



DEPARTMENT OF COMPUTER SCIENCE

Video Diffusion Models for Climate Simulations

George Herbert

A dissertation submitted to the University of Bristol in accordance with the requirements of the degree
of Master of Science in the Faculty of Engineering.

Tuesday 28th March, 2023

Abstract

Dedication and Acknowledgements

Declaration

I declare that the work in this dissertation was carried out in accordance with the requirements of the University's Regulations and Code of Practice for Taught Programmes and that it has not been submitted for any other academic award. Except where indicated by specific reference in the text, this work is my own work. Work done in collaboration with, or with the assistance of others, is indicated as such. I have identified all material in this dissertation which is not my own work through appropriate referencing and acknowledgement. Where I have quoted or otherwise incorporated material which is the work of others, I have included the source in the references. Any views expressed in the dissertation, other than referenced material, are those of the author.

George Herbert, Tuesday 28th March, 2023

Contents

| | | |
|----------|--|----------|
| 1 | Introduction | 1 |
| 2 | Background | 2 |
| 2.1 | Diffusion Models | 2 |
| 2.2 | Climate Simulations | 3 |
| 3 | Results | 4 |
| 4 | Conclusion | 5 |
| A | Diffusion Models | 7 |
| A.1 | Derivation of $q(\mathbf{z}_t \mathbf{z}_s)$ | 7 |
| A.2 | α -Cosine Noise Schedule | 8 |

List of Figures

List of Tables

Ethics Statement

Notation and Acronyms

Chapter 1

Introduction

Chapter 2

Background

2.1 Diffusion Models

Diffusion models are

Given observed datapoints \mathbf{x} , the goal of a generative model is to learn to model its true data distribution $p(\mathbf{x})$.

2.1.1 Forward Diffusion Process

The *forward diffusion process* is a Gaussian diffusion process that defines a sequence of increasingly noisy versions of \mathbf{x} , which we call the *latent variables*:

$$\mathbf{z} = \{\mathbf{z}_t \mid t \in [0, 1]\} \quad (2.1)$$

The forward process forms a conditional joint distribution $q(\mathbf{z}|\mathbf{x})$, whose marginal distributions of latent variables \mathbf{z}_t given $\mathbf{x} \sim p(\mathbf{x})$ are given by:

$$q(\mathbf{z}_t|\mathbf{x}) = \mathcal{N}(\mathbf{z}_t; \alpha_t \mathbf{x}, \sigma_t^2 \mathbf{I}) \quad (2.2)$$

where α_t and σ_t are strictly positive scalar-valued functions of t . The joint distribution of latent variables $\mathbf{z}_r, \mathbf{z}_s, \mathbf{z}_t$ at subsequent timesteps $0 \leq r < s < t \leq 1$ is Markovian:

$$q(\mathbf{z}_t|\mathbf{z}_s, \mathbf{z}_r) = q(\mathbf{z}_t|\mathbf{z}_s) = \mathcal{N}(\mathbf{z}_t; \alpha_{t|s} \mathbf{z}_s, \sigma_{t|s}^2 \mathbf{I}) \quad (2.3)$$

where $\alpha_{t|s} = \alpha_t \alpha_s^{-1}$ and $\sigma_{t|s}^2 = \sigma_t^2 - \alpha_{t|s}^2 \sigma_s^2$. A full derivation of $q(\mathbf{z}_t|\mathbf{z}_s)$ is given in Appendix A.1.

2.1.2 Noise Schedule

We formalise the notion that \mathbf{z}_t is increasingly noisy by defining the log signal-to-noise ratio

$$\lambda_t = \log \left(\frac{\alpha_t^2}{\sigma_t^2} \right) \in [\lambda_{\min}, \lambda_{\max}] \quad (2.4)$$

as a strictly monotonically decreasing function f_λ of time $t \in [0, 1]$, known as the *noise schedule*.

In this work, we use a truncated continuous-time version of the α -cosine schedule [5], introduced in its original discrete-time form by Nichol and Dhariwal [5]. The α -cosine schedule was motivated by the fact that the ‘linear’ schedule introduced in prior work by Ho et al. [1] causes α_t to fall to zero more quickly than is optimal. Nichol and Dhariwal empirically found that this induces too much noise in the latter stages of the forward diffusion process; as such, the latent variables \mathbf{z}_t in these stages contribute little to sample quality. In response, they proposed the original discrete-time α -cosine schedule. In this work, we use a continuous-time diffusion model and therefore use an adapted model described in [3]. More formally, we define:

$$f_\lambda(t) = -2 \log \left(\tan \left(\frac{\pi}{2} (t_0 + t(t_1 - t_0)) \right) \right) \quad (2.5)$$

where t_0 and t_1 truncate the range of $f_\lambda(t)$ to $[\lambda_{\min}, \lambda_{\max}]$ for $t \in [0, 1]$, and are themselves defined as:

$$t_0 = \frac{2}{\pi} \arctan \left(\exp \left(-\frac{1}{2} \lambda_{\max} \right) \right) \quad (2.6)$$

$$t_1 = \frac{2}{\pi} \arctan \left(\exp \left(-\frac{1}{2} \lambda_{\min} \right) \right) \quad (2.7)$$

We compute α_t and σ_t from λ_t via the following equations:

$$\alpha_t = \sqrt{S(\lambda_t)} \quad (2.8)$$

$$\sigma_t = \sqrt{S(-\lambda_t)} \quad (2.9)$$

where S is the sigmoid function. A full derivation of f_λ is given in Appendix A.2.

2.1.3 Generative Model

The *generative model* is a learned hierarchical model that matches the forward process running in reverse-time: we sequentially generate latent variables \mathbf{z}_t starting from $t = 1$ and working backwards to $t = 0$.

In this work, our diffusion model is variance preserving (i.e. $\alpha_t^2 = 1 - \sigma_t^2$) and λ_{\min} is sufficiently small. As such, we can model the marginal distribution of \mathbf{z}_1 as the multivariate standard Gaussian:

$$p(\mathbf{z}_1) \approx \mathcal{N}(\mathbf{z}_1; \mathbf{0}, \mathbf{I}) \quad (2.10)$$

Once we have sampled $\mathbf{z}_1 \sim p_\theta(\mathbf{z}_1) = \mathcal{N}(\mathbf{z}_1; \mathbf{0}, \mathbf{I})$, we use the discrete-time ancestral sampler [1] to sequentially generate each latent variable \mathbf{z}_s from \mathbf{z}_t where $0 \leq s < t \leq 1$. The discrete-time ancestral sampler samples $\mathbf{z}_s \sim p_\theta(\mathbf{z}_s | \mathbf{z}_t)$ via:

$$p_\theta(\mathbf{z}_s | \mathbf{z}_t) = q(\mathbf{z}_s | \mathbf{z}_t, \mathbf{x} = \hat{\mathbf{x}}_\theta(\mathbf{z}_t, \lambda_t)) \quad (2.11)$$

$$= \mathcal{N} \left(\tilde{\boldsymbol{\mu}}_{s|t}(\mathbf{z}_t, \mathbf{x} = \hat{\mathbf{x}}_\theta(\mathbf{z}_t, \lambda_t)), \tilde{\sigma}_{s|t} \mathbf{I} \right) \quad (2.12)$$

where $\hat{\mathbf{x}}_\theta(\mathbf{z}_t, \lambda_t)$ is our denoised estimate of the original data \mathbf{x} given latent \mathbf{z}_t and log signal-to-noise ratio λ_t , and

$$\tilde{\boldsymbol{\mu}}_{s|t}(\mathbf{z}_t, \mathbf{x}) = \frac{\alpha_{t|s} \sigma_s^2}{\sigma_t^2} \mathbf{z}_t + \frac{\alpha_s \sigma_{t|s}^2}{\sigma_t^2} \mathbf{x} \quad (2.13)$$

$$\tilde{\sigma}_{s|t}^2 = \frac{\sigma_{t|s} \sigma_s}{\sigma_t} \quad (2.14)$$

For large enough λ_{\max} , \mathbf{z}_0 is almost noiseless, so learning a model $p(\mathbf{z}_0)$ is practically equivalent to learning a model $p(\mathbf{x})$.

2.1.4 Objective Function

Kingma and Gao [4] discovered that diffusion models in the broader literature are optimised with various objectives that are almost all special cases of a weighted loss per datapoint \mathbf{x} , which is defined as:

$$\mathcal{L}_w = w(\lambda_{\min}) \mathcal{L}(\lambda_{\min}) + \int_{\lambda_{\min}}^{\lambda_{\max}} w(\lambda) \mathcal{L}'(\lambda) d\lambda \quad (2.15)$$

where

$$\mathcal{L}(\lambda) = D_{KL}(q(\mathbf{z}_t, \dots, \mathbf{z}_1 | \mathbf{x}) \| p(\mathbf{z}_t, \dots, \mathbf{z}_1)) \quad (2.16)$$

$$\mathcal{L}'(\lambda) = \frac{d}{d\lambda} \mathcal{L} = \frac{1}{2} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\|\boldsymbol{\epsilon} - \hat{\boldsymbol{\epsilon}}_\theta(\mathbf{z}_t, \lambda_t)\|_2^2] \quad (2.17)$$

and $w(\lambda)$ is a weighting function.

In this work, we use the \mathbf{v} -parameterisation, wherein

$$\mathbf{v} = \quad (2.18)$$

2.1.5 Reconstruction-Guided Sampling

2.2 Climate Simulations

Chapter 3

Results

Chapter 4

Conclusion

Bibliography

- [1] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- [2] Jonathan Ho, Tim Salimans, Alexey A. Gritsenko, William Chan, Mohammad Norouzi, and David J. Fleet. Video diffusion models. *CoRR*, abs/2204.03458, 2022.
- [3] Emiel Hooeboom, Jonathan Heek, and Tim Salimans. simple diffusion: End-to-end diffusion for high resolution images. *CoRR*, abs/2301.11093, 2023.
- [4] Diederik P. Kingma and Ruiqi Gao. Understanding the diffusion objective as a weighted integral of elbos. *CoRR*, abs/2303.00848, 2023.
- [5] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 8162–8171. PMLR, 2021.

Appendix A

Diffusion Models

A.1 Derivation of $q(\mathbf{z}_t|\mathbf{z}_s)$

From Equation 2.2, we know $q(\mathbf{z}_t|\mathbf{x})$ is an isotropic Gaussian probability density function. As such, we can sample $\mathbf{z}_t \sim q(\mathbf{z}_t|\mathbf{x})$ by sampling $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ from the multivariate standard Gaussian distribution and computing:

$$\mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \boldsymbol{\epsilon}_t \quad (\text{A.1})$$

With some algebraic manipulation, we can show that:

$$\mathbf{z}_t = \alpha_t \mathbf{x} + \sqrt{\sigma_t^2} \boldsymbol{\epsilon}_t \quad (\text{A.2})$$

$$= \alpha_t \mathbf{x} + \sqrt{\sigma_t^2 - \frac{\alpha_t^2}{\alpha_s^2} \sigma_s^2 + \frac{\alpha_t^2}{\alpha_s^2} \sigma_s^2} \boldsymbol{\epsilon}_t \quad (\text{A.3})$$

$$= \alpha_t \mathbf{x} + \sqrt{\sigma_t^2 - \frac{\alpha_t^2}{\alpha_s^2} \sigma_s^2 + \left(\frac{\alpha_t}{\alpha_s} \sigma_s\right)^2} \boldsymbol{\epsilon}_t \quad (\text{A.4})$$

The sum of two independent Gaussian random variables with mean μ_1 and μ_2 and variance σ_1^2 and σ_2^2 is a Gaussian random variable with mean $\mu_1 + \mu_2$ and variance $\sigma_1^2 + \sigma_2^2$. As such, we can manipulate the above equation further to show that:

$$\mathbf{z}_t = \alpha_t \mathbf{x} + \sqrt{\sigma_t^2 - \frac{\alpha_t^2}{\alpha_s^2} \sigma_s^2} \boldsymbol{\epsilon}_t^* + \frac{\alpha_t}{\alpha_s} \sigma_s \boldsymbol{\epsilon}_s \quad (\text{A.5})$$

$$= \alpha_t \mathbf{x} + \frac{\alpha_t}{\alpha_s} \sigma_s \boldsymbol{\epsilon}_s + \sqrt{\sigma_t^2 - \frac{\alpha_t^2}{\alpha_s^2} \sigma_s^2} \boldsymbol{\epsilon}_t^* \quad (\text{A.6})$$

$$= \frac{\alpha_s}{\alpha_s} \alpha_t \mathbf{x} + \frac{\alpha_t}{\alpha_s} \sigma_s \boldsymbol{\epsilon}_s + \sqrt{\sigma_t^2 - \frac{\alpha_t^2}{\alpha_s^2} \sigma_s^2} \boldsymbol{\epsilon}_t^* \quad (\text{A.7})$$

$$= \frac{\alpha_t}{\alpha_s} (\alpha_s \mathbf{x} + \sigma_s \boldsymbol{\epsilon}_s) + \sqrt{\sigma_t^2 - \frac{\alpha_t^2}{\alpha_s^2} \sigma_s^2} \boldsymbol{\epsilon}_t^* \quad (\text{A.8})$$

$$(\text{A.9})$$

where $\boldsymbol{\epsilon}_t^*, \boldsymbol{\epsilon}_s \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ are similarly both sampled from the multivariate standard Gaussian distribution. We can substitute $\mathbf{z}_s = \alpha_s \mathbf{x} + \sigma_s \boldsymbol{\epsilon}_s$ into the above equation to show that:

$$\mathbf{z}_t = \frac{\alpha_t}{\alpha_s} \mathbf{z}_s + \sqrt{\sigma_t^2 - \frac{\alpha_t^2}{\alpha_s^2} \sigma_s^2} \boldsymbol{\epsilon}_t^* \quad (\text{A.10})$$

$$= \alpha_{t|s} \mathbf{z}_s + \sigma_{t|s} \boldsymbol{\epsilon}_t^* \quad (\text{A.11})$$

$$\sim \mathcal{N}(\mathbf{z}_t; \alpha_{t|s} \mathbf{z}_s, \sigma_{t|s}^2 \mathbf{I}) \quad (\text{A.12})$$

The subscript $t|s$ relates to the fact that $\alpha_{t|s}$ and $\sigma_{t|s}$ define the parameters of the Gaussian probability density function $q(\mathbf{z}_t|\mathbf{z}_s)$.

A.2 α -Cosine Noise Schedule

Before truncation, the continuous-time version of the α -cosine schedule [5] as described in [3] defines α_t^2 at a given timestep $t \in [0, 1]$ as:

$$\alpha_t^2 = \cos^2\left(\frac{\pi}{2}t\right) \quad (\text{A.13})$$

Since our model is a variance-preserving diffusion model, we can show that:

$$\sigma_t^2 = 1 - \alpha_t^2 \quad (\text{A.14})$$

$$= 1 - \cos^2\left(\frac{\pi}{2}t\right) \quad (\text{A.15})$$

$$= \sin^2\left(\frac{\pi}{2}t\right) \quad (\text{A.16})$$

As such, we define our noise schedule before truncation \tilde{f}_λ for all $t \in [0, 1]$ as:

$$\tilde{f}_\lambda(t) = \log\left(\frac{\alpha_t^2}{\sigma_t^2}\right) \quad (\text{A.17})$$

$$= \log\left(\frac{\cos^2\left(\frac{\pi}{2}t\right)}{\sin^2\left(\frac{\pi}{2}t\right)}\right) \quad (\text{A.18})$$

$$= -2\log\left(\tan\left(\frac{\pi}{2}t\right)\right) \quad (\text{A.19})$$

However, the above noise schedule means that $\tilde{f}_\lambda : [0, 1] \rightarrow [-\infty, \infty]$; in simpler terms, λ_t is unbounded. We follow prior work (e.g. [3, 2]) by truncating λ_t to the desired range $[\lambda_{\min}, \lambda_{\max}]$. To do so, we first need to define the inverse of the unbounded noise schedule:

$$\tilde{f}_\lambda^{-1}(\lambda) = \frac{2}{\pi} \arctan\left(\exp\left(-\frac{1}{2}\lambda\right)\right) \quad (\text{A.20})$$

From this, we define t_0 and t_1 as:

$$t_0 = \tilde{f}_\lambda^{-1}(0) = \frac{2}{\pi} \arctan\left(\exp\left(-\frac{1}{2}\lambda_{\max}\right)\right) \quad (\text{A.21})$$

$$t_1 = \tilde{f}_\lambda^{-1}(1) = \frac{2}{\pi} \arctan\left(\exp\left(-\frac{1}{2}\lambda_{\min}\right)\right) \quad (\text{A.22})$$

The truncated noise schedule used in this work is then defined as:

$$f_\lambda(t) = \tilde{f}_\lambda(t_0 + t(t_1 - t_0)) \quad (\text{A.23})$$

$$= -2\log\left(\tan\left(\frac{\pi}{2}(t_0 + t(t_1 - t_0))\right)\right) \quad (\text{A.24})$$