

## A DERIVATION OF THE LOSS

In this section we derive the objectives given in Eqs. (4) and (9), following a path similar to that presented in Ho et al. [10]. Let  $\mathbf{s}_t = (\mathbf{x}_t, \mathbf{y}_t)$  be the state that combines both the inputs and observations. We wish to find the set of parameters  $\theta$  which maximise the likelihood given the initial state,  $p_\theta(\mathbf{s}_0)$ . While a direct evaluation of the likelihood appears intractable,

$$p_\theta(\mathbf{s}_0) = \int p_\theta(\mathbf{s}_{0:T}) d\mathbf{s}_{1:T},$$

this may be recast in a form which allows for a comparison to be drawn between the forward and reverse trajectories [31]

$$p_\theta(\mathbf{s}_0) = \int p_\theta(\mathbf{s}_T) q(\mathbf{s}_{1:T} | \mathbf{s}_0) \prod_{t=1}^T \frac{p_\theta(\mathbf{s}_{t-1} | t)}{q(\mathbf{s}_t | \mathbf{s}_{t-1})} d\mathbf{s}_{1:T}.$$

The appeal of introducing the reverse process is that it is tractable when conditioned on the first state  $\mathbf{s}_0$ , taking a Gaussian form

$$q(\mathbf{s}_{t-1} | \mathbf{s}_t, \mathbf{s}_0) = \mathcal{N}(\mathbf{s}_{t-1} | \tilde{\boldsymbol{\mu}}(\mathbf{s}_t, \mathbf{s}_0), \tilde{\boldsymbol{\beta}}_t \mathbf{I}).$$

Our loss function reflects a lower bound on the negative log likelihood

$$\mathbb{E}[-\log p_\theta(\mathbf{s}_0)] \leq \mathbb{E}_{q(\mathbf{s}_{0:T})} \left[ \log \frac{q(\mathbf{s}_{1:T} | \mathbf{s}_0)}{p_\theta(\mathbf{s}_{0:T})} \right] := \mathcal{L}_\theta,$$

which may be decomposed into two edge terms and a sum over the intermediate steps, as follows

$$\mathcal{L}_\theta = \mathbb{E}_q \left[ L_T + \sum_{t=2}^T L_{t-1} + L_0 \right],$$

where

$$\begin{aligned} L_0 &= -\log p_\theta(\mathbf{s}_0 | \mathbf{s}_1), \\ L_{t-1} &= \text{KL}(q(\mathbf{s}_{t-1} | \mathbf{s}_t, \mathbf{s}_0) \| p_\theta(\mathbf{s}_{t-1} | \mathbf{s}_t)), \\ L_T &= \text{KL}(q(\mathbf{s}_T | \mathbf{s}_0) \| p_\theta(\mathbf{s}_T)). \end{aligned} \tag{10}$$

We can write

$$L_{t-1} = \mathbb{E}_{\mathbf{s}_0, \boldsymbol{\varepsilon}} \left[ \frac{1}{2\sigma^2} \|\tilde{\boldsymbol{\mu}}(\mathbf{s}_t, \mathbf{s}_0) - \boldsymbol{\mu}_\theta(\mathbf{s}_t, t)\|^2 \right], \tag{11}$$

where the mean is a function of the previous and first state

$$\tilde{\boldsymbol{\mu}}(\mathbf{s}_t, \mathbf{s}_0) = \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{s}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{s}_0$$

and the variance

$$\tilde{\boldsymbol{\beta}}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t.$$

It is helpful to consider the relationship between the initial and final states

$$\mathbf{s}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{s}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_s),$$

where  $\boldsymbol{\varepsilon}_s = (\boldsymbol{\varepsilon}_x, \boldsymbol{\varepsilon}_y)$ , so that we can rewrite the mean as

$$\tilde{\boldsymbol{\mu}}(\mathbf{s}_t, \boldsymbol{\varepsilon}_s) = \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{s}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon}_s). \tag{12}$$

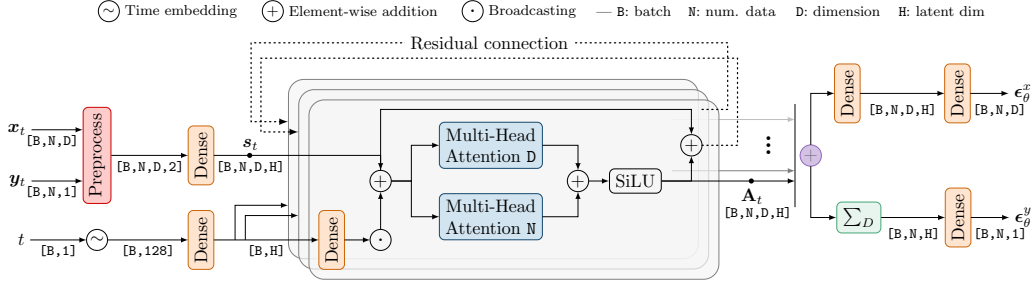


Figure 7: Architecture of the NDP’s neural network noise models  $\epsilon_\theta^x$  and  $\epsilon_\theta^y$ . Compared to Fig. 3 this architecture has two outputs, one to predict the corruption on the inputs  $x_t$  and one to predict the corruption on the function outputs  $y_t$ . Both output share a lot of weights in the network and only bifurcate in the last few layers.

Equation 11 can now be expressed as

$$L_{t-1} = \mathbb{E}_{s_0, \epsilon_s} \left[ \frac{\beta_t^2}{2\sigma^2\alpha_t(1-\bar{\alpha}_t)} \|\epsilon_s - \epsilon_\theta^s\|^2 \right]. \quad (13)$$

Finally, since the variance schedule is fixed, and the edge terms are not found to improve empirical performance, our simplified training objective is given by

$$\mathcal{L}_\theta = \mathbb{E}_{x_0, y_0, \epsilon_x, \epsilon_y, t} \left[ \|\epsilon_x - \epsilon_\theta^x(x_t, y_t, t)\|^2 + \|\epsilon_y - \epsilon_\theta^y(x_t, y_t, t)\|^2 \right]. \quad (14)$$

So far we have derived the objective of the ‘full’ NDP model (i.e., the model which diffuses both  $x$  and  $y$ ) given in Eq. (9) of the main paper. In Fig. 7 we detail the architecture for the noise models  $\epsilon_\theta^x(\cdot)$  and  $\epsilon_\theta^y(\cdot)$ .

#### A.1 DETERMINISTIC INPUT LOCATIONS

In supervised learning the inputs are typically known in advance, which renders the joint generative distribution  $p(x_0, y_0)$  of less importance than the conditional  $p(y_0 | x_0)$ . Therefore, if regression is the use-case of your NDP model, we recommend not to diffuse the inputs, but instead keeping them fixed, i.e.  $x_0 = x_1 = \dots = x_T$ . Mathematically, this is equivalent to setting  $\epsilon_x$  deterministically to zero, which renders its corresponding predictor  $\epsilon_\theta^x(\cdot)$  unnecessary and thus reduces Eq. (14) to

$$\mathcal{L}_\theta = \mathbb{E}_{x_0, y_0, \epsilon_y, t} \left[ \|\epsilon_y - \epsilon_\theta^y(x_0, y_t, t)\|^2 \right], \quad (15)$$

which matches Eq. (4) from the main paper. The noise model  $\epsilon_\theta^y(\cdot)$  now acts on the *uncorrupted* inputs  $x_0$  but still on corrupted observations  $y_t = \sqrt{\bar{\alpha}_t}y_0 + \sqrt{1-\bar{\alpha}_t}\epsilon_y$ .

## B ALGORITHMS

In this section we list pseudo-code for training and sampling NDPs. The code to handle the special case of deterministic input locations (Appendix A.1) is marked by ‘OR’.

### B.1 TRAINING

The training procedure is concerned with fitting the parameters of the NDP’s noise model neural network, with outputs  $\epsilon_\theta^x$  and  $\epsilon_\theta^y$ , such that they can accurately predict the noise that was added to the corrupted function draws. In each iteration, we create a batch of datasets, each originating from a Gaussian process draw. We stress that the number of datapoints and even the input dimensionality of the datasets can vary within the training procedure.

**Algorithm 1** Training

---

**input** Input distribution  $q(\mathbf{x}_0)$  and covariance function  $k_\psi$ , with prior over hyperparameters  $p_\psi$ . Noise schedule  $\beta_t$  for  $t \in \{1, 2, \dots, T\}$ . A loss function  $L$  (e.g., MSE or MAE).

**begin**

Precompute  $\gamma_t = \sqrt{1 - \bar{\alpha}_t}$  and  $\bar{\alpha}_t = \prod_{j=1}^t (1 - \beta_j)$ .

**for**  $i = 1, 2, \dots, N_{\text{iter}}$  **do**

Sample  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ ,  $\psi \sim p_\psi$ ,  $\mathbf{y}_0 \sim \mathcal{N}(\mathbf{0}, k_\psi(\mathbf{x}_0, \mathbf{x}_0) + \sigma^2 \mathbf{I})$ .

Sample  $\epsilon_x, \epsilon_y \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and  $t \sim \mathcal{U}(\{1, \dots, T\})$ .

Compute  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \gamma_t \epsilon_x$  and  $\mathbf{y}_t = \sqrt{\bar{\alpha}_t} \mathbf{y}_0 + \gamma_t \epsilon_y$ .

Update  $\theta$  using gradient  $\nabla_\theta [L(\epsilon_x, \epsilon_\theta^x(\mathbf{x}_t, \mathbf{y}_t, t)) + L(\epsilon_y, \epsilon_\theta^y(\mathbf{x}_t, \mathbf{y}_t, t))]$ . Eq. (9)

OR, in the case that we only diffuse on the observations:

Update  $\theta$  using gradient  $\nabla_\theta L(\epsilon_y, \epsilon_\theta^y(\mathbf{x}_0, \mathbf{y}_t, t))$ . Eq. (4)

---

## B.2 PRIOR AND CONDITIONAL SAMPLING

In practice, the code to sample the prior is a special case of the conditional code with an empty context dataset. However, here we list them both for the clarity of the exposition.

**Algorithm 2** Prior Sampling

---

**input** A pretrained NDP noise model  $\epsilon_\theta^x(\cdot)$  and  $\epsilon_\theta^y(\cdot)$ . Precomputed  $\gamma_t$  and  $\bar{\alpha}_t$  from a given noise schedule  $\beta_t$ .

**begin**

Sample a random initial state  $\mathbf{x}_T$  and  $\mathbf{y}_T$  from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ .

**for**  $t = T, T - 1, \dots, 1$  **do**

Sample using backward kernel:

Eq. (3)

$$\begin{bmatrix} \mathbf{x}_{t-1} \\ \mathbf{y}_{t-1} \end{bmatrix} \sim \mathcal{N}\left(\frac{1}{\sqrt{1 - \beta_t}} \left( \begin{bmatrix} \mathbf{x}_t \\ \mathbf{y}_t \end{bmatrix} - \frac{\beta_t}{\gamma_t} \begin{bmatrix} \epsilon_\theta^x(\mathbf{x}_t, \mathbf{y}_t, t) \\ \epsilon_\theta^y(\mathbf{x}_t, \mathbf{y}_t, t) \end{bmatrix} \right), \frac{\gamma_t^2}{\gamma_{t-1}^2} \beta_t \mathbf{I} \right).$$

OR, in the case that we only diffuse on the observations and the inputs  $\mathbf{x}_0$  are given:

$$\mathbf{y}_{t-1} \sim \mathcal{N}\left(\frac{1}{\sqrt{1 - \beta_t}} \left( \mathbf{y}_t - \frac{\beta_t}{\gamma_t} \epsilon_\theta^y(\mathbf{x}_0, \mathbf{y}_t, t) \right), \frac{\gamma_t^2}{\gamma_{t-1}^2} \beta_t \mathbf{I} \right).$$

**return**  $(\mathbf{x}_0, \mathbf{y}_0)$

---

The conditional sampling algorithm is based on recent advances in image inpainting using diffusion models [19]. In Algorithm 3, for optimal results, we follow Lugmayr et al. [19] and use a re-sampling step which mixes the predictions with the context dataset. In all our experiments  $U$  is set to 5.

**Algorithm 3** Conditional Sampling

**input** A context dataset  $\mathcal{D} = \{([x_0^c]_i \in \mathcal{X}, [y_0^c]_i \in \mathbb{R})\}_{i=1}^N$ . A pretrained NDP noise model  $\epsilon_\theta^x(\cdot)$  and  $\epsilon_\theta^y(\cdot)$ . Precomputed  $\gamma_t$  and  $\bar{\alpha}_t$  from a given noise schedule  $\beta_t$ .

**begin**

Let the augmented states be  $\check{x}_T = [x_T^c, x_T]$  and  $\check{y}_T = [y_T^c, y_T]$ . Sample them according to  $\check{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $\check{y}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .

Create a mask  $m = [0, \dots, 0, 1, \dots, 1]$  indicating the positions of  $x_t^c/y_t^c$  with 0s and  $x_t/y_t$  with 1s.

**for**  $t = T, T-1, \dots, 1$  **do**

**for**  $u = 1, \dots, U$  **do**

    Sample backward from augmented intermediate state:

Eq. (3)

$$\begin{bmatrix} \check{x}_{t-1} \\ \check{y}_{t-1} \end{bmatrix} \sim \mathcal{N}\left(\frac{1}{\sqrt{1-\beta_t}}\left(\begin{bmatrix} \check{x}_t \\ \check{y}_t \end{bmatrix} - \frac{\beta_t}{\gamma_t} \begin{bmatrix} \epsilon_\theta^x(\check{x}_t, \check{y}_t, t) \\ \epsilon_\theta^y(\check{x}_t, \check{y}_t, t) \end{bmatrix}\right), \frac{\gamma_t^2}{\gamma_{t-1}^2} \beta_t \mathbf{I}\right)$$

    Select unknown part from the augmented state using the mask:

$$x_{t-1} = m \odot \check{x}_{t-1} \quad \text{and} \quad y_{t-1} = m \odot \check{y}_{t-1}$$

    Sample context points  $(t-1)$ -steps forward from context dataset:

$$\begin{bmatrix} x_{t-1}^c \\ y_{t-1}^c \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} x_{t-1}^c \\ y_{t-1}^c \end{bmatrix} \mid \sqrt{\bar{\alpha}_{t-1}} \begin{bmatrix} x_0^c \\ y_0^c \end{bmatrix}, \gamma_{t-1}^2 \mathbf{I}_n\right)$$

    Update the augmented intermediate state using the forward sampled context dataset and the new intermediate state:

$$\check{x}_{t-1} = \begin{bmatrix} x_{t-1}^c \\ x_{t-1} \end{bmatrix} \quad \text{and} \quad \check{y}_{t-1} = \begin{bmatrix} y_{t-1}^c \\ y_{t-1} \end{bmatrix}$$

    Diffuse forward by one step

$$\begin{bmatrix} \check{x}_t \\ \check{y}_t \end{bmatrix} \sim \mathcal{N}(\sqrt{1-\beta_t} \begin{bmatrix} \check{x}_{t-1} \\ \check{y}_{t-1} \end{bmatrix}, \beta_t \mathbf{I})$$

**return**  $(x_0, y_0)$

## C PROOFS

In this section, we shall formally demonstrate that NDP’s noise model adheres to the symmetries associated with permutations of the datapoint orderings and the permutations of the input dimensions. We focus our attention to the full noise model of Fig. 7 because Fig. 3 is a simplification of it in which we only keep a single output. We start by proving properties of its main building block: the bi-dimensional attention block. We can then straightforwardly prove the equivariance and invariance properties that hold for the NDP’s noise model. Before that, we start by setting the notation.

### C.1 NOTATION, DEFINITIONS AND PRELIMINARY LEMMAS

**Notation.** Let  $s \in \mathbb{R}^{N \times D \times H}$  be a tensor of rank (or dimension) three, where we refer to each dimension according to the following convention:

$$\text{shape}(s) = [N, D, H].$$

First dim.: sequence length  
 Third dim.: embedding  
 Second dim.: input dimensionality

In the next definitions we will use NumPy-based indexing and slicing notation. We assume the reader is familiar with this convention. Most notably, we use a colon (:) to reference every element in a dimension.

**Definition 1.** Let  $\Pi_N$  be the set of all permutations of indices  $\{1, \dots, N\}$ . Let  $\pi_n \in \Pi_N$  and  $\mathbf{s} \in \mathbb{R}^{N \times D \times H}$ . Then  $(\pi_n \circ \mathbf{s}) \in \mathbb{R}^{N \times D \times H}$  denotes a tensor where the ordering of indices in the first dimension are reshuffled (i.e., permuted) according to  $\pi_n$ . We write

$$\pi_n \circ \mathbf{s} = \mathbf{s}_{\pi_n(1), \pi_n(2), \dots, \pi_n(N); :, :}$$

**Definition 2.** Let  $\Pi_D$  be the set of all permutations of indices  $\{1, \dots, D\}$ . Let  $\pi_d \in \Pi_D$  and  $\mathbf{s} \in \mathbb{R}^{N \times D \times H}$ . Then  $(\pi_d \circ \mathbf{s}) \in \mathbb{R}^{N \times D \times H}$  denotes a tensor where the ordering of indices in the second dimension are reshuffled (i.e., permuted) according to  $\pi_d$ . We write

$$\pi_d \circ \mathbf{s} = \mathbf{s}_{:, \pi_d(1), \pi_d(2), \dots, \pi_d(D); :, \cdot}$$

**Definition 3.** A function  $f : \mathbb{R}^{N \times D \times H} \rightarrow \mathbb{R}^{N \times D \times H}$  is equivariant to  $\Pi$  if for any permutation  $\pi \in \Pi$ , we have

$$f(\pi \circ \mathbf{s}) = \pi \circ f(\mathbf{s}).$$

**Definition 4.** A function  $f : \mathbb{R}^{N \times D \times H} \rightarrow \mathbb{R}^{N \times D \times H}$  is invariant to  $\Pi$  if for any permutation  $\pi \in \Pi$ , we have

$$f(\pi \circ \mathbf{s}) = f(\mathbf{s}).$$

Invariance in layman’s terms means that the output is not affected by a permutation of the inputs.

**Lemma 1.** The composition of equivariant functions is equivariant.

*Proof.* Let  $f$  and  $g$  be equivariant to  $\Pi$ , then for all  $\pi \in \Pi$ :

$$(f \circ g)(\pi \circ \mathbf{s}) = f(g(\pi \circ \mathbf{s})) = f(\pi \circ g(\mathbf{s})) = \pi \circ f(g(\mathbf{s})),$$

which shows that the composition,  $f \circ g$ , is equivariant as well.  $\square$

**Lemma 2.** An element-wise operation between equivariant functions remains equivariant.

*Proof.* Let  $f$  and  $g$  be equivariant to  $\Pi$ , then for all  $\pi \in \Pi$  and an element-wise operation  $\oplus$  (e.g., addition), we have

$$(f \oplus g)(\pi \circ \mathbf{s}) = f(\pi \circ \mathbf{s}) \oplus g(\pi \circ \mathbf{s}) = (\pi \circ f(\mathbf{s})) \oplus (\pi \circ g(\mathbf{s})) = \pi \circ (f \oplus g).$$

which shows that the composition,  $f \oplus g$ , is equivariant as well.  $\square$

**Lemma 3.** A function  $f : \mathbb{R}^{N \times D \times H} \rightarrow \mathbb{R}^{N \times D \times H}$  which applies the same function  $g$  to all its rows, i.e.  $f : \mathbf{s} \mapsto [g(\mathbf{s}_{1,:,:,}), g(\mathbf{s}_{2,:,:,}), \dots, g(\mathbf{s}_{N,:,:,})]$  with  $g : \mathbb{R}^{D \times H} \rightarrow \mathbb{R}^{D \times H}$  is equivariant in the first dimension.

*Proof.* Follows immediately from the structure of  $f$ .  $\square$

## C.2 BI-DIMENSIONAL ATTENTION BLOCK

With the notation, definitions and lemmas in place, we now prove that the bi-dimensional attention block is equivariant in its first and second dimension, respectively to permutations in the set  $\Pi_N$  and  $\Pi_D$ . We start this section by formally defining the bi-dimensional attention block and its main components attention components  $\text{Attn}_N$  and  $\text{Attn}_D$ . Finally, in Appendix C.3 we prove the properties of the noise models  $\epsilon_\theta^x$  and  $\epsilon_\theta^y$  making use of the results in this section.

**Definition 5.** Ignoring the batch dimension  $B$ , let  $\mathbf{A}_t : \mathbb{R}^{N \times D \times H} \rightarrow \mathbb{R}^{N \times D \times H}; \mathbf{s} \mapsto \mathbf{A}_t(\mathbf{s})$  be the bi-dimensional attention block. As illustrated in Fig. 7, it operates on three dimensional tensors in  $\mathbb{R}^{N \times D \times H}$  and applies attention across the first and second dimension using  $\text{Attn}_N$  and  $\text{Attn}_D$ , respectively. The final output  $\mathbf{A}_t(\mathbf{s})$  is obtained by summing the two attention outputs, before applying an element-wise non-linearity.

We now proceed by defining the component  $\text{Attn}_D$  and its properties. Subsequently, in Definition 7 we define  $\text{Attn}_N$  and its properties. Finally, we combine both to prove Proposition 1 in the main paper about the bi-dimensional attention block.

**Definition 6.** Let  $\text{Attn}_D : \mathbb{R}^{N \times D \times H} \rightarrow \mathbb{R}^{N \times D \times H}$  be a self-attention block [35] acting across  $D$ . Let  $\sigma$  be a softmax activation function operating on the last dimension of a tensor. Then,  $\text{Attn}_D$  is defined as

$$\text{Attn}_D(\mathbf{s})[n, d, h] = \sum_{d'=1}^D \Sigma_{n,d,d'}(\mathbf{s}) \mathbf{s}_{n,d',h}^v, \quad \text{where} \quad \Sigma_{n,d,d'}(\mathbf{s}) = \sigma\left(\frac{1}{\sqrt{H}} \sum_l \mathbf{s}_{n,d,\ell}^k \mathbf{s}_{n,d',\ell}^q\right)$$

given a linear projection of the inputs  $\mathbf{s}$  which maps them into keys ( $k$ ), queries ( $q$ ) and values ( $v$ )

$$\mathbf{s}_{n,d,\ell}^k = \sum_j \mathbf{s}_{n,d,j} \mathbf{W}_{j,\ell}^k, \quad \mathbf{s}_{n,d',\ell}^q = \sum_j \mathbf{s}_{n,d',j} \mathbf{W}_{j,\ell}^q, \quad \mathbf{s}_{n,d',h}^v = \sum_j \mathbf{s}_{n,d',j} \mathbf{W}_{j,h}^v.$$

**Proposition 3.**  $\text{Attn}_D$  is equivariant to  $\Pi_N$  and  $\Pi_D$  (i.e., across sequence length and input dimensionality).

*Proof.* We prove the equivariance to  $\Pi_N$  and  $\Pi_D$  separately. First, from the definition we can see that  $\text{Attn}_D$  is a function that acts on each element row of  $\mathbf{s}$  separately. Thus, by Lemma 3,  $\text{Attn}_D$  is equivariant to  $\Pi_N$ .

Next, we want to prove equivariance to  $\Pi_D$ . We want to show that for all  $\pi_d \in \Pi_D$ :

$$\text{Attn}_D(\pi_d \circ \mathbf{s}) = \pi_d \circ \text{Attn}_D(\mathbf{s}).$$

The self-attention mechanism consists of a matrix multiplication of the attention matrix  $\Sigma$  and the projected inputs. We start by showing that the attention matrix is equivariant to permutations in  $\Pi_D$

$$\begin{aligned} \Sigma_{n,d,d'}(\pi_d \circ \mathbf{s}) &= \sigma\left(\frac{1}{\sqrt{H}} \sum_{\ell} (\pi_d \circ \mathbf{s})_{n,d,\ell}^k (\pi_d \circ \mathbf{s})_{n,d',\ell}^q\right) \\ &= \sigma\left(\frac{1}{\sqrt{H}} \sum_{\ell,j} (\mathbf{s}_{n,\pi_d(d),j} \mathbf{W}_{j,\ell}^k) (\mathbf{s}_{n,\pi_d(d'),j} \mathbf{W}_{j,\ell}^q)\right) \\ &= \Sigma_{n,\pi_d(d),\pi_d(d')}(\mathbf{s}) \end{aligned}$$

It remains to show that the final matrix multiplication step restores the row-equivariance

$$\begin{aligned} \text{Attn}_D(\pi_d \circ \mathbf{s}) &= \sum_{d'} \Sigma_{n,d,d'}(\pi_d \circ \mathbf{s}) (\pi_d \circ \mathbf{s})_{n,d',h}^v \\ &= \sum_{d'} \Sigma_{n,\pi_d(d),\pi_d(d')}(\mathbf{s}) \mathbf{s}_{n,\pi_d(d'),h}^v \\ &= \sum_{d'} \Sigma_{n,\pi_d(d),d'}(\mathbf{s}) \mathbf{s}_{n,d',h}^v = \pi_d \circ \text{Attn}_D(\mathbf{s}). \end{aligned}$$

This concludes the proof.  $\square$

We continue by defining and proving the properties of the second main component of the bi-dimensional attention block:  $\text{Attn}_N$ .

**Definition 7.** Let  $\text{Attn}_N : \mathbb{R}^{N \times D \times H} \rightarrow \mathbb{R}^{N \times D \times H}$  be a self-attention block [35] acting across  $N$  (i.e. the sequence length). Let  $\sigma$  be a softmax activation function operating on the last dimension of a tensor. Then  $\text{Attn}_N$  is defined as

$$\text{Attn}_N(\mathbf{s})[n, d, h] = \sum_{n'=1}^N \Sigma_{n,d,n'}(\mathbf{s}) \mathbf{s}_{n',d,h}^v, \quad \text{where} \quad \Sigma_{n,d,n'}(\mathbf{s}) = \sigma\left(\frac{1}{\sqrt{H}} \sum_l \mathbf{s}_{n,d,\ell}^k \mathbf{s}_{n',d,\ell}^q\right)$$

given a linear projection of the inputs  $\mathbf{s}$  which maps them into keys ( $k$ ), queries ( $q$ ) and values ( $v$ )

$$\mathbf{s}_{n,d,\ell}^k = \sum_j \mathbf{s}_{n,d,j} \mathbf{W}_{j,\ell}^k, \quad \mathbf{s}_{n',d,\ell}^q = \sum_j \mathbf{s}_{n',d,j} \mathbf{W}_{j,\ell}^q, \quad \mathbf{s}_{n',d,h}^v = \sum_j \mathbf{s}_{n',d,j} \mathbf{W}_{j,h}^v.$$

**Proposition 4.**  $\text{Attn}_N$  is equivariant to  $\Pi_N$  and  $\Pi_D$  (i.e., across sequence length and input dimensionality).

*Proof.* Follows directly from Proposition 3 after transposing the first and second dimension of the input.  $\square$

Finally, we have the necessary ingredients to prove Proposition 1 from the main paper.

**Proposition 5.** *The bi-dimensional attention block  $\mathbf{A}_t$  is equivariant to  $\Pi_D$  and  $\Pi_N$ .*

*Proof.* The bi-dimensional attention block simply adds the output of  $\text{Attn}_D$  and  $\text{Attn}_N$ , followed by an element-wise non-linearity. Therefore, as a direct application of Lemma 2 the complete bi-dimensional attention block remains equivariant to  $\Pi_N$  and  $\Pi_D$ .  $\square$

### C.3 NDP NOISE MODEL

By building on the equivariant properties of the bi-dimensional attention block, we prove the equivariance and invariance of the NDP’s noise models, denoted by  $\epsilon_\theta^x$  and  $\epsilon_\theta^y$ , respectively. We refer to Fig. 7 for their definition. In short,  $\epsilon_\theta^x$  consists of adding the output of several bi-dimensional attention blocks followed by dense layers operating on the last dimension. Similarly,  $\epsilon_\theta^y$  is constructed by summing the bi-dimensional attention blocks, but is followed by a summation over  $D$  before applying a final dense layer.

The following propositions hold:

**Proposition 6.** *The function  $\epsilon_\theta^x$  is equivariant to  $\Pi_D$  and  $\Pi_N$ .*

*Proof.* The output  $\epsilon_\theta^x$  is formed by element-wise summing the output of bi-dimensional attention layers. Directly applying Lemma 1 and Proposition 5 completes the proof.  $\square$

**Proposition 7.** *The function  $\epsilon_\theta^y$  is equivariant to  $\Pi_N$ .*

*Proof.* The summation over  $D$  does not affect the equivariance over  $\Pi_N$  from the bi-dimensional attention blocks as it can be cast as a row-wise operation.  $\square$

**Proposition 8.** *The function  $\epsilon_\theta^y$  is invariant to  $\Pi_D$ .*

*Proof.* Follows from the equivariance of the bi-dimensional blocks and [41, Thm. 7].  $\square$

### C.4 CONSISTENCY

Consistency<sup>2</sup> is achieved when marginalising over the conditional distribution of a random variable  $y^0$ , we recover its original distribution.

$$p(y^0) = \int p(y^0|y^1)p(y^1)dy^1. \quad (16)$$

The reverse process of a diffusion model begins at  $T$  in a state of white noise, and progresses towards step 0. The features represented at  $T$ , denoted  $\mathbf{y}_T := (y_T^0, y_T^1)$  (for simplicity in this proof we only pick two), are by construction a fully consistent stochastic process. We have  $p(\mathbf{y}_T) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and the consistency is trivially satisfied.

The progression from step  $t$  to  $t - 1$ , involves using a deep neural network to estimate the noise present, and subtracting that noise from the signal. This procedure leaves us with a new probability distribution. Specifically, following Ho et al. [10], this is constructed as follows:

$$p_\theta(\mathbf{y}_{t-1} | \mathbf{y}_t, t) := \mathcal{N}(\mathbf{y}_{t-1}; \mu_\theta(\mathbf{y}_t), \beta_t^2 \mathbf{I}). \quad (17)$$

<sup>2</sup>see also <https://yanndubs.github.io/Neural-Process-Family/text/CNPF.html> for a concise overview.

where the new mean function  $\mu(\mathbf{y}_t)$  is a deterministic function, controlled by a neural network, and  $\mathbf{y}_t$  represents the data vector at step  $t$ . Notice the conditional independence of  $y_{t-1}^0$  and  $y_{t-1}^1$  given  $\mathbf{y}_t$  in the reverse process

$$\mathcal{N}(y_{t-1}^0; \mu_\theta^0(\begin{bmatrix} y_t^0 \\ y_t^1 \end{bmatrix}), \beta_t^2) \cdot \mathcal{N}(y_{t-1}^1; \mu_\theta^1(\begin{bmatrix} y_t^0 \\ y_t^1 \end{bmatrix}), \beta_t^2), \quad (18)$$

where we use  $\mu_\theta^0(\cdot)$  and  $\mu_\theta^1(\cdot)$  to select the first and second entry of the neural network output.

Starting from the RHS of Eq. (16)

$$\int p(y_{t-1}^0 | y_{t-1}^1) p(y_{t-1}^1) dy_{t-1}^1 \quad (19)$$

$$= \int_{y_{t-1}^1} \left( \int_{\mathbf{y}_t} p(\mathbf{y}_t) p(y_{t-1}^0 | y_{t-1}^1, \mathbf{y}_t) d\mathbf{y}_t \right) p(y_{t-1}^1) dy_{t-1}^1 \quad (20)$$

$$= \int_{y_{t-1}^1} \left( \int_{\mathbf{y}_t} p(\mathbf{y}_t) p(y_{t-1}^0 | \mathbf{y}_t) d\mathbf{y}_t \right) p(y_{t-1}^1) dy_{t-1}^1 \quad (21)$$

$$= \int_{\mathbf{y}_t} p(\mathbf{y}_t) p(y_{t-1}^0 | \mathbf{y}_t) d\mathbf{y}_t = p(y_{t-1}^0). \quad (22)$$

In Eq. (20) we marginalised over the previous function values  $\mathbf{y}_t$ . This allowed us to make use of the conditional independence between  $y_{t-1}^0$  and  $y_{t-1}^1$  given  $\mathbf{y}_t$  in Eq. (21). Finally, the random variable  $y_{t-1}^1$  integrates to 1, which leads to the LHS of Eq. (16). This argument can be repeated for each step, until we reach  $\mathbf{y}_0$ . The NDP therefore offers samples drawn from a consistent set of conditional distributions.

**Summarising** Consistency of the NDP samples is ensured because *the neural network plays no role in the conditioning step*. It is therefore infeasible for it to be responsible for generating inconsistent draws. By contrast, if we used a generic neural network to generate the conditional probability distributions itself, there is no such guarantee to this consistency.

## D ADDITIONAL INFORMATION ON THE EXPERIMENTS

### D.1 EXPERIMENTAL SETUP

Here we provide more detailed information describing how the numerical experiments were conducted.

All experiments share the same model architecture illustrated in Figure 3, there are however a number of model parameters that must be chosen. An L1 (i.e., Mean Absolute Error, MAE) loss function was used throughout. We use five bi-dimensional attention blocks, each consisting of multi-head self-attention blocks [35] containing a representation dimensionality of  $H = 64$  and 8 heads. Each experiment used either 500 or 1,000 diffusion steps, where we find larger values produce more accurate samples at the expense of computation time. Following Nichol et al. [24] we use a cosine-based scheduling of  $\beta_t$  during training. The Adam optimiser is used throughout. Our learning rate follows a cosine-decay function, with a 20 epochs linear learning rate warm-up to a maximum learning rate of  $\eta = 0.001$  before decaying. All NDP models were trained for 250 epochs with the exception of the lengthscale marginalisation experiment, which was trained for 500 epochs. Each epoch contained 4096 example training  $(\mathbf{y}_0, \mathbf{x}_0)$  pairs. Training data was provided in batches of 32, with each batch containing data with the same kernel hyperparameters but different realisations of prior GP samples. The complete configuration for each experiment is given in Table 2

Experiments were conducted on a 32-core machine and utilised a single Tesla V100-PCIE-32GB GPU. Training of each model used in the experiments takes no longer than 30 minutes.

**Time embedding** The diffusion step  $t$  is a crucial input of the neural network noise estimator as the model needs to be able to differentiate between noise added at the start or the end of the process.



Table 2: Experiment configuration and training time. Bold values indicate a deviation from the fiducial values used in the regression experiments.

Experiment	Regression	Hyperparameter marginalisation	Step	High dim BO	1D opt
Epochs	250	<b>500</b>	250	250	250
Total samples seen	1024k	<b>2048k</b>	1024k	1024k	1024k
Batch size	32	32	32	32	32
Loss	L1	L1	L1	L1	L1
LR decay	cosine	cosine	cosine	cosine	cosine
LR init	0.001	0.001	0.001	0.001	0.001
LR warmup epochs	20	20	20	20	20
Num blocks	5	5	5	5	5
Representation dim ( $H$ )	64	64	64	64	64
Num heads	8	8	8	8	8
Num timesteps (T)	500	<b>1000</b>	<b>1000</b>	500	<b>2000</b>
Num data per batch (N)	100	100	100	<b>256</b>	100
Deterministic inputs	True	True	True	True	<b>False</b>
Training time (mins)	17	33	16	21	16

Following Vaswani et al. [35] we use a cyclic 128-dimensional encoding vector for each step

$$t \mapsto [\sin(10^{\frac{0 \times 4}{63}} t), \sin(10^{\frac{1 \times 4}{63}} t), \dots, \sin(10^{\frac{64 \times 4}{63}} t), \cos(10^{\frac{0 \times 4}{63}} t), \dots, \cos(10^{\frac{64 \times 4}{63}} t)] \in \mathbb{R}^{128}$$

**Data** In our experiments, we provide training data to the models in the form of unique prior samples from a ground truth GP model. The input for each prior sample,  $\mathbf{x}_0$ , may be randomly sampled or deterministically spaced across the relevant input range. For all but the last experiment we deterministically space  $\mathbf{x}_0$  during training and inference in the hypercube  $[-1, 1]^D$ , where  $D$  is the dimensionality of the input. For  $D \leq 2$  the spacing of points is a simple linear spacing (a grid for  $D = 2$ ), in  $D > 2$  we sample  $\mathbf{x}_0$  according to a Halton sequence in the hypercube. The corresponding output  $\mathbf{y}_0$  is sampled from a GP prior  $\mathbf{y}_0 | \mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, k_\psi(\mathbf{x}_0, \mathbf{x}_0) + \sigma^2)$  where the kernel is a stationary kernel (Matern  $\frac{3}{2}$ , Matérn- $\frac{5}{2}$  or exponentiated quadratic, indicated in each experiment) with hyperparameters  $\psi = [\ell, \sigma_k^2]$  corresponding to the automatic relevance determination (ARD) lengthscales and kernel variance respectively. Kernel hyperparameters are fixed for all experiments except for the hyperparameter marginalisation experiment where we allow the model to observe data from a variety of different lengthscales during training. The noise variance is fixed to  $\sigma^2 = 10^{-6}$  and the kernel variance is fixed to  $\sigma_k^2 = 1.0$  throughout.

## D.2 ONE-DIMENSIONAL CONDITIONAL SAMPLES

In the experiment, we use GPflow [34] for the GP regression model using a kernel that matches the training data: a squared exponential with lengthscale set to 0.2. The other baseline, namely the Attentive Latent NP, is a pretrained model which was trained on a dataset with the same configuration. The ALNP model originates from the reference implementation of Dubois et al. [3]. Figure 8 shows additional conditional samples.

## D.3 HIGH-DIMENSIONAL BAYESIAN OPTIMISATION

In this experiment, we perform Bayesian optimisation on the Hartmann 3D and 6D objectives. We re-scale, without loss of generality, the inputs of the objectives such that the search space is  $[-1, 1]^D$  for both. We refer to <https://www.sfu.ca/~ssurjano/hart3.html> and <https://www.sfu.ca/~ssurjano/hart6.html> for more details about Hartmann 3D and 6D, respectively.

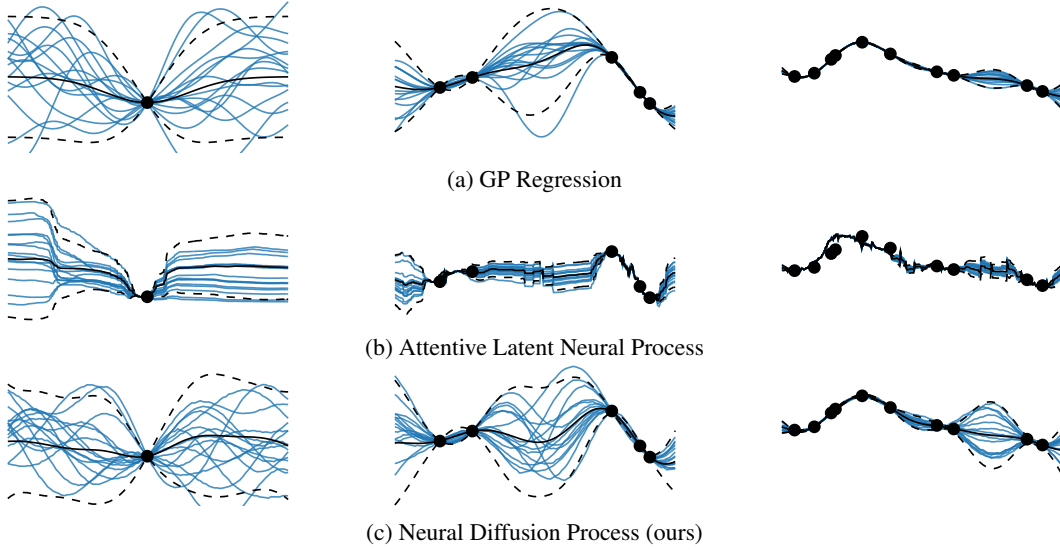


Figure 8: **1D regression:** The blue curves are posterior samples from different probabilistic models. We also plot the empirical mean and two standard deviations of the samples in black. From left to right we increase the number of data points (black dots) and notice how the process gets closer to the true underlying function.

The baseline models, GPR and Random, originate from Trieste [1] —a TensorFlow/GPflow based Bayesian Optimisation Python package. We compare them against two NDP models: Fixed and Marginalised. The Fixed NDP model is trained on Matérn-5/2 samples with a fixed lengthscale set to 0.5 along all dimensions. The Marginalised NDP model is trained on Matérn-5/2 samples originating from different lengthscales, where the lengthscales are drawn from a log-Normal prior  $\log \mathcal{N}(0, 1)$ .

#### D.4 COMPARISON TO NEURAL PROCESSES

Neural processes [8; 7] use an encoder and decoder architecture. The encoder is a neural network which operates on a dataset  $\mathcal{D} = (X, Y)$  to output a dataset representation. Using this representation, a decoder predicts the function output at a test location. For a NP this is given by

$$p(\mathbf{y}^* | \mathbf{x}^*, \mathcal{D}) = \prod_n \mathcal{N}(y_n^* | \text{dec}_\mu(x_n^*, \text{enc}(\mathcal{D})), \text{dec}_{\sigma^2}(x_n^*, \text{enc}(\mathcal{D})))$$

A latent NP works similarly but parameterises a latent variable  $z$ , over which we marginalise to make a prediction:

$$p(\mathbf{y}^* | \mathbf{x}^*, \mathcal{D}) = \int_z \prod_n \mathcal{N}(y_n^* | \text{dec}_\mu(x_n^*, z), \text{dec}_{\sigma^2}(x_n^*, z)) p(z | \text{enc}(\mathcal{D})) dz.$$

An attentive Neural processes (ANPs) [12] uses an attention mechanism [35] to parameterise the encoder and decoder networks. While ANPs perform better than plain (L)NPs, they suffer from two key weaknesses. Firstly, their predictions result in jittery functions as a result of shifting attention patterns. This behaviour is well-known and can also be observed in the top row of Fig. 10 in the supplementary. Secondly, ANP do not encode dimensionality invariance in their architecture, which we show is critical for tabular data applications.

Figure 10 shows example fits for all the different models. We used <https://github.com/wes-selb/neuralprocesses> for the NP baselines and evaluation framework.

#### D.5 GLOBAL OPTIMISATION

### E CODE

Full source code will be made publicly available upon acceptance.

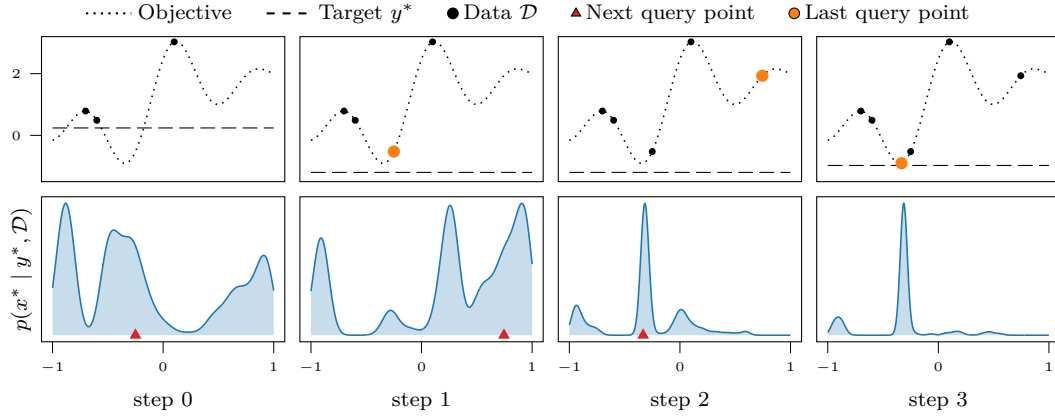


Figure 9: Global optimisation by diffusing the input locations. We condition the model at each step on a target  $y^*$  value (dashed line), and use the NDP to sample from  $p(x^* | y^*, \mathcal{D})$ . The panels in the upper row illustrate the progression of the optimisation. Bottom row shows the distribution  $p(x^* | y^*, \mathcal{D})$  and the red triangles mark the selected query point.

#### PREPROCESSING METHOD

The NDP’s noise model starts by rearranging the inputs  $x_t \in \mathbb{R}^{N \times D}$  and  $y_t \in \mathbb{R}^N$  to be of shape  $[N, D, 2]$ , as highlighted in the red box in Figs. 3 and 7. Below we show the code which is responsible for this preprocessing:

```
import tensorflow as tf

def preprocess(x: tf.Tensor, y: tf.Tensor) -> tf.Tensor:
    """
    Transform inputs to split out the x dimensions for dimension-
    agnostic processing.

    :param x: [B, N, D]
    :param y: [B, N, 1]
    :return: [B, N, D, 2]
    """
    D = tf.shape(x)[-1]
    x = tf.expand_dims(x, axis=-1)
    y = tf.repeat(tf.expand_dims(y, axis=-1), D, axis=2)
    return tf.concat([x, y], axis=-1)
```

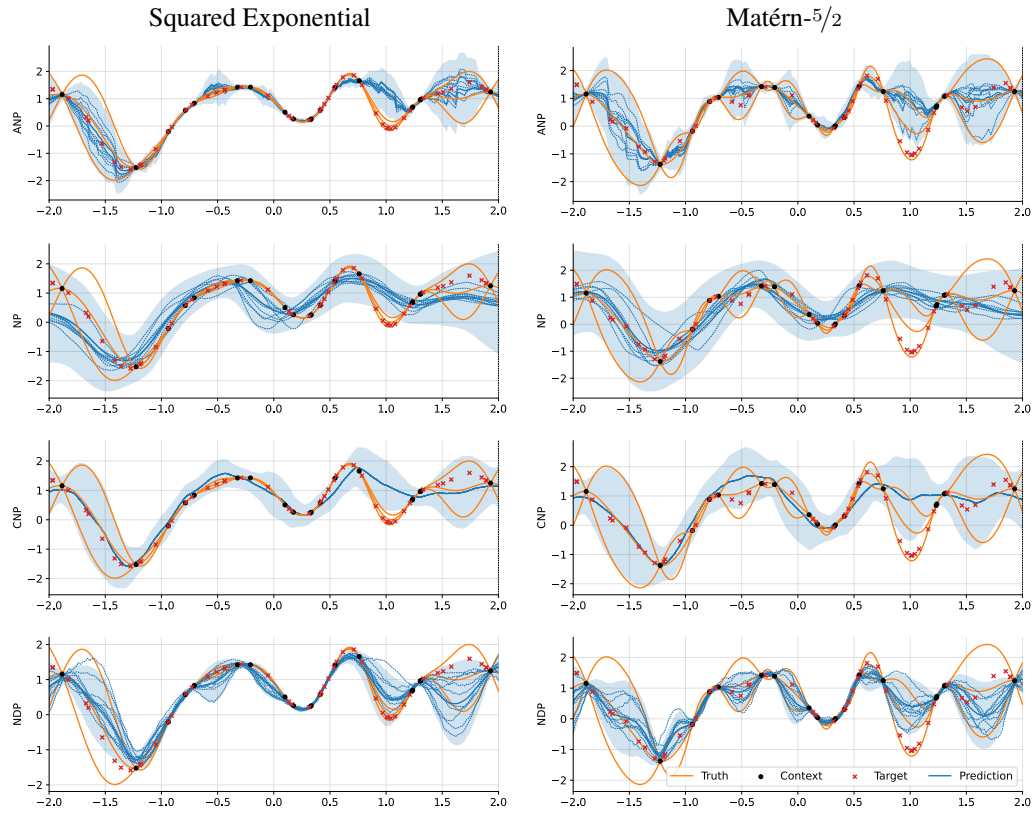


Figure 10: Model predictions on 1D datasets. Black dots: context points (conditioning points). Red crosses: Test/target points. Blue lines: samples of the model, mean and 2 std. dev. Orange: mean and 2 std dev from the true underlying GP model.