

THE SUPERPOSITION OF DIFFUSION MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

The Cambrian explosion of easily accessible pre-trained diffusion models suggests a demand for methods that combine multiple different pre-trained diffusion models without incurring the significant computational burden of re-training a larger combined model. In this paper, we cast the problem of combining multiple pre-trained diffusion models at the generation stage under a novel proposed framework termed superposition. Theoretically, we derive superposition from rigorous first principles stemming from the celebrated continuity equation and design two novel algorithms tailor-made for combining diffusion models in SUPERDIFF. We demonstrate that SUPERDIFF is scalable to large pre-trained diffusion models as superposition is performed *solely through composition during inference*, and also enjoys painless implementation as it combines different pre-trained vector fields through an automated re-weighting scheme. Notably, we show that SUPERDIFF is efficient during inference time, and mimics traditional composition operators such as the logical OR and the logical AND. We empirically demonstrate the utility of using SUPERDIFF for generating more diverse images on CIFAR-10, more faithful prompt conditioned image editing using Stable Diffusion, and improved unconditional *de novo* structure design of proteins.



Figure 1: Concept interpolations via different methods: SUPERDIFF (top row), the averaging of outputs with different prompts (middle row), and joint prompting with standard Stable Diffusion (SD) (bottom row) for six different prompt combinations. Here we use SUPERDIFF with the AND operation (sampling equal densities).

1 INTRODUCTION

The design and application of generative models at scale are arguably one of the fastest-growing use cases of machine learning, with generational leaps in performance that often exceed expert expectations (Steinhardt, 2022). A few of the central facilitators of this rapid progress are the availability of high-quality training data and large computing hardware (Kaplan et al., 2020); which in tandem provide a tried and trusted recipe to scale generative models in a variety of data modalities such as video generation (Brooks et al., 2024), natural language understanding (OpenAI, 2023; Achiam

054 et al., 2023; Dong et al., 2022), and other challenging domains like mathematical reasoning (Trinh
 055 et al., 2024), or code assistance (Bubeck et al., 2023). As a result, it is not surprising that a driving force
 056 behind current generative modeling research is centered around developing open-source tooling (Dao
 057 et al., 2022; Kwon et al., 2023) to enable further scaling and understanding emergent behavior of
 058 such models (Schaeffer et al., 2023), including probing current limitations (Dziri et al., 2024).

059 Indeed, the rapid escalation of generative model development has also induced a democratizing effect,
 060 given the easy access to large pre-trained in the current AI climate (Stability AI, 2023; Midjourney,
 061 2023; Ramesh et al., 2021). Furthermore, with the rise of open-source models, it is now easier than
 062 ever to host and deploy fine-tuned models. However, the current pace of progress also makes it
 063 infeasible to easily scale further models without confronting practical challenges. For instance, for
 064 continuous domains such as natural images current pre-trained diffusion models already exhaust all
 065 public data, with a growing proportion of the web already populated with synthetic data (Schuhmann
 066 et al., 2022). Compounding these challenges is the tremendous cost of pre-training these large
 067 diffusion models, which at present makes it computationally unattractive for individuals to build
 068 large pre-trained models on different datasets without re-training a larger combined model.

069 A compelling alternative to training ever larger models is to consider the efficacy of maximizing the
 070 utility of existing pre-trained models. In particular, it is interesting to consider the compositional
 071 benefits of combining pre-trained models at the generation stage in place of training a single mono-
 072 lithic model (Du & Kaelbling, 2024). For diffusion models, which are the current de facto modeling
 073 paradigm over continuous domains, compositional generation can be framed as modifying the in-
 074 ference mechanism through either guidance (Dhariwal & Nichol, 2021; Ho & Salimans, 2022) or
 075 applying complex MCMC correction schemes (Du et al., 2023). Despite advancements, methods that
 076 utilize guidance lack a firm theoretical underpinning (Bradley & Nakkiran, 2024) and MCMC tech-
 077 niques are prone to scaling issues, and more importantly, remain unproven at combining existing large
 078 pre-trained diffusion models. This motivates the following timely research question: *Can we combine
 pre-trained diffusion models solely at inference in a theoretically sound and efficient manner?*

079 **Present work.** In this paper, we cast the problem of combining and composing multiple pre-trained
 080 diffusion models under a novel joint inference framework we term superposition. Intuitively we
 081 develop our framework by starting from the well-known principle of superposition in physical systems
 082 which summarizes the net response to multiple inputs in a linear system as the sum of individual inputs.
 083 Applying the superposition principle we develop SUPERDIFF, a novel approach to combine pre-
 084 trained diffusion models by superimposing learned vector fields at inference. Moreover, we demon-
 085 strate that the same superposition principle underlies both flow-based methods (Liu et al., 2022b;
 086 Lipman et al., 2022; Albergo et al., 2023) and diffusion models (Ho et al., 2020; Song et al., 2020)
 (see App. B), and thus provides a unifying perspective over combining modern generative models.

087 SUPERDIFF relies on the key insight that following a superposition principle necessitates the efficient
 088 evaluation of the densities of different models during the generation, which we present in Sec. 3.1.
 089 Knowing the current density of the generated sample under different models (or the same model
 090 with different conditioning) we are able to guide the generation to the desired result by controlling
 091 the relative weights of model outputs. Namely, in Sec. 3.2, we propose two different algorithms that
 092 implement: a mixture of densities (sampling from either from one model OR another), generating
 093 samples of the equal densities (samples that are likely under one model AND another, e.g. see Fig. 1).

094 We test the applicability of our approach SUPERDIFF using the two proposed superimposition
 095 strategies for image and protein generation tasks. For images, we first demonstrate the ability to
 096 combine the outputs of two models trained on disjoint datasets such that they yield better performance
 097 than the model trained on both of the datasets (see Sec. 4.1). In addition, we demonstrate the ability
 098 to interpolate between densities which correspond to concept interpolation in the image space (see
 099 Sec. 4.2). For proteins, we demonstrate that combining two different generative models leads to
 100 improvements in designability and novelty generation (see Sec. 4.3). Across all our experimental
 101 settings, we find that combining pre-trained diffusion models using SUPERDIFF leads to higher
 102 fidelity generated samples that better match task specification in text-conditioned image generation
 103 and also produce more diverse generated protein structures than comparable composition strategies.

104 2 PRELIMARIES

105 Generative models learn to approximate a target data distribution $p_{\text{data}} \in \mathbb{P}(\mathbb{R}^d)$ defined over \mathbb{R}^d
 106 using a parametric model q_θ with learnable parameters θ . In the conventional problem definition, the
 107 data distribution is realized as an empirical distribution that is provided as a training set of samples

$\mathcal{D} = \{\mu^i\}_{i=1}^m$. Whilst there are multiple generative model families to choose from, we restrict our attention to diffusion models (Song et al., 2020; Ho et al., 2020) which are arguably the most popular modeling family driving current application domains. We next review the basics of the continuous time formulation of diffusion models before casting them within our superposition framework.

2.1 CONTINUOUS-TIME DIFFUSION MODELS

A diffusion model can be cast as the solution to the Stochastic Differential Equation (Øksendal, 2003),

$$dx_t = f_t(x)dt + g_t dW_t, \quad x_0 \sim q_0(x_0). \quad (1)$$

In the Itô SDE literature, the function $f_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is known as the drift coefficient while $g_t : \mathbb{R} \rightarrow \mathbb{R}$ is a real-valued function called the diffusion coefficient and W_t is the standard Wiener process. The subscript index $t \in [0, 1]$ indicates the time-valued nature of the stochastic process. Specifically, we fix the starting time $t = 0$ to correspond to the data distribution $p_{\text{data}} := q_0(x_0)$ and set $t = 1$ as the terminal time to an easy to sample prior such as a standard Normal distribution $p_{\text{noise}} := q_1(x_1) = \mathcal{N}(x_1|0, I)$. As such the diffusion SDE, also called the forward process, can be seen as progressively corrupting the data distribution and ultimately hitting a terminal distribution devoid of any structure.

To generate samples from the marginal density $q_t(x_t)$ induced by the diffusion SDE in equation 1 we leverage the reverse-time SDE with the same marginal density as demonstrated below.

Proposition 1. [Reverse-time SDEs/ODE] *Marginal densities $q_t(x)$ induced by Eq. (1) correspond to the densities induced by the following SDE that goes back in time ($\tau = 1 - t$) with the corresponding initial condition*

$$dx_\tau = \left(-f_t(x_\tau) + \left(\frac{g_t^2}{2} + \xi_\tau \right) \nabla \log q_t(x_\tau) \right) d\tau + \sqrt{2\xi_\tau} d\bar{W}_\tau, \quad x_{\tau=0} \sim q_1(x_0), \quad (2)$$

where \bar{W}_τ is the standard Wiener process in time τ , and ξ_τ is any positive schedule.

See proof in App. A.1. The reverse SDE flows backward in time $\tau = 1 - t$ and is linked to the diffusion SDE through the score $\nabla_x \log q_\tau(x)$, and $d\bar{W}_\tau$ is another Wiener process. As a result, a parametric model $\nabla \log q_\tau(x; \theta)$ may directly learn to approximate this score function for every point in time and then draw samples by simulating the reverse SDE in equation 2 by plugging back in the learned score. Notably, for $\xi_t \equiv 0$, the SDE becomes an ODE which defines a smooth change of measure corresponding to p_{noise} into the measure corresponding to p_{data} .

In practice, for generative modelling, the forward SDE from Eq. (1) is chosen to be so simple that it can be integrated in time analytically without simulating the SDE itself. This is equivalent to choosing the noising schedule first, and then deriving the SDE that corresponds to this schedule. Namely, for every training sample μ^i , we can define the density of the corrupted μ^i as a normal density with the mean scaled according to α_t and the variance σ_t^2 , then the density of the entire corrupted dataset is simply a mixture over all training samples μ^i , i.e.

$$q_t^i(x) = \mathcal{N}(x | \alpha_t \mu^i, \sigma_t^2 I), \quad q_t(x) = \frac{1}{N} \sum_{i=1}^N q_t^i(x). \quad (3)$$

Clearly, choosing α_t, σ_t such that $\alpha_0 = 1, \sigma_0 = 0$ and $\alpha_1 = 0, \sigma_1 = 1$, we guarantee $p_{\text{data}} := q_0(x_0)$ and $p_{\text{noise}} := q_1(x_1) = \mathcal{N}(x_1|0, I)$. This perturbation of the data distribution using a Gaussian kernel offers specific forms for the drift f_t , and diffusion coefficient g_t as described next.

Proposition 2. [Ornstein–Uhlenbeck SDE] *The time-dependent densities in Eq. (3) correspond to the marginal densities of the following SDE, with the corresponding initial condition*

$$dx_t = \underbrace{\frac{\partial \log \alpha_t}{\partial t} x_t}_{f_t(x_t)} dt + \underbrace{\sqrt{2\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t}}}_{g_t} dW_t, \quad x_0 \sim q_0(x_0). \quad (4)$$

See proof in App. A.2. We highlight the simplicity of the drift term, a linear scaling, that allows us to simulate efficiently the reverse SDE and is crucial for the proposed density estimators in Sec. 3.1. Altogether, the derivations of this section allow us to go from the noise schedules of samples

in Eq. (3) used during the training of a given diffusion model to the corresponding forward SDE in Prop. 2, and finally to the reverse SDEs or ODE used in Prop. 1.

2.2 SUPERPOSITION OF ODES AND SDES

In this section, we introduce the superposition of multiple time-dependent densities that correspond to different stochastic processes. A suggestive view of these densities is as processes corresponding to different training data (e.g. different datasets), different conditions (e.g. different text prompts), or simply differently trained diffusion models. Namely, we consider N forward noising processes $\{q_t^i(x)\}_{i=1}^N$ that possibly start from different initial distributions $\{q_{t=0}^i(x)\}_{i=1}^N$ (e.g. different datasets). Assume that we know the individual vector fields $v_t^i(x)$ that define the change of corresponding densities $q_t^i(x)$ via the state-space ODEs and continuity equations,

$$\frac{dx_t}{dt} = v_t^i(x_t) \implies \frac{\partial}{\partial t} q_t^i(x) = -\langle \nabla_x, q_t^i(x) v_t^i(x) \rangle, \quad \forall i \in [N]. \quad (5)$$

The *superposition* of the noising processes $\{q_t^i(x)\}_{i=1}^N$ is the mixture of corresponding densities:

$$q_t^{\text{mix}}(x) := \sum_{j=1}^N \omega^j q_t^j(x), \quad \sum_{i=1}^N \omega^i = 1, \quad \omega^i \geq 0, \quad (6)$$

where ω^j is a mixing coefficient. Note that superimposed $q_t^{\text{mix}}(x)$ also satisfies the continuity equation for the superposition of the vector fields $v_t^i(x)$ as demonstrated in the following proposition.

Proposition 3. [Superposition of ODEs (Liu, 2022)] *The mixture density in Eq. (6) follows the continuity equation with the superposed vector fields from Eq. (5), i.e.*

$$\frac{\partial}{\partial t} q_t^{\text{mix}}(x) = -\langle \nabla_x, q_t^{\text{mix}}(x) v_t(x) \rangle, \quad v_t(x) = \frac{1}{\sum_{j=1}^N \omega^j q_t^j(x)} \sum_{i=1}^N \omega^i q_t^i(x) v_t^i(x). \quad (7)$$

We reproduce the proof for Prop. 3 in the context of our superposition in App. A.3. The superposition principle is the core principle that allows for efficient simulation-free learning of the flow-based models (Liu et al., 2022b; Lipman et al., 2022; Albergo et al., 2023) and diffusion models (Song et al., 2020). We discuss how these frameworks are derived from the superposition principle in App. B.

The superposition principle straightforwardly extends to the marginal densities of SDEs. That is, consider marginal densities $q_\tau^i(x)$ generated by the following SDEs

$$dx_\tau = u_\tau^i(x_\tau) d\tau + g_\tau d\bar{W}_\tau, \quad x_{\tau=0} = q_{\tau=0}^i(x_0), \quad (8)$$

where one has to note the same diffusion coefficient for all the SDEs. Then the mixture of densities from Eq. (6) can be simulated by the SDE from the following proposition.

Proposition 4. [Superposition of SDEs] *The mixture $q_t^{\text{mix}}(x) := \sum_{i=1}^N \omega^i q_t^i(x)$ of density marginals $\{q_t^i(x)\}_{i=1}^N$ induced by SDEs from Eq. (8) corresponds to the following SDE*

$$dx_\tau = u_\tau(x_\tau) d\tau + g_\tau d\bar{W}_\tau, \quad u_\tau(x) = \frac{1}{\sum_{j=1}^N \omega^j q_\tau^j(x)} \sum_{i=1}^N \omega^i q_\tau^i(x) u_\tau^i(x). \quad (9)$$

See App. A.4 for the proof. Both Prop. 3 and Prop. 4 can be easily extended to the families of densities parameterized with a continuous variable, but this is beyond the scope of the current paper

3 SUPERPOSITION OF DIFFERENT MODELS

We now introduce our method for combining pre-trained diffusion models using the principle of superposition. The result of this is a novel inference time algorithm SUPERDIFF which can be easily applied without further fine-tuning or post-processing of any of the pre-trained diffusion weights. Our proposed approach SUPERDIFF can be instantiated in two distinct ways that allow for the composition diffusion models that can be informally interpreted as logical composition operators in the logical AND and the logical OR. More precisely, given two pre-trained diffusion models that are trained on datasets A and B inference using SUPERDIFF can be done by either sampling from the mixture of the two learned densities, i.e. logical AND Fig. 2c, or sampling from the equal density locus (logical OR Fig. 2d). In such a manner, superposition using AND leads to generated samples that are equally likely

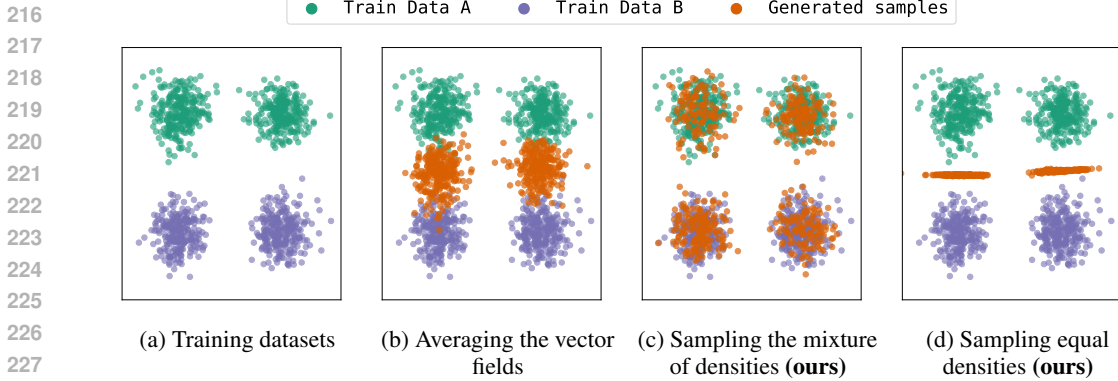


Figure 2: An intuitive illustration of using model superposition for improving inference performance. We show an example of two disjoint datasets and train a model for each set. Each individual model learns to generate samples only from their respective datasets. Using model superposition enables sampling from both densities.

under both pre-trained diffusion models while superposition using OR creates samples that are preferentially generated by either pre-trained model—and thus mimics the empirical distributions A or B.

Method overview. SUPERDIFF is applicable in settings where the modeler has access to M pre-trained diffusion models, along with its learned score function $\nabla_x \log q_t^i(x)$. Each of the pre-trained models follows a marginal density $q_t^i(x)$ and as a result must admit a corresponding vector field $v_t^i(x)$ that satisfies the continuity equation in Equation 5. The key idea of our approach is an adaptive re-weighting scheme of the pre-trained model’s vector fields that relies on the likelihood of a sample under different models. A naive approach to estimating each marginal density during generation immediately presents several technical challenges as it requires the estimation of the divergence of superimposed vector fields. In particular, these challenges can be stated as follows:

- (C1) The marginal superpositioned vector field differs from the vector fields of either of the models.
- (C2) The divergence operation requires backpropagation through the network and is computationally expensive even with Hutchinson’s trace estimator (Hutchinson, 1989).

Our proposed approach SUPERDIFF overcomes these computational challenges by introducing a novel density estimator in Sec. 3.1. Crucially, this new estimator does not require divergence estimation and enjoys having the same variance as the computationally expensive Hutchinson’s trace estimator, making it a favorable choice when generating using large pre-trained diffusion models. In section Sec. 3.2 we exploit this new density estimator to formally present our algorithm SUPERDIFF and derive connections to the composition operators that intuitively resemble logical AND and OR.

3.1 EVALUATING THE DENSITIES ON THE FLY

In this section, we introduce a novel method for evaluating the marginal density of a diffusion model during the inference process. The conventional way to evaluate the density uses the continuity equation and solves the same ODE that is used for generating samples. This, however, is not easily possible in the case of our superposition of vector fields framework as outlined in Prop. 3. To solve this, we present the following proposition that disentangles the vector field generating the sample ($u_t(x)$ in the proposition) and the vector fields corresponding to different generative models $v_t^i(x)$.

Proposition 5. [Smooth density estimator] *For the integral curve $x(t)$ solving $dx/dt = u_t(x_t)$, and the density $q_t^i(x(t))$ satisfying the continuity equation $\frac{\partial}{\partial t} q_t^i(x) = -\langle \nabla_x, q_t^i(x) v_t^i(x) \rangle$, the log-density along the curve changes according to the following ODE*

$$\frac{d}{dt} \log q_t^i(x(t)) = -\langle \nabla_x, v_t^i(x) \rangle - \langle \nabla_x \log q_t^i(x), v_t^i(x) - u_t(x) \rangle. \quad (10)$$

See proof in App. C. We use this proposition for our experiments conducted in Sec. 4.1. However, as outlined previously, evaluating the marginal density via the continuity equation is restricted to small-scale models due to the computational challenges associated with efficiently estimating the divergence of the associated vector field. As a result, a common line of attack assumes constructing

Algorithm 1: SUPERDIFF pseudocode (for **OR** and **AND** operations)

Input : M pre-trained score models $\nabla_x \log q_t^i(x)$, the parameters of the schedule α_t, σ_t , stepsize $d\tau > 0$, temperature parameter T , bias parameter ℓ , and initial noise $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

for $\tau = 0, \dots, 1$ **do**

$t = 1 - \tau, \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

$\kappa_\tau^i \leftarrow \begin{cases} \text{softmax}(T \log q_t^i(x_\tau) + \ell) & // \text{ for OR according to Prop. 3} \\ \text{solve Linear Equations} & // \text{ for AND according to Prop. 6} \end{cases}$

$u_t(x) \leftarrow \sum_{i=1}^M \kappa_\tau^i \nabla \log q_t^i(x_\tau)$

$dx_\tau \leftarrow (-f_{1-\tau}(x_\tau) + g_{1-\tau}^2 u_t(x))d\tau + g_{1-\tau}d\bar{W}_\tau \quad // \text{ using Prop. 1}$

$x_{\tau+d\tau} \leftarrow x_\tau + dx_\tau$

$d \log q_{1-\tau}(x_\tau) = \langle dx_\tau, \nabla \log q_{1-\tau}(x_\tau) \rangle + \left(\langle \nabla, f_{1-\tau}(x_\tau) \rangle + \left\langle f_{1-\tau}(x_\tau) - \frac{g_{1-\tau}^2}{2} \nabla \log q_{1-\tau}(x_\tau), \nabla \log q_{1-\tau}(x_\tau) \right\rangle \right) d\tau \quad // \text{ using Thm. 1}$

return x

a stochastic unbiased estimator that trades for increased speed by introducing a bit of variance into the divergence estimate. This approach is known as Hutchinson’s estimator and requires computing a Jacobian-vector product at every step of the inference.

Instead, we propose a new way to estimate density that allows for efficient computation while integrating the backward SDE from Prop. 1 with a specific choice of the diffusion coefficient.

Theorem 1. [$\hat{\text{O}}_t$ density estimator] *Consider time-dependent density $q_t(x)$ induced by the marginals of the following “forward” process*

$$dx_t = f_t(x_t)dt + g_t dW_t, \quad x_{t=0} \sim q_0(x), \quad t \in [0, 1], \quad (11)$$

where dW_t is the Wiener process. Then, for the following “backward” SDE (with $\tau = 1 - t$)

$$dx_\tau = u_\tau(x_\tau)d\tau + g_{1-\tau}d\bar{W}_\tau, \quad \tau \in [0, 1], \quad (12)$$

the change of the log-density $\log q_\tau(x_\tau)$ follows the following SDE

$$d \log q_{1-\tau}(x_\tau) = \langle dx_\tau, \nabla \log q_{1-\tau}(x_\tau) \rangle + \left(\langle \nabla, f_{1-\tau}(x_\tau) \rangle + \left\langle f_{1-\tau}(x_\tau) - \frac{g_{1-\tau}^2}{2} \nabla \log q_{1-\tau}(x_\tau), \nabla \log q_{1-\tau}(x_\tau) \right\rangle \right) d\tau. \quad (13)$$

We provide the proof in App. C. Notably, the SDE Eq. (13) that keeps track of the change of log-density does not include only the divergence of the forward SDE drift $\langle \nabla, f_{1-\tau}(x_\tau) \rangle$. However, in practice, when using the Ornstein-Uhlenbeck SDE, this divergence is simply a constant because the drift term is a linear scaling. In App. D, we derive the same estimator but in discrete time using the detailed balance condition.

3.2 SUPERDIFF: SUPERPOSING PRE-TRAINED DIFFUSION MODELS

Mixture of densities (logical OR). For a mixture of the densities, superposition of the models follows directly from the propositions in Sec. 2.2. That is, for every $q_t^i(x)$, we assume that it can be generated from SDEs or an ODE from Prop. 1 and we assume the knowledge of scores $\nabla \log q_t^i(x)$.

Then, for the ODE simulation, according to Prop. 3, we can sample from the mixture of densities $q_t^{\text{mix}}(x) := 1/M \sum_{i=1}^M q_t^i(x)$ using the following vector field

$$v_\tau(x) = -f_{1-\tau}(x) + \frac{g_{1-\tau}^2}{2} \sum_{i=1}^M \frac{q_t^i(x)}{\sum_j q_t^j(x)} \nabla \log q_{1-\tau}^i(x), \quad (14)$$

starting from samples $x_{\tau=0} \sim q_1(x_0)$. The densities are estimated along the trajectory using Prop. 5.

We highlight this analogously applies for simulation using the SDE $dx_\tau = u_\tau(x_\tau)d\tau + g_{1-\tau}d\bar{W}_\tau$. Namely, according to [Prop. 4](#), we use the following vector field:

$$u_\tau(x) = -f_{1-\tau}(x) + g_{1-\tau}^2 \sum_{i=1}^M \frac{q_t^i(x)}{\sum_j q_t^j(x)} \nabla \log q_{1-\tau}^i(x), \quad (15)$$

starting from the samples $x_{\tau=0} \sim q_1(x_0)$. The densities are estimated along the trajectory using [Thm. 1](#). We provide the pseudocode in [Algorithm 1](#).

Sampling equal densities (logical AND). To produce the sample that have equal densities under different models we rely on our formula for the density update (see [Thm. 1](#)) to find the optimal weights for the vector fields. Indeed, for M diffusion models, we have a system of M equations: the equal change of density for every model and the normalization constraint for model weights, which is a linear system w.r.t. vector field weights as we show in the following proposition.

Proposition 6. [Density control] *For the SDE*

$$dx_\tau = \sum_{j=1}^M \kappa_j u_\tau^j(x_\tau) d\tau + g_{1-\tau} d\bar{W}_\tau, \quad (16)$$

where κ are the weights of different models and $\sum_j \kappa_j = 1$, one can find κ that satisfies

$$d \log q_{1-\tau}^i(x_\tau) = d \log q_{1-\tau}^j(x_\tau), \quad \forall i, j \in [M], \quad (17)$$

by solving a system of M linear equations w.r.t. κ .

We provide the proof and the formulas for the system of linear equations in [App. C.1](#). This proposition allows us to find κ that controls the densities to stay the same for all the models as described in [Algorithm 1](#). This approach can also be straightforwardly extended to the case of diffusion models for satisfying different prescribed density ratios, i.e.

$$d \log q_{1-\tau}^1(x_\tau) = d \log q_{1-\tau}^i(x_\tau) + \ell^i. \quad (18)$$

4 EXPERIMENTS

4.1 JOINING MODELS WITH DISJOINT TRAINING DATA

We validate the proposed [Algorithm 1](#) for the generation of the mixture of the distributions (OR setting). We split CIFAR-10 into two disjoint training sets of equal size (first 5 labels and last 5 labels), train two diffusion models on each part, and generate the samples jointly using both models. Namely, the stochastic inference is the OR implementation of [Algorithm 1](#), whereas the deterministic setting is the integration of the ODEs (see [Prop. 1](#)) and the estimation of the log-density according to [Prop. 5](#). For the choice of hyperparameters, architecture, and data preprocessing, we follow ([Song et al., 2020](#)).

In [Table 1](#), we demonstrate that the performance of SUPERDIFF drastically outperforms the performance of the individual models and performs even better than the model trained on the union of both parts of the dataset. For comparison, we evaluate conventional image quality metrics (Fréchet inception distance (FID), Inception score (IS), and feature likelihood divergence (FLD) ([Jiralerspong et al., 2023](#)), which takes into account the generalization abilities of the model.¹

4.2 CONCEPT INTERPOLATION AND SELECTION WITH SUPERDIFF AND STABLE DIFFUSION

Next, we evaluate the ability of SUPERDIFF to interpolate (logical AND) or select (logical OR) different concepts using prompt-conditioned Stable Diffusion (SD). In this setup, we generate images from SD by conditioning it on a prompt using classifier-free guidance. We define two models using two separate prompts that represent concepts — e.g. "a sunflower" and "a lemon". We can thus consider each prompt-conditioned model as a separate diffusion model and apply [Algorithm 1](#) to generate images sampled with proportionally equal likelihood with respect to both prompt-conditioned models (logical AND), or to generate from the mixture of distributions (logical OR). We generate 20 images for 20 different concept pairs (tasks) for our model and baselines (c.f. [App. H](#)).

¹In [Table 1](#), we report results for both the SDE and the ODE inference. However, we note that for the ODE setting, evaluating densities takes about 4x longer than our density estimator in the SDE setting and requires 1.5x more memory, and so we carry out remaining experiments only using SDEs.

Table 1: Image generation performance for CIFAR-10 with *unconditional* models trained on two disjoint partitions of the training data (labeled A and B). We compare SUPERDIFF (OR) with the respective models (model_A and model_B) with the model that is trained on the full dataset (model_{A∪B}) and random choice between two models (model_{A OR B}).

	ODE inference			SDE inference		
	FID (↓)	IS (↑)	FLD (↓)	FID (↓)	IS (↑)	FLD (↓)
model _A	14.00	8.73	15.50 ± 0.17	15.33	7.98	15.47 ± 0.18
model _B	13.09	7.89	18.90 ± 0.18	13.50	7.98	18.54 ± 0.23
model _{A∪B}	6.00	8.95	8.06 ± 0.12	3.50	9.14	7.51 ± 0.11
model _{A OR B}	4.28	9.14	5.96 ± 0.11	3.99	9.36	5.29 ± 0.14
SUPERDIFF (OR)	4.41	9.12	6.10 ± 0.11	4.00	9.36	5.33 ± 0.05
SUPERDIFF _{T=100} (OR)	4.11	9.21	5.89 ± 0.11	4.00	9.48	5.20 ± 0.11

Table 2: Quantitative evaluation of SD-generated images. For logical AND, we compare SUPERDIFF (AND), joint prompting, and averaging of outputs for two concept prompts. We report the minimum CLIP, ImageReward, and TIFA scores from the two concept prompts, which gives a measure of how well each method represents both concepts. For logical OR, we compare SUPERDIFF (OR), joint prompting, and an upper bound of randomly selecting a single concept prompt (Prompt_{A OR B}). We report the maximum scores from each prompt pair, as well as the absolute difference between the maximum and minimum scores (|Δ|). These metrics reflect how well a method can select one concept to generate.

	Min. CLIP (↑)	Min. ImageReward (↑)	Min. TIFA (↑)
Joint prompting	23.87	-1.62	27.58
Average of scores (Liu et al., 2022a)	24.23	-1.57	32.48
SUPERDIFF (AND)	24.79	-1.39	39.92
	Δ /Max. CLIP (↑)	Δ /Max. ImageReward (↑)	Δ /Max. TIFA (↑)
Prompt _{A OR B} (uncorrelated random choice)	9.13/29.72	2.87/0.70	88.21/97.58
Joint prompting	7.20/29.80	2.47/0.59	79.46/97.92
SUPERDIFF (OR)	8.58/29.87	2.76/0.64	84.10/95.83

Baselines for concept interpolation (AND). We consider two approaches for composing images as baselines. The first is simple averaging of SD outputs based on the approach in (Liu et al., 2022a); we set $\kappa = 0.5$. The second guides SD generation with a single joint prompt that encourages SD to interpolate concepts. The prompts are constructed by joining two input concepts with the linking term "that looks like". For example, given the concepts "a sunflower" and "a lemon", the prompt is "a sunflower that looks like a lemon". For fairness, we also flip the order of the prompt and keep the image with the higher score for all metrics listed below (Luo et al., 2024).

Baselines for concept selection (OR). For the OR setting, our baseline is prompting SD that prompts the model to select between two concepts: "a sunflower or a lemon". As with AND, we also flip the prompt and keep the better image for each metric. As an upper bound, we generate images from SD by prompting it with a random choice between the two concepts (Prompt_{A OR B}).

SUPERDIFF qualitatively generates images with better concept interpolations and selections. We plot sample generated images for SUPERDIFF (AND) in Fig. 1. For the complete set of generated images, see App. H (Figs. A1–A20). We observe that SUPERDIFF (AND) can interpolate concepts while also maintaining high perceptual quality. In contrast, the averaging baseline either fails to interpolate concepts from both prompts fully, or yields images with lower perceptual quality. We observe that SD using a single prompt struggles to interpolate both concepts.

For concept selection (OR), we find that SUPERDIFF (OR) can faithfully generate images with a single concept. The joint prompting baseline can sometimes generate images that combine fragments of both concepts, but other times it also generates images of a single concept, typically the first concept in the joint prompt (this also underscores why this method struggles with concept interpolation). We show examples of these cases in Fig. A21.

SUPERDIFF outperforms baselines on three image evaluation metrics. To quantitatively evaluate the generated images, we consider three metrics: CLIP Score (Radford et al., 2021), ImageReward (Xu et al., 2024), and TIFA (Hu et al., 2023). CLIP Score measures the cosine similarity between an image embedding and text prompt embedding. ImageReward evaluates generated images by assigning a score that reflects how closely they align with human preferences, including aesthetic quality and prompt adherence. TIFA generates several question-answer pairs using a Large Language Model for a given prompt and assigns a score by answering the questions based on the image with a visual question-answering model. We report these metrics for SUPERDIFF and all baselines in Table 2.

Table 3: Evaluation of SUPERDIFF and baseline models, Proteus & FrameDiff, for unconditional protein generation. We show results for three categories of metrics: designability, novelty, and diversity. We also include a baseline of simple averaging of scores ($\kappa = 0.5$). We use the parameter ℓ (see Eq. (18)) to control and bias model superposition towards Proteus, i.e. ($\ell > 0$) or FrameDiff ($\ell < 0$). We use temperature values $T = 1$ for all variants of SUPERDIFF. For each type of model composition (averaging or SUPERDIFF), we mark each metric with a (†) if it is better than both Proteus and FrameDiff on their own, and with a (*) if it better than either one of them.

	Designability		Novelty			Diversity		
	$< 2\text{\AA}$ sCRMSD (†)	sCRMSD (↓)	< 0.5 sCTM (†)	< 0.3 sCTM (†)	Max. sCTM (↓)	Frac. β (†)	Pairwise sCTM (↓)	Max. cluster (†)
FrameDiff	0.392 ± 0.03	4.315 ± 0.25	0.152 ± 0.02	0.016 ± 0.01	0.570 ± 0.02	0.175 ± 0.01	0.337 ± 0.02	0.326 ± 0.05
Proteus	0.928 ± 0.02	1.014 ± 0.07	0.360 ± 0.03	0.020 ± 0.01	0.536 ± 0.01	0.119 ± 0.01	0.312 ± 0.01	0.217 ± 0.02
Average of scores (Liu et al., 2022a)	0.740 ± 0.03*	1.960 ± 0.14*	0.360 ± 0.03*	0.024 ± 0.01†	0.511 ± 0.01†	0.139 ± 0.01*	0.310 ± 0.01†	0.253 ± 0.01*
SUPERDIFF $_{\ell=0}$ (OR)	0.752 ± 0.03*	1.940 ± 0.14*	0.276 ± 0.03*	0.008 ± 0.01	0.547 ± 0.01*	0.147 ± 0.01*	0.309 ± 0.02†	0.268 ± 0.02*
SUPERDIFF $_{\ell=1}$ (OR)	0.976 ± 0.01†	0.929 ± 0.05†	0.396 ± 0.03†	0.024 ± 0.01†	0.528 ± 0.01†	0.127 ± 0.01*	0.307 ± 0.02†	0.246 ± 0.03*
SUPERDIFF $_{\ell=0}$ (AND)	0.740 ± 0.03*	2.119 ± 0.16*	0.320 ± 0.03*	0.044 ± 0.01†	0.513 ± 0.01†	0.142 ± 0.01*	0.310 ± 0.02†	0.265 ± 0.01†

For the logical **AND** setting, we evaluate the image against each concept prompt separately (i.e., "a sunflower" as one prompt and "a lemon" as the other) and take the minimum score for each metric. This is so that we can measure how well *both* concepts are represented. We find that SUPERDIFF (**AND**) obtains the highest scores across all metrics, indicating that our method can better represent both concepts in the images, while the baseline methods typically only represent one concept or (especially in the case of averaging outputs), generate compositions with lower fidelity. For the logical **OR** setting, we again evaluate the image against each concept prompt separately and take both the maximum score and the absolute difference between both scores for each metric. This is so that we can measure how well *one* concept is represented. The upper bound for this setting is randomly prompting SD with either of the prompts; we find that we are almost able to match this setting across all scores, indicating that our method is able to faithfully select a single concept. SD with joint prompting does not perform as well, as nothing prevents it from combining components from both concepts.

4.3 PROTEIN GENERATION

Finally, we apply our method in the setting of unconditional *de novo* protein generation. Protein generation has critical implications in drug discovery (Abramson et al., 2024). A good understanding of the protein landscape is important to rationally find novel proteins.

We evaluate proteins generated by the superposition of two existing protein diffusion models, Proteus (Wang et al., 2024) and FrameDiff (Yim et al., 2023b), in terms of their *designability*, *novelty*, and *diversity*, following metrics from Bose et al. (2024). We report the results of our best model in Table 3, as well as results for each model individually and simply averaging them. All results are averages over 50 generated proteins at lengths {100, 150, 200, 250, 300} for 500 timesteps.

Metrics for evaluations. We consider *designability*, *novelty*, and *diversity* metrics for evaluation of unconditional generation of protein structures. Protein *designability* refers to the in-silico agreement between generated structures and refolded structures as computed using a purpose-built folding model e.g. ESMFold (Lin et al., 2022), which is positively correlated with the synthesizability of the protein in a wet-lab setting. Generally, if the root-mean-square distance between the generated and refolded proteins (sCRMSD) is less than 2Å, it is considered to be designable. We provide details of all metrics in this section in App. G.2.

We compute several *novelty* metrics based on the similarity of generated proteins to those from the set of known proteins (the training set), which are called sCTM scores; the lower the score, the less similar the generated protein is to the training data. Lastly, we use *diversity* to assess the degree of heterogeneity present in the set of generated proteins. This is done by clustering the generated proteins in terms of sequence overlap, and measuring the fraction of proteins with challenging-to-generate secondary structures such as β sheets (Bose et al., 2024).

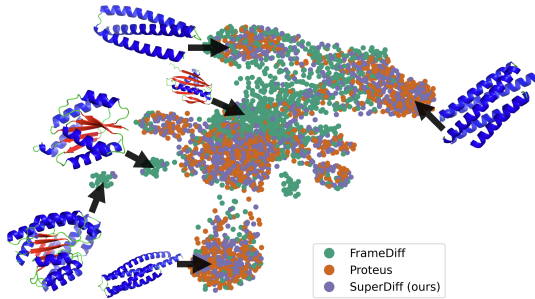


Figure 3: UMAP visualization of protein structures showing cluster archetypes where structure diversity is maintained with SUPERDIFF $_{\ell=1}$ (OR).

SUPERDIFF improves structure generation. By combining two protein diffusion models, SUPERDIFF is able to outperform both of them. This is somewhat surprising as FrameDiff is substantially less designable than Proteus (0.392 vs. 0.928 s_c RMSD) (see Table 3). Nevertheless, by using SUPERDIFF (OR), we can increase designability, novelty, and maintain diversity. We also find that SUPERDIFF (OR) substantially outperforms simple averaging, which simply mixes the performance of the two models, and does not improve over either one.

We also find that SUPERDIFF (AND) can match averaging in terms of designability while improving diversity. Perhaps more impressively, SUPERDIFF (AND) can generate the most proteins that are furthest away from the set of known proteins (s_c TM score < 0.3) by almost two times more than the next best method. This motivates the utility of applying our method in novel discovery settings.

We further investigate the composition made by SUPERDIFF (OR) in Fig. 3. Here, we see a few modes (particularly on the plot’s left-hand side) that Proteus does not generate. We find that SUPERDIFF (OR) can maintain knowledge of these clusters (although to a lesser extent) while maintaining designability.

5 RELATED WORK

Compositional diffusion models. Combining multiple density models into one model with better properties is a classical subject in machine learning (Hinton, 2002). For diffusion models, the most straightforward combination is via averaging their respective learned score functions (Liu et al., 2022a; Kong et al., 2024), which can also be viewed as a form of guidance (Ho et al., 2020; Dhariwal & Nichol, 2021; Bansal et al., 2023). Another way to combine diffusion models comes from their connections to energy-based models (Du et al., 2023; 2020; Nie et al., 2021; Ajay et al., 2024). This, however, comes under a strong assumption that the marginal densities of the noising process are given, which is not the case in most of the modern applications (e.g. Rombach et al. (2022)); our method resolves this. Density estimation, also, is the main bottleneck for the continual learning of diffusion models, e.g. Golatkar et al. (2023) proposes to learn a separate model for the densities in order to simulate the vector fields from Prop. 3. Finally, some works (Zhong et al., 2024; Biggs et al., 2024) propose to combine models by averaging their weights. This does not allow for merging of the models with different architectures (as we do in Sec. 4.3) and does not allow for merging the same model with different conditions (as we do in Sec. 4.2).

Protein generation. Structure-based de novo protein design using deep generative models has recently seen a surge in interest, with a particular emphasis on diffusion-based approaches (Watson et al., 2023; Yim et al., 2023b), and also flow matching methods (Bose et al., 2024; Yim et al., 2023a; 2024; Huguet et al., 2024). Moreover, building on the initial SE(3) equivariant diffusion paradigm multiple recent approaches have sought to increase the performance of the methods through architectural innovations (Wang et al., 2024), conditioning on auxiliary modalities such as sequence or sidechains (Ingraham et al., 2023; Lin et al., 2024). Finally, recent approaches tackle the problem of co-generation which seeks to define a joint inference procedure over both structure and sequences (Campbell et al., 2024; Ren et al., 2024; Lisanza et al., 2023), but remains distinct from the setting of this work which attempts to combine different pre-trained models using superposition.

6 CONCLUSION

Despite the ubiquity of diffusion models, many possible ways of performing generation remain unexplored, with classifier-free guidance being the only practical option (Ho & Salimans, 2022). In this paper, we address this shortcoming by proposing two novel methods for combining different models (or the same model with different condition variables) using SUPERDIFF for joint generation.

Limitations. While computationally efficient SUPERDIFF is still limited by the computational budget required to produce the outputs of each model. In particular, the combination at the level of model outputs cannot be simply done via cheap combinations of pre-trained model weights. This, however, invites us to develop a more principled—architecture and training-agnostic—method that does not require any assumptions prior assumptions which makes SUPERDIFF a general purpose method.

Future Work. Our method unlocks novel research directions by allowing for the principled generation of novel samples that were not previously possible to generate easily Sec. 4.2. Furthermore, we argue that the proposed way to estimate the density during the generation enables numerous new potential ways to control the generation process by providing the information about the likelihoods. We discuss the broader impacts of this work in App. E.

540 REPRODUCIBILITY STATEMENT

541
542 To facilitate reproducibility of our empirical results and findings, we intend to make our code
543 publicly available in the final version. We describe all mathematical and algorithmic details necessary
544 to reproduce our results throughout this paper. In [Sec. 3](#) we outline the theoretical basis and
545 mathematical framework for our method. Furthermore, we provide pseudocode for our method
546 in [Algorithm 1](#). For our theoretical contributions, we provide detailed proofs for all theorems and
547 propositions in [App. A](#), [App. B](#) and [App. C](#). We provide experimental details for the CIFAR-10
548 image experiment results in [Sec. 4.1](#). Details regarding experiments for concept interpolation via
549 Stable Diffusion are discussed in [Sec. 4.2](#) and [App. H](#). Experimental details for unconditional protein
550 generation are described in [Sec. 4.3](#) and [App. G.1](#).

551 REFERENCES

- 552 Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf
553 Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure
554 prediction of biomolecular interactions with alphafold 3. *Nature*, pp. 1–3, 2024.
- 555 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
556 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
557 *arXiv preprint arXiv:2303.08774*, 2023.
- 558 Anurag Ajay, Seungwook Han, Yilun Du, Shuang Li, Abhi Gupta, Tommi Jaakkola, Josh Tenenbaum,
559 Leslie Kaelbling, Akash Srivastava, and Pulkit Agrawal. Compositional foundation models for
560 hierarchical planning. *Advances in Neural Information Processing Systems*, 36, 2024.
- 561 Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying
562 framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.
- 563 Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas
564 Geiping, and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the*
565 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 843–852, 2023.
- 566 Helen M Berman, John Westbrook, Zukang Feng, Gary Gilliland, Talapady N Bhat, Helge Weissig,
567 Ilya N Shindyalov, and Philip E Bourne. The protein data bank. *Nucleic acids research*, 28(1):
568 235–242, 2000.
- 569 Benjamin Biggs, Arjun Seshadri, Yang Zou, Achin Jain, Aditya Golatkar, Yusheng Xie, Alessandro
570 Achille, Ashwin Swaminathan, and Stefano Soatto. Diffusion soup: Model merging for text-to-
571 image diffusion models. *arXiv preprint arXiv:2406.08431*, 2024.
- 572 Avishek Joey Bose, Tara Akhound-Sadegh, Guillaume Huguet, Killian Fatras, Jarrid Rector-Brooks,
573 Cheng-Hao Liu, Andrei Cristian Nica, Maksym Korablyov, Michael Bronstein, and Alexander
574 Tong. Se(3)-stochastic flow matching for protein backbone generation. In *The International*
575 *Conference on Learning Representations (ICLR)*, 2024.
- 576 Arwen Bradley and Preetum Nakkiran. Classifier-free guidance is a predictor-corrector. *arXiv*
577 *preprint arXiv:2408.09000*, 2024.
- 578 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe
579 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video
580 generation models as world simulators. 2024. URL [https://openai.com/research/
581 video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators).
- 582 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrike, Eric Horvitz, Ece Kamar,
583 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence:
584 Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- 585 Andrew Campbell, Jason Yim, Regina Barzilay, Tom Rainforth, and Tommi Jaakkola. Generative
586 flows on discrete state-spaces: Enabling multimodal flows with applications to protein co-design.
587 *arXiv preprint arXiv:2402.04997*, 2024.
- 588 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-
589 efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*,
590 35:16344–16359, 2022.

- 594 Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J Ragotte, Lukas F Milles,
595 Basile IM Wicky, Alexis Courbet, Rob J de Haas, Neville Bethel, et al. Robust deep learning-based
596 protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022.
- 597 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*
598 *in neural information processing systems*, 34:8780–8794, 2021.
- 600 Q. Dong, Lei L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, and Z. Sui. A survey for
601 in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- 602 Yilun Du and Leslie Kaelbling. Compositional generative modeling: A single model is not all you
603 need. *arXiv preprint arXiv:2402.01103*, 2024.
- 605 Yilun Du, Shuang Li, and Igor Mordatch. Compositional visual generation with energy based models.
606 *Advances in Neural Information Processing Systems*, 33:6637–6647, 2020.
- 607 Yilun Du, Conor Durkan, Robin Strudel, Joshua B Tenenbaum, Sander Dieleman, Rob Fergus, Jascha
608 Sohl-Dickstein, Arnaud Doucet, and Will Sussman Grathwohl. Reduce, reuse, recycle: Composi-
609 tional generation with energy-based diffusion models and mcmc. In *International conference on*
610 *machine learning*, pp. 8489–8510. PMLR, 2023.
- 612 Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Sean
613 Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, et al. Faith and fate: Limits of
614 transformers on compositionality. *Advances in Neural Information Processing Systems*, 36, 2024.
- 615 Aditya Golatkar, Alessandro Achille, Ashwin Swaminathan, and Stefano Soatto. Training data
616 protection with compositional diffusion models. *arXiv preprint arXiv:2308.01937*, 2023.
- 617 A Herbert and M Sternberg. Maxcluster: A tool for protein structure comparison and clustering,
618 2008.
- 620 Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural*
621 *computation*, 14(8):1771–1800, 2002.
- 622 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,
623 2022.
- 625 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
626 *Neural Information Processing Systems*, 33:6840–6851, 2020.
- 627 Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A.
628 Smith. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question
629 answering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*,
630 pp. 20406–20417, October 2023.
- 631 Guillaume Huguet, James Vuckovic, Kilian Fatras, Eric Thibodeau-Laufer, Pablo Lemos, Riashat
632 Islam, Cheng-Hao Liu, Jarrid Rector-Brooks, Tara Akhound-Sadegh, Michael Bronstein, et al.
633 Sequence-augmented se (3)-flow matching for conditional protein backbone generation. *arXiv*
634 *preprint arXiv:2405.20313*, 2024.
- 636 Michael F Hutchinson. A stochastic estimator of the trace of the influence matrix for laplacian
637 smoothing splines. *Communications in Statistics-Simulation and Computation*, 18(3):1059–1076,
638 1989.
- 639 John B Ingraham, Max Baranov, Zak Costello, Karl W Barber, Wujie Wang, Ahmed Ismail, Vincent
640 Frappier, Dana M Lord, Christopher Ng-Thow-Hing, Erik R Van Vlack, et al. Illuminating protein
641 space with a programmable generative model. *Nature*, 623(7989):1070–1078, 2023.
- 642 Marco Jiralerspong, Avishek Joey Bose, Ian Gemp, Chongli Qin, Yoram Bachrach, and Gauthier
643 Gidel. Feature likelihood score: Evaluating generalization of generative models using samples,
644 2023.
- 646 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott
647 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.
arXiv preprint arXiv:2001.08361, 2020.

- 648 Lingkai Kong, Yuanqi Du, Wenhao Mu, Kirill Neklyudov, Valentin De Bortol, Haorui Wang, Dongxia
649 Wu, Aaron Ferber, Yi-An Ma, Carla P Gomes, et al. Diffusion models as constrained samplers for
650 optimization with unknown constraints. *arXiv preprint arXiv:2402.18012*, 2024.
- 651
652 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph
653 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
654 serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems
655 Principles*, pp. 611–626, 2023.
- 656 Yeqing Lin, Minji Lee, Zhao Zhang, and Mohammed AlQuraishi. Out of many, one: Designing
657 and scaffolding proteins at the scale of the structural universe with genie 2. *arXiv preprint
658 arXiv:2405.15489*, 2024.
- 659 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Allan
660 dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Sal Candido, et al. Language models of
661 protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv*, 2022.
- 662
663 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching
664 for generative modeling. *International Conference on Learning Representations*, 2022.
- 665
666 Sidney Lyayuga Lisanza, Jake Merle Gershon, Sam Tipps, Lucas Arnoldt, Samuel Hendel,
667 Jeremiah Nelson Sims, Xinting Li, and David Baker. Joint generation of protein sequence and
668 structure with rosettafold sequence space diffusion. *bioRxiv*, pp. 2023–05, 2023.
- 669
670 Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual
671 generation with composable diffusion models. In *European Conference on Computer Vision*, pp.
423–439. Springer, 2022a.
- 672
673 Qiang Liu. Rectified flow: A marginal preserving approach to optimal transport. *arXiv preprint
674 arXiv:2209.14577*, 2022.
- 675
676 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and
transfer data with rectified flow. *International Conference on Learning Representations*, 2022b.
- 677
678 Michael Luo, Justin Wong, Brandon Trabucco, Yanping Huang, Joseph E. Gonzalez, Zhifeng Chen,
679 Ruslan Salakhutdinov, and Ion Stoica. Stylus: Automatic adapter selection for diffusion models.
Advances in Neural Information Processing Systems, 2024.
- 680
681 Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and
682 projection for dimension reduction, 2018.
- 683
684 Midjourney. <https://www.midjourney.com/home/>, 2023. Accessed: 2023-09-09.
- 685
686 Weili Nie, Arash Vahdat, and Anima Anandkumar. Controllable and compositional generation
687 with latent-space energy-based models. *Advances in Neural Information Processing Systems*, 34:
13497–13510, 2021.
- 688
689 Bernt Øksendal. Stochastic differential equations. In *Stochastic differential equations*, pp. 65–84.
Springer, 2003.
- 690
691 OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023.
- 692
693 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
694 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
695 models from natural language supervision. In *International conference on machine learning*, pp.
8748–8763. PMLR, 2021.
- 696
697 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen,
698 and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine
699 Learning*. PMLR, 2021.
- 700
701 Milong Ren, Tian Zhu, and Haicang Zhang. Carbonovo: Joint design of protein structure and
sequence using a unified energy-based model. In *Forty-first International Conference on Machine
Learning*, 2024.

- 702 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
703 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Confer-*
704 *ence on Computer Vision and Pattern Recognition (CVPR)*, pp. 10684–10695, June 2022.
- 705
- 706 Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language
707 models a mirage? arxiv. *arXiv preprint arXiv:2304.15004*, 2023.
- 708 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
709 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. LAION-5B:
710 An open large-scale dataset for training next generation image-text models. *Advances in Neural*
711 *Information Processing Systems*, 2022.
- 712
- 713 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
714 Poole. Score-based generative modeling through stochastic differential equations. In *International*
715 *Conference on Learning Representations*, 2020.
- 716 Stability AI. <https://stability.ai/stablediffusion>, 2023. Accessed: 2023-09-09.
- 717
- 718 Jacob Steinhardt. AI Forecasting: One Year In . [https://bounded-regret.ghost.io/](https://bounded-regret.ghost.io/ai-forecasting-one-year-in/)
719 [ai-forecasting-one-year-in/](https://bounded-regret.ghost.io/ai-forecasting-one-year-in/), 2022. Accessed: 2023-09-09.
- 720 Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. Solving olympiad geometry
721 without human demonstrations. *Nature*, 625(7995):476–482, 2024.
- 722
- 723 Brian L. Trippe, Jason Yim, Doug Tischer, David Baker, Tamara Broderick, Regina Barzilay, and
724 Tommi S. Jaakkola. Diffusion probabilistic modeling of protein backbones in 3d for the motif-
725 scaffolding problem. In *The Eleventh International Conference on Learning Representations*, 2023.
726 URL <https://openreview.net/forum?id=6TxBxqNME1Y>.
- 727 Michel Van Kempen, Stephanie S Kim, Charlotte Tumescheit, Milot Mirdita, Jeongjae Lee,
728 Cameron LM Gilchrist, Johannes Söding, and Martin Steinegger. Fast and accurate protein
729 structure search with foldseek. *Nature biotechnology*, 42(2):243–246, 2024.
- 730
- 731 Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computa-*
732 *tion*, 23(7):1661–1674, 2011.
- 733
- 734 Chentong Wang, Yannan Qu, Zhangzhi Peng, Yukai Wang, Hongli Zhu, Dachuan Chen, and Longxing
735 Cao. Proteus: exploring protein structure generation for enhanced designability and efficiency.
736 *bioRxiv*, pp. 2024–02, 2024.
- 737
- 738 Joseph L. Watson, David Juergens, Nathaniel R. Bennett, Brian L. Trippe, Jason Yim, Helen E.
739 Eisenach, Woody Ahern, Andrew J. Borst, Robert J. Ragotte, Lukas F. Milles, Basile I. M.
740 Wicky, Nikita Hanikel, Samuel J. Pellock, Alexis Courbet, William Sheffler, Jue Wang, Preetham
741 Venkatesh, Isaac Sappington, Susana Vázquez Torres, Anna Lauko, Valentin De Bortoli, Emile
742 Mathieu, Sergey Ovchinnikov, Regina Barzilay, Tommi S. Jaakkola, Frank DiMaio, Minkyung
743 Baek, and David Baker. De novo design of protein structure and function with rfdiffusion. *Nature*,
744 620(7976):1089–1100, 2023. ISSN 1476-4687.
- 745
- 746 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong.
747 Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances*
748 *in Neural Information Processing Systems*, 36, 2024.
- 749
- 750 Jason Yim, Andrew Campbell, Andrew Y. K. Foong, Michael Gastegger, José Jiménez-Luna, Sarah
751 Lewis, Victor Garcia Satorras, Bastiaan S. Veeling, Regina Barzilay, Tommi Jaakkola, and Frank
752 Noé. Fast protein backbone generation with se(3) flow matching, 2023a.
- 753
- 754 Jason Yim, Brian L. Trippe, Valentin De Bortoli, Emile Mathieu, Arnaud Doucet, Regina Barzilay,
755 and Tommi Jaakkola. Se (3) diffusion model with application to protein backbone generation.
arXiv preprint arXiv:2302.02277, 2023b.
- 756
- 757 Jason Yim, Andrew Campbell, Emile Mathieu, Andrew YK Foong, Michael Gastegger, José Jiménez-
758 Luna, Sarah Lewis, Victor Garcia Satorras, Bastiaan S Veeling, Frank Noé, et al. Improved
759 motif-scaffolding with se (3) flow matching. *arXiv preprint arXiv:2401.04082*, 2024.

756 Ming Zhong, Yelong Shen, Shuohang Wang, Yadong Lu, Yizhu Jiao, Siru Ouyang, Donghan Yu,
757 Jiawei Han, and Weizhu Chen. Multi-lora composition for image generation. *arXiv preprint*
758 *arXiv:2402.16843*, 2024.
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

APPENDIX

A PROOFS FOR SEC. 2

A.1 PROOF OF PROP. 1

Proposition 1. [Reverse-time SDEs/ODE] *Marginal densities $q_t(x)$ induced by Eq. (1) correspond to the densities induced by the following SDE that goes back in time ($\tau = 1 - t$) with the corresponding initial condition*

$$dx_\tau = \left(-f_t(x_\tau) + \left(\frac{g_t^2}{2} + \xi_\tau \right) \nabla \log q_t(x_\tau) \right) d\tau + \sqrt{2\xi_\tau} d\bar{W}_\tau, \quad x_{\tau=0} \sim q_1(x_0), \quad (2)$$

where \bar{W}_τ is the standard Wiener process in time τ , and ξ_τ is any positive schedule.

For the forward SDE we can write the corresponding Fokker-Planck equation

$$\frac{\partial}{\partial t} q_t(x) = -\langle \nabla, q_t(x) f_t(x) \rangle + \frac{g_t^2}{2} \Delta q_t(x), \quad (19)$$

then for the inverse time $\tau = 1 - t$, we have

$$\frac{\partial}{\partial \tau} q_{1-\tau}(x) = \langle \nabla, q_{1-\tau}(x) f_{1-\tau}(x) \rangle - \frac{g_{1-\tau}^2}{2} \Delta q_{1-\tau}(x) \quad (20)$$

$$= \left\langle \nabla, q_{1-\tau}(x) \left(f_{1-\tau}(x) - \frac{g_{1-\tau}^2}{2} \nabla \log q_{1-\tau}(x) \right) \right\rangle \quad (21)$$

$$= - \left\langle \nabla, q_{1-\tau}(x) \left(-f_{1-\tau}(x) + \left(\frac{g_{1-\tau}^2}{2} + \xi_\tau \right) \nabla \log q_{1-\tau}(x) \right) \right\rangle + \xi_\tau \Delta q_{1-\tau}(x), \quad (22)$$

which corresponds to the SDE

$$dx_\tau = d\tau \left(-f_{1-\tau}(x_\tau) + \left(\frac{g_{1-\tau}^2}{2} + \xi_\tau \right) \nabla \log q_{1-\tau}(x_\tau) \right) + \sqrt{2\xi_\tau} d\bar{W}_\tau. \quad (23)$$

A.2 PROOF OF PROP. 2

Proposition 2. [Ornstein–Uhlenbeck SDE] *The time-dependent densities in Eq. (3) correspond to the marginal densities of the following SDE, with the corresponding initial condition*

$$dx_t = \underbrace{\frac{\partial \log \alpha_t}{\partial t} x_t}_{f_t(x_t)} dt + \underbrace{\sqrt{2\sigma_t^2} \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t}}_{g_t} dW_t, \quad x_0 \sim q_0(x_0). \quad (4)$$

Proof. For individual components q_t^i from Eq. (3), we have to find the vector field satisfying the continuity equation, which we can rewrite as

$$\frac{\partial}{\partial t} \log q_t^i(x) = -\langle \nabla_x, v_t^i(x) \rangle - \langle \nabla_x \log q_t^i(x), v_t^i(x) \rangle. \quad (24)$$

Using the formula for the density of the normal distribution $q_t^i(x) = \mathcal{N}(x | \alpha_t \mu^i, \sigma_t^2)$, we have

$$\nabla_x \log q_t^i(x) = -\frac{1}{\sigma_t^2} (x - \alpha_t \mu^i), \quad (25)$$

$$\frac{\partial}{\partial t} \log q_t^i(x) = d \frac{\partial}{\partial t} \log \sigma_t + \frac{1}{\sigma_t^3} \frac{\partial \sigma_t}{\partial t} \|x - \alpha_t \mu^i\|^2 + \frac{1}{\sigma_t^2} \langle x - \alpha_t \mu^i, \mu^i \rangle \frac{\partial \alpha_t}{\partial t} \quad (26)$$

$$= \left\langle \nabla_x, \frac{\partial \log \sigma_t}{\partial t} x \right\rangle - \left\langle \nabla_x \log q_t^i(x), \underbrace{\frac{\partial \log \sigma_t}{\partial t} (x - \alpha_t \mu^i) + \frac{\partial \alpha_t}{\partial t} \mu^i}_{v_t^i(x)} \right\rangle \quad (27)$$

$$= \langle \nabla_x, v_t^i(x) \rangle - \langle \nabla_x \log q_t^i(x), v_t^i(x) \rangle. \quad (28)$$

For the mixture of densities q_t from Eq. (3), we can use Prop. 3 and write

$$v_t(x) = \frac{1}{q_t(x)N} \sum_{i=1}^N q_t^i(x) \left[\frac{\partial \log \sigma_t}{\partial t} (x - \alpha_t \mu^i) + \frac{\partial \alpha_t}{\partial t} \mu^i \right] \quad (29)$$

$$= \frac{\partial \log \sigma_t}{\partial t} x + \left[\frac{\partial \alpha_t}{\partial t} - \frac{\partial \log \sigma_t}{\partial t} \alpha_t \right] \frac{1}{q_t(x)N} \sum_{i=1}^N q_t^i(x) \mu^i. \quad (30)$$

At the same time

$$\nabla_x \log q_t(x) = \frac{1}{q_t(x)N} \sum_{i=1}^N q_t^i(x) \left[-\frac{1}{\sigma_t^2} (x - \alpha_t \mu^i) \right], \quad (31)$$

$$\frac{1}{q_t(x)N} \sum_{i=1}^N q_t^i(x) \mu^i = \frac{1}{\alpha_t} [\sigma_t^2 \nabla_x \log q_t(x) + x]. \quad (32)$$

Using this formula in Eq. (30), we have

$$v_t(x) = \left[\frac{\partial \log \sigma_t}{\partial t} + \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} - \frac{\partial \log \sigma_t}{\partial t} \right] x + \left[\frac{\partial \alpha_t}{\partial t} - \frac{\partial \log \sigma_t}{\partial t} \alpha_t \right] \frac{\sigma_t^2}{\alpha_t} \nabla_x \log q_t(x) \quad (33)$$

$$= \frac{\partial \log \alpha_t}{\partial t} x - \sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t} \nabla_x \log q_t(x). \quad (34)$$

Hence, we have

$$\frac{\partial}{\partial t} q_t(x) = -\langle \nabla, q_t(x) v_t(x) \rangle = -\left\langle \nabla, q_t(x) \left(\frac{\partial \log \alpha_t}{\partial t} x \right) \right\rangle + \sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t} \Delta q_t(x), \quad (35)$$

which corresponds to the SDE in the statement. \square

A.3 PROOF OF PROP. 3

Proposition 3. [Superposition of ODEs (Liu, 2022)] *The mixture density in Eq. (6) follows the continuity equation with the superposed vector fields from Eq. (5), i.e.*

$$\frac{\partial}{\partial t} q_t^{\text{mix}}(x) = -\langle \nabla_x, q_t^{\text{mix}}(x) v_t(x) \rangle, \quad v_t(x) = \frac{1}{\sum_{j=1}^N \omega^j q_t^j(x)} \sum_{i=1}^N \omega^i q_t^i(x) v_t^i(x). \quad (7)$$

Proof. By the straightforward substitution, we have

$$\frac{\partial}{\partial t} q_t^{\text{mix}}(x) = \sum_{i=1}^m \omega^i \frac{\partial}{\partial t} q_t^i(x) = -\sum_{i=1}^m \omega^i \langle \nabla_x, q_t^i(x) v_t^i(x) \rangle \quad (36)$$

$$= -\left\langle \nabla_x, \sum_{i=1}^m \omega^i q_t^i(x) v_t^i(x) \right\rangle = -\left\langle \nabla_x, \frac{q_t^{\text{mix}}(x)}{q_t^{\text{mix}}(x)} \sum_{i=1}^m \omega^i q_t^i(x) v_t^i(x) \right\rangle \quad (37)$$

$$= -\left\langle \nabla_x, q_t^{\text{mix}}(x) \underbrace{\frac{1}{\sum_{j=1}^m \omega^j q_t^j(x)} \sum_{i=1}^m \omega^i q_t^i(x) v_t^i(x)}_{v_t(x)} \right\rangle. \quad (38)$$

\square

A.4 PROOF OF PROP. 4

Proposition 4. [Superposition of SDEs] *The mixture $q_t^{\text{mix}}(x) := \sum_{i=1}^N \omega^i q_t^i(x)$ of density marginals $\{q_t^i(x)\}_{i=1}^N$ induced by SDEs from Eq. (8) corresponds to the following SDE*

$$dx_\tau = u_\tau(x_\tau) d\tau + g_\tau d\bar{W}_\tau, \quad u_t(x) = \frac{1}{\sum_{j=1}^N \omega^j q_t^j(x)} \sum_{i=1}^N \omega^i q_t^i(x) u_t^i(x). \quad (9)$$

918 *Proof.* By the straightforward substitution, we have

$$919 \frac{\partial}{\partial t} q_t^{\text{mix}}(x) = \sum_{i=1}^m \omega^i \frac{\partial}{\partial t} q_t^i(x) = \sum_{i=1}^m \omega^i \left(-\langle \nabla_x, q_t^i(x) u_t^i(x) \rangle + \frac{g_t^2}{2} \Delta q_t^i(x) \right). \quad (39)$$

922 The first term is analogous to [Prop. 3](#), i.e.

$$923 \sum_{i=1}^m \omega^i \left(-\langle \nabla_x, q_t^i(x) u_t^i(x) \rangle \right) = - \left\langle \nabla_x, q_t^{\text{mix}}(x) \underbrace{\frac{1}{\sum_{j=1}^m \omega^j q_t^j(x)} \sum_{i=1}^m \omega^i q_t^i(x) u_t^i(x)}_{u_t(x)} \right\rangle. \quad (40)$$

927 The second term is

$$928 \sum_{i=1}^m \omega^i \frac{g_t^2}{2} \Delta q_t^i(x) = \frac{g_t^2}{2} \Delta \sum_{i=1}^m \omega^i q_t^i(x) = \frac{g_t^2}{2} \Delta q_t^{\text{mix}}(x). \quad (41)$$

931 This results in the following PDE

$$932 \frac{\partial}{\partial t} q_t^{\text{mix}}(x) = - \left\langle \nabla_x, q_t^{\text{mix}}(x) \underbrace{\frac{1}{\sum_{j=1}^m \omega^j q_t^j(x)} \sum_{i=1}^m \omega^i q_t^i(x) u_t^i(x)}_{u_t(x)} \right\rangle + \frac{g_t^2}{2} \Delta q_t^{\text{mix}}(x). \quad (42)$$

936 □

938 B FLOWS AND DIFFUSION AS A SUPERPOSITION OF ELEMENTARY VECTOR 939 FIELDS

941 From [Prop. 3](#), one can immediately get the main principle for the simulation-free training of generative
942 flow models, as demonstrated in the following proposition.

944 **Proposition 7.** [Superposition of L_2 -losses] *For the parametric vector field model $v_t(x; \theta)$
945 with parameters θ , the L_2 -loss for the vector field from [Prop. 3](#) can be decomposed into the
946 losses for the vector fields from [Eq. \(5\)](#), i.e.*

$$947 \int_0^1 dt \mathbb{E}_{q_t^{\text{mix}}(x)} \|v_t(x) - v_t(x; \theta)\|^2 = \sum_{i=1}^N \omega_i \int_0^1 dt \mathbb{E}_{q_t^i(x)} \|v_t^i(x) - v_t(x; \theta)\|^2 + \text{constant},$$

950 where the constant does not depend on the parameters θ .

953 *Proof.* By the straightforward calculation, we have

$$954 \int_0^1 dt \mathbb{E}_{q_t^{\text{mix}}(x)} \|v_t(x) - v_t(x; \theta)\|^2 = \int_0^1 dt \mathbb{E}_{q_t^{\text{mix}}(x)} \left[\|v_t(x)\|^2 - 2\langle v_t(x), v_t(x; \theta) \rangle + \|v_t(x; \theta)\|^2 \right],$$

957 where the first term is constant w.r.t. θ and the last term is amenable for a straightforward estimation.

958 The middle term, according to [Prop. 3](#), is

$$959 \int dx q_t^{\text{mix}}(x) \langle v_t(x), v_t(x; \theta) \rangle = \int dx \sum_{i=1}^m \omega^i q_t^i(x) \langle v_t^i(x), v_t(x; \theta) \rangle. \quad (43)$$

962 Thus, we have

$$963 \int_0^1 dt \mathbb{E}_{q_t^{\text{mix}}(x)} \left[\|v_t(x)\|^2 - 2\langle v_t(x), v_t(x; \theta) \rangle + \|v_t(x; \theta)\|^2 \right] \quad (44)$$

$$964 = \int_0^1 dt \left[\underbrace{\mathbb{E}_{q_t^{\text{mix}}(x)} \|v_t(x)\|^2 - \sum_{i=1}^m \omega^i \mathbb{E}_{q_t^i(x)} \|v_t^i(x)\|^2}_{\text{constant}} + \sum_{i=1}^m \omega^i \mathbb{E}_{q_t^i(x)} \|v_t^i(x)\|^2 - \right] \quad (45)$$

$$966 - 2 \sum_{i=1}^m \omega^i \mathbb{E}_{q_t^i(x)} \langle v_t^i(x), v_t(x; \theta) \rangle + \sum_{i=1}^m \omega^i \mathbb{E}_{q_t^i(x)} \|v_t(x; \theta)\|^2 \quad (46)$$

$$= \int_0^1 dt \sum_{i=1}^m \omega^i \mathbb{E}_{q_t^i(x)} \left\| v_t^i(x) - v_t(x; \theta) \right\|^2 + \text{constant}. \quad (47)$$

□

The proof for proposition follows a simple extension of [Liu \(2022\)](#). Consequently, the generative modeling problem explicitly boils down to learning a time-dependent vector field $v_t(x, \theta)$. Moreover, for the generative modeling task, the noising process is required to satisfy the following boundary conditions $q_0^i(x) = \delta(x - \mu^i)$ and $q_1^i(x) = p_{\text{noise}}(x)$, $\forall i \in \mathcal{D}$. Generating samples that resemble $p_{\text{data}}(\mu)$ during inference then simply amounts to solving the reverse-time ODE with the learned model $v_t(x; \theta)$ starting from a noisy sample $x_1 \sim p_{\text{noise}}(x)$.

Diffusion models can be incorporated into the flow-based framework by selecting a Gaussian forward process which gives the marginal density at a timestep t as the mixture of corresponding Gaussians,

$$q_t^i(x) = \mathcal{N}(x | \alpha_t \mu^i, \sigma_t^2), \quad q_t(x) = \frac{1}{N} \sum_{i=1}^N q_t^i(x). \quad (48)$$

Under this Gaussian diffusion framework, we can express the vector fields analytically using the score, which is formalized by the following proposition.

Proposition 8. [Diffusion processes] *The time-dependent densities in [Eq. \(48\)](#) satisfy the following continuity equations*

$$\frac{\partial}{\partial t} q_t^i(x) = - \langle \nabla_x, q_t^i(x) v_t^i(x) \rangle, \quad v_t^i(x) = \frac{\partial \log \sigma_t}{\partial t} (x - \alpha_t \mu^i) + \frac{\partial \alpha_t}{\partial t} \mu^i, \quad (49)$$

$$\frac{\partial}{\partial t} q_t(x) = - \langle \nabla_x, q_t(x) v_t(x) \rangle, \quad v_t(x) = \frac{\partial \log \alpha_t}{\partial t} x - \left(\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t} \right) \nabla_x \log q_t(x). \quad (50)$$

The proof repeats the proof for [Prop. 2](#) in [App. A.2](#) since [Eq. \(50\)](#) corresponds to the Ornstein-Uhlenbeck process:

$$dx_t = \frac{\partial \log \alpha_t}{\partial t} x_t dt + \sqrt{2 \left(\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t} \right)} dW_t. \quad (51)$$

Diffusion models using a Gaussian perturbation kernel enjoy the benefit of giving an exact expression for the Stein score of the perturbed data which leads to the celebrated denoising score matching objective ([Vincent, 2011](#); [Ho et al., 2020](#)). Analogously to [Prop. 7](#) we can write the denoising score matching objective as the superposition of the scores as demonstrated in the following proposition.

Proposition 9. [Denoising Score Matching ([Vincent, 2011](#))] *For the parametric score model $\nabla \log q_t(x; \theta)$ with parameters θ , the score matching objective can be decomposed into the corresponding objectives for the individual scores, i.e.*

$$\begin{aligned} & \int_0^1 dt \mathbb{E}_{q_t(x)} \left\| \nabla \log q_t(x) - \nabla \log q_t(x; \theta) \right\|^2 = \\ & = \frac{1}{N} \sum_{i=1}^N \int_0^1 dt \mathbb{E}_{q_t^i(x)} \left\| \nabla \log q_t^i(x) - \nabla \log q_t(x; \theta) \right\|^2 + \text{constant}, \end{aligned}$$

where the constant does not depend on the parameters θ and $\nabla \log q_t^i(x) = -\frac{1}{\sigma_t^2} (x - \alpha_t \mu^i)$.

Proof. From [Prop. 8](#), we have

$$v_t^i(x) = \frac{\partial \log \sigma_t}{\partial t} (x - \alpha_t \mu^i) + \frac{\partial \alpha_t}{\partial t} \mu^i \quad (52)$$

$$= \frac{\partial \log \sigma_t}{\partial t} (x - (x + \sigma_t^2 \nabla \log q_t^i(x))) + \frac{\partial \alpha_t}{\partial t} \frac{1}{\alpha_t} (x + \sigma_t^2 \nabla \log q_t^i(x)) \quad (53)$$

$$= \frac{\partial \log \alpha_t}{\partial t} x + \sigma_t^2 \frac{\partial}{\partial t} \log \frac{\alpha_t}{\sigma_t} \nabla_x \log q_t^i(x) \quad (54)$$

$$v_t(x) = \frac{\partial \log \alpha_t}{\partial t} x + \sigma_t^2 \frac{\partial}{\partial t} \log \frac{\alpha_t}{\sigma_t} \nabla_x \log q_t(x). \quad (55)$$

Applying the same change of variables to the parametric model

$$v_t(x; \theta) = \frac{\partial \log \alpha_t}{\partial t} x + \sigma_t^2 \frac{\partial}{\partial t} \log \frac{\alpha_t}{\sigma_t} \nabla_x \log q_t(x; \theta), \quad (56)$$

and using Prop. 7, we have

$$\int_0^1 dt \mathbb{E}_{q_t(x)} \|v_t(x) - v_t(x; \theta)\|^2 = \sum_{i=1}^m \frac{1}{N} \int_0^1 dt \mathbb{E}_{q_t^i(x)} \|v_t^i(x) - v_t(x; \theta)\|^2 + \text{constant}, \quad (57)$$

$$\int_0^1 dt \mathbb{E}_{q_t(x)} \|\nabla_x \log q_t(x) - \nabla_x \log q_t(x; \theta)\|^2 \quad (58)$$

$$= \sum_{i=1}^m \frac{1}{N} \int_0^1 dt \mathbb{E}_{q_t^i(x)} \|\nabla_x \log q_t^i(x) - \nabla_x \log q_t(x; \theta)\|^2 + \frac{\text{constant}}{\left(\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\alpha_t}{\sigma_t}\right)}, \quad (59)$$

which concludes the proof. \square

C DENSITY ESTIMATORS

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

Proposition 5. [Smooth density estimator] *For the integral curve $x(t)$ solving $dx/dt = u_t(x_t)$, and the density $q_t^i(x(t))$ satisfying the continuity equation $\frac{\partial}{\partial t} q_t^i(x) = -\langle \nabla_x, q_t^i(x) v_t^i(x) \rangle$, the log-density along the curve changes according to the following ODE*

$$\frac{d}{dt} \log q_t^i(x(t)) = -\langle \nabla_x, v_t^i(x) \rangle - \langle \nabla_x \log q_t^i(x), v_t^i(x) - u_t(x) \rangle. \quad (10)$$

Proof. By straightforward computation, we have

$$\frac{d}{dt} \log q_t^i(x(t)) = \frac{\partial}{\partial t} \log q_t^i(x) + \left\langle \nabla_x \log q_t^i(x), \frac{dx}{dt} \right\rangle, \quad (60)$$

and using the continuity equation $\frac{\partial}{\partial t} q_t^i(x) = -\langle \nabla_x, q_t^i(x) v_t^i(x) \rangle$, we have

$$\frac{\partial}{\partial t} \log q_t^i(x) = -\langle \nabla_x, v_t^i(x) \rangle - \langle \nabla_x \log q_t^i(x), v_t^i(x) \rangle, \quad (61)$$

$$\frac{d}{dt} \log q_t^i(x(t)) = -\langle \nabla_x, v_t^i(x) \rangle - \langle \nabla_x \log q_t^i(x), v_t^i(x) - u_t(x) \rangle. \quad (62)$$

\square

Theorem 1. [$\hat{\text{O}}\text{ti}$ density estimator] *Consider time-dependent density $q_t(x)$ induced by the marginals of the following “forward” process*

$$dx_t = f_t(x_t)dt + g_t dW_t, \quad x_{t=0} \sim q_0(x), \quad t \in [0, 1], \quad (11)$$

where dW_t is the Wiener process. Then, for the following “backward” SDE (with $\tau = 1 - t$)

$$dx_\tau = u_\tau(x_\tau)d\tau + g_{1-\tau} d\bar{W}_\tau, \quad \tau \in [0, 1], \quad (12)$$

the change of the log-density $\log q_\tau(x_\tau)$ follows the following SDE

$$\begin{aligned} d \log q_{1-\tau}(x_\tau) &= \langle dx_\tau, \nabla \log q_{1-\tau}(x_\tau) \rangle + \left(\langle \nabla, f_{1-\tau}(x_\tau) \rangle + \right. \\ &\quad \left. + \left\langle f_{1-\tau}(x_\tau) - \frac{g_{1-\tau}^2}{2} \nabla \log q_{1-\tau}(x_\tau), \nabla \log q_{1-\tau}(x_\tau) \right\rangle \right) d\tau. \end{aligned} \quad (13)$$

1080 *Proof.* The Fokker-Planck equation describing the evolution of the marginals for the forward process
1081 is

$$1082 \frac{\partial q_t(x)}{\partial t} = -\langle \nabla, f_t(x)q_t(x) \rangle + \frac{g_t^2}{2} \Delta q_t(x), \quad (63)$$

1084 hence, for the inverse time $\tau = 1 - t$, we have

$$1085 \frac{\partial}{\partial \tau} q_{1-\tau}(x) = \langle \nabla, f_{1-\tau}(x)q_{1-\tau}(x) \rangle - \frac{g_{1-\tau}^2}{2} \Delta q_{1-\tau}(x) \quad (64)$$

$$1087 \frac{\partial}{\partial \tau} \log q_{1-\tau}(x) = \langle \nabla, f_{1-\tau}(x) \rangle + \langle \nabla \log q_{1-\tau}(x), f_{1-\tau}(x) \rangle - \quad (65)$$

$$1088 - \frac{g_{1-\tau}^2}{2} \Delta \log q_{1-\tau}(x) - \frac{g_{1-\tau}^2}{2} \|\nabla \log q_{1-\tau}(x)\|^2.$$

1091 For the following reverse-time SDE

$$1092 dx_\tau = u_\tau(x_\tau)d\tau + g_{1-\tau}d\bar{W}_\tau, \quad (66)$$

1094 using Itô's lemma for $\log q_{1-\tau}(x_\tau)$, we have

$$1095 d \log q_{1-\tau}(x_\tau) = \left(\frac{\partial}{\partial \tau} \log q_{1-\tau}(x_\tau) + \langle \nabla \log q_{1-\tau}(x_\tau), u_\tau(x_\tau) \rangle + \quad (67)$$

$$1096 + \frac{g_{1-\tau}^2}{2} \Delta \log q_{1-\tau}(x_\tau) \right) d\tau + g_{1-\tau} \langle \nabla \log q_{1-\tau}(x_\tau), d\bar{W}_\tau \rangle. \quad (68)$$

1100 Using Eq. (65), we have

$$1101 d \log q_{1-\tau}(x_\tau) = \left(\left\langle \nabla \log q_{1-\tau}(x_\tau), f_{1-\tau}(x_\tau) - \frac{g_{1-\tau}^2}{2} \nabla \log q_{1-\tau}(x_\tau) \right\rangle + \quad (69)$$

$$1102 + \langle \nabla, f_{1-\tau}(x_\tau) \rangle \right) d\tau + \left\langle \nabla \log q_{1-\tau}(x_\tau), \underbrace{u_\tau(x_\tau)d\tau + g_{1-\tau}d\bar{W}_\tau}_{dx_\tau} \right\rangle. \quad (70)$$

1107 □

1108 C.1 PROOF OF PROP. 6

1110 **Proposition 6.** [Density control] *For the SDE*

$$1111 dx_\tau = \sum_{j=1}^M \kappa_j u_\tau^j(x_\tau) d\tau + g_{1-\tau} d\bar{W}_\tau, \quad (16)$$

1113 where κ are the weights of different models and $\sum_j \kappa_j = 1$, one can find κ that satisfies

$$1114 d \log q_{1-\tau}^i(x_\tau) = d \log q_{1-\tau}^j(x_\tau), \quad \forall i, j \in [M], \quad (17)$$

1115 by solving a system of M linear equations w.r.t. κ .

1118 *Proof.* Using the result of [Thm. 1](#), for the density change of every model, we have

$$1119 d \log q_{1-\tau}^i(x_\tau) = \langle dx_\tau, \nabla \log q_{1-\tau}^i(x_\tau) \rangle + \left(\langle \nabla, f_{1-\tau}(x_\tau) \rangle + \quad (71)$$

$$1120 + \left\langle f_{1-\tau}(x_\tau) - \frac{g_{1-\tau}^2}{2} \nabla \log q_{1-\tau}^i(x_\tau), \nabla \log q_{1-\tau}^i(x_\tau) \right\rangle \right) d\tau \quad (72)$$

$$1121 = \sum_{j=1}^M \kappa_j d\tau \langle u_\tau^j(x_\tau), \nabla \log q_{1-\tau}^i(x_\tau) \rangle + \left(\langle \nabla, f_{1-\tau}(x_\tau) \rangle + \quad (73)$$

$$1122 + \left\langle g_{1-\tau} d\bar{W}_\tau + \left(f_{1-\tau}(x_\tau) - \frac{g_{1-\tau}^2}{2} \nabla \log q_{1-\tau}^i(x_\tau) \right) d\tau, \nabla \log q_{1-\tau}^i(x_\tau) \right\rangle \right). \quad (74)$$

Let us denote

$$a_{ij} = d\tau \langle u_\tau^j(x_\tau), \nabla \log q_{1-\tau}^i(x_\tau) \rangle, \quad (75)$$

$$b_i = \langle \nabla, f_{1-\tau}(x_\tau) \rangle + \left\langle g_{1-\tau} d\bar{W}_\tau + \left(f_{1-\tau}(x_\tau) - \frac{g_{1-\tau}^2}{2} \nabla \log q_{1-\tau}^i(x_\tau) \right) d\tau, \nabla \log q_{1-\tau}^i(x_\tau) \right\rangle. \quad (76)$$

Then the change of densities can be written as the following linear transformation

$$\begin{bmatrix} d \log q_{1-\tau}^1(x_\tau) \\ \vdots \\ d \log q_{1-\tau}^i(x_\tau) \\ \vdots \\ d \log q_{1-\tau}^M(x_\tau) \end{bmatrix} = A\kappa + b. \quad (77)$$

Solving

$$A\kappa + b = \begin{bmatrix} 1 \\ \vdots \\ 1 \\ \vdots \\ 1 \end{bmatrix}, \quad (78)$$

and normalizing kappas so we have $\sum_j \kappa_j = 1$, we get the solution. \square

D DISCRETE-TIME PERSPECTIVE ON THM. 1

In this section we derive [Thm. 1](#) from another perspective by operating with discrete time transition kernels and the detailed balance. Namely, we aim to compute the marginal density by starting from the detailed balance equation, which states the equivalence of the following joint densities:

$$q_t(y)k_{\Delta t}(x|y) = q_{t+\Delta t}(x)r_{\Delta t}(y|x). \quad (79)$$

The detailed balance condition simply states the pair (x, y) can be sampled in two equivalent ways: 1.) sampling y from the marginal $q_t(y)$ and then sampling x via the noising kernel $k_{\Delta t}(x|y)$, or 2.) sampling x from $q_{t+\Delta t}(y)$ and then denoising it via $r_{\Delta t}(y|x)$. Indeed, there are infinitely many valid kernels that may satisfy the detailed balance; we make a specific choice informed by the following principle. We aim to construct a universal noising kernel that is independent of the data distribution or other densities, i.e. $p_{\text{data}}(\cdot), q_t(\cdot), q_{t+\Delta t}(\cdot)$. Remarkably, this principle results in a unique kernel choice, which we formalize in the following proposition.

Proposition 10. [Universal noising kernel] *For any data density $p_{\text{data}}(\mu)$, for the continuous noising process $q_t(x) = \int d\mu \mathcal{N}(x | \alpha_t \mu, \sigma_t^2) p_{\text{data}}(\mu)$, there exists unique noising kernel*

$$k_{\Delta t}(x|y) = \mathcal{N}\left(x \left| \frac{\alpha_{t+\Delta t}}{\alpha_t} y, S_{t+\Delta t}^2\right.\right), q_{t+\Delta t}(x) = \int dy k_{\Delta t}(y|x)q_t(y), \quad (80)$$

where $S_{t+\Delta t}^2 = \sigma_{t+\Delta t}^2 - \sigma_t^2 \frac{\alpha_{t+\Delta t}^2}{\alpha_t^2}$ and it is independent of the densities $p_{\text{data}}(\cdot), q_t(\cdot), q_{t+\Delta t}(\cdot)$.

The proof for this proposition is presented in [App. D.1](#). Note that the proposition holds exactly for any Δt —i.e. there are no any assumptions on the scale of Δt or any approximations.

Once the noising kernel is fixed, the detailed balance [Eq. \(79\)](#) uniquely defines the denoising kernel $r_{\Delta t}(y|x)$ that propagates samples back in time. However, the analytic form of this kernel depends on the densities $q_{t+\Delta t}(x)$ and $q_t(y)$, which are unavailable for the modern large-scale diffusion models that are served using limited API endpoints. Therefore, we consider the Gaussian approximation of the reverse kernel and justify its applicability in the following theorem.

Theorem 2. [Denoising kernel] *For the universal noising kernel $k_{\Delta t}(x|y)$, the Gaussian approximation of the corresponding reverse kernel $r_{\Delta t}(y|x) = k_{\Delta t}(x|y)q_t(y)/q_{t+\Delta t}(x)$ is*

$$\tilde{r}_{\Delta t}(y|x) = \mathcal{N}\left(y \left| \frac{\alpha_t}{\alpha_{t+\Delta t}} x + \frac{\alpha_t}{\alpha_{t+\Delta t}} S_{t+\Delta t}^2 \nabla \log q_{t+\Delta t}(x), \left(\frac{\alpha_t}{\alpha_{t+\Delta t}} S_{t+\Delta t}\right)^2\right.\right), \quad (81)$$

which corresponds to a single step of the Euler discretization of the following SDE

$$dx_\tau = -\left(\frac{\partial}{\partial t} \log \alpha_t x - 2\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t} \nabla_x \log q_t(x)\right) \cdot d\tau + \sqrt{2\sigma_t^2} \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t} dW_\tau, \quad (82)$$

and

$$D_{\text{KL}}(\tilde{r}_{\Delta t}(y|x)||r_{\Delta t}(y|x)) = o(\Delta t). \quad (83)$$

See proof in [App. D.3](#). Given the noising kernel from [Prop. 10](#) and the approximation of the reverse kernel from [Thm. 2](#), we propose to estimate the log marginal density of the current sample using the detailed balance as follows,

$$\log q_t(y) - \log q_{t+\Delta t}(x) = -\log k_{\Delta t}(x|y) + \log \tilde{r}_{\Delta t}(y|x) + \log \frac{r_{\Delta t}(y|x)}{\tilde{r}_{\Delta t}(y|x)}, \quad (84)$$

where the last term involves the unknown reversed kernel. This term can typically be approximated by taking the divergence of the score (see [Lemma 1](#)), but due to the computational challenge of estimating the divergence we opt for a different strategy. Instead, in the following theorem, we study the distribution of the last term for different samples y generated from x and argue that this term can be ignored.

Theorem 3. [Detailed balance density estimator error] For $y = x + \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon$, where $g_{t+\Delta t} = \sqrt{2\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t}}$ and ε is a standard normal random variable, the following estimator is unbiased w.r.t. the generated samples y , i.e.

$$\mathbb{E}_\varepsilon \left[\log \frac{r_{\Delta t}(y|x)}{\tilde{r}_{\Delta t}(y|x)} \right] = o(\Delta t), \quad (85)$$

and the variance is

$$\mathbf{D}_\varepsilon \left[\log \frac{r_{\Delta t}(y|x)}{\tilde{r}_{\Delta t}(y|x)} \right] = \Delta t^2 \frac{g_{t+\Delta t}^4}{4} \mathbf{D}_\varepsilon \left(\varepsilon^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \varepsilon \right) + o(\Delta t^2), \quad (86)$$

where $\frac{\partial^2 \log q_{t+\Delta t}(x)}{\partial x^2}$ is the Hessian matrix of the log-density.

We postpone the proof of this theorem until [App. D.4](#). Note that the last term in [Eq. \(84\)](#) has zero expectation and the same variance as Hutchinson’s trace estimator. Thus, we argue that the following estimator is as efficient as the conventional divergence estimator with Hutchinson’s trick

$$\log q_t(y) - \log q_{t+\Delta t}(x) \approx -\log k_{\Delta t}(x|y) + \log \tilde{r}_{\Delta t}(y|x). \quad (87)$$

We further highlight that since the generation via the inference process happens in many steps with small Δt , we approximate the relation [Eq. \(87\)](#) up to the second order terms in Δt .

Proposition 11. [Recurrence relation for the log-density] For $y = x + \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon$, where $g_{t+\Delta t} = \sqrt{2\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t}}$ and ε is a standard normal random variable, the estimator from [Eq. \(87\)](#) can be expanded as follows

$$\begin{aligned} -\log k_{\Delta t}(x|y) + \log \tilde{r}_{\Delta t}(y|x) &= d\Delta t \frac{\partial}{\partial t} \log \alpha_{t+\Delta t} - \Delta t \frac{g_{t+\Delta t}^2}{2} \|\nabla \log q_{t+\Delta t}(x)\|^2 + \\ &+ \left\langle \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon + \Delta t \frac{\partial}{\partial t} \log \alpha_{t+\Delta t} x, \nabla \log q_{t+\Delta t}(x) \right\rangle + o(\Delta t) \end{aligned} \quad (88)$$

The proof for this proposition is included in [App. D.5](#). For practical use cases we note that [Eq. \(88\)](#) is the final form of our density estimator.

D.1 PROOF OF [PROP. 10](#)

Proposition 10. [Universal noising kernel] For any data density $p_{\text{data}}(\mu)$, for the continuous noising process $q_t(x) = \int d\mu \mathcal{N}(x|\alpha_t \mu, \sigma_t^2) p_{\text{data}}(\mu)$, there exists unique noising kernel

$$k_{\Delta t}(x|y) = \mathcal{N}\left(x \left| \frac{\alpha_{t+\Delta t}}{\alpha_t} y, S_{t+\Delta t}^2 \right.\right), q_{t+\Delta t}(x) = \int dy k_{\Delta t}(y|x) q_t(y), \quad (80)$$

where $S_{t+\Delta t}^2 = \sigma_{t+\Delta t}^2 - \sigma_t^2 \frac{\alpha_{t+\Delta t}^2}{\alpha_t^2}$ and it is independent of the densities $p_{\text{data}}(\cdot)$, $q_t(\cdot)$, $q_{t+\Delta t}(\cdot)$.

Proof.

$$q_t(x) = \int d\mu \mathcal{N}(x | \alpha_t \mu, \sigma_t^2) p_{\text{data}}(\mu) \quad (89)$$

Let's derive the incremental kernel $k_{\Delta t}(x | y)$ from this formula, i.e.

$$q_{t+\Delta t}(x) = \int dy k_{\Delta t}(x | y) q_t(y) \quad (90)$$

$$\int d\mu \mathcal{N}(x | \alpha_{t+\Delta t} \mu, \sigma_{t+\Delta t}^2) p_{\text{data}}(\mu) = \int dy k_{\Delta t}(x | y) \int d\mu \mathcal{N}(y | \alpha_t \mu, \sigma_t^2) p_{\text{data}}(\mu) \quad (91)$$

$$\int d\mu \left(\mathcal{N}(x | \alpha_{t+\Delta t} \mu, \sigma_{t+\Delta t}^2) - \int dy k_{\Delta t}(x | y) \mathcal{N}(y | \alpha_t \mu, \sigma_t^2) \right) p_{\text{data}}(\mu) = 0 \quad (92)$$

$$\mathcal{N}(x | \alpha_{t+\Delta t} \mu, \sigma_{t+\Delta t}^2) = \int dy k_{\Delta t}(x | y) \mathcal{N}(y | \alpha_t \mu, \sigma_t^2) \quad (93)$$

$$x = \alpha_{t+\Delta t} \mu + \frac{y - \alpha_t \mu}{\sigma_t} S + R \varepsilon, \text{ where } y \sim \mathcal{N}(y | \alpha_t \mu, \sigma_t^2), \text{ and } \varepsilon \sim \mathcal{N}(\varepsilon | 0, 1) \quad (94)$$

$$S^2 + R^2 = \sigma_{t+\Delta t}^2 \quad (95)$$

To make the kernel independent of μ we have to choose $R = \sigma_t \frac{\alpha_{t+\Delta t}}{\alpha_t}$, then the kernel is

$$k_{\Delta t}(x | y) = \mathcal{N} \left(x \mid \frac{\alpha_{t+\Delta t}}{\alpha_t} y, \underbrace{\sigma_{t+\Delta t}^2 - \sigma_t^2 \frac{\alpha_{t+\Delta t}^2}{\alpha_t^2}}_{S_{t+\Delta t}^2} \right). \quad (96)$$

□

D.2 LEMMA 1

Lemma 1. [Reverse kernel lemma] *For the universal noising kernel $k_{\Delta t}(x | y)$, the corresponding reverse kernel $r_{\Delta t}(y | x) = k_{\Delta t}(x | y) q_t(y) / q_{t+\Delta t}(x)$, and its Gaussian approximation*

$$\tilde{r}_{\Delta t}(y | x) = \mathcal{N} \left(y \mid \frac{\alpha_t}{\alpha_{t+\Delta t}} x + \frac{\alpha_t}{\alpha_{t+\Delta t}} S_{t+\Delta t}^2 \nabla \log q_{t+\Delta t}(x), \left(\frac{\alpha_t}{\alpha_{t+\Delta t}} S_{t+\Delta t} \right)^2 \right), \quad (97)$$

we have

$$\log \frac{r_{\Delta t}(x + \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon | x)}{\tilde{r}_{\Delta t}(x + \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon | x)} = \quad (98)$$

$$= \Delta t \cdot \frac{g_{t+\Delta t}^2}{2} \left(\Delta \log q_{t+\Delta t}(x) - \varepsilon^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \varepsilon \right) + o(\Delta t), \quad (99)$$

where $g_{t+\Delta t} = \sqrt{2\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}}}$.

Proof. From the detailed balance equation, we have

$$\log r_{\Delta t}(y | x) = \log k_{\Delta t}(x | y) + \log q_t(y) - \log q_{t+\Delta t}(x) \quad (100)$$

$$= \frac{-1}{2S_{t+\Delta t}^2} \left(x - \frac{\alpha_{t+\Delta t}}{\alpha_t} y \right)^2 - \frac{d}{2} \log 2\pi S_{t+\Delta t}^2 + \log q_t(y) - \log q_{t+\Delta t}(x), \quad (101)$$

where we use the definition of the forward kernel. Thus, the difference between kernels is

$$\log \frac{r_{\Delta t}(y|x)}{\tilde{r}_{\Delta t}(y|x)} = \frac{-1}{2S_{t+\Delta t}^2} \left(x - \frac{\alpha_{t+\Delta t}}{\alpha_t} y \right)^2 - \frac{d}{2} \log 2\pi S_{t+\Delta t}^2 + \log q_t(y) - \log q_{t+\Delta t}(x) + \quad (102)$$

$$+ \frac{1}{2(S_{t+\Delta t} \frac{\alpha_t}{\alpha_{t+\Delta t}})^2} \left(y - \frac{\alpha_t}{\alpha_{t+\Delta t}} x - \frac{\alpha_t}{\alpha_{t+\Delta t}} S_{t+\Delta t}^2 \nabla \log q_{t+\Delta t}(x) \right)^2 + \quad (103)$$

$$+ \frac{d}{2} \log 2\pi S_{t+\Delta t}^2 + \frac{d}{2} \log \left(\frac{\alpha_t}{\alpha_{t+\Delta t}} \right)^2. \quad (104)$$

Opening the brackets, we have

$$\log \frac{r_{\Delta t}(y|x)}{\tilde{r}_{\Delta t}(y|x)} = \frac{-1}{2S_{t+\Delta t}^2} \left(x - \frac{\alpha_{t+\Delta t}}{\alpha_t} y \right)^2 + \log q_t(y) - \log q_{t+\Delta t}(x) + \quad (105)$$

$$+ \frac{1}{2S_{t+\Delta t}^2} \left(\frac{\alpha_{t+\Delta t}}{\alpha_t} y - x \right)^2 - \left\langle \frac{\alpha_{t+\Delta t}}{\alpha_t} y - x, \nabla \log q_{t+\Delta t}(x) \right\rangle + \quad (106)$$

$$+ \frac{S_{t+\Delta t}^2}{2} \|\nabla \log q_{t+\Delta t}(x)\|^2 + \frac{d}{2} \log \left(\frac{\alpha_t}{\alpha_{t+\Delta t}} \right)^2. \quad (107)$$

The constants can be estimated as follows

$$S_{t+\Delta t}^2 = \sigma_{t+\Delta t}^2 - \sigma_t^2 \frac{\alpha_{t+\Delta t}^2}{\alpha_t^2} = \sigma_{t+\Delta t}^2 - (\sigma_{t+\Delta t}^2 - 2dt\sigma_{t+\Delta t} \frac{\partial \sigma_{t+\Delta t}}{\partial t} + o(\Delta t)) \frac{\alpha_{t+\Delta t}^2}{\alpha_t^2} \quad (108)$$

$$= \sigma_{t+\Delta t}^2 - (\sigma_{t+\Delta t}^2 - 2dt\sigma_{t+\Delta t} \frac{\partial \sigma_{t+\Delta t}}{\partial t} + o(\Delta t)) \left(1 + 2dt \frac{\alpha_{t+\Delta t}^2}{\alpha_{t+\Delta t}^3} \frac{\partial \alpha_{t+\Delta t}}{\partial t} + o(\Delta t) \right) \quad (109)$$

$$= 2dt\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}} + o(\Delta t), \quad (110)$$

and

$$\frac{\alpha_{t+\Delta t}}{\alpha_t} = 1 - \frac{\alpha_{t+\Delta t}}{\alpha_t^2} \frac{\partial \alpha_{t+\Delta t}}{\partial t} (-\Delta t) + o(\Delta t) = 1 + \Delta t \frac{\partial \log \alpha_{t+\Delta t}}{\partial t} + o(\Delta t) \quad (111)$$

$$\frac{d}{2} \log \left(\frac{\alpha_t}{\alpha_{t+\Delta t}} \right)^2 = -dtd \frac{\partial}{\partial t} \log \alpha_{t+\Delta t} + o(\Delta t). \quad (112)$$

Thus, we have

$$\log \frac{r_{\Delta t}(y|x)}{\tilde{r}_{\Delta t}(y|x)} = \log q_t(y) - \log q_{t+\Delta t}(x) + \Delta t \sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}} \|\nabla \log q_{t+\Delta t}(x)\|^2 - \quad (113)$$

$$-dtd \frac{\partial}{\partial t} \log \alpha_{t+\Delta t} - \left\langle (1 + \Delta t \frac{\partial \log \alpha_{t+\Delta t}}{\partial t}) y - x, \nabla \log q_{t+\Delta t}(x) \right\rangle + o(\Delta t). \quad (114)$$

From Prop. 8, the time-derivative of the density is

$$\frac{\partial}{\partial t} \log q_t(x) = - \langle \nabla, v_t(x) \rangle - \langle \nabla \log q_t(x), v_t(x) \rangle \quad (115)$$

$$= -d \frac{\partial \log \alpha_t}{\partial t} + \left(\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t} \right) \Delta \log q_t(x) - \quad (116)$$

$$- \frac{\partial \log \alpha_t}{\partial t} \langle \nabla \log q_t(x), x \rangle + \left(\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t} \right) \|\nabla \log q_t(x)\|^2. \quad (117)$$

Hence, we have

$$\log q_t(x) - \log q_{t+\Delta t}(x) = -\Delta t \frac{\partial}{\partial t} \log q_{t+\Delta t}(x) + o(\Delta t) \quad (118)$$

$$= \Delta t \frac{\partial \log \alpha_{t+\Delta t}}{\partial t} - \Delta t \left(\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}} \right) \Delta \log q_{t+\Delta t}(x) + \quad (119)$$

$$+ \Delta t \frac{\partial \log \alpha_{t+\Delta t}}{\partial t} \langle \nabla \log q_{t+\Delta t}(x), x \rangle - \Delta t \left(\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}} \right) \|\nabla \log q_{t+\Delta t}(x)\|^2. \quad (120)$$

Using it in Eq. (114), we have

$$\log \frac{r_{\Delta t}(y|x)}{\tilde{r}_{\Delta t}(y|x)} = \log q_t(y) - \log q_t(x) - \Delta t \left(\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}} \right) \Delta \log q_{t+\Delta t}(x) - \quad (121)$$

$$- (1 + \Delta t \frac{\partial \log \alpha_{t+\Delta t}}{\partial t}) \langle y - x, \nabla \log q_{t+\Delta t}(x) \rangle + o(\Delta t). \quad (122)$$

Finally, we consider $y = x + \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon$. Then we can estimate

$$\log q_t(y) = \log q_t(x) + \langle \nabla \log q_t(x), y - x \rangle + \frac{1}{2} (y - x)^T \frac{\partial^2 \log q_t(x)}{\partial x^2} (y - x) \quad (123)$$

$$+ o(\|y - x\|^2) \quad (124)$$

$$= \log q_t(x) + \langle \nabla \log q_t(x), \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon \rangle \quad (125)$$

$$+ \Delta t \frac{g_{t+\Delta t}^2}{2} \varepsilon^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \varepsilon + o(\Delta t), \quad (126)$$

where $\frac{\partial^2 \log q_t(x)}{\partial x^2}$ denotes the Hessian of the log-density. Thus, for $y = x + \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon$, we have

$$\log \frac{r_{\Delta t}(y|x)}{\tilde{r}_{\Delta t}(y|x)} = \Delta t \frac{g_{t+\Delta t}^2}{2} \varepsilon^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \varepsilon - \Delta t \left(\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}} \right) \Delta \log q_{t+\Delta t}(x) + \quad (127)$$

$$+ \langle \nabla \log q_t(x) - \nabla \log q_{t+\Delta t}(x), \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon \rangle + o(\Delta t). \quad (128)$$

The last term is $o(\Delta t)$ by expanding the scores in time, hence, for $\frac{g_{t+\Delta t}^2}{2} = \sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}}$, we have

$$\log \frac{r_{\Delta t}(x + \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon | x)}{\tilde{r}_{\Delta t}(x + \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon | x)} = \quad (129)$$

$$= \Delta t \left(\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}} \right) \left(\Delta \log q_{t+\Delta t}(x) - \varepsilon^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \varepsilon \right) + o(\Delta t). \quad (130)$$

□

D.3 PROOF OF THM. 2

Theorem 2. [Denoising kernel] *For the universal noising kernel $k_{\Delta t}(x|y)$, the Gaussian approximation of the corresponding reverse kernel $r_{\Delta t}(y|x) = k_{\Delta t}(x|y)q_t(y)/q_{t+\Delta t}(x)$ is*

$$\tilde{r}_{\Delta t}(y|x) = \mathcal{N}\left(y \mid \frac{\alpha_t}{\alpha_{t+\Delta t}}x + \frac{\alpha_t}{\alpha_{t+\Delta t}}S_{t+\Delta t}^2 \nabla \log q_{t+\Delta t}(x), \left(\frac{\alpha_t}{\alpha_{t+\Delta t}}S_{t+\Delta t}\right)^2\right), \quad (81)$$

which corresponds to a single step of the Euler discretization of the following SDE

$$dx_\tau = -\left(\frac{\partial}{\partial t} \log \alpha_t x - 2\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t} \nabla_x \log q_t(x)\right) \cdot d\tau + \sqrt{2\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t}} dW_\tau, \quad (82)$$

and

$$D_{\text{KL}}(\tilde{r}_{\Delta t}(y|x) \| r_{\Delta t}(y|x)) = o(\Delta t). \quad (83)$$

1404 *Proof.* To find the reverse the kernel, we first use Tweedie’s formula to find the expectation, i.e.

$$1405 \quad \nabla \log q_{t+\Delta t}(x) = -\frac{1}{q_{t+\Delta t}(x)} \int dy \frac{1}{S_{t+\Delta t}^2} \left(x - \frac{\alpha_{t+\Delta t}}{\alpha_t} y \right) k_{\Delta t}(x|y) q_t(y) \quad (131)$$

$$1408 \quad \frac{\alpha_t}{\alpha_{t+\Delta t}} S_{t+\Delta t}^2 \nabla \log q_{t+\Delta t}(x) = -\frac{1}{q_{t+\Delta t}(x)} \int dy \left(\frac{\alpha_t}{\alpha_{t+\Delta t}} x - y \right) k_{\Delta t}(x|y) q_t(y) \quad (132)$$

$$1410 \quad \frac{\alpha_t}{\alpha_{t+\Delta t}} x + \frac{\alpha_t}{\alpha_{t+\Delta t}} S_{t+\Delta t}^2 \nabla \log q_{t+\Delta t}(x) = \int dy y \frac{k_{\Delta t}(x|y) q_t(y)}{q_{t+\Delta t}(x)} = \int dy y r_{\Delta t}(y|x). \quad (133)$$

1412 Thus, we are going to approximate

$$1414 \quad \tilde{r}_{\Delta t}(y|x) = \mathcal{N} \left(y \mid \frac{\alpha_t}{\alpha_{t+\Delta t}} x + \frac{\alpha_t}{\alpha_{t+\Delta t}} S_{t+\Delta t}^2 \nabla \log q_{t+\Delta t}(x), \left(\frac{\alpha_t}{\alpha_{t+\Delta t}} S_{t+\Delta t} \right)^2 \right) \quad (134)$$

1417 where $\left(\frac{\alpha_t}{\alpha_{t+\Delta t}} S_{t+\Delta t} \right)^2$ is chosen to match the leading term of $k_{\Delta t}(x|y)$. The constants can be
1418 estimated as follows

$$1419 \quad S_{t+\Delta t}^2 = \sigma_{t+\Delta t}^2 - \sigma_t^2 \frac{\alpha_{t+\Delta t}^2}{\alpha_t^2} = \sigma_{t+\Delta t}^2 - (\sigma_{t+\Delta t}^2 - 2dt\sigma_{t+\Delta t} \frac{\partial \sigma_{t+\Delta t}}{\partial t} + o(\Delta t)) \frac{\alpha_{t+\Delta t}^2}{\alpha_t^2} \quad (135)$$

$$1422 \quad = \sigma_{t+\Delta t}^2 - (\sigma_{t+\Delta t}^2 - 2dt\sigma_{t+\Delta t} \frac{\partial \sigma_{t+\Delta t}}{\partial t} + o(\Delta t)) \left(1 + 2dt \frac{\alpha_{t+\Delta t}^2}{\alpha_{t+\Delta t}^3} \frac{\partial \alpha_{t+\Delta t}}{\partial t} + o(\Delta t) \right) \quad (136)$$

$$1425 \quad = 2dt\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}} + o(\Delta t), \quad (137)$$

1427 and

$$1428 \quad \frac{\alpha_t}{\alpha_{t+\Delta t}} = 1 + \frac{1}{\alpha_{t+\Delta t}} \frac{\partial \alpha_{t+\Delta t}}{\partial t} (-\Delta t) + o(\Delta t) = 1 - \Delta t \frac{\partial \log \alpha_{t+\Delta t}}{\partial t} + o(\Delta t). \quad (138)$$

1430 Thus, y can be generated as

$$1431 \quad y = x - \Delta t \frac{\partial \log \alpha_{t+\Delta t}}{\partial t} x + 2dt\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}} \nabla \log q_{t+\Delta t}(x) + \quad (139)$$

$$1434 \quad + \sqrt{2dt\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}}} \varepsilon, \quad (140)$$

1436 where ε is a standard normal random variable. This corresponds to the single step of the Euler
1437 discretization of the following SDE

$$1439 \quad dx_t = - \left(\frac{\partial}{\partial t} \log \alpha_t x - 2\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t} \nabla_x \log q_t(x) \right) \cdot \Delta t + \sqrt{2\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t}} dW_t. \quad (141)$$

1441 For the KL-divergence between the reverse kernel and its Gaussian approximation, we have

$$1442 \quad D_{\text{KL}}(\tilde{r}_{\Delta t}(y|x) \| r_{\Delta t}(y|x)) = \mathbb{E}_\varepsilon \log \frac{\tilde{r}_{\Delta t}(x + \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon | x)}{r_{\Delta t}(x + \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon | x)} + o(\Delta t), \quad (142)$$

1444 where

$$1446 \quad v = - \left(\frac{\partial}{\partial t} \log \alpha_t x - 2\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t} \nabla_x \log q_t(x) \right), \quad \text{and} \quad g_{t+\Delta t} = \sqrt{2\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t}}. \quad (143)$$

1448 From [Lemma 1](#), we have

$$1449 \quad D_{\text{KL}}(\tilde{r}_{\Delta t}(y|x) \| r_{\Delta t}(y|x)) = \quad (144)$$

$$1450 \quad = \Delta t \left(\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}} \right) \mathbb{E}_\varepsilon \left(\varepsilon^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \varepsilon - \Delta \log q_{t+\Delta t}(x) \right) + o(\Delta t) \quad (145)$$

$$1452 \quad = o(\Delta t). \quad (146)$$

1454 \square

1455 D.4 PROOF OF [THM. 3](#)

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511

Theorem 3. [Detailed balance density estimator error] For $y = x + \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon$, where $g_{t+\Delta t} = \sqrt{2\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t}}$ and ε is a standard normal random variable, the following estimator is unbiased w.r.t. the generated samples y , i.e.

$$\mathbb{E}_\varepsilon \left[\log \frac{r_{\Delta t}(y|x)}{\tilde{r}_{\Delta t}(y|x)} \right] = o(\Delta t), \quad (85)$$

and the variance is

$$\mathbf{D}_\varepsilon \left[\log \frac{r_{\Delta t}(y|x)}{\tilde{r}_{\Delta t}(y|x)} \right] = \Delta t^2 \frac{g_{t+\Delta t}^4}{4} \mathbf{D}_\varepsilon \left(\varepsilon^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \varepsilon \right) + o(\Delta t^2), \quad (86)$$

where $\frac{\partial^2 \log q_{t+\Delta t}(x)}{\partial x^2}$ is the Hessian matrix of the log-density.

Proof. From Lemma 1, we have

$$\mathbb{E}_\varepsilon \log \frac{r_{\Delta t}(y|x)}{\tilde{r}_{\Delta t}(y|x)} = \mathbb{E}_\varepsilon \Delta t \frac{g_{t+\Delta t}^2}{2} \left(\Delta \log q_{t+\Delta t}(x) - \varepsilon^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \varepsilon \right) + o(\Delta t) \quad (147)$$

$$= \Delta t \frac{g_{t+\Delta t}^2}{2} \left(\Delta \log q_{t+\Delta t}(x) - \mathbb{E}_\varepsilon \text{Tr} \left[\varepsilon^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \varepsilon \right] \right) + o(\Delta t) \quad (148)$$

$$= \Delta t \frac{g_{t+\Delta t}^2}{2} \left(\Delta \log q_{t+\Delta t}(x) - \mathbb{E}_\varepsilon \text{Tr} \left[\frac{\partial^2 \log q_t(x)}{\partial x^2} \right] \right) + o(\Delta t) \quad (149)$$

$$= o(\Delta t). \quad (150)$$

For the variance, we have

$$\mathbf{D}_\varepsilon \log \frac{r_{\Delta t}(y|x)}{\tilde{r}_{\Delta t}(y|x)} = \Delta t^2 \frac{g_{t+\Delta t}^4}{4} \mathbb{E}_\varepsilon \left(\Delta \log q_{t+\Delta t}(x) - \varepsilon^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \varepsilon \right)^2 + o(\Delta t^2) \quad (151)$$

$$= \Delta t^2 \frac{g_{t+\Delta t}^4}{4} \mathbb{E}_\varepsilon \left(\mathbb{E}_\eta \left[\eta^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \eta \right] - \varepsilon^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \varepsilon \right)^2 + o(\Delta t^2) \quad (152)$$

$$= \Delta t^2 \frac{g_{t+\Delta t}^4}{4} \mathbf{D}_\varepsilon \left(\varepsilon^T \frac{\partial^2 \log q_t(x)}{\partial x^2} \varepsilon \right) + o(\Delta t^2), \quad (153)$$

where we used a standard normal random variable η . \square

D.5 PROOF OF PROP. 11

Proposition 11. [Recurrence relation for the log-density] For $y = x + \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon$, where $g_{t+\Delta t} = \sqrt{2\sigma_t^2 \frac{\partial}{\partial t} \log \frac{\sigma_t}{\alpha_t}}$ and ε is a standard normal random variable, the estimator from Eq. (87) can be expanded as follows

$$\begin{aligned} -\log k_{\Delta t}(x|y) + \log \tilde{r}_{\Delta t}(y|x) &= \Delta t \frac{\partial}{\partial t} \log \alpha_{t+\Delta t} - \Delta t \frac{g_{t+\Delta t}^2}{2} \|\nabla \log q_{t+\Delta t}(x)\|^2 + \\ &+ \left\langle \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon + \Delta t \frac{\partial}{\partial t} \log \alpha_{t+\Delta t} x, \nabla \log q_{t+\Delta t}(x) \right\rangle + o(\Delta t) \end{aligned} \quad (88)$$

1512 *Proof.* Indeed,

1513
1514
$$\log \frac{\tilde{r}_{\Delta t}(y|x)}{k_{\Delta t}(x|y)} = \frac{1}{2B_{t+\Delta t}^2} \left(x - \frac{\alpha_{t+\Delta t}}{\alpha_t} y \right)^2 + \frac{d}{2} \log 2\pi S_{t+\Delta t}^2 - \quad (154)$$

1515
1516
1517
$$- \frac{1}{2(S_{t+\Delta t} \frac{\alpha_t}{\alpha_{t+\Delta t}})^2} \left(y - \frac{\alpha_t}{\alpha_{t+\Delta t}} x - \frac{\alpha_t}{\alpha_{t+\Delta t}} S_{t+\Delta t}^2 \nabla \log q_{t+\Delta t}(x) \right)^2 - \quad (155)$$

1518
1519
$$- \frac{d}{2} \log 2\pi S_{t+\Delta t}^2 - d \log \frac{\alpha_t}{\alpha_{t+\Delta t}} \quad (156)$$

1520
1521
$$= -d \log \frac{\alpha_t}{\alpha_{t+\Delta t}} + \left\langle \frac{\alpha_{t+\Delta t}}{\alpha_t} y - x, \nabla \log q_{t+\Delta t}(x) \right\rangle - \quad (157)$$

1522
1523
1524
$$- \frac{S_{t+\Delta t}^2}{2} \|\nabla \log q_{t+\Delta t}(x)\|^2 \quad (158)$$

1525 Expanding the time around $t + \Delta t$, we have

1526
1527
$$\log \frac{\tilde{r}_{\Delta t}(y|x)}{k_{\Delta t}(x|y)} = d\Delta t \frac{\partial}{\partial t} \log \alpha_{t+\Delta t} + \left\langle (1 + \Delta t \frac{\partial}{\partial t} \log \alpha_{t+\Delta t}) y - x, \nabla \log q_{t+\Delta t}(x) \right\rangle - \quad (159)$$

1528
1529
$$- \Delta t \left(\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}} \right) \|\nabla \log q_{t+\Delta t}(x)\|^2 + o(\Delta t). \quad (160)$$

1530 For $y = x + \Delta t \cdot v + \sqrt{\Delta t} g_{t+\Delta t} \varepsilon$, $g_{t+\Delta t} = \sqrt{2\sigma_{t+\Delta t}^2 \frac{\partial}{\partial t} \log \frac{\sigma_{t+\Delta t}}{\alpha_{t+\Delta t}}}$, we have

1531
1532
1533
$$\log \frac{\tilde{r}_{\Delta t}(y|x)}{k_{\Delta t}(x|y)} = d\Delta t \frac{\partial}{\partial t} \log \alpha_{t+\Delta t} - \Delta t \frac{g_{t+\Delta t}^2}{2} \|\nabla \log q_{t+\Delta t}(x)\|^2 + \quad (161)$$

1534
1535
1536
$$+ \left\langle \Delta t \cdot v + g_{t+\Delta t} \sqrt{\Delta t} \varepsilon + \Delta t \frac{\partial}{\partial t} \log \alpha_{t+\Delta t} x, \nabla \log q_{t+\Delta t}(x) \right\rangle + o(\Delta t). \quad (162)$$

1537
1538 \square

1539 E BROADER IMPACTS

1540
1541 In this paper, we present theoretical results and demonstrate use cases in generation tasks such as
1542 image generation and unconditional protein generation. Because of the theoretical nature of our
1543 contributions, this work carries little societal impact. SUPERDIFF can be used to generate protein
1544 structure from a composition of existing protein diffusion models. Better protein generation methods
1545 can potentially lead to negative use in generating bio-hazardous molecules and proteins. We do not
1546 perceive this as a great risk at the current stage of these models.

1547 F ADDITIONAL RESULTS FOR CIFAR-10

1548
1549 Table A1: Image generation performance for CIFAR-10 with *conditional* models trained on two random
1550 partitions of the training data (labeled A and B). We compare SUPERDIFF (OR) with the respective models
1551 (model_A and model_B) with the model that is trained on the full dataset (model_{A∪B}) and random choice between
1552 two models (model_{A OR B}).

	ODE inference			SDE inference		
	FID (↓)	IS (↑)	FLD (↓)	FID (↓)	IS (↑)	FLD (↓)
model _A	4.75	8.98	6.95 ± 0.12	4.66	9.35	6.39 ± 0.13
model _B	4.78	8.97	6.86 ± 0.15	4.36	9.45	6.20 ± 0.14
model _{A∪B}	5.30	9.04	6.82 ± 0.09	2.83	9.44	6.26 ± 0.11
model _{A OR B}	4.75	8.94	6.86 ± 0.15	4.41	9.40	6.3 ± 0.18
SUPERDIFF (OR)	4.74	9.00	6.98 ± 0.10	4.46	9.40	6.22 ± 0.15
SUPERDIFF _{T=100} (OR)	4.63	8.98	6.81 ± 0.19	4.23	9.42	6.27 ± 0.12

1560 G PROTEIN GENERATION

1561 G.1 EXPERIMENTAL DETAILS

1562
1563 In our setting, we consider two state-of-the-art diffusion models for protein generation: Proteus (Wang
1564 et al., 2024) and FrameDiff (Yim et al., 2023b), which were trained on protein structures from Protein
1565

Data Bank (PDB, (Berman et al., 2000)) to estimate the special Euclidean group (SE(3)) equivariant score over multiple diffusion steps. The models’ outputs predict the coordinates of a monomeric protein backbone.

We use pre-trained checkpoints from Proteus² and FrameDiff³. During protein generation at inference, we separately combine scores for translations and rotations from both models using Algorithm 1.

We investigate the use of different temperature (T) settings to scale κ for controlling the densities. We also found that adding a small bias (ω) towards Proteus log densities improved designability.

G.2 METRICS FOR EVALUATING GENERATED PROTEINS

Designability. We assess designability using the *self-consistency* evaluation from Trippe et al. (2023), where given a generated backbone, we predict its scaffold using a sequence prediction model (we use ProteinMPNN (Dauparas et al., 2022)) and re-fold the sequence using a structure prediction model (we use ESMFold (Lin et al., 2022)). We then compare the re-folded protein to the original generated backbone by computing their template modeling score (s_{cTM}) and root-mean-square-distance (s_{cRMSD}). A protein is considered to be designable if its s_{cRMSD} is $< 2\text{\AA}$ to the refolded structure. For each protein, we repeat this process 8 times and keep the sequence with the lowest s_{cRMSD} . We report the fraction of designable proteins for each method in Table 3 as well as the s_{cRMSD} mean of the resulting *designable* proteins.

Novelty. A significant impetus for generative modelling in biology and chemistry is to propose compounds that have not been previously identified (i.e., different from the training data), but are also possible to make. We compute three novelty metrics: the fraction of designable proteins with a s_{cTM} score < 0.5 (higher is better), the fraction of designable proteins with a s_{cTM} score < 0.3 , and the mean s_{cTM} score between generated proteins and the proteins from PDB that the original diffusion models were trained on, which represents the collective human knowledge of protein structures; a lower score indicates higher distance from the training set and is desirable since it shows generalization ability.

Diversity. To measure how diverse the generated proteins are, we compute their mean pairwise s_{cTM} score (lower is better), and also report the fraction of unique clusters formed after clustering them with MaxCluster (Herbert & Sternberg, 2008) (higher is better). Finally, we report the fraction of proteins that contain β -sheet secondary structures, as it has been found that these structures are typically more rare to generate (Bose et al., 2024).

G.3 PROTEIN DIVERSITY EXPLORATION SETTINGS

To embed and cluster the generated protein backbones we use the foldseek (Van Kempen et al., 2024) package to compute pairwise aligned TM-scores. We then use the UMAP (McInnes et al., 2018) package to compute a 2D embedding. Where proteins are represented as points that are close to each other if they are structurally similar (by aligned TM-Score).

We then clustered the proteins again with the foldseek tool to find representative structures. Finally we used KMeans to explore the space and narrow down which protein structures belong where on the protein structure manifold.

H GENERATING IMAGE COMPOSITIONS WITH STABLE DIFFUSION

For our concept interpolation experiments, we use publicly-available pre-trained weights, models, and schedulers from Stable Diffusion v1-4 <https://huggingface.co/CompVis/stable-diffusion-v1-4>.

In Figs. A1–A20, we show examples of image compositions generated by SUPERDIFF (AND), averaging of outputs, and joint prompting. Prompts are shown in each image caption. We show images generated by the first 6 seeds (uniform sampling), as well as our favourite images generated from 20 seeds. For joint prompting, we generate prompts from two concepts using the linking term "that looks like". For example, given the concepts "a lemon" and "a sunflower", the resulting prompt would be "a sunflower that looks like a lemon". We also generate

²<https://github.com/Wangchentong/Proteus>

³https://github.com/jasonkyuyim/se3_diffusion

images with the reversed prompt ("a lemon that looks like a sunflower"), and keep the images generated by the prompt resulting in the higher mean score for each metric. The order of the concepts in each image caption reflects the ordering that obtained the higher TIFA score.

In Fig. A21, we show examples of image compositions generated by SUPERDIFF (OR), joint prompting, and an upper bound of randomly selecting a prompt of one concept. Again, we keep the prompt order that resulted in higher scores for the baseline.



Figure A1: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "an airplane" and "an eagle".

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727



Figure A2: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "bottle cap" and "chess pawn".

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

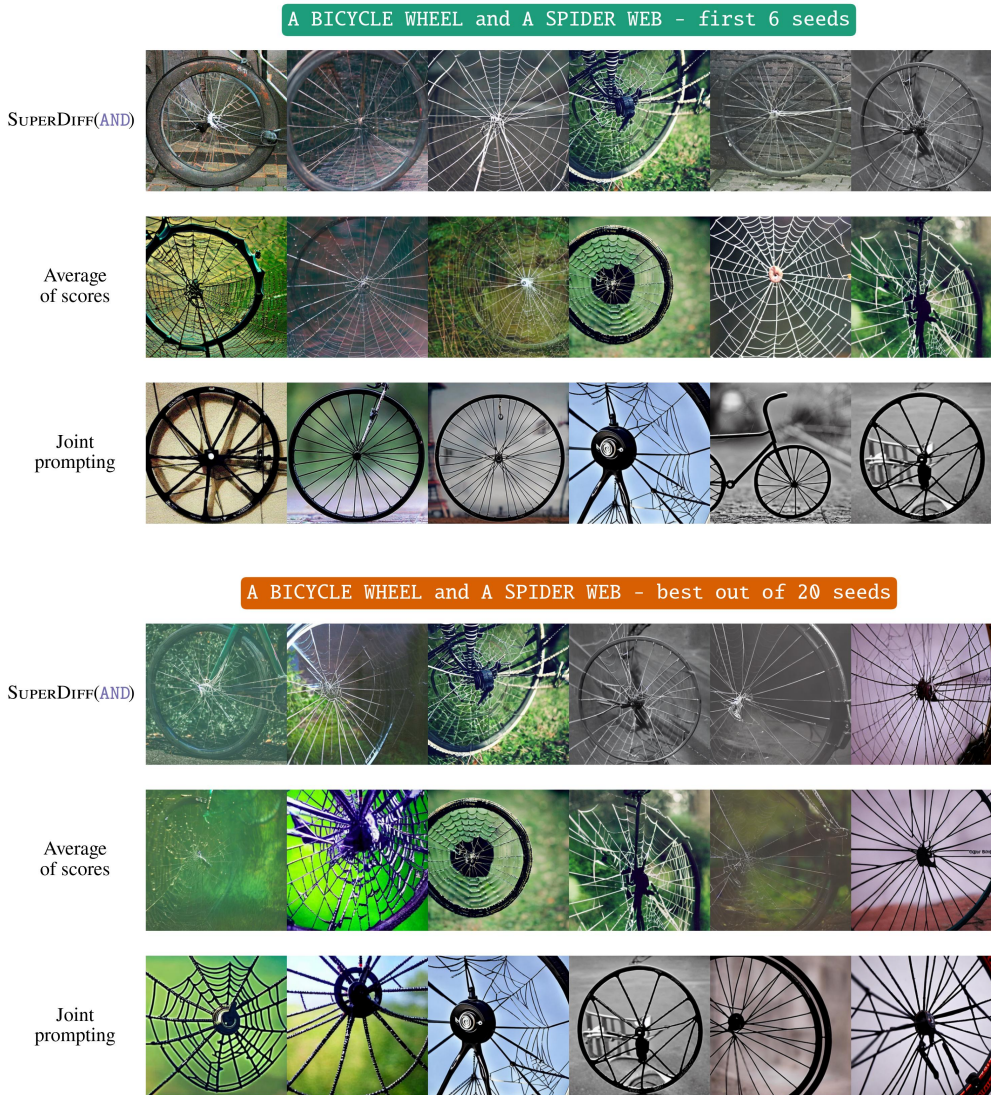


Figure A3: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "a bicycle wheel" and "a spider web".

1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835



Figure A4: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "a chair" and "an avocado".

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889



Figure A5: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "a cinnamon roll" and "a snail".

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943



Figure A6: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "cookie" and "moon".

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997



Figure A7: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "a dog" and "a cat".

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

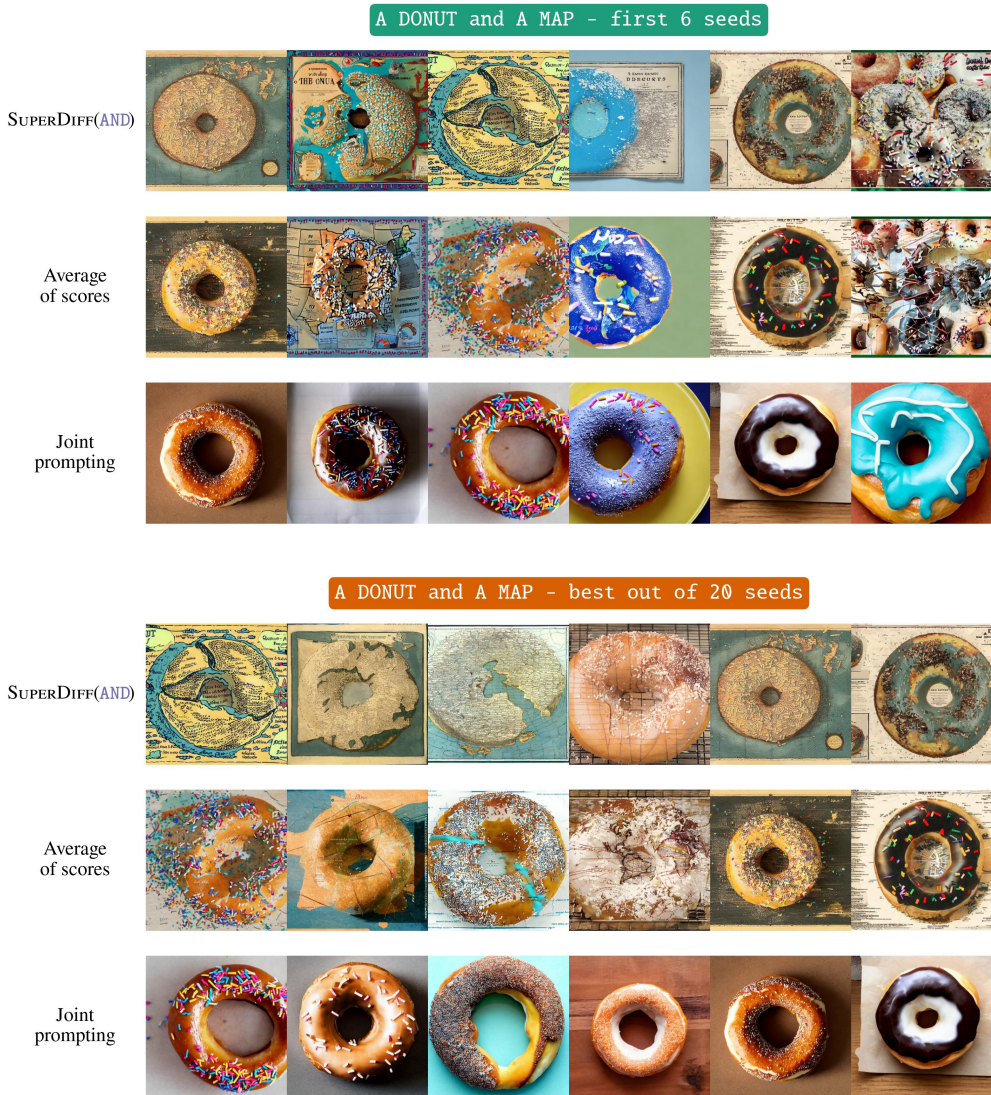


Figure A8: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "a donut" and "a map".

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105



Figure A9: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "duck" and "otter".

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

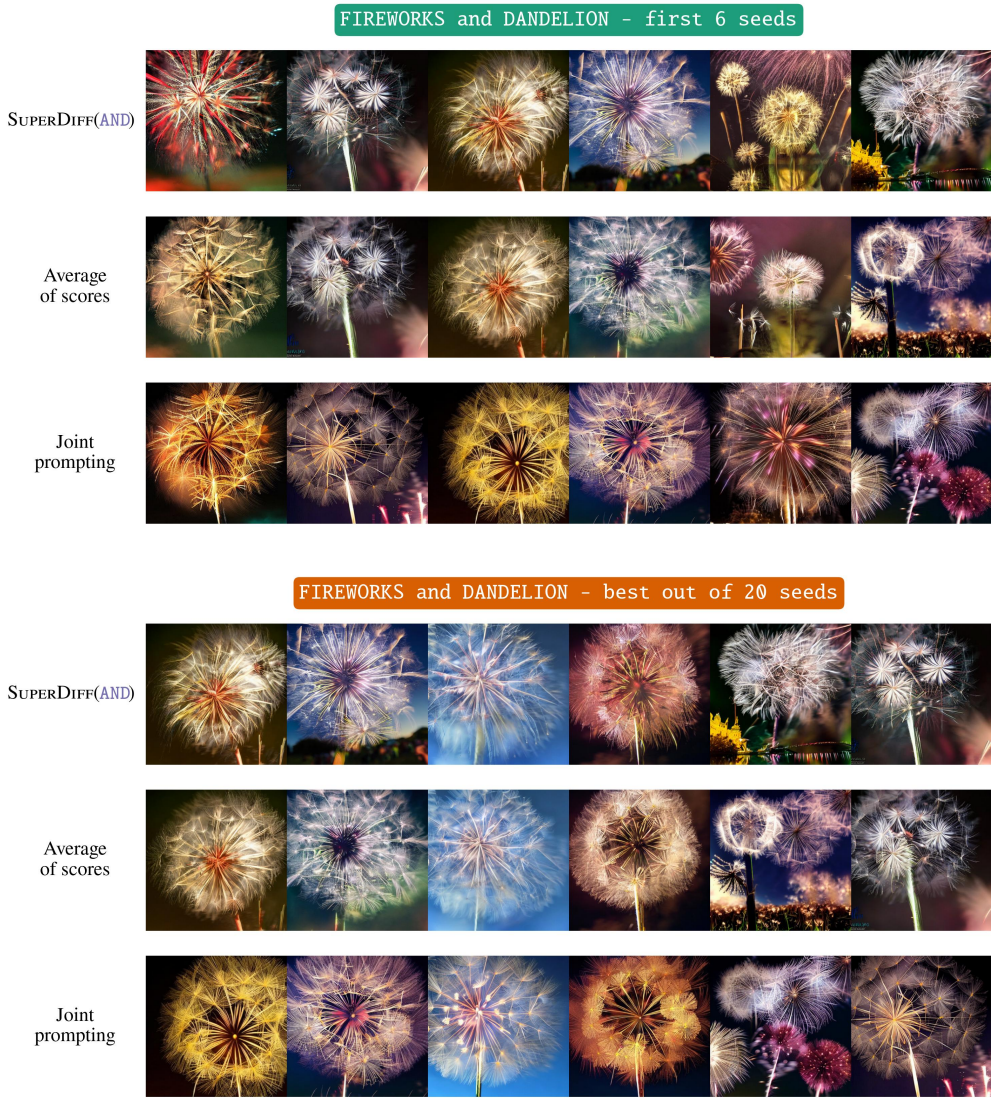


Figure A10: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "fireworks" and "dandelion".

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213



Figure A11: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "a flamingo" and "a candy cane".

2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267



Figure A12: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "a helicopter" and "a dragonfly".

2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321

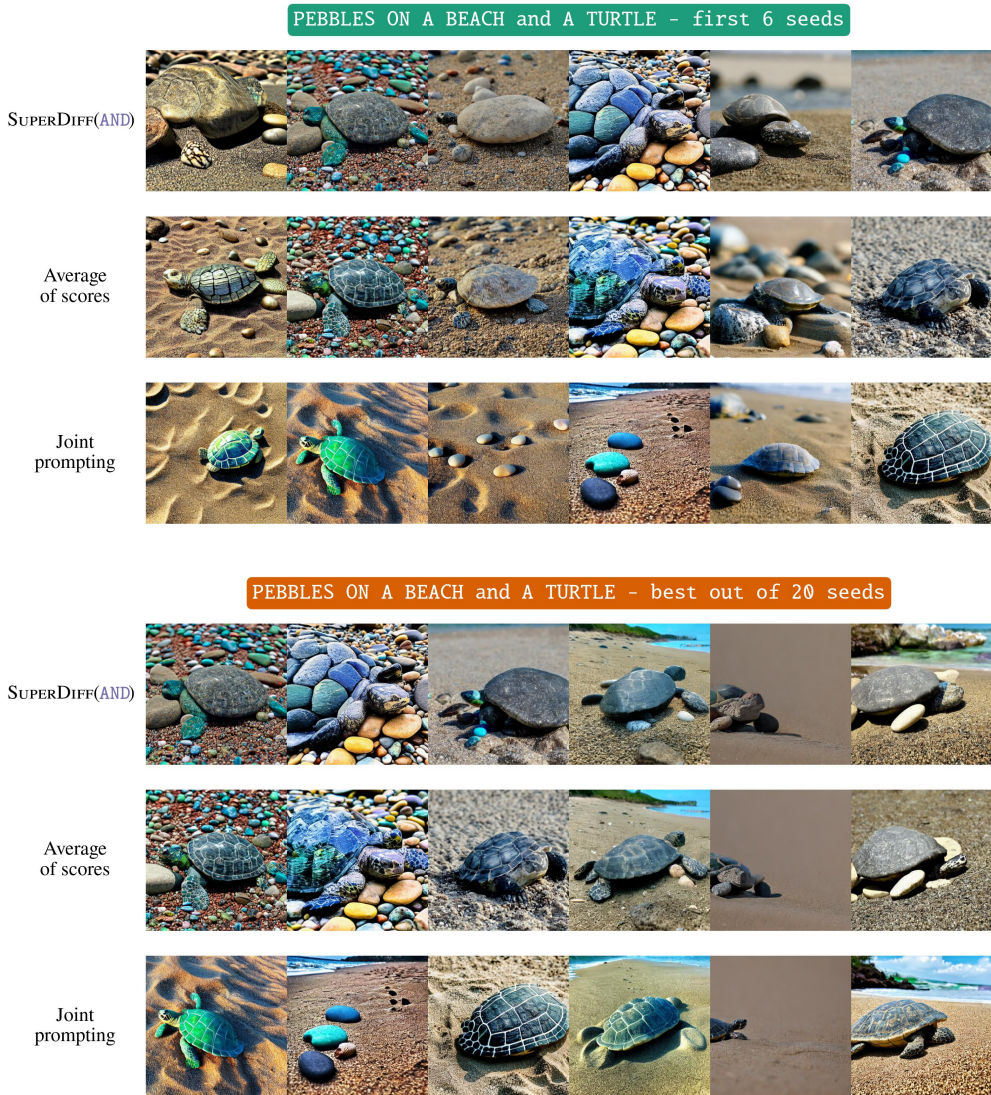


Figure A13: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "pebbles on a beach" and "a turtle".

2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375

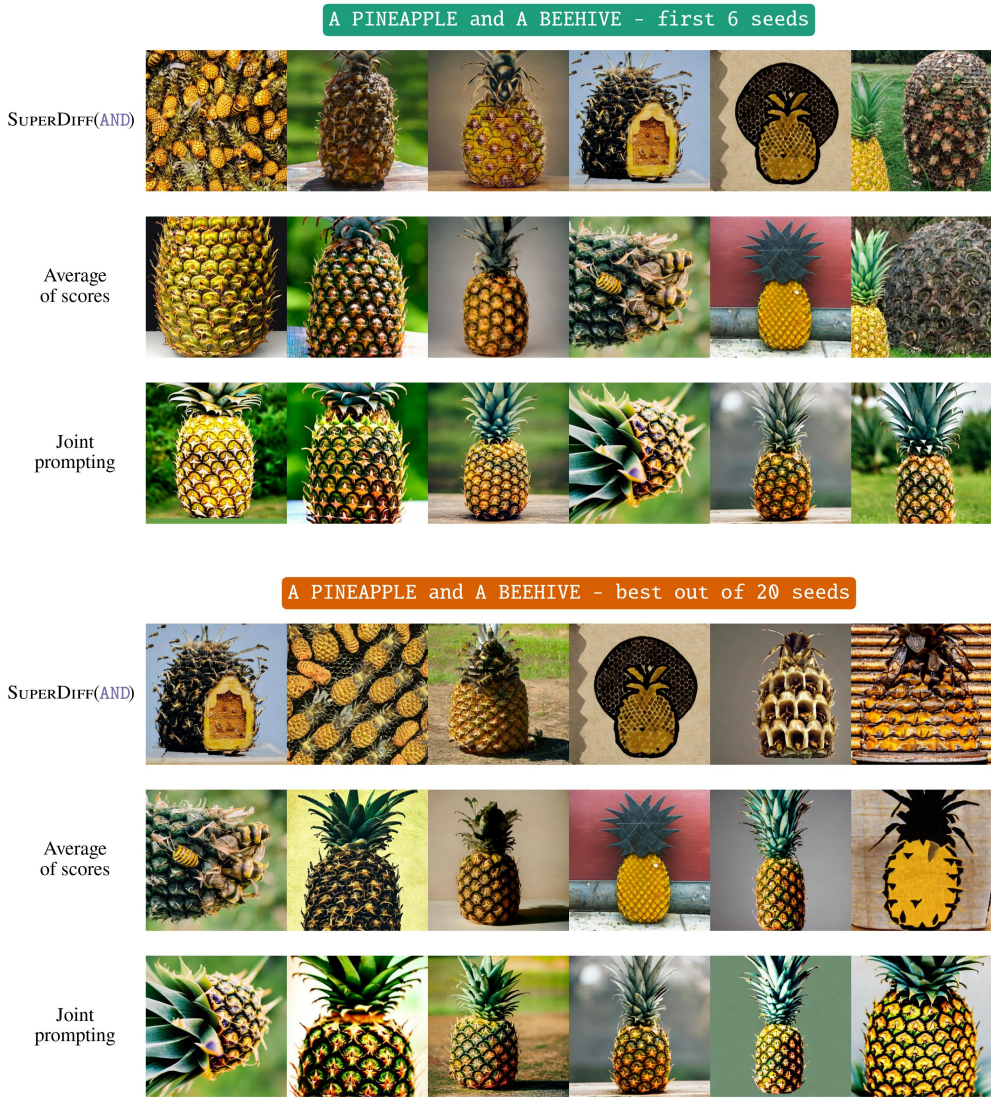


Figure A14: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "a pineapple" and "a beehive".

2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429

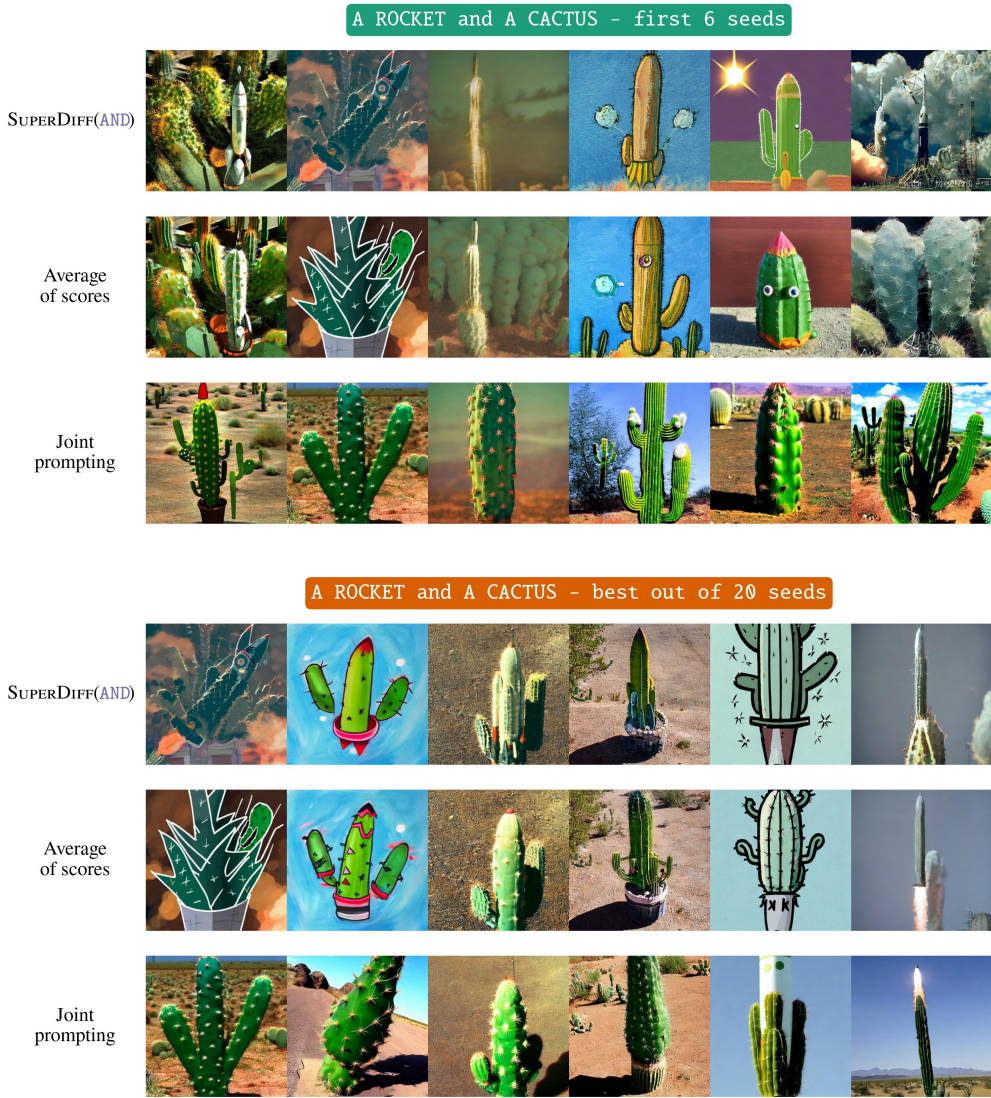


Figure A15: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "a rocket" and "a cactus".

2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483



Figure A16: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "a silhouette of a dog" and "a mountain landscape".

2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537

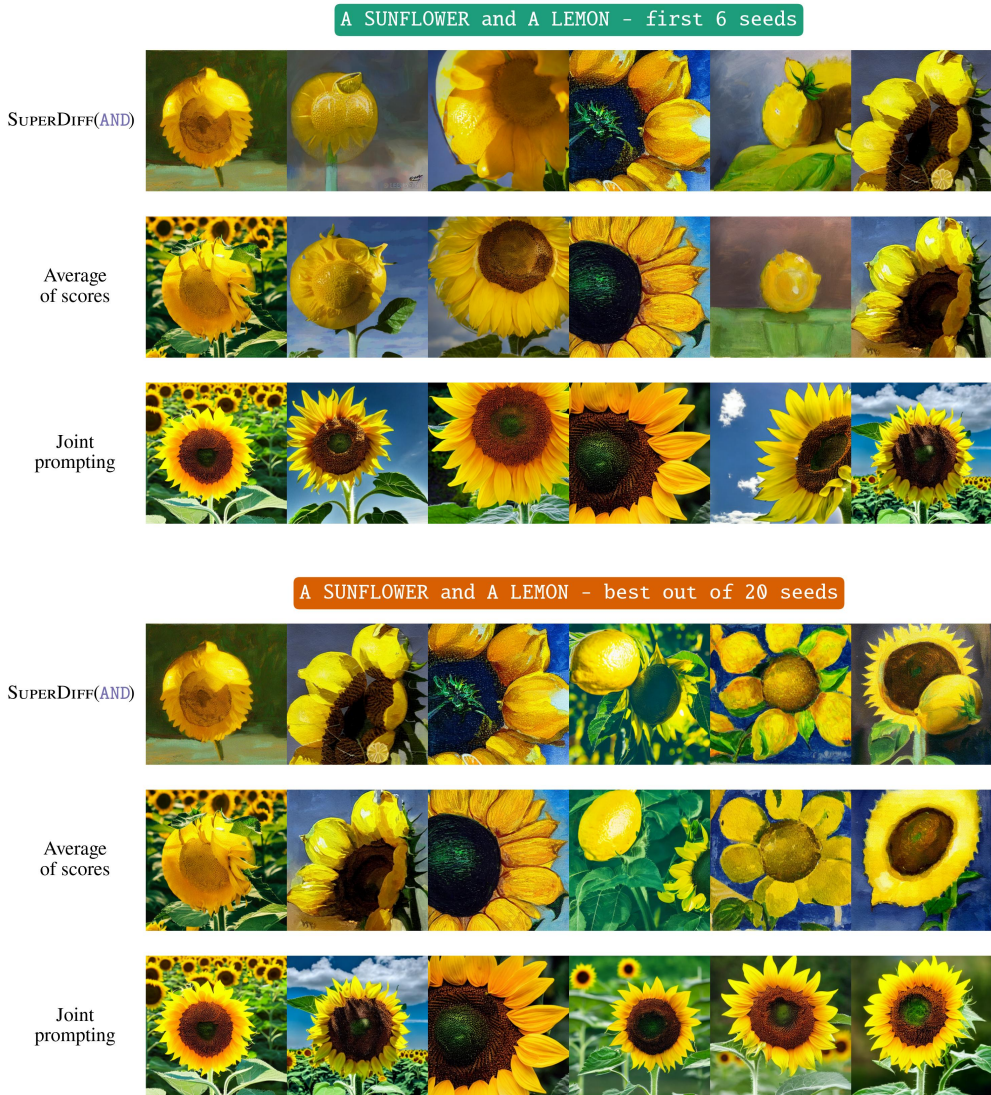


Figure A17: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "a sunflower" and "a lemon".

2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591



Figure A18: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "teddy bear" and "panda".

2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645

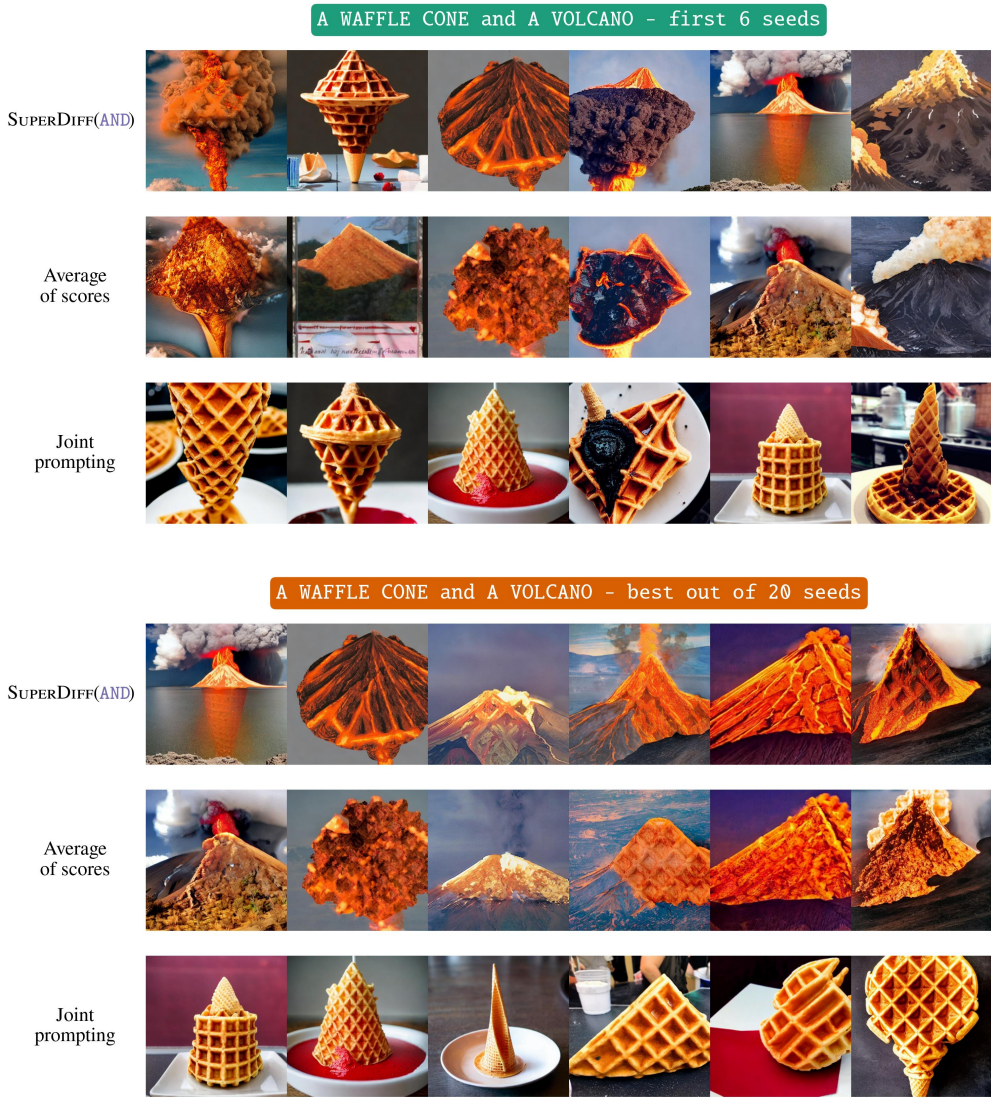


Figure A19: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "a waffle cone" and "a volcano".

2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699

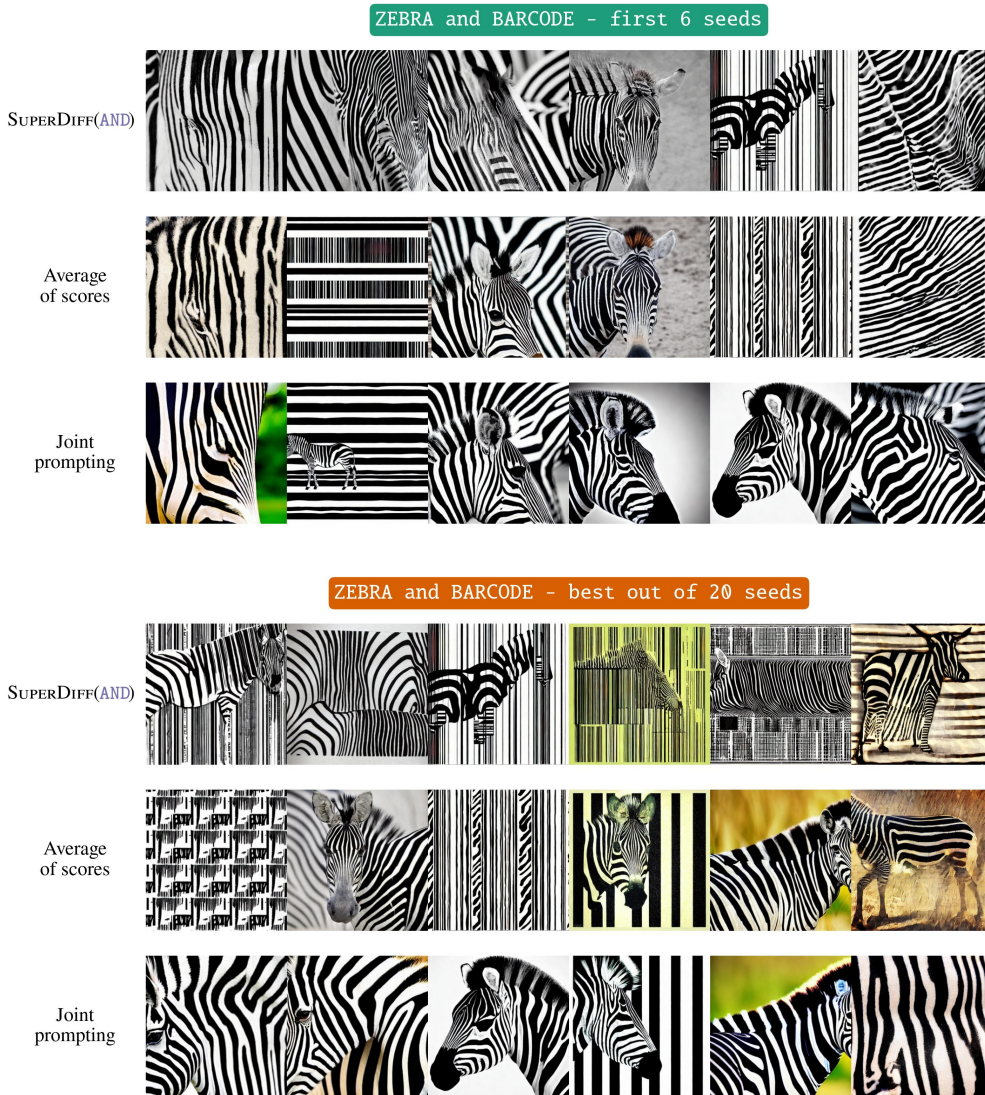


Figure A20: Image compositions generated using SUPERDIFF (AND), averaging of outputs, or joint prompting for the concepts "zebra" and "barcode".

2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753

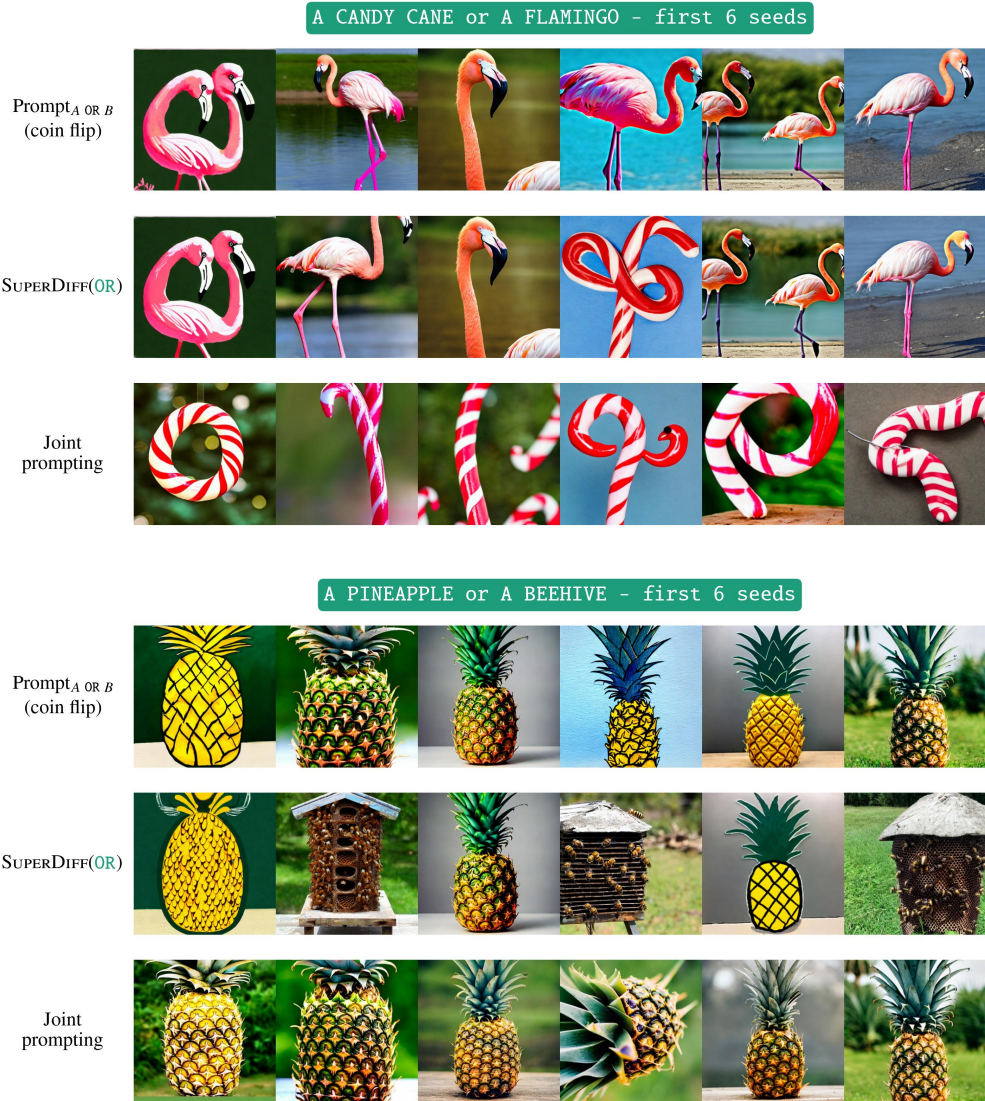


Figure A21: Concept selections using SUPERDIFF (OR), joint prompting, and randomly selecting a prompt of one concept. The top subfigure shows generated images for the concepts "a candy cane" or "a flamingo". The bottom subfigure shows generated images for the concepts "a pineapple" or "a beehive".