



## MOTIVATION

We learn **forward** and **reverse** diffusion processes **jointly** to achieve **SOTA** likelihoods with diffusion models.

### Goal

- Jointly optimize the “forward / noising” process and the “reverse / denoising” process for a **tighter** ELBO.
- Current diffusion models **ONLY** optimize the reverse process.

### Challenge

ELBO is **invariant** to (scalar) noise schedules [1].

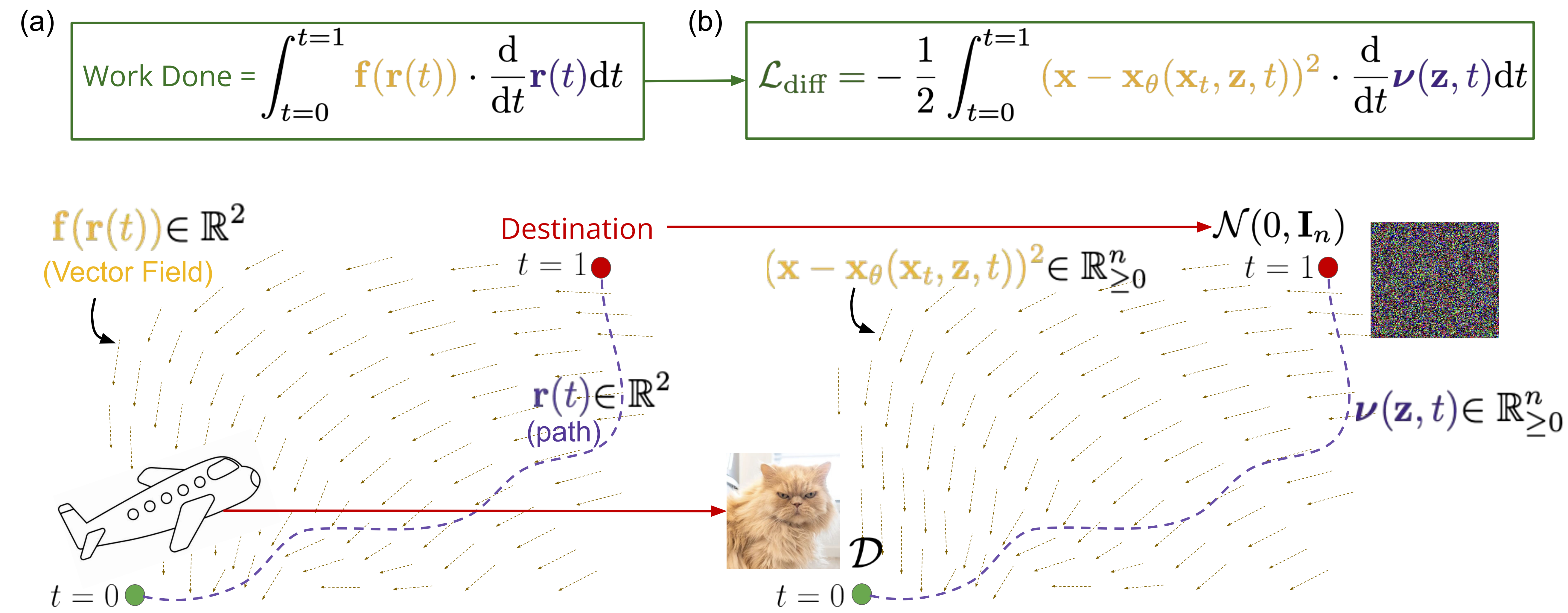
### Solution

- Key Discovery:** Multivariate schedules can alter ELBO.
- Multivariate Learned Adaptive Noise schedule (MuLAN)** **tightens** the ELBO.

| Noise Schedule Properties | MuLAN | Scalar |
|---------------------------|-------|--------|
| Multivariate              | ✓     | ✗      |
| Learned                   | ✓     | ✗      |
| Adaptive                  | ✓     | ✗      |
| Improves ELBO             | ✓     | ✗      |

## INTUITION

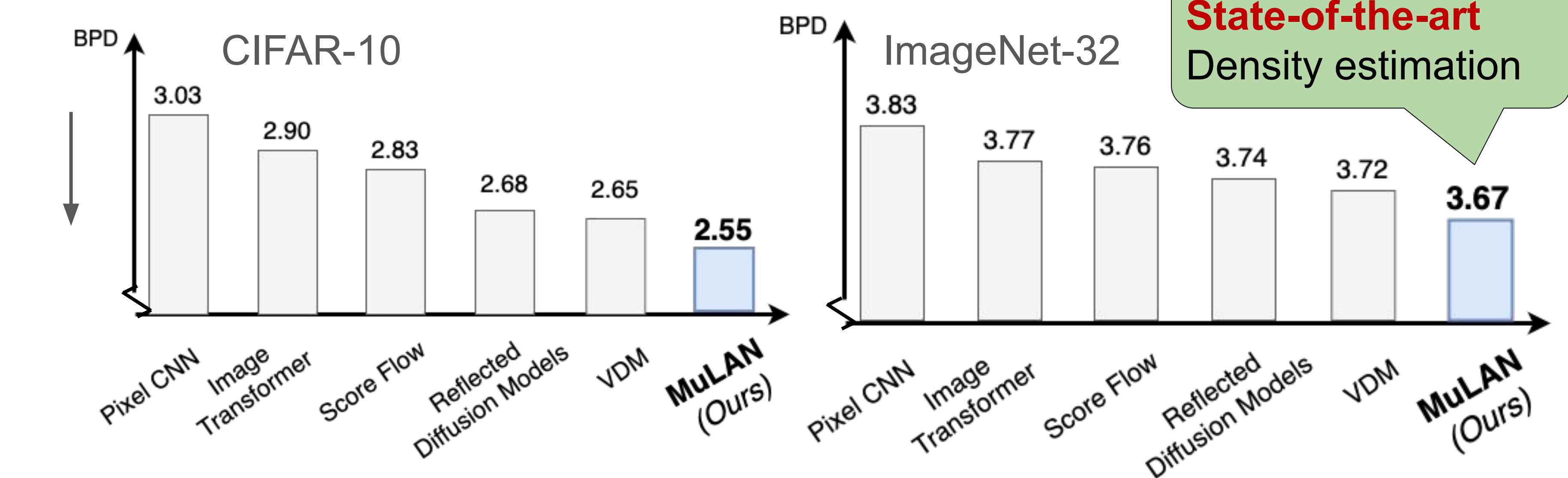
Objects in a **vector field** follow the **path of least resistance**, **NOT** the shortest Euclidean (straight-line) path.



- The same intuition maps to diffusion processes:
  - Objects:** Images
  - Vector field:** per-pixel reconstruction error
  - Path:** noise schedule
  - Work Done:** Diffusion loss  $\mathcal{L}_{\text{diff}}$
- Scalar** noise schedule: **shortest Euclidean path** (i.e., straight line) from the data distribution to prior.
- Since **MuLAN** is multivariate (and adaptive), it learns a **path of least resistance** i.e., a noise schedule corresponding to the **optimal ELBO**.

## RESULTS

### Likelihood estimation on CIFAR-10 and ImageNet-32

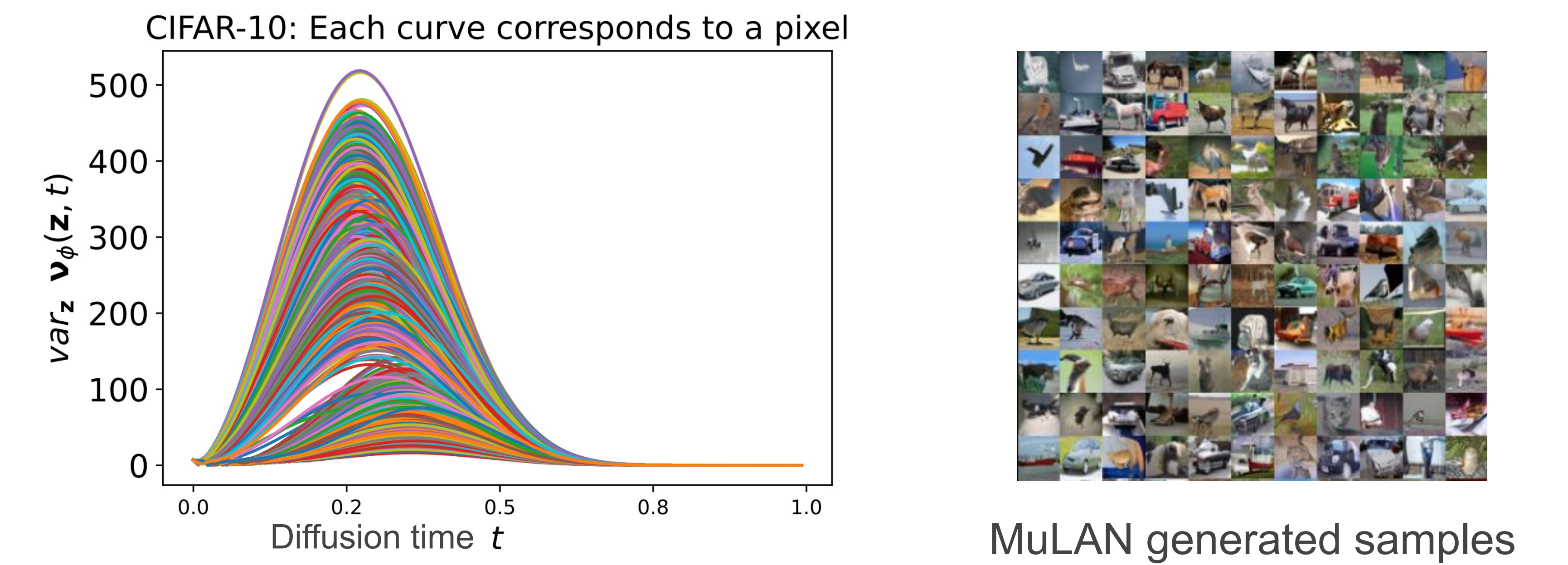


### Faster Training

- MuLAN augmented VDM achieves a BPD of
  - 2.65 in **2M** steps (vs. **10M**) on CIFAR-10.
  - 3.72 in **1M** steps (vs. **2M**) on ImageNet-32.

Upto **5x** faster training!

### Noise Schedule variations across different images



## NOTATION

$\mathbf{x}_0 \sim \mathcal{D}$  : Sample from the Data distribution  
 $\mathbf{x}_{(\cdot)} \in \mathbb{R}^n$   
 $\mathbf{z} \in \{0, 1\}^m$   
 $m < n$   
 $0 \leq s < t \leq 1$   
 $\alpha_s(\mathbf{z}), \alpha_t(\mathbf{z}) \in [0, 1]^n$      $\sigma_s(\mathbf{z}), \sigma_t(\mathbf{z}) \in [0, 1]^n$   
 $\alpha_{t|s}(\mathbf{z}) = \alpha_t(\mathbf{z}) / \alpha_s(\mathbf{z})$      $\sigma_{t|s}(\mathbf{z}) = \sigma_t(\mathbf{z}) / \sigma_s(\mathbf{z})$   
 $\nu_\phi(\mathbf{z}, t) = \alpha_t^2(\mathbf{z}) / \sigma_t^2(\mathbf{z})$   
 $\cdot$  denotes dot product     $\odot$  : element wise mult.  
 $\mathbf{x}_\theta(\mathbf{x}_t, \mathbf{z}, t)$  : Denoising Model with parameters  $\theta$   
 $\phi$  : Parameters of the Encoder network

## REFERENCES

[1] Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. Advances in neural information processing systems, 34:21696–21707, 2021.

[2] Subham S. Sahoo, Anselm Paulus, Marin Vlastelica, Vit Musil, Volodymyr Kuleshov, Georg Martius. Backpropagation through Combinatorial Algorithms: Identity with Projection Works. International Conference on Learning Representations (ICLR - 2023), 2023.

## OUR METHOD: Multivariate Learned Adaptive Noise (MuLAN)

