

---

# A Multi-Resolution Framework for U-Nets with Applications to Hierarchical VAEs

---

Fabian Falck <sup>\*1,3,4</sup> Christopher Williams <sup>\*,1</sup> Dominic Danks <sup>2,4</sup> George Deligiannidis <sup>1</sup>  
Christopher Yau <sup>1,3,4</sup> Chris Holmes <sup>1,3,4</sup> Arnaud Doucet <sup>1</sup> Matthew Willetts <sup>4</sup>

<sup>1</sup>University of Oxford <sup>2</sup>University of Birmingham

<sup>3</sup>Health Data Research UK <sup>4</sup>The Alan Turing Institute

{fabian.falck, williams, deligian, cholmes, doucet}@stats.ox.ac.uk,  
{ddanks, cyau, mwilletts}@turing.ac.uk}

## Abstract

U-Net architectures are ubiquitous in state-of-the-art deep learning, however their regularisation properties and relationship to wavelets are understudied. In this paper, we formulate a multi-resolution framework which identifies U-Nets as finite-dimensional truncations of models on an infinite-dimensional function space. We provide theoretical results which prove that average pooling corresponds to projection within the space of square-integrable functions and show that U-Nets with average pooling implicitly learn a Haar wavelet basis representation of the data. We then leverage our framework to identify state-of-the-art hierarchical VAEs (HVAEs), which have a U-Net architecture, as a type of two-step forward Euler discretisation of multi-resolution diffusion processes which flow from a point mass, introducing sampling instabilities. We also demonstrate that HVAEs learn a representation of time which allows for improved parameter efficiency through weight-sharing. We use this observation to achieve state-of-the-art HVAE performance with half the number of parameters of existing models, exploiting the properties of our continuous-time formulation.

## 1 Introduction

U-Net architectures are extensively utilised in modern deep learning models. First developed for image segmentation in biomedical applications [1], U-Nets have been widely applied for text-to-image models [2], image-to-image translation [3], image restoration [4, 5], super-resolution [6], and multiview learning [7], amongst other tasks [8]. They also form a core building block as the neural architecture of choice in state-of-the-art generative models, particularly for images, such as HVAEs [9, 10, 11, 12] and diffusion models [2, 13, 14, 15, 16, 17, 18, 19, 20]. In spite of their empirical success, it is poorly understood why U-Nets work so well, and what regularisation they impose.

In likelihood-based generative modelling, various model classes are competing for superiority, including normalizing flows [21, 22], autoregressive models [23, 24], diffusion models, and hierarchical variational autoencoders (HVAEs), the latter two of which we focus on in this work. HVAEs form groups of latent variables with a conditional dependence structure, use a U-Net neural architecture, and are trained with the typical VAE ELBO objective (for a detailed introduction to HVAEs, see Appendix B). HVAEs show impressive synthesis results on facial images, and yield competitive likelihood performance, consistently outperforming the previously state-of-the-art autoregressive models, VAEs and flow models on computer vision benchmarks [9, 10]. HVAEs have undergone a journey

---

\*Equal contribution.

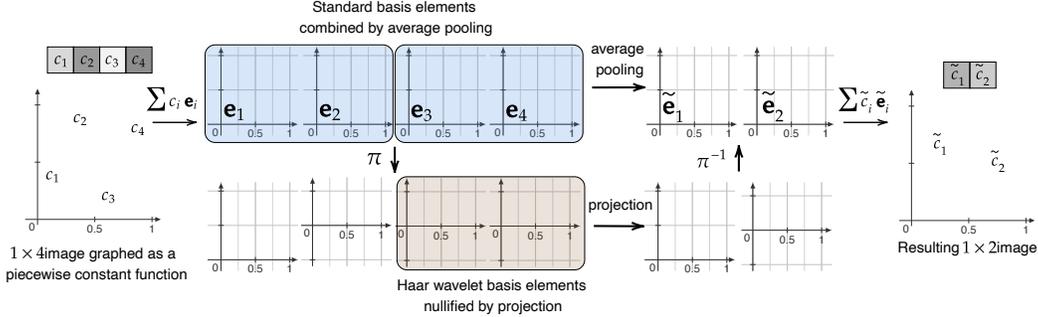


Figure 1: U-Nets with average pooling learn a Haar wavelet basis representation of the data.

of design iterations and architectural improvements in recent years, for example the introduction a deterministic backbone [25, 26, 27] and ResNet elements [28, 29] with shared parameters between the inference and generative model parts. There has also been a massive increase in the number of latent variables and overall stochastic depth, as well as the use of different types of residual cells in the decoder [9, 10] (see §4 and Fig. A.1 for a detailed discussion). However, a theoretical understanding of these choices is lacking. For instance, it has not been shown why a residual backbone may be beneficial, or what the specific cell structures in VDVAE [9] and NVAE [10] correspond to, or how they could be improved.

In this paper we provide a theoretical framework for understanding the latent spaces in U-Nets, and apply this to HVAEs specifically. Doing so allows us to relate HVAEs to diffusion processes, and also to motivate a new type of piecewise time-homogenous model which demonstrates state-of-the-art performance with approximately half the number of parameters of a VDVAE [9]. More formally, our contributions are as follows: **(a)** We provide a multi-resolution framework for U-Nets. We formally define U-Nets as acting over a multi-resolution hierarchy of  $L^2([0, 1]^2)$ . We prove that average pooling is a conjugate operation to projection in the Haar wavelet basis within  $L^2([0, 1]^2)$ . We use this insight to show how U-Nets with average pooling implicitly learn a Haar wavelet basis representation of the data (see Fig. 1), helping to characterise the regularisation within U-Nets. **(b)** We apply this framework to state-of-the-art HVAEs as an example, identifying their residual cell structure as a type of two-step forward Euler discretisation of a multi-resolution diffusion bridge. We uncover that this diffusion process flows from a point mass, which causes instabilities, for instance during sampling, and identify parameter redundancies through our continuous-time formulation. Our framework both allows us to understand the heuristic choices of existing work in HVAEs and enables future work to optimise their design, for instance their residual cell. **(c)** In our experiments, we demonstrate these sampling instabilities and train HVAEs with the largest stochastic depth ever, achieving state-of-the-art performance with half the number of parameters by exploiting our theoretical insights. We explain these results by uncovering that HVAEs secretly represent time in their state and show that they use this information during training. We finally provide extensive ablation studies which, for instance, rule out other potential factors which correlate with stochastic depth, show the empirical gain of multiple resolutions, and find that Fourier features (which discrete-time diffusion models strongly benefit from [19]) do not improve performance in the HVAE setting.

## 2 The Multi-Resolution Framework

A grayscale image with infinite resolution can be thought of as the graph<sup>2</sup> of a two-dimensional function over the unit square. To store these infinitely-detailed images in computers, we project them to some finite resolution. These projections can still be thought of as the graphs of functions with support over the unit square, but they are piecewise constant on finitely many intervals or ‘pixels’, e.g.  $512^2$  pixels, and we store the function values obtained at these pixels in an array or ‘grid’. The relationship between the finite-dimensional version and its infinitely-fine counterpart depends entirely on how we construct this projection to preserve the details we wish to keep. One approach is to prioritise preserving the large-scale details of our images, so unless closely inspected, the projection

<sup>2</sup>For a function  $f(\cdot)$ , its graph is the set  $\bigcup_{x \in [0, 1]^2} \{x, f(x)\}$ .

is indistinguishable from the original. This can be achieved with a multi-resolution projection [30] of the image. In this section we introduce a *multi-resolution framework* for constructing neural network architectures that utilise such projections, prove what regularisation properties they impose, and show as an example how HVAEs with a U-Net [1] architecture can be interpreted in our framework. Proofs of all theorems in the form of an extended exposition of our framework can be found in Appendix A.

## 2.1 Multi-Resolution Framework: Definitions and Intuition

What makes a multi-resolution projection good at prioritising large-scale details can be informally explained through the following thought experiment. Imagine we have an image, represented as the graph of a function, and its finite-dimensional projection drawn on the wall. We look at the wall, start walking away from it and stop when the image and its projection are indistinguishable by eye. The number of steps we took away from the wall can be considered our measure of ‘how far away’ the approximation is from the underlying function. The goal of the multi-resolution projection is therefore to have to take as few steps away as possible. The reader is encouraged to physically conduct this experiment with the examples provided in Appendix B.1. We can formalise the aforementioned intuitions by defining a *multi-resolution hierarchy* [30] of sub-spaces we may project to:

**Definition 1.** [Daubechies (1992) [30]] Given a nested sequence of *approximation spaces*  $\dots \subset V_1 \subset V_0 \subset V_{-1} \subset \dots$ ,  $\{V_{-j}\}_{j \in \mathbb{Z}}$  is a *multi-resolution hierarchy* of the function space  $L^2(\mathbb{R}^m)$  if: **(A1)**  $\bigcup_{j \in \mathbb{Z}} V_{-j} = L^2(\mathbb{R}^m)$ ; **(A2)**  $\bigcap_{j \in \mathbb{Z}} V_{-j} = \{0\}$ ; **(A3)**  $f(\cdot) \in V_{-j} \Leftrightarrow f(2^j \cdot) \in V_0$ ; **(A4)**  $f(\cdot) \in V_0 \Leftrightarrow f(\cdot - n) \in V_0$  for  $n \in \mathbb{Z}$ . For a compact set  $\mathbb{X} \subset \mathbb{R}^m$ , a *multi-resolution hierarchy* of  $L^2(\mathbb{X})$  is  $\{V_{-j}\}_{j \in \mathbb{Z}}$  as defined above, restricting functions in  $V_{-j}$  to be supported on  $\mathbb{X}$ .

In Definition 1, the index  $j$  references how many steps we took in our thought experiment, so negative  $j$  corresponds to ‘zooming in’ on the images. The original image<sup>3</sup> is a member of  $L^2([0, 1]^2)$ , the space of square-integrable functions on the unit square, and its finite projection to  $2^j \cdot 2^j$  many pixels is a member of  $V_{-j}$ . Images can be represented as piecewise continuous functions in the subspaces  $V_{-j} = \{f \in L^2([0, 1]) \mid f|_{[2^{-j} \cdot k, 2^{-j} \cdot (k+1))} = c_k, k \in \{0, \dots, 2^j - 1\}, c_k \in \mathbb{R}\}$ . The nesting property  $V_{-j+1} \subset V_{-j}$  ensures that any image with  $(2^{j-1})^2$  pixels can also be represented by  $(2^j)^2$  pixels, but at a higher resolution. Assumption **(A1)** states that with infinitely many pixels, we can describe any infinitely detailed image. In contrast, **(A2)** says that with no pixels, we cannot approximate any images. Assumptions **(A3)** and **(A4)** allow us to form a basis for images in any  $V_{-j}$  if we know the basis of  $V_0$ . One basis made by extrapolating from  $V_0$  in this way is known as a *wavelet basis* [30]. Wavelets have proven useful for representing images, for instance in the JPEG standard [31], and are constructed to be orthonormal.

Now suppose we have a probability measure  $\nu_\infty$  over infinitely detailed images represented in  $L^2([0, 1]^2)$  and wish to represent it at a lower resolution. Similar to how we did for infinitely detailed images, we want to project the measure  $\nu_\infty$  to a lower dimensional measure  $\nu_j$  on the finite dimensional space  $V_{-j}$ . In extension to this, we want the ability to reverse this projection so that we may sample from the lower dimensional measure and create a generative model for  $\nu_\infty$ . We would like to again prioritise the presence of large-scale features of the original image within the lower dimensional samples. We do this by constructing a *multi-resolution bridge* from  $\nu_\infty$  to  $\nu_j$ , as defined below.

**Definition 2.** Let  $\mathbb{X} \subset \mathbb{R}^m$  be compact,  $\{V_{-j}\}_{j=0}^\infty$  be a multi-resolution hierarchy of scaled so  $L^2(\mathbb{X}) = \bigcup_{j \in \mathbb{N}_0} V_{-j}$  and  $V_0 = \{0\}$ . If  $\mathbb{D}(L^2(\mathbb{X}))$  is the space of probability measures over  $L^2(\mathbb{X})$ , then a family of probability measures  $\{\nu_t\}_{t \in [0, 1]}$  on  $L^2(\mathbb{X})$  is a *multi-resolution bridge* if:

- (i) there exist increasing times  $\mathcal{I} := \{t_j\}_{j \in \mathbb{N}_0}$  where  $t_0 = 0$ ,  $\lim_{j \rightarrow \infty} t_j = 1$ , such that  $s \in [t_j, t_{j+1})$  implies  $\text{supp}(\nu_s) \subset V_{-j}$ , i.e  $\nu_s \in \mathbb{D}(V_{-j})$ ; and,
- (ii) for  $s \in (0, 1)$ , the mapping  $s \mapsto \nu_s$  is continuous for  $s \in (t_j, t_{j+1})$  for some  $j$ .

The continuous time dependence in Definition 2 plays a movie of the measure  $\nu_0$  supported on  $V_0$  growing to  $\nu_\infty$ , a measure on images with infinite resolution. At a time interval  $[t_j, t_{j+1})$ , the space  $V_{-j}$  which the measure is supported on is fixed. We may therefore define a finite-dimensional model

<sup>3</sup>We here focus on grayscale, squared images for simplicity, but note that our framework can be seamlessly extended to colour images with a Cartesian product  $L^2([0, 1]^2) \times L^2([0, 1]^2) \times L^2([0, 1]^2)$ , and other continuous signals such as time series.

transporting probability measures within  $V_{-j}$ , but at  $t_{j+1}$  the support flows over to  $V_{-j-1}$ . Given a multi-resolution hierarchy, we may glue these finite models, each acting on a disjoint time interval, together in a unified fashion. In Theorem 1 we show this for the example of a continuous-time multi-resolution diffusion process truncated up until some time  $t_J = T \in (0, 1)$  and in the *standard basis* discussed in §2.2, which will be useful when viewing HVAEs as discretisations of diffusion processes on functions in §2.3.

**Theorem 1.** *Let  $B_j : [t_j, t_{j+1}] \times \mathbb{D}(V_{-j}) \mapsto \mathbb{D}(V_{-j})$  be a linear operator (such as a diffusion transition kernel, see Appendix A) for  $j < J$  with coefficients  $\mu^{(j)}, \sigma^{(j)} : [t_j, t_{j+1}] \times V_{-j} \mapsto V_{-j}$ , and define the natural extensions within  $V_{-J}$  in bold, i.e.  $\mathbf{B}_j := B_j \oplus \mathbf{I}_{V_{-j}^\perp}$ . Then the operator  $\mathbf{B} : [0, T] \times \mathbb{D}(V_{-J}) \mapsto \mathbb{D}(V_{-J})$  and the coefficients  $\boldsymbol{\mu}, \boldsymbol{\sigma} : [0, T] \times V_{-J} \mapsto V_{-J}$  given by*

$$\mathbf{B} := \sum_{j=0}^J \mathbb{1}_{[t_j, t_{j+1})} \cdot \mathbf{B}_j, \quad \boldsymbol{\mu} := \sum_{j=0}^J \mathbb{1}_{[t_j, t_{j+1})} \cdot \boldsymbol{\mu}^{(j)}, \quad \boldsymbol{\sigma} := \sum_{j=0}^J \mathbb{1}_{[t_j, t_{j+1})} \cdot \boldsymbol{\sigma}^{(j)},$$

*induce a multi-resolution bridge of measures from the dynamics for  $t \in [0, T]$  and on the standard basis as  $dZ_t = \boldsymbol{\mu}_t(Z_t)dt + \boldsymbol{\sigma}_t(Z_t)dW_t$  (see Appendix A.4 for details) for  $Z_t \in V_{-j}$  for  $t \in [t_j, t_{j+1})$ , i.e. a multi-resolution diffusion process.*

The concept of a multi-resolution bridge will become important in Section 2.2 where we will show that current U-Net bottleneck structures used for unconditional sampling impose a multi-resolution bridge on the modelled densities. To preface this, we here provide a description of a U-Net within our framework, illustrated in 2. Consider  $B_{j,\theta}, F_{j,\theta} : \mathbb{D}(V_{-j}) \rightarrow \mathbb{D}(V_{-j})$  as the forwards and backwards passes of a U-Net on resolution  $j$ . Further, let  $P_{-j+1} : \mathbb{D}(V_{-j}) \rightarrow \mathbb{D}(V_{-j+1})$  and  $E_{-j} : \mathbb{D}(V_{-j+1}) \rightarrow \mathbb{D}(V_{-j})$  be the projection (here: average pooling) and embedding maps (e.g. interpolation), respectively. When using an  $L^2$ -reconstruction error, a U-Net [1] architecture implicitly learns a sequence of models  $\mathbf{B}_{j,\phi} : \mathbb{D}(V_{-j+1}) \times \mathbb{D}(V_{-j+1}^\perp) \mapsto \mathbb{D}(V_{-j})$  due to the orthogonal decomposition  $V_{-j} = V_{-j+1} \oplus U_{-j+1}$  where  $U_{-j+1} := V_{-j} \cap V_{-j+1}^\perp$ . The backwards operator for the U-Net has a (bottleneck) input from  $\mathbb{D}(V_{-j+1})$  and a (skip) input yielding information from  $\mathbb{D}(V_{-j+1}^\perp)$ . A simple *bottleneck* map  $U_{j,\theta} : \mathbb{D}(V_{-j}) \rightarrow \mathbb{D}(V_{-j})$  (without skip connection) is given by

$$U_{j,\theta} := B_{j,\theta} \circ E_{-j} \circ P_{-j+1} \circ F_{j,\theta}, \quad (1)$$

and a U-Net bottleneck with skip connection is

$$\mathbf{U}_{j,\phi} := B_{j,\phi}(E_{-j} \circ P_{-j+1} \circ F_{j,\theta}, F_{j,\theta}). \quad (2)$$

In HVAEs, the map  $\mathbf{U}_{j,\phi} : \mathbb{D}(V_{-j}) \rightarrow \mathbb{D}(V_{-j})$  is trained to be the identity by minimising reconstruction error, and further shall approximate  $U_{j,\theta} \approx \mathbf{U}_{j,\phi}$  via a KL divergence. The  $L^2$ -reconstruction error for  $\mathbf{U}_{j,\phi}$  has an orthogonal partition of the inputs from  $V_{-j+1} \times V_{-j}$ , hence the only new subspace added is  $U_{-j+1}$ . As each orthogonal  $U_{-j+1}$  is added sequentially in HVAEs, the skip connections induce a multi-resolution structure of this hierarchical neural network structure. What we will investigate in Theorem 3 is the regularisation imposed on this partitioning by enforcing  $U_{j,\theta} \approx \mathbf{U}_{j,\phi}$ , as is often enforced for generative models with VAEs.

## 2.2 The regularisation property imposed by U-Net architectures with average pooling

Having defined U-Net architectures within our multi-resolution framework, we are now interested in the regularisation they impose. We do so by analysing a U-Net when skip connections are absent, so that we may better understand what information is transferred through each skip connection when they are present. In practice, a pixel representation of images is used when training U-Nets, which we henceforth call the *standard basis* (see A.2, Eq. (A.9)). The standard basis is not convenient to derive theoretical results. It is instead preferable to use a basis natural to the multi-resolution bridge imposed by a U-Net with a corresponding projection operation, which for average pooling is the *Haar (wavelet) basis* [32] (see Appendix A.2). The Haar basis, like a Fourier basis, is an orthonormal basis

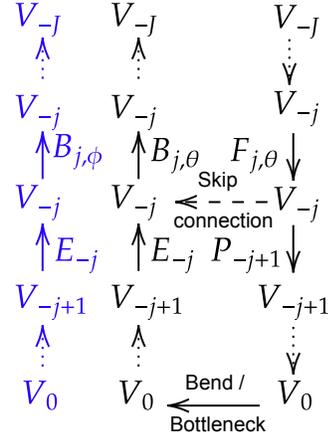


Figure 2: A U-Net in our multi-resolution framework. See Appendix B.2 for details.

of  $L^2(\mathbb{X})$  which has desirable  $L^2$ -approximation properties. We formalise this in Theorem 2 which states that the dimension reduction operation of average pooling in the standard basis is a conjugate operation to co-ordinate projection within the Haar basis (details are provided in Appendix A.2).

**Theorem 2.** *Given  $V_{-j}$  as in Definition 1, let  $x \in V_{-j}$  be represented in the standard basis  $\mathbf{E}_j$  and Haar basis  $\mathbf{\Psi}_j$ . Let  $\pi_j : \mathbf{E}_j \mapsto \mathbf{\Psi}_j$  be the change of basis map illustrated in Fig. 3, then we have the conjugacy  $\pi_{j-1} \circ \text{pool}_{-j,-j+1} = \text{proj}_{V_{-j+1}} \circ \pi_j$ .*

Theorem 2 means that if we project an image from  $V_{-j}$  to  $V_{-j+1}$  in the Haar wavelet basis, we can alternatively view this as changing to the standard basis via  $\pi_j^{-1}$ , performing average pooling, and reverting back via  $\pi_{j-1}$  (see Figure 3). This is important because the Haar basis is orthonormal, which in Theorem 3 allows us to precisely quantify what information is lost with average pooling.

**Theorem 3.** *Let  $\{V_{-j}\}_{j=0}^J$  be a multi-resolution hierarchy of  $V_{-J}$  where  $V_{-j} = V_{-j+1} \oplus U_{-j+1}$ , and further, let  $F_{j,\phi}, B_{j,\theta} : \mathbb{D}(V_{-j}) \mapsto \mathbb{D}(V_{-j})$  be such that  $B_{j,\theta}F_{j,\phi} = I$  with parameters  $\phi$  and  $\theta$ . Define  $\mathbf{F}_{j_1|j_2,\phi} := \mathbf{F}_{j_1,\phi} \circ \dots \circ \mathbf{F}_{j_2,\phi}$  by  $\mathbf{F}_{j,\phi} : \mathbb{D}(V_{-j}) \mapsto \mathbb{D}(V_{-j+1})$  where  $\mathbf{F}_{j,\phi} := \text{proj}_{V_{-j+1}} \circ F_{j,\phi}$ , and analogously define  $\mathbf{B}_{j_1|j_2,\theta}$  with  $\mathbf{B}_{j,\theta} := B_{j,\theta} \circ \text{embd}_{V_{-j}}$ . Then, the sequence  $\{\mathbf{B}_{1|j,\theta}(\mathbf{F}_{1|J,\phi}\nu_J)\}_{j=0}^J$  forms a discrete multi-resolution bridge between  $\mathbf{F}_{1|J,\phi}\nu_J$  and  $\mathbf{B}_{1|J,\theta}\mathbf{F}_{1|J,\phi}\nu_J$  at times  $\{t_j\}_{j=1}^J$ , and*

$$\sum_{j=0}^J \mathbb{E}_{X_{t_j} \sim \nu_j} \left\| \text{proj}_{U_{-j+1}} X_{t_j} \right\|_2^2 / \left\| \mathbf{F}_{j|J,\phi} \right\|_2^2 \leq (\mathcal{W}_2(\mathbf{B}_{1|J,\theta}\mathbf{F}_{1|J,\phi}\nu_J))^2, \quad (3)$$

where  $\mathcal{W}_2$  is the Wasserstein-2 metric and  $\left\| \mathbf{F}_{j|J,\phi} \right\|_2$  is the Lipschitz constant of  $\mathbf{F}_{j|J,\phi}$ .

Theorem 3 states that the bottleneck component of a U-Net pushes the latent data distribution to a finite multi-resolution basis, specifically a Haar basis when average pooling is used. To see this, note that the RHS of Eq. (A.65) is itself upper-bounded by the  $L^2$ -reconstruction error. This is because the Wasserstein-2 distance finds the infimum over all possible couplings between the data and the ‘reconstruction’ measure, hence any coupling (induced by the learned model) bounds it. Note that models using a U-Net, for instance HVAEs or diffusion models, either directly or indirectly optimise for low reconstruction error in their loss function. The LHS of Eq. (A.65) represents what percentage of our data enters the orthogonal subspaces  $\{U_{-j}\}_{j=0}^J$  which are (by Theorem 2) discarded by the bottleneck structure when using a U-Net architecture with average pooling. Theorem 3 thus shows that as we minimise the reconstruction error during training, we minimise the percentage of our data transported to the orthogonal sub-spaces  $\{U_{-j}\}_{j=0}^J$ . Consequently, the bottleneck architecture implicitly decomposes our data into a Haar wavelet decomposition, and when the skip connections are absent (like in a traditional auto-encoder) our network learns to compress the discarded subspaces  $U_{-j}$ . This characterises the regularisation imposed by a U-Net in the absence of skip connections.

These results suggest that U-Nets with average pooling provide a direct alternative to Fourier features [19, 33, 34, 35] which impose a Fourier basis, an alternative orthogonal basis on  $L^2(\mathbb{X})$ , as with skip connections the U-Net adds each subspace  $U_{-j}$  sequentially. However, unlike Fourier bases, there are in fact a multitude of wavelet bases which are all encompassed by the multi-resolution framework, and in particular, Theorem 3 pertains to all of them for the bottleneck structure. This opens the door to exploring conjugacy operations beyond average pooling induced by other wavelet bases optimised for specific data types.

### 2.3 Example: HVAEs as Diffusion Discretisations

To show what practical inferences we can derive from our multi-resolution framework, we apply it to analyse state-of-the-art HVAE architectures (see Appendix B.3 for an introduction), identifying parameter redundancies and instabilities. Here and in our experiments, we focus on VDVAEs [9]. We provide similar results for Markovian HVAEs [36, 37] and NVAEs [10] (see § 4) in Appendix A.5.

We start by inspecting VDVAEs. As we show next, we can tie the computations in VDVAE cells to the (forward and backward) operators  $F_{j,\phi}$  and  $B_{j,\theta}$  within our framework and identify them as a type of two-step forward Euler discretisation of a diffusion process. When used with a U-Net, as is done in VDVAE [9], this creates a *multi-resolution diffusion bridge* by Theorem 4.

$$\begin{array}{ccc} (V_{-j}, \mathbf{E}_j) & \xrightarrow{\text{pool}_{-j,-j+1}} & (V_{-j+1}, \mathbf{E}_{j-1}) \\ \uparrow \pi_j^{-1} & & \pi_{j-1} \downarrow \\ (V_{-j}, \mathbf{\Psi}_j) & \xrightarrow{\text{proj}_{V_{-j+1}}} & (V_{-j+1}, \mathbf{\Psi}_{j-1}) \end{array}$$

Figure 3: The function space  $V_{-j}$  remains the same, but the basis changes under  $\pi_j$ .

**Theorem 4.** Let  $t_j := T \in (0, 1)$  and consider (the  $p_\theta$  backward pass)  $\mathbf{B}_{\theta,1|J} : \mathbb{D}(V_{-j}) \mapsto \mathbb{D}(V_0)$  given in multi-resolution Markov process in the standard basis:

$$dZ_t = (\overleftarrow{\mu}_{1,t}(Z_t) + \overleftarrow{\mu}_{2,t}(Z_t))dt + \overleftarrow{\sigma}_t(Z_t)dW_t, \quad (4)$$

where  $\text{proj}_{U_{-j}} Z_{t_j} = 0$ ,  $\|Z_t\|_2 > \|Z_s\|_2$  with  $0 \leq s < t \leq T$  and for a measure  $\nu_j \in \mathbb{D}(V_{-j})$  we have  $X_T, Z_0 \sim \mathbf{F}_{\phi,J|1} \nu_j = \delta_{\{0\}}$ . Then, VDVAEs approximates this process, and its residual cells are a type of two-step forward Euler discretisation of this Stochastic Differential Equation (SDE).

To better understand Theorem 4, we visualise its residual cell structure of VDVAEs and the corresponding discretisation steps in Fig. 4, and together those of NVAEs and Markovian HVAEs in Appendix A.5, Fig. A.1. Note that this process is Markov and increasing in the  $Z_i$  variables. Similar processes have been empirically observed as efficient first-order approximates to higher-order chains, for example the memory state in LSTMs [38]. Further, VDVAEs and NVAEs are even claimed to be high-order chains (see Eqs. (2,3) in [9] and Eq. (1) in [10]), despite only approximating this with a accumulative process.

To show how VDVAEs impose the growth of the  $Z_t$ , we prove that the bottleneck component of VDVAE’s U-Net enforces  $Z_0 = 0$ . This is done by identifying that the measure  $\nu_0$ , which a VDVAE connects to the data  $\nu_\infty$  via a multi-resolution bridge, is a point mass on the zero function. Consequently the backward pass must grow from this, and the network learns this in a monotonic manner as we later confirm in our experiments (see §3.2).

**Theorem 5.** Consider the SDE in Eq. (A.76), trained through the ELBO in Eq. B.101. Let  $\tilde{\nu}_j$  denote the data measure and  $\nu_0 = \delta_{\{0\}}$  be the initial multi-resolution bridge measure imposed by VDVAEs. If  $q_{\phi,j}$  and  $p_{\theta,j}$  are the densities of  $\mathbf{B}_{\phi,1|j} \mathbf{F}_{J|1} \tilde{\nu}_j$  and  $\mathbf{B}_{\theta,1|j} \nu_0$  respectively, then a VDVAE optimises the boundary condition  $\min_{\theta,\phi} KL(q_{\phi,0,1} || q_{\phi,0} p_{\theta,1})$ , where a double index indicates the joint distribution.

Theorem 5 states that the VDVAE architecture forms multi-resolution bridge with the dynamics of Eq. (A.76), and connects our data distribution to the trivial measure on  $V_0$ : a Dirac mass at 0 as the pooling here cascades completely to  $V_0$ . From this insight, we can draw conclusions on instabilities and on parameter redundancies of this HVAE cell. There are two major instabilities in this discretisation. First, the imposed  $\nu_0$  is disastrously unstable as it enforces a data set, with potentially complicated topology to derive from a point-mass in  $U_{-j}$  at each  $t = t_j$ , and we observe the resulting sampling instability in our experiments in §3.3. We note that similar arguments are applicable in settings without a latent hierarchy imposed by a U-Net, see for instance [39]. The VDVAE architecture does, however, bolster this rate through the  $Z_{i,+}^{(\sigma)}$  term, which is absent in NVAEs [10], in the discretisation steps of the residual cell. We empirically observe this controlled backward error in Fig. 6 [Right]. We refer to Fig. A.1 for a detailed comparison of HVAE cells and their corresponding discretisation of the coupled SDE in Eq. (A.76).

Moreover, the current form of VDVAEs is over-parameterised and not informed by this continuous-time formulation. The continuous time analogue of VDVAEs [9] in Theorem 4 has time dependent coefficients  $\overleftarrow{\mu}_{t,1}$ ,  $\overleftarrow{\mu}_{t,2}$ ,  $\overleftarrow{\sigma}_t$ . We hypothesise that the increasing diffusion process in  $Z_i$  implicitly encodes time. Hence, explicitly representing this in the model, for instance via ResNet blocks with independent parameterisations at every time step, is redundant, and a time-homogeneous model (see Appendix A.6 for a precise formulation)—practically speaking, performing weight-sharing across time steps/layers—has the same expressivity, but requires far fewer parameters than the state-of-the-art VDVAE. It is worth noting that such a time-homogeneous model would make the parameterisation of HVAEs more similar to the recently popular (score-based) diffusion models [40, 41] which perform weight-sharing across all time steps.

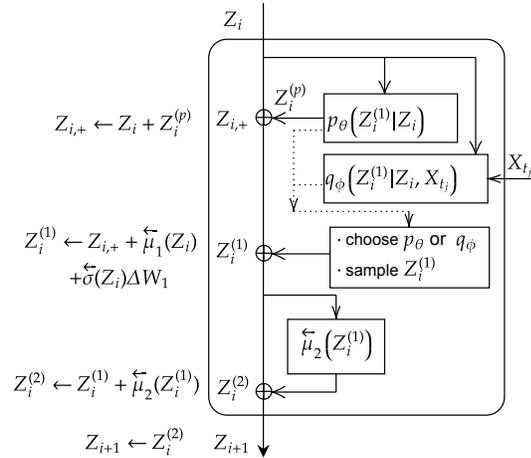


Figure 4: The VDVAE [9] cell is a type of two-step forward Euler discretisations of the continuous-time diffusion process in Eq. A.76. See Fig. A.1 for similar schemas on NVAE [10] and Markovian HVAE [36, 37].

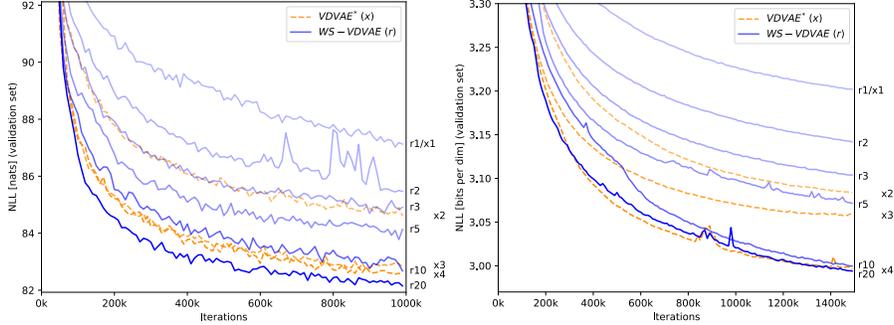


Figure 5: A small-scale study on parameter efficiency of HVAEs. We compare models with with 1,2,3 and 4 parameterised blocks per resolution ( $\{x1, x2, x3, x4\}$ ) against models with a single parameterised block per resolution weight-shared  $\{2, 3, 5, 10, 20\}$  times ( $\{r2, r3, r5, r10, r20\}$ ). We report NLL ( $\downarrow$ ) measured on the validation set of MNIST [left] and CIFAR10 [right]. NLL performance increases with more weight-sharing repetitions and surpasses models without weight-sharing but with more parameters.

### 3 Experiments

In the following we probe the theoretical understanding of HVAEs gained through our framework, demonstrating its utility in four experimental analyses: (a) Improving parameter efficiency in HVAEs, (b) Time representation in HVAEs and how they make use of it, (c) Sampling instabilities in HVAEs, and (d) Ablation studies.

We train HVAEs using VDVAE [9] as the basis model on five datasets: MNIST [42], CIFAR10 [43], two downsampled versions of ImageNet [44, 45], and CelebA [46], splitting each into a training, validation and test set (see Appendix D for details). In general, reported numeric values refer to Negative Log-Likelihood (NLL) in nats (MNIST) or bits per dim (all other datasets) on the test set at model convergence, if not stated otherwise. We note that performance on the validation and test set have similar trends in general. An optional *gradient checkpointing* implementation to trade in GPU memory for compute is discussed in Appendix F. Appendices F and G define the HVAE models we train, i.e.  $p_\theta(\mathbf{z}_L)$ ,  $p_\theta(\mathbf{z}_l|\mathbf{z}_{>l})$ ,  $q_\phi(\mathbf{z}_L|\mathbf{x})$ ,  $q_\phi(\mathbf{z}_l|\mathbf{z}_{>l}, \mathbf{x})$  and  $p_\theta(\mathbf{x}|\bar{\mathbf{z}})$ , and present additional experimental details and results. We provide our PyTorch code base at <https://github.com/FabianFalck/unet-vdvae> (see Appendix C for details).

Table 1: A large-scale study of parameter efficiency in HVAEs. We compare our runs of VDVAE with original hyperparameters [9] (VDVAE\*) against our weight-shared VDVAE (WS-VDVAE). While WS-VDVAEs have improved parameter efficiency by a factor of 2, they reach similar NLL as VDVAE\* with the simple modification inspired by our framework (weight sharing). We note that a parameter count cannot be provided for VDM [19] as the code is not public and the manuscript does not specify it.

Dataset	Method	Type	#Params	NLL $\downarrow$
<b>MNIST</b> 28 × 28	WS-VDVAE (ours)	VAE	<b>232k</b>	$\leq 79.98$
	VDVAE* (ours)	VAE	339k	$\leq 80.14$
	NVAE [10]	VAE	33m	$\leq 78.01$
<b>CIFAR10</b> 32 × 32	WS-VDVAE (ours)	VAE	<b>25m</b>	$\leq 2.88$
	WS-VDVAE (ours)	VAE	39m	$\leq 2.83$
	VDVAE* (ours)	VAE	39m	$\leq 2.87$
	NVAE [10]	VAE	131m	$\leq 2.91$
	VDVAE [9]	VAE	39m	$\leq 2.87$
	VDM [19]	Diff	–	$\leq 2.65$
<b>ImageNet</b> 32 × 32	WS-VDVAE (ours)	VAE	<b>55m</b>	$\leq 3.68$
	WS-VDVAE (ours)	VAE	85m	$\leq 3.65$
	VDVAE* (ours)	VAE	119m	$\leq 3.67$
	NVAE [10]	VAE	268m	$\leq 3.92$
	VDVAE [9]	VAE	119m	$\leq 3.80$
	VDM [19]	Diff	–	$\leq 3.72$
<b>CelebA</b> 64 × 64	WS-VDVAE (ours)	VAE	<b>75m</b>	$\leq 2.02$
	VDVAE* (ours)	VAE	125m	$\leq 2.02$
	NVAE [10]	VAE	153m	$\leq 2.03$

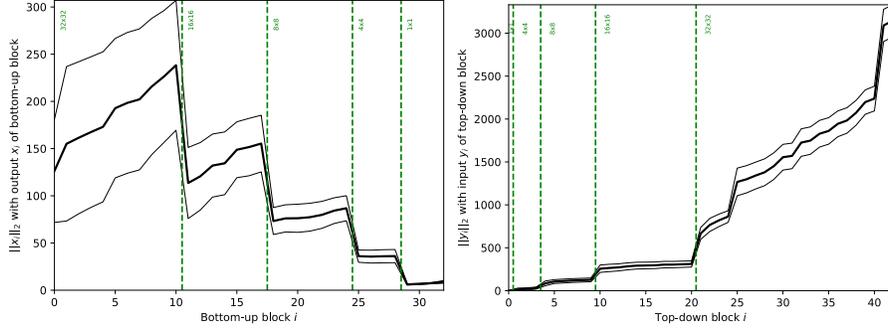


Figure 6: HVAEs secretly represent a notion of time: We measure the  $L_2$ -norm of the residual state for the [Left] forward/bottom-up pass and the [Right] backward/top-down pass over 10 batches with 100 data points each. In both plots, the thick, central line refers to the average and the thin, outer lines refer to  $\pm 2$  standard deviations.

### 3.1 “More from less”: Improving parameter efficiency in HVAEs

In §2.3, we hypothesised that a time-homogeneous model has the same expressivity as a model with time-dependent coefficients, yet uses much less parameters. We start demonstrating this effect by weight-sharing ResNet blocks across time on a small scale. In Fig. 5, we train HVAEs on MNIST and CIFAR10 with  $\{1, 2, 3, 4\}$  ResNet blocks (referred to as  $\{x_1, x_2, x_3, x_4\}$ ) in each resolution with spatial dimensions  $\{32^2, 16^2, 8^2, 4^2, 1^2\}$  (VDVAE\*), and compare their performance when weight-sharing a single parameterised block per resolution  $\{2, 3, 5, 10, 20\}$  times (referred to as  $\{r_2, r_3, r_5, r_{10}, r_{20}\}$ ; WS-VDVAE), excluding projection and embedding blocks. As hypothesised by our framework, yet very surprising in HVAEs, NLL after 1m iterations measured on the validation set gradually increases the more often blocks are repeated even though all weight-sharing models have an identical parameter count to the  $x_1$  model (MNIST: 107k, CIFAR10: 8.7m). Furthermore, the weight-sharing models often outperform or reach equal NLLs compared to  $x_2, x_3, x_4$ , all of which have more parameters (MNIST: 140k; 173k; 206k. CIFAR10: 13.0m; 17.3m; 21.6m), yet fewer activations, latent variables, and number of timesteps at which the coupled SDE in Eq. (A.76) is discretised.

We now scale these findings up to large-scale hyperparameter configurations. We train VDVAE closely following the state-of-the-art hyperparameter configurations in [9], specifically with the same number of parameterised blocks and without weight-sharing (VDVAE\*), and compare them against models with weight-sharing (WS-VDVAE) and fewer parameters, i.e. fewer parameterised blocks, in Table 1. On all four datasets, the weight-shared models achieve similar NLLs with fewer parameters compared to their counterparts without weight-sharing: We use 32%, 36%, 54%, and 40% less parameters on the four datasets reported in Table 1, respectively. For the larger runs, weight-sharing has diminishing returns on NLL as these already have many discretisation steps. To the best of our knowledge, our models achieve a new state-of-the-art performance in terms of NLL compared to any HVAE on CIFAR10, ImageNet32 and CelebA. Furthermore, our WS-VDVAE models have stochastic depths of 57, 105, 235, 125, respectively, the highest ever trained. In spite of these results, it is worth noting that current HVAEs, and VDVAE in particular remains notoriously unstable to train, partly due to the instabilities identified in Theorem 5, and finding the right hyperparameters helps, but cannot solve this.

### 3.2 HVAEs secretly represent time and make use of it

In §3.1, we showed how we can exploit insight on HVAEs through our framework to make HVAEs more parameter efficient. We now want to explain and understand this behavior further. In Fig. 6, we measure  $\|Z_i\|_2$ , the  $L_2$ -norm of the residual state at every backward/top-down block with index  $i$ , over several batches for models trained on MNIST (see Appendix G.2 for the corresponding figure of the forward/bottom-up pass, and similar results on CIFAR10 and ImageNet32). On average, we experience an increase in the state norm across time in every resolution, interleaved by discontinuous ‘jumps’ at the resolution transitions (projection or embedding) where the dimension of the residual state changes. This supports our claim in §2 that HVAEs discretise multi-resolution diffusion

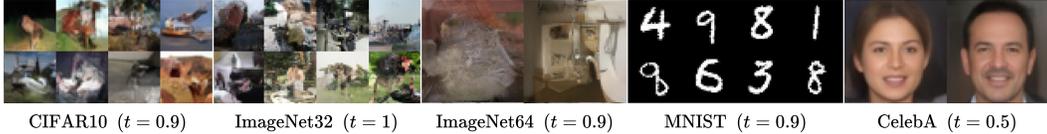


Figure 7: Unconditional samples (not cherry-picked) of VDAE\*. While samples on MNIST and CelebA demonstrate high fidelity and diversity, samples on CIFAR10, ImageNet32 and ImageNet64 are diverse, but are unrecognisable, demonstrating the instabilities identified by Theorem 1. Temperatures  $t$  are tuned for maximum fidelity.

processes which are increasing in the  $Z_i$  variables, and hence learn to represent a notion of time in their residual state.

It is now straightforward to ask how HVAEs benefit from this time representation during training: As we show in Table 2, when normalising the state by its norm at every forward and backward block during training, i.e. forcing a “flat line” in Fig. 6 [Left], learning deteriorates after a short while, resulting in poor NLL results compared to the runs with a regular, non-normalised residual state. This evidence confirms our earlier stated hypothesis: The time representation in ResNet-based HVAEs encodes information which recent HVAEs heavily rely on during learning.

### 3.3 Sampling instabilities in HVAEs

High fidelity unconditional samples of faces, e.g. from models trained on CelebA, cover the front pages of state-of-the-art HVAE papers [9, 10]. Here, we question whether face datasets are an appropriate benchmark for HVAEs. In Theorem 5, we identified the aforementioned state-of-the-art HVAEs as flow from a point mass, hypothesising instabilities during sampling. And indeed, when sampling from our trained VDAE\* with state-of-the-art configurations, we observe high fidelity and diversity samples on MNIST and CelebA, but unrecognisable, yet diverse samples on CIFAR10, ImageNet32 and ImageNet64, in spite of state-of-the-art test set NLLs (see Fig. 7 and Appendix G.3). We argue that MNIST and CelebA, i.e. numbers and faces, have a more uni-modal nature, and are in this sense easier to learn for a discretised multi-resolution process flowing to a point mass, which is uni-modal, than the other “in-the-wild”, multi-modal datasets. Trying to approximate the latter with the, in this case unsuitable, HVAE model leads to the sampling instabilities observed.

Table 2: NLL of HVAEs with and without normalisation of the residual state  $Z_i$ .

Residual state	NLL
<b>MNIST</b>	
Normalised ( $\times$ )	$\leq 464.68$
Non-normalised	$\leq 81.69$
<b>CIFAR10</b>	
Normalised ( $\times$ )	$\leq 6.80$
Non-normalised	$\leq 2.93$
<b>ImageNet</b>	
Normalised	$\leq 6.76$
Non-normalised	$\leq 3.68$

### 3.4 Ablation studies

We conducted several ablation studies which support our experimental results and further probe our multi-resolution framework for HVAEs. In this section we note key findings—a detailed account of all ablations can be found in Appendix G.4. In particular, we find that the number of latent variables, which correlates with stochastic depth, does not explain the performance observed in §3.1, supporting our claims. We further show that Fourier features do not provide a performance gain in HVAEs, in contrast to state-of-the-art diffusion models, where they significantly improve performance [19]. This is consistent with our framework’s finding that a U-Net architecture with pooling is already forced to learn a Haar wavelet basis representation of the data, hence introducing another basis does not add value. We also demonstrate that residual cells are crucial for the performance of HVAEs as they are able to approximate the dynamics of a diffusion process and impose an SDE structure into the model, empirically compare a multi-resolution bridge to a single-resolution model, and investigate synchronous vs. asynchronous processing in time between the forward and backward pass.

## 4 Related work

**U-Nets.** A U-Net [1] is an autoencoding architecture with multiple resolutions where skip connections enable information to pass between matched layers on opposite sides of the autoencoder’s bottleneck. These connections also smooth out the network’s loss landscape [47]. In the literature, U-Nets tend to be convolutional, and a wide range of different approaches have been used for up-sampling and down-sampling between resolutions, with many using average pooling for the down-sampling operation [13, 14, 16, 17, 19]. In this work, we focus on U-Nets as operators on measures interleaved by average pooling as the down-sampling operation (and a corresponding inclusion operation for up-sampling), and we formally characterise U-Nets in Section 2.1 and Appendix B.2. Prior to our work, the dimensionality-reducing bottleneck structure of U-Nets was widely acknowledged as being useful, however it was unclear what regularising properties a U-Net imposes. We provided these in §2.

**HVAEs.** The evolution of HVAEs can be seen as a quest for a parameterisation with more expressiveness than single-latent-layer VAEs [48], while achieving stable training dynamics that avoid common issues such as posterior collapse [36, 49] or exploding gradients. Early HVAEs such as LVAE condition each latent variable directly on only the previous one by taking samples forward [36, 37]. Such VAEs suffer from stability issues even for very small stochastic depths. *Nouveau VAEs (NVAE)* [10] and *Very Deep VAEs (VDVAE)* [9] combine the improvements of several earlier HVAE models (see Appendix B for details), while scaling up to larger stochastic depths. Both use ResNet-based backbones, sharing parameters between the generative and recognition parts of the model. VDVAE is the considerably simpler approach, in particular avoiding common tricks such as a warm-up deterministic autoencoder training phase or data-specific initialisation. VDVAE achieves a stochastic depth of up to 78, improving performance with more ResNet blocks. Worth noting is that while LVAE and NVAE use convolutions with appropriately chosen stride to jump between resolutions, VDVAE use average pooling. In all HVAEs to date, a theoretical underpinning which explains architectural choices, for instance the choice of residual cell, is missing, and we provided this in Section §2.3.

## 5 Conclusion

In this work, we introduced a multi-resolution framework for U-Nets. We provided theoretical results which uncover the regularisation property of the U-Nets bottleneck architecture with average pooling as implicitly learning a Haar wavelet representation of the data. We applied our framework to HVAEs, identifying them as multi-resolution diffusion processes flowing to a point mass. We characterised their backward cell as a type of two-step forward Euler discretisations, providing an alternative to score-matching to approximate a continuous-time diffusion process [16, 18], and observed parameter redundancies and instabilities. We verified the latter theoretical insights in both small- and large-scale experiments, and in doing so trained the deepest ever HVAEs. We explained these results by showing that HVAEs learn a representation of time and performed extensive ablation studies.

An important limitation is that the proven regularisation property of U-Nets is limited to using average pooling as the down-sampling operation. Another limitation is that we only applied our framework to HVAEs, though it is possible to apply it to other model classes. It could also be argued that the lack of exhaustive hyperparameter optimisation performed is a limitation of the work as it may be possible to obtain improved results. We demonstrate, however, that simply adding weight-sharing to the hyperparameter settings given in the original VDVAE paper [9] leads to state-of-the-art performance with improved parameter efficiency, and hence view it as a strength of our results.

## Acknowledgments and Disclosure of Funding

Fabian Falck acknowledges the receipt of studentship awards from the Health Data Research UK-The Alan Turing Institute Wellcome PhD Programme in Health Data Science (Grant Ref: 218529/Z/19/Z), and the Enrichment Scheme of The Alan Turing Institute under the EPSRC Grant EP/N510129/1. Chris Williams acknowledges support from the Defence Science and Technology (DST) Group and from a ESRC DTP Studentship. Dominic Danks is supported by a Doctoral Studentship from The Alan Turing Institute under the EPSRC Grant EP/N510129/1. Christopher Yau is funded by a UKRI Turing AI Fellowship (Ref: EP/V023233/1). Chris Holmes acknowledges support from the Medical Research Council Programme Leaders award MC\_UP\_A390\_1107, The Alan Turing Institute, Health Data Research, U.K., and the U.K. Engineering and Physical Sciences Research Council through the Bayes4Health programme grant. Arnaud Doucet acknowledges support of the UK Defence Science and Technology Laboratory (Dstl) and EPSRC grant EP/R013616/1. This is part of the collaboration between US DOD, UK MOD and UK EPSRC under the Multidisciplinary University Research Initiative. Arnaud Doucet also acknowledges support from the EPSRC grant EP/R034710/1. Matthew Willetts is grateful for the support of UCL Computer Science and The Alan Turing Institute.

The authors report no competing interests.

The three compute clusters used in this work were provided by the Alan Turing Institute, the Oxford Biomedical Research Computing (BMRC) facility, and the Baskerville Tier 2 HPC service (<https://www.baskerville.ac.uk/>) which we detail in the following. First, this research was supported in part through computational resources provided by The Alan Turing Institute under EPSRC grant EP/N510129/1 and with the help of a generous gift from Microsoft Corporation. Second, we used the Oxford BMRC facility, a joint development between the Wellcome Centre for Human Genetics and the Big Data Institute supported by Health Data Research UK and the NIHR Oxford Biomedical Research Centre. The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or the Department of Health. Third, Baskerville was funded by the EPSRC and UKRI through the World Class Labs scheme (EP/T022221/1) and the Digital Research Infrastructure programme (EP/W032244/1) and is operated by Advanced Research Computing at the University of Birmingham.

We thank Tomas Lazauskas, Jim Madge and Oscar Giles from the Alan Turing Institute’s Research Engineering team for their help and support. We thank Adam Huffman, Jonathan Diprose, Geoffrey Ferrari and Colin Freeman from the Biomedical Research Computing team at the University of Oxford for their help and support. We thank Haoting Zhang (University of Cambridge) for valuable comments on the implementation; Huiyu Wang (Johns Hopkins University) for a useful discussion on gradient checkpointing; and Ruining Li and Hanwen Zhu (University of Oxford) for kindly proofreading the manuscript.

## References

- [1] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.
- [2] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- [3] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [4] Xiaojiao Mao, Chunhua Shen, and Yu-Bin Yang. Image restoration using very deep convolutional encoder-decoder networks with symmetric skip connections. In *Advances in Neural Information Processing Systems*, volume 29, 2016.

- [5] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Multi-stage progressive image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021.
- [6] Yuyang Shi, Valentin De Bortoli, George Deligiannidis, and Arnaud Doucet. Conditional simulation using diffusion schrödinger bridges. *arXiv preprint arXiv:2202.13460*, 2022.
- [7] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *Computer Vision – ECCV 2020*. Springer International Publishing, 2020.
- [8] Zoe Landgraf, Fabian Falck, Michael Bloesch, Stefan Leutenegger, and Andrew J. Davison. Comparing view-based and map-based semantic labelling in real-time slam. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6884–6890, 2020.
- [9] Rewon Child. Very Deep VAEs Generalize Autoregressive Models and Can Outperform Them on Images. In *International Conference on Learning Representations*, 2021.
- [10] Arash Vahdat and Jan Kautz. NVAE: A deep hierarchical variational autoencoder. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
- [11] Louay Hazami, Rayhane Mama, and Ragavan Thurairatnam. Efficient-ldvae: Less is more. *arXiv preprint arXiv:2203.13751*, 2022.
- [12] Simon AA Kohl, Bernardino Romera-Paredes, Klaus H Maier-Hein, Danilo Jimenez Rezende, SM Eslami, Pushmeet Kohli, Andrew Zisserman, and Olaf Ronneberger. A hierarchical probabilistic u-net for modeling multi-scale ambiguities. *arXiv preprint arXiv:1905.13077*, 2019.
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
- [14] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021.
- [15] Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. *arXiv preprint arXiv:2102.09672*, 2021.
- [16] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [17] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- [18] Valentin De Bortoli, James Thornton, Jeremy Heng, and Arnaud Doucet. Diffusion schrödinger bridge with applications to score-based generative modeling. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- [19] Diederik P. Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational Diffusion Models. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- [20] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.
- [21] Jonathan Ho, Xi Chen, Aravind Srinivas, Yan Duan, and Pieter Abbeel. Flow++: Improving flow-based generative models with variational dequantization and architecture design. In *International Conference on Machine Learning*, 2019.
- [22] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in Neural Information Processing Systems*, 31, 2018.
- [23] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.

- [24] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. Conditional image generation with pixelcnn decoders. *Advances in Neural Information Processing Systems*, 29, 2016.
- [25] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, 2014.
- [26] Shengjia Zhao, Jiaming Song, and Stefano Ermon. Learning hierarchical features from deep generative models. In *International Conference on Machine Learning*, 2017.
- [27] Fabian Falck, Haoting Zhang, Matthew Willetts, George Nicholson, Christopher Yau, and Chris C Holmes. Multi-facet clustering variational autoencoders. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- [28] Diederik P. Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. In *Advances in Neural Information Processing Systems*, 2016.
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [30] Ingrid Daubechies. *Ten lectures on wavelets*. SIAM, 1992.
- [31] David Taubman and Michael Marcellin. *JPEG2000 image compression fundamentals, standards and practice: image compression fundamentals, standards and practice*, volume 642. Springer Science & Business Media, 2012.
- [32] Alfred Haar. *Zur theorie der orthogonalen funktionensysteme*. Georg-August-Universitat, Gottingen., 1909.
- [33] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
- [34] Haohan Wang, Xindi Wu, Zeyi Huang, and Eric P Xing. High-frequency component helps explain the generalization of convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8684–8694, 2020.
- [35] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, volume 20, 2007.
- [36] Casper Kaae Sønderby, Tapani Raiko, Lars Maaløe, Søren Kaae Sønderby, and Ole Winther. Ladder variational autoencoders. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [37] Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*, 2015.
- [38] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [39] Rob Cornish, Anthony Caterini, George Deligiannidis, and Arnaud Doucet. Relaxing bijectivity constraints with continuously indexed normalising flows. In *International conference on machine learning*, pages 2133–2143. PMLR, 2020.
- [40] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.
- [41] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv preprint arXiv:2006.11239*, 2020.

- [42] Yann LeCun, Corinna Cortes, and C. J. Burges. MNIST handwritten digit database. <http://yann.lecun.com/exdb/mnist/>, 2010.
- [43] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [44] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [45] Patryk Chrabaszcz, Ilya Loshchilov, and Frank Hutter. A downsampled variant of imagenet as an alternative to the cifar datasets. *arXiv preprint arXiv:1707.08819*, 2017.
- [46] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- [47] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- [48] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014.
- [49] Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349*, 2015.
- [50] D La Torre and F Mendivil. The monge–kantorovich metric on multimeasures and self–similar multimeasures. *Set-Valued and Variational Analysis*, 23(2):319–331, 2015.
- [51] Svetlozar T Rachev. *Probability metrics and the stability of stochastic models*, volume 269. Wiley, 1991.
- [52] Stephane G Mallat. Multiresolution approximations and wavelet orthonormal bases of  $L^2(\mathbb{R}^n)$ . *Transactions of the American mathematical society*, 315(1):69–87, 1989.
- [53] Naftali Tishby, Fernando C. Pereira, and William Bialek. The Information Bottleneck Method. 2000.
- [54] Lars Maaløe, Marco Fraccaro, Valentin Liévin, and Ole Winther. BIVA: A very deep hierarchy of latent variables for generative modeling. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [55] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost. *arXiv preprint arXiv:1604.06174*, 2016.
- [56] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019.
- [57] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, et al. Array programming with NumPy. *Nature*, 585(7825):357–362, 2020.
- [58] Lukas Biewald. Experiment tracking with weights and biases, 2020. Software available from wandb.com.
- [59] NVIDIA. Nvidia apex. <https://github.com/NVIDIA/apex>.
- [60] Guido Van Rossum. *The Python Library Reference, release 3.8.2*. Python Software Foundation, 2020.

- [61] J. D. Hunter. Matplotlib: A 2D graphics environment. *Computing in Science & Engineering*, 9(3):90–95, 2007.
- [62] Almar et al. Klein. Imageio. [https://zenodo.org/record/6551868#.Yolo\\_5PMIhg](https://zenodo.org/record/6551868#.Yolo_5PMIhg).
- [63] Lisandro Dalcin and Yao-Lung L Fang. Mpi4py: Status update after 12 years of development. *Computing in Science & Engineering*, 23(4):47–54, 2021.
- [64] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [65] P Umesh. Image processing in python. *CSI Communications*, 23, 2012.
- [66] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. 80 million tiny images. <http://groups.csail.mit.edu/vision/TinyImages/>.
- [67] Vinay Uday Prabhu and Abeba Birhane. Large image datasets: A pyrrhic win for computer vision? *arXiv preprint arXiv:2006.16923*, 2020.
- [68] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Large-scale celebfaces attributes (celeba) dataset. 2022.
- [69] Tomas Mikolov et al. Statistical language models based on neural networks. *Presentation at Google, Mountain View, 2nd April*, 80(26), 2012.
- [70] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *European Conference on Computer Vision*, pages 646–661. Springer, 2016.
- [71] 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. *Advances in neural information processing systems*, 30, 2017.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]** We state our theoretical and experimental contributions in §1. We develop our multi-resolution framework with corresponding Theorems 1 to 5 in §2 (proofs in Appendix A). We provide our experimental results in §3.
  - (b) Did you describe the limitations of your work? **[Yes]** See §5.
  - (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** See Appendix E.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? **[Yes]** We stated the assumptions of all theoretical results in §2 and refer to Appendices A and B for further details.
  - (b) Did you include complete proofs of all theoretical results? **[Yes]** The complete proofs of Theorems 1 to 5, as well as of all additional experimental results are provided in Appendix A.
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]** We provide our code, and instructions on how to download the data and reproduce the main results in the supplementary material. In particular, we refer to the README.md file in the code repository for further details which follow the NeurIPS Code Completeness Checklist.

- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]** See Appendices **F** and **D**.
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[No]** Due to the significant computational cost of training extremely deep HVAEs (multiple Nvidia A100 graphic cards with 40GB of GPU memory each running for 3 weeks per run), we did not perform the multiple runs per hyperparameter setting required to provide error bars for our runs. We note that this is common practice in large-scale HVAE research (compare [9, 10]). Furthermore, in their code base, VDVAE [9], which we directly base our architecture on, reported highly stable test NLL when varying the random seed, varying only in the second decimal place in bits per dim.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** We discuss our computational resources used and an estimate of the total amount of compute required to reproduce our results briefly in §3, and in detail in Appendix **C**.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- (a) If your work uses existing assets, did you cite the creators? **[Yes]** We discuss all existing assets used in our code base in Appendix **C**.
  - (b) Did you mention the license of the assets? **[Yes]** We mention the license of all assets used in Appendix **C**.
  - (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]** We provide our source code with instructions on how to reproduce our results in the supplementary material. Further details on our code are provided in Appendix **C**.
  - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? **[Yes]** We address this question in Appendix **D**.
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[Yes]** We likewise address this question in Appendix **D**.
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[N/A]**
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? **[N/A]**
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **[N/A]**