

# Role of Samplers in Diffusion Based Generative Models

Ernests Lavrinovics  
MSc Medialogy MED10  
Aalborg University, Copenhagen



# Agenda

- Introduction
  - Generative AI and Diffusion-Based Models
  - Applications
- Background
  - Diffusion-Based Speech Enhancement
- Problem Statement
- Methodology
- Results
- Summary

# Introduction

## Generative AI Models

- Discriminative AI – learn the decision boundary
  - Conv.Neural Nets (CNN), Multilayer Percept. (MLP), Transformers
- Generative AI – learn the underlying data distribution
  - Variational Autoencoders (VAE), Generative Adversarial Networks (GAN), Boltzmann machines, Transformers, Diffusion models
- Dependent on the objective and training setup



Image from <https://github.com/NVlabs/stylegan>

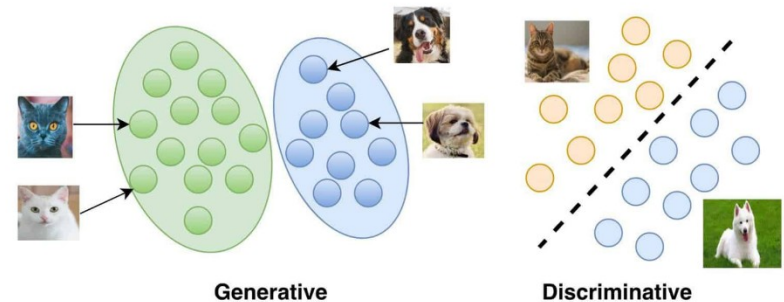


Image from <https://vitalflux.com/generative-vs-discriminative-models-examples/>

# Introduction

## Diffusion-Based Generative AI Models

- Diffusion (iterative process)
  - Forward-diffusion: corrupt a datapoint with noise
  - Reverse-diffusion: reverse the noise into data
- This process can be optionally be guided via conditioning

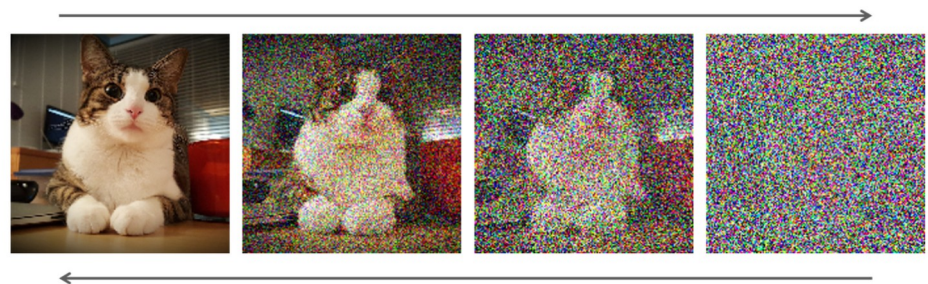
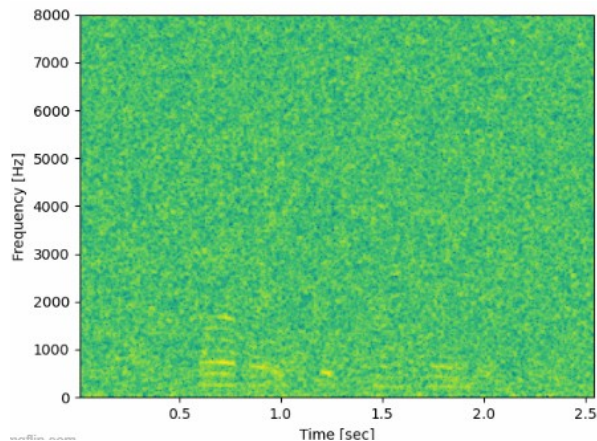


Image from  
<https://developer.nvidia.com/blog/improving-diffusion-models-as-an-alternative-to-gans-part-2/>

# Introduction

## Applications of Diffusion Models

- Image processing
  - Restoration (super resolution, inpainting, colorization)
  - Anomaly detection, semantic segmentation
- Audio processing
  - Speech and music enhancement/generation (communications, restoration)
- Dataset generation



OpenAI's Dall-E 2

Left: "Panda mad scientist mixing sparkling chemicals"  
Right: "A corgi's head depicted as an explosion in nebula"

Noisy / Enhanced



"We'll eat frozen pizzas all day.. All day every day"  
UNIVERSE: <https://serrjoa.github.io/projects/universe/>

# Background

## Diffusion-Based Speech Enhancement

- Task is **conditional generation**, we want to guide the process based on an impaired speech signal
- Ref. contribution\* takes theory from image processing
- Forward-process expressed using an SDE

$$d\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, \mathbf{y})dt + g(t)d\mathbf{w}$$

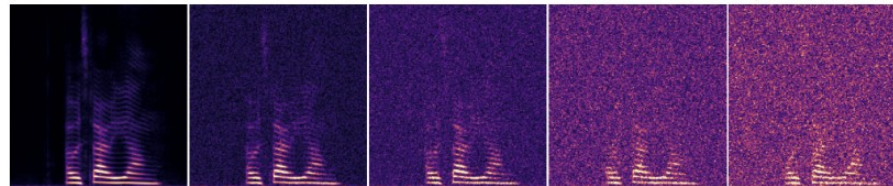
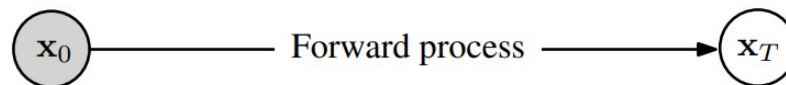


Image from\*

\*Richter, et.al. "Speech Enhancement and Dereverberation with Diffusion-Based Generative Models", IEEE/ACM Transactions on Audio, Speech, and Language Processing

# Background

## Diffusion-Based Speech Enhancement: Forward Process

$$d\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, \mathbf{y})dt + g(t)d\mathbf{w}$$

$\mathbf{f}(\mathbf{x}_t, \mathbf{y}) := \gamma(\mathbf{y} - \mathbf{x}_t)$  //  $\mathbf{y}$  = noisy,  $\mathbf{x}_0$  = clean,  $\gamma$  = *stiffness scalar*

$$g(t) := \sigma_{\min} \left( \frac{\sigma_{\max}}{\sigma_{\min}} \right)^t \sqrt{2 \log \left( \frac{\sigma_{\max}}{\sigma_{\min}} \right)}$$
 // set of  $d\mathbf{w}$  magnitudes

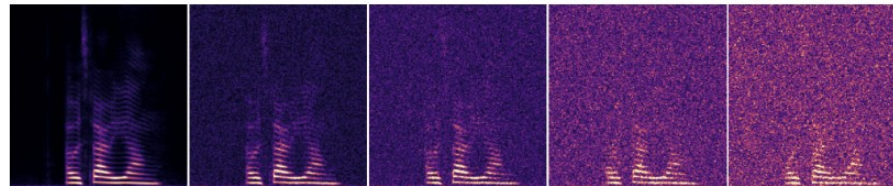
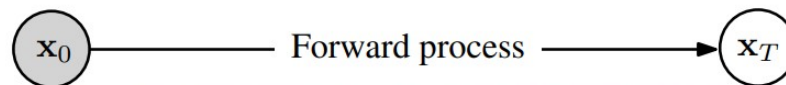


Image from\*

\*Richter, et.al. "Speech Enhancement and Dereverberation with Diffusion-Based Generative Models", IEEE/ACM Transactions on Audio, Speech, and Language Processing

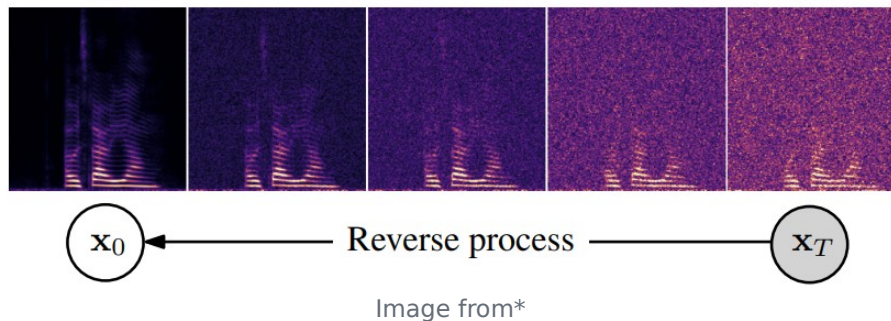
# Background

## Diffusion-Based Speech Enhancement: Reverse Process

- The reverse process is also called **sampling**

$$d\mathbf{x}_t = \left[ -\mathbf{f}(\mathbf{x}_t, \mathbf{y}) + g(t)^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{y}) \right] dt + g(t) d\bar{\mathbf{w}}$$

- Predict by integrating reverse time SDE with an SDE solver
- Correct by numerical optimization



\*Richter, et.al. "Speech Enhancement and Dereverberation with Diffusion-Based Generative Models", IEEE/ACM Transactions on Audio, Speech, and Language Processing



# Background

## Diffusion-Based Speech Enhancement: Reverse Process

$$d\mathbf{x}_t = \left[ -\mathbf{f}(\mathbf{x}_t, \mathbf{y}) + g(t)^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{y}) \right] dt + g(t) d\bar{\mathbf{w}}$$

$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{y})$  // Gradient of log-probability density w/r to data

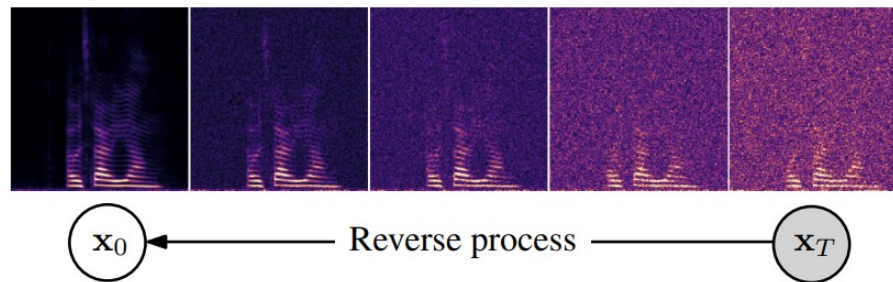


Image from\*

\*Richter, et.al. "Speech Enhancement and Dereverberation with Diffusion-Based Generative Models", IEEE/ACM Transactions on Audio, Speech, and Language Processing

# Background

## Diffusion-Based Speech Enhancement: Reverse Process

$$d\mathbf{x}_t = \left[ -\mathbf{f}(\mathbf{x}_t, \mathbf{y}) + g(t)^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{y}) \right] dt + g(t) d\bar{\mathbf{w}}$$

$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{y})$  // Gradient of log-probability density w/r to data

$d\mathbf{x}_t = [-\mathbf{f}(\mathbf{x}_t, \mathbf{y}) + g(t)^2 s_\theta(\mathbf{x}_t, \mathbf{y}, t)] dt$  // In practice a parametrized neural net  $s_\theta$

---

In the ideal scenario  $s_\theta = \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{y})$

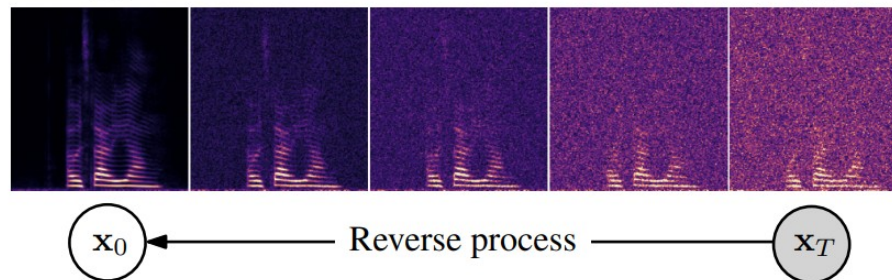


Image from\*

\*Richter, et.al. "Speech Enhancement and Dereverberation with Diffusion-Based Generative Models", IEEE/ACM Transactions on Audio, Speech, and Language Processing

# Background

## Diffusion-Based Speech Enhancement: Noise Schedule

- $\sigma$  controls the noise injections in forward/reverse processes within a defined min/max range

$$\sigma(t)^2 = \frac{\sigma_{\min}^2 \left( (\sigma_{\max}/\sigma_{\min})^{2t} - e^{-2\gamma t} \right) \log(\sigma_{\max}/\sigma_{\min})}{\gamma + \log(\sigma_{\max}/\sigma_{\min})}$$

- Used during

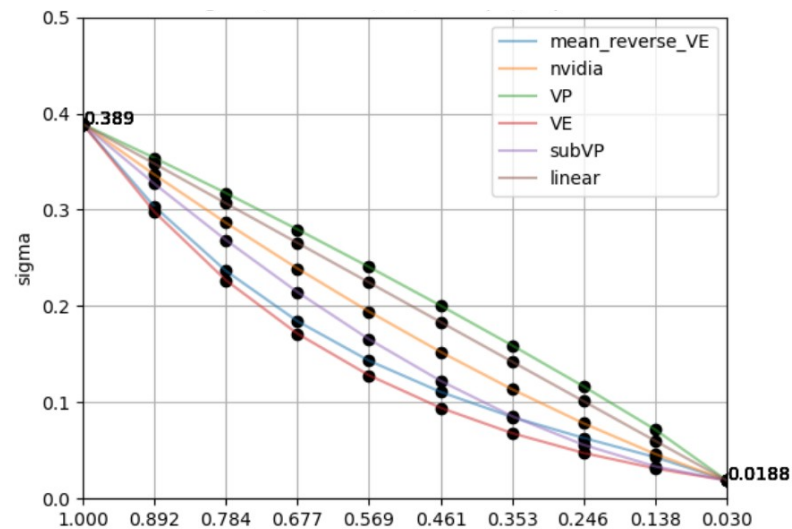
- Training:  $\arg \min_{\theta} \mathbb{E}_{t, (\mathbf{x}_0, \mathbf{y}), \mathbf{z}, \mathbf{x}_t | (\mathbf{x}_0, \mathbf{y})} \left[ \left\| \mathbf{s}_{\theta}(\mathbf{x}_t, \mathbf{y}, t) + \frac{\mathbf{z}}{\sigma(t)} \right\|_2^2 \right]$

- Sampling: Corrector algorithm (score and Gaussian noise scaling)

# Background

## Diffusion-Based Speech Enhancement: Noise Schedule

- Sampling routines in production do not need to follow training [2]
- Schedules can be fine-tuned during inference and is crucial for performance [1, 2]

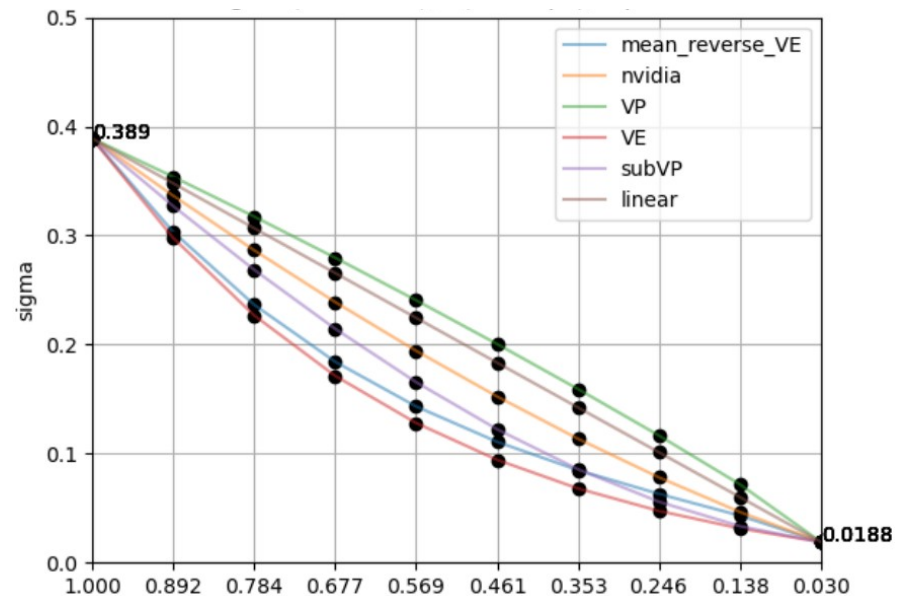


[1] Chen T., On the importance of noise scheduling for diffusion models. *arXiv preprint arXiv:2301.10972*, 2023

[2] Karras T., et al.. *Elucidating the design space of diffusion-based generative models*. *NeurIPS*, 2022

# Problem Statement

- Diffusion models have flexible, tunable moving components
  - Component of the sampling routine – noise schedule
- Investigate the interplay generated output and the noise scheduler
- Empirical study through a set of experiments



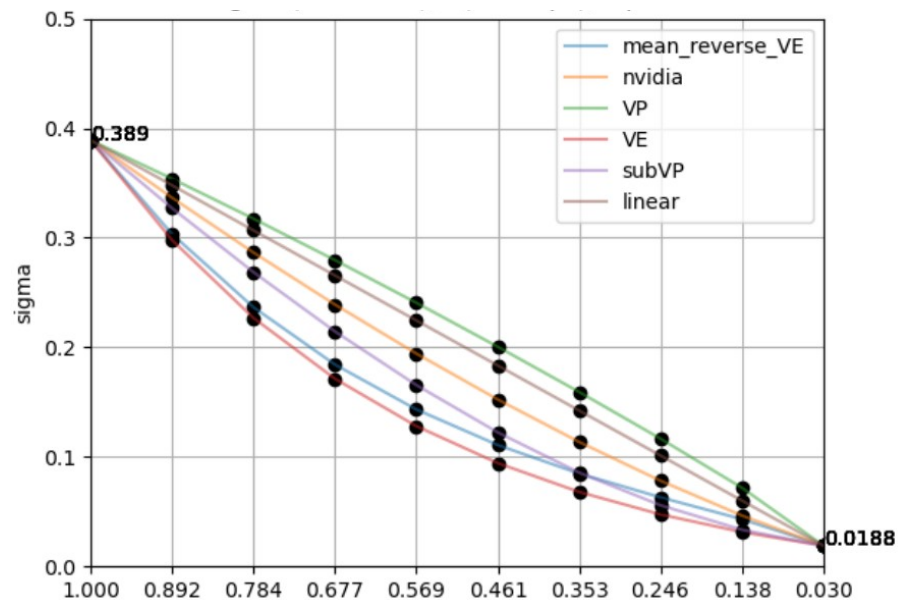
# Methodology

- Four experiments set up
  - Interchanged scheduler functions
  - Baseline scheduler with modified derivative
  - Baseline scheduler with modified timesteps
  - Baseline scheduler with non-uniform timesteps
- Quantified through perceptual metrics
  - STOI (0-1), WARP-Q ( $\dots - 0$ ), DNSMOS and PESQ (1 – 5)
- Results generated with *max* value per each metric
- Tests performed with a pseudo-random 35 datapoint test set
- Fixed model, stochasticity, sampling routine

# Experiment 1

## Interchanged Scheduler Functions

- Baseline function *mean\_reverse\_VE* swapped with 5 other functions
- Given 0.5 – 0.05 range, normalize to effective values of the baseline



# Experiment 1

## Interchanged Scheduler Functions

		PESQ	STOI	WARPQ	DNSMOS
Linear	n=10	<b>2.147</b> ( $\pm 0.617$ )	<b>0.911</b> ( $\pm 0.06$ )	0.776 ( $\pm 0.199$ )	2.909 ( $\pm 0.194$ )
	n=30	2.287 ( $\pm 0.631$ )	<b>0.923</b> ( $\pm 0.057$ )	0.755 ( $\pm 0.187$ )	2.934 ( $\pm 0.186$ )
SubVP	n=10	1.503 ( $\pm 0.261$ )	0.891 ( $\pm 0.062$ )	0.884 ( $\pm 0.142$ )	2.549 ( $\pm 0.252$ )
	n=30	2.202 ( $\pm 0.631$ )	<b>0.923</b> ( $\pm 0.057$ )	0.775 ( $\pm 0.193$ )	2.898 ( $\pm 0.197$ )
VE	n=10	1.62 ( $\pm 0.349$ )	0.885 ( $\pm 0.064$ )	0.907 ( $\pm 0.149$ )	2.619 ( $\pm 0.289$ )
	n=30	2.13 ( $\pm 0.641$ )	0.92 ( $\pm 0.058$ )	0.774 ( $\pm 0.187$ )	2.861 ( $\pm 0.208$ )
VP	n=10	2.122 ( $\pm 0.568$ )	0.907 ( $\pm 0.061$ )	<b>0.774</b> ( $\pm 0.192$ )	<b>2.925</b> ( $\pm 0.184$ )
	n=30	<b>2.31</b> ( $\pm 0.609$ )	0.92 ( $\pm 0.059$ )	<b>0.748</b> ( $\pm 0.179$ )	<b>2.96</b> ( $\pm 0.179$ )
Nvidia	n=10	1.883 ( $\pm 0.493$ )	0.908 ( $\pm 0.062$ )	0.812 ( $\pm 0.177$ )	2.824 ( $\pm 0.212$ )
	n=30	2.249 ( $\pm 0.64$ )	<b>0.923</b> ( $\pm 0.057$ )	0.759 ( $\pm 0.188$ )	2.909 ( $\pm 0.196$ )
Baseline	n=10	1.49 ( $\pm 0.249$ )	0.884 ( $\pm 0.067$ )	0.896 ( $\pm 0.139$ )	2.561 ( $\pm 0.265$ )
	n=30	2.175 ( $\pm 0.671$ )	0.922 ( $\pm 0.056$ )	0.765 ( $\pm 0.198$ )	2.927 ( $\pm 0.187$ )

NOISY



BASELINE

N=10



N=30



VP

N=10



N=30



CLEAN

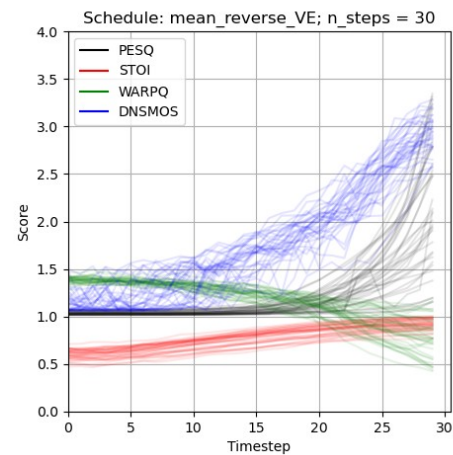
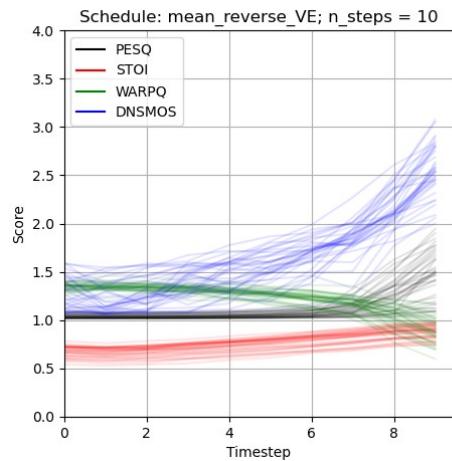
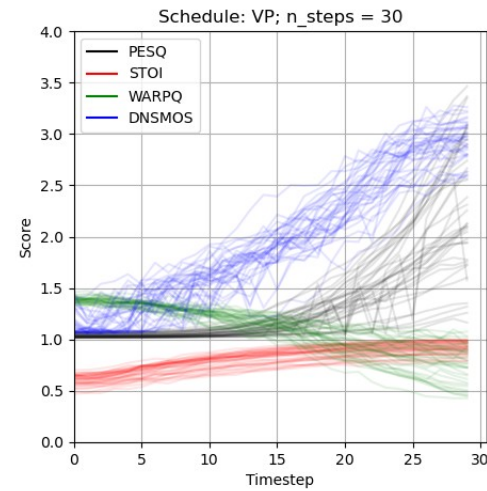
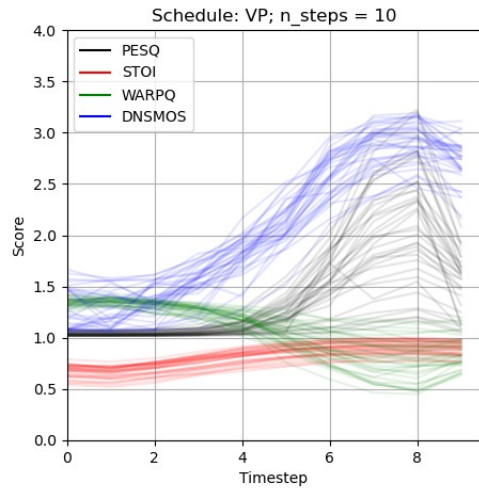


*“Downing street will make the second appointment in the Scotland office today..”*



# Experiment 1

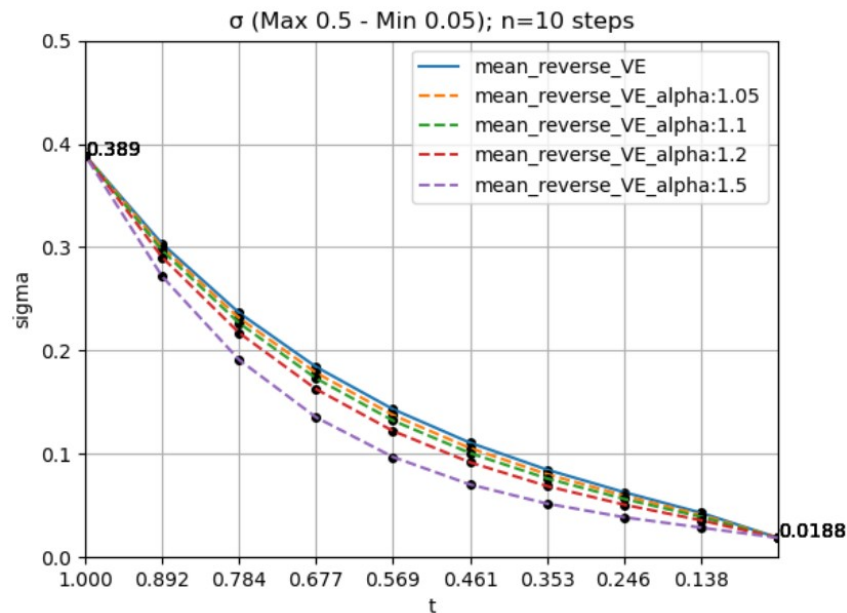
## Interchanged Scheduler Functions



# Experiment 2

## Baseline Function with Modified Derivative

- Introduce an  $\alpha$  modifier, vary the curve steepness
- Idea: significant improvement threshold decrease with potentially higher peak performance



# Experiment 2

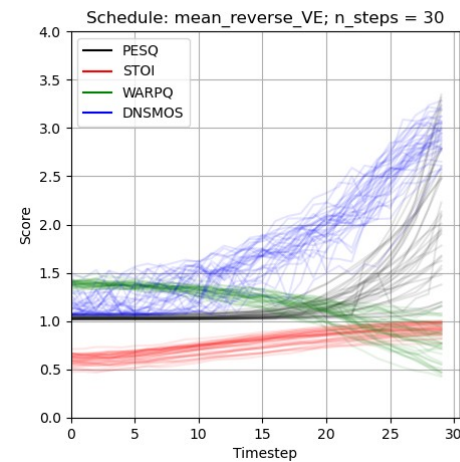
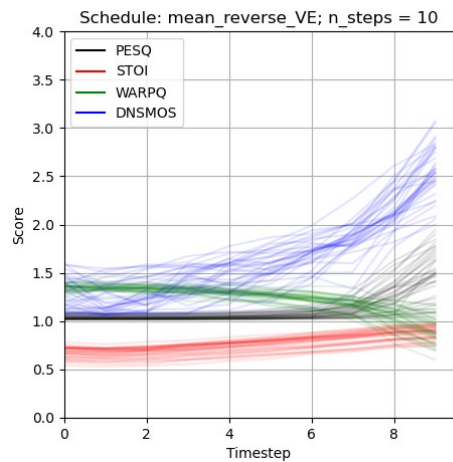
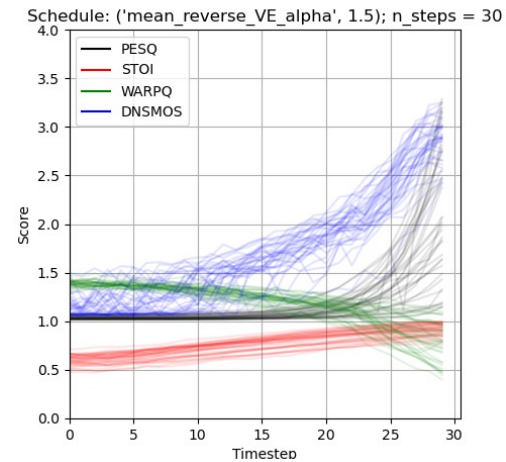
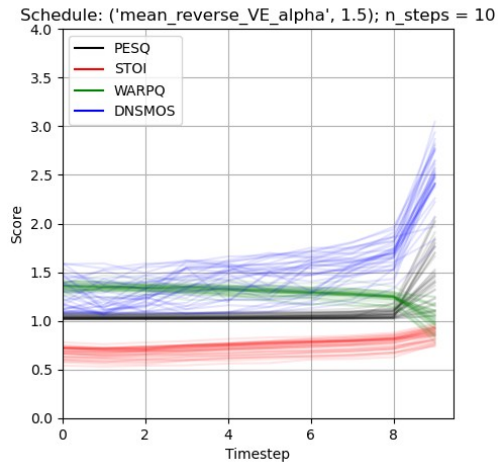
## Baseline Function with Modified Derivative

- N=10 mild-to-none increase in performance
- N=30 mild-to-none decrease of performance
- No clear winner for the modifications

		PESQ	STOI	WARPQ	DNSMOS
Alpha=1.05	n=10	1.498 ( $\pm$ 0.26)	0.882 ( $\pm$ 0.067)	0.905 ( $\pm$ 0.139)	2.563 ( $\pm$ 0.272)
	n=30	2.17 ( $\pm$ 0.671)	0.921 ( $\pm$ 0.056)	<b>0.764</b> ( $\pm$ 0.2)	<b>2.928</b> ( $\pm$ 0.186)
Alpha=1.1	n=10	1.513 ( $\pm$ 0.272)	0.881 ( $\pm$ 0.067)	0.913 ( $\pm$ 0.14)	2.566 ( $\pm$ 0.28)
	n=30	2.167 ( $\pm$ 0.67)	0.921 ( $\pm$ 0.057)	0.768 ( $\pm$ 0.201)	2.912 ( $\pm$ 0.191)
Alpha=1.2	n=10	<b>1.548</b> ( $\pm$ 0.304)	0.879 ( $\pm$ 0.066)	0.928 ( $\pm$ 0.144)	<b>2.586</b> ( $\pm$ 0.287)
	n=30	2.156 ( $\pm$ 0.67)	0.921 ( $\pm$ 0.057)	0.767 ( $\pm$ 0.2)	2.912 ( $\pm$ 0.194)
Alpha=1.5	n=10	1.504 ( $\pm$ 0.296)	0.864 ( $\pm$ 0.064)	1.017 ( $\pm$ 0.114)	2.498 ( $\pm$ 0.284)
	n=30	2.089 ( $\pm$ 0.68)	0.916 ( $\pm$ 0.06)	0.782 ( $\pm$ 0.206)	2.888 ( $\pm$ 0.217)
Baseline	n=10	1.49 ( $\pm$ 0.249)	<b>0.884</b> ( $\pm$ 0.067)	<b>0.896</b> ( $\pm$ 0.139)	2.561 ( $\pm$ 0.265)
	n=30	<b>2.175</b> ( $\pm$ 0.671)	<b>0.922</b> ( $\pm$ 0.056)	0.765 ( $\pm$ 0.198)	2.927 ( $\pm$ 0.187)

# Experiment 2

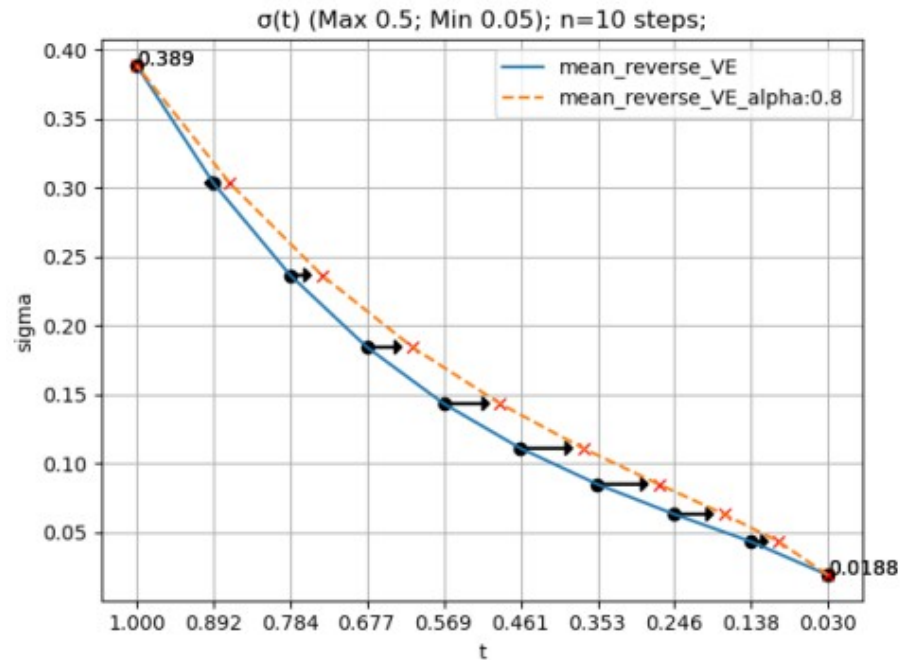
## Baseline Function with Modified Derivative



# Experiment 3

## Baseline Function with Timestep Offset

- Apply an  $\alpha$  modifier and project the unmodified  $\sigma$
- Idea: Compress the sampling points, decrease the amount of noise per step



# Experiment 3

## Baseline Function with Timestep Offset

- Consistent improvement over baseline
- Effect decays with increase of step size

		PESQ	STOI	WARPQ	DNSMOS
Alpha=0.8	n=10	<b>1.945</b> ( $\pm 0.526$ )	<b>0.903</b> ( $\pm 0.064$ )	<b>0.834</b> ( $\pm 0.169$ )	<b>2.809</b> ( $\pm 0.234$ )
	n=15	<b>2.135</b> ( $\pm 0.606$ )	<b>0.915</b> ( $\pm 0.06$ )	<b>0.779</b> ( $\pm 0.189$ )	<b>2.918</b> ( $\pm 0.227$ )
	n=30	<b>2.211</b> ( $\pm 0.669$ )	<b>0.923</b> ( $\pm 0.55$ )	<b>0.759</b> ( $\pm 0.201$ )	<b>2.944</b> ( $\pm 0.181$ )
Baseline	n=10	1.49 ( $\pm 0.249$ )	0.884 ( $\pm 0.067$ )	0.896 ( $\pm 0.139$ )	2.561 ( $\pm 0.265$ )
	n=15	1.93 ( $\pm 0.471$ )	0.911 ( $\pm 0.06$ )	0.805 ( $\pm 0.178$ )	2.858 ( $\pm 0.238$ )
	n=30	2.175 ( $\pm 0.671$ )	0.922 ( $\pm 0.56$ )	0.765 ( $\pm 0.198$ )	2.927 ( $\pm 0.187$ )

NOISY



BASELINE

N=10



N=15

Alpha 0.8

N=10



N=15



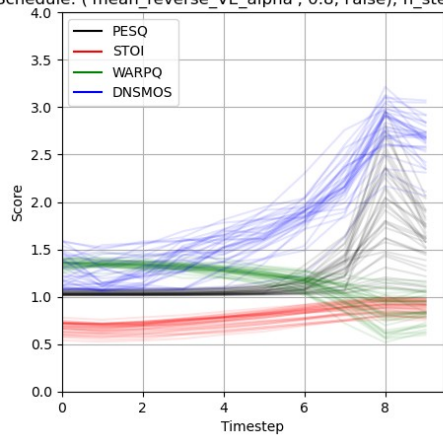
CLEAN



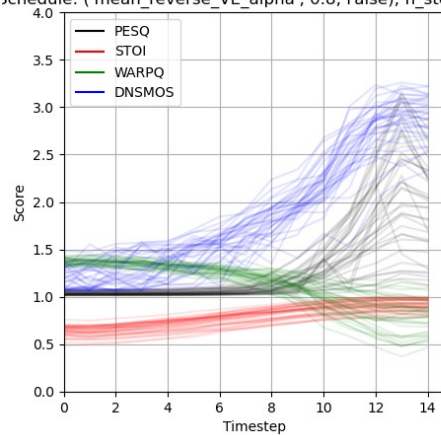
# Experiment 3

## Baseline Function with Timestep Offset

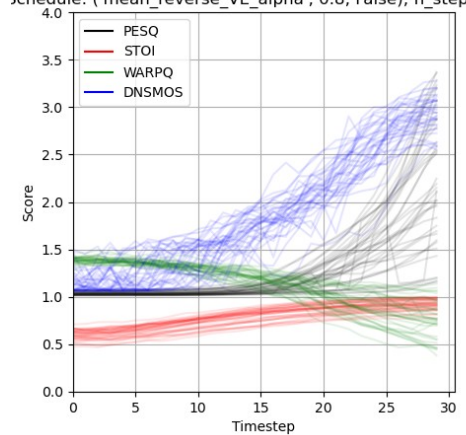
schedule: ('mean\_reverse\_VE\_alpha', 0.8, False); n\_steps = :



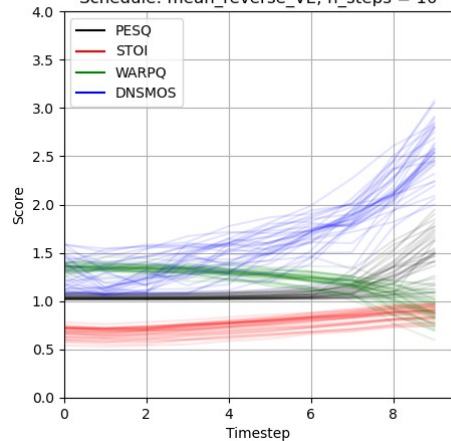
schedule: ('mean\_reverse\_VE\_alpha', 0.8, False); n\_steps = :



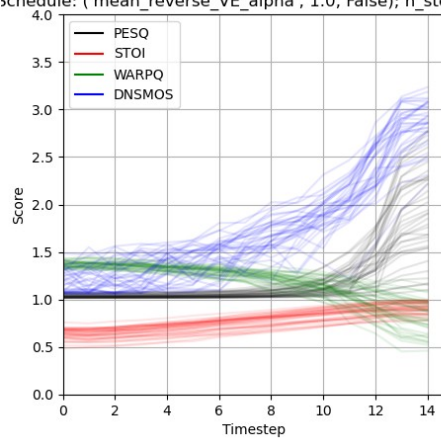
schedule: ('mean\_reverse\_VE\_alpha', 0.8, False); n\_steps = :



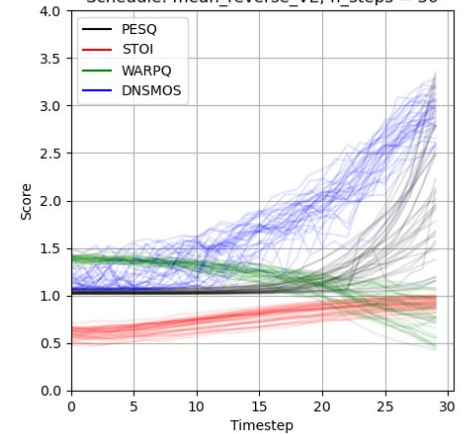
Schedule: mean\_reverse\_VE; n\_steps = 10



schedule: ('mean\_reverse\_VE\_alpha', 1.0, False); n\_steps = :



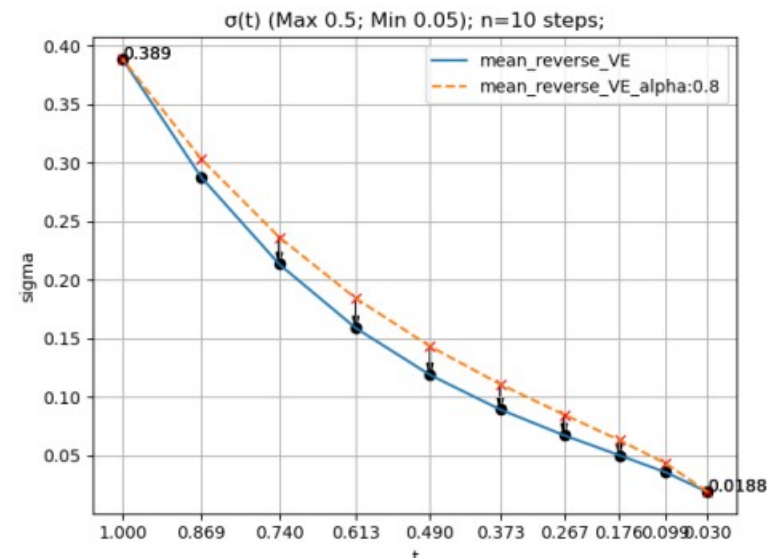
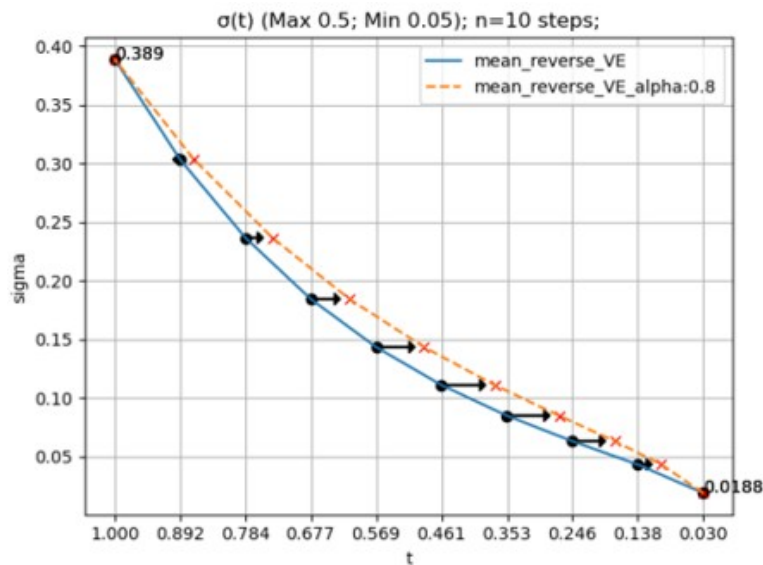
Schedule: mean\_reverse\_VE; n\_steps = 30



# Experiment 4

## Baseline Function with Non-uniform Timesteps

- Exp 3. schedule  $\sigma$  projected back to the baseline
- Idea: Compress sampling points towards the end, move back to familiar  $\sigma$ -t ratios





# Experiment 4

## Baseline Function with Non-uniform Timesteps

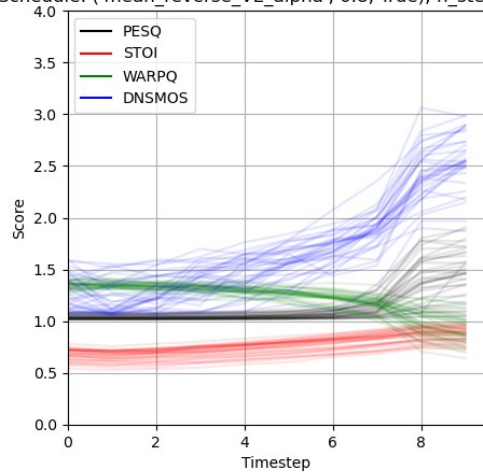
- Minor-to-none improvements in performance
- More performant means – higher std.dev.

		PESQ	STOI	WARPQ	DNSMOS
Alpha=0.8	n=10	1.464 ( $\pm$ 0.238)	0.878 ( $\pm$ 0.067)	0.912 ( $\pm$ 0.134)	<b>2.567</b> ( $\pm$ 0.265)
	n=15	<b>2.1</b> ( $\pm$ 0.603)	0.91 ( $\pm$ 0.061)	<b>0.791</b> ( $\pm$ 0.193)	<b>2.884</b> ( $\pm$ 0.253)
	n=30	<b>2.187</b> ( $\pm$ 0.671)	0.922 ( $\pm$ 0.056)	<b>0.761</b> ( $\pm$ 0.203)	<b>2.936</b> ( $\pm$ 0.192)
Baseline	n=10	<b>1.49</b> ( $\pm$ 0.249)	<b>0.884</b> ( $\pm$ 0.067)	<b>0.896</b> ( $\pm$ 0.139)	2.561 ( $\pm$ 0.265)
	n=15	1.93 ( $\pm$ 0.471)	<b>0.911</b> ( $\pm$ 0.06)	0.805 ( $\pm$ 0.178)	2.858 ( $\pm$ 0.238)
	n=30	2.175 ( $\pm$ 0.671)	0.922 ( $\pm$ 0.56)	0.765 ( $\pm$ 0.198)	2.927 ( $\pm$ 0.187)

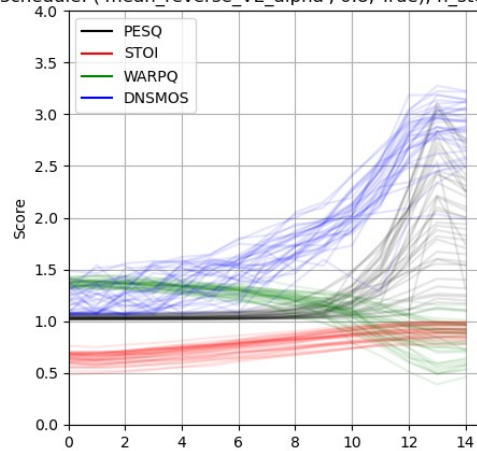
# Experiment 4

## Baseline Function with Non-uniform Timesteps

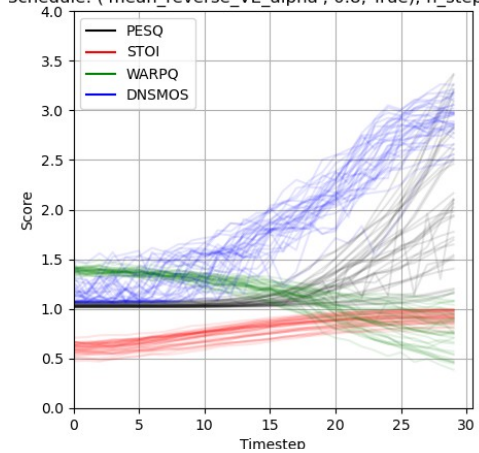
Schedule: ('mean\_reverse\_VE\_alpha', 0.8, True); n\_steps = 1



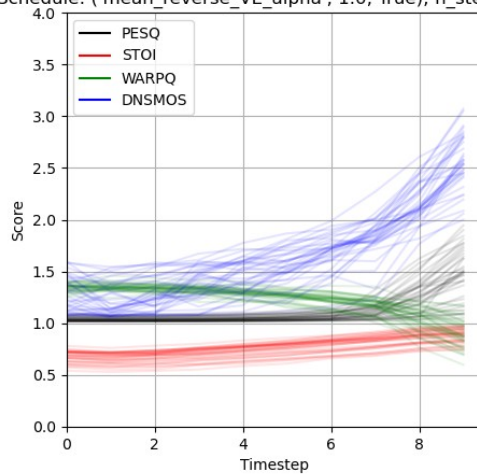
Schedule: ('mean\_reverse\_VE\_alpha', 0.8, True); n\_steps = 1



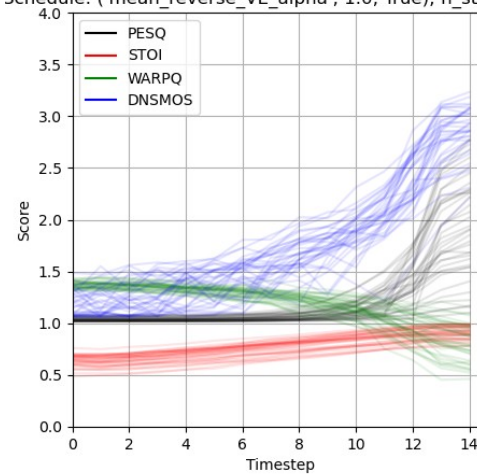
Schedule: ('mean\_reverse\_VE\_alpha', 0.8, True); n\_steps = 3



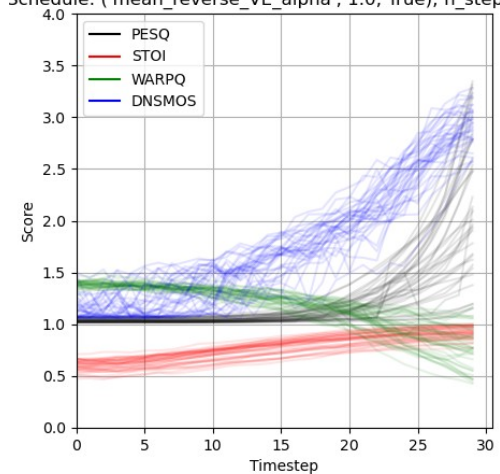
Schedule: ('mean\_reverse\_VE\_alpha', 1.0, True); n\_steps = 1



Schedule: ('mean\_reverse\_VE\_alpha', 1.0, True); n\_steps = 1



Schedule: ('mean\_reverse\_VE\_alpha', 1.0, True); n\_steps = 3



# Summary

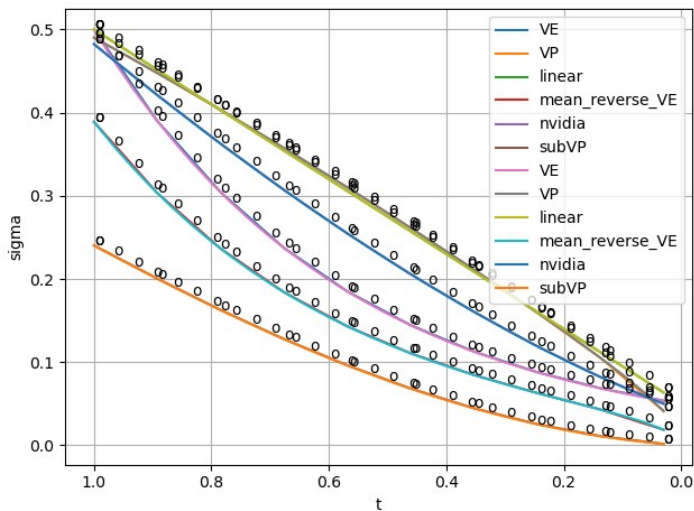
- Interchanged schedule and modified baseline may improve the baseline performance
  - Theory suggests that L2 models tend to remove too much noise
- Effects decline with increased timesteps
  - Robustness of more gradual diffusion
- Quality in may be more dependent on the progression rather than discretization itself
  - Exp 4 compression of timesteps alone does not yield improvements and may degrade results

# Thank you!

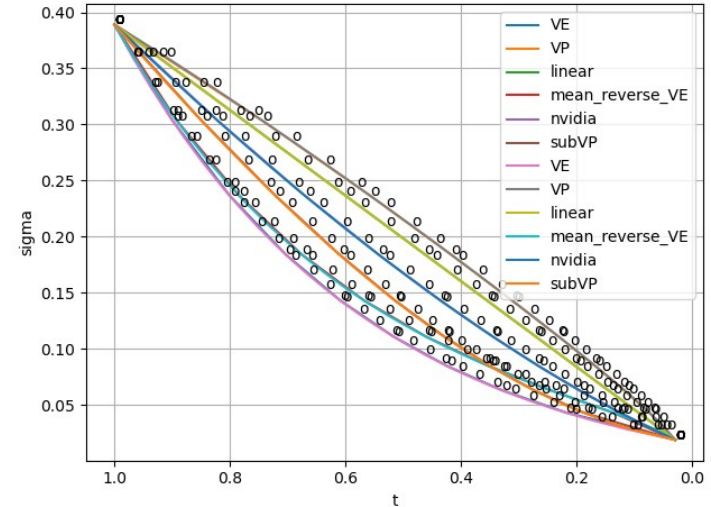
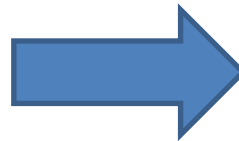
# Backup 1

## Experiment 1: Normalized vs Unnormalized functions

- For range (0.5 – 0.05) different effective values
- Normalize via linear scaling ->  $f(\sigma, x, y) = \frac{(\sigma - \sigma_{min})}{(\sigma_{max} - \sigma_{min})} * (y - x) + x$

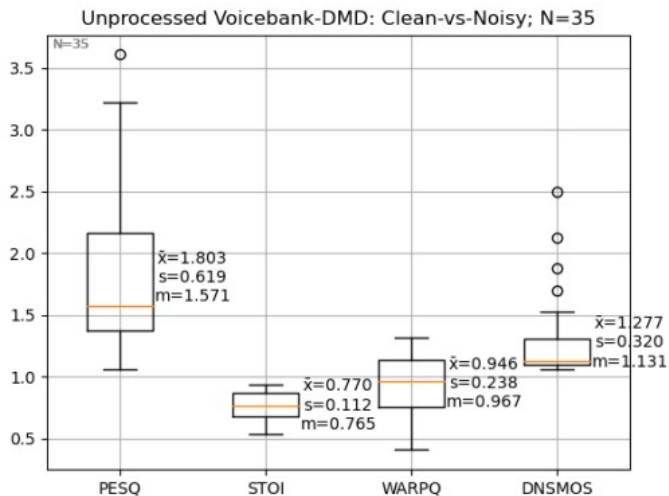


After normalization

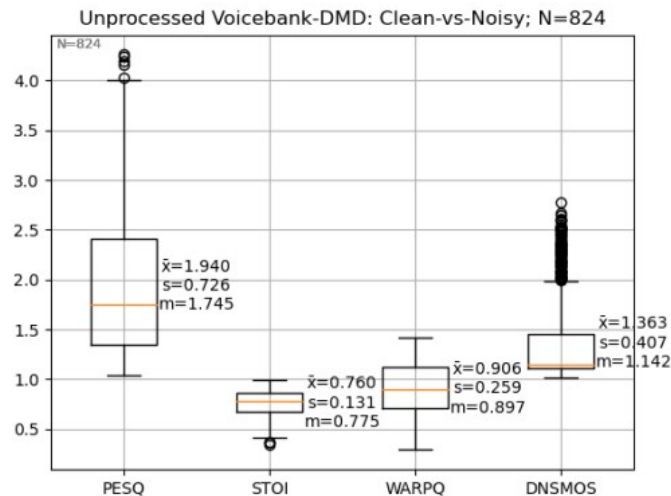


# Backup 2

## Metrics of Unprocessed Dataset



(a) 35 random sample subsplit.



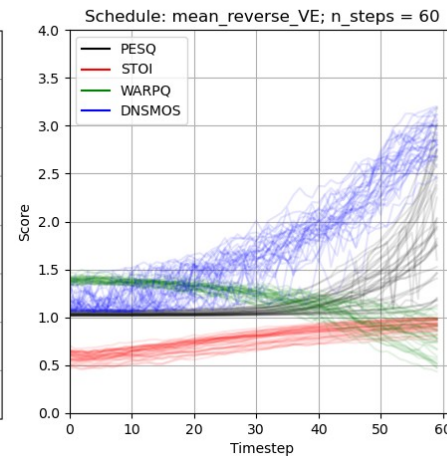
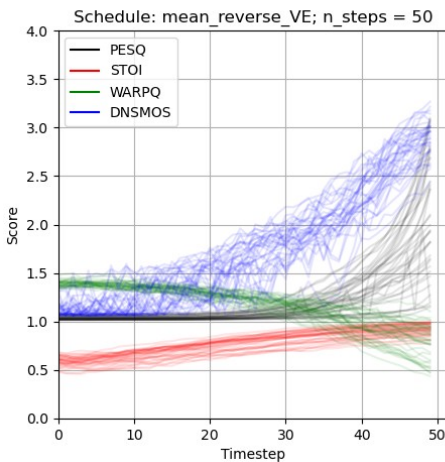
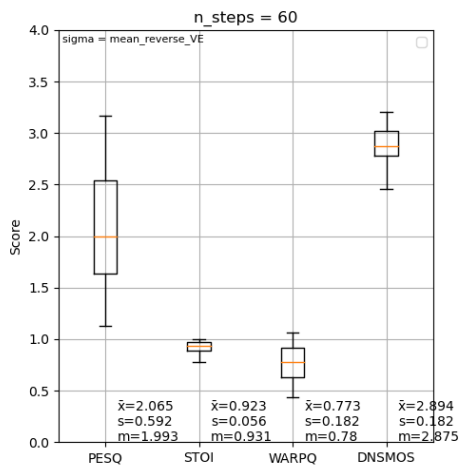
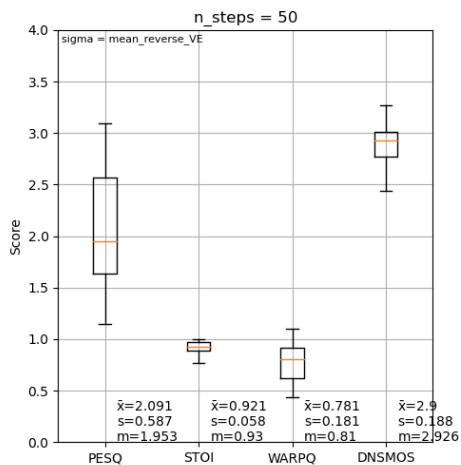
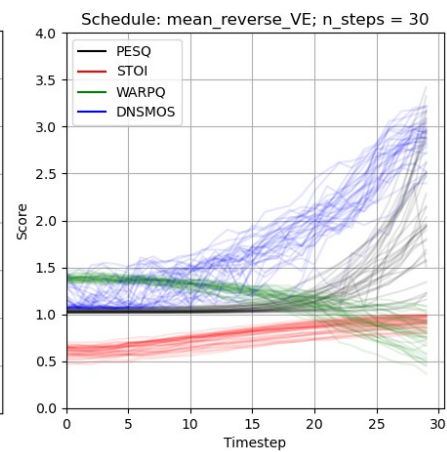
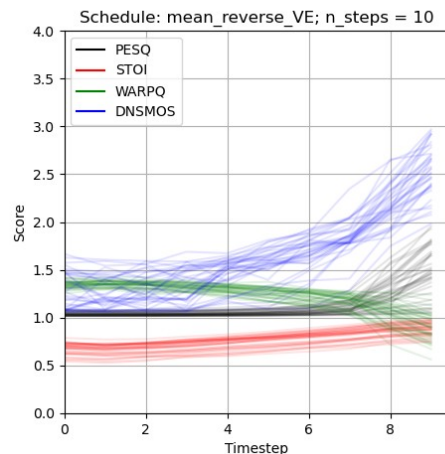
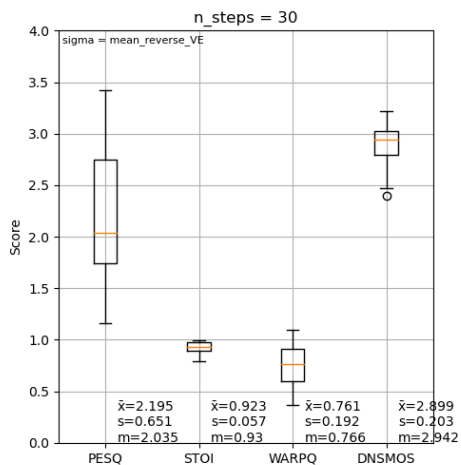
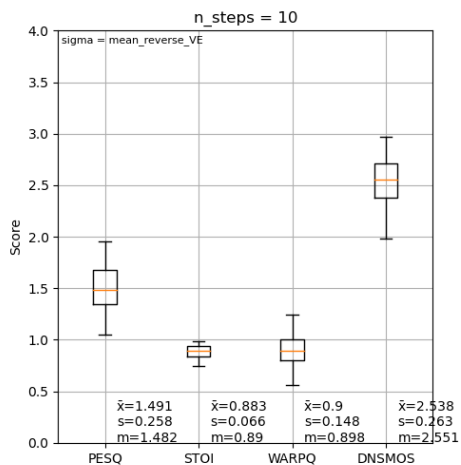
(b) Full dataset.

Baseline	n=10	1.49 ( $\pm 0.249$ )	0.884 ( $\pm 0.067$ )	0.896 ( $\pm 0.139$ )	2.561 ( $\pm 0.265$ )
	n=30	2.175 ( $\pm 0.671$ )	0.922 ( $\pm 0.056$ )	0.765 ( $\pm 0.198$ )	2.927 ( $\pm 0.187$ )

*Baseline model processing results*

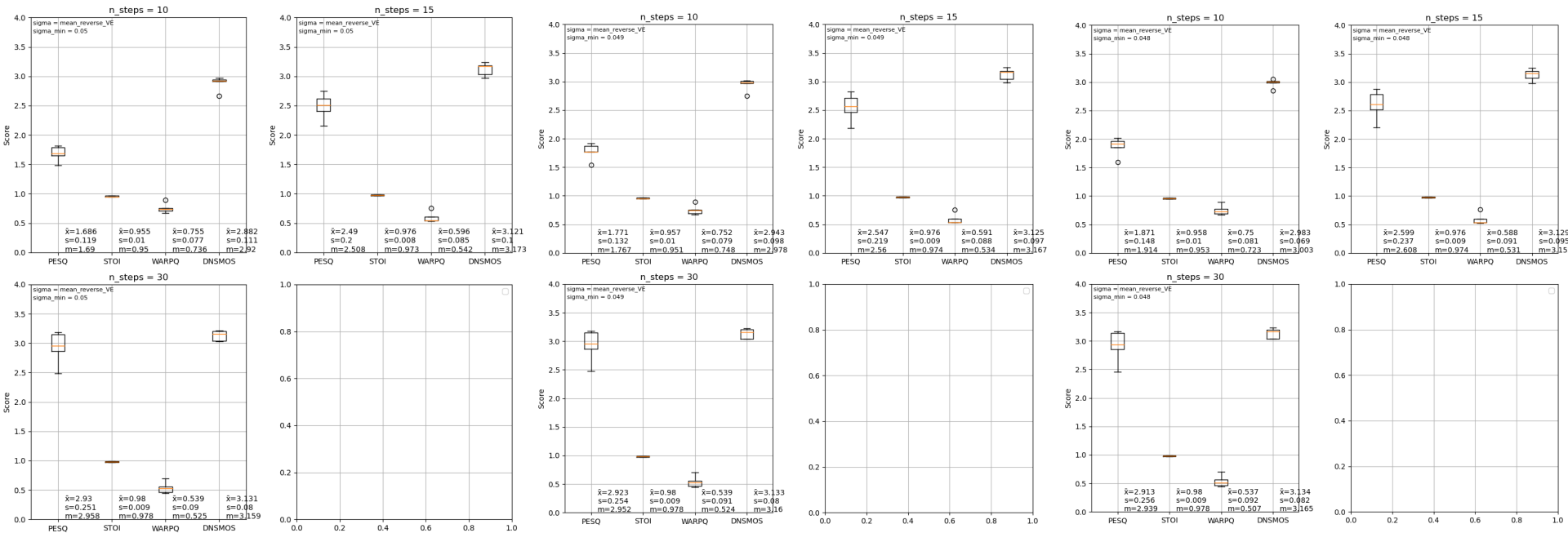
# Backup 3

## Further Increased Timesteps



# Backup 4

## Decreased Sigmas





# Backup 4

## Decreased Sigmas

