

# Relaxing Bijectivity Constraints with Continuously Indexed Normalising Flows

Rob Cornish<sup>1</sup> Anthony Caterini<sup>1</sup> George Deligiannidis<sup>1,2</sup> Arnaud Doucet<sup>1</sup>

## Abstract

We show that normalising flows become pathological when used to model targets whose supports have complicated topologies. In this scenario, we prove that a flow must become arbitrarily numerically noninvertible in order to approximate the target closely. This result has implications for all flow-based models, and especially *residual flows* (ResFlows), which explicitly control the Lipschitz constant of the bijection used. To address this, we propose *continuously indexed flows* (CIFs), which replace the single bijection used by normalising flows with a continuously indexed family of bijections, and which can intuitively “clean up” mass that would otherwise be misplaced by a single bijection. We show theoretically that CIFs are not subject to the same topological limitations as normalising flows, and obtain better empirical performance on a variety of models and benchmarks.

## 1 Introduction

*Normalising flows* (Rezende & Mohamed, 2015) have become popular methods for density estimation (Dinh et al., 2017; Papamakarios et al., 2017; Kingma & Dhariwal, 2018; Chen et al., 2019). These methods model an unknown target distribution  $P_X^*$  on a data space  $\mathcal{X} \subseteq \mathbb{R}^d$  as the marginal of  $X$  obtained by the generative process

$$Z \sim P_Z, \quad X := f(Z), \quad (1)$$

where  $P_Z$  is a *prior* distribution on a space  $\mathcal{Z} \subseteq \mathbb{R}^d$ , and  $f: \mathcal{X} \rightarrow \mathcal{Z}$  is a bijection. The use of a bijection means the density of  $X$  can be computed analytically by the change-of-variables formula, and the parameters of  $f$  can be learned by maximum likelihood using i.i.d. samples from  $P_X^*$ .

To be effective, a normalising flow model must specify an expressive family of bijections with tractable Jacobians.

<sup>1</sup>University of Oxford, Oxford, United Kingdom <sup>2</sup>The Alan Turing Institute, London, United Kingdom. Correspondence to: Rob Cornish <rcornish@robots.ox.ac.uk>.

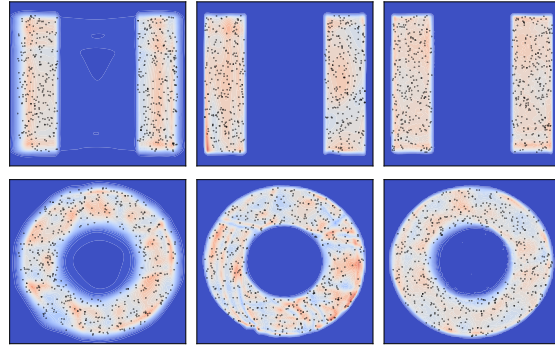


Figure 1: Densities learned by a 10-layer ResFlow (left), 100-layer ResFlow (middle), and 10-layer CIF-ResFlow (right) for two datasets (samples shown in black) that are not homeomorphic to the Gaussian prior. The 10-layer ResFlow visibly leaks mass outside of the support of the target due to its small bi-Lipschitz constant. The larger ResFlow improves on this, but still achieves smaller average log probability than the CIF-ResFlow, as is apparent from the greater homogeneity of the right-hand densities.

Affine coupling layers (Dinh et al., 2015; 2017), autoregressive maps (Germain et al., 2015; Papamakarios et al., 2017), invertible linear transformations (Kingma & Dhariwal, 2018), ODE-based maps (Grathwohl et al., 2019), and invertible ResNet blocks (Behrmann et al., 2019; Chen et al., 2019) are all examples of such bijections that can be composed to produce expressive flows. These models have demonstrated significant promise in their ability to model complex datasets and to synthesise realistic data.

In all these cases,  $f$  and  $f^{-1}$  are both continuous. It follows that  $f$  is a *homeomorphism*, and therefore preserves the topology of its domain (Runde, 2007, Definition 3.3.10). As Dupont et al. (2019) and Dinh et al. (2019) mention, this seems intuitively problematic when  $P_Z$  and  $P_X^*$  are supported on domains with distinct topologies, which occurs for example when the supports differ in their number of connected components or “holes”, or when they are “knotted” differently. This seems inevitable in practice, as  $P_Z$  is usually quite simple (e.g. a Gaussian) while  $P_X^*$  is very complicated (e.g. a distribution over images).

As our first contribution, we make precise the consequences

of using a topologically misspecified prior. We confirm that in this case it is indeed impossible to recover the target perfectly if  $f$  is a homeomorphism. Moreover, in [Theorem 2.1](#) we prove that, in order to *approximate* such a target arbitrarily well, we must have  $\text{BiLip } f \rightarrow \infty$ , where  $\text{BiLip } f$  denotes the *bi-Lipschitz constant* of  $f$  defined as the infimum over  $M \in [1, \infty]$  such that

$$M^{-1}\|z - z'\| \leq \|f(z) - f(z')\| \leq M\|z - z'\| \quad (2)$$

for all  $z, z' \in \mathcal{Z}$ . [Theorem 2.1](#) applies essentially regardless of the training objective, and has implications for the case that  $P_Z$  and  $P_X^*$  both have full support but are heavily concentrated on regions that are not homeomorphic. Since  $\text{BiLip } f$  is a natural measure of the ‘‘invertibility’’ of  $f$  ([Behrmann et al., 2020](#)), this result shows that the goal of designing neural networks with well-conditioned inverses is fundamentally at odds with the goal of designing neural networks that can approximate complicated densities.

[Theorem 2.1](#) also has immediate implications for *residual flows* (ResFlows) ([Behrmann et al., 2019](#); [Chen et al., 2019](#)), which have recently achieved state-of-the-art performance on several large-scale density estimation tasks. Unlike models based on triangular maps ([Jaini et al., 2019](#)), ResFlows have the attractive feature that the structure of their Jacobians is unconstrained, which may explain their greater expressiveness. However, as part of the construction, the bi-Lipschitz constant of  $f$  is bounded, and so these models must be composed many times in order to achieve overall the large bi-Lipschitz constant required for a complex  $P_X^*$ .<sup>1</sup>

To address this problem we introduce *continuously indexed flows* (CIFs), which generalise (1) by replacing the single bijection  $f$  with an indexed family of bijections  $\{F(\cdot; u)\}_{u \in \mathcal{U}}$ , where the index set  $\mathcal{U}$  is continuous. Intuitively, CIFs allow mass that would be erroneously placed by a single bijection to be rerouted into a more optimal location. We show that CIFs can learn the support of a given  $P_X^*$  exactly regardless of the topology of the prior, and without the bi-Lipschitz constant of any  $F(\cdot; u)$  necessarily becoming infinite. CIFs do not specify the form of  $F$ , and can be used in conjunction with any standard normalising flow architecture directly.

Our use of a continuous index overcomes several limitations associated with alternative approaches based on a discrete index ([Dinh et al., 2019](#); [Duan, 2019](#)), which suffer either from a discontinuous loss landscape or an intractable computational complexity. However, as a consequence, we sacrifice the ability to compute the likelihood of our model analytically. To address this, we propose a variational approximation that exploits the bijective structure of the model and is suitable for training large-scale models in practice. We empirically evaluate CIFs applied to ResFlows, neural

spline flows (NSFs) ([Durkan et al., 2019](#)), masked autoregressive flows (MAFs) ([Papamakarios et al., 2017](#)), and RealNVPs ([Dinh et al., 2017](#)), obtaining improved performance in all cases. We observe a particular benefit for ResFlows: with a 10-layer CIF-ResFlow we surpass the performance of a 100-layer baseline ResFlow and achieve state-of-the-art results on several benchmark datasets.

## 2 Bi-Lipschitz Constraints on Pushforwards

Normalising flows fall into a larger class of density estimators based on *pushforwards*. Given a prior measure  $P_Z$  on  $\mathcal{Z}$  and a mapping  $f : \mathcal{Z} \rightarrow \mathcal{X}$ , these models are defined as

$$P_X := f\#P_Z,$$

where the right-hand side denotes a distribution with  $f\#P_Z(B) := P_Z(f^{-1}(B))$  for Borel  $B \subseteq \mathcal{X}$ . Normalising flows take  $f$  to be bijective, which under sufficient regularity yields a closed-form expression for the density<sup>2</sup> of  $P_X$  ([Billingsley, 2008](#), Theorem 17.2).

Intuitively, the pushforward map  $f$  *transports* the mass allocated by  $P_Z$  into  $\mathcal{X}$ -space, thereby defining  $P_X$  based on where each unit of mass ends up. This imposes a global constraint on  $f$  if  $P_X$  is to match perfectly a given target  $P_X^*$ . In particular, denote by  $\text{supp } P_Z$  the *support* of  $P_Z$ . While the precise definition of the support involves topological formalities (see [Section B.1](#) in the Supplement), intuitively this set defines the region of  $\mathcal{Z}$  to which  $P_Z$  assigns mass. It is then straightforward to show that  $P_X = P_X^*$  only if

$$\text{supp } P_X^* = \overline{f(\text{supp } P_Z)}, \quad (3)$$

where  $\overline{A}$  denotes the closure of  $A$  in  $\mathcal{X}$ .<sup>3</sup>

The constraint (3) is especially onerous for normalising flows because of their bijectivity. In practice,  $f$  and  $f^{-1}$  are invariably both continuous, and so  $f$  is a *homeomorphism*. Consequently, for these models (3) entails<sup>4</sup>

$$\text{supp } P_X = \text{supp } P_X^* \text{ only if } \text{supp } P_Z \cong \text{supp } P_X^*, \quad (4)$$

where  $\mathcal{A} \cong \mathcal{B}$  means that  $\mathcal{A}$  and  $\mathcal{B}$  are *homeomorphic*, i.e. isomorphic as topological spaces ([Runde, 2007](#), Definition 3.3.10). This means that  $\text{supp } P_Z$  and  $\text{supp } P_X^*$  must exactly share *all* topological properties, including number of connected components, number of ‘‘holes’’, the way they are ‘‘knotted’’, etc., in order to learn the target perfectly. Condition (4) therefore suggests that normalising flows are not optimally suited to the task of learning complex real-world densities, where such topological mismatch seems inevitable.

<sup>2</sup>Throughout, by ‘‘density’’ we mean Lebesgue density. We will write densities using lowercase, e.g.  $p_X$  for the measure  $P_X$ .

<sup>3</sup>See [Proposition B.3](#) in the Supplement for a proof.

<sup>4</sup>Note that  $f(\overline{\text{supp } P_Z}) = \overline{f(\text{supp } P_Z)}$  here since  $\text{supp } P_Z$  is closed by [Proposition B.2](#) in the Supplement.

<sup>1</sup>[Chen et al. \(2019\)](#) report using 100-200 layers to learn even simple 2D densities.

However, (4) only rules out the limiting case  $P_X = P_X^*$ . In practice it is likely enough to have  $P_X \approx P_X^*$ , and it is therefore relevant to consider the implications of a topologically misspecified prior in this case also. Intuitively, this seems to require  $f$  become *almost* nonbijective as  $P_X$  approaches  $P_X^*$ , but it is not immediately clear what this means, or whether this must occur for all models. Likewise, in practice it might be reasonable to assume the density of  $P_X^*$  is everywhere strictly positive. In this case, even if  $P_X^*$  is concentrated on some very complicated set, the constraint (4) would trivially be met if  $P_Z$  is Gaussian, for example. Nevertheless, it seems that infinitesimal regions of mass should not significantly change the behaviour required of  $f$ , and we would therefore like to extend (4) to apply here also.

The bi-Lipschitz constant (2) naturally quantifies the “invertibility” of  $f$ . Behrmann et al. (2020) recently showed a relationship between the bi-Lipschitz constant and the numerical invertibility of  $f$ . If  $f$  is injective and differentiable,

$$\text{BiLip } f = \max \left( \sup_{z \in \mathcal{Z}} \|Df(z)\|_{\text{op}}, \sup_{x \in f(\mathcal{Z})} \|Df^{-1}(x)\|_{\text{op}} \right),$$

where  $Dg(y)$  is the Jacobian of  $g$  at  $y$  and  $\|\cdot\|_{\text{op}}$  is the operator norm. A large bi-Lipschitz constant thus means  $f$  or  $f^{-1}$  “jumps” somewhere in its domain. More generally, if  $f$  is not injective, then  $\text{BiLip } f = \infty$ , while if  $\text{BiLip } f < \infty$ , then  $f$  is a homeomorphism from  $\mathcal{Z}$  to  $f(\mathcal{Z})$ .<sup>5</sup>

The following theorem shows that if the supports of  $P_Z$  and  $P_X^*$  are not homeomorphic, then the bi-Lipschitz constant of  $f$  must grow arbitrarily large in order to approximate  $P_X^*$ . Here  $\xrightarrow{\mathcal{D}}$  denotes weak convergence.

**Theorem 2.1.** *Suppose  $P_Z$  and  $P_X^*$  are probability measures on  $\mathbb{R}^{d_Z}$  and  $\mathbb{R}^{d_X}$  respectively, and that  $\text{supp } P_Z \not\cong \text{supp } P_X^*$ . Then for any sequence of measurable  $f_n : \mathbb{R}^{d_Z} \rightarrow \mathbb{R}^{d_X}$ , we can have  $f_n \# P_Z \xrightarrow{\mathcal{D}} P_X^*$  only if*

$$\lim_{n \rightarrow \infty} \text{BiLip } f_n = \infty.$$

Weak convergence is implied by the minimisation of all standard statistical divergences used to train generative models, including the KL and Jensen-Shannon divergences and the Wasserstein metric (Arjovsky et al., 2017, Theorem 2). Thus, Theorem 2.1 states that these quantities can vanish only if the bi-Lipschitz constant of the learned mapping becomes arbitrarily large. Likewise, note that we do not assume  $d_Z = d_X$  so that this result also applies to injective flow models (Kumar et al., 2019), as well as other pushforward-based models such as GANs (Goodfellow et al., 2014).<sup>6</sup>

<sup>5</sup>See Section B.2 in the Supplement for proofs.

<sup>6</sup>However, the implications for GANs seem less problematic since a GAN generator is not usually assumed to be bijective.

Theorem 2.1 also applies when  $\text{supp } P_Z$  is *almost* not homeomorphic to  $\text{supp } P_X^*$ , as is made precise by the following corollary. Here  $\rho$  denotes any metric for the weak topology; see Chapter 6 of Villani (2008) for standard examples.

**Corollary 2.2.** *Suppose  $P_Z$  and  $P_X^0$  are probability measures on  $\mathbb{R}^{d_Z}$  and  $\mathbb{R}^{d_X}$  respectively with  $\text{supp } P_Z \not\cong \text{supp } P_X^0$ . Then there exists nonincreasing  $M : [0, \infty) \rightarrow [1, \infty]$  with  $M(\epsilon) \rightarrow \infty$  as  $\epsilon \rightarrow 0$  such that, for any probability measure  $P_X^*$  on  $\mathbb{R}^{d_X}$ , we have  $\text{BiLip } f \geq M(\epsilon)$  whenever  $\rho(P_X^*, P_X^0) \leq \epsilon$  and  $\rho(f \# P_Z, P_X^*) \leq \epsilon$ .*

In other words, if the target is close to a probability measure with non-homeomorphic support to that of the prior (i.e.  $\rho(P_X^*, P_X^0)$  is small), and if the model is a good approximation of the target (i.e.  $\rho(f \# P_Z, P_X^*)$  is small), then the Bi-Lipschitz constant of  $f$  must be large.

Proofs of these results are in Section B.3 of the Supplement.

## 2.1 Practical Implications

The results of this section indicate a limitation of existing flow-based density models. This is most direct for *residual flows* (ResFlows) (Behrmann et al., 2019; Chen et al., 2019), which take  $f = f_L \circ \dots \circ f_1$  with each layer of the form

$$f_\ell^{-1}(x) = x + g_\ell(x), \quad \text{Lip } g_\ell \leq \kappa < 1. \quad (5)$$

Here  $\text{Lip}$  denotes the Lipschitz constant, which is bounded by a fixed constant  $\kappa$  throughout training. The Lipschitz constraint is enforced by spectral normalisation (Miyato et al., 2018; Gouk et al., 2018) and ensures each  $f_\ell$  is bijective. However, it also follows (Behrmann et al., 2019, Lemma 2) that

$$\text{BiLip } f \leq \max(1 + \kappa, (1 - \kappa)^{-1})^L < \infty, \quad (6)$$

and Theorem 2.1 thus restricts how well a ResFlow can approximate  $P_X^*$  with non-homeomorphic support to  $P_Z$ . Figure 1 illustrates this in practice for simple 2-D examples.

It is possible to relax (6) by taking  $\kappa \rightarrow 1$ . However, this can have a detrimental effect on the variance of the Russian roulette estimator (Kahn, 1955) used by Chen et al. (2019) to compute the Jacobian, and in Section B.4 of the Supplement we give a simple example in which the variance is in fact infinite. Alternatively, we can also loosen the bound (6) by taking  $L \rightarrow \infty$ , and Figure 1 shows that this does indeed lead to better performance. However, greater depth means greater computational cost. In the next section we describe an alternative approach that allows relaxing the bi-Lipschitz constraint of Theorem 2.1 without modifying either  $\kappa$  or  $L$ , and thus avoids these potential issues.

Unlike ResFlows, most normalising flows used in practice have an unconstrained bi-Lipschitz constant (Behrmann et al., 2020). As a result, Theorem 2.1 does not prevent

these models from approximating non-homeomorphic targets arbitrarily well, and indeed several architectures have been proposed that can in principle do so (Huang et al., 2018; Jaini et al., 2019). Nevertheless, the constraint (4) shows that these models still face an underlying limitation in practice, and suggests we may improve performance more generally by relaxing the requirement of bijectivity. We verify empirically in Section 5 that, in addition to ResFlows, our proposed method also yields benefits for flows without an explicit bi-Lipschitz constraint.

Finally, Theorem 2.1 has implications for the numerical stability of normalising flows. It was recently pointed out by Behrmann et al. (2020) that, while having a well-defined mathematical inverse, many common flows can become *numerically* noninvertible over the course of training, leading to low-quality reconstructions and calling into question the accuracy of density values output by the change-of-variables formula. Behrmann et al. (2020) suggest explicitly constraining BiLip  $f$  in order to avoid this problem. Theorem 2.1 shows that this involves a fundamental tradeoff against expressivity: if greater numerical stability is required of our normalising flow, then we must necessarily reduce the set of targets we can represent arbitrarily well.

### 3 Continuously Indexed Flows

In this section we propose *continuously indexed flows* (CIFs) for relaxing the bijectivity of standard normalising flows. We begin by defining the model we consider, and then detail our suggested training and inference procedures. In the next section we discuss advantages over related approaches.

#### 3.1 Model Specification

CIFs are obtained by replacing the single bijection  $f$  used by normalising flows with an indexed family  $\{F(\cdot; u)\}_{u \in \mathcal{U}}$ , where  $\mathcal{U} \subseteq \mathbb{R}^{d_u}$  is our index set and each  $F(\cdot; u) : \mathcal{Z} \rightarrow \mathcal{X}$  is a bijection. We then define the model  $P_X$  as the marginal of  $X$  obtained from the following generative process:

$$Z \sim P_Z, \quad U \sim P_{U|Z}(\cdot|Z), \quad X := F(Z; U). \quad (7)$$

Like (1), we assume a prior  $P_Z$  on  $\mathcal{Z}$ , but now also require conditional distributions  $P_{U|Z}(\cdot|z)$  on  $\mathcal{U}$  for each  $z \in \mathcal{Z}$ .

We can increase the complexity of (7) by taking  $P_Z$  itself to have the same form. This is directly analogous to the standard practice of composing simple bijections to obtain a richer class of normalising flows. In our context, stacking  $L$  layers of (7) corresponds to the generative process

$$Z_0 \sim P_{Z_0}, \quad U_\ell \sim P_{U_\ell|Z_{\ell-1}}(\cdot|Z_{\ell-1}), \quad Z_\ell := F_\ell(Z_{\ell-1}; U_\ell), \quad (8)$$

where  $\ell \in \{1, \dots, L\}$ . We then take  $P_X$  to be the marginal of  $X := Z_L$ . We have found this construction

to improve significantly the expressiveness of our models and make extensive use of it in our experiments below. Note that this corresponds to an instance of (7) where, defining  $F^\ell(\cdot; u_1, \dots, u_\ell) := F_\ell(\cdot; u_\ell) \circ \dots \circ F_1(\cdot; u_1)$ , we take  $Z = Z_0$ ,  $U = (U_1, \dots, U_L)$ ,  $P_{U|Z}(du|z) = \prod_\ell P_{U_\ell|Z_{\ell-1}}(du_\ell|F^\ell(z; u_1, \dots, u_\ell))$ , and  $F = F^L$ . We use this to streamline some of the discussion below.

Previous works, most notably RAD (Dinh et al., 2019), have considered related models with a discrete index set  $\mathcal{U}$ . We instead consider a *continuous* index. In particular, our  $\mathcal{U}$  will be an open subset of  $\mathbb{R}^{d_u}$ , with each  $P_{U|Z}(\cdot|z)$  having a density  $p_{U|Z}(\cdot|z)$ . A continuous index confers various advantages that we describe in Section 4. The choice also requires a distinct approach to training and inference that we describe in Section 3.2.

We require choices of  $p_{U|Z}$  and  $F$  for each layer of our model. Straightforward possibilities are

$$F(z; u) = f\left(e^{-s(u)} \odot z - t(u)\right) \quad (9)$$

$$p_{U|Z}(\cdot|z) = \text{Normal}(\mu^p(z), \Sigma^p(z)) \quad (10)$$

for any bijection  $f$  (e.g. a ResFlow step) and appropriately defined neural networks  $s$ ,  $t$ ,  $\mu^p$ , and  $\Sigma^p$ .<sup>7</sup> Here the exponential of a vector is meant elementwise, and  $\odot$  denotes elementwise multiplication. Note that (9) may be used with all existing normalising flow implementations out-of-the-box. These choices yielded strong empirical results despite their simplicity, but more sophisticated alternatives are certainly possible and may bring improvements in some applications.

#### 3.2 Training and Inference

Heuristically,<sup>8</sup> (7) yields the joint ‘‘density’’

$$p_{X,U,Z}(x, u, z) := p_Z(z) p_{U|Z}(u|z) \delta(x - F(z; u)),$$

where  $p_Z$  is the density of  $P_Z$  and  $\delta$  is the Dirac delta. If  $F$  is sufficiently regular, we can marginalise out the dependence on  $z$  by making the change of variable  $x' := F(z; u)$ , which means  $dz = |\det DF^{-1}(x'; u)| dx'$ .<sup>9</sup> This yields a proper density for  $(X, U)$  by integrating over  $x'$ :

$$p_{X,U}(x, u) := p_Z(F^{-1}(x; u)) \times p_{U|Z}(u|F^{-1}(x; u)) |\det DF^{-1}(x; u)|. \quad (11)$$

For an  $L$ -layered model, an extension of this argument also gives the following joint density for each  $(Z_\ell, U_{1:\ell})$ :

$$p_{Z_\ell, U_{1:\ell}}(z_\ell, u_{1:\ell}) := p_{Z_{\ell-1}, U_{1:\ell-1}}(F_\ell^{-1}(z_\ell; u_\ell), u_{1:\ell-1}) \times p_{U_\ell|Z_{\ell-1}}(u_\ell|F_\ell^{-1}(z_\ell; u_\ell)) |\det DF_\ell^{-1}(z_\ell; u_\ell)|. \quad (12)$$

<sup>7</sup>Note this requires  $\mathcal{Z} = \mathcal{X} = \mathbb{R}^d$  and  $\mathcal{U} = \mathbb{R}^{d_u}$ , i.e. these domains are not strict subsets. We assume this in all our experiments.

<sup>8</sup>We make this rigorous in Section B.5 of the Supplement.

<sup>9</sup>Here  $DF(z; u)$  denotes the Jacobian with respect to  $z$  only.

Taking  $X := Z_L$  as before we obtain  $p_{X,U_{1:L}}$  and hence a density for  $P_X$  via

$$p_X(x) := \int p_{X,U_{1:L}}(x, u_{1:L}) du_{1:L}. \quad (13)$$

Since  $\mathcal{U}$  is continuous, this is not analytically tractable. To facilitate likelihood-based training and inference, we make use of a variational scheme that we describe now.

Assuming an  $L$ -layered model (8), we introduce an approximate posterior density  $q_{U_{1:L}|X} \approx p_{U_{1:L}|X}$  and consider the evidence lower bound (ELBO) of  $\log p_X(x)$ :

$$\mathcal{L}(x) := \mathbb{E}_{u_{1:L} \sim q_{U_{1:L}|X}(\cdot|x)} \left[ \log \frac{p_{X,U_{1:L}}(x, u_{1:L})}{q_{U_{1:L}|X}(u_{1:L}|x)} \right]. \quad (14)$$

It is a standard result that  $\mathcal{L}(x) \leq \log p_X(x)$  with equality if and only if  $q_{U_{1:L}|X}$  is the exact posterior  $p_{U_{1:L}|X}$ . This allows learning an approximation to  $P_X^*$  by maximising  $\sum_{i=1}^n \mathcal{L}(x_i)$  jointly in  $p_{X,U_{1:L}}$  and  $q_{U_{1:L}|X}$ , where we assume a dataset of  $n$  i.i.d. samples  $x_i \sim P_X^*$ .

We now consider how to parametrise an effective  $q_{U_{1:L}|X}$ . Standard approaches to designing inference networks for variational autoencoders (VAEs) (Kingma & Welling, 2014; Rezende et al., 2014; Rezende & Mohamed, 2015; Kingma et al., 2016), while mathematically valid, would not exploit the conditional independencies induced by the bijective structure of (8). We therefore propose a novel inference network that is specifically targeted towards our model, which we compare with existing VAE approaches in Section 4.3.

In particular, our  $q_{U_{1:L}|X}$  has the following form:

$$q_{U_{1:L}|X}(u_{1:L}|x) := \prod_{\ell=1}^L q_{U_\ell|Z_\ell}(u_\ell|z_\ell), \quad (15)$$

with  $z_L := x$  and  $z_\ell := F_{\ell+1}^{-1}(z_{\ell+1}; u_{\ell+1})$  for  $\ell \in \{1, \dots, L-1\}$ , and  $q_{U_\ell|Z_\ell}$  can be any parameterised conditional density. We show in Section B.6 of the Supplement that the posterior  $p_{U_{1:L}|X}$  factors in the same way as (15), so that we do not lose any generality. Observe also that this scheme shares parameters between  $q_{U_{1:L}|X}$  and  $p_{X,U_{1:L}}$  in a natural way, since the same  $F_\ell$  are used in both.

We assume each  $q_{U_\ell|Z_\ell}$  can be suitably reparametrised (Kingma & Welling, 2014; Rezende et al., 2014) so that, for some function  $H_\ell$  and some density  $\eta_\ell$  that does not depend on the parameters of  $q_{U_{1:L}|Z_\ell}$  and  $p_{X,U_{1:L}}$ , we have  $H_\ell(\epsilon_\ell, z_\ell) \sim q_{U_\ell|Z_\ell}(\cdot|z_\ell)$  when  $\epsilon_\ell \sim \eta_\ell$ . We can then obtain unbiased estimates of  $\mathcal{L}(x)$  using Algorithm 1, which corresponds to a single-sample approximation to the expectation in (14). It is straightforward to see that Algorithm 1 has  $\Theta(L)$  complexity. Differentiating through this procedure allows maximising  $\sum_{i=1}^n \mathcal{L}(x_i)$  via stochastic gradient descent. At test time, we can also estimate  $\log p_X(x)$  directly using importance sampling as described by Rezende

et al. (2014, (40)). In particular, letting  $\hat{\mathcal{L}}^{(1)}, \dots, \hat{\mathcal{L}}^{(m)}$  denote the result of separate calls to ELBO( $x$ ), we have

$$\text{LogSumExp}(\hat{\mathcal{L}}^{(1)}, \dots, \hat{\mathcal{L}}^{(m)}) - \log m \rightarrow \log p_X(x) \quad (16)$$

almost surely as  $m \rightarrow \infty$ .

---

**Algorithm 1** Unbiased estimation of  $\mathcal{L}(x)$ 


---

**function** ELBO( $x$ )

$z_L \leftarrow x$

$\Delta \leftarrow 0$

**for**  $\ell = L, \dots, 1$  **do**

$\epsilon \sim \eta_\ell$

$u \leftarrow H_\ell(\epsilon, z_\ell)$

$z_{\ell-1} \leftarrow F_\ell^{-1}(z_\ell; u)$

$\Delta \leftarrow \Delta + \log p_{U_\ell|Z_{\ell-1}}(u|z_{\ell-1}) - \log q_{U_\ell|Z_\ell}(u|z_\ell) + \log |\det DF_\ell^{-1}(z_\ell; u)|$

**end for**

**return**  $\Delta + \log p_{Z_0}(z_0)$

**end function**

---

In all our experiments we used

$$q_{U_\ell|Z_\ell}(\cdot|z_\ell) = \text{Normal}(\mu_\ell^q(z_\ell), \Sigma_\ell^q(z_\ell)) \quad (17)$$

for appropriate neural networks  $\mu_\ell^q$  and  $\Sigma_\ell^q$ , which is immediately reparameterisable as described e.g. by Kingma & Welling (2014). We found this gave good enough performance that we did not require alternatives such as IAF (Kingma et al., 2016), but such options may also be useful.

Finally, Algorithm 1 requires an expression for  $\log |\det DF_\ell^{-1}(z_\ell; u_\ell)|$ . For (9) this is

$$\log \left| \det Df_\ell^{-1} \left( e^{s_\ell(u_\ell)} \odot (z_\ell + t_\ell(u_\ell)) \right) \right| + \sum_{i=1}^d [s_\ell(u_\ell)]_i,$$

where  $[x]_i$  denotes the  $i^{\text{th}}$  dimension of  $x$ .

## 4 Comparison with Related Models

### 4.1 Comparison with Normalising Flows

We now compare CIFs with normalising flows, and in particular describe how CIFs relax the constraints of bijectivity identified in Section 2.

#### 4.1.1 ADVANTAGES

Observe that (7) generalises normalising flows: if  $F(\cdot; u)$  does not depend on  $u$ , then we obtain (1). Moreover, training with the ELBO in this case does not reduce performance compared with training a flow directly, as the following result shows. Here the components of our model  $F_\theta$ ,  $p_{U|Z}^\theta$  and  $q_{U|X}^\theta$  are parameterised by  $\theta \in \Theta$ , and for a given

choice of parameters  $\theta$  we will denote by  $P_X^\theta$  and  $\mathcal{L}^\theta$  the corresponding distribution and ELBO (14) respectively.

**Proposition 4.1.** *Suppose there exists  $\phi \in \Theta$  such that, for some bijection  $f : \mathcal{Z} \rightarrow \mathcal{X}$ ,  $F_\phi(\cdot; u) = f(\cdot)$  for all  $u \in \mathcal{U}$ . Likewise, suppose  $p_{U|Z}^\phi$  and  $q_{U|X}^\phi$  are such that, for some density  $r$  on  $\mathcal{U}$ ,  $p_{U|Z}^\phi(\cdot|z) = q_{U|X}^\phi(\cdot|x) = r(\cdot)$  for all  $z \in \mathcal{Z}$  and  $x \in \mathcal{X}$ . If  $\mathbb{E}_{x \sim P_X^*}[\mathcal{L}^\theta(x)] \geq \mathbb{E}_{x \sim P_X^*}[\mathcal{L}^\phi(x)]$ , then*

$$D_{\text{KL}}(P_X^* \parallel P_X^\theta) \leq D_{\text{KL}}(P_X^* \parallel f\#P_Z).$$

Simply stated, in the limit of infinite data, optimising the ELBO will yield at least as performant a model (as measured by the KL) as any normalising flow our model family can express. The proof is in Section B.7 of the Supplement. In practice, our choices (9), (10), and (17) can easily realise the conditions of Proposition 4.1 by zeroing out the output weights of the neural networks (other than  $f$ ) involved. Thus, for a given  $f$ , we have reason to expect a comparative or better performing model (as measured by average log-likelihood) when trained as a CIF rather than as a normalising flow.

We expect this will in fact lead to *improved* performance because, intuitively,  $P_{U|Z}$  can reroute  $z$  that would otherwise map outside of  $\text{supp } P_X^*$ . To illustrate, fix  $f$  in (9) and choose some  $z \in \mathcal{Z}$ . If  $f(z) \in \text{supp } P_X^*$ , then setting  $F(z; u) = f(z)$  for all  $u \in \mathcal{U}$  as described above ensures  $F(z; U) \in \text{supp } P_X^*$  when  $U \sim P_{U|Z}(\cdot|z)$ . If conversely  $f(z) \notin \text{supp } P_X^*$ , then we *still* have  $F(z; U) \in \text{supp } P_X^*$  almost surely if  $P_{U|Z}(\cdot|z)$  is supported on  $\{u \in \mathcal{U} : F(z; u) \in \text{supp } P_X^*\}$ . Of course, if  $f$  is too simple, then  $P_{U|Z}$  must heuristically become very complex in order to obtain this behaviour. This would seem to make inference harder, leading to a looser ELBO (14) and thus overall worse performance after training. We therefore expect CIFs to work well for  $f$  that, like the 10-layer ResFlow in Figure 1, can learn a close approximation to the support of the target but “leak” some mass outside of it due to (4) or Theorem 2.1. A CIF can then use  $P_{U|Z}$  to “clean up” these small extraneous regions of mass.

We provide empirical support for this argument in Section 5. We also summarise our discussion above with the following precise result. Here  $\partial A$  denotes the boundary of a set  $A$ .

**Proposition 4.2.** *If  $P_X^*(\partial \text{supp } P_X^*) = 0$  and  $(z, u) \mapsto F(z; u)$  is jointly continuous with*

$$\overline{F(\text{supp } P_Z \times \mathcal{U})} \supseteq \text{supp } P_X^*, \quad (18)$$

*then there exists  $P_{U|Z}$  such that  $\text{supp } P_X = \text{supp } P_X^*$  if and only if, for all  $z \in \text{supp } P_Z$ , there exists  $u \in \mathcal{U}$  with*

$$F(z; u) \in \text{supp } P_X^*. \quad (19)$$

The assumptions here are fairly minimal: the boundary condition ensures  $P_X^*$  is not pathological, and if (18) does

not hold, then  $D_{\text{KL}}(P_X^* \parallel P_X) = \infty$  for every  $P_{U|Z}$ .<sup>10</sup> Additionally, the following result gives a sufficient condition under which it is possible to learn the target exactly.

**Proposition 4.3.** *If  $F(z; \cdot) : \mathcal{U} \rightarrow \mathcal{X}$  is surjective for each  $z \in \mathcal{Z}$ , then there exists  $P_{U|Z}$  such that  $P_X = P_X^*$ .*

See Section B.8 of the Supplement for proofs. These results do not require  $\text{supp } P_Z \cong \text{supp } P_X^*$ , thereby showing CIFs relax the constraint (4) for standard normalising flows.

Of course, in practice, our parameterisation (9) does not necessarily ensure that  $F$  will satisfy these conditions, and our parameterisation (10) may not be expressive enough to instantiate the  $P_{U|Z}$  that is required. However, these results show that CIFs provide at least a *mechanism* for correcting a topologically misspecified prior. When  $F$  and  $P_{U|Z}$  are sufficiently expressive, we can expect that they will learn to approximate these conditions over the course of training if doing so produces a better density estimate. We therefore anticipate CIFs will improve performance for ResFlows, where Theorem 2.1 applies, and may have benefits more generally, since all flows are ultimately constrained by (4).

#### 4.1.2 DISADVANTAGES

On the other hand, CIFs introduce additional overhead compared with regular normalising flows. It therefore remains to show we obtain better performance on a fixed computational budget, which requires using a smaller model. Empirically this holds for the models and datasets we consider in Section 5, but there are likely cases where it does not, particularly if the topologies of the target and prior are similar.

Likewise, CIFs sacrifice the exactness of normalising flows. We do not see this as a significant problem for the task of density estimation, since the importance sampling estimator (16) means that at test time we can obtain arbitrary accuracy by taking  $m$  to be large. However, the lack of a closed-form density does limit the use of CIFs in some downstream tasks. In particular, CIFs cannot immediately be plugged in to a variational approximation in the manner of Rezende & Mohamed (2015), since this requires exact likelihoods. However, it may be possible to use CIFs in the context of an extended-space variational framework along the lines of Agakov & Barber (2004), and we leave this for future work.

## 4.2 Comparison with Discretely Indexed Models

Similar models to CIFs have been proposed that use a discrete index space. In the context of Bayesian inference, Duan (2019) proposes a single-layer ( $L = 1$ ) model consisting of (7) with  $\mathcal{U} = \{1, \dots, I\}$  and  $F(\cdot; i) = f_i$  for separate normalising flows  $f_1, \dots, f_I$ . A special case of this framework is given by *deep Gaussian mixture models* (Van den

<sup>10</sup>See Proposition B.1 and Proposition B.3 in the Supplement.

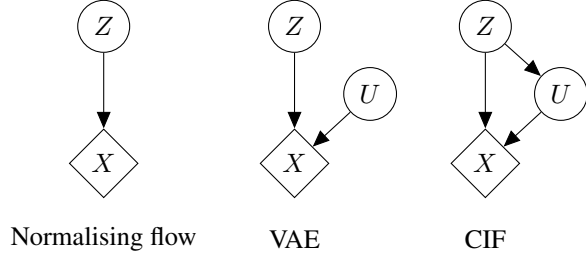


Figure 2: Comparison of related generative models. Circular nodes are random and diamond nodes are deterministic. CIFs generalise both normalising flows and VAEs as shown.

Oord & Schrauwen, 2014; van den Oord & Dambre, 2015), which corresponds to using invertible linear transformations for each  $f_i$ . In this case, (13) becomes a summation that can be computed analytically. However, this quickly becomes intractable as  $L$  grows larger, since the cost to compute this is seen to be  $\Theta(I^L)$ . Unlike for a continuous  $u$ , this cannot easily be reduced to  $\Theta(L)$  using a variational approximation as in Section 3.2, since a discrete  $q_{U|X}$  is not amenable to the reparameterisation trick. In addition, the use of separate bijections also means that the number of parameters of the model grows as  $I$  increases. In contrast, a continuous index allows a natural mechanism for sharing parameters across different  $F(\cdot; u)$  as in (9).

Prior to Duan (2019), Dinh et al. (2019) proposed RAD as a means to mitigate the  $\Theta(I^L)$  cost of naively stacking discrete layers. RAD partitions  $\mathcal{X}$  into  $I$  disjoint subsets  $B_1, \dots, B_I$  and defines bijections  $f_i : \mathcal{Z} \rightarrow B_i$  for each  $i$ . The model is then taken to be the marginal of  $X$  in

$$Z \sim P_Z, \quad U \sim P_{U|Z}(\cdot|Z), \quad X := f_U(Z),$$

where each  $P_{U|Z}(\cdot|z)$  is a discrete distribution on  $\{1, \dots, I\}$ . Note that this is not an instance of our model (7), since we require each  $F(\cdot; u)$  to be surjective onto  $\mathcal{X}$ . The use of partitioning means that (13) is a summation with only a single term, which reduces the cost for  $L$  layers to  $\Theta(L)$ . However, partitioning also makes  $p_X$  discontinuous. This leads to a very difficult optimisation problem and Dinh et al. (2019) only report results for simple 2-D densities. Additionally, partitioning requires ad-hoc architectural changes to existing normalising flows, and does not directly address the increasing parameter cost as  $I$  grows large.

### 4.3 Comparison with Variational Autoencoders

CIFs also generalise a broad family of variational autoencoders (VAEs) (Kingma & Welling, 2014; Rezende et al., 2014). Recall that VAEs take

$$p_X(x) := \int p_U(u) p_{X|U}(x|u) du \quad (20)$$

for some choices of densities  $p_U$  and  $p_{X|U}$ .<sup>11</sup> For instance, a mean-field Gaussian observation density has

$$p_{X|U}(\cdot|u) := \text{Normal} \left( t(u), \text{diag} \left( e^{s(u)} \right) \right),$$

where  $t, s : \mathcal{U} \rightarrow \mathcal{X}$ , and  $\text{diag}(v)$  denotes the matrix with diagonal  $v \in \mathbb{R}^d$  and zeros elsewhere. If  $P_Z$  is a standard Gaussian, if each  $P_{U|Z}(\cdot|z)$  has independent density  $p_U$ , and if  $F$  is (9) with  $f$  the identity, then it follows that (7) has marginal density (20) (modulo the signs of  $s$  and  $t$ ).<sup>12</sup>

More generally, every VAE model (20) with each  $p_{X|U}(\cdot|u)$  strictly positive corresponds to an instance of (7) where  $U$  is sampled independently of  $Z$ . To see this, let  $p_Z$  be any strictly positive density on  $\mathcal{Z}$ , and let each  $F(\cdot; u)$  be the Knothe-Rosenblatt coupling (Villani, 2008) of  $p_Z$  and  $p_{X|U}(\cdot|u)$ . By construction each  $F(\cdot; u)$  is invertible and gives  $F(Z; u) \sim p_{X|U}(\cdot|u)$  when  $Z \sim p_Z$ . As a result, (7) again yields  $X$  with a marginal density defined by (20). Consequently, CIFs generalise the VAE framework by adding an additional edge in the graphical model as shown in Figure 2.

On the other hand, CIFs differ from VAEs in the way they are composed. Whereas CIFs stack by taking  $p_Z$  to be a CIF, VAEs are typically stacked by taking  $p_U$  to be a VAE (Rezende et al., 2014; Kingma et al., 2014; Burda et al., 2016; Sønderby et al., 2016). This has implications for the design of the inference network  $q_{U_{1:L}|X}$ . In particular, a hierarchical VAE obtained in this way is *Markovian*, so that

$$p_{U_{1:L}|X}(x, u_{1:L}) = p_{U_L|X}(u_L|x) \prod_{\ell=1}^L p_{U_\ell|U_{\ell-1}}(u_\ell|u_{\ell-1})$$

where  $L$  is the number of layers. This directly allows specifying  $q_{U_{1:L}|X}$  to be of the same form without any loss of generality (Kingma et al., 2014; Burda et al., 2016; Sønderby et al., 2016). Conversely, CIFs do not factor in this way, which motivates our alternative approach in Section 3.2.

Note finally that CIFs should not be conflated with the large class of methods that use normalising flows to improve the *inference* procedure in VAEs (Rezende & Mohamed, 2015; Kingma et al., 2016; van den Berg et al., 2018). These approaches are orthogonal to ours and indeed may be useful for improving our own inference procedure by replacing (17) with a more expressive model.

### 4.4 Other Related Work

Additional related methods have been proposed. Within a classification context, Dupont et al. (2019) identify topological problems related to ODE-based mappings (Chen et al.,

<sup>11</sup>Note that this notation is nonstandard for VAEs in order to align with the rest of the paper. Here our  $U$  corresponds to  $z$  as used by Kingma & Welling (2014).

<sup>12</sup>Here  $Z$  corresponds to  $\epsilon$  as used by Kingma & Welling (2014).

2018), which like normalising flows are homeomorphisms and hence preserve the topology of their input. To avoid this, Dupont et al. (2019) propose augmenting the data by appending auxiliary dimensions and learning a new mapping on this space. In contrast, CIFs may be understood as augmenting not the data but instead the *model* by considering a family of individual bijections on the *original* space.

In addition, Ho et al. (2019) use a variational scheme to improve on the standard dequantisation method proposed by Theis et al. (2016) for modelling image datasets with normalising flows. This approach is potentially complementary to CIFs, but we do not make use of it in our experiments.

## 5 Experiments

We evaluated the performance of CIFs on several problems of varying difficulty, including synthetic 2-D data, several tabular datasets, and three image datasets. In all cases we took  $\mathcal{Z} = \mathcal{X} = \mathbb{R}^d$  with  $d$  the dimension of the dataset. We used the stacked architecture (8) with the prior  $P_{Z_0}$  a Gaussian. At each layer,  $F$  had form (9) with  $f$  a primitive flow step from a baseline architecture (e.g. a single residual block for ResFlow). Each  $p_{U|Z}$  and  $q_{U|X}$  had form (10) and (17) respectively. We provide an overview of our results for the tabular and image datasets here. Full experimental details, including additional 2-D figures along the lines of Figure 1, are in Section C of the Supplement. See [github.com/jrmcornish/cif](https://github.com/jrmcornish/cif) for our code.

### 5.1 Tabular Datasets

We tested the performance of CIFs on the tabular datasets used by Papamakarios et al. (2017). For each dataset, we trained 10 and 100-layer baseline fully connected ResFlows, and corresponding 10-layer CIF-ResFlows. The CIF-ResFlows had roughly 1.5-4.5% more parameters (depending on the dimension of the dataset) than the otherwise identical 10-layer ResFlows, and roughly 10% of the parameters of the 100-layer ResFlows. Table 1 reports the average log-probability of the test set that we obtained for each model. Observe that in all cases CIF-ResFlows significantly outperform both baseline models. Moreover, for all but GAS, the CIF-ResFlows achieve state-of-the-art performance based on the results reported by Durkan et al. (2019, Table 1). This is particularly noticeable for POWER and BSDS300, where CIF-ResFlow improves on the best results of Durkan et al. (2019) by 0.94 and 2.77 nats respectively.

We additionally tried using *masked autoregressive flows* (MAFs) (Papamakarios et al., 2017) and *neural spline flows* (NSFs) (Durkan et al., 2019) for  $f$ . In each case, we closely match the experimental settings of the baselines and augment using CIFs, controlling for the number of parameters used by the CIF extensions. Table 1 reports the average

log-probability across the test set for each experiment. Here, CIF-NSF-1 is a CIF with the same number of parameters as the baseline, and CIF-NSF-2 is a model using a baseline configuration for  $f$  (but having more parameters overall). We see that CIF-MAFs consistently outperform MAFs across datasets; CIF-NSFs do not improve upon NSFs as dramatically, although we still notice improvements and would expect to improve further with more hyperparameter tuning. Lastly it is important to notice that MAFs and NSFs do not restrict the Lipschitz constant of  $f$ . These results show that CIFs can yield benefits for normalising flows even if Theorem 2.1 is not directly a limitation.

Finally, for ablation purposes we tried taking  $f$  to be the identity. We obtained consistently worse performance than for CIF-ResFlows and CIF-MAF in this case, which aligns with our conjecture in Section 4.1.1 that a performant CIF requires an expressive base flow  $f$ . Details and results are given in Section C.1.4 of the Supplement.

### 5.2 Image Datasets

We also considered CIFs applied to the MNIST (LeCun, 1998), Fashion-MNIST (Xiao et al., 2017), and CIFAR-10 (Krizhevsky & Hinton, 2009) datasets. Following our tabular experiments, we trained a multi-scale convolutional ResFlow and a corresponding CIF-ResFlow, as well as a larger baseline ResFlow to account for the additional parameters and depth introduced by our method. Note that these models were significantly smaller than those used by Chen et al. (2019): e.g. for CIFAR10, the ResFlow used by Chen et al. (2019) had 25M parameters, while our two baseline ResFlows and our CIF-ResFlow had 2.4M, 6.2M, and 5.6M parameters respectively. We likewise considered RealNVPs with the same multi-scale convolutional architecture used by Dinh et al. (2017) for their CIFAR-10 experiments. For these runs we trained baseline RealNVPs, corresponding CIF-RealNVPs, and larger baseline RealNVPs with more depth and parameters.

The results are given in Table 2 and Table 3. Observe CIFs outperformed the baseline models for all datasets, which shows that our approach can scale to high dimensions. For the CIF-ResFlows, we also obtained better performance than Chen et al. (2019) on MNIST and better performance than Glow (Kingma & Dhariwal, 2018) on CIFAR10, despite using a much smaller model. Samples from all models are shown in Section C.2 of the Supplement.

## 6 Conclusion and Future Work

The constraint (4) shows that normalising flows are unable to exactly model targets whose topology differs from that

<sup>13</sup>Only one seed was used per run due to computational limitations. However, the results were not cherry-picked.



Table 1: Mean  $\pm$  standard error (over 3 seeds) of average test set log-likelihood (in nats). Higher is better. Best performing runs for each group are shown in bold. A  $\star$  indicates state-of-the-art performance according to Durkan et al. (2019, Table 1).

	POWER	GAS	HEPMASS	MINIBOONE	BSDS300
RESFLOW ( $L = 10$ )	$-2.73 \pm 0.03$	$4.16 \pm 0.08$	$-20.68 \pm 0.02$	$-14.2 \pm 0.10$	$123.51 \pm 0.09$
RESFLOW ( $L = 100$ )	$0.48 \pm 0.00$	$10.57 \pm 0.17$	$-16.67 \pm 0.05$	$-11.16 \pm 0.04$	$148.05 \pm 0.61$
CIF-RESFLOW ( $L = 10$ )	<b><math>1.60 \pm 0.21^\star</math></b>	<b><math>12.12 \pm 0.10</math></b>	<b><math>-13.74 \pm 0.03^\star</math></b>	<b><math>-8.10 \pm 0.04^\star</math></b>	<b><math>160.50 \pm 0.08^\star</math></b>
MAF	$0.19 \pm 0.02$	$9.23 \pm 0.07$	$-18.33 \pm 0.10$	$-10.98 \pm 0.03$	$156.13 \pm 0.00$
CIF-MAF	<b><math>0.48 \pm 0.01</math></b>	<b><math>12.02 \pm 0.10</math></b>	<b><math>-16.63 \pm 0.09</math></b>	<b><math>-9.93 \pm 0.04</math></b>	<b><math>156.67 \pm 0.02</math></b>
NSF	<b><math>0.69 \pm 0.00</math></b>	$13.01 \pm 0.02$	$-14.30 \pm 0.05$	$-10.68 \pm 0.06$	<b><math>157.59 \pm 0.02</math></b>
CIF-NSF-1	<b><math>0.68 \pm 0.01</math></b>	$12.94 \pm 0.01$	<b><math>-13.83 \pm 0.10</math></b>	<b><math>-9.93 \pm 0.06</math></b>	<b><math>157.60 \pm 0.02</math></b>
CIF-NSF-2	<b><math>0.69 \pm 0.00</math></b>	<b><math>13.08 \pm 0.00</math></b>	$-14.18 \pm 0.09$	$-10.80 \pm 0.01$	$157.56 \pm 0.02$

Table 2: Average test bits per dimension.<sup>13</sup> Lower is better.

	MNIST	CIFAR-10
RESFLOW (SMALL)	1.074	3.474
RESFLOW (BIG)	1.018	3.422
CIF-RESFLOW	<b>0.922</b>	<b>3.334</b>

Table 3: Mean  $\pm$  standard error of average test set bits per dimension over 3 random seeds. Lower is better.

	FASHION-MNIST	CIFAR-10
REALNVP (SMALL)	$2.944 \pm 0.003$	$3.565 \pm 0.001$
REALNVP (BIG)	$2.946 \pm 0.002$	$3.554 \pm 0.001$
CIF-REALNVP	<b><math>2.823 \pm 0.003</math></b>	<b><math>3.477 \pm 0.019</math></b>

of the prior. Moreover, in order to approximate such targets closely, Theorem 2.1 shows that the bi-Lipschitz constant of a flow must become arbitrarily large. To address these problems, we have proposed CIFs, which can “clean up” regions of mass that are placed outside the support of the target by a standard flow. CIFs perform well in practice and outperform baseline flows on several benchmark datasets.

While we have focussed on the use of CIFs for density estimation in this paper, it would also be interesting to apply CIFs in other contexts where normalising flows have been used successfully. As CIFs do not have an analytically available density, this would likely require the modification of existing numerical frameworks, but the expressiveness benefits provided by CIFs might make this additional effort worthwhile. We leave this direction for future work.

## Acknowledgements

Rob Cornish is supported by the EPSRC Centre for Doctoral Training in Autonomous Intelligent Machines & Systems (EP/L015897/1) and NVIDIA. Anthony Caterini is a Commonwealth Scholar supported by the U.K. Government. Arnaud Doucet is partially supported by the U.S. Army Re-

search Laboratory, the U.S. Army Research Office, and by the U.K. Ministry of Defence (MoD) grant EP/R013616/1 and the U.K. EPSRC under grant numbers EP/R034710/1 and EP/R018561/1.

## References

- Agakov, F. V. and Barber, D. An auxiliary variational method. In *International Conference on Neural Information Processing*, pp. 561–566. Springer, 2004.
- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pp. 214–223, 2017.
- Bache, K. and Lichman, M. UCI machine learning repository, 2013. URL <http://archive.ics.uci.edu/ml>.
- Beaton, A. and Adams, R. P. Efficient optimization of loops and limits with randomized telescoping sums. In *International Conference on Machine Learning*, pp. 534–543, 2019.
- Behrmann, J., Grathwohl, W., Chen, R. T., Duvenaud, D., and Jacobsen, J.-H. Invertible residual networks. In *International Conference on Machine Learning*, pp. 573–582, 2019.
- Behrmann, J., Vicol, P., Wang, K.-C., Grosse, R. B., and Jacobsen, J.-H. On the invertibility of invertible neural networks, 2020. URL <https://openreview.net/forum?id=BJlVeyHFwH>.
- Billingsley, P. *Probability and Measure*. John Wiley & Sons, 2008.
- Burda, Y., Grosse, R., and Salakhutdinov, R. Importance weighted autoencoders. In *ICLR*, 2016.
- Chen, T. Q., Rubanova, Y., Bettencourt, J., and Duvenaud, D. K. Neural ordinary differential equations. In *Advances in Neural Information Processing Systems*, pp. 6571–6583, 2018.

- Chen, T. Q., Behrmann, J., Duvenaud, D. K., and Jacobsen, J.-H. Residual flows for invertible generative modeling. In *Advances in Neural Information Processing Systems*, pp. 9913–9923, 2019.
- Dinh, L., Krueger, D., and Bengio, Y. Nice: Non-linear independent components estimation. In *ICLR Workshop*, 2015.
- Dinh, L., Sohl-Dickstein, J., and Bengio, S. Density estimation using real NVP. In *ICLR*, 2017.
- Dinh, L., Sohl-Dickstein, J., Pascanu, R., and Larochelle, H. A RAD approach to deep mixture models. In *ICLR Workshop*, 2019.
- Duan, L. L. Transport Monte Carlo. *arXiv preprint arXiv:1907.10448*, 2019.
- Dudley, R. M. *Real Analysis and Probability*. Cambridge Studies in Advanced Mathematics. Cambridge University Press, 2 edition, 2002.
- Dupont, E., Doucet, A., and Teh, Y. W. Augmented neural ODEs. In *Advances in Neural Information Processing Systems*, pp. 3134–3144, 2019.
- Durkan, C., Bekasov, A., Murray, I., and Papamakarios, G. Neural spline flows. In *Advances in Neural Information Processing Systems*, pp. 7509–7520, 2019.
- Germain, M., Gregor, K., Murray, I., and Larochelle, H. Made: Masked autoencoder for distribution estimation. In *International Conference on Machine Learning*, pp. 881–889, 2015.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pp. 2672–2680, 2014.
- Gouk, H., Frank, E., Pfahringer, B., and Cree, M. Regularisation of neural networks by enforcing lipschitz continuity. *arXiv preprint arXiv:1804.04368*, 2018.
- Grathwohl, W., Chen, R. T., Betterncourt, J., Sutskever, I., and Duvenaud, D. FFJORD: Free-form continuous dynamics for scalable reversible generative models. In *ICLR*, 2019.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016a.
- He, K., Zhang, X., Ren, S., and Sun, J. Identity mappings in deep residual networks. In *European Conference on Computer Vision*, pp. 630–645. Springer, 2016b.
- Ho, J., Chen, X., Srinivas, A., Duan, Y., and Abbeel, P. Flow++: Improving flow-based generative models with variational dequantization and architecture design. In *International Conference on Machine Learning*, pp. 2722–2730, 2019.
- Huang, C.-W., Krueger, D., Lacoste, A., and Courville, A. Neural autoregressive flows. In *International Conference on Machine Learning*, pp. 2083–2092, 2018.
- Hutchinson, M. F. A stochastic estimator of the trace of the influence matrix for laplacian smoothing splines. *Communications in Statistics-Simulation and Computation*, 19(2):433–450, 1990.
- Jaini, P., Selby, K. A., and Yu, Y. Sum-of-squares polynomial flow. In *International Conference on Machine Learning*, pp. 3009–3018, 2019.
- Kahn, H. Use of different Monte Carlo sampling techniques. Technical report, Rand Corporation, 1955.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- Kingma, D. P. and Dhariwal, P. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*, pp. 10215–10224, 2018.
- Kingma, D. P. and Welling, M. Auto-encoding variational Bayes. In *ICLR*, 2014.
- Kingma, D. P., Mohamed, S., Rezende, D. J., and Welling, M. Semi-supervised learning with deep generative models. In *Advances in Neural Information Processing Systems*, pp. 3581–3589, 2014.
- Kingma, D. P., Salimans, T., Jozefowicz, R., Chen, X., Sutskever, I., and Welling, M. Improved variational inference with inverse autoregressive flow. In *Advances in Neural Information Processing Systems*, pp. 4743–4751, 2016.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- Kumar, A., Poole, B., and Murphy, K. Learning generative samplers using relaxed injective flow. In *ICML Workshop on Invertible Neural Nets and Normalizing Flows*, 2019.
- LeCun, Y. The MNIST database of handwritten digits, 1998. URL <http://yann.lecun.com/exdb/mnist/>.
- Lyne, A.-M., Girolami, M., Atchadé, Y., Strathmann, H., Simpson, D., et al. On Russian roulette estimates for Bayesian inference with doubly-intractable likelihoods. *Statistical science*, 30(4):443–467, 2015.

- Martin, D., Fowlkes, C., Tal, D., and Malik, J. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, volume 2, pp. 416–423. IEEE, 2001.
- Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. Spectral normalization for generative adversarial networks. In *ICLR*, 2018.
- Papamakarios, G., Pavlakou, T., and Murray, I. Masked autoregressive flow for density estimation. In *Advances in Neural Information Processing Systems*, pp. 2338–2347, 2017.
- Rezende, D. and Mohamed, S. Variational inference with normalizing flows. In *International Conference on Machine Learning*, pp. 1530–1538, 2015.
- Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, pp. 1278–1286, 2014.
- Rhee, C.-h. and Glynn, P. W. Unbiased estimation with square root convergence for SDE models. *Operations Research*, 63(5):1026–1043, 2015.
- Rudin, W. *Principles of Mathematical Analysis*, volume 3. McGraw-hill New York, 1964.
- Rudin, W. *Real and Complex Analysis*. Tata McGraw-hill education, 2006.
- Runde, V. *A Taste of Topology*. Springer, 2007.
- Skilling, J. The eigenvalues of mega-dimensional matrices. In *Maximum Entropy and Bayesian Methods*, pp. 455–466. Springer, 1989.
- Sønderby, C. K., Raiko, T., Maaløe, L., Sønderby, S. K., and Winther, O. Ladder variational autoencoders. In *Advances in Neural Information Processing Systems*, pp. 3738–3746, 2016.
- Theis, L., Oord, A. v. d., and Bethge, M. A note on the evaluation of generative models. In *ICLR*, 2016.
- van den Berg, R., Hasenclever, L., Tomczak, J. M., and Welling, M. Sylvester normalizing flows for variational inference. In *UAI 2018: The Conference on Uncertainty in Artificial Intelligence (UAI)*, pp. 393–402, 2018.
- van den Oord, A. and Dambre, J. Locally-connected transformations for deep GMMs. In *International Conference on Machine Learning (ICML): Deep learning Workshop*, pp. 1–8, 2015.
- Van den Oord, A. and Schrauwen, B. Factoring variations in natural images with deep Gaussian mixture models. In *Advances in Neural Information Processing Systems*, pp. 3518–3526, 2014.
- Villani, C. *Optimal Transport: Old and New*, volume 338. Springer Science & Business Media, 2008.
- Xiao, H., Rasul, K., and Vollgraf, R. Fashion-MNIST: A novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

---

# Relaxing Bijectivity Constraints with Continuously Indexed Normalising Flows: Supplementary Material

---

## A Guide to Notation

$(a_n)$	A sequence of elements $a_1, a_2, \dots$
$a(n) = \Theta(b(n))$	$a(n)$ differs from $b(n)$ by at most a constant factor as $n \rightarrow \infty$
$u \odot v$	The elementwise product of tensors $u$ and $v$
$\text{LogSumExp}(a_1, \dots, a_m)$	$\log(\sum_{i=1}^m \exp(a_i))$
$e^v$ , where $v \in \mathbb{R}^d$	$(e^{v_1}, \dots, e^{v_d})$
$\ v\ $	The norm of a vector $v \in \mathbb{R}^d$ (our results are agnostic to the specific choice of $\ \cdot\ $ )
$\ A\ _{\text{op}}$	The operator norm of a matrix $A \in \mathbb{R}^{d_1 \times d_2}$ induced by $\ \cdot\ $
$I_d$	The $d \times d$ identity matrix
$\det A$	The determinant of a square matrix $A$
$Df(z)$	The Jacobian matrix of a function $f$ evaluated at $z$
$DF(z; u)$	The Jacobian matrix of a function $DF(\cdot; u)$ (i.e. with $u$ fixed) evaluated at $z$
$\text{Lip } f$	The Lipschitz constant of a function $f$
$\text{BiLip } f$	The bi-Lipschitz constant of a function $f$
$\mathcal{A} \cong \mathcal{B}$	The topological spaces $\mathcal{A}$ and $\mathcal{B}$ are homeomorphic
$\overline{B}$	The topological closure of a set $B$
$\text{int}(B)$	The interior of a set $B$
$\partial B$	The boundary of a set $B$
$\text{supp } \mu$	The support of a measure $\mu$
$f\#\mu$	The pushforward of a measure $\mu$ by a function $f$
$\mu_n \xrightarrow{\mathcal{D}} \mu$	Weak convergence of the measures $\mu_n$ to $\mu$

## B Proofs

### B.1 Preliminaries

We require some basic results that we include here for completeness. We will make use of standard definitions and results from topology and real analysis. A complete background to these topics can be found in [Dudley \(2002\)](#).

#### B.1.1 SUPPORTS OF MEASURES

Recall that for a Borel measure  $\mu$  on a topological space  $\mathcal{Z}$ , the *support* of  $\mu$ , denoted  $\text{supp } \mu$ , is the set of all  $z \in \mathcal{Z}$  such that  $\mu(N_z) > 0$  for every open set  $N_z$  containing  $z$ .

The following is an immediate consequence:

**Proposition B.1.** *Suppose  $\mu$  and  $\nu$  are Borel measures with  $\mu$  absolutely continuous with respect to  $\nu$ . Then*

$$\text{supp } \mu \subseteq \text{supp } \nu.$$

*Proof.* Suppose  $z \notin \text{supp } \nu$ . Then there exists an open set  $N_z$  containing  $z$  such that  $\nu(N_z) = 0$ . By absolute continuity, we have also that  $\mu(N_z) = 0$  and hence  $z \notin \text{supp } \mu$ .  $\square$

In general the converse need not hold. For example, the Dirac measure on 0 has support contained within the Lebesgue measure on  $\mathbb{R}$  (which has full support), but is not absolutely continuous with respect to it.

The following characterisation is useful:

**Proposition B.2.** For any Borel measure  $\mu$ ,

$$(\text{supp } \mu)^c = \bigcup_{\substack{A \text{ open:} \\ \mu(A)=0}} A, \tag{B.1}$$

and hence  $\text{supp } \mu$  is closed.

*Proof.* This follows directly from the definitions, since  $z \notin \text{supp } \mu$  if and only if there exists open  $N_z$  with  $z \in N_z$  and  $\mu(N_z) = 0$ , which is just another way of saying that  $z$  is contained in the right-hand side of (B.1). It follows that  $(\text{supp } \mu)^c$  is open, and hence  $\text{supp } \mu$  is closed.  $\square$

We mainly care about how the support of a measure is transformed by a pushforward function. The following proposition characterises what occurs in this case.

**Proposition B.3.** Suppose  $\mathcal{Z}$  and  $\mathcal{X}$  are topological spaces. If  $\mu$  is a Borel measure on  $\mathcal{Z}$  such that  $\mu((\text{supp } \mu)^c) = 0$ , and if  $f : \mathcal{Z} \rightarrow \mathcal{X}$  is continuous, then

$$\text{supp } f\#\mu = \overline{f(\text{supp } \mu)}.$$

*Proof.* Suppose  $x \notin \overline{f(\text{supp } \mu)}$ . Then  $x$  must have an open neighbourhood  $N_x$  such that

$$N_x \cap f(\text{supp } \mu) = \emptyset.$$

This implies

$$\begin{aligned} f^{-1}(N_x) \cap \text{supp } \mu &\subseteq f^{-1}(N_x) \cap f^{-1}(f(\text{supp } \mu)) \\ &= f^{-1}(N_x \cap f(\text{supp } \mu)) \\ &= f^{-1}(\emptyset) \\ &= \emptyset. \end{aligned}$$

We then have

$$f\#\mu(N_x) = \mu(f^{-1}(N_x)) = \mu(f^{-1}(N_x) \cap \text{supp } \mu) = 0,$$

where the second equality follows since we assumed  $\mu((\text{supp } \mu)^c) = 0$ , and hence  $x \notin \text{supp } f\#\mu$ . Consequently

$$\text{supp } f\#\mu \subseteq \overline{f(\text{supp } \mu)}.$$

In the other direction, suppose  $x \in \overline{f(\text{supp } \mu)}$ , so that  $x = f(z)$  for some  $z \in \text{supp } \mu$ . Given an open neighbourhood  $N_x$  it then follows from continuity that  $f^{-1}(N_x)$  is an open neighbourhood of  $z$ , and so

$$f\#\mu(N_x) = \mu(f^{-1}(N_x)) > 0$$

since  $z \in \text{supp } \mu$ . This entails  $\text{supp } f\#\mu \supseteq f(\text{supp } \mu)$ , which means

$$\text{supp } f\#\mu = \overline{\text{supp } f\#\mu} \supseteq \overline{f(\text{supp } \mu)}$$

by Proposition B.2.  $\square$

Note that in general we need not have  $\text{supp } f\#\mu = f(\text{supp } \mu)$ . For example, if  $\mu$  is Gaussian and  $f = \arctan$ , then

$$f(\text{supp } \mu) = (-1, 1) \neq [-1, 1] = \text{supp } f\#\mu.$$

Likewise, in general we do require the assumption  $\mu((\text{supp } \mu)^c) = 0$ . This is because there exist examples of nontrivial Borel measures  $\mu$  such that  $\text{supp } \mu = \emptyset$ . Taking  $f \equiv x_0$  to be any constant  $x_0 \in \mathcal{X}$  (in which case  $f$  is certainly continuous) then gives

$$\overline{f(\text{supp } \mu)} = \emptyset \neq \{x_0\} = \text{supp } f\#\mu.$$

However, for our purposes, the following proposition shows that this is not a restriction.

**Proposition B.4.** *Suppose  $\mu$  is a Borel measure on a separable metric space  $\mathcal{Z}$ . Then*

$$\mu((\text{supp } \mu)^c) = 0.$$

*Proof.* Throughout the proof, for each  $z$  and  $r > 0$ , we will denote by  $B(z, r)$  an open ball of radius  $r$  centered at  $z$ . Likewise, for each  $z \notin \text{supp } \mu$ , let

$$r^*(z) := \sup\{r > 0 \mid \mu(B(z, r)) = 0\}.$$

Observe that  $r^*$  is well-defined (but possibly infinite) since  $z \notin \text{supp } \mu$  means there must exist some  $r > 0$  such that  $\mu(B(z, r)) = 0$ .

We first show that  $\mu(B(z, r^*(z))) = 0$  for all  $z \notin \text{supp } \mu$ . To this end, fix  $z$  and choose a sequence  $r_m \uparrow r^*(z)$  with  $r_m < r^*(z)$ . We then have

$$B(z, r^*(z)) = \bigcup_{m=1}^{\infty} B(z, r_m),$$

and so

$$\mu(B(z, r^*(z))) = \lim_{m \rightarrow \infty} \mu(B(z, r_m)) = 0$$

by continuity of measure.

Now, by separability, we can choose a countable sequence  $(z_k) \subseteq (\text{supp } \mu)^c$  such that  $\overline{\{z_k\}} = \overline{(\text{supp } \mu)^c}$ . We show that

$$(\text{supp } \mu)^c = \bigcup_{k=1}^{\infty} B(z_k, r^*(z_k)),$$

from which the result follows by countable subadditivity. It is clear from (B.1) that the left-hand side is a superset of the right. In the other direction, let  $z \in (\text{supp } \mu)^c$ . By construction of  $(z_k)$ , there exists a subsequence  $(z_{k'})$  such that  $z_{k'} \rightarrow z$ . For all  $k'$  large enough we then have  $z_{k'} \in B(z, r^*(z)/2)$  and hence

$$B(z_{k'}, r^*(z)/2) \subseteq B(z, r^*(z))$$

by triangle inequality. It follows that for such  $k'$  we have

$$\mu(B(z_{k'}, r^*(z)/2)) \leq \mu(B(z, r^*(z))) = 0,$$

and so  $r^*(z_{k'}) \geq r^*(z)/2$  since  $r^*(z_{k'})$  is the supremum. But then we have

$$z \in B(z_{k'}, r^*(z)/2) \subseteq B(z_{k'}, r^*(z_{k'})),$$

so that

$$z \in \bigcup_{k=1}^{\infty} B(z_k, r^*(z_k))$$

and we are done. □

## B.2 Lipschitz and Bi-Lipschitz Functions

We assume that  $\mathcal{Z} \subseteq \mathbb{R}^{d_z}$ ,  $\mathcal{X} \subseteq \mathbb{R}^{d_x}$ , and  $f : \mathcal{Z} \rightarrow \mathcal{X}$ . Recall that the *Lipschitz* constant of  $f$ , denoted  $\text{Lip } f$ , is defined as the infimum over  $M \in [0, \infty]$  such that

$$\|f(z) - f(z')\| \leq M\|z - z'\|$$

for all  $z, z' \in \mathcal{Z}$ . Likewise the *bi-Lipschitz* constant  $\text{BiLip } f$  is defined as the infimum over  $M \in [1, \infty]$  such that

$$M^{-1}\|z - z'\| \leq \|f(z) - f(z')\| \leq M\|z - z'\|$$

for all  $z, z' \in \mathcal{Z}$ . We prove some basic properties that follow from this definition.

**Proposition B.5.** *BiLip  $f < \infty$  if and only if  $f$  is injective and  $\max(\text{Lip } f, \text{Lip } f^{-1}) < \infty$ , where  $f^{-1} : f(\mathcal{Z}) \rightarrow \mathcal{Z}$ . For all injective  $f$ , we then have  $\text{BiLip } f = \max(\text{Lip } f, \text{Lip } f^{-1})$ .*

*Proof.* For the first statement, suppose  $\text{BiLip } f < \infty$ . It is immediate that  $\text{BiLip } f \geq \text{Lip } f$ . To see that  $f$  is injective, note that for  $z \neq z'$  we have

$$\|f(z) - f(z')\| \geq (\text{BiLip } f)^{-1} \|z - z'\| > 0$$

and so  $f(z) \neq f(z')$ . On the other hand, for  $x, x' \in f(\mathcal{Z})$ , we have

$$(\text{BiLip } f)^{-1} \|f^{-1}(x) - f^{-1}(x')\| \leq \|f(f^{-1}(x)) - f(f^{-1}(x'))\| = \|x - x'\|,$$

which gives that  $\text{BiLip } f \geq \text{Lip } f^{-1}$ . Altogether we have

$$\max(\text{Lip } f, \text{Lip } f^{-1}) \leq \text{BiLip } f < \infty, \tag{B.2}$$

which gives the forward direction.

Next suppose  $f$  is injective and that

$$M := \max(\text{Lip } f, \text{Lip } f^{-1}) < \infty.$$

For  $z, z' \in \mathcal{Z}$ , we certainly have

$$\|f(z) - f(z')\| \leq M \|z - z'\|.$$

Likewise, since  $f(z), f(z') \in f(\mathcal{Z})$ ,

$$\|z - z'\| = \|f^{-1}(f(z)) - f^{-1}(f(z'))\| \leq M \|f(z) - f(z')\|,$$

so that

$$M^{-1} \|z - z'\| \leq \|f(z) - f(z')\|$$

because injectivity of  $f$  means that  $M > 0$ . From this it follows that

$$\text{BiLip } f \leq M < \infty, \tag{B.3}$$

which gives the reverse direction, proving the first statement.

For the second statement, suppose  $f$  is injective. Then if  $\text{BiLip } f < \infty$ , (B.2) and (B.3) together give

$$\text{BiLip } f = \max(\text{Lip } f, \text{Lip } f^{-1}).$$

On the other hand, if  $\text{BiLip } f = \infty$  then  $\max(\text{Lip } f, \text{Lip } f^{-1}) = \infty$  since we would otherwise obtain a contradiction by the first statement of the proposition. This completes the proof.  $\square$

It follows directly that if  $\text{BiLip } f < \infty$ , then  $f$  is a homeomorphism from  $\mathcal{Z}$  to  $f(\mathcal{Z})$ .<sup>14</sup> Moreover, in this case  $f$  maps closed sets to closed sets, as the following result shows:

**Proposition B.6.** *If  $\text{BiLip } f < \infty$  and  $\mathcal{Z}$  is closed in  $\mathbb{R}^{d_{\mathcal{Z}}}$ , then  $f(\mathcal{Z})$  is closed in  $\mathbb{R}^{d_x}$ .*

*Proof.* It is a straightforward consequence of Proposition B.5 that if  $(x_n) \subseteq f(\mathcal{Z})$  is Cauchy, then  $(f^{-1}(x_n))$  is Cauchy. Consequently  $(f^{-1}(x_n))$  converges to some  $z_\infty \in \mathcal{Z}$ , since  $\mathcal{Z}$  is a closed subset of a complete space and therefore complete. But then

$$\begin{aligned} \|x_n - f(z_\infty)\| &= \|f(f^{-1}(x_n)) - f(z_\infty)\| \\ &\leq M \|f^{-1}(x_n) - z_\infty\| \\ &\rightarrow 0 \end{aligned}$$

as  $n \rightarrow \infty$ . Consequently  $f(\mathcal{Z})$  is complete, and so  $f(\mathcal{Z})$  is closed as desired since the ambient space  $\mathbb{R}^{d_x}$  is complete.  $\square$

The Lipschitz constant can be computed from the operator norm  $\|\cdot\|_{\text{op}}$  of the Jacobian of  $f$ . Recall that  $\|\cdot\|_{\text{op}}$  is defined as for a matrix  $A \in \mathbb{R}^{d_x \times d_{\mathcal{Z}}}$  as

$$\|A\|_{\text{op}} := \sup_{\substack{v \in \mathbb{R}^{d_{\mathcal{Z}}}: \\ \|v\|=1}} \|Av\|$$

where we think of elements of  $\mathbb{R}^{d_{\mathcal{Z}}}$  as column vectors.

<sup>14</sup>Note however that the converse is not true in general: for example,  $\exp$  is a homeomorphism from  $\mathbb{R}$  to  $(0, \infty)$ , but  $\text{BiLip } \exp = \infty$ .

**Proposition B.7.** *If  $\mathcal{Z} = \mathbb{R}^{d_{\mathcal{Z}}}$ ,  $\mathcal{X} = \mathbb{R}^{d_{\mathcal{X}}}$ , and  $f$  is everywhere differentiable, then*

$$\text{Lip } f = \sup_{z \in \mathcal{Z}} \|\text{D}f(z)\|_{\text{op}}.$$

*Proof.* If  $v \in \mathcal{Z}$  with  $\|v\| = 1$ , then

$$\begin{aligned} \|[\text{D}f(z)]v\| &= \lim_{t \rightarrow 0} \frac{\|f(z + tv) - f(z)\|}{|t|} \\ &\leq \lim_{t \rightarrow 0} \frac{(\text{Lip } f)\|(z + tv) - z\|}{|t|} \\ &= \text{Lip } f. \end{aligned}$$

It follows directly that

$$\|\text{D}f(z)\|_{\text{op}} \leq \text{Lip } f.$$

On the other hand, suppose  $\text{Lip } f > M$ . Then there exists  $z, z' \in \mathcal{Z}$  such that

$$\|f(z) - f(z')\| > M\|z - z'\|.$$

Since  $f$  is differentiable, so too is the map  $\varphi : [0, 1] \rightarrow \mathcal{X}$  defined by

$$\varphi(t) := f(tz + (1 - t)z').$$

By Theorem 5.19 of Rudin (1964), there exists  $t_0 \in (0, 1)$  such that the derivative  $\varphi'$  satisfies

$$\|\varphi'(t_0)\| \geq \|f(z') - f(z)\| > M\|z - z'\|.$$

But, letting  $z_0 := t_0z + (1 - t_0)z'$ , observe that

$$\begin{aligned} \varphi'(t_0) &= \lim_{t \rightarrow 0} \frac{f(z_0 + t(z - z')) - f(z_0)}{t} \\ &= [\text{D}f(z_0)](z - z'), \end{aligned}$$

where we think of  $z, z'$  as column vectors. As such,

$$\begin{aligned} \|\text{D}f(z_0)\|_{\text{op}}\|z - z'\| &\geq \|[\text{D}f(z_0)](z - z')\| \\ &= \|\varphi'(t_0)\| \\ &> M\|z - z'\| \end{aligned}$$

and so

$$\sup_{z \in \mathcal{Z}} \|\text{D}f(z)\|_{\text{op}} > M.$$

Since  $M$  was arbitrary this means that

$$\text{Lip } f \leq \sup_{z \in \mathcal{Z}} \|\text{D}f(z)\|_{\text{op}}$$

which gives the result. □

**Proposition B.5** and **Proposition B.7** then immediately entail the following:

**Corollary B.8.** *Suppose  $\mathcal{Z} = \mathbb{R}^{d_{\mathcal{Z}}}$  and  $\mathcal{X} = \mathbb{R}^{d_{\mathcal{X}}}$ . If  $f$  is injective, and if  $f$  and  $f^{-1} : f(\mathcal{Z}) \rightarrow \mathcal{Z}$  are everywhere differentiable, then*

$$\text{BiLip } f = \max \left( \sup_{z \in \mathcal{Z}} \|\text{D}f(z)\|_{\text{op}}, \sup_{x \in f(\mathcal{Z})} \|\text{D}f^{-1}(x)\|_{\text{op}} \right).$$



## B.2.1 ARZELÀ-ASCOLI

Our proof of [Theorem 2.1](#) makes use of the Arzelà-Ascoli theorem. This is a standard and foundational result in analysis, but we include a statement here for completeness. To this end, suppose we have a sequence of functions  $f_n : \mathcal{Z} \subseteq \mathbb{R}^{d_x} \rightarrow \mathcal{X} \subseteq \mathbb{R}^{d_x}$ . We say that  $(f_n)$  is *pointwise bounded* if, for all  $z \in \mathcal{Z}$ ,

$$\sup_n \|f_n(z)\| < \infty.$$

Likewise,  $(f_n)$  is *uniformly equicontinuous* if for every  $\epsilon > 0$  there exists  $\delta > 0$  such that, for all  $n$ ,

$$\|f_n(z) - f_n(z')\| < \epsilon$$

whenever  $\|z - z'\| < \delta$ .

**Theorem B.9** (Arzelà-Ascoli). *If a sequence of functions  $f_n : \mathcal{Z} \subseteq \mathbb{R}^{d_z} \rightarrow \mathcal{X} \subseteq \mathbb{R}^{d_x}$  is pointwise bounded and uniformly equicontinuous, then there exists a subsequence of  $(f_n)$  that converges uniformly on every compact subset of  $\mathcal{Z}$ .*

*Proof.* The case  $d = 1$  is proven for example by [Rudin \(2006, Theorem 11.28\)](#). This can be extended to the case  $d > 1$  by a standard argument. In particular, write

$$f_n =: (f_{n,1}, \dots, f_{n,d}),$$

where  $f_{n,i} : \mathcal{Z} \rightarrow \mathbb{R}$ . Then extract a subsequence  $(f_{n_1})$  of  $(f_n)$  such that  $f_{n_1,1}$  converges uniformly on every compact subset of  $\mathcal{Z}$ . Then extract a subsequence of  $(f_{n_1})$  such that the same holds for  $f_{n_1,2}$ , and so on. The result is a subsequence  $(f_{n'})$  such that each  $f_{n',i}$  converges uniformly on compact subsets of  $\mathcal{Z}$ , from which the same holds for  $f_{n'}$  also by the triangle inequality.  $\square$

## B.3 Pushforward Maps Require Unbounded Bi-Lipschitz Constants

**Theorem 2.1.** *Suppose  $P_Z$  and  $P_X^*$  are probability measures on  $\mathbb{R}^{d_z}$  and  $\mathbb{R}^{d_x}$  respectively, and that  $\text{supp } P_Z \not\cong \text{supp } P_X^*$ . Then for any sequence of measurable  $f_n : \mathbb{R}^{d_z} \rightarrow \mathbb{R}^{d_x}$ , we can have  $f_n \# P_Z \xrightarrow{\mathcal{D}} P_X^*$  only if*

$$\lim_{n \rightarrow \infty} \text{BiLip } f_n = \infty.$$

*Proof.* We suppose that  $f_n \# P_Z \xrightarrow{\mathcal{D}} P_X^*$  and prove the contrapositive. That is, without loss of generality (pass to a subsequence if necessary) we assume

$$M := \sup_n \text{BiLip } f_n < \infty, \tag{B.4}$$

and prove that  $\text{supp } P_Z \cong \text{supp } P_X^*$ .

We first show that  $(f_n)$  is pointwise bounded. To this end, observe that Prokhorov's theorem ([Dudley, 2002, Proposition 9.3.4](#)) means that  $P_Z$  is tight and that the sequence  $(f_n \# P_Z)$  is uniformly tight. As such, there exists compact  $K \subseteq \mathbb{R}^{d_z}$  such that  $P_Z(K) > 0$ , and compact  $K' \subseteq \mathbb{R}^{d_x}$  such that

$$\inf_n f_n \# P_Z(K') > 1 - P_Z(K).$$

For each  $n$ , we must then have some  $z_n \in K$  such that  $f_n(z_n) \in K'$ ; otherwise  $K' \subseteq f_n(K)^c$  and so

$$\begin{aligned} f_n \# P_Z(K') &\leq f_n \# P_Z(f_n(K)^c) \\ &= 1 - f_n \# P_Z(f_n(K)) \\ &= 1 - P_Z(f_n^{-1}(f_n(K))) \\ &= 1 - P_Z(K) \end{aligned}$$

since  $f_n$  is injective by [Proposition B.5](#). But for any fixed  $z \in \mathbb{R}^{d_Z}$ , this entails

$$\begin{aligned} \sup_n \|f_n(z)\| &\leq \sup_n \|f_n(z_n)\| + \|f_n(z) - f_n(z_n)\| \\ &\leq \sup_{x \in K'} \|x\| + \sup_{z \in K} M \|z - z_n\| \\ &\leq \sup_{x \in K'} \|x\| + 2M \sup_{z \in K} \|z\| \\ &< \infty \end{aligned}$$

since  $K$  and  $K'$  are compact.

Next, observe that [\(B.4\)](#) easily means  $(f_n)$  is uniformly equicontinuous. In particular, for  $\epsilon > 0$ , choosing  $\delta := \epsilon/M$  gives

$$\|f_n(z) - f_n(z')\| \leq M \|z - z'\| < \epsilon$$

for all  $n$  whenever  $\|z - z'\| < \delta$

[Theorem B.9](#) now entails the existence of a subsequence  $(f_{n'})$  that converges uniformly on every compact subset of  $\mathbb{R}^{d_Z}$ . In particular,  $(f_{n'})$  converges pointwise to a limit that we denote by  $f_\infty$ . Moreover,  $f_\infty$  is bi-Lipschitz. To see this, recall that for all  $n'$  and  $z, z' \in \mathbb{R}^{d_Z}$  we have

$$\frac{1}{M} \|z - z'\| \leq \|f_{n'}(z) - f_{n'}(z')\| \leq M \|z - z'\|,$$

by our assumption [\(B.4\)](#). Taking  $n' \rightarrow \infty$  shows that  $\text{BiLip } f_\infty \leq M < \infty$ .

We also have that

$$f_{n'} \# P_Z \xrightarrow{\mathcal{D}} f_\infty \# P_Z. \tag{B.5}$$

This follows from the Portmanteau theorem ([Dudley, 2002](#), Theorem 11.3.3). In particular, suppose  $h$  is a bounded Lipschitz function, and let  $B_r \subseteq \mathbb{R}^{d_Z}$  denote a ball of radius  $r > 0$  at the origin. Then

$$\begin{aligned} \left| \int h(x) f_{n'} \# P_Z(dx) - \int h(x) f_\infty \# P_Z(dx) \right| &= \left| \int h(f_{n'}(z)) - h(f_\infty(z)) P_Z(dz) \right| \\ &\leq \int_{B_r} |h(f_{n'}(z)) - h(f_\infty(z))| P_Z(dz) \\ &\quad + \int_{B_r^c} |h(f_{n'}(z))| + |h(f_\infty(z))| P_Z(dz) \\ &\leq P_Z(B_r) (\text{Lip } h) \sup_{z \in B_r} \|f_n(z) - f_\infty(z)\| + 2P_Z(B_r^c) \sup_{z \in \mathbb{R}^{d_Z}} |h(z)|. \end{aligned}$$

Hence

$$\limsup_{n' \rightarrow \infty} \left| \int h(x) f_{n'} \# P_Z(dx) - \int h(x) f_\infty \# P_Z(dx) \right| \leq 2P_Z(B_r^c) \sup_{z \in \mathbb{R}^{d_Z}} |h(z)|$$

by the uniform convergence of  $f_{n'}$  to  $f_\infty$  on compact subsets, and since  $\text{Lip } h < \infty$ . Taking  $r \rightarrow \infty$ , the right-hand side vanishes since  $h$  is bounded, and we obtain [\(B.5\)](#).

We are now ready to complete the proof. Since  $f_\infty$  is bi-Lipschitz, [Proposition B.5](#) means that  $f_\infty$  is a homeomorphism from  $\mathbb{R}^{d_Z}$  to  $f_\infty(\mathbb{R}^{d_Z})$ . This certainly gives

$$\text{supp } P_Z \cong f_\infty(\text{supp } P_Z).$$

But now [Proposition B.6](#) means

$$f_\infty(\text{supp } P_Z) = \overline{f_\infty(\text{supp } P_Z)}$$

where the closure is taken in  $\mathbb{R}^{d_X}$ . However, from [\(B.5\)](#) we have

$$P_X^* = f_\infty \# P_Z,$$

which by [Proposition B.3](#) means that

$$\text{supp } P_X^* = \text{supp } f_\infty \# P_Z = \overline{f_\infty(\text{supp } P_Z)}.$$

Consequently

$$\text{supp } P_X^* = f_\infty(\text{supp } P_Z) \cong \text{supp } P_Z$$

as desired.  $\square$

The following corollary extends the above result to the case where  $\text{supp } P_X^*$  may be homeomorphic to  $\text{supp } P_Z$ , but  $P_X^*$  is very *close* to a probability measure with non-homeomorphic support to  $P_Z$ . Here  $\rho$  denotes any metric for the weak topology. In other words,  $\rho$  must be a metric on the space of distributions that satisfies  $\rho(P_n, P) \rightarrow 0$  as  $n \rightarrow \infty$  if and only if  $P_n \xrightarrow{\mathcal{D}} P$ . The Lévy-Prokhorov and bounded Lipschitz metrics provide standard examples of such  $\rho$  ([Villani, 2008](#), Definition 3.3.10).

**Corollary 2.2.** *Suppose  $P_Z$  and  $P_X^0$  are probability measures on  $\mathbb{R}^{d_Z}$  and  $\mathbb{R}^{d_X}$  respectively with  $\text{supp } P_Z \not\cong \text{supp } P_X^0$ . Then there exists nonincreasing  $M : [0, \infty) \rightarrow [1, \infty]$  with  $M(\epsilon) \rightarrow \infty$  as  $\epsilon \rightarrow 0$  such that, for any probability measure  $P_X^*$  on  $\mathbb{R}^{d_X}$ , we have  $\text{BiLip } f \geq M(\epsilon)$  whenever  $\rho(P_X^*, P_X^0) \leq \epsilon$  and  $\rho(f \# P_Z, P_X^*) \leq \epsilon$ .*

*Proof.* Define  $M : [0, \infty) \rightarrow [1, \infty]$  by

$$M(\epsilon) := \inf \{ \text{BiLip } f \mid f : \mathbb{R}^{d_Z} \rightarrow \mathbb{R}^{d_X}, \rho(f \# P_Z, P_X^0) \leq 2\epsilon \},$$

with  $M(\epsilon) := \infty$  if the infimum is taken over the empty set. Certainly  $M$  is nonincreasing. If we have both  $\rho(P_X^*, P_X^0) \leq \epsilon$  and  $\rho(f \# P_Z, P_X^*) \leq \epsilon$ , then the triangle inequality gives

$$\rho(f \# P_Z, P_X^0) \leq \rho(f \# P_Z, P_X^*) + \rho(P_X^*, P_X^0) \leq 2\epsilon$$

and so  $\text{BiLip } f \geq M(\epsilon)$  since the right-hand side is an infimum. It remains only to show that  $M(\epsilon) \rightarrow \infty$  as  $\epsilon \rightarrow 0$ . For contradiction, suppose there exists  $\epsilon_n \rightarrow 0$  such that  $\sup_n M(\epsilon_n) < \infty$ . From the definition of  $M$ , this means that for each  $n$  there exists  $f_n : \mathbb{R}^{d_Z} \rightarrow \mathbb{R}^{d_X}$  such that  $\rho(f_n \# P_Z, P_X^0) \leq 2\epsilon_n$  and  $\text{BiLip } f_n \leq M(\epsilon_n) + 1$ . It follows directly that  $\rho(f_n \# P_Z, P_X^0) \rightarrow 0$  as  $n \rightarrow \infty$ , which in turn means  $f_n \# P_Z \xrightarrow{\mathcal{D}} P_X^0$  since  $\rho$  is a metric for the weak topology. At the same time we have

$$\sup_n \text{BiLip } f_n \leq \sup_n M(\epsilon_n) + 1 < \infty,$$

which contradicts [Theorem 2.1](#), since we assumed  $\text{supp } P_Z \not\cong \text{supp } P_X^0$ .  $\square$

## B.4 Variance of the Russian Roulette Estimator

In this section we briefly review the Russian roulette estimator used in [Chen et al. \(2019\)](#), and then discuss some scenarios in which we expect the variance of this estimator to increase unboundedly.

### B.4.1 RUSSIAN ROULETTE ESTIMATOR

Residual Flows (ResFlows, ([Chen et al., 2019](#))), building off of Invertible Residual Networks (iResNets, ([Behrmann et al., 2019](#))), model the data by repeatedly stacking bijections of the form  $f_\ell^{-1}(x) = x + g_\ell(x)$ , where  $\text{Lip } g_\ell =: \kappa < 1$ , as mentioned in (5). The change-of-variable formula for one layer of flow reads as, for  $x \in \mathbb{R}^d$ ,

$$\log p_X(x) = \log p_Z(f_\ell^{-1}(x)) + \text{tr} \left( \sum_{j=1}^{\infty} \frac{(-1)^{j+1}}{j} \text{D}g_\ell(x)^j \right). \quad (\text{B.6})$$

To deal with this infinite series, iResNets truncate after a fixed number of terms – this provides a biased estimate of the log-likelihood of a point  $x$  under the model. ResFlows rely on an alternative method of estimating (B.6), first using a Russian roulette procedure to rewrite the series as follows:

$$\sum_{j=1}^{\infty} \frac{(-1)^{j+1}}{j} \text{tr} (\text{D}g_\ell(x)^j) = \mathbb{E}_N \left[ \sum_{j=1}^N \frac{(-1)^{j+1}}{j} \frac{\text{tr} (\text{D}g_\ell(x)^j)}{p_j} \right] =: S(x),$$

where  $N \sim \text{Geom}(p)$  is a geometric random variable, and  $p_k := \mathbb{P}(N \geq k)$ . Then, taking a single sample  $N \sim \text{Geom}(p)$ , an unbiased estimator of  $S$  is given as  $S_N$ , where  $S_n$  is defined for any  $n \in \mathbb{N}$  and  $x \in \mathbb{R}$  as

$$S_n(x) := \sum_{j=1}^n \frac{(-1)^{j+1}}{j} \frac{\text{tr}(\text{D}g_\ell(x)^j)}{p_j} \quad (\text{B.7})$$

for any  $x \in \mathbb{R}^d$ . We will study the variance of  $S_N$  in this section.<sup>15</sup>

First, however, define the quantity  $\alpha_j(x)$  for  $j \in \mathbb{N}$ ,  $x \in \mathbb{R}^d$  as

$$\alpha_j(x) := \frac{(-1)^{j+1}}{j} \text{tr}(\text{D}g(x)^j), \quad (\text{B.8})$$

where we now drop the dependence of  $g$  on  $\ell$ . Then,  $S(x) = \sum_{j=1}^{\infty} \alpha_j(x)$ , and  $S_N(x) = \sum_{j=1}^N \alpha_j(x)/p_j$ .

#### B.4.2 WHAT MIGHT HAPPEN WHEN $\kappa \rightarrow 1$ ?

We begin with an informal discussion on the variance of  $S_N$  as  $\kappa \rightarrow 1$ . First of all we know that, as  $\kappa \rightarrow 1$ , the mapping  $f^{-1}$  gets arbitrarily close to a non-invertible mapping: consider e.g.  $g(x) = -\kappa x$ , then  $f^{-1} = (1 - \kappa)\text{Id} \rightarrow 0$  as  $\kappa \rightarrow 1$ . This near non-invertibility has implications for the speed of convergence of both  $S(x)$  and its gradient,<sup>16</sup> as noted in these two results from [Behrmann et al. \(2019\)](#):

1. **Theorem 3:**  $\left| \sum_{j=1}^n \alpha_j(x) - \log \det(I + \text{D}g(x)) \right| \leq -d \left( \log(1 - \kappa) + \sum_{j=1}^n \frac{\kappa^j}{j} \right)$ ,
2. **Theorem 4:**  $\|\nabla_\theta (\alpha_j(x) - \log \det(I + \text{D}g(x)))\|_\infty = \mathcal{O}(\kappa^n)$ .

We can see that both bounds become very loose as  $\kappa \rightarrow 1$ , implying we cannot guarantee the fast convergence of either series. It then follows that we cannot invoke the results from [Rhee & Glynn \(2015\)](#) and [Beatson & Adams \(2019\)](#) to argue that the variance of the Russian roulette estimator  $S_N$  will be small. Indeed, in the next section, we will look at a specific example where this variance becomes *infinite*.

#### B.4.3 A SPECIFIC EXAMPLE OF INFINITE VARIANCE

Now consider the case where  $d = 1$ . We will show that when  $\kappa^2 > 1 - p$ , there is a set of  $x$  having positive Lebesgue measure such that  $S_N(x)$  from (B.7) has infinite variance.

We note that here we have  $\text{tr}(\text{D}g(x)^j) = (g'(x))^j$  for any  $j \in \mathbb{N}$ . We can thus rewrite  $\alpha_j$  from (B.8) as

$$\alpha_j(x) := \frac{(-1)^{j+1}}{j} (g'(x))^j. \quad (\text{B.9})$$

Also recall that  $N \sim \text{Geom}(p)$  and  $p_j := \mathbb{P}(N \geq j)$  for all  $j \in \mathbb{N}$ .

**Proposition B.10.** *For any  $x \in \mathbb{R}$  and random variable  $N$  satisfying  $\text{supp } N = \mathbb{N}$ ,  $S_N(x)$  has finite expectation if  $\kappa < 1$ .*

*Proof.* Refer to [Lyne et al. \(2015, Proposition A.1\)](#). □

**Proposition B.11.** *Under the same conditions as [Proposition B.10](#),*

$$\text{Var} S_N(x) \geq \lim_{n \rightarrow \infty} 2 \sum_{j=1}^n \alpha_j(x) S_{j-1}(x) - \mathbb{E}[S_N(x)]^2.$$

<sup>15</sup>[Chen et al. \(2019\)](#) additionally approximate  $\text{tr}(\text{D}g_\ell(x)^j)$  by the Hutchinson's trace estimator  $v^T \text{D}g_\ell(x)^j v$  for  $v \sim \mathcal{N}(0, I)$ . Since  $v$  is independent of  $N$ , their estimator has strictly higher variance than (B.7).

<sup>16</sup>With respect to the flow parameters  $\theta$

*Proof.* This proof is taken from Lyne et al. (2015, Proposition A.2); we mostly rewrite the proof but adapt it to our specific setting and notation. Note that we will drop the dependence of  $S_j$  and  $\alpha_j$  on  $x$  throughout the proof.

We know from Proposition B.10 that  $\mathbb{E}[S_N(x)]$  is finite. Thus we will simply lower-bound  $\mathbb{E}[S_N(x)^2]$ .

We will first use induction to show the following holds for any  $n \in \mathbb{N}$ :

$$\sum_{j=1}^n S_j^2(p_j - p_{j+1}) = \alpha_1^2 + \sum_{j=2}^n \frac{\alpha_j^2}{p_j} + 2 \sum_{j=2}^n \alpha_j S_{j-1} - S_n^2 p_{n+1}. \quad (\text{B.10})$$

The base case is

$$S_1^2(p_1 - p_2) = \frac{\alpha_1^2}{p_1} p_1 - S_1^2 p_2 = \alpha_1^2 - S_1^2 p_2$$

since  $p_1 = 1$ . Now, assume (B.10) holds for some  $m \in \mathbb{N}$ . Then, for  $n = m + 1$ ,

$$\begin{aligned} \sum_{j=1}^{m+1} S_j^2(p_j - p_{j+1}) &= \sum_{j=1}^m S_j^2(p_j - p_{j+1}) + S_{m+1}^2(p_{m+1} - p_{m+2}) \\ &= \alpha_1^2 + \sum_{j=2}^m \frac{\alpha_j^2}{p_j} + 2 \sum_{j=2}^m \alpha_j S_{j-1} - S_m^2 p_{m+1} \\ &\quad + S_{m+1}^2(p_{m+1} - p_{m+2}) \end{aligned} \quad (\text{B.11})$$

by the inductive hypothesis. We also have

$$\begin{aligned} p_{m+1}(S_m^2 - S_{m+1}^2) &= p_{m+1}(S_m - S_{m+1})(S_m + S_{m+1}) \\ &= p_{m+1} \frac{\alpha_{m+1}}{p_{m+1}} \left( 2S_m + \frac{\alpha_{m+1}}{p_{m+1}} \right) \\ &= \frac{\alpha_{m+1}^2}{p_{m+1}} + 2\alpha_{m+1} S_m. \end{aligned}$$

Substituting this result into (B.11) completes the induction and proves (B.10) for all  $n \in \mathbb{N}$ .

Now, by Jensen's inequality,

$$S_n^2 = \left( \sum_{j=1}^n \frac{p_j \frac{\alpha_j}{p_j}}{p_j} \right)^2 \leq \frac{\sum_{j=1}^n \frac{\alpha_j^2}{p_j}}{\sum_{j=1}^n p_j}.$$

This implies

$$p_{n+1} S_n^2 \leq p_n S_n^2 \leq \frac{p_n}{\sum_{j=1}^n p_j} \sum_{j=1}^n \frac{\alpha_j^2}{p_j} \leq \sum_{j=1}^n \frac{\alpha_j^2}{p_j}$$

since  $(p_n)$  is a positive sequence.

This finally implies the following lower bound for any  $n \in \mathbb{N}$ :

$$\begin{aligned}
 \sum_{j=1}^n S_j^2 \mathbb{P}(N = j) &= \sum_{j=1}^n S_j^2 (p_j - p_{j+1}) \\
 &= \alpha_1^2 + \sum_{j=2}^n \frac{\alpha_j^2}{p_j} + 2 \sum_{j=2}^n \alpha_j S_{j-1} - S_n^2 p_{n+1} \\
 &\geq \alpha_1^2 + \sum_{j=2}^n \frac{\alpha_j^2}{p_j} + 2 \sum_{j=2}^n \alpha_j S_{j-1} - \sum_{j=1}^n \frac{\alpha_j^2}{p_j} \\
 &= \alpha_1^2 (1 - p_1^{-1}) + 2 \sum_{j=2}^n \alpha_j S_{j-1} \\
 &= 2 \sum_{j=2}^n \alpha_j S_{j-1},
 \end{aligned}$$

where the final line follows because  $p_1 = 1$ .

Since  $\mathbb{E}[S_N^2] = \lim_{n \rightarrow \infty} \sum_{j=1}^n S_j^2 \mathbb{P}(N = j)$ , the proof is complete.  $\square$

We are about ready to prove the main result but require one more auxiliary result first.

**Proposition B.12.** *Suppose  $|b| > 1$ . Then,*

$$\lim_{n \rightarrow \infty} \frac{n}{b^n} \sum_{j=1}^{n-1} \frac{b^j}{j} = \frac{1}{b-1}.$$

*Proof.* We will first show that the limit exists, and then show that it equals  $(b-1)^{-1}$ . Let

$$c_n = \frac{n}{b^n} \sum_{j=1}^{n-1} \frac{b^j}{j}.$$

We can rewrite this as follows:

$$c_n = \sum_{j=1}^{n-1} \frac{n}{b^{n-j} j} = \sum_{j=1}^{n-1} \frac{n}{b^j (n-j)} = \sum_{j=1}^{n-1} \frac{1}{b^j} + \sum_{j=1}^{n-1} \frac{j}{b^j (n-j)}.$$

Since  $b > 1$ , the first sum is a convergent geometric series as  $n \rightarrow \infty$ . We can decompose the second sum into its positive and negative terms:

$$\sum_{j=1}^{n-1} \frac{j}{b^j (n-j)} = \sum_{j \geq 1: b^j > 0} \frac{j}{b^j (n-j)} + \sum_{j \geq 1: b^j < 0} \frac{j}{b^j (n-j)} \equiv \textcircled{1}_n + \textcircled{2}_n.$$

We can see, for all  $n \in \mathbb{N}$ ,

$$\textcircled{1}_n \geq - \sum_{j=1}^{n-1} \frac{j}{|b|^j (n-j)} \quad \text{and} \quad \textcircled{2}_n \leq \sum_{j=1}^{n-1} \frac{j}{|b|^j (n-j)}.$$

Furthermore, for all  $j \in \{1, \dots, n-1\}$ , we have

$$\frac{j}{n-j} \leq j.$$

Now notice that the series  $\sum_{j=1}^{\infty} \frac{j}{|b|^j}$  converges by the ratio test:

$$\lim_{j \rightarrow \infty} \left| \frac{\frac{j+1}{|b|^{j+1}}}{\frac{j}{|b|^j}} \right| = \lim_{j \rightarrow \infty} \frac{j+1}{j|b|} = \frac{1}{|b|} < 1.$$

This implies the existence of  $\lim_{n \rightarrow \infty} \sum_{j=1}^{n-1} \frac{j}{|b|^j (n-j)}$ . Since the sequence  $(\textcircled{1}_n)$  (resp.  $(\textcircled{2}_n)$ ) is negative, non-increasing, and bounded below (resp. positive, non-decreasing, and bounded above), this implies the existence of  $\lim_{n \rightarrow \infty} \textcircled{1}_n$  (resp.  $\lim_{n \rightarrow \infty} \textcircled{2}_n$ ). Altogether, this implies the existence of

$$\lim_{n \rightarrow \infty} \left( \sum_{j=1}^{n-1} \frac{1}{b^j} + \sum_{j=1}^{n-1} \frac{j}{b^j (n-j)} \right) = \lim_{n \rightarrow \infty} c_n =: c_\infty.$$

Now we will determine its precise value. Note the following recurrence for all  $n \in \mathbb{N}$ :

$$c_{n+1} = \frac{n+1}{bn} (1 + c_n).$$

Taking the limit of both sides as  $n \rightarrow \infty$  gives

$$c_\infty = \frac{1}{b} (1 + c_\infty).$$

Solving this gives us  $c_\infty = \frac{1}{b-1}$ , which completes the proof.  $\square$

**Proposition B.13.** *Suppose  $N \sim \text{Geom}(p)$ ,  $g$  is continuously differentiable, and  $1 - p < \kappa^2 < 1$ . Then*

$$\{x \in \mathbb{R} \mid \text{Var}S_N(x) = \infty\}$$

*has positive Lebesgue measure.*

*Proof.* From [Proposition B.11](#), for a given  $x \in \mathbb{R}$ , we can see that showing  $\sum_{n=2}^{\infty} \alpha_n(x) S_{n-1}(x)$  diverges is sufficient to prove  $\text{Var}S_N(x)$  is infinite.

Consider using the ratio test to assess the convergence of the above series, with terms defined as  $a_n(x) := \alpha_n(x) S_{n-1}(x)$ . We have the following for any  $n \geq 2$ :

$$\begin{aligned} \left| \frac{a_{n+1}(x)}{a_n(x)} \right| &= \left| \frac{\alpha_{n+1}(x) S_n(x)}{\alpha_n(x) S_{n-1}(x)} \right| \\ &= \frac{|(g'(x))^{n+1}|}{\frac{|(g'(x))^n|}{n}} \cdot \left| \frac{\sum_{j=1}^n \frac{\alpha_j(x)}{p_j}}{\sum_{j=1}^{n-1} \frac{\alpha_j(x)}{p_j}} \right| \\ &= \frac{n|g'(x)|}{n+1} \cdot \left| \frac{(-1)^{n+1} \cdot (g'(x))^n}{np_n} \left( \sum_{j=1}^{n-1} \frac{(-1)^{j+1} \cdot (g'(x))^j}{jp_j} \right)^{-1} + 1 \right|. \end{aligned}$$

Recall  $p_j = (1-p)^{j-1} \equiv q^{j-1}$ . Then, writing  $b = -\frac{g'(x)}{q}$ , we have

$$\frac{(-1)^{n+1} \cdot (g'(x))^n}{np_n} \left( \sum_{j=1}^{n-1} \frac{(-1)^{j+1} \cdot (g'(x))^j}{jp_j} \right)^{-1} = \frac{1}{n} b^n \left( \sum_{j=1}^{n-1} \frac{1}{j} b^j \right)^{-1}.$$

Now let us assume that  $|g'(x)|^2 > q$ . We can see that  $|g'(x)|^2 > q \implies |g'(x)| > q$  since  $q \in (0, 1)$ , which then entails  $|b| > 1$ . Therefore, by [Proposition B.12](#),

$$\lim_{n \rightarrow \infty} \frac{n}{b^n} \sum_{j=1}^{n-1} \frac{b^j}{j} = \frac{1}{b-1}.$$

This then implies

$$\begin{aligned} \lim_{n \rightarrow \infty} \left| \frac{a_{n+1}(x)}{a_n(x)} \right| &= \lim_{n \rightarrow \infty} \frac{n|g'(x)|}{n+1} \left| \frac{1}{\frac{n}{b^n} \sum_{j=1}^{n-1} \frac{b^j}{j}} + 1 \right| \\ &= |g'(x)| \left| \frac{1}{\frac{1}{b-1}} + 1 \right| = \frac{|g'(x)|^2}{q} > 1 \end{aligned}$$

since we have assumed that  $|g'(x)|^2 > q$ . Thus, for all  $x$  in the set

$$V_{g,q} := \{x \in \mathbb{R} \mid |g'(x)|^2 > q\},$$

the series  $\sum_{n=2}^{\infty} \alpha_n(x) S_{n-1}(x)$  diverges by the ratio test. This means that  $\text{Var} S_N(x) = \infty$  for all  $x \in V_{g,q}$ .

Finally, we will prove the set  $\{x \in \mathbb{R} \mid \text{Var} S_N(x) = \infty\}$  has positive Lebesgue measure. Recall that  $\text{Lip } g = \kappa$ , which directly implies  $\sup_{x \in \mathbb{R}} |g'(x)| = \kappa$  from [Proposition B.7](#) and thus  $\sup_{x \in \mathbb{R}} |g'(x)|^2 = \kappa^2$ . Then, since  $\kappa^2 > q$ , there exists  $x_0 \in \mathbb{R}$  such that  $|g'(x_0)|^2 \in (q, \kappa^2)$ . By the continuity of  $|g'|$ , there is open ball of nonzero radius around  $x_0$ , denoted  $\mathcal{B}(x_0)$ , such that  $|g'(x)| > q$  for all  $x \in \mathcal{B}(x_0)$ . Since  $\mathcal{B}(x_0)$  is open and non-empty, it has positive Lebesgue measure. The inclusions

$$\mathcal{B}(x_0) \subseteq V_{g,q} \subseteq \{x \in \mathbb{R} \mid \text{Var} S_N(x) = \infty\}$$

thus conclude the proof. □

#### B.4.4 DISCUSSION

**Changing  $p$  as  $\kappa$  increases** An obvious strategy to avoid satisfying the conditions of [Proposition B.13](#) is to set  $p$  such that  $1 - \kappa^2 > p$ . However, lowering  $p$  in this way incurs additional computational cost: the average number of iterations per training step is equal to  $p^{-1}$ , or is lower-bounded by  $(1 - \kappa^2)^{-1}$  if  $p < 1 - \kappa^2$ . Thus, if we send  $\kappa \rightarrow 1$  to mitigate the bi-Lipschitz constraint (6), we will either incur an infinite computational cost or run the risk of encountering infinite variance.

**Higher dimensions** Although [Proposition B.13](#) only applies for  $d = 1$ , it is conceivable that similar results can be derived for  $d > 1$ , especially when considering the discussion in [Section B.4.2](#). We leave a deeper investigation for future work.

#### B.5 Density of a CIF

We make precise our heuristic derivation of the density (11) via the following result.

**Proposition B.14.** *Suppose  $\mathcal{Z}, \mathcal{X} \subseteq \mathbb{R}^d$  are open, and that  $F(\cdot; u) : \mathcal{Z} \rightarrow \mathcal{X}$  is a continuously differentiable bijection with everywhere invertible Jacobian for each  $u \in \mathcal{U}$ . Under the generative model (8),  $(X, U)$  has joint density*

$$p_Z(F^{-1}(x; u)) p_{U|Z}(u|F^{-1}(x; u)) |\det DF^{-1}(x; u)|.$$

*Proof.* Suppose  $h : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$  is a bounded measurable test function. Then

$$\begin{aligned} \mathbb{E}[h(X, U)] &= \mathbb{E}[h(F(Z; U), U)] \\ &= \int \left[ \int h(F(z; u), u) p_Z(z) p_{U|Z}(u|z) dz \right] du \\ &= \int h(x, u) p_Z(F^{-1}(x; u)) p_{U|Z}(u|F^{-1}(x; u)) |\det DF^{-1}(x; u)| dz du, \end{aligned}$$

where in the third line we substitute  $x := F(z; u)$  on the inner integral, which is valid by Theorem 17.2 of [Billingsley \(2008\)](#). Now for  $A \subseteq \mathcal{X} \times \mathcal{U}$ , let  $h := \mathbb{I}_A$ . It follows that

$$\mathbb{P}((X, U) \in A) = \mathbb{E}[\mathbb{I}_A(X, U)] = \int_A p_Z(F^{-1}(x; u)) p_{U|Z}(u|F^{-1}(x; u)) |\det DF^{-1}(x; u)| dz du,$$

which gives the result since  $A$  was arbitrary. □

#### B.6 Our Approximate Posterior Does Not Sacrifice Generality

The following result shows that our parameterisation of the approximate posterior  $q_{U_{1:L}|X}$  in (15) does not lose generality. In particular, provided each  $q_{U_\ell|Z_\ell}$  is sufficiently expressive, we can always recover the exact posterior.



**Proposition B.15.** *Under the generative model (8), the posterior factors like*

$$p_{U_{1:L}|X}(u_{1:L}|x) = \prod_{\ell=1}^L p_{U_\ell|Z_\ell}(u_\ell|z_\ell),$$

where  $z_L := x$  and  $z_\ell := F_{\ell+1}^{-1}(z_{\ell+1}; u_{\ell+1})$  for  $\ell \in \{1, \dots, L-1\}$ .

*Proof.* Writing  $p_{U_{1:L}|X}$  autoregressively gives

$$p_{U_{1:L}|X}(u_{1:L}|x) = \prod_{\ell=1}^L p_{U_\ell|U_{\ell+1:L},X}(u_\ell|u_{\ell+1:L}, x).$$

But now it is clear from the generative model (8) that  $U_\ell$  is conditionally independent of  $(U_{\ell+1:L}, X)$  given  $Z_\ell$ , and as such

$$p_{U_\ell|U_{\ell+1:L},X}(u_\ell|u_{\ell+1:L}, x) = p_{U_\ell|Z_\ell}(u_\ell|z_\ell).$$

Substituting this into the above expression then gives the result.  $\square$

### B.7 Conditions for a CIF to Outperform an Underlying Normalising Flow

For this result, the components of our model are assumed to be parameterised by  $\theta \in \Theta$ , which we will indicate by by  $F_\theta$ ,  $p_{U|Z}^\theta$ , and  $q_{U|X}^\theta$ . We will also use  $\theta$  to indicate quantities that result from the choice of parameters  $\theta$  (e.g.  $P_X^\theta$  for the distribution obtained), and will denote by  $\mathcal{L}^\theta$  the corresponding ELBO (14).

**Proposition 4.1.** *Suppose there exists  $\phi \in \Theta$  such that, for some bijection  $f : \mathcal{Z} \rightarrow \mathcal{X}$ ,  $F_\phi(\cdot; u) = f(\cdot)$  for all  $u \in \mathcal{U}$ . Likewise, suppose  $p_{U|Z}^\phi$  and  $q_{U|X}^\phi$  are such that, for some density  $r$  on  $\mathcal{U}$ ,  $p_{U|Z}^\phi(\cdot|z) = q_{U|X}^\phi(\cdot|x) = r(\cdot)$  for all  $z \in \mathcal{Z}$  and  $x \in \mathcal{X}$ . If  $\mathbb{E}_{x \sim P_X^*}[\mathcal{L}^\theta(x)] \geq \mathbb{E}_{x \sim P_X^*}[\mathcal{L}^\phi(x)]$ , then*

$$D_{\text{KL}}(P_X^* \parallel P_X^\theta) \leq D_{\text{KL}}(P_X^* \parallel f\#P_Z).$$

*Proof.* Observe from (11) that

$$p_{X,U}^\phi(x, u) = p_Z(f^{-1}(x)) |\det Df^{-1}(x)| p_{U|Z}(u|f^{-1}(x)).$$

It then follows from (13) that, under  $\phi$ , the model has density

$$p_X^\phi(x) = p_Z(f^{-1}(x)) |\det Df^{-1}(x)| \underbrace{\int p_{U|Z}(u|f^{-1}(x)) du}_{=1}$$

which is exactly the density of the normalising flow  $f\#P_Z$ . We also obtain the posterior

$$\begin{aligned} p_{U|X}^\phi(u|x) &= \frac{p_{X,U}^\phi(x, u)}{p_X^\phi(x)} \\ &= p_{U|Z}(u|f^{-1}(x)) \\ &= r(u). \end{aligned}$$

Since each  $q_{U|X}^\phi(\cdot|x) = r(\cdot)$  also, it follows that  $\mathcal{L}^\phi$  is tight, so that  $\mathcal{L}^\phi(x) = \log p_X^\phi(x)$  for all  $x \in \mathcal{X}$ .

Now suppose some  $\theta \in \Theta$  has

$$\mathbb{E}_{x \sim P_X^*}[\mathcal{L}^\theta(x)] \geq \mathbb{E}_{x \sim P_X^*}[\mathcal{L}^\phi(x)].$$

It follows that

$$\mathbb{E}_{x \sim P_X^*}[\log p_X^\theta(x)] \geq \mathbb{E}_{x \sim P_X^*}[\mathcal{L}^\phi(x)] = \mathbb{E}_{x \sim P_X^*}[\log p_X^\phi(x)].$$

Subtracting  $\mathbb{E}_{x \sim P_X^*}[\log p_X^*(x)]$  from both sides and negating gives

$$D_{\text{KL}}(P_X^* \parallel P_X^\theta) \leq D_{\text{KL}}(P_X^* \parallel P_X^\phi) = D_{\text{KL}}(P_X^* \parallel f\#P_Z).$$

$\square$

### B.8 CIFs Can Learn Target Supports Exactly

In this section we give necessary and sufficient conditions for a CIF to learn the support of a target distribution exactly, without needing changes to  $F$ . However, our argument applies more generally and does not make specific use of the bijective structure of  $F$ . To make this clear, we formulate our result here in terms of a generalisation of the model (7). In particular, we will take  $P_X$  as the marginal in  $X$  of

$$Z \sim P_Z, \quad U \sim P_{U|Z}(\cdot|Z), \quad X := G(Z, U), \quad (\text{B.12})$$

where  $G : \mathcal{Z} \times \mathcal{U} \rightarrow \mathcal{X}$ . We will assume that

- $\mathcal{Z} \subseteq \mathbb{R}^{d_Z}$ ,  $\mathcal{U} \subseteq \mathbb{R}^{d_U}$ , and  $\mathcal{X} \subseteq \mathbb{R}^{d_X}$  are equipped with the subspace topology;
- $P_Z$  and each  $P_{U|Z}(\cdot|z)$  are Borel probability measures on  $\mathcal{Z}$  and  $\mathcal{U}$  respectively;
- $G$  is continuous with respect to the product topology  $\mathcal{Z} \times \mathcal{U}$ .

We then have the following formula for  $\text{supp } P_X$ :

**Lemma B.16.** *Under the model (B.12),*

$$\text{supp } P_X = \overline{\bigcup_{z \in \text{supp } P_Z} G(\{z\} \times \text{supp } P_{U|Z}(\cdot|z))}.$$

*Proof.* Denote the joint distribution of  $(Z, U)$  by  $P_{Z,U}$ . Observe from [Proposition B.3](#) that

$$\text{supp } P_X = \overline{G(\text{supp } P_{Z,U})}.$$

Let

$$B := \bigcup_{z \in \text{supp } Z} \{z\} \times \text{supp } P_{U|Z}(\cdot|z).$$

The result follows if we can show that

$$\text{supp } P_{Z,U} = \overline{B},$$

since  $\overline{G(\overline{B})} = \overline{G(B)}$  because  $G$  is continuous.

We first show that  $\text{supp } P_{Z,U} \supseteq \overline{B}$ . Suppose  $(z, u) \in B$ , and let  $N_{(z,u)} \subseteq \mathcal{Z} \times \mathcal{U}$  be an open set containing  $(z, u)$ . Then there exists open  $N_z$  and  $N_u$  containing  $z$  and  $u$  respectively such that  $N_z \times N_u \subseteq N_{(z,u)}$ , since the open rectangles form a base for the product topology. It follows that

$$\begin{aligned} P_{Z,U}(N_{(z,u)}) &\geq P_{Z,U}(N_z \times N_u) \\ &= \int_{N_z} P_{U|Z}(N_u|z') P_Z(dz') \\ &> 0, \end{aligned}$$

since by the definition of  $B$  we have  $P_Z(N_z) > 0$  and  $P_{U|Z}(N_u|z) > 0$  for each  $u \in N_z$ . From this we have  $\text{supp } P_{Z,U} \supseteq B$ , and taking the closure of each side gives  $\text{supp } P_{Z,U} \supseteq \overline{B}$ .

In the other direction, suppose that  $(z, u) \notin \overline{B}$ . Then there exist open sets  $N_z$  and  $N_u$  containing  $z$  and  $u$  respectively such that

$$(N_z \times N_u) \cap \overline{B} = \emptyset.$$

By the definition of  $B$ , it follows that if  $(z', u') \in N_z \times N_u$  and  $z' \in \text{supp } P_Z$ , then  $u' \notin \text{supp } P_{U|Z}(\cdot|z')$ . Otherwise stated, if  $z' \in N_z \cap \text{supp } P_Z$ , then

$$N_u \cap \text{supp } P_{U|Z}(\cdot|z') = \emptyset.$$

Thus

$$\begin{aligned}
 P_{Z,U}(N_z \times N_u) &= \int_{N_z} \left[ \int_{N_u} P_{U|Z}(du'|z') \right] P_Z(dz') \\
 &= \int_{N_z \cap \text{supp } P_Z} \left[ \int_{N_u \cap \text{supp } P_{U|Z}(\cdot|z')} P_{U|Z}(du'|z') \right] P_Z(dz') \\
 &= 0,
 \end{aligned}$$

where the second line follows from [Proposition B.4](#). Consequently  $(z, u) \notin \text{supp } P_{Z,U}$ , which gives  $\text{supp } P_{Z,U} \subseteq \overline{B}$ .  $\square$

We now give necessary and sufficient conditions for the model (B.12) to learn a given target support exactly.

**Proposition B.17.** *Suppose  $P_X^*(\partial \text{supp } P_X^*) = 0$  and that*

$$\overline{G(\text{supp } P_Z \times \mathcal{U})} \supseteq \text{supp } P_X^*. \quad (\text{B.13})$$

*Then there exists  $P_{U|Z}$  such that  $\text{supp } P_X = \text{supp } P_X^*$  if and only if, for all  $z \in \text{supp } P_Z$ , there exists  $u \in \mathcal{U}$  with*

$$G(z, u) \in \text{supp } P_X^*.$$

*Proof.* ( $\Rightarrow$ ) Choose  $P_{U|Z}$  such that  $\text{supp } P_X = \text{supp } P_X^*$ . [Lemma B.16](#) gives

$$\bigcup_{z \in \text{supp } P_Z} G(\{z\} \times \text{supp } P_{U|Z}(\cdot|z)) \subseteq \text{supp } P_X^*.$$

Suppose  $z \in \text{supp } P_Z$ . Then for indeed all  $u \in \text{supp } P_{U|Z}(\cdot|z)$  we must have  $G(z, u) \in \text{supp } P_X^*$ , which proves this direction since  $\text{supp } P_{U|Z}(\cdot|z) \neq \emptyset$  by [Proposition B.4](#).

( $\Leftarrow$ ) For  $z \in \text{supp } P_Z$ , let

$$A_z := \{u \in \mathcal{U} : G(z, u) \in \text{int}(\text{supp } P_X^*)\}, \quad (\text{B.14})$$

where  $\text{int}$  denotes the *interior* operator. If  $A_z = \emptyset$ , define  $P_{U|Z}(\cdot|z)$  to be Dirac on some  $u$  such that  $G(z, u) \in \text{supp } P_X^*$ , which exists by assumption. Otherwise, we let  $P_{U|Z}(\cdot|z)$  be a probability measure with support  $\overline{A_z}$ . To show that such a measure exists, observe that  $A_z$  is open since  $G$  is continuous. Since  $\mathcal{U}$  is separable, we can therefore write

$$A_z = \bigcup_{n=1}^{\infty} B_n$$

for a countable collection of open sets  $B_n \subseteq A_z$ . We can then define a probability measure  $\mu$  by

$$\mu(C) := \sum_{n=1}^{\infty} 2^{-n} \mathbb{I}(C \cap B_n \neq \emptyset)$$

for measurable  $C \subseteq \mathcal{U}$ . Since  $A_z \neq \emptyset$ , it is straightforward to see that this is a probability measure with  $\mu(A_z) = 1$ . Consequently  $\text{supp } \mu = \overline{A_z}$  by [Proposition B.2](#), since  $A_z$  is open and  $\text{supp } \mu$  is the smallest closed set with  $\mu$ -probability 1.

We show this construction gives  $\text{supp } P_X \subseteq \text{supp } P_X^*$ . To this end, we first prove that if  $z \in \text{supp } P_Z$  and  $u \in \text{supp } P_{U|Z}(\cdot|z)$  then

$$G(z, u) \in \text{supp } P_X^*.$$

If  $A_z = \emptyset$  this is immediate. Otherwise, since  $\text{supp } P_{U|Z}(\cdot|z) = \overline{A_z}$ , there exists  $(u_n) \subseteq A_z$  such that  $u_n \rightarrow u$ . By (B.14), each  $G(z, u_n) \in \text{supp } P_X^*$ . By continuity we then have

$$G(z, u_n) \rightarrow G(z, u) \in \text{supp } P_X^*$$

since  $\text{supp } P_X^*$  is closed. It follows that

$$\bigcup_{z \in \text{supp } P_Z} G(\{z\} \times \text{supp } P_{U|Z}(\cdot|z)) \subseteq \text{supp } P_X^*,$$

which gives  $\text{supp } P_X \subseteq \text{supp } P_X^*$  from Lemma B.16 since  $\text{supp } P_X^*$  is closed.

We now show  $\text{supp } P_X \supseteq \text{supp } P_X^*$ . Since  $P_X^*(\partial \text{supp } P_X^*) = 0$  we have

$$\text{supp } P_X^* = \overline{\text{int}(\text{supp } P_X^*)}$$

by Proposition B.2, so that  $\text{supp } P_X \supseteq \text{supp } P_X^*$  if  $\text{supp } P_X \supseteq \text{int}(\text{supp } P_X^*)$ . Now suppose  $x \in \text{int}(\text{supp } P_X^*)$ . Then there exists  $(z_n) \subseteq \text{supp } P_Z$  and  $(u_n) \subseteq \mathcal{U}$  such that  $G(z_n, u_n) \rightarrow x$  by (B.13). But then we must have  $G(z_n, u_n) \in \text{int}(\text{supp } P_X^*)$  for  $n$  large enough because  $x$  lies in the interior. Consequently, for  $n$  large enough,

$$u_n \in A_{z_n} \subseteq \text{supp } P_{U|Z}(\cdot|z_n)$$

and hence  $G(z_n, u_n) \in \text{supp } P_X$  by Lemma B.16. This means  $x \in \text{supp } P_X$  since  $\text{supp } P_X$  is closed.  $\square$

The following proposition then gives a straightforward condition under which it is additionally possible to recover the *target* exactly (i.e. not just its support). In our experiments we do not enforce this condition explicitly. However, since we learn the parameters of  $G$  here, we can expect our model will approximate this behaviour if doing so produces a better density estimator.

**Proposition B.18.** *If  $G(z, \cdot)$  is surjective for each  $z \in \mathcal{Z}$ , then there exists  $P_{U|Z}$  such that  $P_X = P_X^*$ .*

*Proof.* Fix  $z \in \mathcal{Z}$ . Surjectivity of  $G(z, \cdot)$  means that, for  $x \in \mathcal{X}$ , there exists  $u \in \mathcal{U}$  such that  $G(z, u) = x$ . Thus we can define  $H_z : \mathcal{X} \rightarrow \mathcal{U}$  such that

$$G(z, H_z(x)) = x$$

for all  $x \in \mathcal{X}$ . We then define each

$$P_{U|Z}(\cdot|z) := H_z \# P_X^*.$$

From this it follows that  $P_X = P_X^*$ . For, letting  $B \subseteq \mathcal{X}$  be measurable,

$$\begin{aligned} P_X(B) &= \int_{G^{-1}(B)} P_{U|Z}(du|z) P_Z(dz) \\ &= \int \left[ \int \mathbb{I}_B(G(z, u)) H_z \# P_X^*(du) \right] P_Z(dz) \\ &= \int \left[ \int \mathbb{I}_B(G(z, H_z(x))) P_X^*(dx) \right] P_Z(dz) \\ &= \int \left[ \int \mathbb{I}_B(x) P_X^*(dx) \right] P_Z(dz) \\ &= P_X^*(B), \end{aligned}$$

which gives the result.  $\square$

## C Experimental Details

Our choices (10) and (17) required parameterising  $s$ ,  $t$ ,  $\mu^p$ ,  $\Sigma^p$ ,  $\mu^q$ , and  $\Sigma^q$ . Since these terms are naturally paired, at each layer of our model we set

$$\begin{aligned} [s(u), t(u)] &:= \text{NN}_F(u), \\ [\mu^p(z), \varsigma^p(z)] &:= \text{NN}_p(z), \\ \Sigma^p(z) &:= \text{diag}(e^{\varsigma^p(z)}), \\ [\mu^q(x), \varsigma^q(x)] &:= \text{NN}_q(x), \\ \Sigma^q(x) &:= \text{diag}(e^{\varsigma^q(x)}), \end{aligned}$$

where NN denotes a separate neural network and  $\varsigma^p(z), \varsigma^q(x) \in \mathbb{R}^d$ .

In all experiments we trained our models to maximise either the log-likelihood (for the baseline flows) or the ELBO (for the CIFs) using the ADAM optimiser (Kingma & Ba, 2015) with default hyperparameters and no weight decay. The ELBO was estimated using a single sample per datapoint (i.e. a single call to Algorithm 1). We used a held-out validation set and trained each model until its validation score stopped improving, except for the NSF tabular data experiments where we train for a fixed number of epochs as specified in Durkan et al. (2019). After training, we used validation performance to select the best parameters found during training for use at test time (again except for the NSF experiments, where we just test with the final model). Both validation and test scores were computed using the exact log-likelihood for the baseline and the importance sampling estimate (16) for the CIFs, with  $m = 5$  samples for validation and  $m = 100$  for testing.

## C.1 Tabular Data Experiments

Following Papamakarios et al. (2017), we experimented with the POWER, GAS, HEPMASS, and MINIBOONE datasets from the UCI repository (Bache & Lichman, 2013), as well as a dataset of  $8 \times 8$  image patches extracted from the BSDS300 dataset (Martin et al., 2001). We preprocessed these datasets identically to Papamakarios et al. (2017), and used the same train/validation/test splits. For all CIF-ResFlow models, we used a batch size of 1000 and a learning rate of  $10^{-3}$ . For the MAF experiments, we used a batch size of 1000 and a learning rate of  $10^{-3}$ , except for BSDS300 where we used a learning rate of  $10^{-4}$  to control the instability of the baseline. For the NSF experiments, we used batch sizes and learning rates as dictated by Durkan et al. (2019, Table 5), along with their cosine learning rate annealing scheme.

Also, for all CIF models, each  $U_\ell$  had the same dimension  $d_{U_\ell}$ , which we took to be roughly a quarter of the dimensionality of the data (except in Section C.1.4 for which  $d_{U_\ell} = d_{\mathcal{X}}$ ). In particular, we set  $d_{U_\ell} := 2$  for POWER and GAS,  $d_{U_\ell} := 5$  for HEPMASS,  $d_{U_\ell} := 10$  for MINIBOONE, and  $d_{U_\ell} := 15$  for BSDS300.

### C.1.1 RESIDUAL FLOWS

The residual blocks in all ResFlow models used multilayer perceptrons (MLPs) with 4 hidden layers of 128 hidden units (denoted  $4 \times 128$ ), LipSwish nonlinearities (Chen et al., 2019, (10)) before each linear layer, and a residual connection from the input to the output. We did not use any kind of normalisation (e.g. ActNorm or BatchNorm) for these experiments. For all models we set  $\kappa = 0.9$  in (5) to match the value for the 2-D experiments in the codebase of Chen et al. (2019). Other design choices followed Chen et al. (2019). In particular:

- We always exactly computed several terms at the beginning of the series expansion of the log Jacobian, and then used Russian Roulette sampling (Kahn, 1955) to estimate the sum of the remaining terms. In particular, at training time we computed 2 exact terms, while at test time we computed 20 exact terms;
- We used a geometric distribution with parameter 0.5 for the number of terms to compute in our Russian Roulette estimators;
- We used the Skilling-Hutchinson trace estimator (Skilling, 1989; Hutchinson, 1990) to estimate the trace in the log Jacobian term;
- At both training and test time, we used a single Monte Carlo sample of  $(n, v)$  to estimate (6) of Chen et al. (2019);

However, note that for these experiments, for the sake of simplicity, we did not use the memory-saving techniques in (8) and (9) of Chen et al. (2019), nor the adaptive power iteration scheme described in their Appendix E.

For  $\text{NN}_F$ ,  $\text{NN}_p$ , and  $\text{NN}_q$  we used  $2 \times 10$  MLPs with tanh nonlinearities. These networks were much smaller than  $4 \times 128$ , and hence the CIF-ResFlows had only roughly 1.5-4.5% more parameters (depending on the dimension of the dataset) than the otherwise identical 10-layer ResFlows, and roughly 10% of the parameters of the 100-layer ResFlows.

The 100-layer ResFlows were significantly slower to train than the 10-layer models, and for POWER, GAS, and BSDS300 we were forced to stop these before their validation loss had converged. However, to ensure a fair comparison, we allocated more total computing power to these models than to the 10-layer models, which were terminated properly. In particular, we trained each 100-layer ResFlow on POWER and GAS for a total of 10 days on a single NVIDIA GeForce GTX 1080 Ti, and on BSDS300 for a total of 7 days. In contrast, the 10-layer ResFlows converged after around 1 day on POWER, 4.5 days on GAS, and around 3 days on BSDS300. Likewise, the 10-layer CIF-ResFlows converged after around 1 day on POWER, 6 days on GAS, and 2 days on BSDS300.

Table C.4: MAF and CIF-MAF parameter configurations for POWER and GAS.

	LAYERS ( $L$ )	AUTOREGRESSIVE NETWORK SIZE	$NN_p$ SIZE	$NN_q$ SIZE	$NN_F$ SIZE
MAF	5, 10, 20	$2 \times 100, 2 \times 200, 2 \times 400$	-	-	-
CIF-MAF	5, 10	$2 \times 128$	$2 \times 100, 2 \times 200$	$2 \times 100, 2 \times 200$	$2 \times 128$

Table C.5: MAF and CIF-MAF parameter configurations for HEPMASS and MINIBOONE

	LAYERS ( $L$ )	AUTOREGRESSIVE NETWORK SIZE	$NN_p$ SIZE	$NN_q$ SIZE	$NN_F$ SIZE
MAF	5, 10, 20	$2 \times 128, 2 \times 512, 2 \times 1024$	-	-	-
CIF-MAF	5, 10	$2 \times 128$	$2 \times 128, 2 \times 512$	$2 \times 128, 2 \times 512$	$2 \times 128$

### C.1.2 MASKED AUTOREGRESSIVE FLOWS

The experiment comparing MAF baselines to CIF-MAFs was inspired by the experimental setup in Papamakarios et al. (2017). For each dataset, we specified a set of hyperparameters over which to search for both the baselines and the CIFs; these hyperparameters are provided in Table C.4, Table C.5, and Table C.6. Then, we trained each model until no validation improvement had been observed for 50 epochs. We then evaluated the model with the best validation score among all candidate models on the test dataset to obtain a log-likelihood score. We performed this procedure with three separate random seeds, and report the average and standard error across the runs in Table 1.

We searched over all combinations of parameters listed in Table C.4, Table C.5, and Table C.6. For example, on HEPMASS or MINIBOONE, our set of candidate MAF models included: for  $L = 5$ , an autoregressive network of size of either  $2 \times 128$ ,  $2 \times 512$ , or  $2 \times 1024$ ; for  $L = 10$ , an autoregressive network size of either  $2 \times 128$ ,  $2 \times 512$ , or  $2 \times 1024$ ; and for  $L = 20$ , again an autoregressive network size of either  $2 \times 128$ ,  $2 \times 512$ , or  $2 \times 1024$ ; this gave us a total of 9 candidate MAF models for each seed. The set of candidate CIF-MAF models can similarly be determined via the table and gave us a total of 8 candidate models for each seed. We maintained this split of 9 candidates for MAF and 8 candidates for CIF-MAF across datasets to fairly compare against the baseline by allowing them more configurations. We also considered deeper and wider MAF models to compensate for the additional parameters introduced by  $NN_F$ ,  $NN_p$ , and  $NN_q$  in the CIF-MAFs. Finally, we allowed the baseline MAF models to use batch normalization between MADE layers as recommended by Papamakarios et al. (2017), but we do not use them within CIF-MAFs as the structure of our  $F$  generalises this transformation.

We should note that our evaluation of models is slightly different from Papamakarios et al. (2017). For the model which scores best on the validation set, Papamakarios et al. (2017) report the average and standard deviation of log-likelihood across the points in the test dataset. However, our error bars emerge as the error in average test-set log-likelihood across *multiple* runs of the same experiment; this style of evaluation is often employed in other works as well (e.g. FFJORD (Grathwohl et al., 2019), NAF (Huang et al., 2018), and SOS (Jaini et al., 2019) as noted in Durkan et al. (2019, Table 1)).

### C.1.3 NEURAL SPLINE FLOWS

The experiment comparing NSF baselines to CIF-NSFs mirrors the experimental setup in Durkan et al. (2019). Specifically, we constructed baseline NSFs that exactly copied the settings in Durkan et al. (2019, Table 5). We also built CIF-NSFs using these baseline settings, although for the CIF-NSF-1 model we lowered the number of hidden channels in the autoregressive networks so that the total number of trainable parameters matched that of the baseline. Our parameter settings are provided in Table C.7; note that parameter settings are homogeneous across datasets, besides MINIBOONE for which we reduced

Table C.6: MAF and CIF-MAF parameter configurations for BSDS300

	LAYERS ( $L$ )	AUTOREGRESSIVE NETWORK SIZE	$NN_p$ SIZE	$NN_q$ SIZE	$NN_F$ SIZE
MAF	5, 10, 20	$2 \times 512, 2 \times 1024, 2 \times 2048$	-	-	-
CIF-MAF	5, 10	$2 \times 512$	$2 \times 128, 2 \times 512$	$2 \times 128, 2 \times 512$	$2 \times 128$

Table C.7: CIF-NSF configurations for all tabular datasets. The number of hidden features in the autoregressive network is referred to as  $n_h$ .

	NN <sub>p</sub> SIZE	NN <sub>q</sub> SIZE	NN <sub>F</sub> SIZE	$n_h$ VS. BASELINE
CIF-NSF-1 (MINIBOONE)	3 × 50	2 × 10	3 × 25	FEWER
CIF-NSF-1 (NON-MINIBOONE)	3 × 200	2 × 10	3 × 100	FEWER
CIF-NSF-2	3 × 200	2 × 10	3 × 100	SAME

 Table C.8: Mean ± standard error of average test set log-likelihood (higher is better). Best performing runs are shown in bold. CIF-Id-1 had  $s \equiv 0$  and  $t = \text{Id}$ . CIF-Id-2 had  $s \equiv 0$  and  $t = \text{NN}_F$ . CIF-Id-3 had  $(s, t) = \text{NN}_F$ .

	POWER	GAS	HEPMASS	MINIBOONE
CIF-ID-1 (NN <sub>q</sub> = 10 × 2)	0.43 ± 0.01	10.92 ± 0.10	-17.06 ± 0.05	-11.26 ± 0.03
CIF-ID-1 (NN <sub>q</sub> = 100 × 4)	0.42 ± 0.01	10.86 ± 0.16	-17.44 ± 0.09	-10.91 ± 0.04
CIF-ID-2 (NN <sub>q</sub> = 10 × 2)	0.45 ± 0.01	10.43 ± 0.08	-17.63 ± 0.10	-11.13 ± 0.08
CIF-ID-2 (NN <sub>q</sub> = 100 × 4)	0.47 ± 0.01	10.89 ± 0.18	-17.51 ± 0.09	-10.75 ± 0.07
CIF-ID-3 (NN <sub>q</sub> = 10 × 2)	<b>0.50 ± 0.01</b>	<b>11.32 ± 0.14</b>	-17.08 ± 0.02	-10.45 ± 0.04
CIF-ID-3 (NN <sub>q</sub> = 100 × 4)	<b>0.50 ± 0.01</b>	<b>11.58 ± 0.12</b>	<b>-16.68 ± 0.07</b>	<b>-10.01 ± 0.04</b>

the size NN<sub>p</sub> and NN<sub>F</sub> by a factor of 4 as per Durkan et al. (2019)<sup>17</sup>. We trained both NSF and CIF-NSFs for a number of training epochs corresponding to the number of training steps divided by the number of batches in the training set, i.e.

$$n_e = \lceil n_s / (n_t / n_b) \rceil,$$

where  $n_e$  is the number of epochs,  $n_s$  is the number of training steps,  $n_b$  is the batch size, and  $n_t$  is the number of training data points. Note that  $n_s$  and  $n_b$  are from Durkan et al. (2019, Table 5), and  $n_t$  is fixed by the pre-processing steps from Papamakarios et al. (2017). We then evaluated the test-set performance of each model after the pre-specified number of epochs, averaging across three seeds, and put the results in Table 1. We again average randomness across seeds, rather than across points in the test set, as discussed in the previous section.

We quickly note here that we selected our parameters after trying a few settings on various UCI datasets. There were other settings which performed better for individual datasets that are not included here, as we would like the proposed configurations to be as homogeneous as possible. It appeared as though the NSF models were already fairly good at modelling the data, which allowed us to make NN<sub>q</sub> much smaller while still achieving good inference.

We also should note that we wrapped our code around the NSF bijection code from <https://github.com/bayesiains/nsf>. We also disabled weight decay in all of these experiments without observing any problems with convergence.

#### C.1.4 ABLATING $f$

We ran ablation experiments to gain some insight into the relative importance of  $f$  in (9). In particular, we considered a 10 layer model ( $L = 10$ ) where at each layer  $U_\ell$  had the same dimension as the data and  $f = \text{Id}$  was the identity. We refer to this model as CIF-Id.

We considered three parameterisations of CIF-Id. The first had  $s \equiv 0$  and  $t = \text{Id}$ , which from our choice (10) of  $p_{U|Z}$  corresponds to stacking the following generative process:

$$\begin{aligned} Z &\sim P_Z \\ \epsilon &\sim \text{Normal}(0, I_d) \\ X &:= Z - \mu^p(Z) - e^{s^p(Z)} \odot \epsilon. \end{aligned} \tag{C.1}$$

<sup>17</sup>Indeed, there was no choice of  $n_h$  which would allow us to achieve the same number of parameters as the baseline for the models noted in row 2 of Table C.7.

Observe this generalise ResFlows, since (5) can be realised by sending  $\zeta^p \rightarrow -\infty$  and having  $\mu^p < 1$ . Accordingly, we took  $\text{NN}_p$  to be a  $4 \times 128$  MLP to match the size of the residual blocks used in our tabular ResFlow experiments.

The second CIF-Id parameterisation had  $s \equiv 0$  and  $t = \text{NN}_F$ , which amounts to replacing (C.1) with

$$X := Z - t \left( \mu^p(Z) + e^{s^p(Z)} \odot \epsilon \right).$$

To align with the first CIF-Id, we took  $\text{NN}_F$  and  $\text{NN}_p$  to be  $2 \times 128$  MLPs, and zeroed out the  $s$  output of  $\text{NN}_F$  to obtain  $s \equiv 0$ . The third parameterisation had  $(s, t) = \text{NN}_F$ , which replaces (C.1) with

$$X := \exp \left( -s \left( \mu^p(Z) + e^{s^p(Z)} \odot \epsilon \right) \right) \odot Z - t \left( \mu^p(Z) + e^{s^p(Z)} \odot \epsilon \right).$$

Again, we took  $\text{NN}_F$  and  $\text{NN}_p$  to be  $2 \times 128$  MLPs in this case.

We ran all configurations with two different choices of  $\text{NN}_q$ : a  $2 \times 10$  MLP as in our tabular ResFlow experiments, as well as a  $4 \times 100$  MLP. The results are given in Table C.8.<sup>18</sup> Observe that these models performed comparably or better than the 100-layer ResFlows, but worse than the CIF-ResFlows and CIF-MAFs in Table 1. As discussed in Section 4.1.1, we conjecture this occurs because a CIF-Id requires greater complexity from  $p_{U|Z}$  to make up for its simple choice of  $f$ , which in turn makes inference harder and hence the ELBO (14) looser, resulting in a poorer model that is learned overall. Likewise, note that the best performance in all cases was obtained when  $(s, t) = \text{NN}_F$ . This provides some justification for the generality of our choice of (9), as opposed to simpler alternatives that omit  $s$  or  $t$ .

## C.2 Image Experiments

In all our image experiments we applied the same uniform dequantisation scheme as Theis et al. (2016), after which we applied the logit transform of Dinh et al. (2017) with  $\alpha = 10^{-5}$  for Fashion-MNIST and  $\alpha = 0.05$  for CIFAR10.

### C.2.1 RESFLOW

For our baseline ResFlow experiments we used the same architecture as Chen et al. (2019). In particular, our convolutional residual blocks (denoted Conv-ResBlock) had the form

$$\text{LipSwish} \rightarrow 3 \times 3 \text{ Conv} \rightarrow \text{LipSwish} \rightarrow 1 \times 1 \text{ Conv} \rightarrow \text{LipSwish} \rightarrow 3 \times 3 \text{ Conv},$$

while our fully connected residual blocks (denoted FC-ResBlock) had the form

$$\text{LipSwish} \rightarrow \text{Linear} \rightarrow \text{LipSwish} \rightarrow \text{Linear},$$

with a residual connection from the input to the output in both cases. The overall architecture of the flow in all cases was:

$$\text{Image} \rightarrow \text{LogitTransform}(\alpha) \rightarrow k \times \text{Conv-ResBlock} \rightarrow [\text{Squeeze} \rightarrow k \times \text{Conv-ResBlock}] \times 2 \rightarrow 4 \times \text{FC-ResBlock},$$

where the Squeeze operation was as defined by Dinh et al. (2017). Like Chen et al. (2019), we used ActNorm layers (Kingma & Dhariwal, 2018) before and after each residual block.

Due to computational constraints, the models we considered were smaller than those used by Chen et al. (2019). In particular, our smaller ResFlow models used 128 hidden channels in their Conv-ResBlocks, 64 hidden channels in the linear layers of their FC-ResBlocks, and had  $k = 4$ . Our larger ResFlow models used 256 hidden channels in their Conv-ResBlocks, 128 hidden channels in the linear layers of their FC-ResBlocks, and had  $k = 6$ . In contrast, Chen et al. (2019) used 512 hidden channels in their Conv-ResBlocks, 128 hidden channels in their FC-ResBlocks, and had  $k = 16$ .

As described for our tabular experiments, we used the same estimation scheme as Chen et al. (2019). Additionally:

- We took  $\kappa = 0.98$ ;
- We used the Neumann gradient series expression for the log Jacobian (Chen et al., 2019, (8)) and computed gradients in the forward pass (Chen et al., 2019, (9)) to reduce memory overhead;

<sup>18</sup>Due to computational constraints we did not run these experiments on BSDS300.



- We used an adaptive rather than a fixed number of power iterations for spectral normalisation (Gouk et al., 2018), with a tolerance of 0.001;

For the CIF-ResFlows, we augmented the smaller baseline ResFlow by treating each composition of ActNorm  $\rightarrow$  ResBlock, as well as the final ActNorm, as an instance of  $f$  in (9). Each  $\text{NN}_F$ ,  $\text{NN}_p$ , and  $\text{NN}_q$  was a ResNet (He et al., 2016a;b) consisting of 2 residual blocks with 32 hidden channels (denoted  $2 \times 32$ ). We gave each  $U_\ell$  the same shape as a single channel of  $Z_\ell$ , and upsampled to the dimension of  $Z_\ell$  by adding channels at the output of each  $\text{NN}_F$ . Note that we did not experiment with using the larger baseline ResFlow model as the basis for a CIF.

For all models we used a learning rate of  $10^{-3}$  and a batch size of 64.

Figure C.3 through to Figure C.8 show samples synthesised from the ResFlow and CIF-ResFlow density models trained on MNIST and CIFAR-10.

### C.2.2 REALNVP

For our RealNVP-based image experiments, we took the baseline to be a RealNVP with the same architecture used by Dinh et al. (2017) for their CIFAR-10 experiments. In particular, we used 10 affine coupling layers with the corresponding alternating channelwise and checkerboard masks. Each coupling layer used a ResNet (He et al., 2016a;b) consisting of 8 residual blocks of 64 channels (denoted  $8 \times 64$ ). We replicated the multi-scale architecture of Dinh et al. (2017), squeezing the channel dimension after the first 3 coupling layers, and splitting off half the dimensions after the first 6. This model had 5.94M parameters for Fashion-MNIST and 6.01M parameters for CIFAR-10.

For the CIF-RealNVP, we considered each affine coupling layer to be an instance of  $f$  in (9). When choosing the size of our networks, we sought to maintain roughly the same depth over which gradients were propagated as in the baseline. To this end, our coupling networks were  $4 \times 64$  ResNets, each  $\text{NN}_p$  and  $\text{NN}_q$  were  $2 \times 64$  ResNets, and each  $\text{NN}_F$  was a  $2 \times 8$  ResNet. We gave each  $U_\ell$  the same shape as a single channel of  $Z_\ell$ , and upsampled to the dimension of  $Z_\ell$  by adding channels at the output of  $\text{NN}_F$ . Our model had 5.99M parameters for Fashion-MNIST and 6.07M parameters for CIFAR-10.

For completeness, we also trained a RealNVP model with coupler networks of size  $4 \times 64$  to match our CIF-RealNVP configuration. This model had 2.99M parameters for Fashion-MNIST and 3.05M for CIFAR-10.

In all cases for these experiments we used a learning rate of  $10^{-4}$  and a batch size of 100.

Figure C.9 through to Figure C.14 show samples synthesised from the RealNVP and CIF-RealNVP density models trained on Fashion-MNIST and CIFAR-10.

## C.3 2-D Experiments

To gain intuition about our model, we ran experiments on some simple 2-D datasets. For the datasets in Figure 1, we used a 10-layer ResFlow, a 100-layer ResFlow, and 10-layer CIF-ResFlow. For the CIF-ResFlows we took  $d_U = 1$ . Other architectural and training details were the same as for the tabular experiments described in Section 5.1 and Section C.1.1. The resulting average test set log-likelihoods for the top dataset were:

- -1.501 for the 10-layer ResFlow
- -1.419 for the 100-layer ResFlow
- -1.409 for the 10-layer CIF-ResFlow

The final average test set log likelihoods for the bottom dataset were:

- -2.357 for the 10-layer ResFlow
- -2.287 for the 100-layer ResFlow
- -2.275 for the 10-layer CIF-ResFlow

Note that in both cases the CIF-ResFlow slightly outperformed the 100-layer ResFlow.

We additionally ran several experiments comparing a baseline MAF against a CIF-MAF on the 2-D datasets shown in [Figure C.15](#). The baseline MAFs had 20 autoregressive layers, while the CIF-MAFs had 5. The network used at each layer had 4 hidden layers of 50 hidden units (denoted  $4 \times 50$ ). For the CIF-MAF, we took  $d_u = 1$ , and used  $2 \times 10$  MLPs for  $\text{NN}_F$  and  $4 \times 50$  MLPs for  $\text{NN}_p$  and  $\text{NN}_q$ . In total the baseline MAF had 160160 parameters, while our model had 119910 parameters.

The results of these experiments are shown in [Figure C.15](#). Observe that CIF-MAF consistently produces a more faithful representation of the target distribution than the baseline, and in all cases achieved higher average test set log probability. A failure mode of our approach is exhibited in the spiral dataset, where our model still lacks the power to fully capture the topology of the target. However, we did not find it difficult to improve on this: by increasing the size of  $\text{NN}_p$  to  $8 \times 50$  (and keeping all other parameters fixed), we were able to obtain the result shown in [Figure C.16](#). This model had a total of 221910 parameters. We also tried a larger MAF model with autoregressive networks of size  $8 \times 50$ , (obtaining 364160 parameters total). This model diverged after approximately 160 epochs. The result after 150 epochs is shown in [Figure C.16](#).



Figure C.3: Synthetic MNIST samples generated by the small baseline ResFlow model



Figure C.4: Synthetic MNIST samples generated by the large baseline ResFlow model

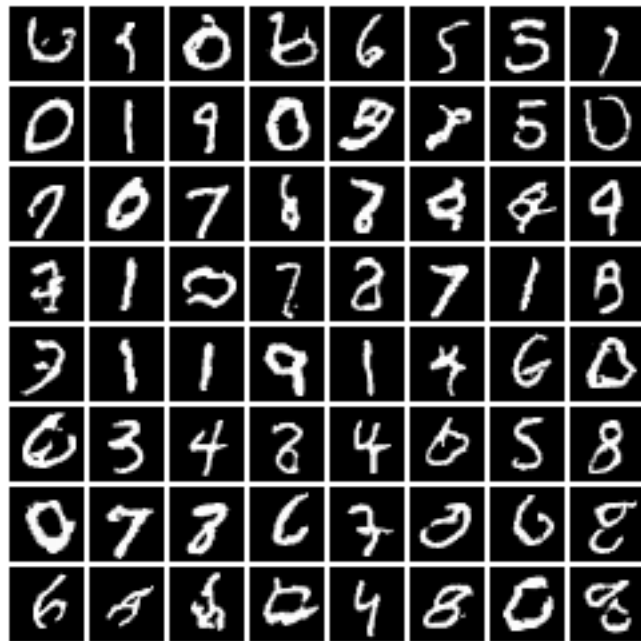


Figure C.5: Synthetic MNIST samples generated by the CIF-ResFlow model

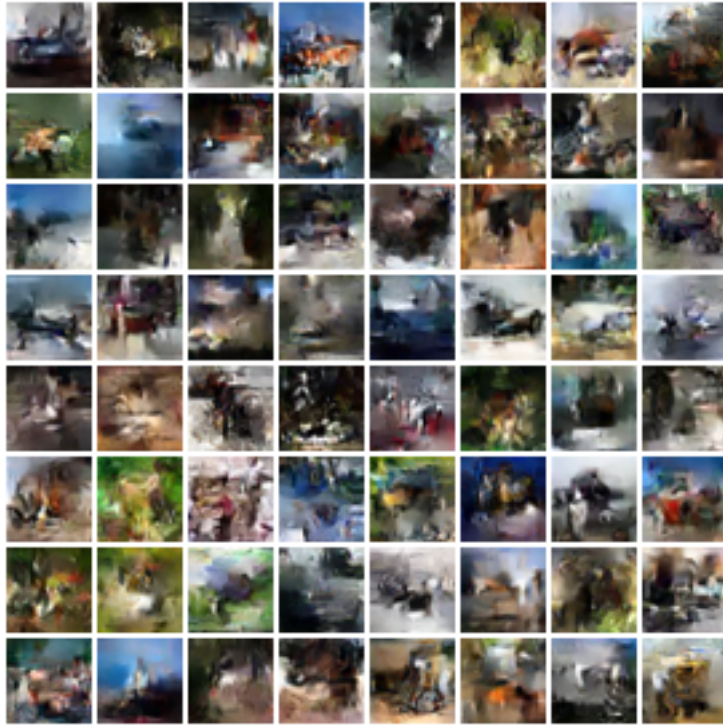


Figure C.6: Synthetic CIFAR-10 samples generated by the small baseline ResFlow model

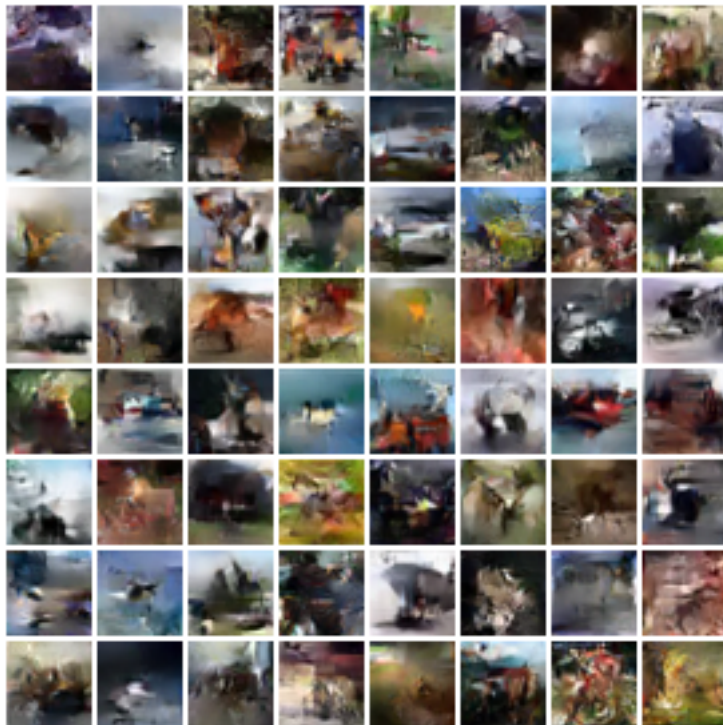


Figure C.7: Synthetic CIFAR-10 samples generated by the large baseline ResFlow model

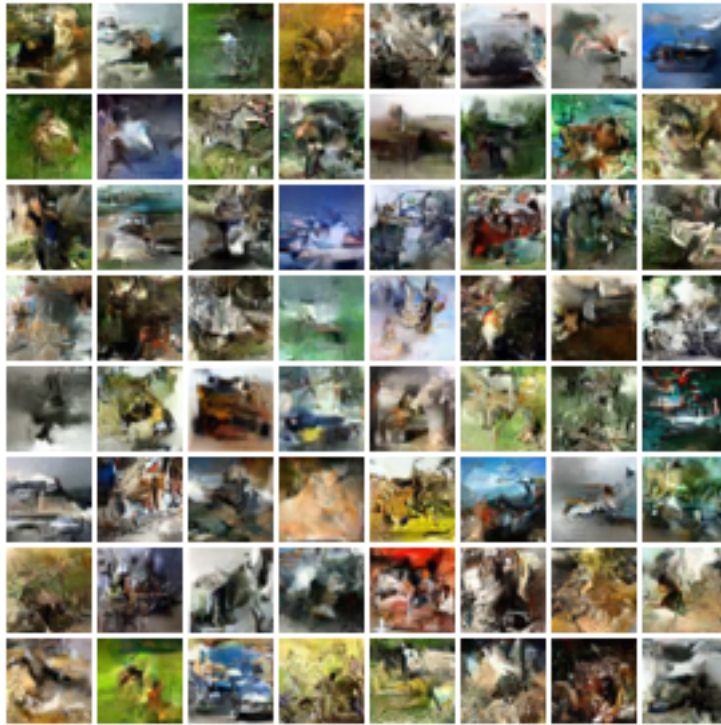


Figure C.8: Synthetic CIFAR-10 samples generated by the CIF-ResFlow model



Figure C.9: Synthetic Fashion-MNIST samples generated by RealNVP with coupling networks of size  $4 \times 64$



Figure C.10: Synthetic Fashion-MNIST samples generated by RealNVP with coupling networks of size  $8 \times 64$



Figure C.11: Synthetic Fashion-MNIST samples generated by CIF-RealNVP



Figure C.12: Synthetic CIFAR-10 samples generated by RealNVP with coupling networks of size  $4 \times 64$



Figure C.13: Synthetic CIFAR-10 samples generated by RealNVP with coupling networks of size  $8 \times 64$



Figure C.14: Synthetic CIFAR-10 samples generated by CIF-RealNVP



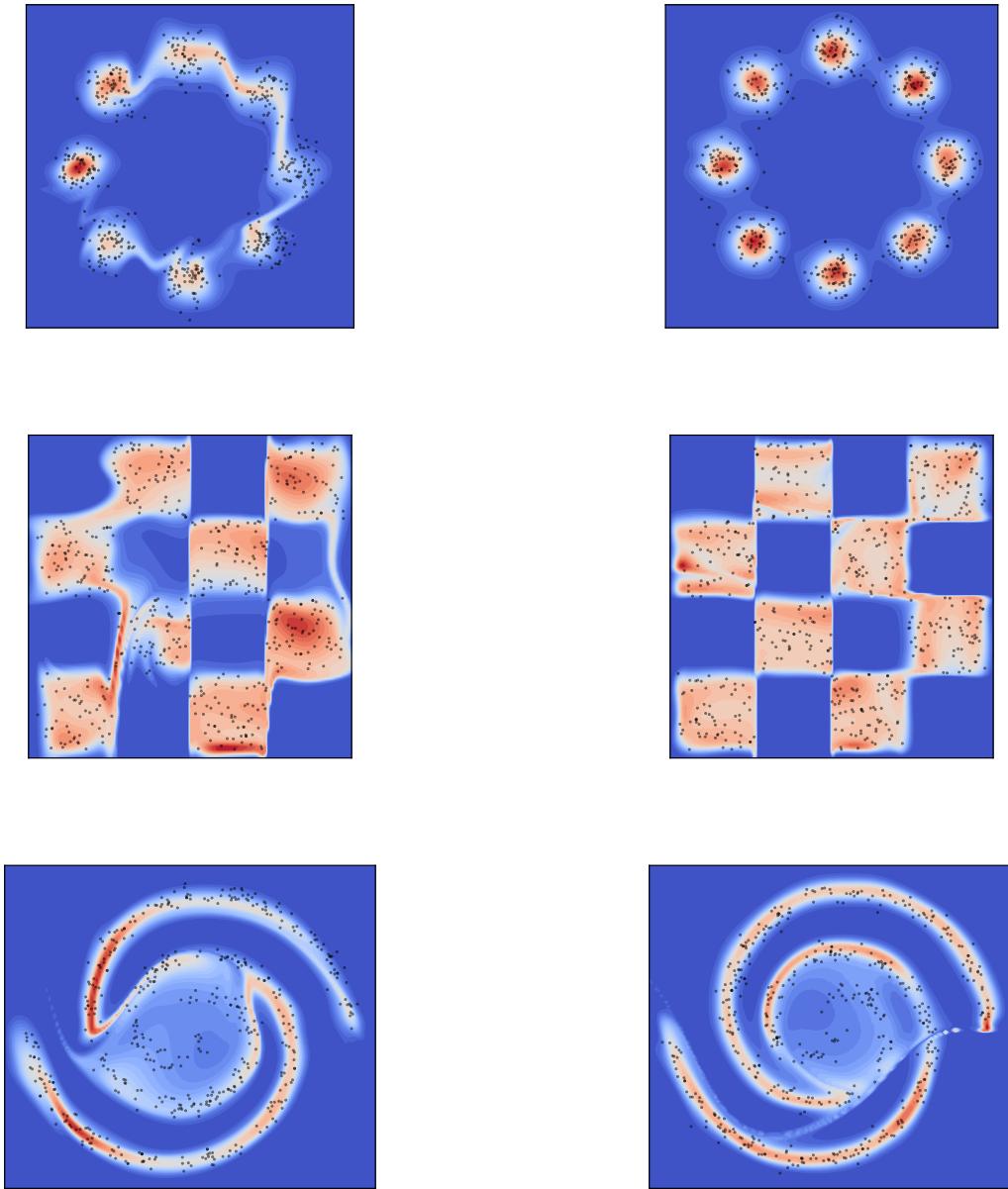


Figure C.15: Density models learned by a standard 20 layer MAF (left) and by a 5 layer CIF-MAF (right) for a variety of 2-D target distributions. Samples from the target are shown in black.

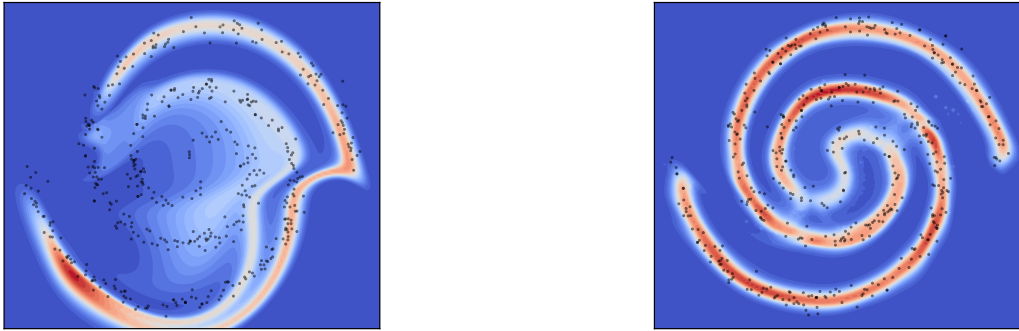


Figure C.16: Density models learned by a larger 20 layer MAF (left) and a larger 5 layer CIF-MAF (right) for the spirals dataset.