

Normalizing Flows with Multi-Scale Autoregressive Priors

Shweta Mahajan*¹ Apratim Bhattacharyya*² Mario Fritz³ Bernt Schiele² Stefan Roth¹

¹Department of Computer Science, TU Darmstadt

²Max Planck Institute for Informatics, Saarland Informatics Campus

³CISPA Helmholtz Center for Information Security, Saarland Informatics Campus

Abstract

Flow-based generative models are an important class of exact inference models that admit efficient inference and sampling for image synthesis. Owing to the efficiency constraints on the design of the flow layers, e.g. split coupling flow layers in which approximately half the pixels do not undergo further transformations, they have limited expressiveness for modeling long-range data dependencies compared to autoregressive models that rely on conditional pixel-wise generation. In this work, we improve the representational power of flow-based models by introducing channel-wise dependencies in their latent space through multi-scale autoregressive priors (mAR). Our mAR prior for models with split coupling flow layers (mAR-SCF) can better capture dependencies in complex multimodal data. The resulting model achieves state-of-the-art density estimation results on MNIST, CIFAR-10, and ImageNet. Furthermore, we show that mAR-SCF allows for improved image generation quality, with gains in FID and Inception scores compared to state-of-the-art flow-based models.

1. Introduction

Deep generative models aim to learn complex dependencies within very high-dimensional input data, e.g. natural images [3, 28] or audio data [7], and enable generating new samples that are representative of the true data distribution. These generative models find application in various downstream tasks like image synthesis [11, 20, 38] or speech synthesis [7, 39]. Since it is not feasible to learn the exact distribution, generative models generally approximate the underlying true distribution. Popular generative models for capturing complex data distributions are Generative Adversarial Networks (GANs) [11], which model the distribution implicitly and generate (high-dimensional) samples by transforming a noise distribution into the desired space with complex dependencies; however, they may not cover

* Authors contributed equally.

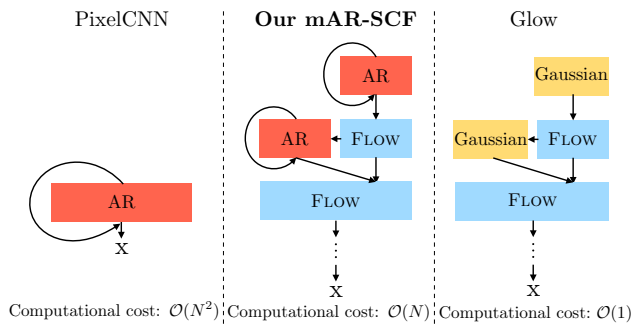


Figure 1. Our mAR-SCF model combines normalizing flows with autoregressive (AR) priors to improve modeling power while ensuring that the computational cost grows linearly with the spatial image resolution $N \times N$.

all modes of the underlying data distribution. Variational Autoencoders (VAEs) [20] optimize a lower bound on the log-likelihood of the data. This implies that VAEs can only approximately optimize the log-likelihood [30].

Autoregressive models [10, 37, 38] and normalizing flow-based generative models [8, 9, 18] are exact inference models that optimize the exact log-likelihood of the data. Autoregressive models can capture complex and long-range dependencies between the dimensions of a distribution, e.g. in case of images, as the value of a pixel is conditioned on a large context of neighboring pixels. The main limitation of this approach is that image synthesis is sequential and thus difficult to parallelize. Recently proposed normalizing flow-based models, such as NICE [8], RealNVP [9], and Glow [18], allow exact inference by mapping the input data to a known base distribution, e.g. a Gaussian, through a series of invertible transformations. These models leverage invertible split coupling flow (SCF) layers in which certain dimensions are left unchanged by the invertible transformation as well as SPLIT operations following which certain dimensions do not undergo subsequent transformations. This allows for considerably easier parallelization of both inference and generation processes. However, these models lag behind autoregressive models for density estimation.

In this work, we (i) propose multi-scale autoregressive priors for invertible flow models with split coupling flow layers, termed *mAR-SCF*, to address the limited modeling power of non-autoregressive invertible flow models [9, 14, 18, 28] (Fig. 1); (ii) we apply our multi-scale autoregressive prior after every SPLIT operation such that the computational cost of sampling grows linearly in the spatial dimensions of the image compared to the quadratic cost of traditional autoregressive models (given sufficient parallel resources); (iii) our experiments show that we achieve state-of-the-art density estimation results on MNIST [22], CIFAR10 [21], and ImageNet [31] compared to prior invertible flow-based approaches; and finally (iv) we show that our multi-scale autoregressive prior leads to better sample quality as measured by the FID metric [13] and the Inception score [32], significantly lowering the gap to GAN approaches [27, 40].

2. Related Work

In this work, we combine the expressiveness of autoregressive models with the efficiency of non-autoregressive flow-based models for exact inference.

Autoregressive models [6, 10, 12, 15, 26, 37, 38] are a class of exact inference models that factorize the joint probability distribution over the input space as a product of conditional distributions, where each dimension is conditioned on the previous ones in a pre-defined order. Recent autoregressive models, such as PixelCNN and PixelRNN [37, 38], can generate high-quality images but are difficult to parallelize since synthesis is sequential. It is worth noting that autoregressive image models, such as that of Domke *et al.* [10], significantly pre-date their recent popularity. Various extensions have been proposed to improve the performance of the PixelCNN model. For example, Multiscale-PixelCNN [29] extends PixelCNN to improve the sampling runtime from linear to logarithmic in the number of pixels, exploiting conditional independence between the pixels. Chen *et al.* [6] introduce self-attention in PixelCNN models to improve the modeling power. Salimans *et al.* [33] introduce skip connections and a discrete logistic likelihood model. WaveRNN [17] leverages customized GPU kernels to improve the sampling speed for audio synthesis. Menick *et al.* [23] synthesize images by sequential conditioning on sub-images within an image. These methods, however, still suffer from slow sampling speed and are difficult to parallelize.

Flow-based generative models, first introduced in [8], also allow for exact inference. These models are composed of a series of invertible transformations, each with a tractable Jacobian and inverse, which maps the input distribution to a known base density, *e.g.* a Gaussian. Papamakarios *et al.* [25] proposed autoregressive invertible transformations using masked decoders. However, these are difficult to parallelize just like PixelCNN-based approaches. Kingma *et al.*

[19] propose inverse autoregressive flow (IAF), where the means and variances of pixels depend on random variables and not on previous pixels, making it easier to parallelize. However, the approach offers limited generalization [39].

Recent extensions. Recent work [1, 9, 18] extends normalizing flows [8] to multi-scale architectures with split couplings, which allow for efficient inference and sampling. For example, Kingma *et al.* [18] introduce additional invertible 1×1 convolutions to capture non-linearities in the data distribution. Hooeboom *et al.* [16] extend this to $d \times d$ convolutions, increasing the receptive field. [4] improve the residual blocks of flow layers with memory efficient gradients based on the choice of activation functions. A key advantage of flow-based generative models is that they can be parallelized for inference and synthesis. Ho *et al.* [14] propose Flow++ with various modifications in the architecture of the flows in [9], including attention and a variational quantization method to improve the data likelihood. The resulting model is computationally expensive as non-linearities are applied along all the dimensions of the data at every step of the flow, *i.e.* all the dimensions are instantiated with the prior distribution at the last layer of the flow. While comparatively efficient, such flow-based models have limited expressiveness compared to autoregressive models, which is reflected in their lower data log-likelihood. It is thus desirable to develop models that have the expressiveness of autoregressive models and the efficiency of flow-based models. This is our goal here.

Methods with complex priors. Recent work [5] develops complex priors to improve the data likelihoods. VQ-VAE2 integrates autoregressive models as priors [28] with discrete latent variables [5] for high-quality image synthesis and proposes latent graph-based models in a VAE framework. Tomczak *et al.* [36] propose mixtures of Gaussians with pre-defined clusters, and [5] use neural autoregressive model priors in the latent space, which improves results for image synthesis. Ziegler *et al.* [41] learn a prior based on normalizing flows to capture multimodal discrete distributions of character-level texts in the latent spaces with nonlinear flow layers. However, this invertible layer is difficult to be optimized in both directions. Moreover, these models do not allow for exact inference. In this work, we propose complex autoregressive priors to improve the power of invertible split coupling-based normalizing flows [9, 14, 18].

3. Overview and Background

In this work, we propose multi-scale autoregressive priors for split coupling-based flow models, termed *mAR-SCF*, where we leverage autoregressive models to improve the modeling flexibility of invertible normalizing flow models without sacrificing sampling efficiency. As we build upon normalizing flows and autoregressive models, we first provide an overview of both.

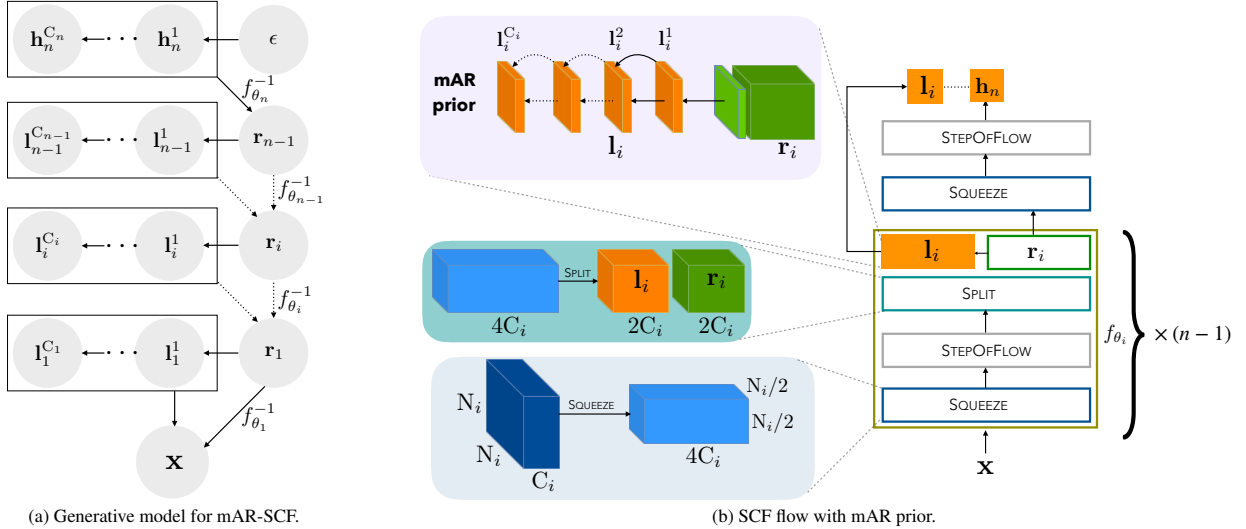


Figure 2. Flow-based generative models with multi-scale autoregressive priors (*mAR-SCF*). The generative model (*left*) shows the multi-scale autoregressive sampling of the channel dimensions of \mathbf{l}_i at each level. The spatial dimensions of each channel are sampled in parallel. \mathbf{r}_i are computed with invertible transformations. The *mAR-SCF* model (*right*) shows the complete multi-scale architecture with the mAR prior applied along the channels of \mathbf{l}_i , *i.e.* at each level i after the SPLIT operation.

Normalizing flows. Normalizing flows [8] are a class of exact inference generative models. Here, we consider invertible flows, which allow for both efficient exact inference and sampling. Specifically, invertible flows consist of a sequence of n invertible functions f_{θ_i} , which transform a density on the data \mathbf{x} to a density on latent variables \mathbf{z} ,

$$\mathbf{x} \xleftrightarrow{f_{\theta_1}} \mathbf{h}_1 \xleftrightarrow{f_{\theta_2}} \mathbf{h}_2 \cdots \xleftrightarrow{f_{\theta_n}} \mathbf{z}. \quad (1)$$

Given that we can compute the likelihood of $p(\mathbf{z})$, the likelihood of the data \mathbf{x} under the transformation f can be computed using the change of variables formula,

$$\log p_{\theta}(\mathbf{x}) = \log p(\mathbf{z}) + \sum_{i=1}^n \log |\det J_{\theta_i}|, \quad (2)$$

where $J_{\theta_i} = \partial \mathbf{h}_i / \partial \mathbf{h}_{i-1}$ is the Jacobian of the invertible transformation f_{θ_i} going from \mathbf{h}_{i-1} to \mathbf{h}_i with $\mathbf{h}_0 \equiv \mathbf{x}$. Note that most prior work [4, 8, 9, 14, 18] considers i.i.d. Gaussian likelihood models of \mathbf{z} , *e.g.* $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mu, \sigma)$.

These models, however, have limitations. First, the requirement of invertibility constrains the class of functions f_{θ_i} to be monotonically increasing (or decreasing), thus limiting expressiveness. Second, of the three possible variants of f_{θ_i} to date [41], MAF (masked autoregressive flows), IAF (inverse autoregressive flows), and SCF (split coupling flows), MAFs are difficult to parallelize due to sequential dependencies between dimensions and IAFs do not perform well in practice. SCFs strike the right balance with respect to parallelization and modeling power. In detail, SCFs partition the dimensions into two equal halves *and* transform one of

the halves \mathbf{r}_i conditioned on \mathbf{l}_i , leaving \mathbf{l}_i unchanged and thus not introducing any sequential dependencies (making parallelization easier). Examples of SCFs include the affine couplings of RealNVP [9] and MixLogCDF couplings of Flow++ [14].

In practice, SCFs are organized into blocks [8, 9, 18] to maximize efficiency such that each f_{θ_i} typically consists of SQUEEZE, STEPOFFLOW, and SPLIT operations. SQUEEZE trades off spatial resolution for channel depth. Suppose an intermediate layer \mathbf{h}_i is of size $[C_i, N_i, N_i]$, then the SQUEEZE operation transforms it into size $[4C_i, N_i/2, N_i/2]$ by re-shaping 2×2 neighborhoods into 4 channels. STEPOFFLOW is a series of SCF (possibly several) coupling layers and invertible 1×1 convolutions [8, 9, 18].

The SPLIT operation (distinct from the split couplings) splits an intermediate layer \mathbf{h}_i into two halves $\{\mathbf{l}_i, \mathbf{r}_i\}$ of size $[2C_i, N_i/2, N_i/2]$ each. Subsequent invertible layers $f_{\theta_j > i}$ operate only on \mathbf{r}_i , leaving \mathbf{l}_i unchanged. In other words, the SPLIT operation fixes some dimensions of the latent representation \mathbf{z} to \mathbf{l}_i as they are not transformed any further. This leads to a significant reduction in the amount of computation and memory needed. In the following, we denote the spatial resolutions at the n different levels as $\mathbf{N} = \{N_0, \dots, N_n\}$, with $N = N_0$ being the input resolution. Similarly, $\mathbf{C} = \{C_0, \dots, C_n\}$ denotes the number of feature channels, with $C = C_0$ being the number of input channels.

In practice, due to limited modeling flexibility, prior SCF-based models [9, 14, 18] require many SCF coupling layers in f_{θ_i} to model complex distributions, *e.g.* images. This in turn leads to high memory requirements and also leads to less efficient sampling procedures.

Autoregressive models. Autoregressive generative models are another class of powerful and highly flexible exact inference models. They factorize complex target distributions by decomposing them into a product of conditional distributions, *e.g.* images with $N \times N$ spatial resolution as $p(\mathbf{x}) = \prod_{i=1}^N \prod_{j=1}^N p(\mathbf{x}_{i,j} | \mathbf{x}_{\text{pred}(i,j)})$ [10, 12, 25, 37, 38]. Here, $\text{pred}(i, j)$ denotes the set of predecessors of pixel (i, j) . The functional form of these conditionals can be highly flexible, and allows such models to capture complex multimodal distributions. However, such a dependency structure only allows for image synthesis via ancestral sampling by generating each pixel sequentially, conditioned on the previous pixels [37, 38], making parallelization difficult. This is also inefficient since autoregressive models, including PixelCNN and PixelRNN, require $\mathcal{O}(N^2)$ time steps for sampling.

4. Multi-scale Autoregressive Flow Priors

We propose to leverage the strengths of autoregressive models to improve invertible normalizing flow models such as [9, 18]. Specifically, we propose novel *multi-scale autoregressive priors for split coupling flows (mAR-SCF)*. Using them allows us to learn complex multimodal latent priors $p(\mathbf{z})$ in multi-scale SCF models, *cf.* Eq. (2). This is unlike [9, 14, 18, 28], which rely on Gaussian priors in the latent space. Additionally, we also propose a scheme for interpolation in the latent space of our *mAR-SCF* models.

The use of our novel autoregressive *mAR* priors for invertible flow models has two distinct advantages over both vanilla SCF and autoregressive models. First, the powerful autoregressive prior helps mitigate the limited modeling capacity of the vanilla SCF flow models. Second, as only the prior is autoregressive, this makes flow models with our *mAR* prior an order of magnitude faster with respect to sampling time than fully autoregressive models. Next, we describe our multi-scale autoregressive prior in detail.

Our *mAR-SCF* model uses an efficient invertible split coupling flow $f_{\theta_i}(\mathbf{x})$ to map the distribution over the data \mathbf{x} to a latent variable \mathbf{z} and then models an autoregressive *mAR* prior over \mathbf{z} , parameterized by ϕ . The likelihood of a data point \mathbf{x} of dimensionality $[C, N, N]$ can be expressed as

$$\log p_{\theta, \phi}(\mathbf{x}) = \log p_{\phi}(\mathbf{z}) + \sum_{i=1}^n \log |\det J_{\theta_i}|. \quad (3)$$

Here, J_{θ_i} is the Jacobian of the invertible transformations f_{θ_i} . Note that, as $f_{\theta_i}(\mathbf{x})$ is an invertible function, \mathbf{z} has the same total dimensionality as the input data point \mathbf{x} .

Formulation of the *mAR* prior. We now introduce our *mAR* prior $p_{\phi}(\mathbf{z})$ along with our *mAR-SCF* model, which combines the split coupling flows f_{θ_i} with an *mAR* prior. As shown in Fig. 2, our *mAR* prior is applied after every SPLIT operation of the invertible flow layers as well as at the smallest spatial resolution. Let $\mathbf{l}_i = \{\mathbf{l}_i^1, \dots, \mathbf{l}_i^{C_i}\}$ be the C_i

channels of size $[C_i, N_i, N_i]$, which do not undergo further transformation f_{θ_i} after the SPLIT at level i . Following the SPLIT at level i , our *mAR* prior is modeled as a conditional distribution, $p_{\phi}(\mathbf{l}_i | \mathbf{r}_i)$; at the coarsest spatial resolution it is an unconditional distribution, $p_{\phi}(\mathbf{h}_n)$. Thereby, we assume that our *mAR* prior at each level i autoregressively factorizes along the channel dimension as

$$p_{\phi}(\mathbf{l}_i | \mathbf{r}_i) = \prod_{j=1}^{C_i} p_{\phi}(\mathbf{l}_i^j | \mathbf{l}_i^1, \dots, \mathbf{l}_i^{j-1}, \mathbf{r}_i). \quad (4)$$

Furthermore, the distribution at each spatial location (m, n) within a channel \mathbf{l}_i^j is modeled as a conditional Gaussian,

$$p_{\phi}(\mathbf{l}_{i(m,n)}^j | \mathbf{l}_i^1, \dots, \mathbf{l}_i^{j-1}, \mathbf{r}_i) = \mathcal{N}(\mu_{i(m,n)}^j, \sigma_{i(m,n)}^j). \quad (5)$$

Thus, the mean, $\mu_{i(m,n)}^j$ and variance, $\sigma_{i(m,n)}^j$ at each spatial location are autoregressively modeled along the channels. This allows the distribution at each spatial location to be highly flexible and capture multimodality in the latent space. Moreover from Eq. (4), our *mAR* prior can model long-range correlations in the latent space as the distribution of each channel is dependent on all previous channels.

This autoregressive factorization allows us to employ Conv-LSTMs [34] to model the distributions $p_{\phi}(\mathbf{l}_i^j | \mathbf{l}_i^1, \dots, \mathbf{l}_i^{j-1}, \mathbf{r}_i)$ and $p_{\phi}(\mathbf{h}_n)$. Conv-LSTMs can model long-range dependencies across channels in their internal state. Additionally, long-range spatial dependencies within channels can be modeled by stacking multiple Conv-LSTM layers with a wide receptive field. This formulation allows all pixels within a channel to be sampled in parallel, while the channels are sampled in a sequential manner,

$$\hat{\mathbf{l}}_i^j \sim p_{\phi}(\mathbf{l}_i^j | \mathbf{l}_i^1, \dots, \mathbf{l}_i^{j-1}, \mathbf{r}_i). \quad (6)$$

This is in contrast to PixelCNN/RNN-based models, which sample one pixel at a time.

The *mAR-SCF* model. We illustrate our *mAR-SCF* model architecture in Fig. 2b. Our *mAR-SCF* model leverages the SQUEEZE and SPLIT operations for invertible flows introduced in [8, 9] for efficient parallelization. Following [8, 9, 18], we use several SQUEEZE and SPLIT operations in a multi-scale setup at n scales (Fig. 2b) until the spatial resolution at \mathbf{h}_n is reasonably small, typically 4×4 . Note that there is no SPLIT operation at the smallest spatial resolution. Therefore, the latent space is the concatenation of $\mathbf{z} = \{\mathbf{l}_1, \dots, \mathbf{l}_{n-1}, \mathbf{h}_n\}$. The split coupling flows (SCF) f_{θ_i} in the *mAR-SCF* model remain invertible by construction. We consider different SCF couplings for f_{θ_i} , including the affine couplings of [9, 18] and MixLogCDF couplings [14].

Given the parameters ϕ of our multimodal *mAR* prior modeled by the Conv-LSTMs, we can compute $p_{\phi}(\mathbf{z})$ using the formulation in Eqs. (4) and (5). We can thus express

Algorithm 1: MARPS: Multi-scale Autoregressive Prior Sampling for our *mAR-SCF* models

```

1 Sample  $\hat{\mathbf{h}}_n \sim p_\phi(\mathbf{h}_n)$ ;
2 for  $i \leftarrow n - 1$  to 1 do
    /* SPLITINVERSE */
3    $\hat{\mathbf{r}}_i \leftarrow \hat{\mathbf{h}}_{i+1}$ ;           // Assign previous
4    $\hat{\mathbf{l}}_i \sim p_\phi(\mathbf{l}_i | \mathbf{r}_i)$ ;       // Sample mAR prior
5    $\hat{\mathbf{h}}_i \leftarrow \{\hat{\mathbf{l}}_i, \hat{\mathbf{r}}_i\}$ ;       // Concatenate
    /* STEPOFFLOWINVERSE */
6   Apply  $f_i^{-1}(\hat{\mathbf{h}}_i)$ ;           // SCF coupling
    /* SQUEEZEINVERSE */
7   Reshape  $\hat{\mathbf{h}}_i$ ;           // Depth to Space
8 end
9  $\mathbf{x} \leftarrow \hat{\mathbf{h}}_1$ ;

```

Eq. (3) in *closed form* and directly maximize the likelihood of the data under the multimodal *mAR* prior distribution learned by the Conv-LSTMs.

Next, we show that the computational cost of our *mAR-SCF* model is $\mathcal{O}(N)$ for sampling an image of size $[C, N, N]$; this is in contrast to the standard $\mathcal{O}(N^2)$ computational cost required by purely autoregressive models.

Analysis of sampling time. We now formally analyze the computational cost in the number of steps T required for sampling with our *mAR-SCF* model. First, we describe the sampling process in detail in Algorithm 1 (the forward training process follows the sampling process in reverse order). Next, we derive the worst-case number of steps T required by MARPS, given sufficient parallel resources to sample a channel in parallel. Here, the number of steps T can be seen as the length of the critical path while sampling.

Lemma 4.1. *Let the sampled image \mathbf{x} be of resolution $[C, N, N]$, then the worst-case number of steps T (length of the critical path) required by MARPS is $\mathcal{O}(N)$.*

Proof. At the first sampling step (Fig. 2a) at layer f_{θ_n} , our *mAR* prior is applied to generate \mathbf{h}_n , which is of shape $[2^{n+1}C, N/2^n, N/2^n]$. Therefore, the number of sequential steps required at the *last flow layer* \mathbf{h}_n is

$$T_n = C \cdot 2^{n+1}. \quad (7)$$

Here, we are assuming that each channel can be sampled in parallel in one time-step.

From $f_{\theta_{n-1}}$ to f_{θ_1} , f_{θ_i} always contains a SPLIT operation. Therefore, at each f_{θ_i} we use our *mAR* prior to sample \mathbf{l}_i , which has shape $[2^i C, N/2^i, N/2^i]$. Therefore, the number of sequential steps required for sampling at layers \mathbf{h}_i , $1 \leq i < n$ of our *mAR-SCF* model is

$$T_i = C \cdot 2^i. \quad (8)$$

Therefore, the total number of sequential steps (length of the critical path) required for sampling is

$$\begin{aligned} T &= T_n + T_{n-1} + \dots + T_i + \dots + T_1 \\ &= C \cdot (2^{n+1} + 2^{n-1} + \dots + 2^i + \dots + 2^1) \\ &= C \cdot (3 \cdot 2^n - 2). \end{aligned} \quad (9)$$

Now, the total number of layers in our *mAR-SCF* model is $n \leq \log(N)$. This is because each layer reduces the spatial resolution by a factor of two. Therefore, the total number of time-steps required is

$$T \leq 3 \cdot C \cdot N. \quad (10)$$

In practice, $C \ll N$, with $C = C_0 = 3$ for RGB images. Therefore, the total number of sequential steps required for sampling in our *mAR-SCF* model is $T = \mathcal{O}(N)$. \square

It follows that with our multi-scale autoregressive *mAR* priors in our *mAR-SCF* model, sampling can be performed in a linear number of time-steps in contrast to fully autoregressive models like PixelCNN, which require a quadratic number of time-steps [38].

Interpolation. A major advantage of invertible flow-based models is that they allow for latent spaces, which are useful for downstream tasks like interpolation – smoothly transforming one data point into another. Interpolation is simple in case of typical invertible flow-based models, because the latent space is modeled as a unimodal i.i.d. Gaussian. To allow interpolation in the space of our multimodal *mAR* priors, we develop a simple method based on [2].

Let \mathbf{x}_A and \mathbf{x}_B be the two images (points) to be interpolated and \mathbf{z}_A and \mathbf{z}_B be the corresponding points in the latent space. We begin with an initial linear interpolation between the two latent points, $\{\mathbf{z}_A, \mathbf{z}_{A,B}^1, \dots, \mathbf{z}_{A,B}^k, \mathbf{z}_B\}$, such that, $\mathbf{z}_{A,B}^i = (1 - \alpha^i) \mathbf{z}_A + \alpha^i \mathbf{z}_B$. The initial linearly interpolated points $\mathbf{z}_{A,B}^i$ may not lie in a high-density region under our multimodal prior, leading to non-smooth transformations. Therefore, we next project the interpolated points $\mathbf{z}_{A,B}^i$ to a high-density region, without deviating too much from their initial position. This is possible because our *mAR* prior allows for exact inference. However, the image corresponding to the projected $\bar{\mathbf{z}}_{A,B}^i$ must also not deviate too far from either \mathbf{x}_A and \mathbf{x}_B either to allow for smooth transitions. To that end, we define the projection operation as

$$\begin{aligned} \bar{\mathbf{z}}_{A,B}^i &= \arg \min \left(\|\bar{\mathbf{z}}_{A,B}^i - \mathbf{z}_{A,B}^i\| - \lambda_1 \log p_\phi(\bar{\mathbf{z}}_{A,B}^i) \right. \\ &\quad \left. + \lambda_2 \min \left(\|f^{-1}(\bar{\mathbf{z}}_{A,B}^i) - \mathbf{x}_A\|, \|f^{-1}(\bar{\mathbf{z}}_{A,B}^i) - \mathbf{x}_B\| \right) \right), \end{aligned} \quad (11)$$

where λ_1, λ_2 are the regularization parameters. The term controlled by λ_1 pulls the interpolated $\mathbf{z}_{A,B}^i$ back to high-density regions, while the term controlled by λ_2 keeps the result close to the two images \mathbf{x}_A and \mathbf{x}_B . Note that this reduces to linear interpolation when $\lambda_1 = \lambda_2 = 0$.

Method	Coupling	MNIST				CIFAR10			
		Levels	SCF	Channels	bits/dim (\downarrow)	Levels	SCF	Channels	bits/dim (\downarrow)
PixelCNN [38]	Autoregressive	–	–	–	–	–	–	–	3.00
PixelCNN++ [37]	Autoregressive	–	–	–	–	–	–	–	2.92
Glow [18]	Affine	3	32	512	1.05	3	32	512	3.35
Flow++ [14]	MixLogCDF	–	–	–	–	3	–	96	3.29
Residual Flow [4]	Residual	3	16	–	0.97	3	16	–	3.28
<i>mAR-SCF</i> (Ours)	Affine	3	32	256	1.04	3	32	256	3.33
<i>mAR-SCF</i> (Ours)	Affine	3	32	512	1.03	3	32	512	3.31
<i>mAR-SCF</i> (Ours)	MixLogCDF	3	4	96	0.88	3	4	96	3.27
<i>mAR-SCF</i> (Ours)	MixLogCDF	–	–	–	–	3	4	256	3.24

Table 1. Evaluation of our *mAR-SCF* model on MNIST and CIFAR10 (using uniform dequantization for fair comparison with [4, 18]).

5. Experiments

We evaluate our approach on the MNIST [22], CIFAR10 [21], and ImageNet [38] datasets. In comparison to datasets like CelebA, CIFAR10 and ImageNet are highly multimodal and the performance of invertible SCF models has lagged behind autoregressive models in density estimation and behind GAN-based generative models regarding image quality.

5.1. MNIST and CIFAR10

Architecture details. Our *mAR* prior at each level f_{θ_i} consists of three convolutional LSTM layers, each of which uses 32 convolutional filters to compute the input-to-state and state-to-state components. Keeping the *mAR* prior architecture constant, we experiment with different SCF couplings in f_{θ_i} to highlight the effectiveness of our *mAR* prior. We experiment with affine couplings of [9, 18] and MixLogCDF couplings [14]. Affine couplings have limited modeling flexibility. The more expressive MixLogCDF applies the cumulative distribution function of a mixture of logistics. In the following, we include experiments varying the number couplings and the number of channels in the convolutional blocks of the neural networks used to predict the affine/MixLogCDF transformation parameters.[†]

Hyperparameters. We use Adamax (as in [18]) with a learning rate of 8×10^{-4} . We use a batch size of 128 with affine and 64 with MixLogCDF couplings (following [14]).

Density estimation. We report density estimation results on MNIST and CIFAR10 in Table 1 using the per-pixel log-likelihood metric in bits/dim. We also include the architecture details (# of levels, coupling type, # of channels). We compare to the state-of-the-art Flow++ [14] method with SCF couplings and Residual Flows [4]. Note that in terms of architecture, our *mAR-SCF* model with affine couplings is closest to that of Glow [18]. Therefore, the comparison with Glow serves as an ideal ablation to assess the effectiveness of our *mAR* prior. Flow++ [14], on the other hand, uses the more powerful MixLogCDF transformations and their

model architecture does not include SPLIT operations. Because of this, Flow++ has higher computational and memory requirements for a given batch size compared to Glow. Furthermore, for fair comparison with Glow [18] and Residual flows [4], we use uniform dequantization unlike Flow++, which proposes to use variational dequantization.

In comparison to Glow, we achieve improved density estimation results on both MNIST and CIFAR10. In detail, we outperform Glow (*e.g.* 1.05 vs. 1.04 bits/dim on MNIST and 3.35 vs. 3.33 bits/dim on CIFAR10) with |SCF|= 32 affine couplings and 3 levels, while using parameter prediction networks with only half (256 vs. 512) the number of channels. We observe that increasing the capacity of our parameter prediction networks to 512 channels boosts the log-likelihood further to 1.03 bits/dim on MNIST and 3.31 bits/dim on CIFAR10. As this setting with 512 channels is identical to setting reported in [18], this shows that our *mAR* prior boosts the accuracy by ~ 0.04 bits/dim in case of CIFAR10. To place this performance gain in context, it is competitive with the ~ 0.03 bits/dim boost reported in [18] (*cf.* Fig. 3 in [18]) with the introduction of the 1×1 convolution. We train our model for ~ 3000 epochs, similar to [18]. Also note that we only require a batch size of 128 to achieve state-of-the-art likelihoods, whereas Glow uses batches of size 512. Thus our *mAR-SCF* model improves density estimates and requires significantly lower computational resources (~ 48 vs. ~ 128 GB memory). Overall, we also observe competitive sampling speed (see also Table 2). This firmly establishes the utility of our *mAR-SCF* model.

For fair comparison with Flow++ [14] and Residual Flows [4], we employ the more powerful MixLogCDF couplings. Our *mAR-SCF* model uses 4 MixLogCDF couplings at each level with 96 channels but includes SPLIT operations unlike Flow++. Here, we outperform Flow++ and Residual Flows (3.27 vs. 3.29 and 3.28 bits/dim on CIFAR10) while being equally fast to sample as Flow++ (Table 2). A baseline model without our *mAR* prior has performance comparable to Flow++ (3.29 bits/dim). Similarly on MNIST, our *mAR-SCF* model again outperforms Residual Flows (0.88 vs.

[†]Code available at: <https://github.com/visinf/mar-scf>



Figure 3. Comparison of random samples from our *mAR-SCF* model with state-of-the-art models.

0.97 bits/dim). Finally, we train a more powerful *mAR-SCF* model with 256 channels with sampling speed competitive with [4], which achieves state-of-the-art 3.24 bits/dim on CIFAR10 [‡]. This is attained after ~ 400 training epochs (comparable to ~ 350 epochs required by [4] to achieve 3.28 bits/dim). Next, we compare the sampling speed of our *mAR-SCF* model with that of Flow++ and Residual Flow.

Method	Coupling	Levels	SCF	Ch.	Speed (ms, ↓)
Glow [18]	Affine	3	32	512	13
Flow++ [14]	MixLogCDF	3	–	96	19
Residual Flow [4]	Residual	3	16	–	34
PixelCNN++ [33]	Autoregressive	–	–	–	5×10^3
<i>mAR-SCF</i> (Ours)	Affine	3	32	256	6
<i>mAR-SCF</i> (Ours)	Affine	3	32	512	17
<i>mAR-SCF</i> (Ours)	MixLogCDF	3	4	96	19
<i>mAR-SCF</i> (Ours)	MixLogCDF	3	4	256	32

Table 2. Evaluation of sampling speed with batches of size 32.

Sampling speed. We report the sampling speed of our *mAR-SCF* model in Table 2 in terms of sampling one image on CIFAR10. We report the average over 1000 runs using a batch size of 32. We performed all tests on a single Nvidia V100 GPU with 32GB of memory. First, note that our *mAR-*

[‡]*mAR-SCF* with 256 channels trained on CIFAR10 (~ 1000 epochs) trained to convergence achieves 3.22 bits/dim

SCF model with affine coupling layers in 3 levels with 512 channels needs 17 ms on average to sample an image. This is comparable with Glow, which requires 13 ms. This shows that our *mAR* prior causes only a slight increase in sampling time – particularly because our *mAR-SCF* requires only $\mathcal{O}(N)$ steps to sample and the prior has far fewer parameters compared to the invertible flow network. Moreover, our *mAR-SCF* model with affine coupling layers with 256 channels is considerably faster (6 vs. 13 ms) with an accuracy advantage. Similarly, our *mAR-SCF* with MixLogCDF and 96 channels is competitive in speed with [14] with an accuracy advantage and considerably faster than [4] (19 vs. 34 ms). This is because Residual Flows are slower to invert (sample) as there is no closed-form expression of the inverse. Furthermore, our *mAR-SCF* with MixLogCDF and 256 channels is competitive with respect to [4] in terms of sampling speed while having a large accuracy advantage. Finally, note that these sampling speeds are two orders of magnitude faster than state-of-the-art fully autoregressive approaches, e.g. PixelCNN++ [33].

Sample quality. Next, we analyze the sample quality of our *mAR-SCF* model in Table 3 using the FID metric [13] and Inception scores [32]. The analysis of sample quality is important as it is well-known that visual fidelity and test log-likelihoods are not necessarily indicative of each other [35].



Figure 4. Interpolations of our *mAR-SCF* model on CIFAR10.



(a) Residual Flows [4] (3.75 bits/dim)



(b) Our *mAR-SCF* (Affine, 3.80 bits/dim)

Figure 5. Random samples on ImageNet (64×64).

We achieve an FID of 41.0 and an Inception score of 5.7 with our *mAR-SCF* model with affine couplings, significantly better than Glow with the same specifications and Residual Flows. While our *mAR-SCF* model with MixLogCDF couplings also performs comparably, empirically we find affine couplings to lead to better image quality as in [4]. We show random samples from our *mAR-SCF* model with both affine and MixLogCDF couplings in Fig. 3. Here, we compare to the version of Flow++ with MixLogCDF couplings and variational dequantization (which gives even better log-likelihoods) and Residual Flows. Our *mAR-SCF* model achieves better sample quality with more clearly defined objects. Furthermore, we also obtain improved sample quality over both PixelCNN and PixelIQN and close the gap in comparison to adversarial approaches like DCGAN [27] and WGAN-GP [40]. This highlights that our *mAR-SCF* model is able to better capture long-range correlations.

Interpolation. We show interpolations on CIFAR10 in Fig. 4, obtained using Eq. (11). We observe smooth interpolation between images belonging to distinct classes. This shows that the latent space of our *mAR* prior can be potentially used for downstream tasks similarly to Glow [18]. We include additional analyses in Appendix E.

5.2. ImageNet

Finally, we evaluate our *mAR-SCF* model on ImageNet (32×32 and 64×64) against the best performing models on

Method	Coupling	FID (\downarrow)	Inception Score (\uparrow)
PixelCNN [38]	Autoregressive	65.9	4.6
PixelIQN [24]	Autoregressive	49.4	–
Glow [18]	Affine	46.9	–
Residual Flow [4]	Residual	46.3	5.2
<i>mAR-SCF</i> (Ours)	MixLogCDF	41.9	5.7
<i>mAR-SCF</i> (Ours)	Affine	41.0	5.7
DCGAN [27]	Adversarial	37.1	6.4
WGAN-GP [40]	Adversarial	36.4	6.5

Table 3. Evaluation of sample quality on CIFAR10. Other results are quoted from [4, 24].

Method	Coupling	SCF	Ch.	bits/dim (\downarrow)	Mem (GB, \downarrow)
Glow [18]	Affine	32	512	4.09	~ 128
Residual Flow [4]	Residual	32	–	4.01	–
<i>mAR-SCF</i> (Ours)	Affine	32	256	4.07	~ 48
<i>mAR-SCF</i> (Ours)	MixLogCDF	4	460	3.99	~ 80

Table 4. Evaluation on ImageNet (32×32).

MNIST and CIFAR10 in Table 4, *i.e.* Glow [18] and Residual Flows [4]. Our model with affine couplings outperforms Glow while using fewer channels (4.07 vs. 4.09 bits/dim). For comparison with the more powerful Residual Flow models, we use four MixLogCDF couplings at each layer f_{θ_i} with 460 channels. We again outperform Residual Flows [4] (3.99 vs. 4.01 bits/dim). These results are consistent with the findings in Table 1, highlighting the advantage of our *mAR* prior. Finally, we also evaluate on the ImageNet (64×64) dataset. Our *mAR-SCF* model with affine flows achieves 3.80 vs. 3.81 bits/dim in comparison to Glow [18]. We show qualitative examples in Fig. 5 and compare to Residual Flows. We see that although the powerful Residual Flows obtain better log-likelihoods (3.75 bits/dim), our *mAR-SCF* model achieves better visual fidelity. This again highlights that our *mAR* is able to better capture long-range correlations.

6. Conclusion

We presented *mAR-SCF*, a flow-based generative model with novel multi-scale autoregressive priors for modeling long-range dependencies in the latent space of flow models. Our *mAR* prior considerably improves the accuracy of flow-based models with split coupling layers. Our experiments show that not only does our *mAR-SCF* model improve density estimation (in terms of bits/dim), but also considerably improves the sample quality of the generated images compared to previous state-of-the-art exact inference models. We believe the combination of complex priors with flow-based models, as demonstrated by our *mAR-SCF* model, provides a path toward efficient models for exact inference that approach the fidelity of GAN-based approaches.

Acknowledgement. SM and SR acknowledge the support by the German Research Foundation as part of the Research Training Group *Adaptive Preparation of Information from Heterogeneous Sources (AIPHES)* under grant No. GRK 1994/1.

References

- [1] Jens Behrmann, Will Grathwohl, Ricky T. Q. Chen, David Duvenaud, and Jörn-Henrik Jacobsen. Invertible residual networks. In *ICML*, pages 573–582, 2019. 2
- [2] Christoph Breger and Stephen M. Omohundro. Nonlinear image interpolation using manifold learning. In *NIPS*, pages 973–980, 1994. 5
- [3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *ICLR*, 2019. 1
- [4] Ricky T. Q. Chen, Jens Behrmann, David Duvenaud, and Jörn-Henrik Jacobsen. Residual flows for invertible generative modeling. In *NeurIPS*, pages 9913–9923, 2019. 2, 3, 6, 7, 8, 11, 13, 14
- [5] Xi Chen, Diederik P. Kingma, Tim Salimans, Yan Duan, Prafulla Dhariwal, John Schulman, Ilya Sutskever, and Pieter Abbeel. Variational lossy autoencoder. In *ICLR*, 2017. 2
- [6] Xi Chen, Nikhil Mishra, Mostafa Rohaninejad, and Pieter Abbeel. Pixelsnail: An improved autoregressive generative model. In *ICML*, pages 864–872, 2018. 2
- [7] Sander Dieleman, Aäron van den Oord, and Karen Simonyan. The challenge of realistic music generation: Modelling raw audio at scale. In *NeurIPS*, pages 7989–7999, 2018. 1
- [8] Laurent Dinh, David Krueger, and Yoshua Bengio. NICE: Non-linear independent components estimation. In *ICLR Workshop*, 2015. 1, 2, 3, 4
- [9] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In *ICLR*, 2017. 1, 2, 3, 4, 6
- [10] Justin Domke, Alap Karapurkar, and Yiannis Aloimonos. Who killed the directed model? In *CVPR*, 2008. 1, 2, 4
- [11] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014. 1
- [12] Alex Graves. Generating sequences with recurrent neural networks. *arXiv:1308.0850*, 2013. 2, 4
- [13] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In *NeurIPS*, pages 6626–6637, 2017. 2, 7
- [14] Jonathan Ho, Xi Chen, Aravind Srinivas, Yan Duan, and Pieter Abbeel. Flow++: Improving flow-based generative models with variational dequantization and architecture design. In *ICML*, pages 2722–2730, 2019. 2, 3, 4, 6, 7, 11, 14
- [15] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997. 2
- [16] Emiel Hoogeboom, Rianne van den Berg, and Max Welling. Emerging convolutions for generative normalizing flows. In *ICML*, pages 2771–2780, 2019. 2
- [17] Nal Kalchbrenner, Erich Elsen, Karen Simonyan, Seb Noury, Norman Casagrande, Edward Lockhart, Florian Stimberg, Aäron van den Oord, Sander Dieleman, and Koray Kavukcuoglu. Efficient neural audio synthesis. In *ICML*, pages 2410–2419, 2018. 2
- [18] Diederik P. Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *NeurIPS*, pages 10215–10224, 2018. 1, 2, 3, 4, 6, 7, 8
- [19] Diederik P. Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. In *NIPS*, pages 4743–4751, 2016. 2
- [20] Diederik P. Kingma and Max Welling. Auto-encoding variational Bayes. In *ICLR*, 2014. 1
- [21] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, U. of Toronto, 2009. 2, 6
- [22] Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86(11):2278–2324, 1998. 2, 6
- [23] Jacob Menick and Nal Kalchbrenner. Generating high fidelity images with subscale pixel networks and multidimensional upscaling. In *ICLR*, 2019. 2
- [24] Georg Ostrovski, Will Dabney, and Rémi Munos. Autoregressive quantile networks for generative modeling. In *ICML*, pages 3936–3945, 2018. 8
- [25] George Papamakarios, Iain Murray, and Theo Pavlakou. Masked autoregressive flow for density estimation. In *NeurIPS*, pages 2338–2347, 2017. 2, 4
- [26] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In *ICML*, pages 4055–4064, 2018. 2
- [27] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016. 2, 8
- [28] Ali Razavi, Aäron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with VQ-VAE-2. In *NeurIPS*, pages 14837–14847, 2019. 1, 2, 4
- [29] Scott E. Reed, Aäron van den Oord, Nal Kalchbrenner, Sergio Gomez Colmenarejo, Ziyu Wang, Yutian Chen, Dan Belov, and Nando de Freitas. Parallel multiscale autoregressive density estimation. In *ICML*, pages 2912–2921, 2017. 2
- [30] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *ICML*, pages 1278–1286, 2014. 1
- [31] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(13):211–252, 2015. 2
- [32] Tim Salimans, Ian J. Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training GANs. In *NIPS*, pages 2234–2242, 2016. 2, 7

- [33] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P. Kingma. PixelCNN++: Improving the PixelCNN with discretized logistic mixture likelihood and other modifications. In *ICLR*, 2017. [2](#), [7](#)
- [34] Xingjian Shi, Zhihan Gao, Leonard Lausen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Deep learning for precipitation nowcasting: A benchmark and a new model. In *NeurIPS*, pages 5617–5627, 2017. [4](#)
- [35] Lucas Theis, Aäron van den Oord, and Matthias Bethge. A note on the evaluation of generative models. In *ICLR*, 2016. [7](#)
- [36] Jakob M. Tomczak and Max Welling. VAE with a VampPrior. In *AISTATS*, pages 1214–1223, 2018. [2](#)
- [37] Aäron van den Oord, Nal Kalchbrenner, Lasse Espeholt, Koray Kavukcuoglu, Oriol Vinyals, and Alex Graves. Conditional image generation with PixelCNN decoders. In *NIPS*, pages 4790–4798, 2016. [1](#), [2](#), [4](#), [6](#)
- [38] Aäron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *ICML*, pages 1747–1756, 2016. [1](#), [2](#), [4](#), [5](#), [6](#), [8](#), [11](#), [13](#)
- [39] Aäron van den Oord, Yazhe Li, Igor Babuschkin, Karen Simonyan, Oriol Vinyals, Koray Kavukcuoglu, George van den Driessche, Edward Lockhart, Luis C. Cobo, Florian Stimberg, Norman Casagrande, Dominik Grewe, Seb Noury, Sander Dieleman, Erich Elsen, Nal Kalchbrenner, Heiga Zen, Alex Graves, Helen King, Tom Walters, Dan Belov, and Demis Hassabis. Parallel WaveNet: Fast high-fidelity speech synthesis. In *ICML*, pages 3918–3926, 2018. [1](#), [2](#)
- [40] Xiang Wei, Boqing Gong, Zixia Liu, Wei Lu, and Liqiang Wang. Improving the improved training of Wasserstein GANs: A consistency term and its dual effect. In *ICLR*, 2018. [2](#), [8](#)
- [41] Zachary M. Ziegler and Alexander M. Rush. Latent normalizing flows for discrete sequences. In *ICML*, pages 7673–7682, 2019. [2](#), [3](#)

Appendices

We provide additional details of Lemma 4.1 in the main paper, additional details of our *mAR-SCF* model architecture, as well as additional results and qualitative examples.

A. Channel Dimensions in *mAR-SCF*

We begin by providing additional details of Lemma 4.1 in the main paper. We formally prove the claim that the number of channels at the last layer n (which does not have a SPLIT operation) is $C_n = 2^{n+1} \cdot C$ for an image of size $[C, N, N]$. Thereby, we also show that the number of channels at any layer i (with a SPLIT operation) is $C_i = 2^i \cdot C$.

Note that the number of time steps required for sampling depends on the number of channels for *mAR-SCF* (flow model) based on *MARPS* (Algorithm 1).

Lemma A.1. *Let the size of the sampled image be $[C, N, N]$. The number of channels at the last layer \mathbf{h}_n is $C_n = 2^{n+1} \cdot C$.*

Proof. The shape of image to be sampled is $[C, N, N]$. Since the network is invertible, the size of the image at level \mathbf{h}_0 is the size of the image sampled, *i.e.* $[C, N, N]$.

First, let the number of layers n be 1. This implies that during the forward pass the network applies a SQUEEZE operation, which reshapes the image to $[4C, N/2, N/2]$ followed by STEPOFFLOW, which does not modify the shape of the input, *i.e.* the shape remains $[4C, N/2, N/2]$. Thus the number of channels at the last layer \mathbf{h}_1 is $C_1 = 4 \cdot C = 2^{1+1} \cdot C$ when $n = 1$.

Next, let us assume that the number of channels at the last layer \mathbf{h}_{k-1} for a flow with $k - 1$ layers is

$$C_{k-1} = 2^k \cdot C. \quad (12)$$

We aim to prove by induction that the number of channels at the last layer \mathbf{h}_k for a flow with k layers then is

$$C_k = 2^{k+1} \cdot C. \quad (13)$$

To that end, we note that the dimensionality of the output at layer $k - 1$ after SQUEEZE and STEPOFFLOW from Eq. (12) by assumption is $C_{k-1} = 2^k \cdot C$. For a flow with k layers, the $(k - 1)$ st layer has a SPLIT operation resulting in $\{\mathbf{l}_{k-1}, \mathbf{r}_{k-1}\}$, each with size $[2^{k-1}C, N/2^{k-1}, N/2^{k-1}]$. The number of channels for \mathbf{r}_{k-1} at layer $k - 1$ is thus $2^k \cdot C/2 = 2^{k-1} \cdot C$. At layer k the input with $2^{k-1} \cdot C$ channels is transformed by SQUEEZE to $2^{k-1} \cdot 4 \cdot C = 2^{(k-1)+2} \cdot C = 2^{k+1} \cdot C$ channels.

Therefore, by induction Eq. (13) holds for $k = n$. Thus the number of channels at the last layer \mathbf{h}_n is given as $C_n = 2^{n+1} \cdot C$. \square

B. *mAR-SCF* Architecture

In Fig. 6, we show the architecture of our *mAR* prior in detail for a layer i of *mAR-SCF*. The network is a convolutional LSTM with three LSTM layers. The input at each time-step is the previous channel of \mathbf{l}_i , concatenated with the output after convolution on \mathbf{r}_i . Our *mAR* prior autoregressively outputs the probability of each channel \mathbf{l}_i^j given by the distribution $p_\phi(\mathbf{l}_i^j | \mathbf{l}_i^{j-1}, \dots, \mathbf{l}_i^1, \mathbf{r}_i)$ modeled as $\mathcal{N}(\mu_i^j, \sigma_i^j)$ during inference. Because of the internal state of the Convolutional LSTM (memory), our *mAR* prior can learn long-range dependencies.

In Fig. 2 in the main paper, the STEPOFFLOW operation consists of an activation normalization layer, an invertible 1×1 convolution layer followed by split coupling layers (32 layers for affine couplings or 4 layers of MixLogCDF at each level).

C. Empirical Analysis of Sampling Speed

In Fig. 7, we analyze the real-world sampling time of our state-of-the-art *mAR-SCF* model with MixLogCDF couplings using varying image sizes $[C, N, N]$. In particular, we vary the input spatial resolution N . We report the mean and variance of 1000 runs with batch size of 8 on an Nvidia V100 GPU with 32GB memory. Using smaller batch sizes of 8 ensures that we always have enough parallel resources for all image sizes. Under these conditions, we see that the sampling time increases linearly with the image size. This is because the number of sampling steps of our *mAR-SCF* model scales as $\mathcal{O}(N)$, *cf.* Lemma 4.1.

D. Additional Qualitative Examples

We include additional qualitative examples for training our generative model on the MNIST, CIFAR10, and ImageNet (32×32) datasets. In Fig. 9, we see that the visual sample quality of our *mAR-SCF* model with MixLogCDF couplings is competitive with those of the state-of-the-art Residual Flows [4]. Moreover, note that we achieve better test log-likelihoods (0.88 *vs.* 1.00 bits/dim).

We additionally provide qualitative examples for ImageNet (32×32) in Figure 11. We observe that in comparison to the state-of-the-art Flow++ [14] and Residual Flows [4], the images generated by our *mAR-SCF* model are detailed with a competitive visual quality.

Finally, on CIFAR10 we compare with the fully autoregressive PixelCNN model [38] in Figure 10. We see that although the fully autoregressive PixelCNN achieves significantly better test log-likelihoods (3.00 *vs.* 3.24 bits/dim), our *mAR-SCF* models achieves better visual sample quality (also shown by the FID and Inception metrics in Table 3 of the main paper).

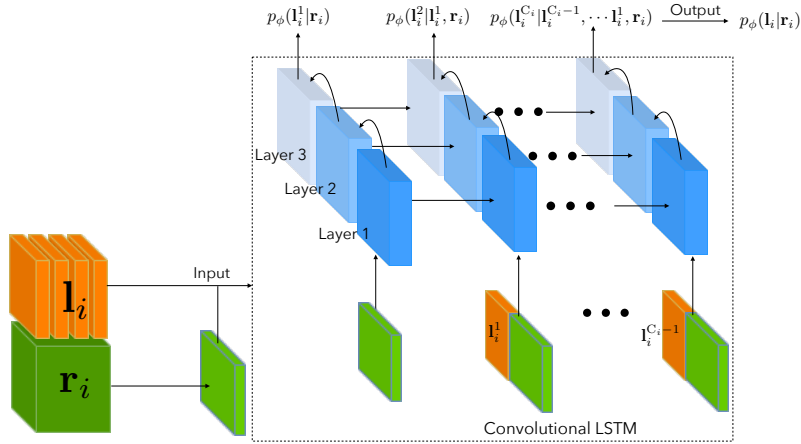


Figure 6. Architecture of our multi-scale autoregressive (*mAR*) prior.

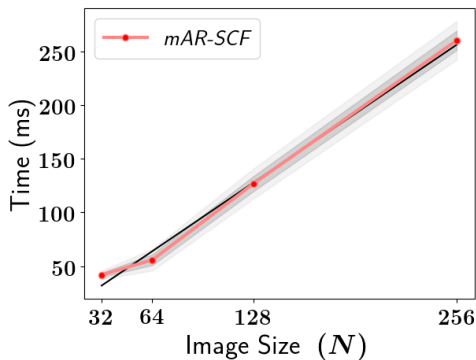


Figure 7. Analysis of real-world sampling time of our *mAR-SCF* model with varying image size $[C, N, N]$. For reference, we show the line $y = 2x$ in black.

E. Additional Interpolation Results

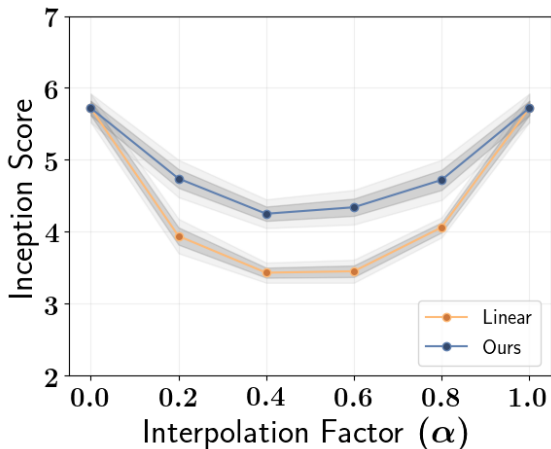
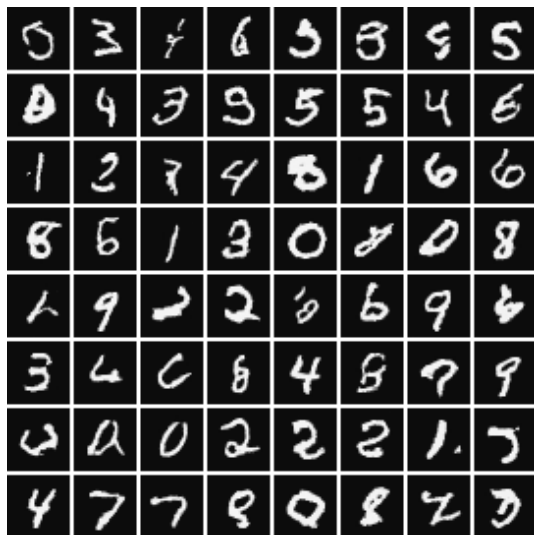


Figure 8. Inception scores of random samples generated with our interpolation scheme versus linear interpolation.

Finally, in Fig. 12 we show qualitative examples to compare the effect using our proposed interpolation method (Eq. 11) versus simple linear interpolation in the multimodal latent space of our *mAR* prior. We use the Adamax optimizer to optimize Eq. (11) with a learning rate of 5×10^{-2} for ~ 100 iterations. The computational requirement per iteration is approximately equivalent to a standard training iteration of our *mAR-SCF* model, *i.e.* ~ 1 sec for a batch of size 128. We observe that interpolated images obtained from our *mAR-SCF* affine model using our proposed interpolation method (Eq. 11) have a better visual quality, especially for the interpolations in the middle of the interpolating path. Note that we find our scheme to be stable in practice in a reasonably wide range of the hyperparameters $\lambda_2 \approx \lambda_1 \in [0.2, 0.5]$ as it interpolates in high-density regions between $\mathbf{x}_A, \mathbf{x}_B$. We support the better visual quality of the interpolations of our scheme compared to the linear method with Inception scores computed for interpolations obtained from both methods in Fig. 8.

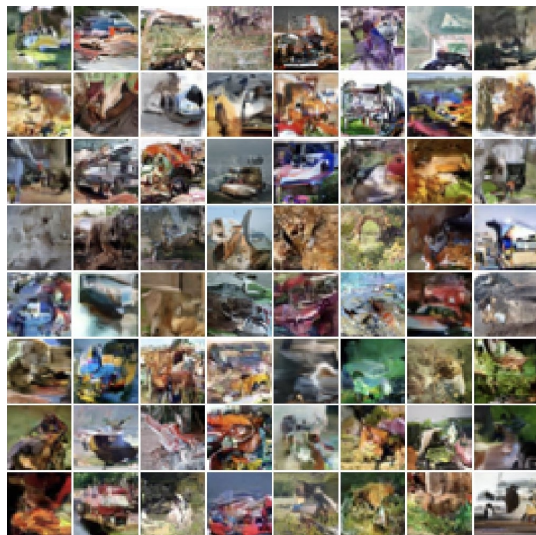


(a) Residual Flows [4]



(b) Our *mAR-SCF* (MixLogCDF)

Figure 9. Random samples when trained on MNIST.

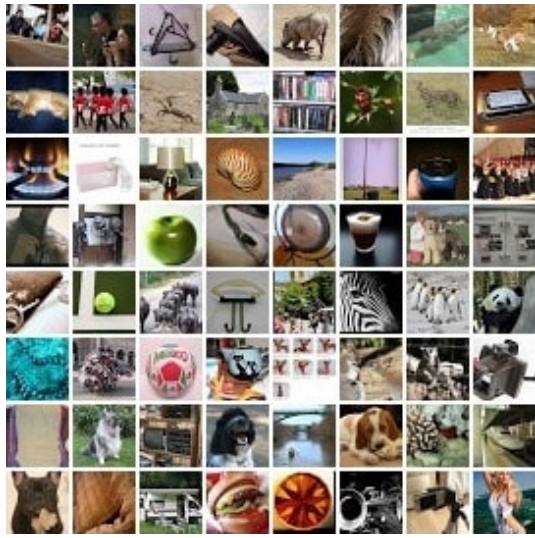


(a) PixelCNN [38]

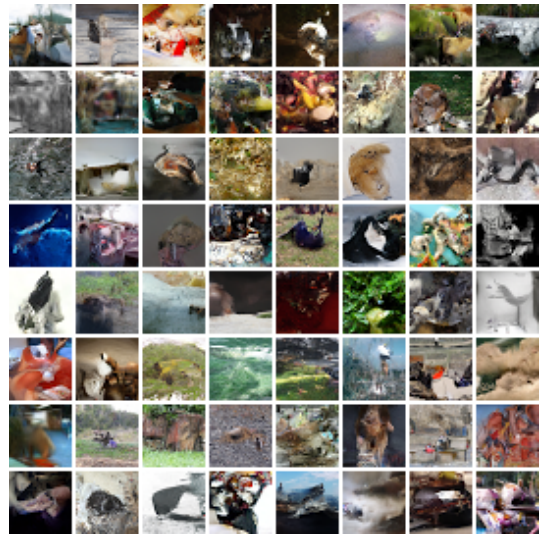


(b) Our *mAR-SCF* (MixLogCDF)

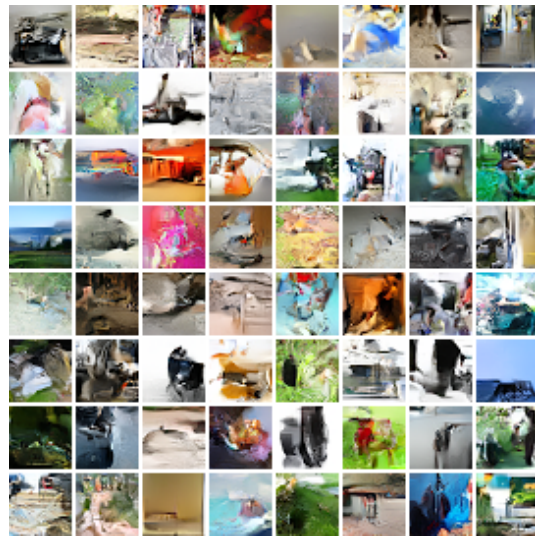
Figure 10. Random samples when trained on CIFAR10 (32×32).



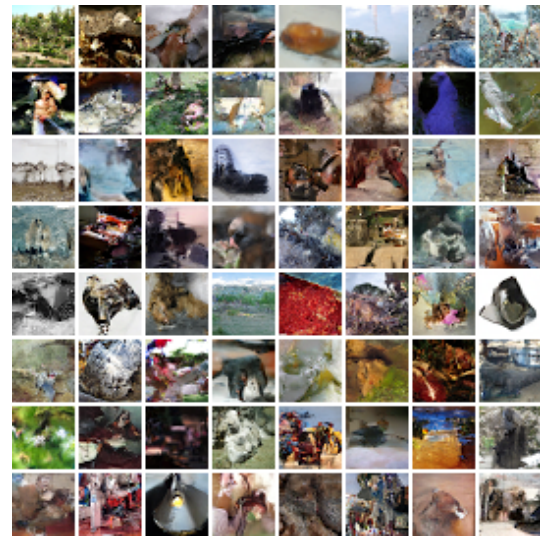
(a) Real Data



(b) Flow++ [14]



(c) Residual Flows [4]



(d) Our *mAR-SCF* (MixLogCDF)

Figure 11. Random samples when trained on ImageNet (32×32).

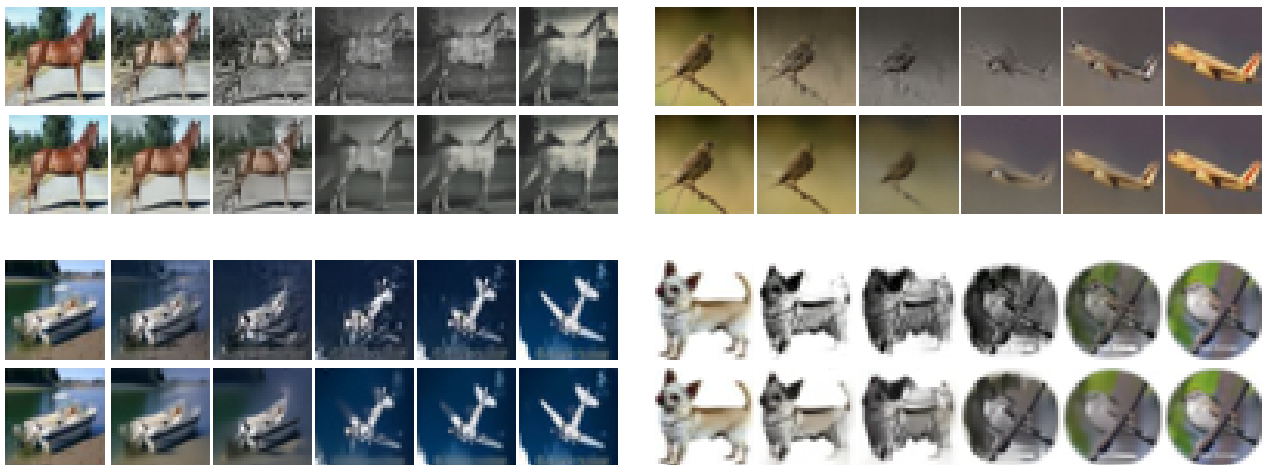


Figure 12. Comparison of interpolation quality of our interpolation scheme (Eq. 11) with a standard linear interpolation scheme. The top row of each result shows the interpolations between real images for the linear method and the corresponding bottom rows are the interpolations with our proposed scheme.