# ScaleDreamer: Scalable Text-to-3D Synthesis with Asynchronous Score Distillation

Zhiyuan Ma, Yuxiang Wei, Yabin Zhang, Xiangyu Zhu,
Zhen Lei, and Lei Zhang
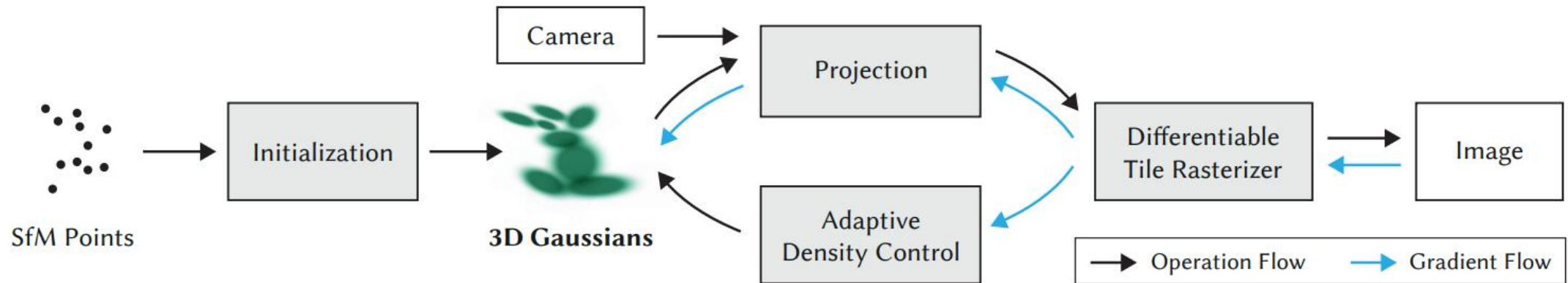
ECCV 2024

STRUCT Group Seminar
Presenter: Yifan Li
2024.7.20

# Outline

- Background

- Method

- Experiments

- Conclusion

# Background

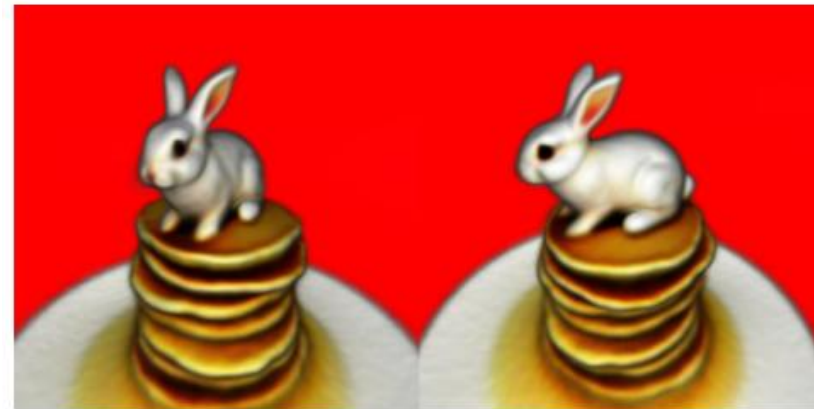## 3D Reconstruction



Explicit access of constraints and prior knowledge
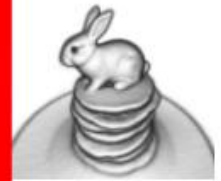
# Background

## 3D Generation



a tiger dressed as a doctor*

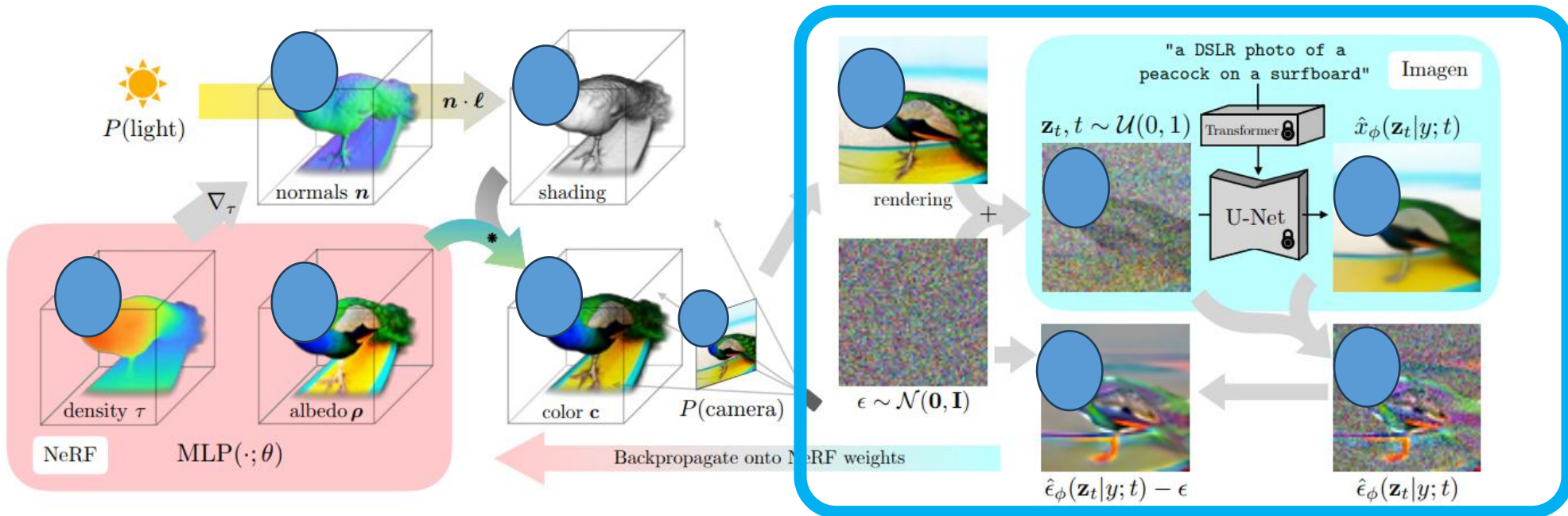a baby bunny sitting on top of a stack of pancakes†

Lack of enough external priors to generate a high quality 3D object

High quality huge scale 3D dataset is hard to collect

How to utilize strong generative power of **2D Diffusion Models** to 3D?

# Background

## DreamFusion: Diffusion Model as a loss



## Score Distillation Sampling (SDS)

"DreamFusion: Text-to-3D using 2D Diffusion", Ben Poole, Ajay Jain, Jonathan T. Barron, Ben Mildenhall, ICLR23 outstanding paper
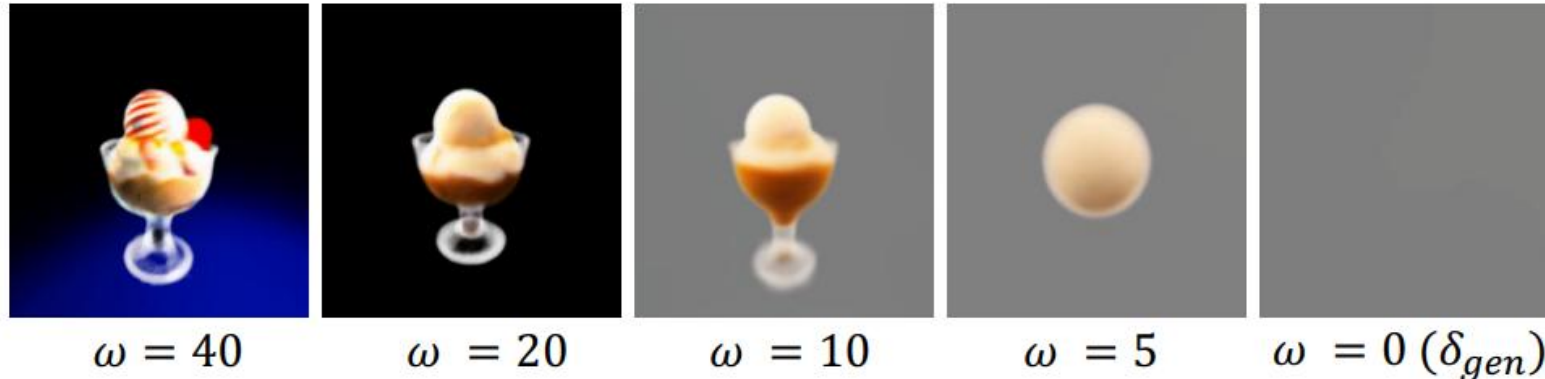
# Background: DreamFusion

- Pros:

  - Do not need to backpropagate through the diffusion model

  - DM simply acts like an efficient, frozen critic predicts image-space edits

  - Effectively insert 2D DM's generative priors to produce 3D objects

- Cons:

  - Need to set the Classifier Free Guidance as high as 100 for convergence

  - Produce excessively large gradients and lead to unstable optimization

  - High-saturation results

# Background

## Classifier Score Distillation (CSD)

Classifier score is the true essential component that drives the optimization

$$\delta_x(\mathbf{x}_t; y, t) = \underbrace{[\epsilon_\phi(\mathbf{x}_t; y, t) - \epsilon]}_{\delta_x^{\text{gen}}} + \omega \cdot \underbrace{[\epsilon_\phi(\mathbf{x}_t; y, t) - \epsilon_\phi(\mathbf{x}_t; t)]}_{\delta_x^{\text{cls}}}$$



$\omega = 40$      $\omega = 20$      $\omega = 10$      $\omega = 5$      $\omega = 0 \, (\delta_{gen})$

$\omega$: classifier free guidance intensity

"Text-to-3D with Classifier Score Distillation", Xin Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Song-Hai Zhang, Xiaojuan Qi,  arXiv 23.10

# Background

## Classifier Score Distillation (CSD)

Replace "ground truth noise"

$$\nabla_\theta \mathcal{L}_{\text{SDS}} = \mathbb{E}_{t,\epsilon,\mathbf{c}} \left[ w(t)(\epsilon_\phi(\mathbf{x}_t; y, t) - \boxed{\epsilon)} \frac{\partial \mathbf{x}}{\partial \theta} \right]$$

$$\nabla_\theta \mathcal{L}_{\text{CSD}} = \mathbb{E}_{t,\epsilon,\mathbf{c}} \left[ w(t)(\epsilon_\phi(\mathbf{x}_t; y, t) - \boxed{\epsilon_\phi(\mathbf{x}_t; t))} \frac{\partial \mathbf{x}}{\partial \theta} \right]$$

$\theta$: parameters of 3D model (NeRF, …), used to generate a rendered 2D image
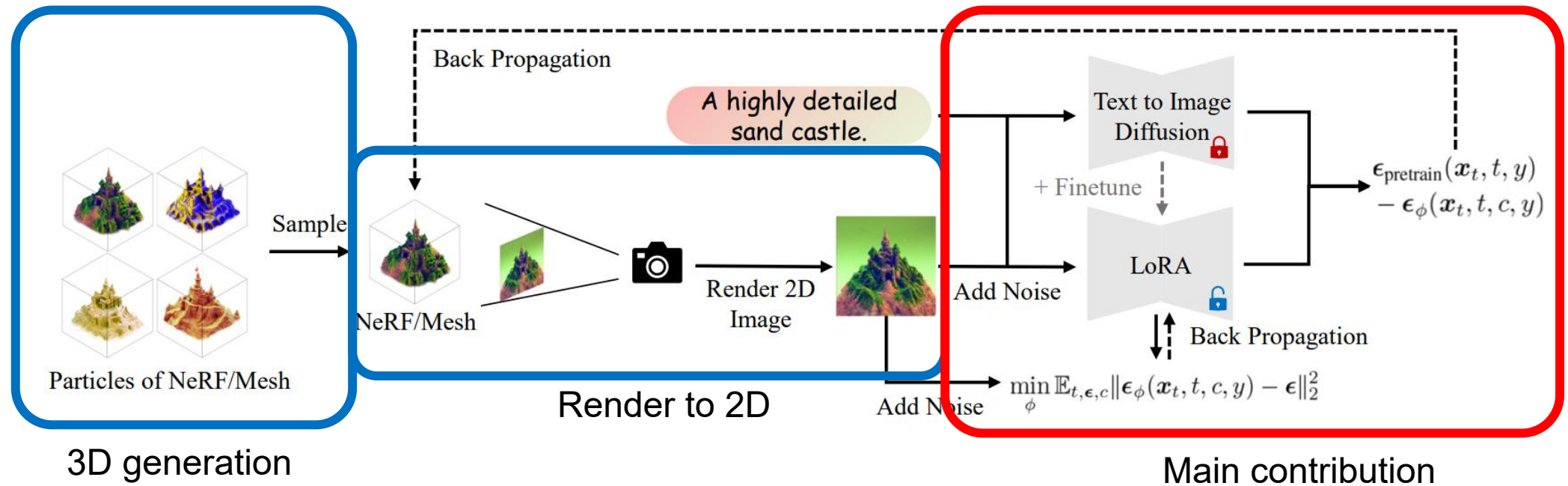$\phi$: parameters of diffusion model
$y$: text prompt

# Background

## Variational Score Distillation (VSD)

Predict noise adaptively and more accurately

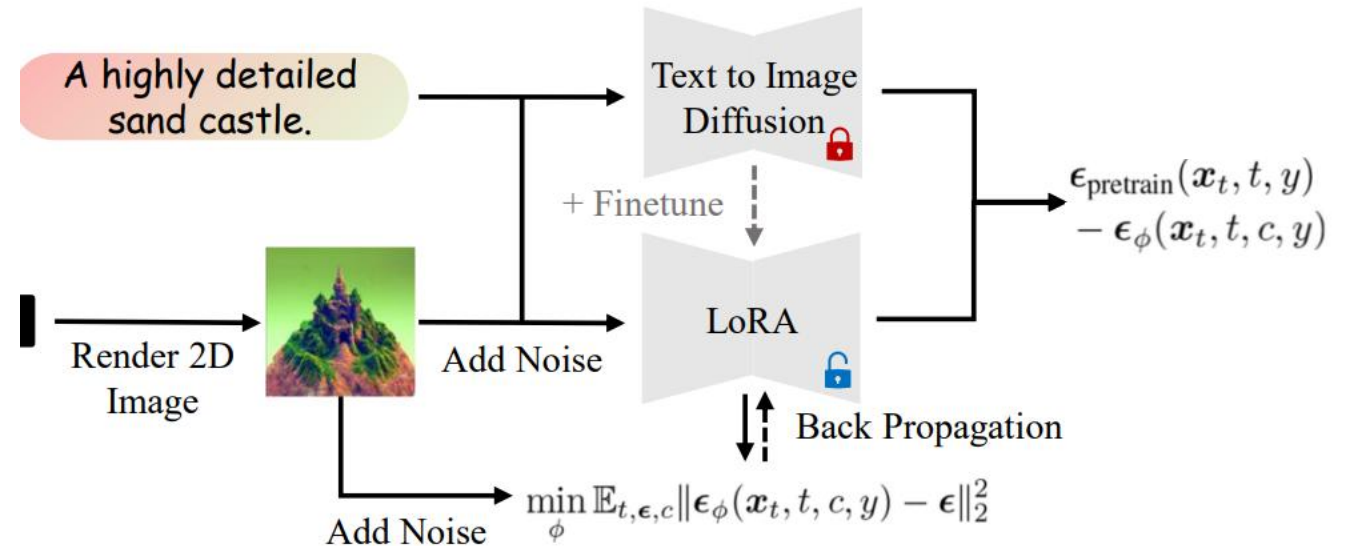Find a better alignment to rendered images distribution



"ProlificDreamer: High-Fidelity and Diverse Text-to-3D Generation with Variational Score Distillation", Zhengyi Wang et al., NeurIPS 23.

# Background

## Variational Score Distillation (VSD)

A better alignment achieves more accurate noise prediction

- Train a LoRA to predict noise

- Form a bi-level optimization

  - Finetune LoRA first, then predict noise to optimize 3D generation model iteratively



$$\epsilon_{\text{pretrain}}(\boldsymbol{x}_t, t, y) - \boldsymbol{\epsilon}_\phi(\boldsymbol{x}_t, t, c, y)$$

$$\min_\phi \mathbb{E}_{t,\boldsymbol{\epsilon},c}\|\boldsymbol{\epsilon}_\phi(\boldsymbol{x}_t, t, c, y) - \boldsymbol{\epsilon}\|_2^2$$

# Background

## Variational Score Distillation (VSD)

**Algorithm 1** Variational Score Distillation

**Input:** Number of particles $n$ ($\geq 1$). Large text-to-image diffusion model $\epsilon_{\text{pretrain}}$. Learning rate $\eta_1$ and $\eta_2$ for 3D structures and diffusion model parameters, respectively. A prompt $y$.

1: **initialize** $n$ 3D structures $\{\theta^{(i)}\}_{i=1}^{n}$, a noise prediction model $\epsilon_\phi$ parameterized by $\phi$.
2: **while** not converged **do**
3:   Randomly sample $\theta \sim \{\theta^{(i)}\}_{i=1}^{n}$ and a camera pose $c$.
4:   Render the 3D structure $\theta$ at pose $c$ to get a 2D image $\boldsymbol{x}_0 = \boldsymbol{g}(\theta, c)$.
5:   $\theta \leftarrow \theta - \eta_1 \mathbb{E}_{t,\epsilon,c} \left[ \omega(t) \left( \boldsymbol{\epsilon}_{\text{pretrain}}(\boldsymbol{x}_t, t, y^c) - \boldsymbol{\epsilon}_\phi(\boldsymbol{x}_t, t, c, y) \right) \frac{\partial \boldsymbol{g}(\theta, c)}{\partial \theta} \right]$
6:   $\phi \leftarrow \phi - \eta_2 \nabla_\phi \mathbb{E}_{t,\epsilon} || \boldsymbol{\epsilon}_\phi(\boldsymbol{x}_t, t, c, y) - \boldsymbol{\epsilon} ||_2^2.$
7: **end while**
8: **return**

# Outline

- Background

- Method

- Experiments

- Conclusion

# Method

Asynchronous Score Distillation (ASD)

Improve VSD which limited to:
- problematic optimization
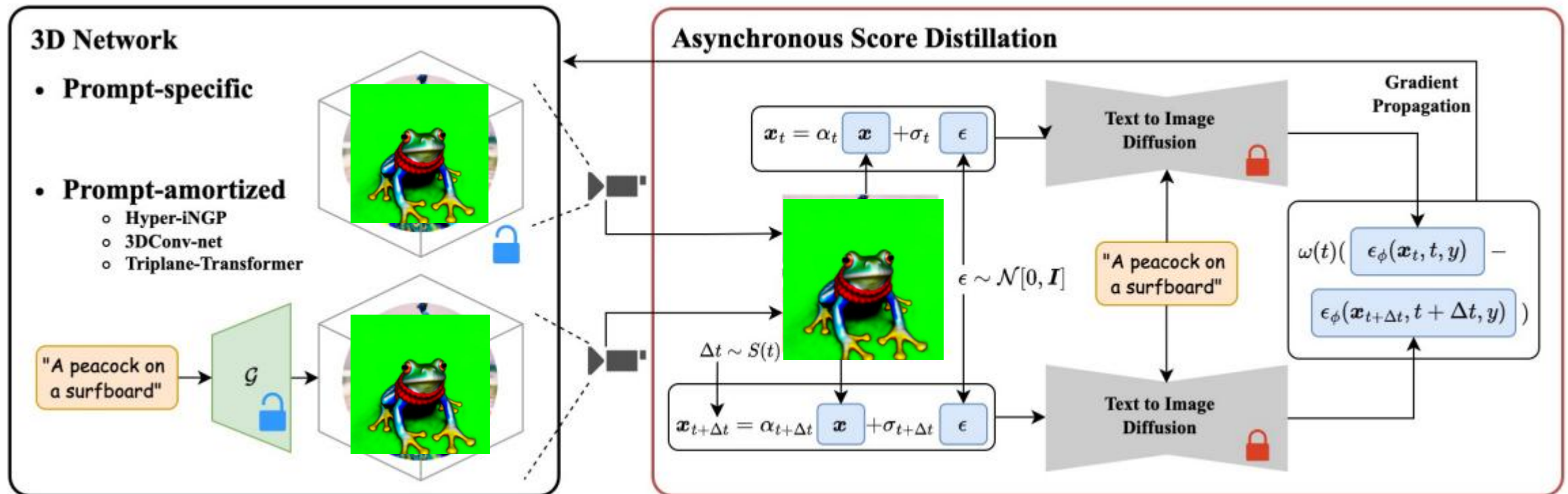- sacrificed comprehension ability to diverse prompts

Assumption:
better alignment with rendered image distribution, will lead to:
- more accurate noise prediction
- more effective loss gradient for optimization
- better 3D generation results
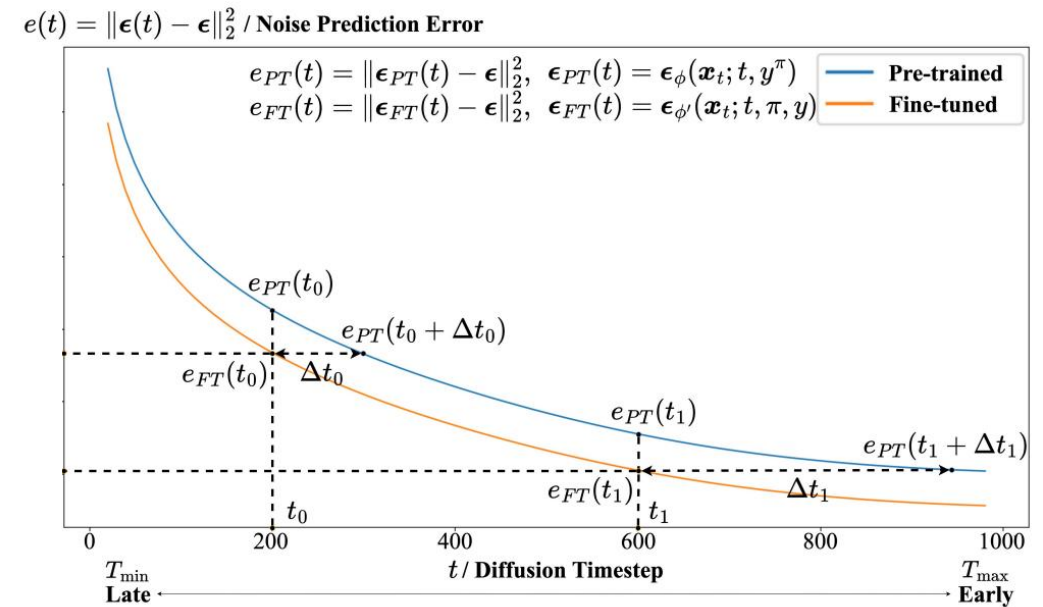
# Method

## Asynchronous Score Distillation (ASD)

- Predict noises on different timesteps, use discrepancy as gradient

# Method

## Observation

- Finetuned Diffusion model predict more accurate noise

- Noise prediction error will decrease as timestep increase, both on original UNet and finetuned UNet

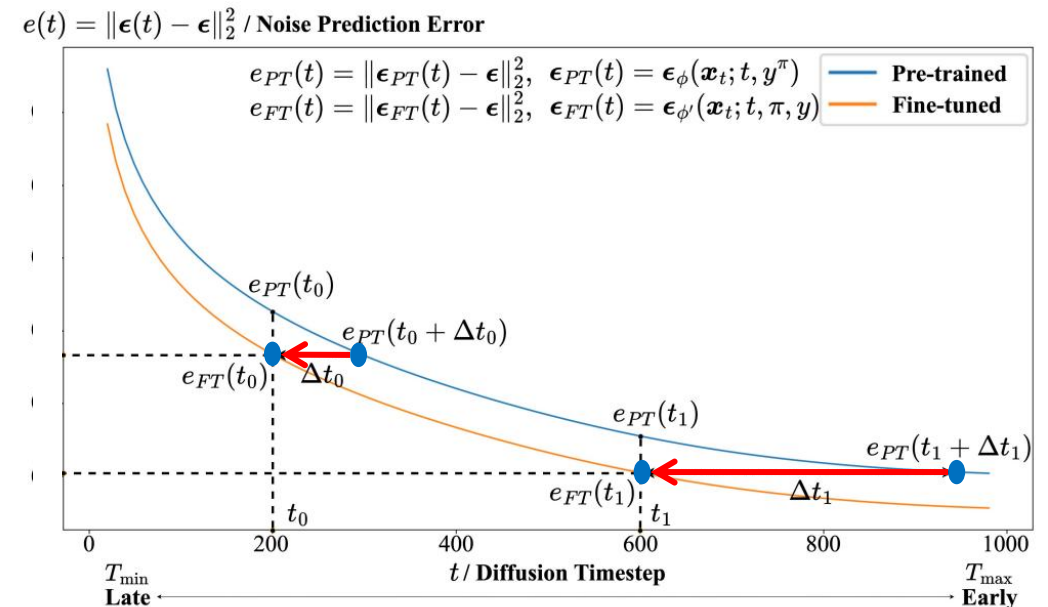- The speed of decrease becomes slower and slower as timestep increase



$e(t) = \|\boldsymbol{\epsilon}(t) - \boldsymbol{\epsilon}\|_2^2$ / Noise Prediction Error

$e_{PT}(t) = \|\boldsymbol{\epsilon}_{PT}(t) - \boldsymbol{\epsilon}\|_2^2, \quad \boldsymbol{\epsilon}_{PT}(t) = \boldsymbol{\epsilon}_\phi(\boldsymbol{x}_t; t, y^\pi)$

$e_{FT}(t) = \|\boldsymbol{\epsilon}_{FT}(t) - \boldsymbol{\epsilon}\|_2^2, \quad \boldsymbol{\epsilon}_{FT}(t) = \boldsymbol{\epsilon}_{\phi'}(\boldsymbol{x}_t; t, \pi, y)$

Pre-trained
Fine-tuned

# Method

Goal: Find a more accurate noise prediction on rendered images

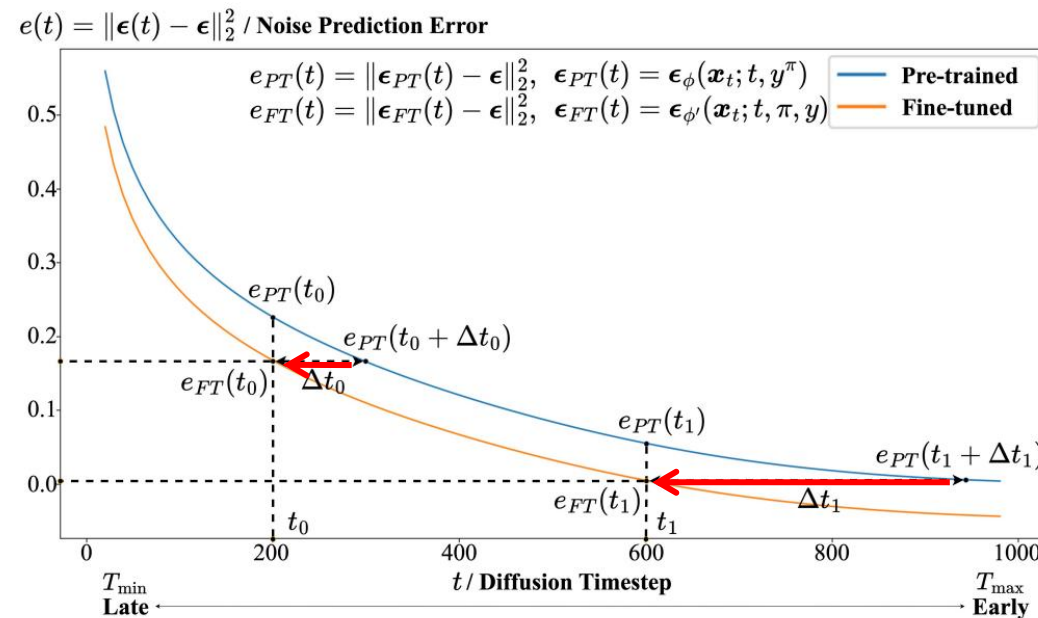Replace "ground truth noise" with noise prediction at <span style="color:red">larger timestep</span>

- More accurate prediction: $e_{PT}(t) \rightarrow e_{FT}(t)$

- Use $e_{PT}(t + \Delta t)$ approximate $e_{FT}(t)$



$e(t) = \|\boldsymbol{\epsilon}(t) - \boldsymbol{\epsilon}\|_2^2$ / **Noise Prediction Error**

$e_{PT}(t) = \|\boldsymbol{\epsilon}_{PT}(t) - \boldsymbol{\epsilon}\|_2^2, \quad \boldsymbol{\epsilon}_{PT}(t) = \boldsymbol{\epsilon}_\phi(\boldsymbol{x}_t; t, y^\pi)$
$e_{FT}(t) = \|\boldsymbol{\epsilon}_{FT}(t) - \boldsymbol{\epsilon}\|_2^2, \quad \boldsymbol{\epsilon}_{FT}(t) = \boldsymbol{\epsilon}_{\phi'}(\boldsymbol{x}_t; t, \pi, y)$

Pre-trained
Fine-tuned

# Method

A heuristic strategy to increase $\Delta t$ with larger timestep:

- if $t_0 < t_1$, then $\Delta t_0 < \Delta t_1$

- Sample a timestep shift $\Delta t \sim S(t) = \mathcal{U}\left[0, \eta\left(t - T_{\min}\right)\right]$

# Method

Comparison with VSD:

- No LoRA need to be trained, a single training objective

- Maintain most generative prior of pretrained diffusion model

**Algorithm 1** Variational Score Distillation

**Input:** Number of particles $n$ ($\geq 1$). Large text-to-image diffusion model $\epsilon_{\text{pretrain}}$. Learning rate $\eta_1$ and $\eta_2$ for 3D structures and diffusion model parameters, respectively. A prompt $y$.

1: **initialize** $n$ 3D structures $\{\theta^{(i)}\}_{i=1}^n$, a noise prediction model $\epsilon_\phi$ parameterized by $\phi$.
2: **while** not converged **do**
3:      Randomly sample $\theta \sim \{\theta^{(i)}\}_{i=1}^n$ and a camera pose $c$.
4:      Render the 3D structure $\theta$ at pose $c$ to get a 2D image $x_0 = g(\theta, c)$.
5:      $\theta \leftarrow \theta - \eta_1 \mathbb{E}_{t,\epsilon,c} \left[ \omega(t) \left( \epsilon_{\text{pretrain}}(x_t, t, y^c) - \epsilon_\phi(x_t, t, c, y) \right) \frac{\partial g(\theta, c)}{\partial \theta} \right]$
6:      $\phi \leftarrow \phi - \eta_2 \nabla_\phi \mathbb{E}_{t,\epsilon} \| \epsilon_\phi(x_t, t, c, y) - \epsilon \|_2^2$.
7: **end while**
8: **return**

**Algorithm 1** Asynchronous Score Distillation (ASD)

**Input:** 3D representation $\theta$; Text prompt $y$; Hyperparamter $\eta$; 2D diffusion prior $\epsilon_\phi$
**while** *not converged* **do**
     Sample a camera pose $\pi$
     Render an image $x = g(\theta, \pi)$
     Sample a timestep $t \sim \mathcal{U}[T_{\min}, T_{\max}]$, Gaussian noise $\epsilon \sim \mathcal{N}(0, I)$
     Sample a timestep shift $\Delta t \sim S(t) = \mathcal{U}[0, \eta(t - T_{\min})]$
     $x_t \leftarrow \alpha_t x + \sigma_t \epsilon$, $x_{t+\Delta t} \leftarrow \alpha_{t+\Delta t} x + \sigma_{t+\Delta t} \epsilon$
     Update $\theta$ with $\Delta \theta \leftarrow \omega(t) \left( \epsilon_\phi(x_t; t, y^\pi) - \epsilon_\phi(x_{t+\Delta t}; t + \Delta t, y^\pi) \right) \frac{\partial x}{\partial \theta}$
**end**
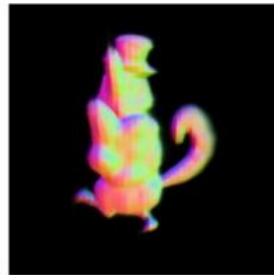
# Outline

- Background

- Method

- <span style="color:red">Experiments</span>
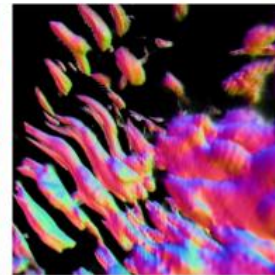
- Conclusion

# Experiments

Prompt amortized: optimize a general 3D generator to produce 3D objects given different prompts
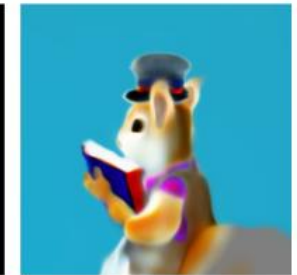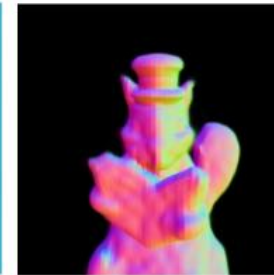
- More strict to prompt comprehensive ability
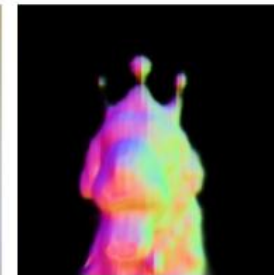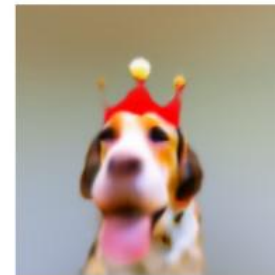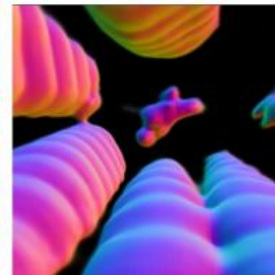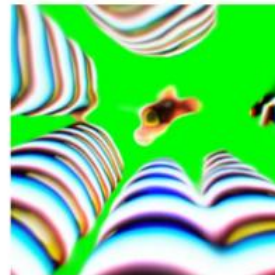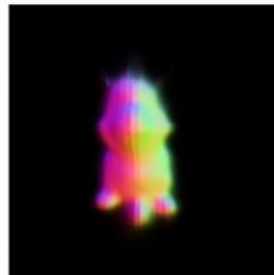


**Classifier Score Distillation (CSD)**     **Variational Score Distillation (VSD)**     **Asynchronous Score Distillation (ASD, ours)**

"A squirrel holding a book wearing a sweater wearing a tophat" in AT2520

"A DSLR photo of a cocker spaniel wearing a crown" in DF415

# Experiments

Scalability: train a 3D generator on 100k prompts
- VSD will be crashed
- ASD can generator more vivid results compared with CSD
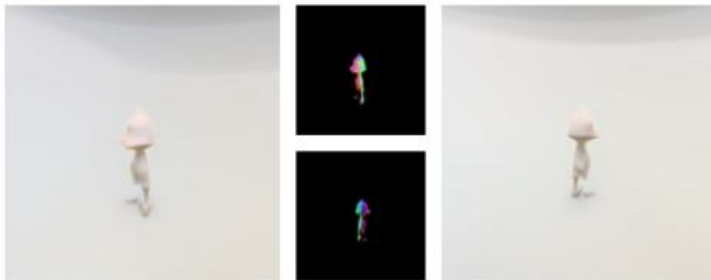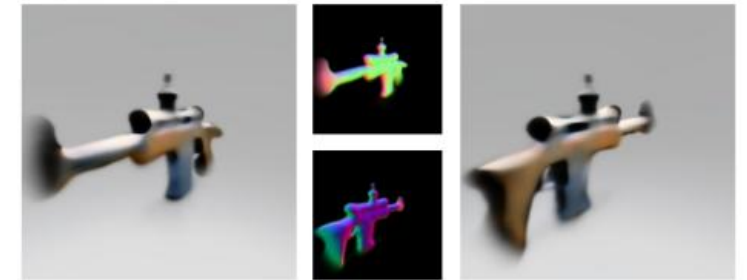


Classifier Score Distillation (CSD)  Variational Score Distillation (VSD)  Asynchronous Score Distillation (ASD, ours)

"An AR-15 rifle (M4 Carbine)"

# Experiments: Ablation

## Different setting of $\Delta t$

- No random sampling: not work
- Too big $\eta$: $\epsilon_\phi(\boldsymbol{x}_t; t, y^\pi) \approx \boldsymbol{\epsilon}$ , degrade to SDS, which is not suitable with CFG=7.5
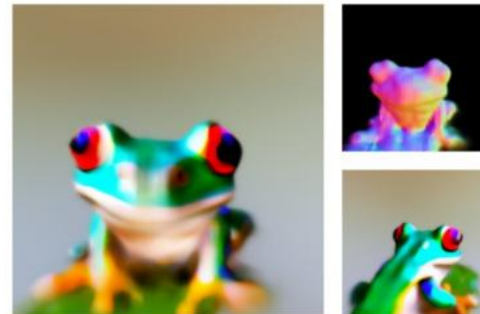


$$\Delta t = \eta(t - T_{\min})$$
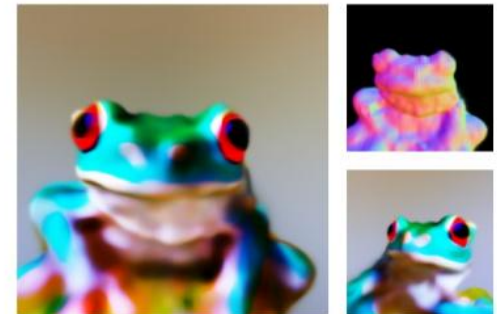$$\eta = 0.1$$

$$\eta = 0$$

$$\Delta t \sim \mathcal{U}[0, \eta(t - T_{\min})]$$
$$\eta = 0.1 \ \textbf{(Default)}$$
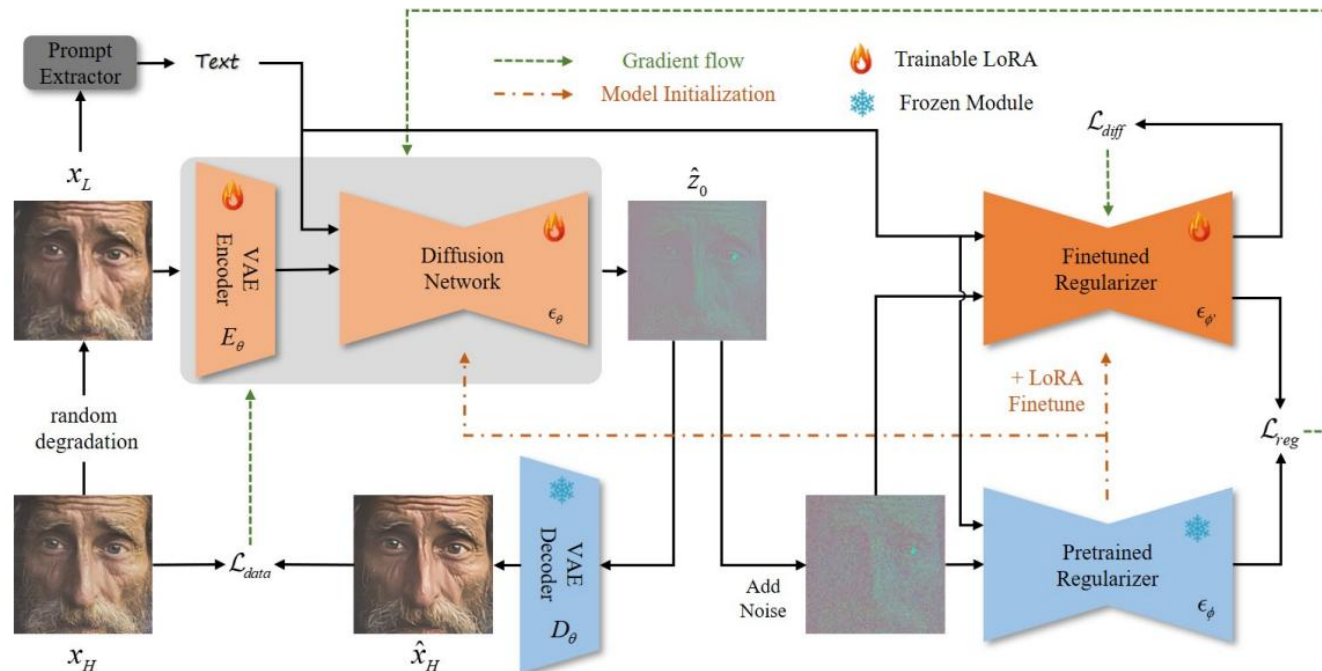
$$\eta = 0.2$$

"A DSLR photo of a red-eyed tree frog" in DF415

# Outline

- Background

- Method

- Experiments

- Conclusion

# Conclusion: Potential Extension

OSEDiff (arXiv 24.06): utilize VSD Loss on **one step** RealSR task
(Comes from the same team)



No noise input: stable, high fidelity

# Conclusion

- Regularize 3D generation by 2D pretrained Diffusion models effectively and efficiently

- Improve comprehension ability with scalable prompts

- Potential applications to more vision tasks

# Thanks for listening!