1 \sim P^{\text{tgt}}$ , and  $\epsilon \sim \mathcal{N}(0, I)$ , and form the noisy linear interpolation  
 938

$$x_t = (1 - t)x_0 + tx_1 + \sigma\epsilon, \quad (10)$$

940 with  $\sigma = 0.5$  unless otherwise stated. We train with the conditional flow matching objective  
 941

$$\mathcal{L}_{\text{FM}} = \mathbb{E} \left[ \|v_\theta(x_t, t, z_{\text{src}}, z_{\text{tgt}}) - (x_1 - x_0)\|_2^2 \right]. \quad (11)$$

942 At inference, for each source sample  $x(0) = x_0$  we solve the ODE  
 943

$$\frac{dx}{dt} = v_\theta(x, t, z_{\text{src}}, z_{\text{tgt}}), \quad t \in [0, 1], \quad (12)$$

944 and take  $\hat{y} = x(1)$ . We use an adaptive Dormand–Prince solver (dopri5) with absolute and relative tolerances  $10^{-4}$  unless  
 945 otherwise stated.  
 946

947 **D. Standard evaluation metrics**

948 Across experiments we report three distributional distances between a generated sample set  $\hat{S} = \{\hat{y}_i\}_{i=1}^n$  and a ground-truth  
 949 target set  $S = \{y_j\}_{j=1}^n$ : sliced Wasserstein distance (**SWD**), energy distance (**Energy**), and MMD with an RBF kernel  
 950 (**MMD-RBF**). These metrics are computed in the data space unless stated otherwise (e.g. in the TCR experiments, they are  
 951 computed in ESM2 embedding space).  
 952

953 **D.1. Energy distance (MMD with energy kernel)**

954 The energy distance between distributions  $P$  and  $Q$  can be written as  
 955

$$D_{\text{energy}}(P, Q) = 2\mathbb{E}\|X - Y\| - \mathbb{E}\|X - X'\| - \mathbb{E}\|Y - Y'\|, \quad (13)$$

956 where  $X, X' \sim P$  i.i.d. and  $Y, Y' \sim Q$  i.i.d. We estimate this quantity empirically from the two finite sample sets. This is  
 957 equivalent to MMD with kernel  $k(x, y) = -\|x - y\|$ . The energy distance is closely related to the *energy score*, a strictly  
 958 proper scoring rule for probabilistic prediction (Gneiting & Raftery, 2007).  
 959

960 **D.2. Sliced 2-Wasserstein distance (SWD)**

961 We use the standard sliced 2-Wasserstein distance (Kolouri et al., 2019) estimated with  $L = 100$  random projections  
 962  $\{\omega_\ell\}_{\ell=1}^L$  drawn uniformly from the unit sphere. For each projection, we compute and sort the scalar projections of both sets,  
 963 match order statistics, and average the squared differences across projections:  
 964

$$\widehat{\text{SWD}}_2^2(\hat{S}, S) = \frac{1}{L} \sum_{\ell=1}^L \frac{1}{n} \sum_{j=1}^n (\text{sort}(\{\langle \omega_\ell, \hat{y}_i \rangle\}_{i=1}^n)_j - \text{sort}(\{\langle \omega_\ell, y_i \rangle\}_{i=1}^n)_j)^2. \quad (14)$$

965 We report  $\widehat{\text{SWD}}_2(\hat{S}, S) = \sqrt{\widehat{\text{SWD}}_2^2(\hat{S}, S)}$ .  
 966

967 **D.3. MMD with RBF kernel (MMD-RBF)**

968 Given a positive definite kernel  $k$ , the squared MMD is  
 969

$$\text{MMD}^2(P, Q) = \mathbb{E}[k(X, X')] + \mathbb{E}[k(Y, Y')] - 2\mathbb{E}[k(X, Y)], \quad (15)$$

970 with  $X, X' \sim P$  and  $Y, Y' \sim Q$  i.i.d., estimated empirically from samples. MMD is a standard kernel two-sample statistic  
 971 (Gretton et al., 2012). We use an RBF kernel  
 972

$$k_{\text{RBF}}(x, y) = \exp\left(-\frac{\|x-y\|_2^2}{2\sigma^2}\right), \quad (16)$$

973 with bandwidth  $\sigma$  set by the median heuristic on the pooled samples from  $\hat{S} \cup S$ .  
 974

---

## 990 E. Gaussian and Gaussian Mixture Experiments

991 This appendix provides experimental details for the Gaussian experiments in Sections 5.1 and 6.1. We describe the data-  
 992 generating processes, model architectures, training configurations, and evaluation protocols. Generator families and metrics  
 993 are defined in Appendices C and D.

### 995 E.1. Data-Generating Process

#### 997 E.1.1. MULTIVARIATE NORMAL (MVN) DISTRIBUTIONS

998 We construct a dataset of  $n = 50,000$  sample sets, each drawn from a distinct bivariate normal distribution. The parameters  
 1000 for each distribution  $i$  are sampled as follows:

1002 **Mean Prior.** Each component of the mean vector is drawn independently from a uniform distribution:

$$1003 \quad \mu_i \sim \text{Uniform}([0, 5]^2) \quad (17)$$

1005 **Covariance Prior.** The covariance matrix is sampled from an inverse-Wishart distribution:

$$1006 \quad \Sigma_i \sim \text{InverseWishart}(\nu, \Psi) \quad (18)$$

1008 where the degrees of freedom  $\nu = 10 \cdot (d - 1) = 10$  for  $d = 2$  dimensions, and the scale matrix  $\Psi = I_d$  (identity matrix).  
 1009 This parameterization ensures well-conditioned covariance matrices with moderate variability.

1011 **Sample Generation.** For each distribution  $i$ , we draw  $n_{\text{set}} = 100$  i.i.d. samples:

$$1012 \quad x_{ij} \sim \mathcal{N}(\mu_i, \Sigma_i), \quad j = 1, \dots, n_{\text{set}} \quad (19)$$

1013 Sampling is performed via the Cholesky decomposition  $\Sigma_i = L_i L_i^\top$ :

$$1015 \quad x_{ij} = \mu_i + L_i z_{ij}, \quad z_{ij} \sim \mathcal{N}(0, I_d) \quad (20)$$

1017 **Unique Distribution Subsets.** For experiments studying generalization across distribution density (varying  $K$ ), we sample  
 1018  $K \in \{10, 100, 1,000, 10,000\}$  unique parameter sets  $(\mu_k, \Sigma_k)$  and repeat each parameter set  $n/K$  times. This creates  
 1019 multiple sample sets from the same underlying distribution, enabling the study of how models perform when training  
 1020 distributions are seen multiple times versus once.

#### 1022 E.1.2. GAUSSIAN MIXTURE MODELS (GMM)

1023 We extend the MVN setting to Gaussian mixture models with  $C = 3$  components. The parameters for each GMM  $i$  are  
 1024 sampled as follows:

1027 **Mixture Weight Prior.** Component weights are drawn from a symmetric Dirichlet distribution:

$$1028 \quad \pi_i \sim \text{Dirichlet}(\alpha \cdot \mathbf{1}_C) \quad (21)$$

1029 where  $\alpha = 1.0$  yields a uniform prior over the simplex.

1031 **Component Mean Prior.** Each component mean is drawn independently:

$$1033 \quad \mu_{ic} \sim \text{Uniform}([0, 5]^2), \quad c = 1, \dots, C \quad (22)$$

1035 **Component Covariance Prior.** Each component covariance is drawn independently:

$$1036 \quad \Sigma_{ic} \sim \text{InverseWishart}(\nu, \Psi), \quad c = 1, \dots, C \quad (23)$$

1037 with the same hyperparameters as the MVN case ( $\nu = 10$ ,  $\Psi = I_d$ ).

1039 **Sample Generation.** For each GMM  $i$ , we draw  $n_{\text{set}} = 1,000$  samples via ancestral sampling:

$$1041 \quad c_{ij} \sim \text{Categorical}(\pi_i) \quad (24)$$

$$1042 \quad x_{ij} \sim \mathcal{N}(\mu_{i,c_{ij}}, \Sigma_{i,c_{ij}}) \quad (25)$$

1043 The larger set size (1,000 vs. 100 for MVN) provides sufficient samples to capture the multimodal structure.  
 1044

## 1045 E.2. Supervised Gaussian Experiments

1046 For the semi-supervised experiments, we construct paired source-target datasets with known transformations.

## 1047 E.2.1. SUPERVISED MVN DATASET

1048 The supervised MVN dataset creates paired distributions via a fixed shift transformation:

1049 
$$Y = X + b, \quad b = (1, 1)^\top \quad (26)$$

 1050 **Source Distribution.** To study the generalization beyond the training distribution and evaluate the value of any-to-any  
 1051 models for generalization, source means are restricted:

1052 
$$\mu_i^{\text{src}} \sim \text{Uniform}([0, 2.5]^2) \quad (27)$$

 1053 **Target Generation.** Given source samples  $\{x_{ij}\}$ , target samples are:

1054 
$$y_{ij} = x_{ij} + b \quad (28)$$

1055 The target samples are then randomly permuted to remove point-wise correspondence, creating an unpaired but distributionally related dataset.

## 1056 E.2.2. SUPERVISED GMM DATASET

1057 The supervised GMM dataset creates a more challenging bimodal target structure:

1058 
$$Y = X + b \pm b_{\text{off}} \quad (29)$$

 1059 where  $b = (1, 1)^\top$  and  $b_{\text{off}} = 0.1 \cdot (-1, 1)^\top$ .

 1060 **Target Generation.** For each source sample, we create two shifted versions:

1061 
$$y_{ij}^+ = x_{ij} + b + b_{\text{off}} \quad (30)$$

1062 
$$y_{ij}^- = x_{ij} + b - b_{\text{off}} \quad (31)$$

 1063 These are concatenated and subsampled to maintain the original set size, then randomly permuted. This creates a bimodal  
 1064 target distribution from a unimodal (per-component) source, testing the model's ability to learn complex distributional  
 1065 transformations.

## 1066 E.3. Model Architecture

## 1067 E.3.1. DISTRIBUTION ENCODER

 1068 We use the standard distribution encoder described in Appendix C.2 (mean-pooled DeepSets/GNN encoder). For MVN  
 1069 experiments, we use hidden dimension  $d_h = 64$  and latent dimension  $d_z = 16$ ; for GMM, we increase these to  $d_h = 256$   
 1070 and  $d_z = 128$  to accommodate the more complex multimodal structure.

## 1071 E.3.2. EMBEDDING ENCODER (BASELINE)

 1072 For the  $K$ -to- $K$  baseline that treats each distribution as a discrete label, we use a learnable embedding table:

1073 
$$z_i = \text{Embedding}(i) \in \mathbb{R}^{d_z} \quad (32)$$

 1074 where  $i \in \{1, \dots, K\}$  is the distribution index. This encoder cannot generalize to unseen distributions at test time.

## 1075 E.3.3. TRANSPORT MAP

 1076 The transport map  $\mathcal{T}(x | z_{\text{src}}, z_{\text{tgt}})$  is a 4-layer MLP with SELU activations that takes the concatenation of a source sample  
 1077  $x$  with both the source and target distribution embeddings:  $\hat{y} = \text{MLP}([x; z_{\text{src}}; z_{\text{tgt}}])$ . The hidden dimension matches the  
 1078 encoder ( $d_h = 64$  for MVN,  $d_h = 256$  for GMM).

 1079 We compare the three standard continuous-data generator families defined in Appendix C: sliced Wasserstein regression  
 1080 (SWD), energy/MMD regression (Energy), and flow matching (FM).

1100 **E.4. Training Protocol**

1101 All models are trained using Adam with a fixed learning rate of  $2 \times 10^{-4}$  for 200 epochs. We use a batch size of 256  
 1102 distribution pairs, where each pair consists of source and target sample sets. No learning rate scheduling or early stopping is  
 1103 applied.  
 1104

1105  
 1106 **Pair Sampling.** During training, we sample source-target pairs uniformly at random. For the unsupervised any-to-any  
 1107 setting, source and target distributions are sampled independently from the full dataset, so the model sees arbitrary pairings.  
 1108 For the supervised setting, each target is deterministically derived from its corresponding source via the fixed transformation  
 1109 (shift for MVN, shift plus off-axis displacement for GMM).  
 1110

1111  
 1112 **Bidirectional vs. Unidirectional Loss.** The unsupervised any-to-any objective trains transport in both directions, summing  
 1113 the losses:  
 1114

$$\mathcal{L}_{\text{unsup}} = \mathcal{L}(\text{src} \rightarrow \text{tgt}) + \mathcal{L}(\text{tgt} \rightarrow \text{src}) \quad (33)$$

1115 This symmetric objective encourages the encoder to learn representations that support transport regardless of which  
 1116 distribution serves as source or target. For supervised training, we use only the forward direction  $\mathcal{L}_{\text{sup}} = \mathcal{L}(\text{src} \rightarrow \text{tgt})$ , since  
 1117 the transformation is inherently asymmetric.  
 1118

1119 **E.5. Evaluation Protocol**

## 1120 E.5.1. UNSUPERVISED EVALUATION

1121 We evaluate transport quality in two settings. For *in-distribution* evaluation, we test on held-out source-target pairs where  
 1122 both distributions appeared during training, though not necessarily paired together. For *out-of-distribution* evaluation, we fix  
 1123 a source distribution from the training set and transport to a dense grid of novel target distributions with means  $\mu \in [0, 5]^2$   
 1124 and randomly sampled covariances. This tests zero-shot generalization to target distributions never seen during training.  
 1125

1126 We measure transport quality using the energy distance (Appendix D), estimated empirically from the finite sample sets.  
 1127

## 1128 E.5.2. SEMI-SUPERVISED EVALUATION

1129 For the semi-supervised experiments, we restrict the supervised training data to distributions with source means satisfying  
 1130  $\|\mu\|_\infty \leq 2.5$ , then evaluate across the full range  $\|\mu\|_\infty \leq 5$ . We partition test results by the  $L^\infty$  norm of the source mean to  
 1131 assess how performance degrades outside the supervised training support.  
 1132

1133 The semi-supervised approach first trains an any-to-any encoder on the full unsupervised dataset (all 50,000 distributions),  
 1134 then fits a ridge regression predictor  $\hat{z}_{\text{tgt}} = W z_{\text{src}} + b$  on the paired subset within the restricted support. At test time, we  
 1135 encode the source distribution, predict the target embedding via the linear map, and apply the learned transport.  
 1136

1137 We compare against two baselines: a *supervised* model that conditions only on the source embedding and is trained  
 1138 exclusively on paired data within the support, and an *oracle* that uses the true target embedding computed by the encoder  
 1139 (providing an upper bound on semi-supervised performance).  
 1140

1141 **E.6. Hyperparameter Summary**

1142 We summarize the key hyperparameters for reproducibility. All experiments use the same random seed protocol and data  
 1143 splits.  
 1144

## 1145 E.6.1. DATASET AND MODEL HYPERPARAMETERS

1146 Table 6 provides a consolidated view of the dataset and model configurations for both experimental settings.  
 1147

1148 For the supervised experiments, we restrict the mean prior to  $[0, 2.5]^2$  so that the shifted targets remain within the original  
 1149  $[0, 5]^2$  range. The shift vector is  $b = (1, 1)^\top$  for both MVN and GMM; GMM additionally uses an off-axis displacement of  
 1150 magnitude 0.1 to create bimodal targets.  
 1151

1155  
 1156 *Table 6.* Dataset and model hyperparameters for Gaussian experiments.  
 1157

Hyperparameter	MVN	GMM
<i>Dataset</i>		
Number of sample sets	50,000	50,000
Samples per set	100	1,000
Data dimension	2	2
Mean prior range	[0, 5] <sup>2</sup>	[0, 5] <sup>2</sup>
Covariance prior (inv-Wishart $\nu$ )	10	10
Mixture components	—	3
Dirichlet concentration $\alpha$	—	1.0
<i>Encoder</i>		
Latent dimension $d_z$	16	128
Hidden dimension $d_h$	64	256
Mean-pooled GNN layers	2	2
FC layers per block	2	2
<i>Transport Map</i>		
MLP layers	4	4
Hidden dimension	64	256
Activation	SELU	SELU
<i>Training</i>		
Batch size	256	256
Learning rate	$2 \times 10^{-4}$	$2 \times 10^{-4}$
Epochs	200	200
Optimizer	Adam	Adam

## E.7. Full Results

We present comprehensive results from our Gaussian experiments. All tables report energy distance (lower is better), averaged across three random seeds.

### E.7.1. UNSUPERVISED TRANSPORT: SCALING WITH $K$

Tables 7 and 8 report the full breakdown of transport quality as a function of  $K$ , the number of unique training distributions. The key observation is that the  $K$ -to- $K$  baseline performs well on in-distribution targets when  $K$  is small (it can memorize each distribution), but degrades significantly on out-of-distribution targets. In contrast, the any-to-any encoder maintains consistent performance across both IID and OOD targets, demonstrating genuine distributional generalization.

### E.7.2. SEMI-SUPERVISED GENERALIZATION

Tables 9 and 10 compare supervised, semi-supervised, and oracle methods on the paired Gaussian tasks. The supervised model is trained only on paired data within the restricted support ( $\|\mu\|_\infty \leq 2.5$ ), while the semi-supervised approach uses the any-to-any encoder with a ridge predictor for the target embedding. The oracle uses the true target embedding.

The results show that the semi-supervised approach matches or exceeds supervised performance within the training support (IID) while maintaining substantially lower error outside it (OOD). This demonstrates that the any-to-any encoder learns distributional structure from the unsupervised objective that enables extrapolation even when the latent predictor must generalize beyond its training data.

**Flow Matching Anomaly.** We observe that supervised flow matching generalizes surprisingly well outside its training support in this synthetic setting (low OOD error even for the supervised method). This behavior is not replicated in real-world experiments, suggesting it may be an artifact of the simple linear transformation and low-dimensional setting where the learned flow field happens to extrapolate correctly.

1210  
 1211 *Table 7.* Distributional distances ( $\downarrow$ ) for MVN experiments across  $K$  unique distributions. We report MMD, SWD, and Energy distance.  
 1212

Model	$K = 10$			$K = 100$			$K = 1,000$			$K = 10,000$		
	MMD	SWD	Energy	MMD	SWD	Energy	MMD	SWD	Energy	MMD	SWD	Energy
SWD ( $K$ )	0.003	0.044	0.002	0.002	0.037	0.001	0.011	0.083	0.010	0.016	0.099	0.014
SWD (Any)	0.045	0.146	0.044	0.003	0.044	0.002	0.003	0.038	0.002	0.001	0.031	0.000
Energy ( $K$ )	0.002	0.036	0.001	0.004	0.052	0.002	0.009	0.092	0.008	0.016	0.101	0.015
Energy (Any)	0.077	0.173	0.076	0.003	0.047	0.002	0.001	0.031	-0.000	0.001	0.033	-0.000
FM ( $K$ )	0.008	0.072	0.006	0.009	0.077	0.008	0.012	0.092	0.011	0.021	0.125	0.020
FM (Any)	0.034	0.120	0.032	0.014	0.098	0.012	0.012	0.094	0.011	0.013	0.094	0.012

OOD Targets												
Model	$K = 10$			$K = 100$			$K = 1,000$			$K = 10,000$		
	MMD	SWD	Energy	MMD	SWD	Energy	MMD	SWD	Energy	MMD	SWD	Energy
SWD ( $K$ )	0.911	0.642	0.910	0.147	0.242	0.145	0.037	0.138	0.036	0.020	0.110	0.019
SWD (Any)	0.147	0.219	0.146	0.004	0.046	0.003	0.003	0.039	0.001	0.001	0.031	0.000
Energy ( $K$ )	0.942	0.651	0.941	0.127	0.227	0.126	0.036	0.143	0.035	0.020	0.111	0.019
Energy (Any)	0.171	0.237	0.169	0.004	0.048	0.002	0.001	0.035	0.000	0.001	0.033	-0.000
FM ( $K$ )	0.930	0.649	0.929	0.123	0.221	0.121	0.034	0.134	0.032	0.027	0.134	0.025
FM (Any)	0.071	0.168	0.070	0.015	0.102	0.014	0.013	0.096	0.012	0.012	0.090	0.011

## F. MNIST-Colors image benchmark

This appendix describes the MNIST-Colors image experiment, which serves as a sanity check that the qualitative behavior observed in the MVN/GMM scaling experiments (Appendix E.7) persists on simple image data. Generator families and standard distributional metrics are defined in Appendices C and D.

### F.1. Dataset

We construct a dataset of image distributions by taking grayscale MNIST digits (LeCun et al., 2010) and applying an RGB “style” via elementwise multiplication with a color vector  $c \in [0, 1]^3$ . Each distribution is represented by a set of  $n = 64$  colored images of shape  $3 \times 28 \times 28$ . We study scaling with  $K$  by restricting training to  $K \in \{10, 100, 1,000, 10,000\}$  unique colors (with repetition to keep the total number of training sets fixed). We evaluate both in-distribution (IID) targets that use colors from the training set and out-of-distribution (OOD) targets that use held-out colors.

### F.2. Models

**Encoders.** For any-to-any (STC) models, we use a convolutional distribution encoder that processes each image with a shared CNN and aggregates across the set via mean-pooled message passing, outputting  $z \in \mathbb{R}^{64}$ . For the  $K$ -to- $K$  baseline, we use a learned embedding lookup table with  $K$  embeddings.

**Generators.** All MNIST-Colors generators use a context-conditioned U-Net backbone. For the deterministic SWD and Energy regressors (Appendix C), the U-Net is conditioned on the concatenated embeddings  $(z_{\text{src}}, z_{\text{tgt}})$  and predicts transported images  $\hat{y} = \mathcal{T}_\theta(x | z_{\text{src}}, z_{\text{tgt}})$ . For flow matching (FM), we use the same backbone as a time-conditional velocity field  $v_\theta(x_t, t, z_{\text{src}}, z_{\text{tgt}})$  with  $\sigma = 0.5$ .

**Training.** Unless stated otherwise, we use latent dimension 64, set size 64, batch size 32, learning rate  $2 \times 10^{-4}$ , and train for 200 epochs.

### F.3. Metrics and results

We report SWD and MMD-RBF (Appendix D) computed on flattened pixels, as well as a color error metric defined as the mean squared error between the per-image mean RGB vectors of generated and target samples (**Color MSE**). Full results are shown in Table 11.

1265  
1266 **Table 8.** Distributional distances ( $\downarrow$ ) for GMM experiments across  $K$  unique distributions. We report MMD, SWD, and Energy distance.  
1267

Model	$K = 10$			$K = 100$			$K = 1,000$			$K = 10,000$		
	MMD	SWD	Energy	MMD	SWD	Energy	MMD	SWD	Energy	MMD	SWD	Energy
SWD ( $K$ )	0.004	0.115	0.001	0.006	0.127	0.003	0.011	0.148	0.008	0.020	0.186	0.017
SWD (Any)	0.948	0.866	0.945	0.056	0.271	0.053	0.019	0.163	0.016	0.009	0.127	0.006
Energy ( $K$ )	0.004	0.114	0.000	0.005	0.131	0.002	0.009	0.154	0.005	0.020	0.208	0.017
Energy (Any)	0.937	0.947	0.934	0.043	0.276	0.040	0.010	0.159	0.007	0.009	0.149	0.005
FM ( $K$ )	0.045	0.302	0.042	0.054	0.352	0.051	0.062	0.366	0.059	0.055	0.339	0.052
FM (Any)	0.284	0.564	0.280	0.060	0.308	0.057	0.044	0.280	0.041	0.037	0.269	0.034

#### OOD Targets

Model	$K = 10$			$K = 100$			$K = 1,000$			$K = 10,000$		
	MMD	SWD	Energy	MMD	SWD	Energy	MMD	SWD	Energy	MMD	SWD	Energy
SWD ( $K$ )	0.787	0.882	0.784	0.337	0.567	0.334	0.169	0.408	0.166	0.113	0.341	0.110
SWD (Any)	0.764	0.851	0.760	0.058	0.271	0.054	0.015	0.149	0.011	0.010	0.129	0.006
Energy ( $K$ )	0.786	0.884	0.783	0.335	0.564	0.332	0.171	0.409	0.168	0.109	0.345	0.106
Energy (Any)	1.395	1.109	1.392	0.041	0.269	0.038	0.010	0.154	0.007	0.008	0.142	0.005
FM ( $K$ )	0.782	0.858	0.779	0.350	0.592	0.348	0.187	0.465	0.184	0.124	0.401	0.121
FM (Any)	0.267	0.607	0.264	0.059	0.304	0.056	0.046	0.284	0.043	0.039	0.278	0.036

1288 **Table 9.** Distributional distances ( $\downarrow$ ) for MVN semi-supervised experiments. IID:  $\|\mu\|_\infty \leq 2.5$ , OOD:  $\|\mu\|_\infty > 2.5$ .

Generator	Method	IID			OOD		
		MMD	SWD	Energy	MMD	SWD	Energy
Flow Matching	Supervised	0.000	0.003	-0.000	0.000	0.006	-0.000
	Semi-sup.	0.001	0.019	0.001	0.001	0.021	0.001
	Oracle	0.001	0.020	0.001	0.003	0.028	0.003
Energy	Supervised	0.000	0.003	-0.000	0.020	0.062	0.020
	Semi-sup.	0.001	0.021	0.000	0.003	0.033	0.002
	Oracle	0.001	0.021	0.000	0.002	0.030	0.002
SWD	Supervised	0.000	0.003	-0.000	0.017	0.052	0.017
	Semi-sup.	0.001	0.024	0.001	0.002	0.034	0.002
	Oracle	0.001	0.024	0.001	0.003	0.036	0.003

1302 The MNIST-Colors results mirror the Gaussian experiments: at small  $K$ ,  $K$ -to- $K$  baselines can perform competitively on  
1303 IID targets but degrade on OOD targets, while any-to-any models improve OOD performance consistently across generator  
1304 families. As  $K$  increases, the gap between the two regimes narrows, reflecting the reduced need to interpolate beyond the  
1305 training set.

## G. Experimental details for scRNA-seq batch effect transfer

1309 This appendix provides experimental details for the scRNA-seq batch effect transfer experiment in Section 5.2. We describe  
1310 the preprocessing pipeline, model configurations, and evaluation protocol. Generator families and metrics are defined in  
1311 Appendices C and D.

### G.1. Data preprocessing

1315 We use the murine pancreas scRNA-seq atlas data collected in Hrovatin et al. (2023). We use a preprocessed version of the  
1316 dataset distributed by cellxgene<sup>1</sup>. We use the top  $d = 10$  principal components for all downstream analysis. These 10  
1317 PCs are standardized to have unit variance and zero mean.

1<sup>1</sup><https://cellxgene.cziscience.com/collections/296237e2-393d-4e31-b590-b03f74ac5070>

1320  
 1321 *Table 10.* Distributional distances ( $\downarrow$ ) for GMM semi-supervised experiments. IID:  $\|\mu\|_\infty \leq 2.5$ , OOD:  $\|\mu\|_\infty > 2.5$ .

Generator	Method	IID			OOD		
		MMD	SWD	Energy	MMD	SWD	Energy
Flow Matching	Supervised	0.000	0.015	0.000	0.001	0.022	0.001
	Semi-sup.	0.003	0.032	0.002	0.011	0.052	0.011
	Oracle	0.003	0.032	0.002	0.005	0.041	0.005
Energy	Supervised	0.000	0.008	0.000	0.037	0.094	0.037
	Semi-sup.	0.001	0.028	0.001	0.009	0.061	0.008
	Oracle	0.002	0.029	0.001	0.005	0.044	0.005
SWD	Supervised	0.000	0.008	0.000	0.018	0.065	0.018
	Semi-sup.	0.003	0.036	0.003	0.013	0.070	0.013
	Oracle	0.003	0.036	0.003	0.009	0.059	0.009

 1333  
 1334 *Table 11.* Distributional distances ( $\downarrow$ ) for MNIST-Colors across  $K$  unique training colors. We report SWD, MMD-RBF, and Color MSE.  
 1335 **IID Targets**

Model	$K = 10$			$K = 100$			$K = 1,000$			$K = \infty$	
	SWD	MMD	Color MSE	SWD	MMD	Color MSE	SWD	MMD	Color MSE	SWD	MMD
SWD ( $K$ -to- $K$ )	0.0472	0.3746	0.0011	0.0425	0.3395	0.0009	0.0444	0.3549	0.0010	0.0442	0.3549
SWD (Any-to-Any)	0.0415	0.2987	0.0014	0.0400	0.3138	0.0011	0.0404	0.2956	0.0012	0.0414	0.3194
Energy ( $K$ -to- $K$ )	0.0426	0.3225	0.0013	0.0387	0.2905	0.0011	0.0408	0.3104	0.0012	0.0406	0.3104
Energy (Any-to-Any)	0.0392	0.2549	0.0013	0.0378	0.2630	0.0011	0.0400	0.2833	0.0012	0.0404	0.2900
Flow Matching ( $K$ -to- $K$ )	0.0763	1.2057	0.0017	0.0634	0.8942	0.0012	0.0656	0.9301	0.0014	0.0657	0.9301
Flow Matching (Any-to-Any)	0.0678	0.9243	0.0015	0.0575	0.7122	0.0012	0.0595	0.7879	0.0014	0.0563	0.6537

**OOD Targets**

Model	$K = 10$			$K = 100$			$K = 1,000$			$K = \infty$	
	SWD	MMD	Color MSE	SWD	MMD	Color MSE	SWD	MMD	Color MSE	SWD	MMD
SWD ( $K$ -to- $K$ )	0.0646	0.9201	0.0021	0.0493	0.4626	0.0012	0.0459	0.3825	0.0011	0.0459	0.3795
SWD (Any-to-Any)	0.0612	0.8022	0.0016	0.0430	0.3481	0.0012	0.0409	0.3012	0.0012	0.0427	0.3337
Energy ( $K$ -to- $K$ )	0.0638	0.8685	0.0019	0.0463	0.4114	0.0013	0.0424	0.3352	0.0012	0.0422	0.3274
Energy (Any-to-Any)	0.0601	0.7550	0.0018	0.0413	0.3037	0.0012	0.0408	0.2912	0.0012	0.0417	0.3000
Flow Matching ( $K$ -to- $K$ )	0.0793	1.3419	0.0018	0.0692	1.0138	0.0014	0.0677	0.9777	0.0014	0.0659	0.9077
Flow Matching (Any-to-Any)	0.0748	1.2180	0.0018	0.0601	0.7541	0.0013	0.0610	0.8206	0.0014	0.0560	0.6300

 1355  
 1356 **Train-test split.** In order to induce a distribution shift at test-time, we hold out all mice of one condition. We hold out  
 1357 three donor mice at the 2 year age, leaving 53 remaining donor mice for the training dataset.

 1360 **Dataset sampling.** For each donor in the split, we sample sets of  $n = 128$  cells uniformly at random. At training time, for  
 1361 each source donor, a random target donor is selected from the training set, and sets are sampled from both, yielding random  
 1362 source-target sets  $(S_i^{\text{src}}, S_i^{\text{tgt}})$ . This corresponds to STC (any-to-any) training in the terminology of Appendix C.

**G.2. Model configurations**

 1363 We evaluate the three standard continuous-data generator families described in Appendix C: flow matching (FM), energy/  
 1364 MMD regression (Energy), and sliced Wasserstein regression (SWD). For each generator type, we compare two encoder  
 1365 architectures: (i) a *distribution encoder* (as in DCT) that maps each sampled cell set to an embedding, and (ii) a *one-hot*  
 1366 *encoder* that learns a fixed embedding vector for each donor in the training set (a  $K$ -conditioned baseline). Unless stated  
 1367 otherwise, the distribution encoder uses the standard architecture described in Appendix C.2.

 1371 For the one-hot encoder at test time, we assign each test donor to its nearest training donor based on Euclidean distance  
 1372 between donor centroids (computed as the mean of all cells in PCA space for that donor). The model then uses the learned  
 1373 embedding of this nearest training donor.

**Distribution-Conditioned Transport**

<b>Generator</b>	<b>Encoder</b>	<b>Energy ↓</b>	<b>SWD ↓</b>	<b>MMD-RBF ↓</b>
SWD	Any-to-any	$0.1167 \pm 0.0119$	$0.3161 \pm 0.0326$	$0.2334 \pm 0.0239$
	$K$ -to- $K$	$0.3411 \pm 0.0351$	$0.4865 \pm 0.0317$	$0.6822 \pm 0.0702$
Energy	Any-to-any	$0.0729 \pm 0.0090$	$0.3551 \pm 0.0354$	$0.1459 \pm 0.0181$
	$K$ -to- $K$	$0.3623 \pm 0.0327$	$0.5998 \pm 0.0344$	$0.7245 \pm 0.0655$
FM	Any-to-any	$0.1480 \pm 0.0157$	$0.4532 \pm 0.0361$	$0.2960 \pm 0.0315$
	$K$ -to- $K$	$0.4235 \pm 0.0540$	$0.5448 \pm 0.0425$	$0.8470 \pm 0.1080$
scVI		$0.8056 \pm 0.0643$	$0.6190 \pm 0.0630$	$1.6112 \pm 0.1286$

Table 12. Distributional metrics for batch effect transfer between held-out donors (mean  $\pm$  standard error).

**Baselines.** We compare against scVI (Lopez et al., 2018), a deep generative model for scRNA-seq that learns a latent representation while correcting for batch effects. We train scVI on the three held-out donors, using the donor identity as the batch condition. We use a model with 10 latent dimensions and 2 hidden layers, and 128 hidden units per layer. We train for 1000 epochs. We evaluate its ability to transform source distributions to match target distributions via the `transform_batch` functionality, which changes the batch condition passed to the decoder. The resulting normalized expression is then matched to its nearest neighbor in the target dataset.

**Evaluation protocol.** For each test donor pair (source, target), we sample  $n = 10$  independent source-target set pairs. For our models, we generate predictions by sampling from the learned generator conditioned on the source distribution and target embedding. For scVI, we transform the source cells to the target batch and use the observed nearest neighbor for each prediction. We then compute metrics between the predicted and ground truth target distributions.

**Metrics.** We report the standard distributional metrics defined in Appendix D: energy distance, sliced Wasserstein distance, and MMD with an RBF kernel. Metrics are computed between generated and ground-truth target distributions for each sampled pair, aggregated across all test donor combinations and reported as mean  $\pm$  standard error.

### G.3. Full results across metrics

In Table 1, we show MMD-RBF for all evaluated models. In Table 12, we show all three evaluated metrics.

## H. Drug perturbation prediction on organoids (Trellis)

This appendix provides additional details for the Trellis organoid drug-perturbation benchmarks in Section 6.2. We use the Trellis single-cell mass cytometry dataset (Zapatero et al., 2023), representing each cell by 43 protein abundance features and forming 927 matched control–treatment population pairs across 10 patients and 11 treatments. Each population is represented at training time by a set of  $n = 256$  cells sampled uniformly at random.

We evaluate on two replicate-holdout splits (`replicas-1`, `replicas-2`) and three patient-holdout splits (`pdo21`, `pdo27`, `pdo75`). Models use the standard continuous-data generators (SWD/Energy/FM; Appendix C) and are evaluated with the standard metrics (Appendix D). We report three conditioning regimes: a supervised source-conditioned (SC) baseline (MFM) that conditions on treatment identity, a semi-supervised source–target-conditioned (STC) model with a post-hoc ridge predictor for the target embedding, and an oracle STC evaluation that conditions on the true target embedding. Unless stated otherwise, the distribution encoder uses the standard architecture described in Appendix C.2.

**Key hyperparameters.** Unless otherwise stated, we use embedding dimension 256, hidden dimension 256, batch size 4, learning rate  $2 \times 10^{-4}$ , and train for 20,000 epochs.

1430  
 1431 *Table 13.* Distributional metrics ( $\downarrow$ ). Mean  $\pm$  standard deviation reported for IID (replicate holdout) and OOD (patient holdout) settings.  
 1432

Generator	Regime	Energy	SWD	MMD-RBF
<i>IID (Replicate Holdout)</i>				
SWD	Supervised (SC)	$0.0897 \pm 0.0332$	$0.1266 \pm 0.0256$	$0.0064 \pm 0.0022$
	Semi-supervised (STC)	$0.1305 \pm 0.0305$	$0.1524 \pm 0.0185$	$0.0094 \pm 0.0018$
	Oracle (STC)	$0.0154 \pm 0.0005$	$0.0633 \pm 0.0009$	$0.0010 \pm 0.0000$
Energy	Supervised (SC)	$0.0779 \pm 0.0290$	$0.1199 \pm 0.0209$	$0.0054 \pm 0.0019$
	Semi-supervised (STC)	$0.1442 \pm 0.0520$	$0.1588 \pm 0.0297$	$0.0103 \pm 0.0032$
	Oracle (STC)	$0.0169 \pm 0.0026$	$0.0659 \pm 0.0037$	$0.0011 \pm 0.0002$
FM	Supervised (SC)	$0.1431 \pm 0.0377$	$0.1621 \pm 0.0187$	$0.0101 \pm 0.0023$
	Semi-supervised (STC)	$0.1450 \pm 0.0380$	$0.1602 \pm 0.0168$	$0.0106 \pm 0.0025$
	Oracle (STC)	$0.0415 \pm 0.0034$	$0.0929 \pm 0.0054$	$0.0030 \pm 0.0003$
<i>OOD (Patient Holdout)</i>				
SWD	Supervised (SC)	$0.3391 \pm 0.0093$	$0.2452 \pm 0.0078$	$0.0259 \pm 0.0008$
	Semi-supervised (STC)	$0.2580 \pm 0.0403$	$0.2103 \pm 0.0196$	$0.0197 \pm 0.0022$
	Oracle (STC)	$0.0346 \pm 0.0098$	$0.0819 \pm 0.0066$	$0.0026 \pm 0.0009$
Energy	Supervised (SC)	$0.3932 \pm 0.0258$	$0.2622 \pm 0.0074$	$0.0298 \pm 0.0019$
	Semi-supervised (STC)	$0.2503 \pm 0.0556$	$0.2091 \pm 0.0246$	$0.0191 \pm 0.0049$
	Oracle (STC)	$0.0315 \pm 0.0016$	$0.0808 \pm 0.0028$	$0.0022 \pm 0.0002$
FM	Supervised (SC)	$0.2969 \pm 0.0209$	$0.2247 \pm 0.0089$	$0.0227 \pm 0.0006$
	Semi-supervised (STC)	$0.2668 \pm 0.0347$	$0.2138 \pm 0.0167$	$0.0208 \pm 0.0018$
	Oracle (STC)	$0.0705 \pm 0.0214$	$0.1185 \pm 0.0200$	$0.0053 \pm 0.0013$

 1458 **H.1. Full results across metrics**

 1459 **H.2. scGen baseline**

	Energy	SWD	RBF
IID	$0.2924 \pm 0.0881$	$0.2189 \pm 0.0383$	NA
OOD	$0.3278 \pm 0.0345$	$0.2272 \pm 0.0162$	$0.0265 \pm 0.0022$

 1466 **I. Experimental details for lineage-traced scRNA-seq**

 1467 This appendix provides experimental details for the lineage-traced scRNA-seq forecasting experiment in Section 6.3.  
 1468 We describe dataset construction, model configurations, and evaluation. Generator families and metrics are defined in  
 1469 Appendices C and D.

 1472 **I.1. Data preprocessing**

 1474 We use the *in vitro* lineage-traced scRNA-seq dataset from Weinreb et al. (2020), which tracks hematopoietic stem cell  
 1475 differentiation over time points  $t \in \{2, 4, 6\}$  days.

 1476 We download normalized count matrices, metadata, gene names, and clone assignment matrices from the paper data  
 1477 repository<sup>2</sup>. We discard all cells which have not been assigned to a clone.

 1479 **Gene expression preprocessing.** Following standard scRNA-seq preprocessing pipelines (Heumos et al., 2023), we apply:

- 1482 1. Normalization to
- $10^4$
- counts per cell

 1483 <sup>2</sup>[https://kleintools.hms.harvard.edu/paper\\_websites/state\\_fate2020/](https://kleintools.hms.harvard.edu/paper_websites/state_fate2020/)

- 1485 2.  $\log(x + 1)$  transform  
 1486 3. Feature selection of top 10,000 highly variable genes  
 1487 4. Unit variance, zero mean rescaling per gene with maximum value clipped at 10  
 1488 5. PCA to  $d = 50$  principal components  
 1489

1490 Downstream analysis is done on these 50 principal components.  
 1491  
 1492

1493 **Clone set construction.** For each (clone, timepoint) pair, we partition cells into sets of size  $n = 100$  with minimum  
 1494  $n_{\min} = 3$  cells per set. Sets smaller than  $n$  are padded by sampling with replacement from the same (clone, timepoint)  
 1495 group. This yields  $N$  sets  $\{S_i\}_{i=1}^N$ .  
 1496

1497 **Dataset sampling (supervised).** We construct source-target pairs  $(S_i^{\text{src}}, S_i^{\text{tgt}})$  where both sets share the same clone identity  
 1498 but differ by  $\Delta t = 2$  days (i.e.,  $t \in \{2, 4\}$  paired with  $t + 2 \in \{4, 6\}$ ). This results in 1256 clones observed at consecutive  
 1499 timepoints. These clones are split into train/test with ratio 1:1. This split is at the level of clones, not cells – each clone  
 1500 appears in only part of the split.  
 1501

1502 **Dataset sampling (semi-supervised).** For the semi-supervised setting, we additionally construct unstructured pairs by  
 1503 matching source sets from clones without complete targets to random target sets at the correct timepoint. During training,  
 1504 each batch samples from the true paired distribution with probability  $p = 0.25$  and from the random pairing with probability  
 1505  $1 - p = 0.75$ . Clone splits are identical.  
 1506

## 1507 I.2. Model configurations

1508 We evaluate the three standard continuous-data generator families described in Appendix C: flow matching (FM), energy/MMD  
 1509 regression (Energy), and sliced Wasserstein regression (SWD). For each generator type, we compare three regimes  
 1510 (as in the main text): a supervised source-conditioned (SC) baseline trained only on paired clones; a semi-supervised  
 1511 source–target-conditioned (STC) model trained on a mixture of true pairs and random early→late pairings and evaluated  
 1512 using a post-hoc ridge predictor for the target embedding; and an oracle STC evaluation that conditions on the true target  
 1513 embedding. Unless stated otherwise, the distribution encoder uses the standard architecture described in Appendix C.2.  
 1514

1515 For the semi-supervised setting, we fit a linear predictor post-hoc to map source embeddings to target embeddings.  
 1516 Specifically, we encode all training source and target distributions using the semi-supervised encoder, then fit a ridge  
 1517 regression model with 5-fold cross-validation over regularization values  $\alpha \in [10^{-6}, 10^6]$  to predict target embeddings from  
 1518 source embeddings.  
 1519

1520 **Metrics.** We report the standard distributional metrics defined in Appendix D: energy distance (Energy), sliced Wasserstein  
 1521 distance (SWD), and MMD with an RBF kernel (MMD-RBF). Metrics are computed between generated and ground-truth  
 1522 target distributions for each test clone, with results reported as mean  $\pm$  standard error across test clones.  
 1523

## 1524 I.3. Full results across metrics

1525 In Table 6.3 we show MMD-RBF for transport across all evaluated models. In Table 14, we show results across all three  
 1526 evaluated metrics. Regardless of the distributional metric, results remain qualitatively similar.  
 1527

## 1528 J. TCR Repertoire Forecasting Experiments

1529 This appendix provides detailed descriptions of the dataset, model architectures, and evaluation protocols for the T-  
 1530 cell receptor repertoire forecasting experiments presented in Section 6.4. We describe the data processing pipeline, the  
 1531 two discrete sequence generators we employ, and explain architectural differences that may account for their divergent  
 1532 performance. Distributional evaluation metrics are defined in Appendix D and are computed here in ESM2 embedding  
 1533 space.  
 1534

Generator	Regime	Energy ↓	SWD ↓	MMD-RBF ↓
SWD	Supervised (SC)	4.936 ± 0.133	1.679 ± 0.023	9.872 ± 0.267
	Semi-supervised (STC)	4.333 ± 0.108	1.575 ± 0.020	8.666 ± 0.215
	Oracle (STC)	1.494 ± 0.025	0.953 ± 0.011	2.988 ± 0.049
Energy	Supervised (SC)	4.709 ± 0.127	1.650 ± 0.023	9.419 ± 0.255
	Semi-supervised (STC)	4.269 ± 0.103	1.575 ± 0.020	8.538 ± 0.205
	Oracle (STC)	1.440 ± 0.022	0.951 ± 0.011	2.881 ± 0.045
FM	Supervised (SC)	7.049 ± 0.139	2.239 ± 0.025	14.099 ± 0.278
	Semi-supervised (STC)	4.834 ± 0.123	1.747 ± 0.023	9.668 ± 0.247
	Oracle (STC)	2.539 ± 0.052	1.309 ± 0.017	5.077 ± 0.104

Table 14. Distributional metrics for clonal distribution forecasting in held-out clones (mean ± standard error).

## J.1. Data preprocessing

### J.1.1. SOURCE DATA

We use longitudinal TCR repertoire sequencing data from COVID-19 patients collected by Schultheiß et al. (2020). The dataset contains TRB (T-cell receptor beta chain) CDR3 amino acid sequences from 37 patients, of whom 10 were profiled at multiple timepoints during disease progression and recovery. Each repertoire comprises thousands of unique CDR3 junction sequences, typically 10–20 amino acids in length.

### J.1.2. SEQUENCE FILTERING

Raw CDR3 junction sequences are filtered to retain only those composed of the 20 standard amino acids (ACDEFGHIKLM-NPQRSTVWY). Sequences containing ambiguous or non-standard residues are discarded. We truncate or pad sequences to a maximum length of 30 tokens (including special tokens), which accommodates the vast majority of TCR CDR3 sequences.

### J.1.3. TRAIN-TEST SPLIT

The train-test split is performed at the patient level to ensure that test evaluation reflects true generalization to unseen individuals. We partition the 10 multi-timepoint patients, holding out 3 patients (30%) for evaluation and training on the remaining 7. Single-timepoint patients are included only in training for the semi-supervised (any-to-any) model, where they serve as additional unpaired marginals that enrich the learned embedding space.

For the supervised (within-patient) model, only consecutive timepoint pairs  $(t, t+1)$  from multi-timepoint patients are used for training. The semi-supervised model additionally forms cross-patient pairs, training on transitions from any earlier timepoint to any later timepoint across different patients. This pairing strategy encodes the inductive bias that immune repertoires evolve forward in time, while allowing the model to leverage the larger pool of unpaired data.

### J.1.4. TOKENIZATION

Each sequence is tokenized for both ESM2 and ProGen2 vocabularies and cached for efficient training. ESM2 tokenization follows the standard protocol with `<cls>` and `<eos>` special tokens, producing input sequences of the form `<cls>SEQUENCE<eos><pad>...`. ProGen2 tokenization uses `<|bos|>` and `<|eos|>` delimiters with `<|pad|>` for padding.

## J.2. Embedding Architecture

All models share a common distribution encoder that maps repertoires to fixed-dimensional embeddings. The encoder operates in two stages: a pretrained protein language model extracts per-sequence representations, which are then aggregated into a distribution embedding.

## 1595 J.2.1. SEQUENCE ENCODER

1596 We use ESM2 (Lin et al., 2023) with 8 million parameters as the backbone sequence encoder. For each sequence, we extract  
 1597 the final hidden states and apply mean pooling over non-padding positions to obtain a 320-dimensional sequence embedding.  
 1598 The ESM2 backbone is frozen during training to preserve its pretrained representations.  
 1599

 1600 J.2.2. DISTRIBUTION ENCODER  
 1601

1602 Given a set of  $n$  sequence embeddings  $\{h_1, \dots, h_n\} \subset \mathbb{R}^{320}$ , we apply a transformer-based aggregator following the  
 1603 architecture of Zaheer et al. (2017). The aggregator consists of:  
 1604

- 1605 1. An input projection  $\text{Linear}(320, 256) + \text{SELU}$
- 1606 2. Two self-attention layers with 4 heads each
- 1607 3. Mean pooling across sequences
- 1608 4. A latent projection  $\text{Linear}(256, 128) + \text{SELU}$
- 1609 5.  $\ell_2$  normalization of the final embedding

1613 The output is a 128-dimensional distribution embedding  $z \in \mathbb{R}^{128}$  that summarizes the entire repertoire.  
 1614

 1615 J.3. Generator Architectures  
 1616

1617 We compare two discrete sequence generators that share the same conditioning interface but differ fundamentally in their  
 1618 generative mechanisms.  
 1619

## 1620 J.3.1. ESM2 DISCRETE FLOW MATCHING (DFM)

1621 The DFM generator adapts the discrete flow matching framework of Gat & Lipman (2024) to protein sequence generation  
 1622 by repurposing ESM2 as a denoising backbone. Rather than training an autoregressive model, DFM learns to iteratively  
 1623 refine sequences from source to target through a learned flow over discrete tokens.  
 1624

1625 **Architecture.** We extend EsmForMaskedLM with conditioning modules that inject source and target distribution  
 1626 embeddings at two points:  
 1627

- 1628 1. **Input-level conditioning:** The concatenated latent  $[z_{\text{src}}; z_{\text{tgt}}] \in \mathbb{R}^{64}$  is projected through a two-layer MLP  
 1629 ( $\text{Linear}(64, 256) \rightarrow \text{GELU} \rightarrow \text{Linear}(256, 320)$ ) and added to the token embeddings before the transformer  
 1630 encoder. This ensures all 6 transformer layers process conditioned representations.
- 1632 2. **Output-level conditioning:** A parallel projection is added to the hidden states after the transformer, before the language  
 1633 modeling head.

1635 Time conditioning uses sinusoidal embeddings following the original transformer formulation, with the continuous time  
 1636  $t \in [0, 1]$  embedded into the model’s hidden dimension and added to token embeddings.  
 1637

1638 **Training.** Given source sequence  $x^{\text{src}}$  and target sequence  $x^{\text{tgt}}$ , we sample  $t \sim \text{Uniform}(0, 1)$  and construct an interpolated  
 1639 sequence  $x_t$  via Bernoulli masking: at each position,  $x_t$  equals  $x^{\text{tgt}}$  with probability  $t$  and  $x^{\text{src}}$  otherwise. The model predicts  
 1640  $x^{\text{tgt}}$  from  $x_t$  using cross-entropy loss over all positions, including padding tokens. This allows the model to learn length  
 1641 changes (insertions and deletions) between source and target.  
 1642

1643 **Sampling.** Generation proceeds by discrete flow integration from  $t = 0$  to  $t = 1$  with step size  $\Delta t = 0.05$ . At each step,  
 1644 the model predicts a distribution over tokens at each position. The update uses a mixture formula:  
 1645

$$p(x_{t+\Delta t}) = \mathbf{1}_{x_t} + \frac{\Delta t}{1-t} (p_\theta(x^{\text{tgt}} | x_t, t) - \mathbf{1}_{x_t}) \quad (34)$$

1647 where  $\mathbf{1}_{x_t}$  is the one-hot encoding of the current sequence. Vocabulary constraints enforce that position 0 contains only  
 1648  $\langle \text{cls} \rangle$  and subsequent positions contain only amino acids,  $\langle \text{eos} \rangle$ , or  $\langle \text{pad} \rangle$ .  
 1649

1650 J.3.2. PROGEN2 WITH PREFIX CONDITIONING  
16511652 The ProGen2 generator uses the pretrained ProGen2 language model (Nijkamp et al., 2023) with prefix-based conditioning  
1653 for source-to-target generation.1654  
1655 **Architecture.** We use `progen2-medium` as the backbone. Conditioning is implemented by inserting a learned prefix  
1656 embedding between source and target sequences:

- 1657
- 
- 1658 1. The concatenated latent
- $[z_{\text{src}}; z_{\text{tgt}}]$
- is projected through
- $\text{Linear}(64, 256) \rightarrow \text{GELU} \rightarrow \text{Linear}(256, d_{\text{model}})$
- 
- 1659
- 
- 1660 2. The resulting vector is inserted as a single prefix token between the source and target portions of the input
- 
- 1661
- 
- 1662 3. The attention mask is adjusted to allow the prefix to attend to source tokens and target tokens to attend to the prefix
- 
- 1663

1664  
1665 **Training.** The input sequence is formed by concatenating source and target: `<bos>SOURCE<eos><bos>TARGET<eos>`.  
1666 The prefix is inserted at the boundary. The model is trained with causal language modeling loss computed only on target  
1667 tokens, ignoring source and prefix positions.1668  
1669 **Sampling.** Generation proceeds autoregressively: given the source sequence and predicted target embedding, the model  
1670 generates target tokens one at a time until `<eos>` or maximum length. To encourage diversity when generating multiple  
1671 samples, we add Gaussian noise with scale 0.1 to the latent embeddings.1672 J.4. Evaluation Protocol  
16731674 All models are evaluated on the same task: given a patient’s repertoire at time  $t$ , forecast the repertoire at time  $t+1$ . We  
1675 evaluate only on consecutive timepoint pairs from held-out patients.  
16761677 J.4.1. EMBEDDING COMPUTATION  
16781679 For each repertoire, we sample sequences uniformly at random and compute the distribution embedding using the trained  
1680 encoder. At training time, we use 2048 sequences per repertoire to compute embeddings for the predictor.  
16811682 J.4.2. TARGET EMBEDDING PREDICTION  
16831684 For the semi-supervised model, we train a ridge regression predictor on the delta  $z_{t+1} = z_t + \Delta z$ ,  $\Delta z = W z_t + b$  on  
1685 consecutive pairs from training patients, using 2048 sequences per repertoire for embedding computation. Regularization  
1686 strength is selected via leave-one-out cross-validation over a grid of 20 log-spaced values from  $10^{-4}$  to  $10^4$ .1687 J.4.3. SEQUENCE GENERATION AND COMPARISON  
1688

1689 For each test pair:

- 1690
- 
- 1691 1. Encode the current timepoint repertoire to obtain
- $z_t$
- 
- 1692
- 
- 1693 2. Form the target embedding
- $\hat{z}_{t+1}$
- using either (i) a supervised within-patient baseline that sets
- $\hat{z}_{t+1} = 0$
- (source-
- 
- 1694 conditioned), (ii) a semi-supervised ridge predictor
- $\hat{z}_{t+1}$
- (source–target-conditioned), or (iii) an oracle that uses the true
- 
- 1695 target embedding
- $z_{t+1}$
- 
- 1696
- 
- 1697 3. Generate 2048 sequences conditioned on
- $z_t$
- and
- $\hat{z}_{t+1}$
- 
- 1698
- 
- 1699 4. Re-tokenize generated sequences for ESM2 and compute mean-pooled embeddings
- 
- 1700
- 
- 1701 5. Sample 2048 sequences from the actual next timepoint and compute their ESM2 embeddings
- 
- 1702
- 
- 1703 6. Compute distributional distances between generated and actual sequence embeddings
- 
- 1704

## 1705 J.4.4. METRICS

1706 We report the standard distributional metrics defined in Appendix D (energy distance, sliced Wasserstein distance, and  
 1707 MMD-RBF), computed in the ESM2 embedding space.

1708 Results are averaged over all test pairs, with standard errors computed across pairs.

1709

## 1710 J.5. Discussion: Why ProGen2 Underperforms

1711 The results in Table 4 show that the DFM generator substantially outperforms ProGen2, particularly in the semi-supervised  
 1712 setting. We hypothesize several factors contribute to this gap.

1713

1714 **Autoregressive vs. iterative refinement.** ProGen2 generates sequences left-to-right, committing to each token before  
 1715 seeing subsequent context. The DFM generator, by contrast, refines all positions simultaneously through iterative denoising.  
 1716 For the task of transforming one sequence into another (rather than generating de novo), iterative refinement may better  
 1717 preserve structural features of the source while adapting to the target distribution.

1718

1719 **Conditioning mechanism.** The DFM generator injects conditioning at both input and output levels of every transformer  
 1720 layer, ensuring the distributional signal permeates the entire network. ProGen2’s prefix conditioning provides a single  
 1721 conditioning vector that the model must propagate through causal attention. This architectural difference may limit  
 1722 ProGen2’s ability to leverage the target embedding for precise distributional control.

1723

## 1724 J.6. Hyperparameter Summary

1725 Table 15 summarizes the key hyperparameters for reproducibility.

1726

1727

1728 *Table 15.* Hyperparameters for TCR repertoire forecasting experiments.

Hyperparameter	DFM	ProGen2
<i>Dataset</i>		
Sequences per repertoire (training)	2048	2048
Maximum sequence length	30	30
Test fraction (multi-TP patients)	0.3	0.3
<i>Encoder</i>		
ESM2 backbone	esm2_t6_8M	esm2_t6_8M
Distribution embedding dim	128	128
Self-attention layers	2	2
Attention heads	4	4
<i>Generator</i>		
Backbone	esm2_t6_8M	progen2-medium
Conditioning dimension	256	256
Latent input dimension	32	32
Temperature	1.0	1.0
Sample step size ( $\Delta t$ )	0.05	—
<i>Training</i>		
Batch size	1	1
Learning rate	$2 \times 10^{-5}$	$2 \times 10^{-5}$
Optimizer	Adam	Adam
<i>Evaluation</i>		
Sequences per repertoire	2048	2048
Predictor training samples	2048	2048
Predictor regularization	RidgeCV	RidgeCV

## 1760 J.6.1. RESULTS

1761 Table 16 records results across all metrics.

1762

 1763
 

1764 Table 16. TCR repertoire prediction performance across conditioning strategies (mean  $\pm$  standard error). DFM shows meaningful  
 1765 conditioning (Oracle  $\gg$  Predictor  $\gg$  Within-patient), while ProGen2 shows little conditioning effect due to encoder collapse (all  
 1766 embeddings have cosine similarity  $> 0.998$ ).

Generator	Regime	Energy $\downarrow$	SWD $\downarrow$	MMD-RBF $\downarrow$
ProGen2	Within-patient (SC)	0.0682 $\pm$ 0.0103	0.4015 $\pm$ 0.0118	0.0588 $\pm$ 0.0085
	Predictor (STC)	0.0681 $\pm$ 0.0104	0.4012 $\pm$ 0.0118	0.0587 $\pm$ 0.0085
	Oracle (STC)	0.0676 $\pm$ 0.0114	0.3856 $\pm$ 0.0148	0.0567 $\pm$ 0.0093
DFM	Within-patient (SC)	0.0684 $\pm$ 0.0120	0.3895 $\pm$ 0.0147	0.0579 $\pm$ 0.0100
	Predictor (STC)	0.0262 $\pm$ 0.0050	0.2778 $\pm$ 0.0152	0.0215 $\pm$ 0.0043
	Oracle (STC)	0.0090 $\pm$ 0.0015	0.2199 $\pm$ 0.0217	0.0078 $\pm$ 0.0013

1775

1776

## 1777 J.7. Can we close the semi-supervised–oracle gap?

 1778 In the TCR setting, the gap between the *predictor* and *oracle* regimes in Table 16 reflects how well we can predict the next-  
 1779 timepoint embedding  $z_{t+1}$  from the current embedding  $z_t$  using only the limited paired longitudinal data. The main results  
 1780 use a ridge predictor (chosen for stability in low-data regimes). We also evaluated more expressive nonlinear predictors,  
 1781 but on the patient-holdout split these did not improve performance and often degraded it, consistent with overfitting and  
 1782 distribution shift across individuals.  
 1783

Latent predictor	Energy $\downarrow$	SWD $\downarrow$	MMD-RBF $\downarrow$
Ridge (linear)	0.2902 $\pm$ .0090	1.6447 $\pm$ .0184	0.2619 $\pm$ .0104
MLP	0.2902 $\pm$ .0091	1.6455 $\pm$ .0205	0.2652 $\pm$ .0103
Random Forest	0.2901 $\pm$ .0090	1.6461 $\pm$ .0182	0.2623 $\pm$ .0104

 1784
 

1785 Table 17. Closing the semi-supervised–oracle gap on TCR forecasting. Each row uses the same trained STC generator but varies how the  
 1786 target embedding is predicted.

1787

1788

1789

1790

1791

1792

1793

1794

1795

1796

1797

1798

1799

1800

1801

1802

1803

1804

1805

1806

1807

1808

1809

1810

1811

1812

1813

1814