

---

# A framework for conditional diffusion modelling with applications in motif scaffolding for protein design

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Many protein design applications, such as binder or enzyme design, require scaffold-  
2 folding a structural motif with high precision. Generative modelling paradigms  
3 based on denoising diffusion processes emerged as a leading candidate to address  
4 this *motif scaffolding* problem and have shown early experimental success in some  
5 cases. In the diffusion paradigm, motif scaffolding is treated as a conditional  
6 generation task, and several conditional generation protocols were proposed or  
7 imported from the Computer Vision literature. However, most of these protocols  
8 are motivated heuristically, e.g. via analogies to Langevin dynamics, and lack a  
9 unifying framework, obscuring connections between the different approaches. In  
10 this work, we unify conditional training and conditional sampling procedures under  
11 one common framework based on the mathematically well-understood *Doob's*  
12 *h-transform*. This new perspective allows us to draw connections between existing  
13 methods and propose a new conditional training protocol. We illustrate the effective-  
14 ness of this new protocol in both, *image outpainting* and *motif scaffolding* and  
15 find that it outperforms standard methods.

## 16 1 Introduction

17 Denoising diffusion models are a powerful class of generative models where noise is gradually  
18 added to data samples until they converge to pure noise. The time-reversal of this noising process  
19 then allows to transform noise into samples. This process has been widely successful in generating  
20 high-quality images [Ho et al., 2020] and has more recently shown promise in designing protein  
21 backbones that were validated in experimental protein design workflows [Watson et al., 2023].

22 Importantly for protein design, diffusion models allow to subject this time-reversed sampling process  
23 to a target condition. For proteins, a key condition is the inclusion of a structural motif that grants the  
24 protein a particular function, such as binding to a specific target or forming an enzyme active site.  
25 However, for these motifs to be foldable and stable, they often need to be integrated into a larger  
26 protein structure. While there have been notable successes in scaffolding some motifs experimentally,  
27 many still prove challenging to scaffold [Watson et al., 2023]. This makes the development of  
28 better conditional generation methods for diffusion models an active area of research, with several  
29 contributions from the computer vision, molecular and protein design communities in recent times.

30 For instance, several methods cast the conditional sampling problem as an inverse (posterior sampling)  
31 problem and propose adding a *guidance* term to the time-reversal's drift (Fig. 1c) [e.g. Ho et al., 2022,  
32 Chung et al., 2022a]. Another line of work, focusing on 'inpainting', suggests *replacing* the observed  
33 variable in the diffused state (Fig. 1b) [e.g. Song et al., 2021c, Dutordoir et al., 2023, Mathieu et al.,  
34 2023]. Yet other work performs heuristic conditional training with the target variables in place  
35 [Watson et al., 2023, Torge et al., 2023].

METHOD	STAGE	CONSTRAINT			FRAMEWORK
		OPERATOR Leveraged	Soft	Hard	
Amortised $h$ -transform (ours)	Training	✓	✓	✓	Amortised trained $h$
Classifier free [Ho and Salimans, 2022]	Training	×	×	✓	Amortised trained $h$
Replacement [Song et al., 2021b]	Sampling	✓	×	✓	?
w/ particles: SMCDiff [Trippe et al., 2022]	Sampling	✓	✓	✓	?
RFDiffusion [Watson et al., 2023]	Training	✓	×	✓	Marginal of $h$
Classifier guidance [Dhariwal and Nichol, 2021]	Finetuning	×	×	✓	Trained separate $p(\mathbf{y} \mathbf{H}_t)$
Reconstruction guidance [Chung et al., 2022a,b]	Sampling	✓	✓	✓	Moment matching $h$
w/ particles: TDS [Wu et al., 2023]	Sampling	✓	✓	✓	Moment matching $h$

Table 1: Taxonomy of conditional methods. **STAGE** indicates when the conditional information is acquired. **OPERATOR** indicates whether the measurement operator  $\mathcal{A}$  is assumed to be known and thus leveraged by methods. **CONSTRAINT** classifies the likelihood as either *hard* or *soft*, as detailed in the main text. **FRAMEWORK** specifies the mechanism by which conditioning is accomplished. The ‘?’ means that it is unclear how this method fits into the  $h$ -transform framework.

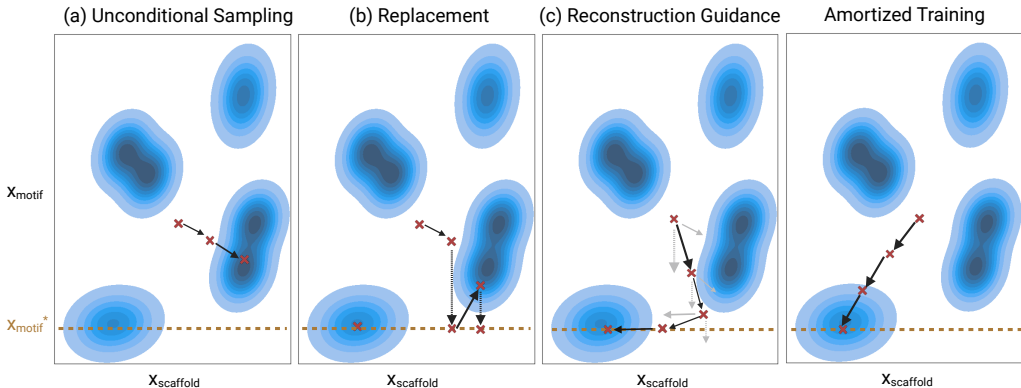


Figure 1: Schematic illustration of several common approaches to (conditionally) sample from a diffusion model. The sampling space is partitioned into motif coordinates (vertical) and scaffold coordinates (horizontal). The target motif is marked as  $x_{\text{motif}}^*$  and regions with plausible scaffolds are illustrated as blue blobs. A clear definition of each approach as pseudo-algorithm is given in App. B.

36 In this work, we reinterpret the conditioning problem leveraging Doob’s  $h$ -transform. This new  
37 perspective provides theoretical backing to existing approaches and naturally leads us to propose a  
38 novel method, which we call *amortised training* (Fig. 1d, Alg. 5). We highlight the straightforward  
39 implementation and practical use of our theoretical framework by applying it to problems, first as  
40 a proof of concept in image generation. We then study the merits and shortcomings of our newly  
41 proposed *amortised training* method in more detail for the *motif scaffolding* problem in protein  
42 design. We do so by comparing an amortised training implementation of the small-scale diffusion  
43 model Genie [Lin and AlQuraishi, 2023] on the RFDiffusion benchmark as well as a newly proposed  
44 benchmark dataset based on the SCOPE classification [Chandonia et al., 2022].

45 Our main contributions are as follows: *i*) We derive a formal framework for conditioning diffusion  
46 processes using Doob’s  $h$ -transform (Sec. 2). *ii*) We use our framework to create a taxonomy of  
47 existing methods (Table 1). *iii*) Our taxonomy elucidated the absence of a specific method within  
48 the current literature, prompting us to develop and implement this novel approach (Sec. 2.3). *iv*) We  
49 empirically assess these different approaches on image generation and protein design (Sec. 3).  
50 *v*) Finally, we present plug-and-play algorithms to implement various conditioning schemes (App. B).

## 51 2 Theory: Conditioning diffusions via the $h$ -transform, a new perspective

52 We first show how Doob’s  $h$ -transform enables diffusion models to satisfy *hard* equality constraints  
53 and then generalise this result to handle *soft* constraints in the context of noisy observations.

54 **2.1 Doob’s  $h$ -transform with hard constraint**

55 Doob’s transform provides a formal mechanism for conditioning a stochastic differential equation  
56 (SDE) to hit an event at a given time. Formally:

**Proposition 2.1.** (Doob’s  $h$ -transform *Rogers and Williams [2000]*) Consider the reverse SDE:

$$d\mathbf{X}_t = \bar{b}_t(\mathbf{X}_t) dt + \sigma_t \overline{d\mathbf{W}}_t, \quad \mathbf{X}_T \sim \mathcal{P}_T, \quad (1)$$

where time flows backwards and with transition densities  $\bar{p}_{t|s}$ . It then follows that the conditioned process  $\mathbf{X}_t | \mathbf{X}_0 \in B$  is a solution of

$$d\mathbf{H}_t = \left( \bar{b}_t(\mathbf{H}_t) - \sigma_t^2 \nabla_{\mathbf{H}_t} \ln \bar{P}_{0|t}(\mathbf{X}_0 \in B | \mathbf{H}_t) \right) dt + \sigma_t \overline{d\mathbf{W}}_t, \quad \mathbf{X}_T \sim \mathcal{P}_T, \quad (2)$$

such that  $\text{Law}(\mathbf{H}_s | \mathbf{H}_t) = \bar{p}_{s|t,0}(\mathbf{h}_s | \mathbf{h}_t, \mathbf{x}_0 \in B)$  and  $\mathbb{P}(\mathbf{H}_0 \in B) = 1$ .

57 This says that by conditioning a diffusion process to hit a particular event  $\mathbf{X}_0 \in B$  at a boundary  
58 time (e.g.  $t = 0$ ), the resulting conditional process is itself an SDE with an additional drift term.  
59 Furthermore, the resulting SDE will hit the specified event within a finite time  $T$ . The function  
60  $h(t, \mathbf{H}_t) \triangleq \bar{P}_{0|t}(\mathbf{X}_0 \in B | \mathbf{H}_t)$  is referred to as the  $h$ -transform [Rogers and Williams, 2000,  
61 De Bortoli et al., 2021a]. The  $h$ -transform drift decomposes into two terms via Bayes rule, a  
62 conditional and a prior score:

$$\nabla_{\mathbf{H}_t} \ln \bar{P}_{0|t}(\mathbf{X}_0 \in B | \mathbf{H}_t) = \nabla_{\mathbf{H}_t} \ln \bar{P}_{t|0}(\mathbf{H}_t | \mathbf{X}_0 \in B) - \nabla_{\mathbf{H}_t} \ln P_t(\mathbf{H}_t), \quad (3)$$

63 whereby the conditional score ensures that the event is hit at the specified boundary time, while the  
64 prior score ensures it is time-reversal of the correct forward process [De Bortoli et al., 2021a] (see  
65 App. A.3).

66 **Hard constraint** We now consider events of the form  $\mathbf{X}_0 \in B$  which are described by an equality  
67 constraint  $\mathcal{A}(\mathbf{X}_0) = \mathbf{y}$  with  $\mathcal{A}$  a known *measurement* (or *forward*) operator and  $\mathbf{y}$  an observation.  
68 We will see concrete examples of  $\mathcal{A}$  in Sec. 3.

**Corollary 2.2.** Consider the reverse SDE (1), then it follows that

$$d\mathbf{H}_t = \left( \bar{b}_t(\mathbf{H}_t) - \sigma_t^2 \nabla_{\mathbf{H}_t} \ln \bar{P}_{0|t}(\mathcal{A}(\mathbf{X}_0) = \mathbf{y} | \mathbf{H}_t) \right) dt + \sigma_t \overline{d\mathbf{W}}_t, \quad (4)$$

satisfies  $\text{Law}(\mathbf{H}_s | \mathbf{H}_t) = \text{Law}(\mathbf{X}_s | \mathbf{X}_t, \mathcal{A}(\mathbf{X}_0) = \mathbf{y})$  thus  $\text{Law}(\mathbf{H}_0) = \text{Law}(\mathbf{X}_0 | \mathcal{A}(\mathbf{X}_0) = \mathbf{y})$ .

69 Sampling (4) directly provides samples  $\mathbf{x} \sim p_{\text{data}}$  which also satisfy  $\mathcal{A}(\mathbf{x}) = \mathbf{y}$ . Crucially, this SDE  
70 is guaranteed to hit the conditioning in finite time, unlike prior equilibrium-motivated approaches  
71 [Chung et al., 2022a, Meng and Kabashima, 2022, Finzi et al., 2023, Song et al., 2022, Han et al.,  
72 2022, Dutordoir et al., 2023].

73 **Reconstruction guidance** To get better insight into the challenge of sampling from Doob’s  $h$ -  
74 transform (4) let us re-express the  $h$ -transform as

$$\bar{P}_{0|t}(\mathcal{A}(\mathbf{X}_0) = \mathbf{y} | \mathbf{H}_t) = \int \mathbb{1}_{\mathcal{A}(\mathbf{x}_0) = \mathbf{y}}(\mathbf{x}_0) \bar{p}_{0|t}(\mathbf{x}_0 | \mathbf{H}_t) d\mathbf{x}_0 \quad (5)$$

75 where in the case of denoising diffusion models  $\bar{p}_{0|t}(\mathbf{x}_0 | \cdot)$  is the transition density of the reverse SDE  
76 (1). In practice, one does not have access to this transition density – i.e. we can sample from this  
77 distribution, but we cannot easily get its value at a certain point. This makes it difficult to approximate  
78 the integral. To alleviate this, a strand of recent works [Finzi et al., 2023, Song et al., 2022, Rozet  
79 and Louppe, 2023] have proposed Gaussian approximation of  $\bar{p}_{0|t}(\mathbf{x}_0 | \cdot) \approx \mathcal{N}(\mathbf{x}_0 | \mathbb{E}[\mathbf{X}_0 | \mathbf{X}_t =$   
80  $\cdot], \Gamma_t)$  leveraging Tweedie’s formula and the already trained score network. This line of work is  
81 referred as *reconstruction guidance*. We note that whilst proposing to approximate the quantity  
82  $\bar{P}_{0|t}(\mathcal{A}(\mathbf{X}_0) = \mathbf{y} | \cdot)$ , they do not make the connection to Doob’s transform and thus are unable to  
83 provide guarantees on the conditional sampling that Cor. 2.2 provides. Overall, the Gaussian-based  
84 approximations of Doob’s  $h$ -transform lead to reconstruction guidance-based approaches [Finzi  
85 et al., 2023, Rozet and Louppe, 2023, Chung et al., 2022a, Han et al., 2022, Song et al., 2022]  
86  $d\mathbf{H}_t = \left( \bar{b}_t(\mathbf{H}_t) + \sigma_t^2 \nabla_{\mathbf{H}_t} \|\mathbf{y} - \mathbb{A}\mathbb{E}[\mathbf{X}_0 | \mathbf{X}_t = \mathbf{H}_t]\|_{\Gamma_t}^2 \right) dt + \sigma_t \overline{d\mathbf{W}}_t$ ,  $\mathbf{X}_T \sim \mathcal{P}_T$ , where  $\Gamma_t$  acts  
87 as a guidance scale [Simon V et al., 2023, Rozet and Louppe, 2023], and  $\mathbb{A}$  is a matrix if  $\mathcal{A}$  is linear  
88 otherwise  $\mathbb{A} = d\mathcal{A}(\mathbb{E}[\mathbf{X}_0 | \mathbf{X}_t = \mathbf{H}_t])$ .

89 **2.2 Generalised  $h$ -transform for soft constraints**

90 In the previous Sec. 2.1, we showed how the  $h$ -transform allows for conditioning on *hard* constraints,  
 91 correcting the reverse process to satisfy some observation  $P(\mathbf{y}|\mathbf{x}_0) \propto \mathbb{1}_{\mathcal{A}(\mathbf{X}_0)=\mathbf{y}}(\mathbf{x}_0)$ . Yet, many sce-  
 92 narios deal with *soft* constraints, modelling noisy observation  $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \eta$  with a density  $p(\mathbf{y}|\mathbf{x}_0)$ ,  
 93 typically with the goal of sampling from the posterior  $p(\mathbf{x}_0|\mathbf{Y} = \mathbf{y}) = p(\mathbf{y}|\mathbf{x}_0)p_{\text{data}}(\mathbf{x}_0)/p(\mathbf{y})$  as  
 94 in noisy inverse problems [Song et al., 2021a, Chung et al., 2022a,b]. In this section, we present a  
 95 generalisation of the  $h$ -transform applicable to denoising diffusion models that build on results in  
 96 [Vargas et al., 2023]:

**Proposition 2.3.** (*Noisy conditioning*) Given the following forward SDE:

$$d\mathbf{X}_t = f_t(\mathbf{X}_t) dt + \sigma_t \overline{\mathbf{W}}_t, \quad \mathbf{X}_0 \sim \mathcal{P}_{\text{data}} \quad (6)$$

it follows that the following reverse SDE with marginals  $p_t$

$$\begin{aligned} \mathbf{H}_T &\sim \text{Law}(\mathbf{X}_T|\mathbf{X}_0) \\ d\mathbf{H}_t &= (f_t(\mathbf{H}_t) + \sigma_t^2(\nabla_{\mathbf{H}_t} \ln p_t(\mathbf{H}_t) + \nabla_{\mathbf{H}_t} \ln p_{y|t}(\mathbf{Y} = \mathbf{y}|\mathbf{H}_t))) dt + \sigma_t d\overline{\mathbf{W}}_t, \end{aligned} \quad (7)$$

satisfies  $\text{Law}(\mathbf{H}_0) = p(\mathbf{x}_0|\mathbf{Y} = \mathbf{y})$  where  $p_{y|t}(\mathbf{Y} = \mathbf{y}|\cdot) = \int p(\mathbf{Y} = \mathbf{y}|\mathbf{x}_0)\overline{p}_{0|t}(\mathbf{x}_0|\cdot)d\mathbf{x}_0$ .

97 In short, the above results give a variant of the  $h$ -transform that allows to sample from noisy posteriors.  
 98 This provides theoretical backing to methodologies such as DPS [Chung et al., 2022a], in which the  
 99 SDE (8) is used to solve noisy inverse problems.

**Corollary 2.4.** *Furthermore, for an Ornstein-Uhlenbeck (OU) forward process, i.e. with drift*  
 $f_t(\mathbf{x}) = -\beta_t\mathbf{x}$  *and diffusion*  $\sigma_t = \sqrt{2\beta_t}$ , *we have that*

$$d\mathbf{H}_t = -\beta_t(\mathbf{H}_t + 2\nabla_{\mathbf{H}_t} \ln p_t(\mathbf{H}_t) + 2\nabla_{\mathbf{H}_t} \ln p_{y|t}(\mathbf{Y} = \mathbf{y}|\mathbf{H}_t))dt + \sqrt{2\beta_t} d\overline{\mathbf{W}}_t, \quad \mathbf{H}_T \sim \mathcal{N}(0, I) \quad (8)$$

satisfies  $\text{Law}(\mathbf{H}_0) \approx p(\mathbf{x}_0|\mathbf{Y} = \mathbf{y})$ . As such,  $\mathbf{H}_T$  inherits the rapid convergence guarantees of the  
 OU process [De Bortoli, 2022, De Bortoli et al., 2021b], in particular  $\|\text{Law}(\mathbf{H}_T) - \mathcal{N}(0, I)\|_{\text{TV}} \leq$   
 $\mathcal{O}(e^{-T/\bar{\beta}})$  for some  $\bar{\beta} > 0$ .

100 **2.3 Amortised training of  $h$ -transform**

101 In this section, we propose an objective for learning Doob’s  $h$ -transform at training time in an  
 102 amortised fashion. Note that since  $\overline{P}_{0|t}(\mathcal{A}(\mathbf{X}_0) = \mathbf{y}|\mathbf{X}_t = \mathbf{h}) = \overline{P}_{t|0}(\mathbf{h}|\mathcal{A}(\mathbf{X}_0) = \mathbf{y})p_0(\mathcal{A}(\mathbf{X}_0) =$   
 103  $\mathbf{y})/p_t(\mathbf{X}_t = \mathbf{h})$  we can re-express the Doob’s transformed SDE of a reversed OU process as:

$$d\mathbf{H}_t = -\beta_t \left( \mathbf{H}_t + 2\nabla_{\mathbf{H}_t} \ln \overline{P}_{t|0}(\mathbf{H}_t|\mathcal{A}(\mathbf{X}_0) = \mathbf{y}) \right) dt + \sqrt{2\beta_t} d\overline{\mathbf{W}}_t, \quad \mathbf{H}_T \sim \text{Law}(\mathbf{X}_T).$$

104

**Proposition 2.5.** *The minimiser of*

$$f^* = \arg \min_f \mathbb{E}_{\mathbf{Y} \sim p|\mathcal{A}, \mathbf{X}_0, \mathcal{A} \sim p, \mathbf{X}_0 \sim p_{\text{data}}} \left[ \int_0^T \|f(t, \mathbf{X}_t, \mathbf{Y}, \mathcal{A}) - \nabla_{\mathbf{X}_t} \ln \overline{p}_{t|0}(\mathbf{X}_t|\mathbf{X}_0)\|^2 dt \right] \quad (9)$$

is given by the conditional score  $f_t^*(\mathbf{h}, \mathbf{y}, \mathcal{A}) = \nabla_{\mathbf{h}} \ln \overline{p}_{t|0}(\mathbf{h}|\mathbf{Y} = \mathbf{y})$ .

105 This is referred as *amortised* learning for conditional sampling, since practically the neural network  
 106 approximating the (conditional) score is amortised over  $\mathcal{A}$  and  $\mathbf{y}$ , instead of learning a separate  
 107 network for each condition. This approach is reminiscent of ‘classifier free guidance’ [Ho and  
 108 Salimans, 2022] where the score network is amortised over some auxiliary variable (e.g. as in text-to-  
 109 image models [Ramesh et al., 2021]), or of RFDiffusion [Watson et al., 2023] where proteins are  
 110 designed given a specific subset motif. Our framework is different to ‘classifier free guidance’ as  $\mathcal{A}$   
 111 is assumed to be known (e.g. an inpainting mask), and to RFDiffusion since the conditioning variable  
 112  $\mathbf{Y}$  being a subset of  $\mathbf{X}$ , is also being noised during training and denoised when sampling (see Alg. 5),  
 113 also note due to its formulation classifier guidance would be unable to noise a subset of  $\mathbf{X}$  (the motif)  
 114 as we do.

115 **2.4 Conceptual comparison with RFDiffusion**

116 As highlighted in Alg. 4 and in contrast to our approach, RFDiffusion [Watson et al., 2023] does not  
 117 noise the motif coordinates  $\mathbf{X}_0^{[M]}$  with the forward OU-Process, instead it directly aims to sample  
 118 from  $p(\mathbf{X}_t^{[\setminus M]}|\mathbf{X}_0^{[M]})$  and estimate this score while keeping the motif fixed.

119 We can relate this approach to our amortised learning of Doob’s  $h$ -transform, by noting that RF  
 120 diffusion can be understood as learning the marginal conditional score:

$$p(\mathbf{X}_t^{[\setminus M]}|\mathbf{X}_0^{[M]}) = \int \overbrace{p(\mathbf{X}_t|\mathbf{X}_0^{[M]})}^{\propto h(t, \mathbf{X}_t)p_t(\mathbf{X}_t)} d\mathbf{X}_t^{[M]}. \quad (10)$$

121 This can be viewed as RFDiffusion estimating a marginal counterpart of our amortised  $h$ -transform  
 122 approach. See Algs. 4 and 5 for more details on how these approaches differ in a pseudo-code  
 123 implementation.

124 **3 Experimental results**

125 To compare the various conditional generation methods, we first highlight our results from initial tests  
 126 in the image setting and then discuss the motif scaffolding problem in protein design in more detail.

127 **3.1 Conditional image generation.**

128 The task of ‘image outpainting’ mimics the motif scaffolding problem in protein design and amounts  
 129 to conditioning the diffusion model on a central patch of an image. The measurement model  
 130  $\mathcal{A} \in \{0, 1\}^{n \times d}$  will select  $n$  central pixels out of an image in  $\mathbb{R}^d$ . We consider noise-free conditions  
 131 (i.e. hard constraints). We focus on the CELEBA [Liu et al., 2015] and FLOWERS [Nilsback and  
 132 Zisserman, 2008] image datasets. We empirically evaluate the AMORTISED approach where the mask  
 133 is provided at training time as an extra channel, along with RECONSTRUCTION GUIDANCE (Alg. 8)  
 134 and REPLACEMENT (Alg. 9) methods for which the score network is trained without access to the  
 135 mask, and are then queried at sampling time. The quality of conditional samples is measured by the  
 136 mean squared error (MSE) and LPIPS perceptual metric [Zhang et al., 2018]. See App. D for further  
 137 details. We empirically observe from Table 3 that the AMORTISED approach slightly outperforms  
 138 sampling-based methods, which are on par with each other.

Condition	Amortised	R. Guidance	Replacement	Metric	AMORTISED	R. GUIDANCE	REPLACEMENT
<b>FLOWERS</b>				MSE (↓)	0.34±0.01	0.27±0.01	0.28±0.01
				LPIPS (↓)	0.25±0.00	0.29±0.01	0.33±0.01
<b>CELEBA</b>				MSE (↓)	0.26±0.01	0.30±0.01	0.34±0.00
				LPIPS (↓)	0.14±0.00	0.15±0.01	0.17±0.00

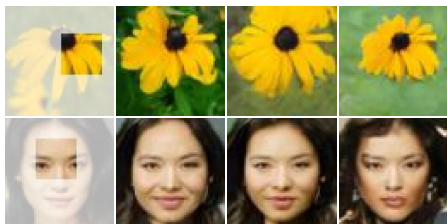


Figure 2: Some conditional samples.

Table 2: Quantitative assessment of conditional samples w.r.t to ground-truth.

139 **3.2 Conditional protein design: motif scaffolding**

140 The task of motif scaffolding in our protein setting amounts to sampling protein C alpha atom  
 141 coordinates  $\mathbf{x} \in \mathbb{R}^d$  such that it contains a given subset of C alpha coordinates  $\mathbf{y} \in \mathbb{R}^n$ , i.e.  
 142  $\mathbf{y} = \mathcal{A}(\mathbf{x}) = \mathbf{A}\mathbf{x}$ , where  $\mathcal{A} \in \{0, 1\}^{n \times d}$  is a masking matrix which selects  $n$  observed C alpha  
 143 coordinates. We perform two sets of motif scaffolding experiments. We firstly compare our proposed  
 144 AMORTISED approach to REPLACEMENT and RECONSTRUCTION GUIDANCE as we did in the image  
 145 case. Upon observing that AMORTISED performs significantly better, we then dive into a more  
 146 detailed analysis of this method on the RFDiffusion benchmark, as well as a new SCOPE based  
 147 benchmark that is created from a hierarchical structure and sequence-based split described below.

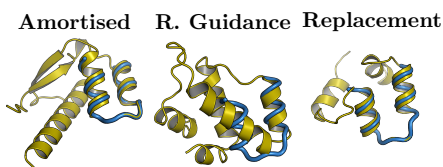


Figure 3: Conditional protein designs in yellow with target motif 3IXT in blue.

METRIC	AMORTISED	R. GUIDANCE	REPLACEMENT
% Success ( $\uparrow$ )	20.0	1.5	0.5
% scRMSD < 2 Å ( $\uparrow$ )	35.1	45.1	0.6
% mRMSD < 1 Å ( $\uparrow$ )	50.0	4.2	24.3

Table 3: RFDIFF benchmark metrics. Success: pAE < 5, scRMSD < 2Å, motifRMSD < 1Å, pLDDT > 70, scTM > 0.5. Details in Sec. 3.2.

148 **Data** We evaluate on the RFDiffusion motif benchmark [Watson et al., 2023] and on a self-curated  
 149 SCOPe benchmark based on a hierarchical structure-based split, jointly with a sequence-similarity  
 150 based split. For the RFDiffusion benchmark, we tested all sequence-contiguous motifs, resulting in  
 151 11 different motif design tasks. Our method readily extends to the non-contiguous motif setting and  
 152 future work will address this in more detail. The performance on each of these targets is depicted  
 153 in Fig. 4. For the SCOPe dataset, we leverage the hierarchical structure classification scheme of the  
 154 SCOPe database [Chandonia et al., 2022] to create train-test splits that allow us to investigate how  
 155 well the model can scaffold motifs from unseen folds, families and superfamilies and how difficult  
 156 these tasks are with respect to each other. In particular, for training, we hold out four clusters of  
 157 protein structures at the fold level, four at the family level and four at the superfamily level (Fig. 5a)  
 158 and evaluate the motif-scaffolding performance of the model on this structure-based hold-out set  
 159 (Fig. 5b-d).

160 **Diffusion process** We use a discrete-time DDPM [Ho et al., 2020] formulation for the diffusion  
 161 model with  $N = 1000$  steps and cosine  $\beta$ -schedule [Dhariwal and Nichol, 2021].

162 **Noise model** The denoising model  $\varepsilon_\theta$  is adapted from the Genie diffusion model [Lin and  
 163 AlQuraishi, 2023]. In Genie, the denoiser architecture consists of an SE(3)-invariant encoder  
 164 and an SE(3)-equivariant decoder. While the network uses Frenet-Serret frames as intermediate  
 165 representations, the diffusion process itself is defined in Euclidean space over the C alpha coordinates.  
 166 Similar to AlphaFold2, the denoiser network consists of a single representation track that is initialised  
 167 via a single feature network and a pair representation track that is initialised via a pair feature network.  
 168 These two representations are further transformed via a pair transform network and are used in the  
 169 decoder for noise prediction via IPA Jumper et al. [2021].

170 To evaluate unconditional sampling-based methods, we retrained the Genie denoising network for  
 171 4000 epochs on 4 A100 GPUs ( $\sim 300$  A100 hours in total). We stopped training at this point, as we  
 172 observed an almost comparable performance to the publicly available model weights (which were  
 173 obtained after training for 50'000 epochs).

174 To evaluate the AMORTISED approach (Alg. 5), we perform a minor modification to the unconditional  
 175 Genie model by adding an additional conditional pair feature network that takes the motif frames as  
 176 input with the ground truth coordinates for the motif and 0 as values for all other coordinates that  
 177 are not part of the motif. The output of this motif-conditional pair feature network is concatenated  
 178 with the output of the unconditional pair feature network to form an intermediate dimension of  
 179 twice the channel size compared to the unconditional model, before being linearly projected down  
 180 to the channel size of the unconditional model. From then onward the output is processed by the  
 181 remaining Genie components as in the unconditional model. The implementation is therefore similar  
 182 to the image case, where the motif features are presented as additional input and the model learns  
 183 to use these for reconstructing the motif. This minor alteration of the Genie architecture means our  
 184 amortised network has 4.162M parameters while the unconditional Genie networks have 4.087M  
 185 parameters ( $\sim 1.8\%$  fewer).

186 **Methods** In the amortised setting we follow the pseudo-code definition given in Alg. 5. In 80%  
 187 of the training steps, we pass a condition to the network. The other 20% contains an empty mask  
 188 consisting of only 0's. For the reconstruction guidance method (Alg. 8), we use a time-dependent  
 189 guidance term of  $\gamma_t = \alpha_t(1 - \alpha_t)$ .

190 **Metrics** We measure the performance of the methods across two axes: designability and success  
 191 rate.



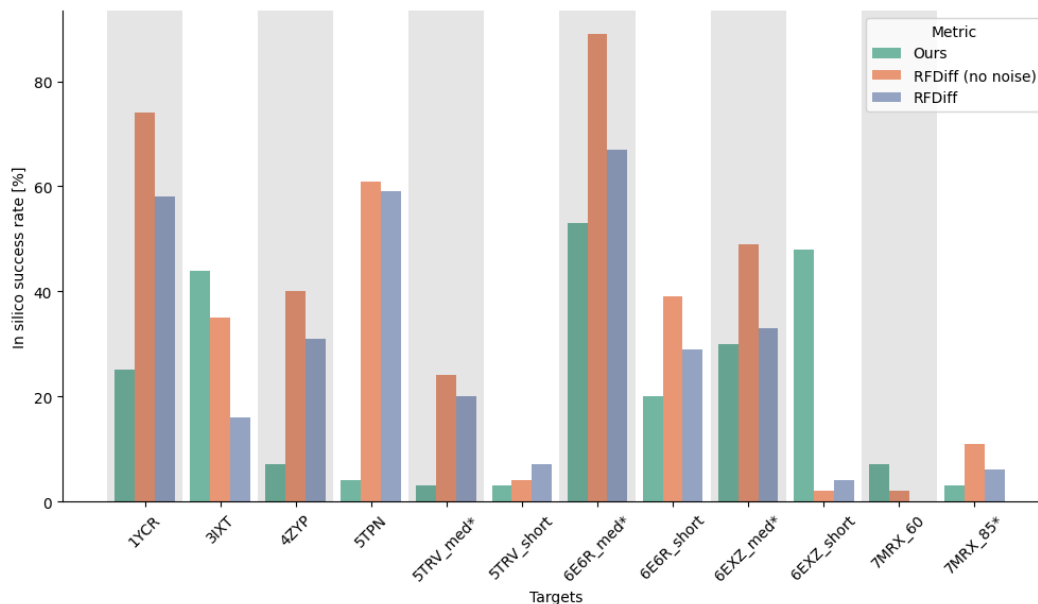


Figure 4: Comparison of our method to RFDiffusion for motif scaffolding for 12 continuous targets. Note that we trained our 4.1M parameter model for only 4000 epochs ( $\sim 300$  A100 hours in total), which is significantly less both in compute and parameter size than RFDiffusion ( $\sim 26'000+$  A100 hours, 59.8M parameters). For the motifs marked with \*, we had to shorten the sampled scaffold ranges on both sides of the motif from 0-65 (0-63 for TMRX80) to 0-50 since we trained our version of Genie only on protein generation up to a length of 128 residues. Performance numbers from RFDiffusion are taken from the original publication [Watson et al. \[2023\]](#) and our designs were created with the same design specifications as described there. We note that our folding step uses ESMFold instead of AlphaFold2, but we have future plans to use AlphaFold2 for a more direct comparison.

192 To assess whether a particular protein scaffold is *designable*, we run the same pipeline as [Lin and](#)  
 193 [AlQuraishi \[2023\]](#), consisting of an inverse folding generated  $C_\alpha$  backbones with ProteinMPNN and  
 194 then re-folding the designed sequences via ESMFold. The considered metrics and their corresponding  
 195 thresholds are the following:

- 196 •  $scTM > 0.5$ : This refers to the TM-score between the structure that's been designed and  
 197 the predicted structure based on self-consistency as previously described. The  $scTM$ -score  
 198 ranges from 0 to 1. Higher scores indicate a higher likelihood that the input structure can be  
 199 designed.
- 200 •  $scRMSD < 2 \text{ \AA}$ : The  $scRMSD$  metric is akin to the  $scTM$  metric. However, it uses the  
 201 RMSD (Root Mean Square Deviation) to measure the difference between the designed and  
 202 predicted structures, instead of the TM-score. This metric is more stringent than  $scTM$  as  
 203 RMSD, being a local metric, is more sensitive to minor structural variances.
- 204 •  $pLDDT > 70$  and  $pAE < 5$ : Both  $scTM$  and  $scRMSD$  metrics depend on a structure prediction  
 205 method like AlphaFold2 or ESMFold to be reliable. Hence, additional confidence metrics  
 206 such as  $pLDDT$  and  $pAE$  are employed to ascertain the reliability of the self-consistency  
 207 metrics.

208 In addition, we want to judge whether the motif scaffolding was successful or not. Therefore, similar  
 209 to previous work by [Watson et al. \[2023\]](#), we calculate the motifRMSD between the predicted design  
 210 structure and the original input motif and judge samples with  $< 1 \text{ \AA}$  motifRMSD as a successful motif  
 211 scaffold.

212 **Results** We evaluate all three approaches on the continuous motifs from the RFDIFFusion motif  
 213 benchmark [[Watson et al., 2023](#)]. For the AMORTISED approach we retrain the Genie model [[Lin and](#)  
 214 [AlQuraishi, 2023](#)] in an amortised fashion (Alg. 5), while for the R. GUIDANCE and REPLACEMENT

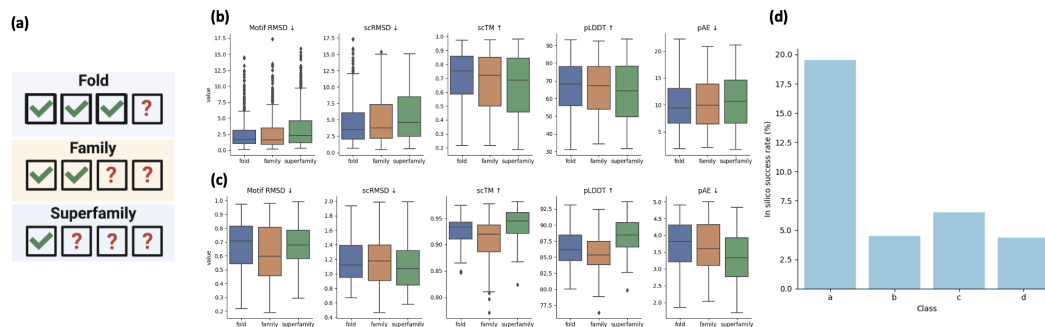


Figure 5: Data ablation study on a newly curated SCOPe benchmark dataset with our amortised training model. (a) We utilise the hierarchical structural clustering of SCOPe to create hold-out sets at three different levels of structural hierarchy: the fold, the family and the superfamily level. (b) We test the motif scaffolding performance on these splits and see decreasing scaffolding success for structurally dissimilar samples. (c) The same metrics as in (b), but only for samples that fulfill the definition of in silico success. (d) Scaffolding success by SCOPe class. Alpha helices can be scaffolded successfully, whereas other classes are more challenging.

215 methods we used the publicly available unconditional model. We observe that amortised training  
 216 outperforms the other approaches, especially replacement sampling (Fig. 3.2).

217 To better understand how well the AMORTISED conditioning approach works, we break down our  
 218 model performance on the different targets and compare it to the performance of RFDiffusion (Fig. 4).  
 219 Despite having trained a smaller model with fewer computing resources, we obtained competitive  
 220 performance on several targets.

221 We also ablate our model performance w.r.t. structural dissimilarity of the motif compared to  
 222 the training set via our previously described SCOPe benchmark. Testing the motif-scaffolding  
 223 performance of the amortised model on this data, we see that the scaffolding success decreases  
 224 from fold-over family to superfamily, indicating that scaffolding a motif from a protein that is  
 225 more different to the training set is harder (Fig. 5b-c). We also quantitatively observe an anecdotal  
 226 phenomenon in protein design: while alpha helices are relatively easy to scaffold, domains from other  
 227 classes have significantly lower success rates (Fig. 5d). We hope that this benchmark set will help to  
 228 address these issues in future modelling efforts.

## 229 4 Conclusion

230 We presented a unified framework, based on Doob’s  $h$ -transform, to better understand and clas-  
 231 sify different conditional diffusion methods. Based on the gained insights, we developed a novel  
 232 AMORTISED conditional sampling scheme (Alg. 5) which differs from existing approaches in that it  
 233 takes into account the measurement operator. For the motif scaffolding task this means we denoise  
 234 both the scaffold and the motif. We evaluated the AMORTISED approach on image outpainting  
 235 and motif scaffolding in protein design and outperform standard methods. We further investigated  
 236 the performance of the AMORTISED approach by comparing to RFDiffusion on contiguous motifs.  
 237 Surprisingly, our AMORTISED implementation of Genie [Lin and AlQuraishi, 2023] achieves notable  
 238 *in-silico* success rates of between 3 – 50% (as per the RFDiffusion definition) across the targets.  
 239 Though it lags behind RFDiffusion in 9/12 targets, it achieves this without low-temperature sampling,  
 240 a mere 10% of RFDiffusion’s parameter count, and after being trained for just 1.2% of the time. This  
 241 positions the AMORTISED approach as a promising candidate for further improving motif scaffolding,  
 242 potentially opening up new applications in protein engineering for drug discovery and enzyme design.



## 243 **References**

- 244 John-Marc Chandonia, Lindsey Guan, Shiangyi Lin, Changhua Yu, Naomi K Fox, and Steven E  
245 Brenner. Scope: improvements to the structural classification of proteins—extended database to  
246 facilitate variant interpretation and machine learning. *Nucleic acids research*, 50(D1):D553–D559,  
247 2022.
- 248 Hyungjin Chung, Jeongsol Kim, Michael T Mccann, Marc L Klasky, and Jong Chul Ye. Diffusion  
249 posterior sampling for general noisy inverse problems. *arXiv preprint arXiv:2209.14687*, 2022a.
- 250 Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Improving diffusion models for  
251 inverse problems using manifold constraints. *Advances in Neural Information Processing Systems*,  
252 35:25683–25696, 2022b.
- 253 Valentin De Bortoli. Convergence of denoising diffusion models under the manifold hypothesis.  
254 *arXiv preprint arXiv:2208.05314*, 2022.
- 255 Valentin De Bortoli, Arnaud Doucet, Jeremy Heng, and James Thornton. Simulating diffusion bridges  
256 with score matching. *arXiv preprint arXiv:2111.07243*, 2021a.
- 257 Valentin De Bortoli, James Thornton, Jeremy Heng, and Arnaud Doucet. Diffusion Schrödinger  
258 bridge with applications to score-based generative modeling. *Advances in Neural Information  
259 Processing Systems*, 34:17695–17709, 2021b.
- 260 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances  
261 in neural information processing systems*, 34:8780–8794, 2021.
- 262 Vincent Dutordoir, Alan Saul, Zoubin Ghahramani, and Fergus Simpson. Neural diffusion processes.  
263 In *International Conference on Machine Learning*, pages 8990–9012. PMLR, 2023.
- 264 Marc Anton Finzi, Anudhyan Boral, Andrew Gordon Wilson, Fei Sha, and Leonardo Zepeda-Núñez.  
265 User-defined event sampling and uncertainty quantification in diffusion models for physical  
266 dynamical systems. In *International Conference on Machine Learning*, pages 10136–10152.  
267 PMLR, 2023.
- 268 Xizewen Han, Huangjie Zheng, and Mingyuan Zhou. Card: Classification and regression diffusion  
269 models. *Advances in Neural Information Processing Systems*, 35:18100–18115, 2022.
- 270 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,  
271 2022.
- 272 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in  
273 Neural Information Processing Systems*, 33:6840–6851, 2020.
- 274 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J.  
275 Fleet. Video Diffusion Models, June 2022.
- 276 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,  
277 Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate  
278 protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- 279 Yeqing Lin and Mohammed AlQuraishi. Generating novel, designable, and diverse protein structures  
280 by equivariantly diffusing oriented residue clouds. *arXiv preprint arXiv:2301.12485*, 2023.
- 281 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In  
282 *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- 283 Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool.  
284 Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the  
285 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022.
- 286 Emile Mathieu, Vincent Dutordoir, Michael J Hutchinson, Valentin De Bortoli, Yee Whye Teh, and  
287 Richard E Turner. Geometric neural diffusion processes. *arXiv preprint arXiv:2307.05431*, 2023.

- 288 Xiangming Meng and Yoshiyuki Kabashima. Diffusion model based posterior sampling for noisy  
289 linear inverse problems. *arXiv preprint arXiv:2211.12343*, 2022.
- 290 Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number  
291 of classes. In *Indian Conference on Computer Vision, Graphics and Image Processing*, Dec 2008.
- 292 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark  
293 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In Marina Meila and Tong Zhang,  
294 editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of  
295 *Proceedings of Machine Learning Research*, pages 8821–8831. PMLR, 18–24 Jul 2021. URL  
296 <https://proceedings.mlr.press/v139/ramesh21a.html>.
- 297 L Chris G Rogers and David Williams. *Diffusions, Markov processes and martingales: Volume 2, Itô*  
298 *calculus*, volume 2. Cambridge university press, 2000.
- 299 Francois Rozet and Gilles Louppe. Score-based data assimilation. *arXiv preprint arXiv:2306.10574*,  
300 2023.
- 301 Mathis Simon V, Julia Komorowska Urszula, Jamnik Mateja, and Lio Pietro. Normal mode diffusion:  
302 Towards dynamics-informed protein design. *The 2023 ICML Workshop on Computational Biology.*  
303 *Baltimore, Maryland, USA, 2023. C*, 2023.
- 304 Jiaming Song, Arash Vahdat, Morteza Mardani, and Jan Kautz. Pseudoinverse-guided diffusion  
305 models for inverse problems. In *International Conference on Learning Representations*, 2022.
- 306 Yang Song, Liyue Shen, Lei Xing, and Stefano Ermon. Solving inverse problems in medical imaging  
307 with score-based generative models. *arXiv preprint arXiv:2111.08005*, 2021a.
- 308 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
309 Poole. Score-based generative modeling through stochastic differential equations. In *International*  
310 *Conference on Learning Representations*, 2021b. URL <https://openreview.net/forum?id=PXTIG12RRHS>.
- 312 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
313 Poole. Score-based generative modeling through stochastic differential equations. In *International*  
314 *Conference on Learning Representations*, 2021c. URL <https://openreview.net/forum?id=PXTIG12RRHS>.
- 316 Jos Torge, Charles Harris, Simon V. Mathis, and Pietro Lio. Diffhopp: A graph diffusion model for  
317 novel drug design via scaffold hopping. *arXiv preprint arXiv:2308.07416*, 2023.
- 318 Brian L Trippe, Jason Yim, Doug Tischer, David Baker, Tamara Broderick, Regina Barzilay, and  
319 Tommi Jaakkola. Diffusion probabilistic modeling of protein backbones in 3d for the motif-  
320 scaffolding problem. *arXiv preprint arXiv:2206.04119*, 2022.
- 321 Francisco Vargas, Andrius Ovsianas, David Fernandes, Mark Girolami, Neil D Lawrence, and Nikolas  
322 Nüsken. Bayesian learning via neural Schrödinger–Föllmer flows. *Statistics and Computing*, 33  
323 (1):3, 2023.
- 324 Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Trippe, Jason Yim, Helen E Eisenach,  
325 Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, et al. De novo design of protein  
326 structure and function with rfdiffusion. *Nature*, pages 1–3, 2023.
- 327 Luhuan Wu, Brian L Trippe, Christian A Naesseth, David M Blei, and John P Cunningham.  
328 Practical and asymptotically exact conditional sampling in diffusion models. *arXiv preprint*  
329 *arXiv:2306.17775*, 2023.
- 330 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable  
331 effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.

## 332 A Background on diffusion formulations

### 333 A.1 Continuous and discrete diffusion formulations

334 The discretised DDPM versions with various discrete time schedules amount to the time-dependent  
335 OU process

$$d\mathbf{X}_t = -\frac{\beta(t)}{2}\mathbf{x}_t dt + \sqrt{\beta(t)} d\overline{\mathbf{W}}_t \quad (11)$$

336 where choosing different time schedules amounts to choosing different functions  $\beta(t)$ . This process  
337 gives rise to the Green's function for transition probabilities

$$p(\mathbf{x}, t | \mathbf{x}_0, 0) = \vec{p}_{t|0}(\mathbf{x} | \mathbf{x}_0) \quad (12)$$

$$= \mathcal{N}\left(\mathbf{x}_0 e^{-\int_0^T \frac{\beta(s)}{2} ds}, \int_0^T \beta(t) e^{-\int_0^{T-t} \beta(s) ds} dt\right) \quad (13)$$

$$= \mathcal{N}\left(\mathbf{x}_0 e^{-\int_0^T \frac{\beta(s)}{2} ds}, \left(1 - e^{-\int_0^T \beta(s) ds}\right)\right). \quad (14)$$

338 With  $\bar{\alpha}(t) = e^{-\int_0^T \beta(s) ds}$ , this is the familiar form (Ho et al. [2020]):

$$p(\mathbf{x}, t | \mathbf{x}_0, 0) = \vec{p}_{t|0}(\mathbf{x} | \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_0 \sqrt{\bar{\alpha}(t)}, (1 - \bar{\alpha}(t))\right), \quad (15)$$

339 with  $\bar{\alpha}(t)$  time-dependent and we can therefore choose different functional forms for the noise  
340 schedule by either choosing the transition parameters  $\beta(t)$  or the cumulative parameters  $\alpha(t)$ .

341 If we define the noise schedule in terms of  $\beta(t)$ , the time-dependent OU process is immediately  
342 apparent (see (11)).

343 If we define the noise schedule in terms of  $\bar{\alpha}(t)$ , the mean and variance of the corresponding OU  
344 process can simply be obtained from

$$\beta(t) = -\frac{d}{dt} [\ln \bar{\alpha}(t)]. \quad (16)$$

### 345 A.2 Score, noise and mean diffusion formulations

346 The score-based model used for generation at inference time can be parametrised to model different  
347 quantities. The three most common one are the score, the noise and the mean.

348 When starting from the DDPM formulation of describing the diffusion process as a Gaussian linear  
349 Markov chain, it is natural to let the network predict the mean of this Gaussian, with the covariance  
350 being a fixed parameter:

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}(t)}} \boldsymbol{\varepsilon}_t \right) \quad (17)$$

351 However, we have access to the input  $\mathbf{x}_t$  at training time and can therefore reparameterize the  
352 Gaussian in order to make our network predict the noise  $\boldsymbol{\varepsilon}_t$  instead of the mean  $\mu_t$ :

$$\mathbf{x}_{t-1} = \mathcal{N}\left(\mathbf{x}_{t-1}; \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}(t)}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \right)\right) \quad (18)$$

353 When starting from the score-based SDE formulation, one can instead let the network predict the  
354 score term in order to minimise the following score matching loss:

$$\mathcal{L} = \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}_t} \left\| s_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \ln \vec{p}_{t|0}(\mathbf{x}_t | \mathbf{x}_0) \right\|^2 / \sigma_t^2 \quad (19)$$

355 **A.3 Doob’s  $h$ -transform intuition**

356 As mentioned before Doob’s  $h$ -transform adds a new drift to the SDE which amounts to two terms  
 357 (via Bayes Theorem), a conditional and an unconditional score:

$$\nabla \ln \bar{P}_{0|t}(\mathbf{X}_0 \in B|\cdot) = \nabla \ln \bar{P}_{t|0}(\cdot|\mathbf{X}_0 \in B) - \nabla \ln P_t(\cdot) \quad (20)$$

358 Interestingly, these two terms provide for a unique intuition: the Doob’s transform SDE is the time  
 359 reversal of the forward SDE corresponding to (1), that is the time reversal of the forward SDE

$$d\mathbf{X}_t = \bar{b}_t(\mathbf{X}_t) dt + \sigma_t d\bar{\mathbf{W}}_t, \quad \mathbf{X}_0 \sim \bar{P}_0(\cdot|\mathbf{X}_0 \in B), \quad (21)$$

360 coincides with the Doob transformed SDE (2) [De Bortoli et al., 2021a].

361 Thus we can view Doob’s transform as the following series of steps:

- 362 1. Time reverse the SDE we want to condition ((2) to (21)).
- 363 2. Impose the condition via ancestral sampling from the conditioned distribution/posterior.
- 364 3. Time reverse once more to be in the same time direction as we started.

365 **A.4 Examples**

366 **Truncated normal** Here for illustrative purposes we frame the problem of sampling from a  
 367 truncated normal distribution as simulating an SDE that is given by Doob’s  $h$ -transform.

368 Let’s remind that a 1d truncated normal distribution had a density  $p(x|a, b) \propto \mathbb{1}_{x \in (a, b)}(x) \mathcal{N}(x|\mu, \sigma^2)$ .  
 369 Now, let’s assume a data distribution  $p_0(x) = \mathcal{N}(x|\mu, \sigma^2)$  which is noised with an OU process (11).  
 370 Thus we have that  $p(x_0|x_t) = \mathcal{N}(x_0|\hat{\mu}_{0|t}(x_t), \hat{\sigma}_{0|t}(x_t)^2)$  is Gaussian, and so is  $p(x_t) = \mathcal{N}(x_t|\hat{\mu}_t, \hat{\sigma}_t^2)$ .  
 371 Let’s add the constraint that the process hit at time  $t = 0$  the event  $\mathbf{X}_0 \in (a, b)$ .

$$d\mathbf{H}_t = \beta(t) \left( \frac{\mathbf{H}_t}{2} + \nabla_{\mathbf{H}_t} \ln \bar{P}_t(\mathbf{H}_t) - \nabla_{\mathbf{H}_t} \ln \bar{P}_{0|t}(\mathbf{X}_0 \in (a, b) | \mathbf{H}_t) \right) dt + \sqrt{\beta(t)} d\bar{\mathbf{W}}_t, \quad (22)$$

372 We have that the  $h$ -transform is given by

$$\begin{aligned} h(t, \mathbf{H}_t) &= \bar{P}_{0|t}(\mathbf{X}_0 \in (a, b) | \mathbf{H}_t) = \int \mathbb{1}_{x \in (a, b)}(\mathbf{x}_0) \bar{p}_{0|t}(\mathbf{x}_0 | \mathbf{H}_t) d\mathbf{x}_0 \\ &= \int \mathbb{1}_{x \in (a, b)}(\mathbf{x}_0) \mathcal{N}(x | \hat{\mu}_{0|t}(\mathbf{H}_t), \hat{\sigma}_{0|t}(\mathbf{H}_t)^2) d\mathbf{x}_0 \\ &= \frac{1}{\hat{\sigma}_{0|t}(\mathbf{H}_t)} \frac{\phi\left(\frac{\mathbf{H}_t - \hat{\mu}_{0|t}(\mathbf{H}_t)}{\hat{\sigma}_{0|t}(\mathbf{H}_t)}\right)}{\Phi\left(\frac{b - \hat{\mu}_{0|t}(\mathbf{H}_t)}{\hat{\sigma}_{0|t}(\mathbf{H}_t)}\right) - \Phi\left(\frac{a - \hat{\mu}_{0|t}(\mathbf{H}_t)}{\hat{\sigma}_{0|t}(\mathbf{H}_t)}\right)} \end{aligned} \quad (23)$$

373 where  $\phi(\xi) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}\xi^2)$  is the pdf of a standard normal distribution,  $\Phi(\xi) =$   
 374  $\frac{1}{2} (1 + \operatorname{erf}(\xi/\sqrt{2}))$  its cumulative function. The corrective drift term due to the  $h$ -transform can then  
 375 be computed via autograd. The unconditional score term can be computed in closed form.

376 **B Algorithms**

377 In this section, we reformulate multiple algorithms from the literature under our common framework  
 378 as a reference for practitioners. In these algorithms, we use the following conventions: our dataset  
 379 is drawn from the law  $\mathcal{P}_{\text{data}}$ , but we can only sample from the simpler law  $\mathcal{P}_{\text{sampling}}$  at inference  
 380 time, which is often chosen as multivariate standard normal  $\mathcal{P}_{\text{sampling}} = \mathcal{N}(0, \mathbf{I})$ . Therefore, we  
 381 construct a forward noising process  $\mathcal{P}_{\text{data}} \rightarrow \mathcal{P}_{\text{sampling}}$  that is parametrised via the noise schedule  
 382  $\beta_t = \beta(t)$ ,  $\bar{\alpha}_t = \bar{\alpha}(t)$  and try to learn the reverse denoising process  $\mathcal{P}_{\text{sampling}} \rightarrow \mathcal{P}_{\text{data}}$ . Due to this  
 383 notion of "forward", and to keep consistency with the literature on denoising diffusion models, we  
 384 explicate the nomenclature  $\mathcal{P}_{\text{data}} = \mathcal{P}_0$  and  $\mathcal{P}_{\text{sampling}} = \mathcal{P}_T$ .

385 There is an additional law  $\mathcal{P}_{\text{noise}}$  that is sometimes confused with  $\mathcal{P}_{\text{sampling}}$  since in practice both are  
 386 often chosen as  $\mathcal{N}(0, \mathbf{I})$ , but they are two distinct laws that could in principle be different.  $\mathcal{P}_{\text{noise}}$  is  
 387 the law from which the noise added during the forward noising process as well as the during the  
 388 reverse diffusion process is drawn from.

389 For reference, we first reiterate the unconditional DDPM training and sampling algorithms, followed  
 390 by the various conditional methods. For each method, we highlight the differences to standard DDPM  
 391 sampling or training in gray boxes for clarity.

## 392 B.1 Unconditional algorithms

---

### Algorithm 1 | Unconditional training of denoising diffusion models [Ho et al., 2020]

---

**Require:** Dataset drawn from law  $\mathcal{P}_{\text{data}} = \mathcal{P}_0$  ▷ Dataset law  $\mathcal{P}_{\text{data}}$   
**Require:** Noise schedule  $\beta_t = \beta(t)$ ,  $\bar{\alpha}_t = \bar{\alpha}(t)$ , parametrising process  $\mathcal{P}_{\text{data}} \rightarrow \mathcal{P}_{\text{sampling}}$   
**Require:** Untrained noise predictor function  $\mathbf{f}_\theta(\mathbf{x}, t)$  with parameters  $\theta$

- 1: **repeat**
- 2:      $\mathbf{x}_0 \sim \mathcal{P}_0 = \mathcal{P}_{\text{data}}$
- 3:      $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4:     ▷ Forward noise sample,  $\mathbf{x}_t \sim \bar{p}_{t|0}(\mathbf{x}_0)$  ◁
- 5:      $\boldsymbol{\varepsilon}_t \sim \mathcal{P}_{\text{noise}}$  ▷ Often Brownian motion,  $\mathcal{P}_{\text{noise}} = \mathcal{N}(0, \mathbf{I})$
- 6:      $\mathbf{x}_t \leftarrow \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_t$
- 7:     ▷ Estimate noise of noised sample ◁
- 8:      $\hat{\boldsymbol{\varepsilon}}_\theta \leftarrow \mathbf{f}_\theta(\mathbf{x}_t, t)$
- 9:     Take gradient descent step on ▷ Typically, loss  $L(x_{\text{true}}, x_{\text{pred}}) = \|x_{\text{true}} - x_{\text{pred}}\|^2$   
        $\nabla_\theta L(\boldsymbol{\varepsilon}_t, \hat{\boldsymbol{\varepsilon}}_\theta)$
- 10: **until** converged or max epoch reached

---



---

### Algorithm 2 | Unconditional sampling with denoising diffusion models [Ho et al., 2020]

---

**Require:** Unconditionally trained noise predictor  $\mathbf{f}_\theta(\mathbf{x}_t, t)$   
**Require:** Noise schedule  $\beta_t = \beta(t)$ ,  $\bar{\alpha}_t = \bar{\alpha}(t)$ , parametrising process  $\mathcal{P}_{\text{data}} \rightarrow \mathcal{P}_{\text{sampling}}$

- 1: ▷ Sample a starting point  $\mathbf{x}_T$  ◁
- 2:  $\mathbf{x}_T \sim \mathcal{P}_T = \mathcal{P}_{\text{sampling}}$  ▷ Often  $\mathcal{P}_T = \mathcal{N}(0, \mathbf{I})$
- 3: ▷ Iteratively denoise for  $T$  steps ◁
- 4: **for**  $t$  in  $(T, T - 1, \dots, 1)$  **do** ◁
- 5:     ▷ Predict noise with learned network ◁
- 6:      $\hat{\boldsymbol{\varepsilon}}_\theta = \mathbf{f}_\theta(\mathbf{x}_t, t)$
- 7:     ▷ Denoise sample with learned reverse process  $\mathbf{x}_{t-1} \sim \bar{p}_{t-1|t}(\mathbf{x}_t)$  ◁
- 8:     ▷ Perform reverse drift ◁
- 9:      $\mathbf{x}_{t-1} \leftarrow \frac{1}{\sqrt{1 - \beta_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \hat{\boldsymbol{\varepsilon}}_\theta \right)$
- 10:     ▷ Perform reverse diffusion, which is often Brownian motion in  $\mathbb{R}^n$ , i.e.  $\mathcal{P}_{\text{noise}} = \mathcal{N}(0, \mathbf{I})$  ◁
- 11:      $\boldsymbol{\varepsilon}_t \sim \mathcal{P}_{\text{noise}}$  if  $t > 1$  else  $\boldsymbol{\varepsilon}_t \leftarrow 0$
- 12:      $\mathbf{x}_{t-1} \leftarrow \mathbf{x}_{t-1} + \sigma_t \boldsymbol{\varepsilon}_t$  ▷ A common choice is  $\sigma_t = \beta(t)$
- 13: **return**  $\mathbf{x}_0$

---

## B.2 Conditional training

### Algorithm 3 | Classifier-free conditional training [Ho and Salimans, 2022]

**Require:** Dataset drawn from  $\mathcal{P}_{\text{data}}$   $\triangleright$  Dataset law  $\mathcal{P}_{\text{data}}$  over data and auxiliary variable  
**Require:** Noise schedule  $\beta_t = \beta(t)$ ,  $\bar{\alpha}_t = \bar{\alpha}(t)$ , parametrising process  $\mathcal{P}_{\text{data}} \rightarrow \mathcal{P}_{\text{sampling}}$   
**Require:** Untrained noise predictor function  $\mathbf{f}_\theta(\mathbf{x}, t)$  with parameters  $\theta$

- 1: **repeat**
- 2:  $\mathbf{x}_0, \mathbf{y} \sim \mathcal{P}_0 = \mathcal{P}_{\text{data}}$
- 3:  $\boldsymbol{\varepsilon}_t \sim \mathcal{P}_{\text{noise}}$   $\triangleright$  Often Brownian motion,  $\mathcal{P}_{\text{noise}} = \mathcal{N}(0, \mathbf{I})$
- 4:  $t \sim \text{Uniform}(\{1, \dots, T\})$
- 5:  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_t$
- 6:  $\hat{\boldsymbol{\varepsilon}}_\theta = \mathbf{f}_\theta(\mathbf{x}_t, t, \mathbf{y})$
- 7: Take gradient descent step on  $\nabla_\theta L(\boldsymbol{\varepsilon}_t, \hat{\boldsymbol{\varepsilon}}_\theta)$   $\triangleright$  Typically,  $L(x_{\text{true}}, x_{\text{pred}}) = \|x_{\text{true}} - x_{\text{pred}}\|^2$
- 8: **until** converged or max epoch reached

### Algorithm 4 | RFDiffusion conditional training [Watson et al., 2023]

**Require:** Dataset drawn from  $\mathcal{P}_{\text{data}}$   $\triangleright$  Dataset law  $\mathcal{P}_{\text{data}}$   
**Require:** Noise schedule  $\beta_t = \beta(t)$ ,  $\bar{\alpha}_t = \bar{\alpha}(t)$ , parametrising process  $\mathcal{P}_{\text{data}} \rightarrow \mathcal{P}_{\text{sampling}}$   
**Require:** Untrained noise predictor function  $\mathbf{f}_\theta(\mathbf{x}, t, M)$  with parameters  $\theta$

- 1: **repeat**
- 2:  $\mathbf{x}_0 \sim \mathcal{P}_0 = \mathcal{P}_{\text{data}}$
- 3:  $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4:  $\mathbf{x}_0^{[M]} \cup \mathbf{x}_0^{[\setminus M]} \leftarrow \mathbf{x}_0$   $\triangleright$  Randomly partition data point into motif and rest
- 5:  $\triangleright$  Forward noise the non-motif rest via sampling from  $\bar{p}_{0|t}(\mathbf{x}_0)$   $\triangleleft$
- 6:  $\boldsymbol{\varepsilon}_t \sim \mathcal{P}_{\text{noise}}$
- 7:  $\mathbf{x}_t^{[\setminus M]} \leftarrow \sqrt{\bar{\alpha}_t} \mathbf{x}_0^{[\setminus M]} + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_t^{[\setminus M]}$
- 8:  $\triangleright$  Combine unnoised motif with noised rest and set timestep of motif part to 0  $\triangleleft$
- 9:  $\mathbf{x}_t \leftarrow \mathbf{x}_0^{[M]} \cup \mathbf{x}_t^{[\setminus M]}$
- 10:  $t^{[M]} \leftarrow 0$
- 11:  $\hat{\boldsymbol{\varepsilon}}_\theta \leftarrow \mathbf{f}_\theta(\mathbf{x}_t, t, M)$   $\triangleright$  Estimate noise of sample with noised rest
- 12: Take gradient descent step on  $\nabla_\theta L(\boldsymbol{\varepsilon}, \hat{\boldsymbol{\varepsilon}}_\theta)$   $\triangleright$  Typically,  $L(x_{\text{true}}, x_{\text{pred}}) = \|x_{\text{true}} - x_{\text{pred}}\|^2$
- 13: **until** converged or max epoch reached

### Algorithm 5 | Amortised training – i.e. Doob’s $h$ -transform conditional training (new)

**Require:** Dataset drawn from  $\mathcal{P}_{\text{data}}$   $\triangleright$  Dataset law  $\mathcal{P}_{\text{data}}$   
**Require:** Noise schedule  $\beta_t = \beta(t)$ ,  $\bar{\alpha}_t = \bar{\alpha}(t)$ , parametrising process  $\mathcal{P}_{\text{data}} \rightarrow \mathcal{P}_{\text{sampling}}$   
**Require:** Untrained noise predictor function  $\mathbf{f}_\theta(\mathbf{x}, t, \mathbf{x}_0^{[M]}, M)$  with parameters  $\theta$

- 1: **repeat**
- 2:  $\mathbf{x}_0 \sim \mathcal{P}_0 = \mathcal{P}_{\text{data}}$
- 3:  $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4:  $\mathbf{x}_0^{[M]} \cup \mathbf{x}_0^{[\setminus M]} \leftarrow \mathbf{x}_0$   $\triangleright$  Randomly partition data point into motif and rest
- 5:  $\triangleright$  Forward noise full sample via sampling from  $\bar{p}_{0|t}(\mathbf{x}_0)$   $\triangleleft$
- 6:  $\boldsymbol{\varepsilon}_t \sim \mathcal{P}_{\text{noise}}$
- 7:  $\mathbf{x}_t \leftarrow \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_t$
- 8:  $\triangleright$  Estimate noise of sample with original motif as additional input  $\triangleleft$
- 9:  $\hat{\boldsymbol{\varepsilon}}_\theta \leftarrow \mathbf{f}_\theta(\mathbf{x}_t, t, \mathbf{x}_0^{[M]}, M)$
- 10: Take gradient descent step on  $\nabla_\theta L(\boldsymbol{\varepsilon}, \hat{\boldsymbol{\varepsilon}}_\theta)$   $\triangleright$  Typically,  $L(x_{\text{true}}, x_{\text{pred}}) = \|x_{\text{true}} - x_{\text{pred}}\|^2$
- 11: **until** converged or max epoch reached



**Algorithm 6** | RFDiffusion conditional sampling [Watson et al., 2023]

---

**Require:** Conditionally trained noise predictor  $\mathbf{f}_\theta(\mathbf{x}, t, M)$

**Require:** Target motif/context  $\mathbf{x}_0^{[M]}$

**Require:** Noise schedule  $\beta_t = \beta(t), \bar{\alpha}_t = \bar{\alpha}(t)$ , parametrising process  $\mathcal{P}_{\text{data}} \rightarrow \mathcal{P}_{\text{sampling}}$

- 1:  $\triangleright$  Sample a starting point  $\mathbf{x}_T$   $\triangleleft$
- 2:  $\mathbf{x}_T \sim \mathcal{P}_T = \mathcal{P}_{\text{sampling}}$
- 3:  $\triangleright$  Iteratively denoise for  $T$  steps  $\triangleright$  Often  $\mathcal{P}_T = \mathcal{N}(0, \mathbf{I})$   $\triangleleft$
- 4: **for**  $t$  in  $(T, T-1, \dots, 1)$  **do**
- 5:  $\triangleright$  Overwrite motif variables with target motif and reset their time parameter  $\triangleleft$
- 6:  $\triangleright$  Note: Original RFDiffusion zero-centers  $\mathbf{x}_t$  and  $\mathbf{x}_0^{[M]}$  individually for equivariance.  $\triangleleft$
- 7:  $\mathbf{x}_t^{[M]} \leftarrow \mathbf{x}_0^{[M]}$   $\triangleright$  Set noisy motif to unnoised motif
- 8:  $t^{[M]} \leftarrow 0$   $\triangleright$  Set timesteps for motif to 0
- 9:  $\hat{\boldsymbol{\epsilon}}_\theta = \mathbf{f}_\theta(\mathbf{x}_t, t, M)$   $\triangleright$  Predict noise with learned network
- 10:  $\triangleright$  Denoise sample with learned reverse process  $\mathbf{x}_{t-1} \sim \bar{p}_{t-1|t}(\mathbf{x}_t)$   $\triangleleft$
- 11:  $\triangleright$  Perform reverse drift  $\triangleleft$
- 12: 
$$\mathbf{x}_{t-1} \leftarrow \frac{1}{\sqrt{1-\beta_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \hat{\boldsymbol{\epsilon}}_\theta \right)$$
- 13:  $\triangleright$  Perform reverse diffusion, which is often Brownian motion in  $\mathbb{R}^n$ , i.e.  $\mathcal{P}_{\text{noise}} = \mathcal{N}(0, \mathbf{I})$   $\triangleleft$
- 14:  $\boldsymbol{\epsilon}_t \sim \mathcal{P}_{\text{noise}}$  if  $t > 1$  else  $\boldsymbol{\epsilon}_t \leftarrow 0$
- 15:  $\mathbf{x}_{t-1} \leftarrow \mathbf{x}_{t-1} + \sigma_t \boldsymbol{\epsilon}_t$   $\triangleright$  A common choice is  $\sigma_t = \beta(t)$
- 16: **return**  $\mathbf{x}_0$

---

**Algorithm 7** | Replacement conditional sampling

---

**Require:** Unconditionally trained noise predictor  $\mathbf{f}_\theta(\mathbf{x}_t, t)$

**Require:** Noise schedule  $\beta_t = \beta(t), \bar{\alpha}_t = \bar{\alpha}(t)$ , parametrising process  $\mathcal{P}_{\text{data}} \rightarrow \mathcal{P}_{\text{sampling}}$

**Require:** Target motif  $\mathbf{x}_0^{[M]}$

- 1:  $\triangleright$  Sample a starting point  $\mathbf{x}_T$   $\triangleleft$
- 2:  $\mathbf{x}_T \sim \mathcal{P}_T = \mathcal{P}_{\text{sampling}}$
- 3:  $\triangleright$  Iteratively denoise for  $T$  steps  $\triangleright$  Often  $\mathcal{P}_T = \mathcal{N}(0, \mathbf{I})$   $\triangleleft$
- 4: **for**  $t$  in  $(T, T-1, \dots, 1)$  **do**
- 5:  $\triangleright$  Predict noise with learned network  $\triangleleft$
- 6:  $\hat{\boldsymbol{\epsilon}}_\theta \leftarrow \mathbf{f}_\theta(\mathbf{x}_t, t)$
- 7:  $\triangleright$  Denoise sample with learned reverse process  $\mathbf{x}_{t-1} \sim \bar{p}_{t-1|t}(\mathbf{x}_t)$   $\triangleleft$
- 8:  $\triangleright$  Perform reverse drift  $\triangleleft$
- 9: 
$$\mathbf{x}_{t-1} \leftarrow \frac{1}{\sqrt{1-\beta_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \hat{\boldsymbol{\epsilon}}_\theta \right)$$
- 10:  $\triangleright$  Perform reverse diffusion, which is often Brownian motion in  $\mathbb{R}^n$ , i.e.  $\mathcal{P}_{\text{noise}} = \mathcal{N}(0, \mathbf{I})$   $\triangleleft$
- 11:  $\boldsymbol{\epsilon}_t \sim \mathcal{P}_{\text{noise}}$  if  $t > 1$  else  $\boldsymbol{\epsilon}_t \leftarrow 0$
- 12:  $\mathbf{x}_{t-1} \leftarrow \mathbf{x}_{t-1} + \sigma_t \boldsymbol{\epsilon}_t$   $\triangleright$  A common choice is  $\sigma_t = \beta(t)$
- 13:  $\triangleright$  Forward noise the target motif  $\mathbf{x}_{t-1}^{[M]} \sim \bar{p}_{0|t-1}(\mathbf{x}_0^{[M]})$   $\triangleleft$
- 14:  $\boldsymbol{\eta}_{t-1} \sim \mathcal{P}_{\text{noise}}$  if  $t > 1$  else  $\boldsymbol{\eta}_{t-1} \leftarrow 0$
- 15:  $\mathbf{x}_{t-1}^{[M]} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0^{[M]} + \sqrt{1-\bar{\alpha}_{t-1}} \boldsymbol{\eta}_{t-1}$
- 16:  $\mathbf{x}_{t-1} \leftarrow \mathbf{x}_{t-1}^{[M]} \cup \mathbf{x}_{t-1}$   $\triangleright$  Insert noised motif into current sample
- 17: **return**  $\mathbf{x}_0$

---

---

**Algorithm 8** | Reconstruction Guidance (i.e. Moment Matching (MM) Approximation to  $h$ -transform)

---

**Require:** Unconditionally trained noise predictor  $\mathbf{f}_\theta(\mathbf{x}_t, t)$ , target motif/context  $\mathbf{x}_0^{[M]}$ .

**Require:** Noise schedule  $\beta_t = \beta(t)$ ,  $\bar{\alpha}_t = \bar{\alpha}(t)$ , parameterising process  $\mathcal{P}_{\text{data}} \rightarrow \mathcal{P}_{\text{sampling}}$

**Require:** Guidance scale (schedule)  $\gamma_t = \gamma(t)$

**Require:** Conditioning loss  $l(x_{\text{true}}, x_{\text{pred}})$ . e.g, Gaussian MM  $l(x_{\text{true}}, x_{\text{pred}}) = \|x_{\text{true}} - x_{\text{pred}}\|^2$

```

1: ▷ Sample a starting point  $\mathbf{x}_T$  ◁
2:  $\mathbf{x}_T \sim \mathcal{P}_T = \mathcal{P}_{\text{sampling}}$  ▷ Often  $\mathcal{P}_T = \mathcal{N}(0, \mathbf{I})$ 

3: ▷ Iteratively denoise and condition for  $T$  steps ◁
4: for  $t$  in  $(T, T - 1, \dots, 1)$  do
5:    $\hat{\boldsymbol{\epsilon}}_\theta = \mathbf{f}_\theta(\mathbf{x}_t, t)$  ▷ Predict noise with learned network
6:   ▷ Estimate current denoised estimate via Tweedie's formula ◁
7:    $\hat{\mathbf{x}}_0(\mathbf{x}_t, \hat{\boldsymbol{\epsilon}}_\theta) \leftarrow \frac{1}{\sqrt{\bar{\alpha}_t}}(\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\hat{\boldsymbol{\epsilon}}_\theta)$  ▷ c.f. also eq. 15 in Ho et al. [2020]
8:   ▷ Perform gradient descent step towards condition on motif dimensions  $M$  ◁
9:    $\mathbf{x}_t \leftarrow \mathbf{x}_t - \gamma_t \nabla_{\theta} l(\mathbf{x}_0^{[M]}, \hat{\mathbf{x}}_0^{[M]}(\mathbf{x}_t, \hat{\boldsymbol{\epsilon}}_\theta))$  ▷ Requires backprop through  $\mathbf{f}_\theta$ 
10:  ▷ Denoise sample with learned reverse process  $\mathbf{x}_{t-1} \sim \bar{p}_{t-1|t}(\mathbf{x}_t)$  ◁
11:   $\mathbf{x}_{t-1} \leftarrow (1 - \beta_t)^{-1/2}(\mathbf{x}_t - \beta_t(1 - \bar{\alpha}_t)^{-1/2}\hat{\boldsymbol{\epsilon}}_\theta)$  ▷ Perform reverse drift
12:  ▷ Perform reverse diffusion, which is often Brownian motion in  $\mathbb{R}^n$ , i.e.  $\mathcal{P}_{\text{noise}} = \mathcal{N}(0, \mathbf{I})$  ◁
13:   $\boldsymbol{\epsilon}_t \sim \mathcal{P}_{\text{noise}}$  if  $t > 1$  else  $\boldsymbol{\epsilon}_t \leftarrow 0$ 
14:   $\mathbf{x}_{t-1} \leftarrow \mathbf{x}_{t-1} + \sigma_t \boldsymbol{\epsilon}_t$  ▷ A common choice is  $\sigma_t = \beta(t)$ 
15: return  $\mathbf{x}_0$ 

```

---

**Algorithm 9** | Replacement conditional Sampling [Lugmayr et al., 2022]

---

**Require:** Unconditionally trained noise predictor  $\mathbf{f}_\theta(\mathbf{x}_t, t)$

**Require:** Noise schedule  $\beta_t = \beta(t)$ ,  $\bar{\alpha}_t = \bar{\alpha}(t)$ , parametrising process  $\mathcal{P}_{\text{data}} \rightarrow \mathcal{P}_{\text{sampling}}$

**Require:** Target motif  $\mathbf{x}_0^{[M]}$

```

1: ▷ Sample a starting point  $\mathbf{x}_T$  ◁
2:  $\mathbf{x}_T \sim \mathcal{P}_T = \mathcal{P}_{\text{sampling}}$ 

3: ▷ Iteratively denoise for  $T$  steps ▷ Often  $\mathcal{P}_T = \mathcal{N}(0, \mathbf{I})$  ◁
4: for  $t$  in  $(T, T - 1, \dots, 1)$  do ▷  $T$  time steps
5:   for  $r$  in  $1, \dots, R$  do ▷  $R$  repaint steps
6:     ▷ Predict noise with learned network ◁
7:      $\hat{\boldsymbol{\epsilon}}_\theta \leftarrow \mathbf{f}_\theta(\mathbf{x}_t, t)$ 
8:     ▷ Denoise sample with learned reverse process  $\mathbf{x}_{t-1} \sim \bar{p}_{t-1|t}(\mathbf{x}_t)$  ◁
9:      $\mathbf{x}_{t-1} \leftarrow (1 - \beta_t)^{-1/2}(\mathbf{x}_t - \beta_t(1 - \bar{\alpha}_t)^{-1/2}\hat{\boldsymbol{\epsilon}}_\theta)$  ▷ Perform reverse drift
10:    ▷ Perform reverse diffusion, often Brownian motion in  $\mathbb{R}^n$ , i.e.  $\mathcal{P}_{\text{noise}} = \mathcal{N}(0, \mathbf{I})$  ◁
11:     $\boldsymbol{\epsilon}_t \sim \mathcal{P}_{\text{noise}}$  if  $t > 1$  else  $\boldsymbol{\epsilon}_t \leftarrow 0$ 
12:     $\mathbf{x}_{t-1} \leftarrow \mathbf{x}_{t-1} + \sigma_t \boldsymbol{\epsilon}_t$  ▷ A common choice is  $\sigma_t = \beta(t)$ 
13:    ▷ Forward noise the target motif  $\mathbf{x}_{t-1}^{[M]} \sim \bar{p}_{0|t-1}(\mathbf{x}_0^{[M]})$  ◁
14:     $\boldsymbol{\eta}_{t-1} \sim \mathcal{P}_{\text{noise}}$  if  $t > 1$  else  $\boldsymbol{\eta}_{t-1} \leftarrow 0$ 
15:     $\mathbf{x}_{t-1}^{[M]} \leftarrow \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0^{[M]} + \sqrt{1 - \bar{\alpha}_{t-1}}\boldsymbol{\eta}_{t-1}$ 
16:     $\mathbf{x}_{t-1} \leftarrow \mathbf{x}_{t-1}^{[M]} \cup \mathbf{x}_{t-1}$  ▷ Insert noised motif into current sample
17:    if  $r < R$  and  $t > 1$  then ▷ Forward noise sample from  $t - 1$  to  $t$ ,  $\mathbf{x}_t \sim \bar{p}_{t|t-1}(\mathbf{x}_{t-1})$ 
18:       $\boldsymbol{\zeta}_{t-1} \sim \mathcal{P}_{\text{noise}}$ 
19:       $\mathbf{x}_t \leftarrow \sqrt{1 - \beta_{t-1}}\mathbf{x}_{t-1} + \sqrt{\beta_{t-1}}\boldsymbol{\zeta}_{t-1}$ 
20: return  $\mathbf{x}_0$ 

```

---

395 **C Amortised learning of Doob’s transform**

396 **C.1 Proof of proposition 2.5**

397 *Proof.* (Informal) Via the mean squared error property of the conditional expectation the minimiser  
 398 is given by:

$$f_t^*(\mathbf{h}, \mathbf{y}, \mathcal{A}) = \mathbb{E} \left[ \nabla_{\mathbf{X}_t} \ln \bar{p}_{t|0}(\mathbf{X}_t | \mathbf{X}_0) | \mathbf{Y} = \mathbf{y}, \mathbf{X}_t = \mathbf{h} \right] \quad (24)$$

399 Then:

$$\begin{aligned} f_t^*(\mathbf{h}, \mathbf{y}, \mathcal{A}) &= \int \nabla_{\mathbf{h}} \ln \bar{p}_{t|0}(\mathbf{h} | \mathbf{X}_0) \bar{p}_{0|t}(\mathbf{X}_0 | \mathbf{X}_t = \mathbf{h}, \mathbf{Y} = \mathbf{y}) d\mathbf{X}_0 \\ &= \int \frac{\nabla_{\mathbf{h}} \bar{p}_{t|0}(\mathbf{h} | \mathbf{X}_0)}{\bar{p}_{0|t}(\mathbf{h} | \mathbf{X}_0)} \frac{\bar{p}_{t|0}(\mathbf{X}_t = \mathbf{h} | \mathbf{X}_0, \mathbf{Y} = \mathbf{y}) p(\mathbf{X}_0 | \mathbf{Y} = \mathbf{y})}{p(\mathbf{X}_t = \mathbf{h} | \mathbf{Y} = \mathbf{y})} d\mathbf{X}_0 \\ &= \frac{1}{p(\mathbf{X}_t = \mathbf{h} | \mathbf{Y} = \mathbf{y})} \int \frac{\nabla_{\mathbf{h}} \bar{p}_{t|0}(\mathbf{h} | \mathbf{X}_0)}{\bar{p}_{0|t}(\mathbf{h} | \mathbf{X}_0)} \bar{p}_{t|0}(\mathbf{X}_t = \mathbf{h} | \mathbf{X}_0) p(\mathbf{X}_0 | \mathbf{Y} = \mathbf{y}) d\mathbf{X}_0 \\ &= \frac{1}{p(\mathbf{X}_t = \mathbf{h} | \mathbf{Y} = \mathbf{y})} \nabla_{\mathbf{h}} \int \bar{p}_{t|0}(\mathbf{h} | \mathbf{X}_0) p(\mathbf{X}_0 | \mathbf{Y} = \mathbf{y}) d\mathbf{X}_0 \\ &= \frac{1}{\bar{p}(\mathbf{X}_t = \mathbf{h} | \mathbf{Y} = \mathbf{y})} \nabla_{\mathbf{h}} \bar{p}(\mathbf{X}_t = \mathbf{h} | \mathcal{A} \mathbf{X}_0 = \mathbf{y}) = \nabla_{\mathbf{h}} \ln \bar{p}(\mathbf{X}_t = \mathbf{h} | \mathbf{Y} = \mathbf{y}), \end{aligned}$$

400

□

401 **D Experimental details: image experiments**

402 In the image experiment, we use the DDPM [Ho et al., 2020] formulation for the diffusion model  
 403 with  $N = 1000$  steps, a linear  $\beta$ -schedule with  $\beta_0 = 10^{-4}$  and  $\beta_N = 2 \cdot 10^{-2}$ .

404 **Data** We focus on the CELEBA [Liu et al., 2015] and FLOWERS [Nilsback and Zisserman, 2008]  
 405 image datasets. For each of these datasets, we follow the same preprocessing procedure consisting  
 406 of centrally cropping the image to size  $64 \times 64$ , and rescaling to pixel values  $[-1, 1]$ . We use this  
 407 information to also clip our model’s prediction.

408 **Noise model** The noise model  $\epsilon_\theta$  consists of a UNET architecture with four downsampling blocks  
 409 consisting of 2d convolutional layers of dimensionality 128, 256, 384 and 512, respectively. We  
 410 apply attention in the middle layers of the UNet with four heads. Throughout the network, we use the  
 411 SiLU activation function, no dropout and group normalisation layers. The amortised network differs  
 412 from the unconditional network in the fact that it accepts as input twice the number of channels (six  
 413 instead of only three RGB channels). The unconditional models operate directly on the three RGB  
 414 channels while the amortised network operates on the RBG channels, the mask and the condition.  
 415 We can represent the mask and the condition information, however, into a single input with the same  
 416 dimension as the image. The values of this input will be equal to the condition when the mask is 1  
 417 and set to a padding value of  $-2$  where the mask is 0. We concatenate the image  $\mathbb{R}^{3 \times H \times W}$  with  
 418 the condition and mask input of size  $\mathbb{R}^{3 \times H \times W}$  into an image with six channels. Due to this minor  
 419 difference, our amortised network has 68.159M parameters while the unconditional networks have  
 420 68.156M parameters (roughly 0.005% fewer).

421 **Methods** In the amortised setting we follow Alg. 5. In 90% of the training steps, we pass a condition  
 422 to the network. The other 10% contains a mask consisting of only 0’s. For the reconstruction guidance  
 423 method, we use a guidance term of  $\gamma = 10.0$ .

424 **Metrics** We measure the performance of the methods using mean squared error (MSE) and the  
 425 perceptual metric LPIPS. Both these metrics compare the similarity between the original image (from  
 426 which a patch was taken) and the conditional sample. For each metric, we compute the mean across  
 427 64 test images and repeat the experiment 5 times to get error estimates.