# SurVAE Flows: Surjections to Bridge the Gap between VAEs and Flows

**Didrik Nielsen**[1], **Priyank Jaini**[2], **Emiel Hoogeboom**[2], **Ole Winther**[1], **Max Welling**[2]

Technical University of Denmark[1], UvA-Bosch Delta Lab, University of Amsterdam[2]

didni@dtu.dk, p.jaini@uva.nl, e.hoogeboom@uva.nl
olwi@dtu.dk, m.welling@uva.nl

## Abstract

Normalizing flows and variational autoencoders are powerful generative models that can represent complicated density functions. However, they both impose constraints on the models: Normalizing flows use bijective transformations to model densities whereas VAEs learn stochastic transformations that are non-invertible and thus typically do not provide tractable estimates of the marginal likelihood. In this paper, we introduce SurVAE Flows: A modular framework of composable transformations that encompasses VAEs and normalizing flows. SurVAE Flows bridge the gap between normalizing flows and VAEs with *surjective transformations*, wherein the transformations are deterministic in one direction – thereby allowing exact likelihood computation, and stochastic in the reverse direction – hence providing a lower bound on the corresponding likelihood. We show that several recently proposed methods, including dequantization and augmented normalizing flows, can be expressed as SurVAE Flows. Finally, we introduce common operations such as the *max value*, the *absolute value*, *sorting* and *stochastic permutation* as composable layers in SurVAE Flows.

## 1 Introduction

Normalizing flows (Tabak and Vanden-Eijnden, 2010; Tabak and Turner, 2013; Rezende and Mohamed, 2015) provide a powerful *modular* and *composable* framework for representing expressive probability densities via differentiable bijections (with a differentiable inverse). These composable bijective transformations accord significant advantages due to their ability to be implemented using a modular software framework with a general interface consisting of three important components: (i) a forward transform, (ii) an inverse transform, and (iii) a log-likelihood contribution through the Jacobian determinant. Thus, significant advances have been made in recent years to develop novel flow modules that are easily invertible, expressive and computationally cheap (Dinh et al., 2015, 2017; Kingma et al., 2016; Papamakarios et al., 2017; Huang et al., 2018; Jaini et al., 2019a,b; Kingma and Dhariwal, 2018; Hoogeboom et al., 2019b, 2020; Durkan et al., 2019; van den Berg et al., 2018).

However, the bijective nature of the transformations used for building normalizing flows limit their ability to alter dimensionality, model discrete data and distributions with discrete structure or disconnected components. Specialized solutions have been developed to address these limitations independently. Uria et al. (2013); Ho et al. (2019) use dequantization to model discrete distributions using continuous densities, while Tran et al. (2019); Hoogeboom et al. (2019a) propose a discrete analog of normalizing flows. Cornish et al. (2019) use an augmented space to model an infinite mixtures of normalizing flows to address the problem of disconnected components whereas Huang et al. (2020); Chen et al. (2020) use a similar idea of augmentation of the observation space to model expressive distributions. VAEs (Kingma and Welling, 2014; Rezende et al., 2014), on the other hand, have no such limitations, but only provide lower bound estimates of the tractable estimates for the
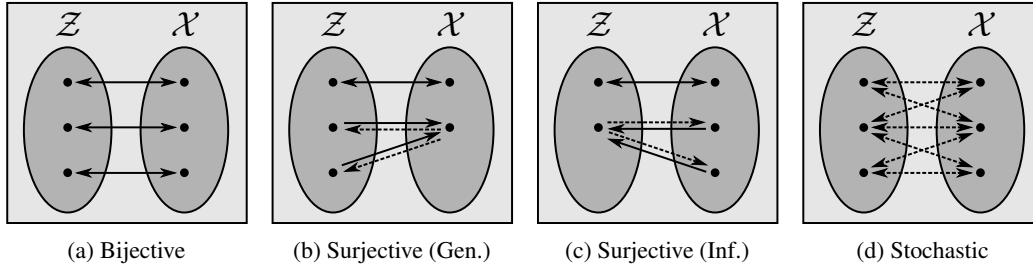
|           |                  |                  |                |
|-----------|------------------|------------------|----------------|
| (a) Bijective | (b) Surjective (Gen.) | (c) Surjective (Inf.) | (d) Stochastic |

Figure 1: Classes of SurVAE transformations $\mathcal{Z} \to \mathcal{X}$ and their inverses $\mathcal{X} \to \mathcal{Z}$. Solid lines indicate deterministic transformations, while dashed lines indicate stochastic transformations.

exact marginal density. These shortcomings motivate the question: *Is it possible to have composable and modular architectures that are expressive, model discrete and disconnected structure, and allow altering dimensions with exact likelihood evaluation?*

In this paper, we answer this affirmatively by introducing SurVAE Flows that use surjections to provide a unified, composable, and modular framework for probabilistic modeling. We introduce our unifying framework in §2 by identifying the components necessary to build composable architectures with modular software implementation for probabilistic modeling. We then introduce surjections for probabilistic modeling in §3 and show that these transformations lie at the interface between VAEs (stochastics maps) and normalizing flows (bijective maps). We unify these transformations (bijections, surjections, and stochastic transformations) in a composable and modular framework that we call SurVAE Flows. Subsequently, in §3.1, we propose novel SurVAE Flow layers like *max value* used for max pooling layers, *absolute value* for modelling symmetries in the data, and *sorting* and *stochastic permutations* that can be used for modelling exchangeable data and order statistics. Finally, in §3.2 we connect SurVAE Flows to several aforementioned specialised models by expressing them using SurVAE Flow layers which can now be implemented easily using our modular implementation. We demonstrate the efficacy of SurVAE Flows with experiments on synthetic datasets, point cloud data, and images. Code to implement SurVAE Flows and reproduce results is publicly available[1].

## 2 Preliminaries and Setup

In this section, we set up our main problem, provide key notations and definitions, and formulate a unifying framework for using different kinds of transformations to model distributions.

Let $\boldsymbol{x} \in \mathcal{X}$ and $\boldsymbol{z} \in \mathcal{Z}$ be two variables with distributions $p(\boldsymbol{x})$ and $p(\boldsymbol{z})$. We call a deterministic mapping $f : \mathcal{Z} \to \mathcal{X}$ *bijective* if it is both *surjective* and *injective*. A mapping is surjective if $\forall \boldsymbol{x} \in \mathcal{X}$, $\exists \boldsymbol{z} \in \mathcal{Z}$ such that $\boldsymbol{x} = f(\boldsymbol{z})$. A mapping is injective if $\forall \boldsymbol{z}_1, \boldsymbol{z}_2 \in \mathcal{Z}, f(\boldsymbol{z}_1) = f(\boldsymbol{z}_2) \implies \boldsymbol{z}_1 = \boldsymbol{z}_2$. If the mapping is not deterministic, we refer to it as a stochastic mapping, and denote it as $\boldsymbol{x} \sim p(\boldsymbol{x}|\boldsymbol{z})$.

Normalizing flows (Tabak and Vanden-Eijnden, 2010; Tabak and Turner, 2013; Rezende and Mohamed, 2015) make use of bijective transformations $f$ to transform a simple base density $p(\boldsymbol{z})$ to a more expressive density $p(\boldsymbol{x})$, making using the change-of-variables formula $p(\boldsymbol{x}) = p(\boldsymbol{z})|\det \nabla_{\boldsymbol{x}} f^{-1}(\boldsymbol{x})|$. VAEs (Kingma and Welling, 2014; Rezende et al., 2014), on the other hand, define a probabilistic graphical model where each observed variable $\boldsymbol{x}$ has an associated latent variable $\boldsymbol{z}$ with the generative process as $\boldsymbol{z} \sim p(\boldsymbol{z})$, $\boldsymbol{x} \sim p(\boldsymbol{x}|\boldsymbol{z})$, where $p(\boldsymbol{x}|\boldsymbol{z})$ may be viewed as a stochastic transformation. VAEs use variational inference with an amortized variational distribution $q(\boldsymbol{z}|\boldsymbol{x})$ to approximate the intractable posterior $p(\boldsymbol{z}|\boldsymbol{x})$ which facilitates computation of a lower bound of $p(\boldsymbol{x})$ known as the evidence lower bound (ELBO) i.e., $\mathcal{L} := \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x})}[\log p(\boldsymbol{x}|\boldsymbol{z})] - \mathbb{D}_{\mathsf{KL}}[q(\boldsymbol{z}|\boldsymbol{x})\|p(\boldsymbol{z})]$.

In the following, we introduce a framework to connect flows and VAEs[2] by showing that bijective and stochastic transformations are *composable* and require three important components for use in probabilistic modeling: (i) a forward transformation, $f : \mathcal{Z} \to \mathcal{X}$ with an associated conditional probability $p(\boldsymbol{x}|\boldsymbol{z})$, (ii) an inverse transformation, $f^{-1} : \mathcal{X} \to \mathcal{Z}$ with an associated distribution $q(\boldsymbol{z}|\boldsymbol{x})$, and (iii) a *likelihood contribution* term used for log-likelihood computation.

---

[1]The code is available at `https://github.com/didriknielsen/survae_flows`
[2]We note that Wu et al. (2020) also considered stochastic maps in flows using MCMC transition kernels.

Table 1: Composable building blocks of SurVAE Flows.

| Transformation | Forward $\boldsymbol{x} \leftarrow \boldsymbol{z}$ | Inverse $\boldsymbol{z} \leftarrow \boldsymbol{x}$ | Likelihood Contribution $\mathcal{V}(\boldsymbol{x}, \boldsymbol{z})$ | Bound Gap $\mathcal{E}(\boldsymbol{x}, \boldsymbol{z})$ |
|---|---|---|---|---|
| Bijective | $\boldsymbol{x} = f(\boldsymbol{z})$ | $\boldsymbol{z} = f^{-1}(\boldsymbol{x})$ | $\log \lvert \det \nabla_{\boldsymbol{x}} \boldsymbol{z} \rvert$ | $0$ |
| Stochastic | $\boldsymbol{x} \sim p(\boldsymbol{x}\mid\boldsymbol{z})$ | $\boldsymbol{z} \sim q(\boldsymbol{z}\mid\boldsymbol{x})$ | $\log \frac{p(\boldsymbol{x}\mid\boldsymbol{z})}{q(\boldsymbol{z}\mid\boldsymbol{x})}$ | $\log \frac{q(\boldsymbol{z}\mid\boldsymbol{x})}{p(\boldsymbol{z}\mid\boldsymbol{x})}$ |
| Surjective (Gen.) | $\boldsymbol{x} = f(\boldsymbol{z})$ | $\boldsymbol{z} \sim q(\boldsymbol{z}\mid\boldsymbol{x})$ | $\log \frac{p(\boldsymbol{x}\mid\boldsymbol{z})}{q(\boldsymbol{z}\mid\boldsymbol{x})}$ as $\begin{smallmatrix} p(\boldsymbol{x}\mid\boldsymbol{z}) \to \\ \delta(\boldsymbol{x} - f(\boldsymbol{z})) \end{smallmatrix}$ | $\log \frac{q(\boldsymbol{z}\mid\boldsymbol{x})}{p(\boldsymbol{z}\mid\boldsymbol{x})}$ |
| Surjective (Inf.) | $\boldsymbol{x} \sim p(\boldsymbol{x}\mid\boldsymbol{z})$ | $\boldsymbol{z} = f^{-1}(\boldsymbol{x})$ | $\log \frac{p(\boldsymbol{x}\mid\boldsymbol{z})}{q(\boldsymbol{z}\mid\boldsymbol{x})}$ as $\begin{smallmatrix} q(\boldsymbol{z}\mid\boldsymbol{x}) \to \\ \delta(\boldsymbol{z} - f^{-1}(\boldsymbol{x})) \end{smallmatrix}$ | $0$ |

**Forward Transformation:** For a stochastic transformation, the forward transformation is the conditional distribution $p(\boldsymbol{x}\mid\boldsymbol{z})$. For a bijective transformation, on the other hand, the forward transformation is deterministic and therefore, $p(\boldsymbol{x}\mid\boldsymbol{z}) = \delta\big(\boldsymbol{x} - f(\boldsymbol{z})\big)$ or simply $\boldsymbol{x} = f(\boldsymbol{z})$.

**Inverse Transformation:** For a bijective transformation, the inverse is also deterministic and given by $\boldsymbol{z} = f^{-1}(\boldsymbol{x})$. For a stochastic transformation, the inverse is also stochastic and is defined by Bayes theorem $p(\boldsymbol{z}\mid\boldsymbol{x}) = p(\boldsymbol{x}\mid\boldsymbol{z})p(\boldsymbol{z})/p(\boldsymbol{x})$. Computing $p(\boldsymbol{z}\mid\boldsymbol{x})$ is typically intractable and we thus resort to a variational approximation $q(\boldsymbol{z}\mid\boldsymbol{x})$.

**Likelihood Contribution:** For bijections, the density $p(\boldsymbol{x})$ can be computed from $p(\boldsymbol{z})$ and the mapping $f$ using the change-of-variables formula as:

$$\log p(\boldsymbol{x}) = \log p(\boldsymbol{z}) + \log \lvert \det \boldsymbol{J} \rvert, \qquad \boldsymbol{z} = f^{-1}(\boldsymbol{x}) \tag{1}$$

where $\lvert \det \boldsymbol{J} \rvert = \lvert \det \nabla_{\boldsymbol{x}} f^{-1}(\boldsymbol{x}) \rvert$ is the absolute value of the determinant of the Jacobian matrix of $f^{-1}$ which defines the likelihood contribution term for a bijective transformation $f$. For stochastic transformations, we can rewrite the marginal density $p(\boldsymbol{x})$ as:

$$\log p(\boldsymbol{x}) = \underbrace{\mathbb{E}_{q(\boldsymbol{z}\mid\boldsymbol{x})}[\log p(\boldsymbol{x}\mid\boldsymbol{z})] - \mathbb{D}_{\mathsf{KL}}[q(\boldsymbol{z}\mid\boldsymbol{x})\|p(\boldsymbol{z})]}_{\text{ELBO}} + \underbrace{\mathbb{D}_{\mathsf{KL}}[q(\boldsymbol{z}\mid\boldsymbol{x})\|p(\boldsymbol{z}\mid\boldsymbol{x})]}_{\text{Gap in Lower Bound}} \tag{2}$$

The ELBO $\mathcal{L}$ in Eq. 2 can then be evaluated using a single Monte Carlo sample: $\mathcal{L} \approx \log p(\boldsymbol{z}) + \log \frac{p(\boldsymbol{x}\mid\boldsymbol{z})}{q(\boldsymbol{z}\mid\boldsymbol{x})}$, $\boldsymbol{z} \sim q(\boldsymbol{z}\mid\boldsymbol{x})$. Therefore, the likelihood contribution term for a stochastic transformation is defined as $\log \frac{p(\boldsymbol{x}\mid\boldsymbol{z})}{q(\boldsymbol{z}\mid\boldsymbol{x})}$. Furthermore, we show in App. A that Eq. 2 allows us to recover the change-of-variables formula given in Eq. 1 by using Dirac delta functions, thereby drawing a precise connection between VAEs and normalizing flows. Crucially, Eq. 2 helps us to reveal a unified modular framework to model a density $p(\boldsymbol{x})$ under any transformation by restating it as:

$$\log p(\boldsymbol{x}) \simeq \log p(\boldsymbol{z}) + \mathcal{V}(\boldsymbol{x}, \boldsymbol{z}) + \mathcal{E}(\boldsymbol{x}, \boldsymbol{z}), \quad \boldsymbol{z} \sim q(\boldsymbol{z}\mid\boldsymbol{x}) \tag{3}$$

where $\mathcal{V}(\boldsymbol{x}, \boldsymbol{z})$ and $\mathcal{E}(\boldsymbol{x}, \boldsymbol{z})$ are the *likelihood contribution* and *bound looseness* terms, respectively. The likelihood contribution is $\mathcal{V}(\boldsymbol{x}, \boldsymbol{z}) = \log \lvert \det \boldsymbol{J} \rvert$ for bijections and $\log \frac{p(\boldsymbol{x}\mid\boldsymbol{z})}{q(\boldsymbol{z}\mid\boldsymbol{x})}$ for stochastic transformations. For bijections, likelihood evaluation is deterministic and *exact* with $\mathcal{E}(\boldsymbol{x}, \boldsymbol{z}) = 0$, while for stochastic maps it is stochastic and unbiased with a bound looseness of $\mathcal{E}(\boldsymbol{x}, \boldsymbol{z}) = \log \frac{q(\boldsymbol{z}\mid\boldsymbol{x})}{p(\boldsymbol{z}\mid\boldsymbol{x})}$. This is summarized in Table 1. The first term in Eq. 3, $\log p(\boldsymbol{z})$, reveals the compositional nature of the transformations, since it can be modeled by further transformations. While the compositional structure has been used widely for bijective transformations, Eq. 3 demonstrates its viability for stochastic maps as well. We demonstrate this unified compositional structure in Alg. 1.

---

**Algorithm 1:** $\log - \mathrm{likelihood}(\boldsymbol{x})$

**Data:** $\boldsymbol{x}, p(\boldsymbol{z})$ & $\{f_t\}_{t=1}^T$
**Result:** $\mathcal{L}(\boldsymbol{x})$
**for** $t$ *in* $\mathrm{range}(T)$, **do**
  **if** $f_t$ *is bijective* **then**
    $\boldsymbol{z} = f_t^{-1}(\boldsymbol{x})$ ;
    $\mathcal{V}_t = \log \lvert \det \frac{\partial \boldsymbol{z}}{\partial \boldsymbol{x}} \rvert$ ;
  **else if** $f_t$ *is stochastic* **then**
    $\boldsymbol{z} \sim q_t(\boldsymbol{z}\mid\boldsymbol{x})$ ;
    $\mathcal{V}_t = \log \frac{p_t(\boldsymbol{x}\mid\boldsymbol{z})}{q_t(\boldsymbol{z}\mid\boldsymbol{x})}$ ;
  $\boldsymbol{x} = \boldsymbol{z}$ ;
**end**
**return** $\log p(\boldsymbol{z}) + \sum_{t=1}^T \mathcal{V}_t$

# 3 SurVAE Flows

As explained in Section 2, bijective and stochastic transformations provide a modular framework for constructing expressive generative models. However, they both impose constraints on the model: bijective transformations are deterministic and allow exact likelihood computation, but they are required to preserve dimensionality. On the other hand, stochastic transformations are capable of altering the dimensionality of the random variables but only provide a stochastic lower bound estimate of the likelihood. *Is it possible to have composable transformations that can alter dimensionality and allow exact likelihood evaluation?* In this section, we answer this question affirmatively by introducing surjective transformations as SurVAE Flows that bridge the gap between bijective and stochastic transformations.

In the following, we will define composable deterministic transformations that are surjective and non-injective. For brevity, we will refer to them as *surjections* or *surjective transformations*. Note that for surjections, multiple inputs can map to a single output, resulting in a *loss of information* since the input is not guaranteed to be recovered through inversion. Similar to bijective and stochastic transformations, the three important components of composable surjective transformations are:

**Forward Transformation:** Like bijections, surjective transformations have a deterministic forward transformation $p(\boldsymbol{x}|\boldsymbol{z}) = \delta(\boldsymbol{x} - f(\boldsymbol{z}))$ or $\boldsymbol{x} = f(\boldsymbol{z})$.

**Inverse Transformation:** In contrast with bijections, surjections $f : \mathcal{Z} \to \mathcal{X}$ are not invertible since multiple inputs can map to the same output. However, they have *right inverses*, i.e. functions $g : \mathcal{X} \to \mathcal{Z}$ such that $f \circ g(\boldsymbol{x}) = \boldsymbol{x}$, but not necessarily $g \circ f(\boldsymbol{z}) = \boldsymbol{z}$. We will use a stochastic right inverse $q(\boldsymbol{z}|\boldsymbol{x})$ which can be thought of as passing $\boldsymbol{x}$ through a random right inverse $g$. Importantly, $q(\boldsymbol{z}|\boldsymbol{x})$ only has support over the preimage of $\boldsymbol{x}$, i.e. the set of $\boldsymbol{z}$ that map to $\boldsymbol{x}$, $\mathcal{B}(\boldsymbol{x}) = \{\boldsymbol{z}|\boldsymbol{x} = f(\boldsymbol{z})\}$.

So far, we have described what we will term *generative surjections*, i.e. transformations that are surjective in the generative direction $\mathcal{Z} \to \mathcal{X}$. We will refer to a transformation which is surjective in the inference direction $\mathcal{X} \to \mathcal{Z}$ as an *inference surjection*. These are illustrated in Fig.1. Generative surjections have stochastic inverse transformations $q(\boldsymbol{z}|\boldsymbol{x})$, while inference surjections have stochastic forward $p(\boldsymbol{x}|\boldsymbol{z})$ transformations.

**Likelihood Contribution:** For continuous surjections, the likelihood contribution term is:

$$\mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x})}\left[\log \frac{p(\boldsymbol{x}|\boldsymbol{z})}{q(\boldsymbol{z}|\boldsymbol{x})}\right], \quad \text{as} \quad \begin{cases} p(\boldsymbol{x}|\boldsymbol{z}) \to \delta(\boldsymbol{x} - f(\boldsymbol{z})), & \text{for gen. surjections.} \\ q(\boldsymbol{z}|\boldsymbol{x}) \to \delta(\boldsymbol{z} - f^{-1}(\boldsymbol{x})), & \text{for inf. surjections.} \end{cases}$$

While generative surjections generally give rise to stochastic estimates of the likelihood contribution and introduce lower bound likelihood estimates, inference surjections allow *exact likelihood computation* (see App. B). Before proceeding further, we give a few examples to better understand the construction of a surjective transformation for probabilistic modeling.

**Example 1 (Tensor slicing)** *Let $f$ be a tensor slicing surjection that takes input $\boldsymbol{z} = (\boldsymbol{z}_1, \boldsymbol{z}_2) \in \mathbb{R}^{d_z}$ and returns a subset of the elements, i.e. $\boldsymbol{x} = f(\boldsymbol{z}) = \boldsymbol{z}_1$. To develop this operation as a SurVAE layer, we first specify the stochastic forward and inverse transformations as:*

$$p(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}|\boldsymbol{z}_1, \sigma^2 \boldsymbol{I}), \quad \text{and} \quad q(\boldsymbol{z}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{z}_1|\boldsymbol{x}, \sigma^2 \boldsymbol{I})q(\boldsymbol{z}_2|\boldsymbol{x})$$

*We next compute the likelihood contribution term in the limit that $p(\boldsymbol{x}|\boldsymbol{z}) \to \delta(\boldsymbol{x} - f(\boldsymbol{z}))$. Here, this corresponds to $\sigma \to 0$. Thus,*

$$\mathcal{V}(\boldsymbol{x}, \boldsymbol{z}) = \lim_{\sigma^2 \to 0} \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x})}\left[\log \frac{p(\boldsymbol{x}|\boldsymbol{z})}{q(\boldsymbol{z}|\boldsymbol{x})}\right] = \mathbb{E}_{q(\boldsymbol{z}_2|\boldsymbol{x})}\left[-\log q(\boldsymbol{z}_2|\boldsymbol{x})\right],$$

*which corresponds to the entropy of $q(\boldsymbol{z}_2|\boldsymbol{x})$ that is used to infer the sliced elements $\boldsymbol{z}_2$. We illustrate the slicing surjection for both the generative and inference directions in Fig. 2.*

**Example 2 (Rounding)** *Let $f$ be a rounding surjection that takes an input $\boldsymbol{z} \in \mathbb{R}^{d_z}$ and returns the rounded $\boldsymbol{x} := \lfloor \boldsymbol{z} \rfloor$. The forward transformation is a discrete surjection $P(\boldsymbol{x}|\boldsymbol{z}) = \mathbb{I}(\boldsymbol{z} \in \mathcal{B}(\boldsymbol{x}))$, for $\mathcal{B}(\boldsymbol{x}) = \{\boldsymbol{x} + \boldsymbol{u}|\boldsymbol{u} \in [0, 1)^d\}$. The inverse transformation $q(\boldsymbol{z}|\boldsymbol{x})$ is stochastic with support in $\mathcal{B}(\boldsymbol{x})$. Inserting this in the likelihood contribution term and simplifying, we find*

$$\mathcal{V}(\boldsymbol{x}, \boldsymbol{z}) = \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x})}\left[-\log q(\boldsymbol{z}|\boldsymbol{x})\right].$$

*This generative rounding surjection gives rise to dequantization (Uria et al., 2013; Ho et al., 2019) which is a method commonly used to train continuous flows on discrete data such as images.*

Table 2: Summary of selected inference surjection layers. See App. C for more SurVAE layers.

| Surjection | Forward | Inverse | $\mathcal{V}(\boldsymbol{x}, \boldsymbol{z})$ |
|---|---|---|---|
| Abs | $s \sim \text{Bern}(\pi(z))$ <br> $x = s \cdot z, \ \ s \in \{-1, 1\}$ | $s = \text{sign } x$ <br> $z = \|x\|$ | $\log p(s\|z)$ |
| Max | $k \sim \text{Cat}(\boldsymbol{\pi}(z))$ <br> $x_k = z, \boldsymbol{x}_{-k} \sim p(\boldsymbol{x}_{-k}\|z, k)$ | $k = \arg \max \boldsymbol{x}$ <br> $z = \max \boldsymbol{x}$ | $\log p(k\|z) + \log p(\boldsymbol{x}_{-k}\|z, k)$ |
| Sort | $\mathcal{I} \sim \text{Cat}(\boldsymbol{\pi}(\boldsymbol{z}))$ <br> $\boldsymbol{x} = \boldsymbol{z}_{\mathcal{I}}$ | $\mathcal{I} = \text{argsort } \boldsymbol{x}$ <br> $\boldsymbol{z} = \text{sort } \boldsymbol{x}$ | $\log p(\mathcal{I}\|\boldsymbol{z})$ |

The preceding discussion shows that surjective transformations can be composed to construct expressive transformations for density modelling. We call a single surjective transformation a SurVAE layer and a composition of bijective, surjective, and/or stochastic transformations a SurVAE Flow. The unified framework of SurVAE Flows allows us to construct generative models learned using the likelihood (or its lower bound) of the data, utilizing Eq. 3, and Table 1.

## 3.1 Novel SurVAE Layers

We developed the tensor slicing and rounding surjections in Examples 1 and 2. In this section, we introduce additional novel SurVAE layers including the *absolute value*, the *maximum value* and *sorting* as surjective layers and *stochastic permutation* as a stochastic layer. We provide a summary of these in Table 2. Due to space constraints, we defer the derivations and details on each of these surjections to Appendix D-G along with detailed tables on generative and inference surjections in Table 6 and 7.

**Abs Surjection** (App. D). The abs surjection returns the the magnitude of its input, $z = |x|$. As a SurVAE layer, we can represent the inference surjection with the forward and inverse transformations as:

$$p(x|z) = \sum_{s \in \{-1,1\}} p(\boldsymbol{x}|z, s)p(s|z) = \sum_{s \in \{-1,1\}} \delta(x - sz)p(s|z),$$

$$q(z|x) = \sum_{s \in \{-1,1\}} q(z|\boldsymbol{x}, s)p(s|x) = \sum_{s \in \{-1,1\}} \delta(z - sx)\delta_{s,\text{sign}(x)}$$

where $q(z|x)$ is deterministic corresponding to $z = |x|$. The forward transformation $p(x|z)$ involves the following steps: (i) sample the sign $s$, conditioned on $z$, and (ii) apply the sign to $z$ to obtain $x = sz$. Abs surjections are useful for modelling data with symmetries which we demonstrate in our experiments.

**Max Surjection** (App. E, Fig. 2). The max operator returns the largest element of an input vector, $z = \max \boldsymbol{x}$. We can represent this transformation as



(a) Gen. slicing



(b) Inf. slicing



(c) Inf. max

Figure 2: Surjections.

$$p(\boldsymbol{x}|z) = \sum_{k=1}^{K} p(\boldsymbol{x}|z, k)p(k|z) = \sum_{k=1}^{K} \delta(x_k - z)p(\boldsymbol{x}_{-k}|z, k)p(k|z),$$

$$q(z|\boldsymbol{x}) = \sum_{k=1}^{K} q(z|\boldsymbol{x}, k)q(k|x) = \sum_{k=1}^{K} \delta(z - x_k)\delta_{k,\arg\max(\boldsymbol{x})},$$

where $q(z|x)$ is deterministic and corresponds to $z = \max \boldsymbol{x}$. While the inverse is deterministic The stochastic forward proceeds by (i) sampling an index $k$ and setting $x_k = z$, and (ii) imputing the remaining values $\boldsymbol{x}_{-k}$ of $\boldsymbol{x}$ such that they are all smaller than $x_k$. Max surjections are useful in implementing the *max pooling* layer commonly used in convolutional architectures for downsampling. In our experiments, we demonstrate the use of max surjections for probabilistic modelling of images.

**Sort Surjection** (App. F). Sorting, $\boldsymbol{z} = \text{sort}(\boldsymbol{x})$ returns a vector in sorted order. It is a surjective (and non-injective) transformation since the original order of the vector is lost in the operation even
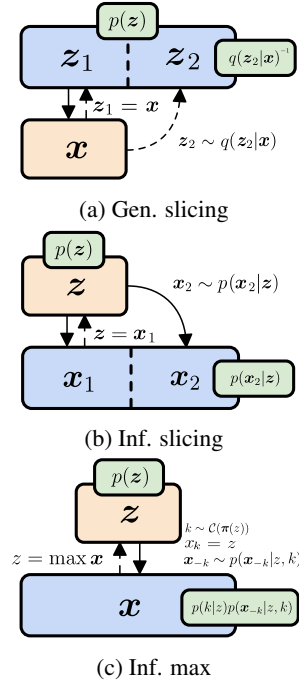
though the dimensions remain the same. Sort surjections are useful in modelling naturally sorted data, learning order statistics, and learning an exchangeable model using flows.

**Stochastic Permutation** (App. G). A stochastic permutation transforms the input vector by shuffling the elements randomly. The inverse pass for a permutation is the same as the forward pass with the likelihood contribution term equal to zero, $\mathcal{V} = 0$. Stochastic permutations helps to enforce permutation invariance i.e. any flow can be made permutation invariant by adding a final permutation SurVAE layer. In our experiments, we compare sorting surjections and stochastic permutations to enforce permutation invariance for modelling exchangeable data.

**Stochastic Inverse Parameterization**. For surjections, the stochastic inverses have to be defined so that they form a distribution over the possible right-inverses. Different right-inverse distributions do not have to align over subsets of the space. Consequently, for more sophisticated choices of right-inverses, the log-likelihood may be discontinuous across boundaries of these subsets. Since these points have measure zero, this does not influence the validity of the log-likelihood. However, it may impede optimization using gradient-based methods. In our experiments, we did not encounter any specific issues, but for a more thorough discussion see (Dinh et al., 2019).

## 3.2 Connection to Previous Work

The results above provide a unified framework based on SurVAE Flows for estimating probability densities. We now connect this general approach to several recent works on generative modelling.

The differentiable and bijective nature of transformations used in normalizing flows limit their ability to alter dimensionality, model discrete data, and distributions with disconnected components. Specialized solutions have been proposed in recent years to address these individually. We now show that these works can be expressed using SurVAE Flow layers, as summarized in Table 3.

### 3.2.1 Using Stochastic Transformations

As discussed in Sec. 2, VAEs (Kingma and Welling, 2014; Rezende et al., 2014) may be formulated as composable stochastic transformations. Probabilistic PCA Tipping and Bishop (1999) can be considered a simple special case of VAEs wherein the forward transformation is linear-Gaussian, i.e. $p(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z}, \sigma^2 \boldsymbol{I})$. Due to the linear-Gaussian formulation, the posterior $p(\boldsymbol{z}|\boldsymbol{x})$ is tractable and we can thus perform exact stochastic inversion for this model. Diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020) are another class of models closely related to VAEs. For diffusion models, the inverse $q(\boldsymbol{z}|\boldsymbol{x})$ implements a diffusion step, while the forward transformation $p(\boldsymbol{x}|\boldsymbol{z})$ learns to reverse the diffusion process. Wu et al. (2020) propose an extended flow framework consisting of bijective and stochastic transformations using MCMC transition kernels. Their method utilizes the same computation as in the general formulation in Algorithm 1, but does not consider surjective maps or an explicit connection to VAEs. Their work shows that MCMC kernels may also be implemented as stochastic transformations in SurVAE Flows.

Table 3: SurVAE Flows as a unifying framework.

| Model | SurVAE Flow architecture |
|---|---|
| Probabilistic PCA (Tipping and Bishop, 1999) VAE (Kingma and Welling, 2014; Rezende et al., 2014) Diffusion Models (Sohl-Dickstein et al., 2015; Ho et al., 2020) | $\mathcal{Z} \xrightarrow{\text{stochastic}} \mathcal{X}$ |
| Dequantization (Uria et al., 2013; Ho et al., 2019) | $\mathcal{Z} \xrightarrow{\text{round}} \mathcal{X}$ |
| ANFs, VFlow (Huang et al., 2020; Chen et al., 2020) | $\mathcal{X} \xrightarrow{\text{augment}} \mathcal{X} \times \mathcal{E} \xrightarrow{\text{bijection}} \mathcal{Z}$ |
| Multi-scale Architectures (Dinh et al., 2017) | $\mathcal{X} \xrightarrow{\text{bijection}} \mathcal{Y} \times \mathcal{E} \xrightarrow{\text{slice}} \mathcal{Y} \xrightarrow{\text{bijection}} \mathcal{Z}$ |
| CIFs, Discretely Indexed Flows, DeepGMMs (Cornish et al., 2019; Duan, 2019; Oord and Dambre, 2015) | $\mathcal{X} \xrightarrow{\text{augment}} \mathcal{X} \times \mathcal{E} \xrightarrow{\text{bijection}} \mathcal{Z} \times \mathcal{E} \xrightarrow{\text{slice}} \mathcal{Z}$ |
| RAD Flows (Dinh et al., 2019) | $\mathcal{X} \xrightarrow{\text{partition}} \mathcal{X}_{\mathcal{E}} \times \mathcal{E} \xrightarrow{\text{bijection}} \mathcal{Z} \times \mathcal{E} \xrightarrow{\text{slice}} \mathcal{Z}$ |

### 3.2.2   Using Surjective Transformations

Dequantization (Uria et al., 2013; Ho et al., 2019) is used for training continuous flow models on ordinal discrete data such as images and audio. Dequantization fits into the SurVAE Flow framework as a composable generative rounding surjection (cf. Example 2) and thus simplifies implementation. When the inverse $q(\boldsymbol{z}|\boldsymbol{x})$ is a standard uniform distributon, *uniform dequantization* is obtained, while a more flexible flow-based distribution $q(\boldsymbol{z}|\boldsymbol{x})$ yields *variational dequantization* (Ho et al., 2019).

VFlow (Chen et al., 2020) and ANFs (Huang et al., 2020) aim to build expressive generative models by augmenting the data space and jointly learning a normalizing flow for the augmented data space as well as the distribution of augmented dimensions. This strategy was also adopted by Dupont et al. (2019) for continuous-time flows. VFlow and ANFs can be obtained as SurVAE Flows by composing a bijection with a generative tensor slicing surjection (cf. Example 1 and Figure 3a). The reverse transformation, i.e. inference slicing, results in the *multi-scale architecture* of Dinh et al. (2017).

CIFs (Cornish et al., 2019) use an indexed family of bijective transformations $g(\cdot; \varepsilon) : \mathcal{Z} \to \mathcal{X}$ where $\mathcal{Z} = \mathcal{X} \subseteq \mathbb{R}^d$, and $\varepsilon \in \mathcal{E} \subseteq \mathbb{R}^{d_\varepsilon}$ with the generative process as: $\boldsymbol{z} \sim p(\boldsymbol{z})$, $\epsilon \sim p(\epsilon|\boldsymbol{z})$ and $\boldsymbol{x} = g(\boldsymbol{z}; \varepsilon)$ and requires specifying $p(\boldsymbol{z})$ and $p(\varepsilon|\boldsymbol{z})$. CIFs are akin to modeling densities using an infinite mixture of normalizing flows since $g$ is a surjection from an augmented space $\mathcal{Z} \times \mathcal{E}$ to the data space $\mathcal{X}$. Consequently, CIFs can be expressed as a SurVAE flow using a augment surjection composed with a bijection and tensor slicing (cf. Figure 3b). Similarly, Duan (2019) used a *finite* mixture of normalizing flows to model densities by using a discrete index set $\mathcal{E} = \{1, 2, 3, \cdots, K\}$ with bijections $g(\cdot; \varepsilon)$. Deep Gaussian mixture models (Oord and Dambre, 2015) form special case wherein the bijections $g(\cdot; \varepsilon)$ are linear transformations.

RAD flows (Dinh et al., 2019) are also "similar" to CIFs but it partitions the data space into finite disjoint subsets $\{\mathcal{X}_i\}_{i=1}^K \subseteq \mathcal{X}$ and defines bijections $g_i : \mathcal{Z} \to \mathcal{X}_i, \forall\, i \in \{1, 2, ..., K\}$ with the generative process as $\boldsymbol{z} \sim p(\boldsymbol{z}), i \sim p(i|\boldsymbol{z})$ and $\boldsymbol{x} = g_i(\boldsymbol{z})$. Interestingly, RAD can be seen to implement a class of inference surjections that rely on partitioning of the data space. The partitioning is learned during training, thus allowing learning of expressive inference surjections. However, careful parameterization is required for stable gradient-based training. We note that the abs and sort inference surjections introduced earlier may be expressed using a static (non-learned) partitioning of the data space $\mathcal{X}$ and thus have close ties to RAD. However, RAD does not express generative surjections or more general inference surjections that do not rely on partitioning, such as dimensional changes.

Finally, we note that apart from providing a general method for modelling densities, SurVAE Flows provide a modular framework for easy implementation of the methods described here. We discuss these important software perspectives using code snippets in App. H.
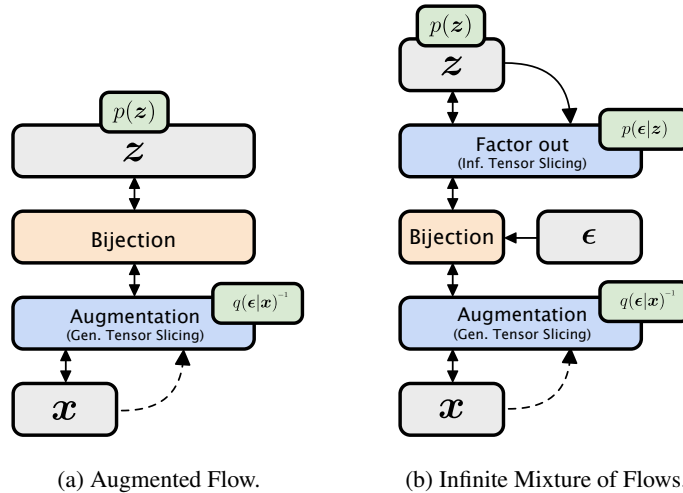


(a) Augmented Flow.          (b) Infinite Mixture of Flows.

Figure 3: Flow architectures making use of tensor slicing.

# 4 Experiments

We investigate the ability of SurVAE flows to model data that is difficult to model with normalizing flows. We show that the absolute value surjection is useful in modelling data where certain symmetries are known to exist. Next, we demonstrate that SurVAE flows allow straightforward modelling of exchangeable data by simply composing *any* flow together with either a sorting surjection or a stochastic permutation layer. Furthermore, we investigate the use of max pooling – which is commonly used for downsampling in convolutional neural networks – as a surjective downsampling layer in SurVAE flows for image data.

**Synthetic Data.** We first consider modelling data where certain symmetries are known to exist. We make use of 3 symmetric and 1 anti-symmetric synthetic 2D datasets. The absolute value inference surjection can be seen to fold the input space across the origin and can thus be useful in modelling such data. The baseline uses 4 coupling bijections, while our *AbsFlow* adds an extra `abs` surjection. For the anti-symmetric data, AbsFlow uses only a single `abs` surjection with a classifier (*i.e.* for $P(s|z)$) which learns the unfolding. For further details, see App. I.1. The results are shown in Fig. 4.

**Point Cloud Data.** We now consider modelling exchangeable data. We use the `SpatialMNIST` dataset (Edwards and Storkey, 2017), where each MNIST digit is represented as a 2D point cloud of 50 points. A point cloud is a set, *i.e.* it is permutation invariant. Using SurVAE flows, we can enforce permutation invariance on *any* flow using either 1) a sorting surjection – forcing a canonical order on the inputs, or 2) a stochastic permutation – forcing a random order on the inputs.



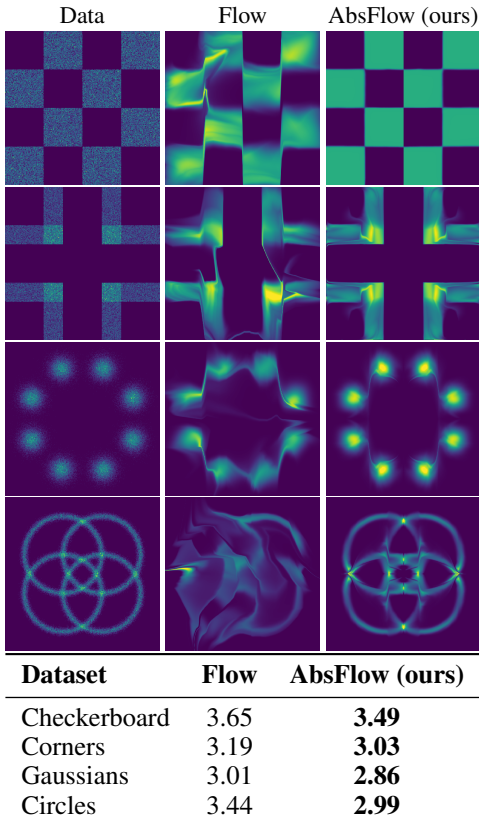| Dataset | Flow | AbsFlow (ours) |
|---|---|---|
| Checkerboard | 3.65 | **3.49** |
| Corners | 3.19 | **3.03** |
| Gaussians | 3.01 | **2.86** |
| Circles | 3.44 | **2.99** |

Figure 4: Comparison of flows with and without absolute value surjections modelling anti-symmetric (top row) and symmetric (3 bottom rows) 2-dimensional distributions.
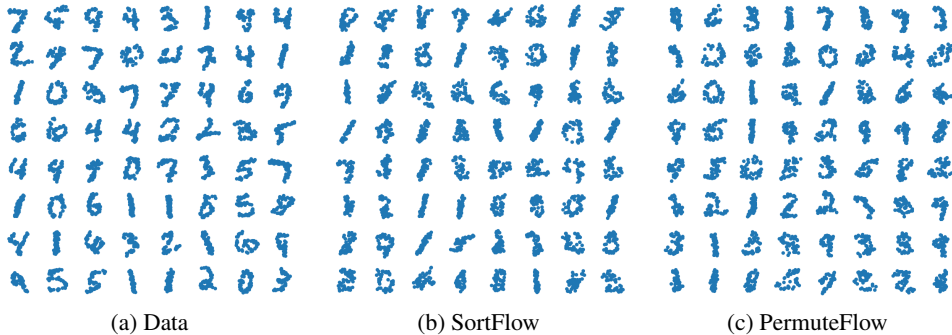
We compare 2 SurVAE flows, *SortFlow* and *PermuteFlow*, both using 64 layers of coupling flows parameterized by Transformer networks (Vaswani et al., 2017). Transformers are – when *not* using positional encoding – permutation equivariant. PermuteFlow uses stochastic permutation in-between the coupling layers. SortFlow, on the other hand, uses and initial sorting surjection, which introduces an ordering, and fixed permutations after. The Transformers thus make use of learned positional encodings for SortFlow, but not for PermuteFlow. See App. I.2 for further details, and Fig. 5 for



|         (a) Data          |        (b) SortFlow         |       (c) PermuteFlow        |

Figure 5: Point cloud samples from permutation-invariant SurVAE flows trained on `SpatialMNIST`.

Table 4: Unconditional image modeling results in bits/dim.

| Model | CIFAR-10 | ImageNet32 | ImageNet64 |
|---|---|---|---|
| RealNVP (Dinh et al., 2017) | 3.49 | 4.28 | - |
| Glow (Kingma and Dhariwal, 2018) | 3.35 | 4.09 | 3.81 |
| Flow++ (Ho et al., 2019) | 3.08 | 3.86 | 3.69 |
| Baseline (Ours) | **3.08** | **4.00** | **3.70** |
| MaxPoolFlow (Ours) | 3.09 | 4.01 | 3.74 |



Figure 7: Samples from CIFAR-10 models. Top: MaxPoolFlow, Bottom: Baseline.

| Model | Inception ↑ | FID ↓ |
|---|---|---|
| DCGAN* | 6.4 | 37.1 |
| WGAN-GP* | 6.5 | 36.4 |
| PixelCNN* | 4.60 | 65.93 |
| PixelIQN* | 5.29 | 49.46 |
| Baseline (Ours) | 5.08 | 49.56 |
| MaxPoolFlow (Ours) | **5.18** | **49.03** |

Table 5: Inception score and FID for CIFAR-10. *Results taken from Ostrovski et al. (2018).

model samples. Interestingly, PermuteFlow outperforms SortFlow, with -5.30 vs. -5.53 PPLL (per-point log-likelihood), even though it only allows computation of lower bound likelihood estimates. For comparison, BRUNO (Korshunova et al., 2018) and FlowScan (Bender et al., 2020) obtain -5.68 and -5.26 PPLL, but make use of autoregressive components. Neural Statistican (Edwards and Storkey, 2017) utilizes hierarchical latent variables without autoregressive parts and obtains -5.37 PPL. PermuteFlow thus obtains *state-of-the-art* performance among non-autoregressive models.

**Image Data.** Max pooling layers are commonly used for downsampling in convolutional neural networks. We investigate their use as surjective downsampling transformations in flow models for image data here.

We train a flow using 2 scales with 12 steps/scale for `CIFAR-10` and `ImageNet 32×32` and 3 scales with 8 steps/scale for `ImageNet 64×64`. Each step consists of an affine coupling bijection and a $1 \times 1$ convolution (Kingma and Dhariwal, 2018). We implement a max pooling surjection for downscaling and compare it to a baseline model with tensor slicing which corresponds to a *multi-scale* architecture (Dinh et al., 2017). We report results for the log-likelihood in Table 4 and the inception and FID scores in Table 5 with bolding indicating best among the baseline and MaxPoolFlow. The results show that compared to slicing surjections, the max pooling surjections yield marginally worse log-likelihoods, but better visual sample quality as measured by the Inception score and FID. We also provide the generated samples from our models in Fig. 7 and App. J. Due to space constrains, we refer the reader to App. I.3 for more details on the experiment.
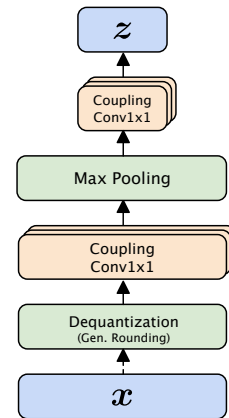


Figure 6: Flow architecture with max pooling. Surjections in green.

## 5 Conclusion

We introduced SurVAE flows, a modular framework for constructing likelihood-based models using composable bijective, surjective and stochastic transformations. We showed how this encompasses normalizing flows, which rely on bijections, as well as VAEs, which rely on stochastic transformations. We further showed that several recently proposed methods such as dequantization and augmented normalizing flows may be obtained as SurVAE flows using surjective transformations. One interesting direction for further research is development of novel non-bijective transformations that might be beneficial as composable layers in SurVAE flows.

## Acknowledgements

## Broader Impact

This work constitutes foundational research on generative models/unsupervised learning by providing a unified view on several lines of work and further by introducing new modules that expand the generative modelling toolkit. This work further suggests how to build software libraries to that allows more rapid implementation of a wider range of deep unsupervised models. Unsupervised learning has the potential to greatly reduce the need for labeled data and thus improve models in applications such as medical imaging where a lack of data can be a limitation. However, it may also potentially be used to improve deep fakes with potentially malicious applications.

## References

Bender, C. M., O'Connor, K., Li, Y., Garcia, J. J., Oliva, J., and Zaheer, M. (2020). Exchangeable generative models with flow scans. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence*, pages 10053–10060.

Bingham, E., Chen, J. P., Jankowiak, M., Obermeyer, F., Pradhan, N., Karaletsos, T., Singh, R., Szerlip, P., Horsfall, P., and Goodman, N. D. (2018). Pyro: Deep Universal Probabilistic Programming. *Journal of Machine Learning Research*.

Burda, Y., Grosse, R. B., and Salakhutdinov, R. (2016). Importance weighted autoencoders. In *4th International Conference on Learning Representations*.

Chen, J., Lu, C., Chenli, B., Zhu, J., and Tian, T. (2020). Vflow: More expressive generative flows with variational data augmentation. *CoRR*, abs/2002.09741.

Cornish, R., Caterini, A. L., Deligiannidis, G., and Doucet, A. (2019). Localised generative flows. *CoRR*, abs/1909.13833.

Dillon, J. V., Langmore, I., Tran, D., Brevdo, E., Vasudevan, S., Moore, D., Patton, B., Alemi, A., Hoffman, M. D., and Saurous, R. A. (2017). Tensorflow distributions. *CoRR*, abs/1711.10604.

Dinh, L., Krueger, D., and Bengio, Y. (2015). NICE: Non-linear independent components estimation. In *3rd International Conference on Learning Representations, ICLR, Workshop Track Proceedings*.

Dinh, L., Sohl-Dickstein, J., and Bengio, S. (2017). Density estimation using real NVP. In *5th International Conference on Learning Representations*.

Dinh, L., Sohl-Dickstein, J., Pascanu, R., and Larochelle, H. (2019). A rad approach to deep mixture models. *arXiv preprint arXiv:1903.07714*.

Duan, L. L. (2019). Transport monte carlo. *arXiv preprint arXiv:1907.10448*.

Dupont, E., Doucet, A., and Teh, Y. W. (2019). Augmented neural odes. In *Advances in Neural Information Processing Systems*, pages 3134–3144.

Durkan, C., Bekasov, A., Murray, I., and Papamakarios, G. (2019). Neural spline flows. In *Advances in Neural Information Processing Systems*, pages 7509–7520.

Edwards, H. and Storkey, A. J. (2017). Towards a neural statistician. In *5th International Conference on Learning Representations*.

Ho, J., Chen, X., Srinivas, A., Duan, Y., and Abbeel, P. (2019). Flow++: Improving flow-based generative models with variational dequantization and architecture design. In *Proceedings of the 36th International Conference on Machine Learning*.

Ho, J., Jain, A., and Abbeel, P. (2020). Denoising diffusion probabilistic models. *CoRR*, abs/2006.11239.

Hoogeboom, E., Peters, J. W. T., van den Berg, R., and Welling, M. (2019a). Integer discrete flows and lossless compression. In *Advances in Neural Information Processing Systems*, pages 12134–12144.

Hoogeboom, E., Satorras, V. G., Tomczak, J. M., and Welling, M. (2020). The convolution exponential and generalized sylvester flows. *CoRR*, abs/2006.01910.

Hoogeboom, E., van den Berg, R., and Welling, M. (2019b). Emerging convolutions for generative normalizing flows. In *Proceedings of the 36th International Conference on Machine Learning*.

Huang, C., Dinh, L., and Courville, A. C. (2020). Augmented normalizing flows: Bridging the gap between generative flows and latent variable models. *CoRR*, abs/2002.07101.

Huang, C., Krueger, D., Lacoste, A., and Courville, A. C. (2018). Neural autoregressive flows. In *Proceedings of the 35th International Conference on Machine Learning*, pages 2083–2092.

Huang, G., Liu, Z., van der Maaten, L., and Weinberger, K. Q. (2017). Densely connected convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2261–2269.

Jaini, P., Kobyzev, I., Brubaker, M., and Yu, Y. (2019a). Tails of triangular flows. *arXiv preprint arXiv:1907.04481*.

Jaini, P., Selby, K. A., and Yu, Y. (2019b). Sum-of-squares polynomial flow. In *Proceedings of the 36th International Conference on Machine Learning*, pages 3009–3018.

Kingma, D. P. and Ba, J. (2015). Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations*.

Kingma, D. P. and Dhariwal, P. (2018). Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*, pages 10236–10245.

Kingma, D. P., Salimans, T., Józefowicz, R., Chen, X., Sutskever, I., and Welling, M. (2016). Improving variational autoencoders with inverse autoregressive flow. In *Advances in Neural Information Processing Systems*, pages 4736–4744.

Kingma, D. P. and Welling, M. (2014). Auto-encoding variational bayes. In *2nd International Conference on Learning Representations*.

Korshunova, I., Degrave, J., Huszar, F., Gal, Y., Gretton, A., and Dambre, J. (2018). BRUNO: A deep recurrent model for exchangeable data. In *Advances in Neural Information Processing Systems*, pages 7190–7198.

Oord, A. v. d. and Dambre, J. (2015). Locally-connected transformations for deep gmms. In *International Conference on Machine Learning (ICML): Deep learning Workshop*, pages 1–8.

Ostrovski, G., Dabney, W., and Munos, R. (2018). Autoregressive quantile networks for generative modeling. In *Proceedings of the 35th International Conference on Machine Learning*, pages 3933–3942.

Papamakarios, G., Murray, I., and Pavlakou, T. (2017). Masked autoregressive flow for density estimation. In *Advances in Neural Information Processing Systems*, pages 2338–2347.

Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035.

Rezende, D. J. and Mohamed, S. (2015). Variational inference with normalizing flows. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pages 1530–1538.

Rezende, D. J., Mohamed, S., and Wierstra, D. (2014). Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31th International Conference on Machine Learning*, pages 1278–1286.

Sohl-Dickstein, J., Weiss, E. A., Maheswaranathan, N., and Ganguli, S. (2015). Deep unsupervised learning using nonequilibrium thermodynamics. In Bach, F. R. and Blei, D. M., editors, *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 2256–2265. JMLR.org.

Tabak, E. G. and Turner, C. V. (2013). A family of nonparametric density estimation algorithms. *Communications on Pure and Applied Mathematics*, 66(2):145–164.

Tabak, E. G. and Vanden-Eijnden, E. (2010). Density estimation by dual ascent of the log-likelihood. *Communications in Mathematical Sciences*, 8(1):217–233.

Theis, L., van den Oord, A., and Bethge, M. (2016). A note on the evaluation of generative models. In *International Conference on Learning Representations*.

Tipping, M. E. and Bishop, C. M. (1999). Probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(3):611–622.

Tran, D., Vafa, K., Agrawal, K. K., Dinh, L., and Poole, B. (2019). Discrete flows: Invertible generative models of discrete data. In *Advances in Neural Information Processing Systems*, pages 14692–14701.

Uria, B., Murray, I., and Larochelle, H. (2013). Rnade: The real-valued neural autoregressive density-estimator. In *Advances in Neural Information Processing Systems*, pages 2175–2183.

Uria, B., Murray, I., and Larochelle, H. (2014). A deep and tractable density estimator. In *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32*.

van den Berg, R., Hasenclever, L., Tomczak, J. M., and Welling, M. (2018). Sylvester normalizing flows for variational inference. In *Proceedings of the Thirty-Fourth Conference on Uncertainty in Artificial Intelligence*, pages 393–402.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems 30*, pages 5998–6008.

Wu, H., Köhler, J., and Noé, F. (2020). Stochastic normalizing flows. In *Advances in Neural Information Processing Systems*.

## A    A Connection Between VAEs and Flows

Variational Autoencoders (VAEs) can be seen as a composable stochastic transformations. From this viewpoint, the log-likelihood resulting from a single transformation can be written as

$$\log p(\boldsymbol{x}) = \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x})}\left[\log p(\boldsymbol{z})\right] + \underbrace{\mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x})}\left[\log \frac{p(\boldsymbol{x}|\boldsymbol{z})}{q(\boldsymbol{z}|\boldsymbol{x})}\right]}_{\text{Lik. contrib. } \mathcal{V}(\boldsymbol{x},\boldsymbol{z})} + \underbrace{\mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x})}\left[\log \frac{q(\boldsymbol{z}|\boldsymbol{x})}{p(\boldsymbol{z}|\boldsymbol{x})}\right]}_{\text{Bound looseness } \mathcal{E}(\boldsymbol{x},\boldsymbol{z})}, \tag{4}$$

which consists of *1)* the log-likelihood of $\boldsymbol{z} \sim q(\boldsymbol{z}|\boldsymbol{x})$ under the remaining layers $p(\boldsymbol{z})$, *2)* the likelihood contribution term $\mathcal{V}(\boldsymbol{x}, \boldsymbol{z})$ and *3)* the looseness of the bound $\mathcal{E}(\boldsymbol{x}, \boldsymbol{z})$.

Normalizing flows, on the other hand, make use of deterministic transformations. Specifically, using a diffeomorphism $f : \mathcal{Z} \to \mathcal{X}$, the log-likelihood can be computed as

$$\log p(\boldsymbol{x}) = \log p(\boldsymbol{z}) + \underbrace{\log\left|\det \frac{\partial \boldsymbol{z}}{\partial \boldsymbol{x}}\right|}_{\text{Lik. contrib. } \mathcal{V}(\boldsymbol{x},\boldsymbol{z})}, \quad \boldsymbol{z} = f^{-1}(\boldsymbol{x}), \tag{5}$$

which consists *1)* the log-likelihood of $\boldsymbol{z} = f^{-1}(\boldsymbol{x})$ under $p(\boldsymbol{z})$ (possibly another flow) and *2)* the likelihood contribution term, which here corresponds to the log Jacobian determinant. Notice that for normalizing flows, the likelihood is exact and hence the *bound looseness* term $\mathcal{E}(\boldsymbol{x}, \boldsymbol{z}) = 0$.

In the remainder of this section we show that the change-of-variables formula (Eq. 5) can be obtained from the ELBO (Eq. 4).

**Proof.** We can use a composition of a function $g$ with a Dirac $\delta$-function:

$$\int \delta(g(\boldsymbol{z}))f(g(\boldsymbol{z}))\left|\det \frac{\partial g(\boldsymbol{z})}{\partial \boldsymbol{z}}\right|d\boldsymbol{z} = \int \delta(\boldsymbol{u})f(\boldsymbol{u})d\boldsymbol{u} \tag{6}$$

to conclude that

$$\delta(g(\boldsymbol{z})) = \left|\det \frac{\partial g(\boldsymbol{z})}{\partial \boldsymbol{z}}\right|^{-1}_{\boldsymbol{z}=\boldsymbol{z}_0} \delta(\boldsymbol{z} - \boldsymbol{z}_0) \tag{7}$$

with $\boldsymbol{z}_0$ being the root of $g(\boldsymbol{z})$. This results assumes that $g$ is smooth (derivative exists), $f$ has compact support, the root is unique and the Jacobian is non-singular.

Let $f : \mathcal{Z} \to \mathcal{X}$ be a diffeomorphism and define a pair of deterministic conditionals

$$p(\boldsymbol{x}|\boldsymbol{z}) = \delta(\boldsymbol{x} - f(\boldsymbol{z})) \tag{8}$$
$$p(\boldsymbol{z}|\boldsymbol{x}) = \delta(\boldsymbol{z} - f^{-1}(\boldsymbol{x})). \tag{9}$$

Applying the above result to $p(\boldsymbol{x}|\boldsymbol{z})$, we set $g(\boldsymbol{z}) = \boldsymbol{x} - f(\boldsymbol{z})$ and find $\boldsymbol{z}_0 = f^{-1}(\boldsymbol{x})$ and

$$p(\boldsymbol{x}|\boldsymbol{z}) = \delta(\boldsymbol{z} - f^{-1}(\boldsymbol{x}))|\det \boldsymbol{J}| = p(\boldsymbol{z}|\boldsymbol{x})|\det \boldsymbol{J}|, \tag{10}$$

where

$$\boldsymbol{J}^{-1} = \left.\frac{\partial f(\boldsymbol{z})}{\partial \boldsymbol{z}}\right|_{\boldsymbol{z}=f^{-1}(\boldsymbol{x})}.$$

Let further $q(\boldsymbol{z}|\boldsymbol{x}) = p(\boldsymbol{z}|\boldsymbol{x}) = \delta(\boldsymbol{z} - f^{-1}(\boldsymbol{x}))$. The resulting ELBO gives rise to the change-of-variables formula,

$$\log p(\boldsymbol{x}) = \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x})}\left[\log p(\boldsymbol{z}) + \log \frac{p(\boldsymbol{x}|\boldsymbol{z})}{q(\boldsymbol{z}|\boldsymbol{x})} + \log \frac{q(\boldsymbol{z}|\boldsymbol{x})}{p(\boldsymbol{z}|\boldsymbol{x})}\right] \tag{11}$$

$$= \log p(\boldsymbol{z}) + \log|\det \boldsymbol{J}|, \quad \text{for } \boldsymbol{z} = f^{-1}(\boldsymbol{x}), \tag{12}$$

where the likelihood contribution $\mathcal{V}(\boldsymbol{x}, \boldsymbol{z}) = \log \frac{p(\boldsymbol{x}|\boldsymbol{z})}{q(\boldsymbol{z}|\boldsymbol{x})} = \log|\det \boldsymbol{J}|$, while the bound looseness term $\mathcal{E}(\boldsymbol{x}, \boldsymbol{z}) = \log \frac{q(\boldsymbol{z}|\boldsymbol{x})}{p(\boldsymbol{z}|\boldsymbol{x})} = 0$, trivially.

## B  The Bound Looseness for Inference Surjections

For inference surjections $f : \mathcal{X} \to \mathcal{Z}$, the bound looseness term $\mathcal{E}(\boldsymbol{x}, \boldsymbol{z}) = 0$, given that the *stochastic right inverse condition* is satified. The stochastic right inverse condition requires that $p(\boldsymbol{x}|\boldsymbol{z})$ defines a distribution over the possible right inverses of the surjection $f$.

A right inverse function $g : \mathcal{Z} \to \mathcal{X}$ to a function $f : \mathcal{X} \to \mathcal{Z}$ satisfies $f \circ g = \mathrm{id}_{\mathcal{Z}}$, but not necessarily $g \circ f = \mathrm{id}_{\mathcal{X}}$. Here $\mathrm{id}_{\mathcal{S}}$ denotes an identity map defined on the space $\mathcal{S}$.

We satisfy the stochastic right inverse condition by requiring that $p(\boldsymbol{x}|\boldsymbol{z})$ only has support over the *fiber* of $\boldsymbol{z}$, i.e. the set of elements $\mathcal{B}(\boldsymbol{z})$ in the domain $\mathcal{X}$ that are mapped to $\boldsymbol{z}$, $\mathcal{B}(\boldsymbol{z}) := \{\boldsymbol{x}|\boldsymbol{z} = f(\boldsymbol{x})\}$. A simple check for stochastic right invertibility is thus: For any $\boldsymbol{z}$, computing $\boldsymbol{z} = f(\boldsymbol{x})$, for $\boldsymbol{x} \sim p(\boldsymbol{x}|\boldsymbol{z})$ should return the original $\boldsymbol{z}$.

Given that the distribution $p(\boldsymbol{z})$ has full support over $\mathcal{Z}$ and the stochastic right inverse condition is satisfied, we have that, for any observed $\boldsymbol{x}$, only one $\boldsymbol{z}$ could have given rise to the observation $\boldsymbol{x}$. Consequently, the posterior distribution $p(\boldsymbol{z}|\boldsymbol{x}) = \delta(\boldsymbol{z} - f(\boldsymbol{x}))$ is deterministic. By defining $q(\boldsymbol{z}|\boldsymbol{x}) = p(\boldsymbol{z}|\boldsymbol{x})$, the bound looseness is thus $\mathcal{E}(\boldsymbol{x}, \boldsymbol{z}) = 0$.

## C  List of SurVAE Layers

See Table 6 and Table 7 for lists of generative and inference surjection layers, respectively.

Table 6: Summary of some generative surjection layers.

| Surjection | Forward | Inverse | $\mathcal{V}(\boldsymbol{x}, \boldsymbol{z})$ |
|---|---|---|---|
| Rounding | $x = \lfloor z \rfloor$ | $z \sim q(z|x)$ where $z \in [x, x+1)$ | $-\log q(z|x)$ |
| Slicing | $\boldsymbol{x} = \boldsymbol{z}_1$ | $\boldsymbol{z}_1 = \boldsymbol{x}, \boldsymbol{z}_2 \sim q(\boldsymbol{z}_2|\boldsymbol{x})$ | $-\log q(\boldsymbol{z}_2|\boldsymbol{x})$ |
| Abs | $s = \mathrm{sign}\, z$ <br> $x = |z|$ | $s \sim \mathrm{Bern}(\pi(x))$ <br> $z = s \cdot x,\ \ s \in \{1, -1\}$ | $-\log q(s|x)$ |
| Max | $k = \arg\max \boldsymbol{z}$ <br> $x = \max \boldsymbol{z}$ | $k \sim \mathrm{Cat}(\boldsymbol{\pi}(x))$ <br> $z_k = x, \boldsymbol{z}_{-k} \sim q(\boldsymbol{z}_{-k}|x, k)$ | $-\log q(k|x) - \log q(\boldsymbol{z}_{-k}|x, k)$ |
| Sort | $\mathcal{I} = \mathrm{argsort}\, \boldsymbol{z}$ <br> $\boldsymbol{x} = \mathrm{sort}\, \boldsymbol{z}$ | $\mathcal{I} \sim \mathrm{Cat}(\boldsymbol{\pi}(\boldsymbol{x}))$ <br> $\boldsymbol{z} = \boldsymbol{x}_{\mathcal{I}}$ | $-\log q(\mathcal{I}|\boldsymbol{x})$ |
| ReLU | $x = \max(z, 0)$ | if $x = 0 : z \sim q(z)$, else $: z = x$ | $\mathbb{I}(x = 0)[-\log q(z)]$ |

Table 7: Summary of some inference surjection layers.

| Surjection | Forward | Inverse | $\mathcal{V}(\boldsymbol{x}, \boldsymbol{z})$ |
|---|---|---|---|
| Rounding | $x \sim p(x|z)$ where $x \in [z, z+1)$ | $z = \lfloor x \rfloor$ | $\log p(z|x)$ |
| Slicing | $\boldsymbol{x}_1 = \boldsymbol{z}, \boldsymbol{x}_2 \sim p(\boldsymbol{x}_2|\boldsymbol{z})$ | $\boldsymbol{z} = \boldsymbol{x}_1$ | $\log p(\boldsymbol{x}_2|\boldsymbol{z})$ |
| Abs | $s \sim \mathrm{Bern}(\pi(z))$ <br> $x = s \cdot z,\ \ s \in \{-1, 1\}$ | $s = \mathrm{sign}\, x$ <br> $z = |x|$ | $\log p(s|z)$ |
| Max | $k \sim \mathrm{Cat}(\boldsymbol{\pi}(z))$ <br> $x_k = z, \boldsymbol{x}_{-k} \sim p(\boldsymbol{x}_{-k}|z, k)$ | $k = \arg\max \boldsymbol{x}$ <br> $z = \max \boldsymbol{x}$ | $\log p(k|z) + \log p(\boldsymbol{x}_{-k}|z, k)$ |
| Sort | $\mathcal{I} \sim \mathrm{Cat}(\boldsymbol{\pi}(\boldsymbol{z}))$ <br> $\boldsymbol{x} = \boldsymbol{z}_{\mathcal{I}}$ | $\mathcal{I} = \mathrm{argsort}\, \boldsymbol{x}$ <br> $\boldsymbol{z} = \mathrm{sort}\, \boldsymbol{x}$ | $\log p(\mathcal{I}|\boldsymbol{z})$ |
| ReLU | if $z = 0 : x \sim p(x)$, else $: x = z$ | $z = \max(x, 0)$ | $\mathbb{I}(z = 0) \log p(x)$ |

# D The Absolute Value Surjection

We here develop the absolute value surjections, both in the generative direction $x = |z|$ and in the inference direction $z = |x|$. We will make use of Dirac delta functions to develop the likelihood contributions, but we could equivalently develop them using Gaussian distributions where $\sigma \to 0$.

## D.1 Generative Direction

**Forward and Inverse.** We define the forward and inverse transformations as

$$p(x|z) = \sum_{s \in \{-1,1\}} p(x|z,s)p(s|z) = \sum_{s \in \{-1,1\}} \delta(x - sz)\delta_{s,\text{sign}(z)}, \tag{13}$$

$$q(z|x) = \sum_{s \in \{-1,1\}} q(z|x,s)q(s|x) = \sum_{s \in \{-1,1\}} \delta(z - sx)q(s|x), \tag{14}$$

where the forward transformation $p(x|z)$ is fully deterministic and corresponds to $x = |z|$. The inference direction involves two steps, 1) sample the sign $s$ of $z$ conditioned of $x$, and 2) deterministically map $x$ to $z = sx$. Note that $q(s|x)$ may either be trained as a classifier or fixed to e.g. $q(s|x) = 1/2$. The last choice especially makes sense when $p(z)$ is symmetric.

**Likelihood Contribution.** We may develop the likelihood contribution by computing

$$\mathcal{V} = \mathbb{E}_{q(z|x,s)q(s|x)} \left[ \log \frac{p(x|z,s)p(s|z)}{q(z|x,s)q(s|x)} \right] \tag{15}$$

$$= \mathbb{E}_{\delta(z-sx)q(s|x)} \left[ \log \frac{\delta(x - sz)\delta_{s,\text{sign}(z)}}{\delta(z - sx)q(s|x)} \right] \tag{16}$$

$$\approx -\log q(s|x), \quad \text{where } z = sx, \ s \sim q(s|x). \tag{17}$$

Here, $\delta(x - sz)$ and $\delta(z - sx)$ cancel since $\delta(x - sz) = \delta(z - x/s)|1/s| = \delta(z - sx)$.

## D.2 Inference Direction

**Forward and Inverse.** We define the forward and inverse transformations as

$$p(x|z) = \sum_{s \in \{-1,1\}} p(x|z,s)p(s|z) = \sum_{s \in \{-1,1\}} \delta(x - sz)p(s|z), \tag{18}$$

$$q(z|x) = \sum_{s \in \{-1,1\}} q(z|x,s)q(s|x) = \sum_{s \in \{-1,1\}} \delta(z - sx)\delta_{s,\text{sign}(x)}, \tag{19}$$

where the inverse transformation $q(z|x)$ is fully deterministic and corresponds to $z = |x|$. The generative direction involves two steps, 1) sample the sign of $x$ conditioned of $z$, and 2) deterministically map $z$ to $x = sz$. Note that $p(s|z)$ may either be trained as a classifier or fixed to e.g. $p(s|z) = 1/2$. The last choice gives rise to an absolute value surjection which may be used to enforce exact symmetry across the origin.

**Likelihood Contribution.** We may develop the likelihood contribution by computing

$$\mathcal{V} = \mathbb{E}_{q(z|x,s)q(s|x)} \left[ \log \frac{p(x|z,s)p(s|z)}{q(z|x,s)q(s|x)} \right] \tag{20}$$

$$= \mathbb{E}_{\delta(z-sx)\delta_{s,\text{sign}(x)}} \left[ \log \frac{\delta(x - sz)p(s|z)}{\delta(z - sx)\delta_{s,\text{sign}(x)}} \right] \tag{21}$$

$$= \log p(s|z), \quad \text{where } z = sx = |x|, \ s = \text{sign}(x). \tag{22}$$

Here, $\delta(x - sz)$ and $\delta(z - sx)$ cancel since $\delta(x - sz) = \delta(z - x/s)|1/s| = \delta(z - sx)$.

# E  The Maximum Value Surjection

We here develop the maximum value surjections, both in the generative direction $x = \max z$ and in the inference direction $z = \max x$. We will make use of Dirac delta functions to develop the likelihood contributions, but we could equivalently develop them using Gaussian distributions where $\sigma \to 0$.

## E.1  Generative Direction

**Forward and Inverse.** We define the forward and inverse transformations as

$$p(x|z) = \sum_{k=1}^{K} p(x|z, k)p(k|z) = \sum_{k=1}^{K} \delta(x - z_k)\delta_{k,\arg\max(z)}, \tag{23}$$

$$q(z|x) = \sum_{k=1}^{K} q(z|x, k)q(k|x) = \sum_{k=1}^{K} \delta(z_k - x)q(z_{-k}|x, k)q(k|x), \tag{24}$$

where $k$ refers to the indices of $z$, $K$ is the number of elements in $z$ and $z_{-k}$ is $z$ excluding element $k$. The forward transformation $p(x|z)$ is fully deterministic and corresponds to $x = \max z$. The inference direction involves three steps, 1) sample the index $k$ for the argmax of $z$ conditioned of $x$, 2) deterministically map $x$ to $z_k = x$, and 3) infer the remaining elements $z_{-k}$ of $z$. Note that $q(k|x)$ may either be trained as a classifier or fixed to e.g. $q(k|x) = 1/K$.

For $q$ to define a right-inverse of $p$, we require that $q(z_{-k}|x, k)$ only has support in $(-\infty, x)^{K-1}$ such that $z_k$ will be the maximum value.

**Likelihood Contribution.** We may develop the likelihood contribution by computing

$$\mathcal{V} = \mathbb{E}_{q(z|x,k)q(k|x)} \left[ \log \frac{p(x|z, k)p(k|z)}{q(z|x, k)q(k|x)} \right] \tag{25}$$

$$= \mathbb{E}_{\delta(z_k-x)q(z_{-k}|x,k)q(k|x)} \left[ \log \frac{\delta(x - z_k)\delta_{k,\arg\max(z)}}{\delta(z_k - x)q(z_{-k}|x, k)q(k|x)} \right] \tag{26}$$

$$\approx -\log q(k|x) - \log q(z_{-k}|x, k), \quad \text{where } z_k = x, \ z_{-k} \sim q(z_{-k}|x, k), \ k \sim q(k|x). \tag{27}$$

## E.2  Inference Direction

**Forward and Inverse.** We define the forward and inverse transformations as

$$p(x|z) = \sum_{k=1}^{K} p(x|z, k)p(k|z) = \sum_{k=1}^{K} \delta(x_k - z)p(x_{-k}|z, k)p(k|z), \tag{28}$$

$$q(z|x) = \sum_{k=1}^{K} q(z|x, k)q(k|x) = \sum_{k=1}^{K} \delta(z - x_k)\delta_{k,\arg\max(x)}, \tag{29}$$

where $k$ refers to the indices of $x$, $K$ is the number of elements in $x$ and $x_{-k}$ is $x$ excluding element $k$. The inverse transformation $q(z|x)$ is fully deterministic and corresponds to $z = \max x$. The inference direction involves three steps, 1) sample the index $k$ for the argmax of $x$ conditioned of $z$, 2) deterministically map $z$ to $x_k = z$, and 3) infer the remaining elements $x_{-k}$ of $x$. Note that $p(k|z)$ may either be trained as a classifier or fixed to e.g. $p(k|z) = 1/K$.

For $p$ to define a right-inverse of $q$, we require that $p(x_{-k}|z, k)$ only has support in $(-\infty, z)^{K-1}$ such that $x_k$ will be the maximum value.

**Likelihood Contribution.** We may develop the likelihood contribution by computing

$$\mathcal{V} = \mathbb{E}_{q(z|x,k)q(k|x)} \left[ \log \frac{p(x|z, k)p(k|z)}{q(z|x, k)q(k|x)} \right] \tag{30}$$

$$= \mathbb{E}_{\delta(z-x_k)\delta_{k,\arg\max(x)}} \left[ \log \frac{\delta(x_k - z)p(x_{-k}|z, k)p(k|z)}{\delta(z - x_k)\delta_{k,\arg\max(x)}} \right] \tag{31}$$

$$= \log p(k|z) + \log p(x_{-k}|z, k), \quad \text{where } z = x_k = \max x, \ k = \arg\max x. \tag{32}$$

## F  The Sort Surjection

We here develop the sorting surjections, both in the generative direction $x = \text{sort} \, z$ and in the inference direction $z = \text{sort} \, x$. We will make use of Dirac delta functions to develop the likelihood contributions, but we could equivalently develop them using Gaussian distributions where $\sigma \to 0$.

### F.1  Generative Direction

**Forward and Inverse.** We define the forward and inverse transformations as

$$p(\boldsymbol{x}|\boldsymbol{z}) = \sum_{\mathcal{I}} p(\boldsymbol{x}|\boldsymbol{z}, \mathcal{I}) p(\mathcal{I}|\boldsymbol{z}) = \sum_{\mathcal{I}} \delta(\boldsymbol{x} - \boldsymbol{z}_{\mathcal{I}}) \delta_{\mathcal{I}, \text{argsort}(\boldsymbol{z})}, \tag{33}$$

$$q(\boldsymbol{z}|\boldsymbol{x}) = \sum_{\mathcal{I}} q(\boldsymbol{z}|\boldsymbol{x}, \mathcal{I}) q(\mathcal{I}|\boldsymbol{x}) = \sum_{\mathcal{I}} \delta(\boldsymbol{z} - \boldsymbol{x}_{\mathcal{I}^{-1}}) q(\mathcal{I}|\boldsymbol{x}), \tag{34}$$

where $\mathcal{I}$ refers to a set of permutation indices, $\mathcal{I}^{-1}$ refers to the inverse permutation indices and $\boldsymbol{z}_{\mathcal{I}}$ refers to the elements of $\boldsymbol{z}$ permuted according to the indices $\mathcal{I}$. Note that there are $D!$ possible permutations.

The forward transformation $p(\boldsymbol{x}|\boldsymbol{z})$ is fully deterministic and corresponds to $x = \text{sort} \, z$. The inference direction involves two steps, 1) sample permutation indices $\mathcal{I}$ conditioned of $\boldsymbol{x}$, and 2) deterministically permute $\boldsymbol{x}$ according to the inverse permutation $\mathcal{I}^{-1}$ to obtain $z = x_{\mathcal{I}^{-1}}$. Note that $q(\mathcal{I}|\boldsymbol{x})$ may either be trained as a classifier or fixed to e.g. $q(\mathcal{I}|\boldsymbol{x}) = 1/D!$.

**Likelihood Contribution.** We may develop the likelihood contribution by computing

$$\mathcal{V} = \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x}, \mathcal{I}) q(\mathcal{I}|\boldsymbol{x})} \left[ \log \frac{p(\boldsymbol{x}|\boldsymbol{z}, \mathcal{I}) p(\mathcal{I}|\boldsymbol{z})}{q(\boldsymbol{z}|\boldsymbol{x}, \mathcal{I}) q(\mathcal{I}|\boldsymbol{x})} \right] \tag{35}$$

$$= \mathbb{E}_{\delta(\boldsymbol{z} - \boldsymbol{x}_{\mathcal{I}^{-1}}) q(\mathcal{I}|\boldsymbol{x})} \left[ \log \frac{\delta(\boldsymbol{x} - \boldsymbol{z}_{\mathcal{I}}) \delta_{\mathcal{I}, \text{argsort}(\boldsymbol{z})}}{\delta(\boldsymbol{z} - \boldsymbol{x}_{\mathcal{I}^{-1}}) q(\mathcal{I}|\boldsymbol{x})} \right] \tag{36}$$

$$\approx -\log q(\mathcal{I}|\boldsymbol{x}), \quad \text{where } \mathcal{I} \sim q(\mathcal{I}|\boldsymbol{x}). \tag{37}$$

### F.2  Inference Direction

**Forward and Inverse.** We define the forward and inverse transformations as

$$p(\boldsymbol{x}|\boldsymbol{z}) = \sum_{\mathcal{I}} p(\boldsymbol{x}|\boldsymbol{z}, \mathcal{I}) p(\mathcal{I}|\boldsymbol{z}) = \sum_{\mathcal{I}} \delta(\boldsymbol{x} - \boldsymbol{z}_{\mathcal{I}^{-1}}) p(\mathcal{I}|\boldsymbol{z}), \tag{38}$$

$$q(\boldsymbol{z}|\boldsymbol{x}) = \sum_{\mathcal{I}} q(\boldsymbol{z}|\boldsymbol{x}, \mathcal{I}) q(\mathcal{I}|\boldsymbol{x}) = \sum_{\mathcal{I}} \delta(\boldsymbol{z} - \boldsymbol{x}_{\mathcal{I}}) \delta_{\mathcal{I}, \text{argsort}(\boldsymbol{x})}, \tag{39}$$

where $\mathcal{I}$ refers to a set of permutation indices, $\mathcal{I}^{-1}$ refers to the inverse permutation indices and $\boldsymbol{x}_{\mathcal{I}}$ refers to the elements of $\boldsymbol{x}$ permuted according to the indices $\mathcal{I}$. Note that there are $D!$ possible permutations.

The inverse transformation $q(\boldsymbol{z}|\boldsymbol{x})$ is fully deterministic and corresponds to $z = \text{sort} \, x$. The generative direction involves two steps, 1) sample permutation indices $\mathcal{I}$ conditioned of $\boldsymbol{z}$, and 2) deterministically permute $\boldsymbol{z}$ according to the inverse permutation $\mathcal{I}^{-1}$ to obtain $x = z_{\mathcal{I}^{-1}}$. Note that $p(\mathcal{I}|\boldsymbol{z})$ may either be trained as a classifier or fixed to e.g. $p(\mathcal{I}|\boldsymbol{z}) = 1/D!$.

**Likelihood Contribution.** We may develop the likelihood contribution by computing

$$\mathcal{V} = \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x}, \mathcal{I}) q(\mathcal{I}|\boldsymbol{x})} \left[ \log \frac{p(\boldsymbol{x}|\boldsymbol{z}, \mathcal{I}) p(\mathcal{I}|\boldsymbol{z})}{q(\boldsymbol{z}|\boldsymbol{x}, \mathcal{I}) q(\mathcal{I}|\boldsymbol{x})} \right] \tag{40}$$

$$= \mathbb{E}_{\delta(\boldsymbol{z} - \boldsymbol{x}_{\mathcal{I}}) \delta_{\mathcal{I}, \text{argsort}(\boldsymbol{x})}} \left[ \log \frac{\delta(\boldsymbol{x} - \boldsymbol{z}_{\mathcal{I}^{-1}}) p(\mathcal{I}|\boldsymbol{z})}{\delta(\boldsymbol{z} - \boldsymbol{x}_{\mathcal{I}}) \delta_{\mathcal{I}, \text{argsort}(\boldsymbol{x})}} \right] \tag{41}$$

$$= \log p(\mathcal{I}|\boldsymbol{z}), \quad \text{where } \boldsymbol{z} = \boldsymbol{x}_{\mathcal{I}} = \text{sort} \, \boldsymbol{x}, \ \mathcal{I} = \text{argsort} \, \boldsymbol{x}. \tag{42}$$

# G  The Stochastic Permutation

We here develop the stochastic permutation layer which randomly permutes its input. The inverse pass mirrors the forward pass. Note that stochastic permutation is *not* a surjection, but rather a stochastic transform. We will make use of Dirac delta functions to develop the likelihood contributions, but we could equivalently develop them using Gaussian distributions where $\sigma \to 0$.

**Forward and Inverse.** We define the forward and inverse transformations as

$$p(\boldsymbol{x}|\boldsymbol{z}) = \sum_{\mathcal{I}} p(\boldsymbol{x}|\boldsymbol{z}, \mathcal{I}) p(\mathcal{I}) = \sum_{\mathcal{I}} \delta(\boldsymbol{x} - \boldsymbol{z}_{\mathcal{I}}) \operatorname{Unif}(\mathcal{I}), \tag{43}$$

$$q(\boldsymbol{z}|\boldsymbol{x}) = \sum_{\mathcal{I}} q(\boldsymbol{z}|\boldsymbol{x}, \mathcal{I}) q(\mathcal{I}) = \sum_{\mathcal{I}} \delta(\boldsymbol{z} - \boldsymbol{x}_{\mathcal{I}^{-1}}) \operatorname{Unif}(\mathcal{I}), \tag{44}$$

where $\mathcal{I}$ refers to a set of permutation indices, $\mathcal{I}^{-1}$ refers to the inverse permutation indices and $\boldsymbol{z}_{\mathcal{I}}$ refers to the elements of $\boldsymbol{z}$ permuted according to the indices $\mathcal{I}$. Note that there are $D!$ possible permutations.

The transformation is stochastic and involves the same two steps in both directions: 1) Sample permutation indices $\mathcal{I}$ uniformly at random, and 2) deterministically permute the input according to the samples indices $\mathcal{I}$.

**Likelihood Contribution.** We may develop the likelihood contribution by computing

$$\mathcal{V} = \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x}, \mathcal{I}) q(\mathcal{I})} \left[ \log \frac{p(\boldsymbol{x}|\boldsymbol{z}, \mathcal{I}) p(\mathcal{I})}{q(\boldsymbol{z}|\boldsymbol{x}, \mathcal{I}) q(\mathcal{I})} \right] \tag{45}$$

$$= \mathbb{E}_{\delta(\boldsymbol{z} - \boldsymbol{x}_{\mathcal{I}^{-1}}) \operatorname{Unif}(\mathcal{I})} \left[ \log \frac{\delta(\boldsymbol{x} - \boldsymbol{z}_{\mathcal{I}}) \operatorname{Unif}(\mathcal{I})}{\delta(\boldsymbol{z} - \boldsymbol{x}_{\mathcal{I}^{-1}}) \operatorname{Unif}(\mathcal{I})} \right] \tag{46}$$

$$= 0. \tag{47}$$

This layer thus takes the simple form: Both during the forward and inverse passes, shuffle the input uniformly at random. The resulting likelihood contribution is zero.

# H  The Software Perspective

Normalizing flows provide a powerful modular framework where flexible densities may be specified using a composition of bijective transformations. Each bijection may be implemented as a module contained 3 important components: 1) A forward transformation $x = f(z)$, 2) an inverse transformation $z = f^{-1}(x)$, and 3) a Jacobian determinant $\log|\det J|$. Several software libraries for normalizing flows have been built using this modular design principle (Dillon et al., 2017; Bingham et al., 2018).

SurVAE flows suggest that such software frameworks may be directly extended since the modules follow the exact same design principles – each module has 3 important components:

1. A forward transformation $\mathcal{Z} \to \mathcal{X}$.
2. An inverse transformation $\mathcal{X} \to \mathcal{Z}$.
3. A likelihood contribution $\mathcal{V}(x, z)$.

SurVAE flows allow compositions of not only bijective transformations, but also surjective and stochastic transformations. This allows us to obtain methods such as dequantization (Uria et al., 2014; Theis et al., 2016; Ho et al., 2019), variational data augmentation (Huang et al., 2020; Chen et al., 2020), multi-scale architectures (Dinh et al., 2017) as composable surjective transformations and VAEs (Kingma and Welling, 2014; Rezende et al., 2014) as composable stochastic transformations.

In our code[3], we provide a library of SurVAE flows that may serve as a prototype for a more extensive library. In the next subsections, we show some selected code snippets from our library. The code is based on PyTorch (Paszke et al., 2019), but can easily be ported to other frameworks. Note that in the implementation, the `forward` method implements the inverse transformation $\mathcal{X} \to \mathcal{Z}$ and the likelihood contribution $\mathcal{V}(x, z)$, since this is what is needed during the forward pass of backpropagation used for training.

In Sec. H.1 we show an implementation of a VAE as a stochastic transformation, while in Sec. H.2 and Sec. H.3 we show implementations of dequantization and variational data augmentation as surjective transformations. Finally, in Sec. H.4, we show an example of how to construct an augmented normalizing flow through composition of SurVAE layers.

## H.1  VAE

We implement VAEs as a composable stochastic transformation.

```python
class VAE(StochasticTransform):
    '''A variational autoencoder layer.'''

    def __init__(self, decoder, encoder):
        super(VAE, self).__init__()
        self.decoder = decoder
        self.encoder = encoder

    def forward(self, x):
        z, log_qz = self.encoder.sample_with_log_prob(context=x)
        log_px = self.decoder.log_prob(x, context=z)
        ldj = log_px - log_qz
        return z, ldj

    def inverse(self, z):
        x = self.decoder.sample(context=z)
        return x
```

---

[3]The code is available at `https://github.com/didriknielsen/survae_flows`

## H.2 Dequantization

We implement `UniformDequantization`, which may be used to convert between discrete and continuous variables, as a generative rounding surjection.

```python
class UniformDequantization(Surjection):
    '''A uniform dequantization layer.'''

    def forward(self, x):
        z = x.float() + torch.rand_like(x)
        ldj = torch.zeros(x.shape[0])
        return z, ldj

    def inverse(self, z):
        x = z.floor().long()
        return x
```

## H.3 Augmentation

We implement `Augment`, a generative tensor slicing surjection, which may be used to construct e.g. augmented normalizing flows (Huang et al., 2020; Chen et al., 2020).

```python
class Augment(Surjection):
    '''An augmentation layer.'''

    def __init__(self, encoder, split_size):
        super(Augment, self).__init__()
        self.encoder = encoder
        self.split_size = split_size

    def forward(self, x):
        z2, log_qz2 = self.encoder.sample_with_log_prob(context=x)
        z = torch.cat([x, z2], dim=1)
        ldj = -log_qz2
        return z, ldj

    def inverse(self, z):
        x, z2 = torch.split(z, self.split_size, dim=1)
        return x
```

## H.4 Example: Augmented Normalizing Flows

We showcase here the simplicity of implementing an augmented normalizing flow using the SurVAE flow framework. In Listing 1, a simple normalizing flow consisting of 2 coupling layers is constructed. In Listing 2, this is extended by adding an `Augment` surjection, resulting in an augmented flow.

| Listing 1: A basic flow. | Listing 2: An augmented flow. |
|---|---|
| ```python Flow(base_dist=Normal((2,)),     transforms=[      CouplingBijection(),     Reverse(),     CouplingBijection(), ]) ``` | ```python Flow(base_dist=Normal((4,)),     transforms=[         Augment(Normal((2,)), (2,2)),         CouplingBijection(),         Reverse(),         CouplingBijection(), ]) ``` |

Using the models in Listing 1 and Listing 2, we compare a standard coupling flow with a simple extension using an additional `Augment` layer. We use 4 coupling layers instead of 2 and train models both using identical setups 10000 iterations each. Augmented flows have improved capabilities of modelling data with disconnected components. In Fig. 8, we observe that the augmented flows tend to place their mass more out in a more "clean" fashion and thus demonstrate improved ability to model complicated 2D densities.
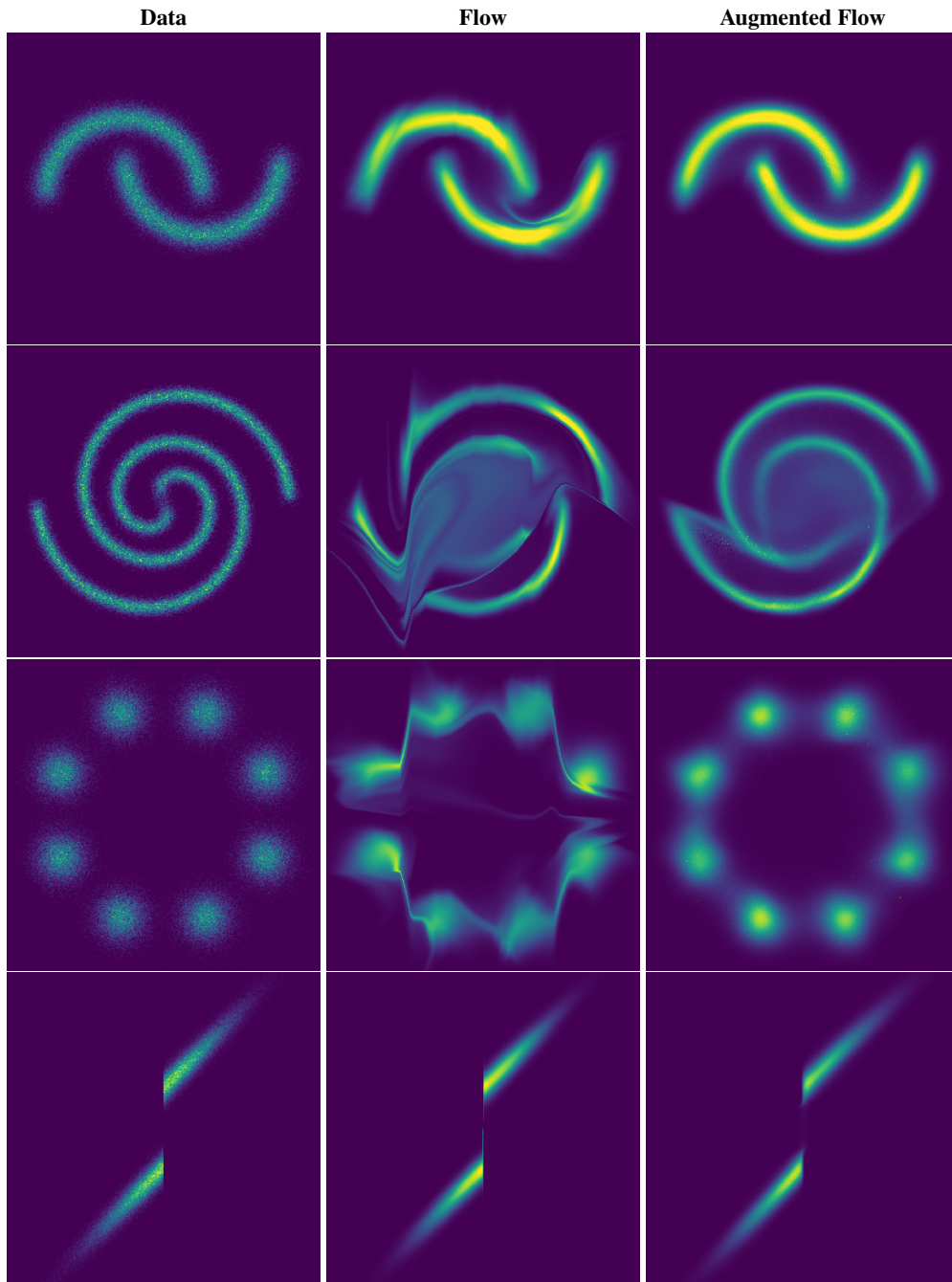
Figure 8: Augmented flows show improved capabilities over flows at modelling 2D densities, especially where there are disconnected components. With SurVAE flows, augmented flows are implemented by adding a single surjective augmentation layer to the flow as shown in Listing 2.

# I   Experimental Details

We here give more details on the experiments. For further details, see our open-source code[4].

## I.1   Synthetic Data

**Data.** We used 4 synthetic datasets, `checkerboard`, `corners`, `gaussians` and `circles`. For each syntheric dataset, 128000 samples were used as a training set and 128000 more samples as a test set. The `checkerboard` dataset is anti-symmetric, while the 3 others are symmetric.

**Training.** We used the Adam optimizer (Kingma and Ba, 2015) with a learning rate of $10^{-3}$. All models were trained for 10000 iterations (10 epochs) using a batch size of 128.

**Baseline.** The baseline flow is a composition of 4 affine coupling bijections with the ordering reversed in-between. The coupling layers are parameterized by MLPs with hidden units (200,100) and ReLU activations. The base distribution is a standard Gaussian.

**Symmetric AbsFlow.** For the symmetric datasets, AbsFlow uses all the same layers as the baseline. In addition, an `abs` surjection is added, followed by and inverse softplus (`gaussians`) or logit (`checkerboard` and `corners`). In the generative direction, the `abs` surjection randomly samples the sign with equal probabilities. The extra layers contain no parameters, and the AbsFlow thus have the exact same number of parameters as the baseline.

**Anti-Symmetric AbsFlow.** For the anti-symmetric dataset, AbsFlow uses only a single `abs` surjection and a uniform base distribution. In the generative direction, a classifier network learns the probabilities of sampling the sign conditioned on $z$. This classifier network is, like the coupling layer networks, an MLP with (200,100) hidden units and ReLU activations. In this case, the AbsFlow thus has ~1/4 the number of parameters.

## I.2   Point Cloud Data

**Data.** We used the `SpatialMNIST` dataset (Edwards and Storkey, 2017). This dataset was constructed by, for each digit in the `MNIST` dataset, sampling 50 points according to the normalized pixel intensities. We used the official code[5] to construct the dataset. We split the dataset into parts of 50000-10000-10000 for training, validation and test (without shuffling). Each data example is a set of 50 2D points which we represent as a tensor of shape `(2,50)`.

**Training.** Both models were trained for 500 epochs using a batch size of 128. We used the Adam optimizer (Kingma and Ba, 2015) with an initial learning rate of $10^{-3}$. The learning rate was warmed up linearly for 2000 iterations and the decayed by 0.995 every epoch. All models were trained using a single GPU for about 40 hours.

**Evaluation.** SortFlow allows exact computation of the likelihood, while the PermuteFlow only allows computation of lower bounds. We evaluated PermuteFlow using the IWBO (importance weighted bound) (Burda et al., 2016) using $k = 1000$ importance samples. PermuteFlow obtains an ELBO of -5.32 PPLL and an IWBO of -5.30 PPLL, while SortFlow obtains an exact log-likelihood of -5.53 PPLL.

**Hyperparameters.** We tuned the dropout rate using the validation set. We considered dropout rates of $\{0.0, 0.05, 0.1, 0.2, 0.3\}$ for both models. We found 0.1 to work best for PermuteFlow, while 0.2 worked best for SortFlow.

**PermuteFlow.** We used a flow of an initial stochastic permutation layer followed by 32 steps with ActNorm layers (Kingma and Dhariwal, 2018) in-between. Each step consisted of 1) an affine coupling bijection which transforms a the first half tensor `(1,50)` conditioned on the other half `(1,50)`, 2) reversing the order along the spatial dimension, 3) an affine coupling bijection which transforms the first half tensor `(2,25)` conditioned on the other half `(2,25)`, 4) a stochastic permutation along the point dimension. Each of the coupling bijections are parametersized by Transformer networks (Vaswani et al., 2017) *without* positional encoding. The Transformers used 2 blocks, with $d_{\text{model}} = 64$, $d_{\text{ff}} = 256$ and 8 attention heads.

---

[4] https://github.com/didriknielsen/survae_flows
[5] https://github.com/conormdurkan/neural-statistician

**SortFlow.** This follows the setup of PermuteFlow, with the following changes: 1) The initial stochastic permutation is replaced by a sorting layer. 2) The stochastic permutations in the flow are swapped with fixed permutations (sampled at random once, before training). 3) The Transformers make use of a learned positional encoding, since the sorting layer enforces a canonical ordering of the points.

### I.3 Image Data

**Data.** We used the CIFAR-10, ImageNet $32 \times 32$ and ImageNet $64 \times 64$ datasets. The CIFAR-10 dataset comes pre-split in 50000 training examples and 10000 test examples. The ImageNet datasets also come pre-split in 1,281,149 training examples and 49,999 validation examples. We use these splits and report results for the test set of CIFAR-10 and the validations sets of ImageNet $32 \times 32$ and ImageNet $64 \times 64$.

**Training.** We used the Adamax optimizer (Kingma and Ba, 2015) with an initial learning rate of $10^{-3}$ and a batch size of 32. The learning rate was linearly warmed up for 5000 iterations. For CIFAR-10, the models were first trained for 500 epochs with the learning rate decayed by 0.995 every epoch. Next, the models were "cooled down" for an additional 50 epochs with a smaller learning rate of $2 \cdot 10^{-5}$. For the ImageNet datasets, the models were first trained for 25 epochs (ImageNet $32 \times 32$) and 20 epochs (ImageNet $64 \times 64$) with the learning rate decayed by 0.95 every epoch. Next, the models were "cooled down" for an additional 2 epochs with a smaller learning rate of $5 \cdot 10^{-5}$. The CIFAR-10 and ImageNet $32 \times 32$ models were trained on a single GPU for about 2 weeks, while the ImageNet $64 \times 64$ models were trained using 4 GPUs for about 3 weeks. We provide pre-trained model checkpoints in our open-source code. Note that data augmentation was applied during training of the CIFAR-10 models, including random flipping and rotations. See code for more details.

**Evaluation.** The CIFAR-10 models were evaluated using the IWBO (importance weighted bound) (Burda et al., 2016) using $k = 1000$ importance samples. The ImageNet models were evaluated using the ELBO (which corresponds to the IWBO with $k = 1$ importance sample).

**Baseline.** For CIFAR-10 and ImageNet $32 \times 32$, the flow uses 2 scales with 12 steps/scale. For ImageNet $64 \times 64$, the flow uses 3 scales with 8 steps/scale. Each step consists of an affine coupling bijection (Dinh et al., 2017) and an invertible $1 \times 1$ convolution (Kingma and Dhariwal, 2018). All models are trained with variational dequantization (Ho et al., 2019) and an initial squeezing layer (Dinh et al., 2017) to increase the number of channels from 3 to 12. The coupling bijections are parameterized by DenseNets (Huang et al., 2017).

**MaxPoolFlow.** The MaxPoolFlow uses the exact same setup as the baseline, but replaces the tensor slicing surjection with a max pooling surjection. In the generative direction, we used the simplest possible choice: Each input pixel is equally likely to be copied to any of the pixels in its corresponding $2 \times 2$ patch. The remaining 3 elements are sampled such that the copied value remains the largest: They are set equal to this maximum value minus noise from a standard half-normal distribution (i.e. Gaussian distribution with only positive values). We used simple choices containing *no extra parameters* in order to facilitate more fair comparison. Note that the max pooling layer could be potentially be improved by using more sophisticated choices for the distribution for sampling the remaining elements, $p(\boldsymbol{x}_{-k}|\boldsymbol{z})$, and/or by using a classifier, $p(\boldsymbol{k}|\boldsymbol{z})$, to predict the indices $\boldsymbol{k}$.

# J Additional Samples

Samples from SurVAE flows trained on `CIFAR-10`, `ImageNet` $32 \times 32$ and `ImageNet` $64 \times 64$ using either max pooling or tensor slicing for downsampling are shown in Fig. 9, Fig. 10 and Fig. 11, respectively.
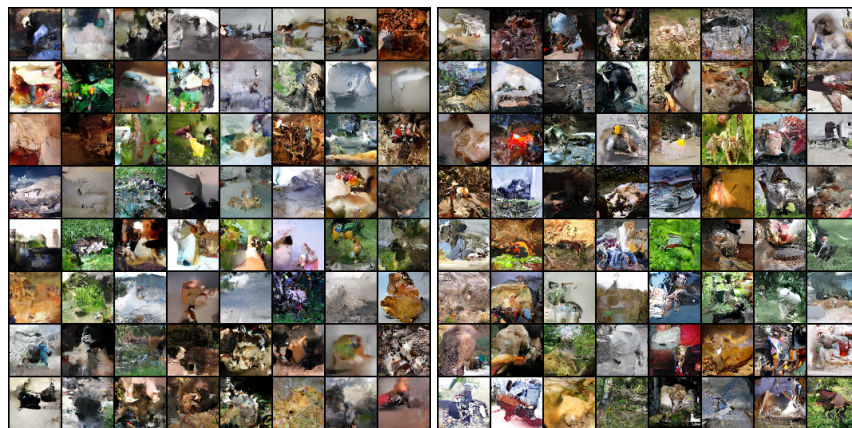


(a) Max Pooling.  (b) No Pooling.

Figure 9: Unconditional samples from SurVAE flows trained on `CIFAR-10`.



(a) Max Pooling.  (b) No Pooling.

Figure 10: Unconditional samples from SurVAE flows trained on `ImageNet` $32 \times 32$.

(a) Max Pooling.

(b) No Pooling.

Figure 11: Unconditional samples from SurVAE flows trained on ImageNet $64 \times 64$.