

# Text-to-3D using 2D Diffusion

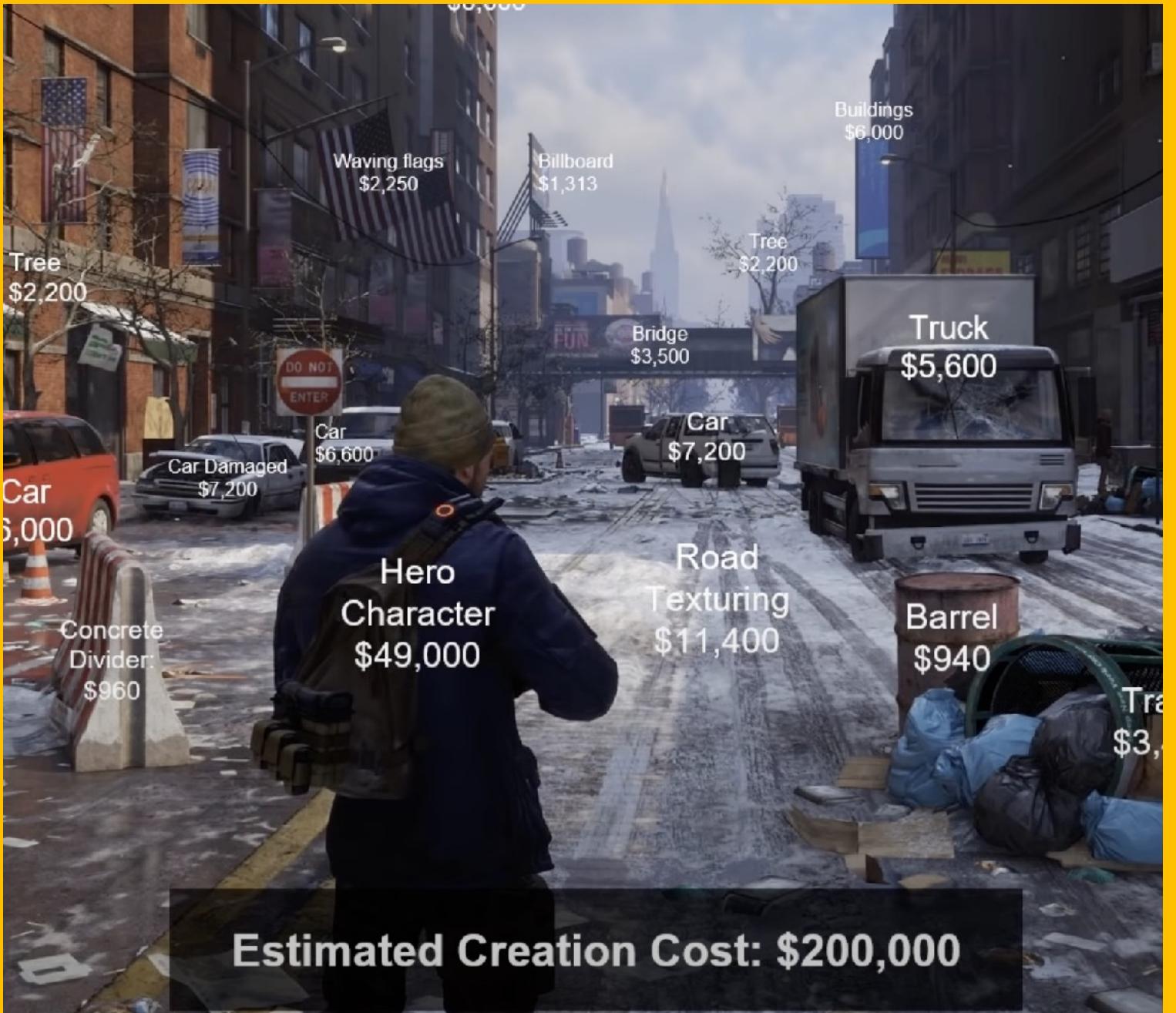


**Kirill Struminsky**

Research Fellow @ HSE Centre for Deep Learning and Bayesian Methods



# Motivation

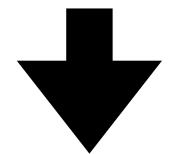
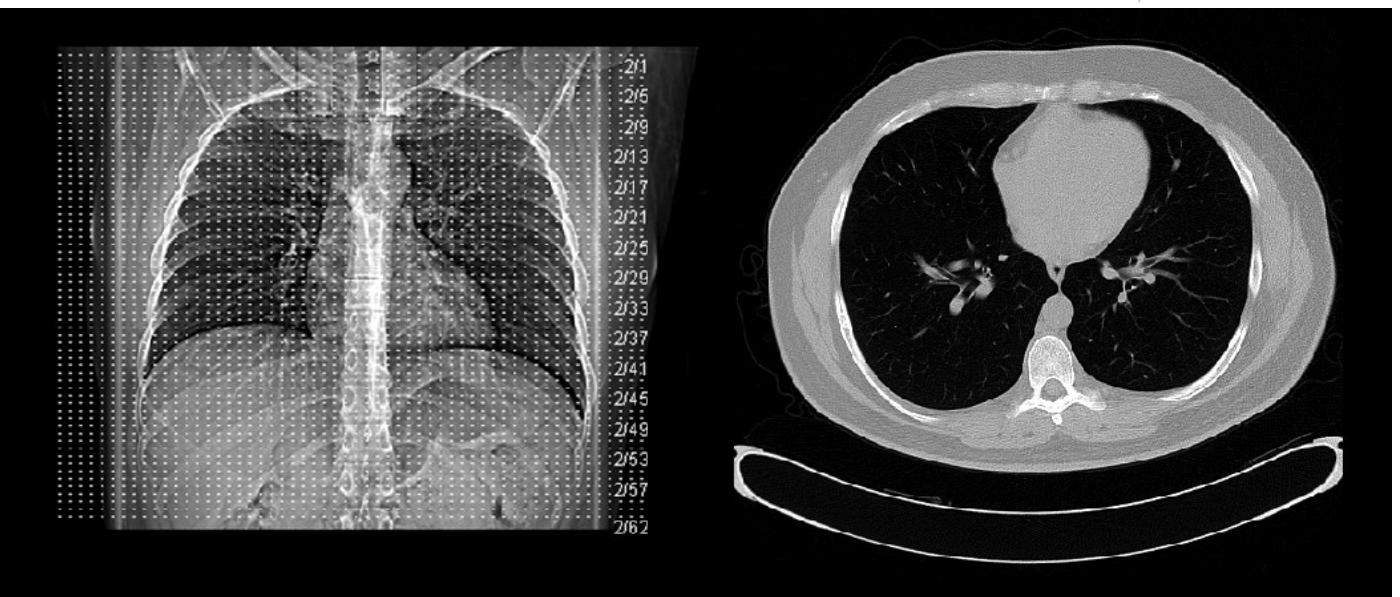


- 3D assets are expensive
- Demand goes beyond game industry

Image taken from “The Next Leap: How A.I. will change the 3D industry”  
talk by Andrew Price

# Motivation

- Inverse problems
  - Observe
  - Find
- Generative model
  - Find
  - Such that are plausible samples
- Can we use as an objective?



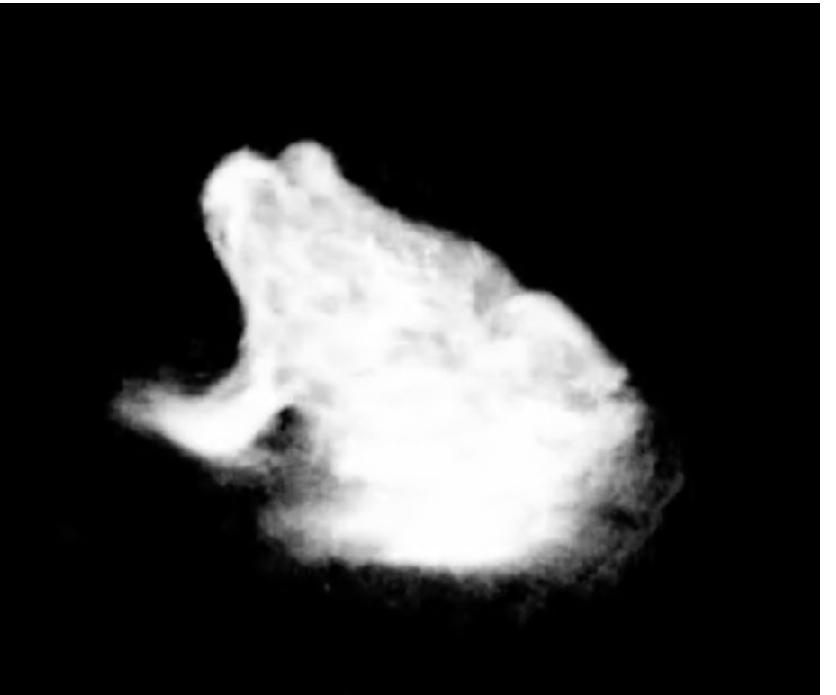
# Setup

- Pre-trained text-to-2D diffusion model
- Input: text prompt
  - A DSLR photo of a squirrel wearing a kimono reading a book*
- Output: 3D model



# Parameterising & Rendering 3D Models

- Neural Radiance Fields
- Volumetric density
- Radiance field
- Often, and are neural networks



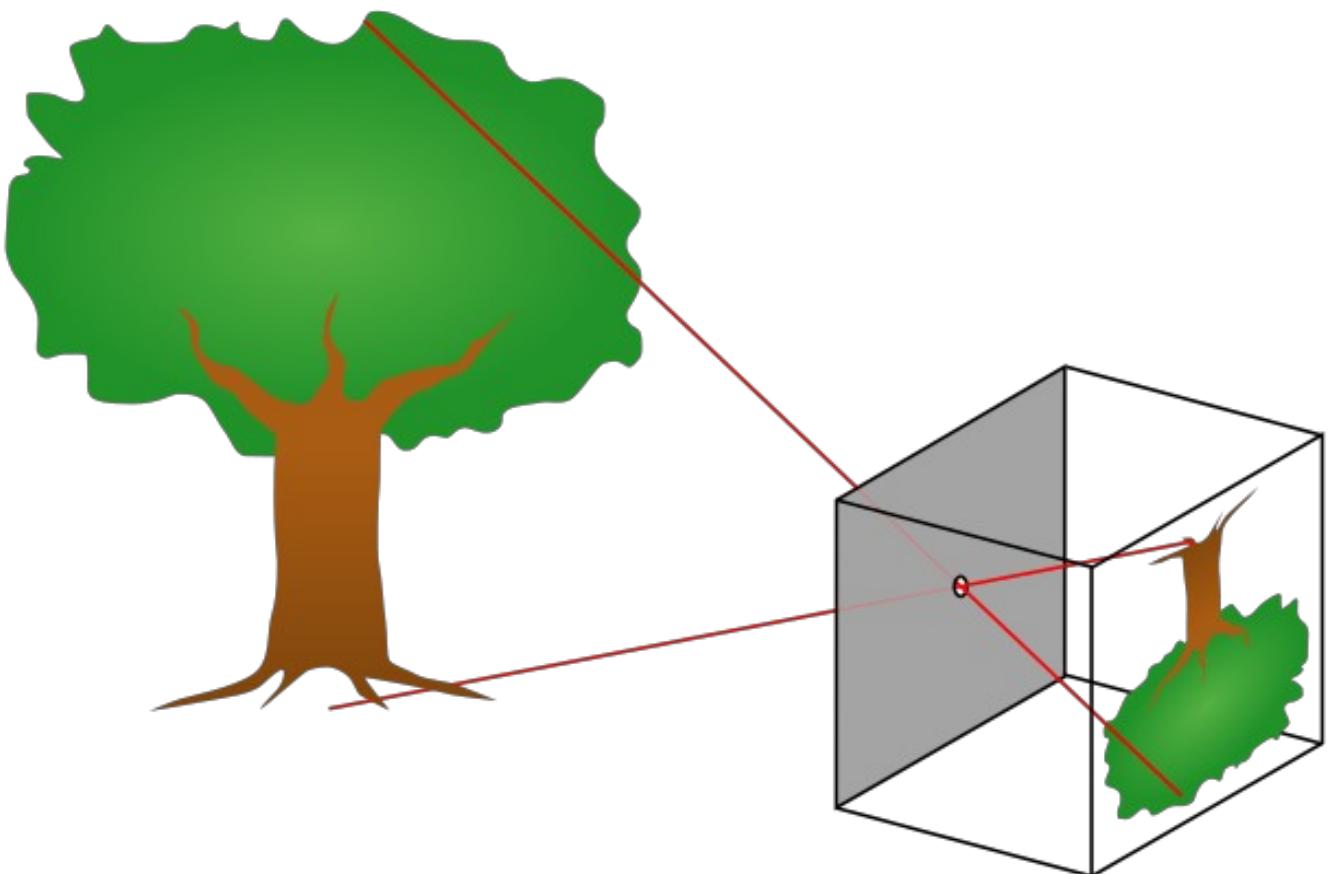
# Parameterising & Rendering 3D Models

- Fields and parameterise the scene
- Render image pixel by pixel

$$r = (r_o, r_d) t_i$$

$$\sigma_i = \sigma(r_o + t_i r_d)$$
$$C_i = c(r_o + t_i r_d)$$

- Pixel color is



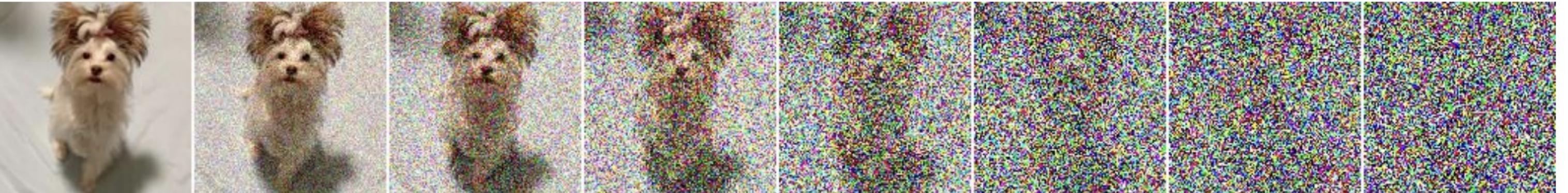
Course Iter.: 1  
Eps. time: 00:00

Course Iter.: 1  
Eps. time: 00:00

Course Iter.: 1  
Eps. time: 00:00

# Diffusion Models Recap

Mapping data to noise



$$x_t = \alpha_t x_0 + \sigma_t \varepsilon$$

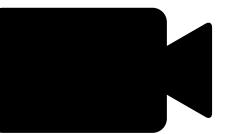
Trained reverse map

$$\hat{x}_\phi(x_t, t) = \frac{(x_t - \sigma_t \varepsilon_\phi(x_t, t))}{\alpha_t}$$

$$\nabla_{x_t} \log p_t(x_t) = s(x_t) \approx -\frac{\varepsilon_\phi(x_t, t)}{\sigma_t}$$

# DreamFusion: Text-to-3D using 2D Diffusion

- Consider rendered image
  - Scene parameterisation
  - View angle, camera parameters, etc.
- How do we find ?
- Optimize ?
  - Training optimises
  - Can we optimise for
- Minimize



# Score Distillation Sampling

- Loss
  - Image and
  - Compute gradient
- 
- Omitting score Jacobian helps



# Score Distillation Sampling Loss

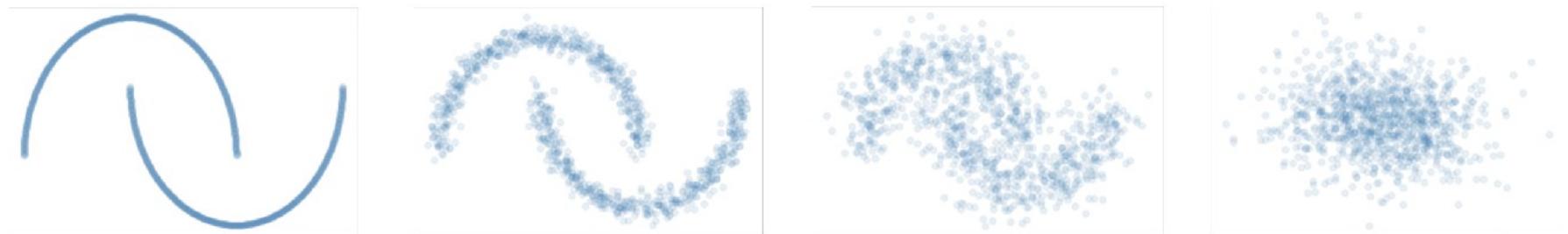
- Let
- Can we interpret ?

# Score Distillation Sampling Loss

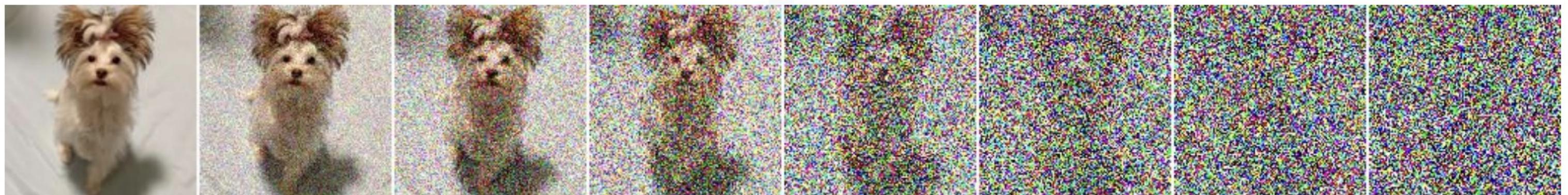
- Let
  - Can we interpret ?
  - Recall and rewrite
- 
- Finally

# Score Distillation Sampling Loss

- In sum,



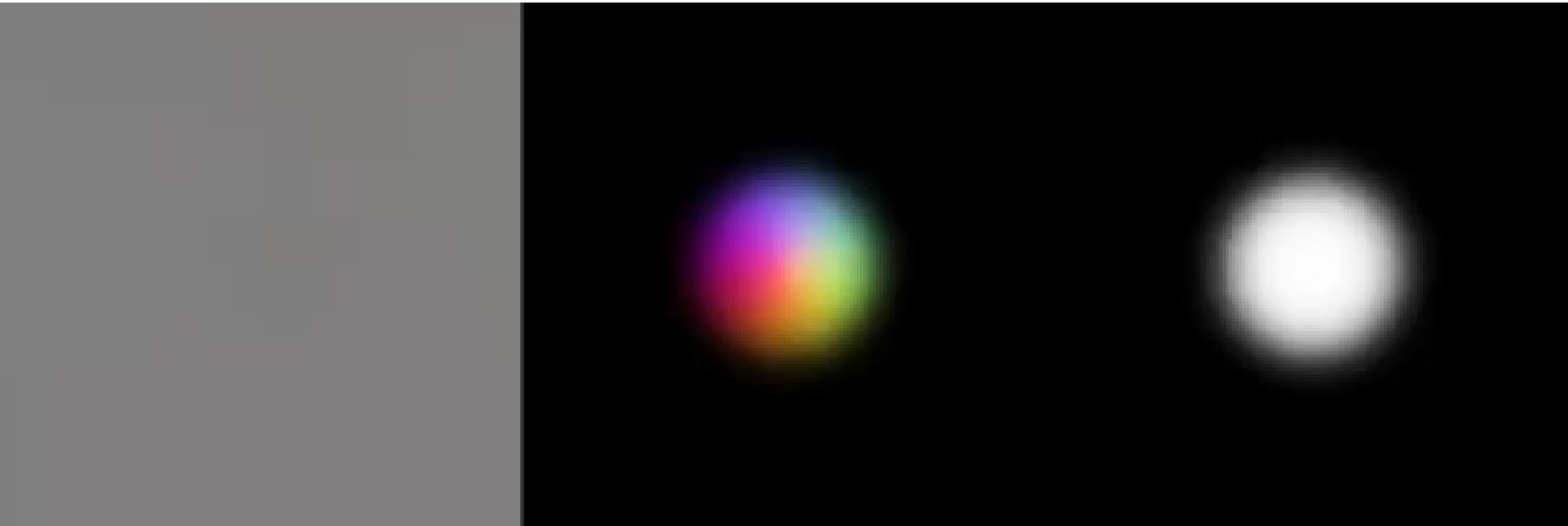
$$L_{Diff}$$



$$L_{SDS}$$

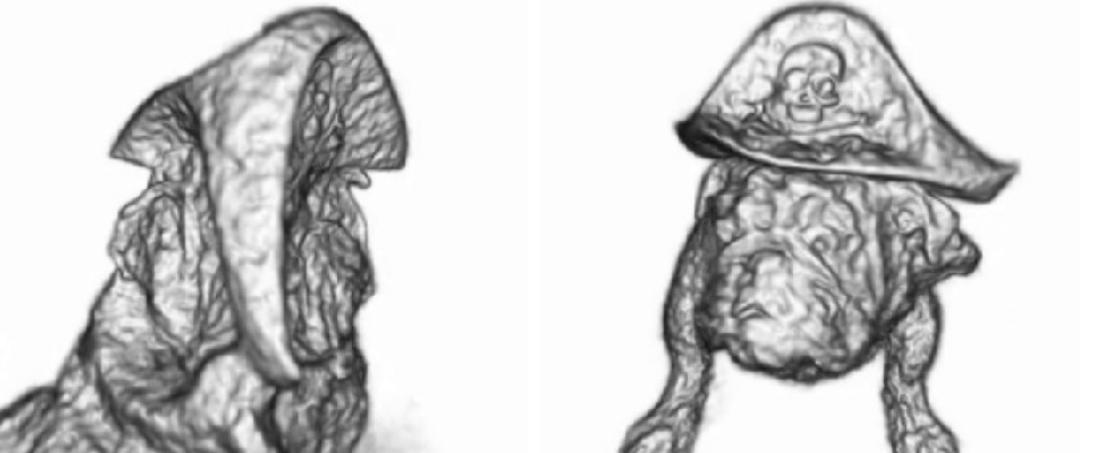
# Training Dreamfusion

- Choose prompt (omitted for clarity)
- Diffusion model (Imagen)
- Training loop:
  - Sample camera position
  - Render image
  - Sample and
  - Update with



# Additional Heuristics (I)

- Janus problem
- Remedy: specify camera angle in prompts

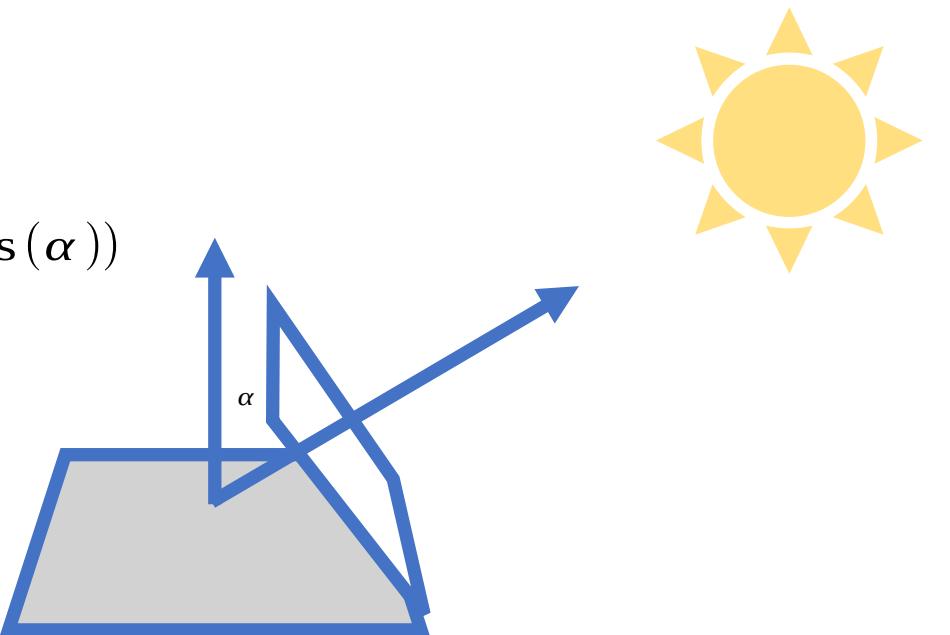


# Additional Heuristics (II)

- Poor geometry
- Remedy (1/2): simulate shadows



$$I \cdot C \max(0, \cos(\alpha))$$



- Remedy (2/2): process colourless images



# Problems to Fix

- Poor model quality
- Seemingly ad-hoc objective
- Low sample diversity
- Janus problem

# Adapting to Latent Diffusion

- Imagen is not public
- Naive adaptation works fine
  - Render image
  - Encode
  - Compute SDS loss in latent space



*a blue poison-dart frog  
sitting on a water lily*



*neuschwanstein castle, aerial view*



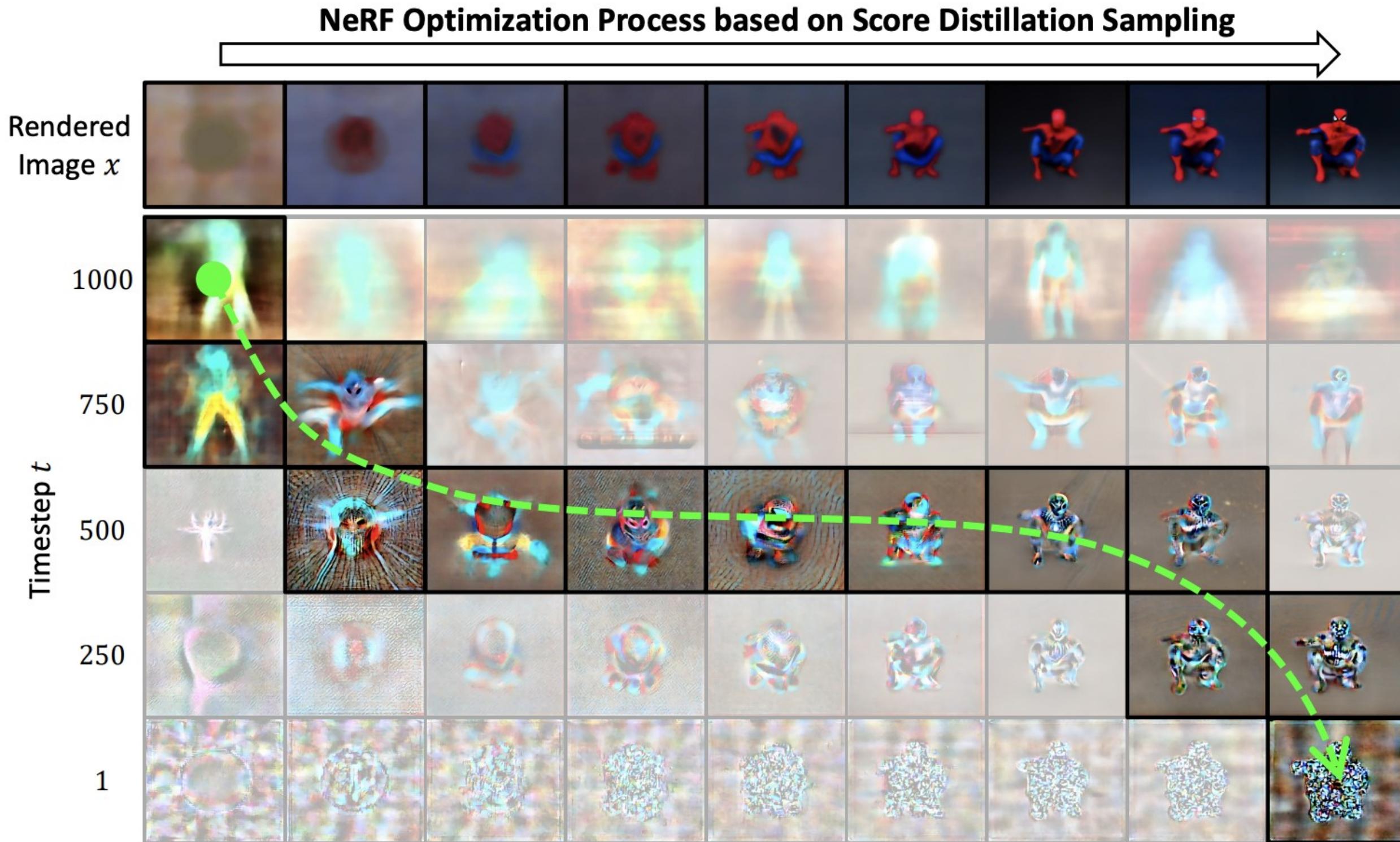
(1) 64 rendering

(2) + 512 rendering

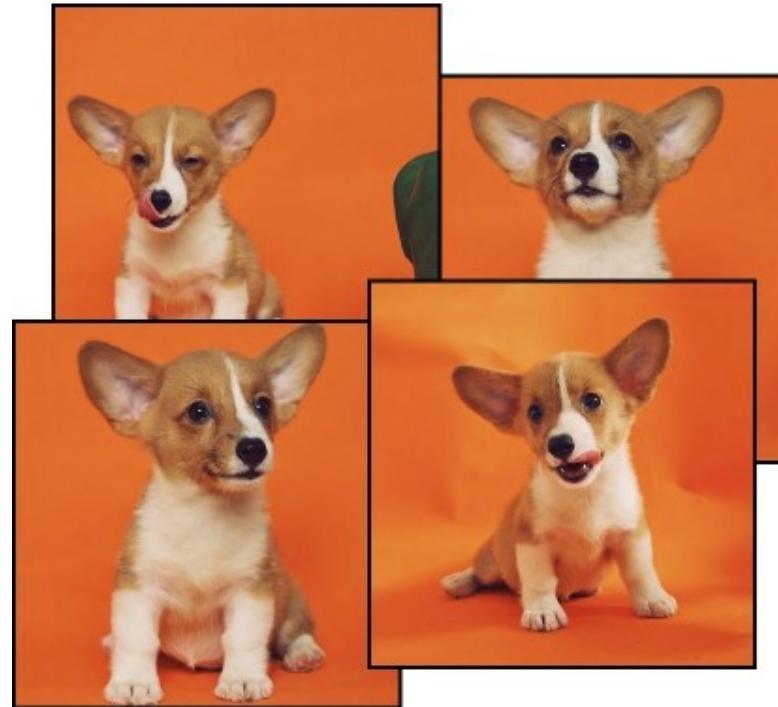
(3) + Annealed  $t$

(4) + VSD

# Time Annealing



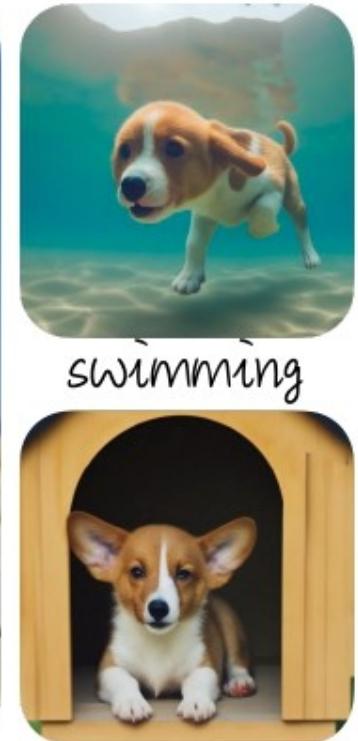
# Tackling Low Sample Diversity



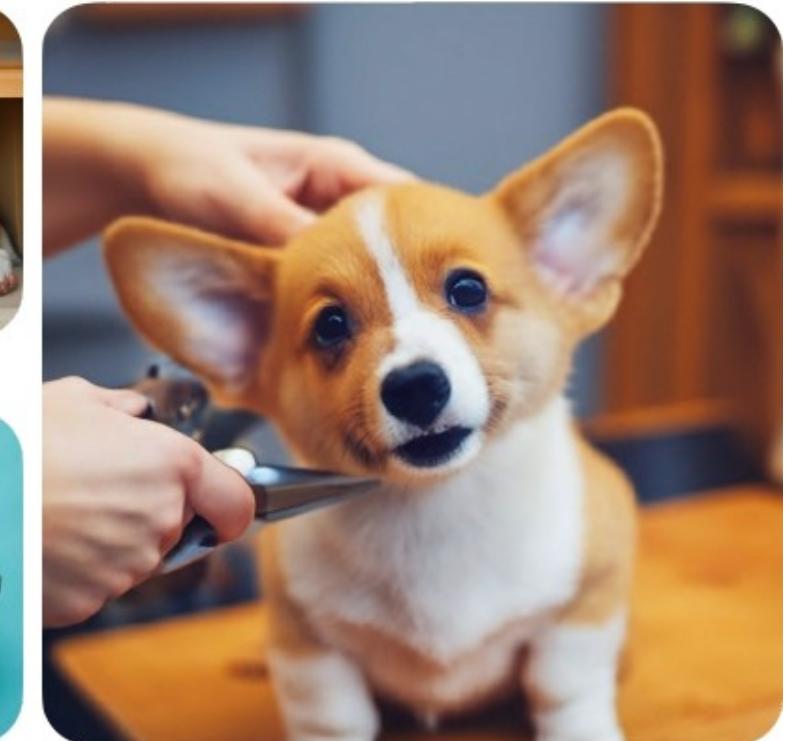
Input images



in the Acropolis

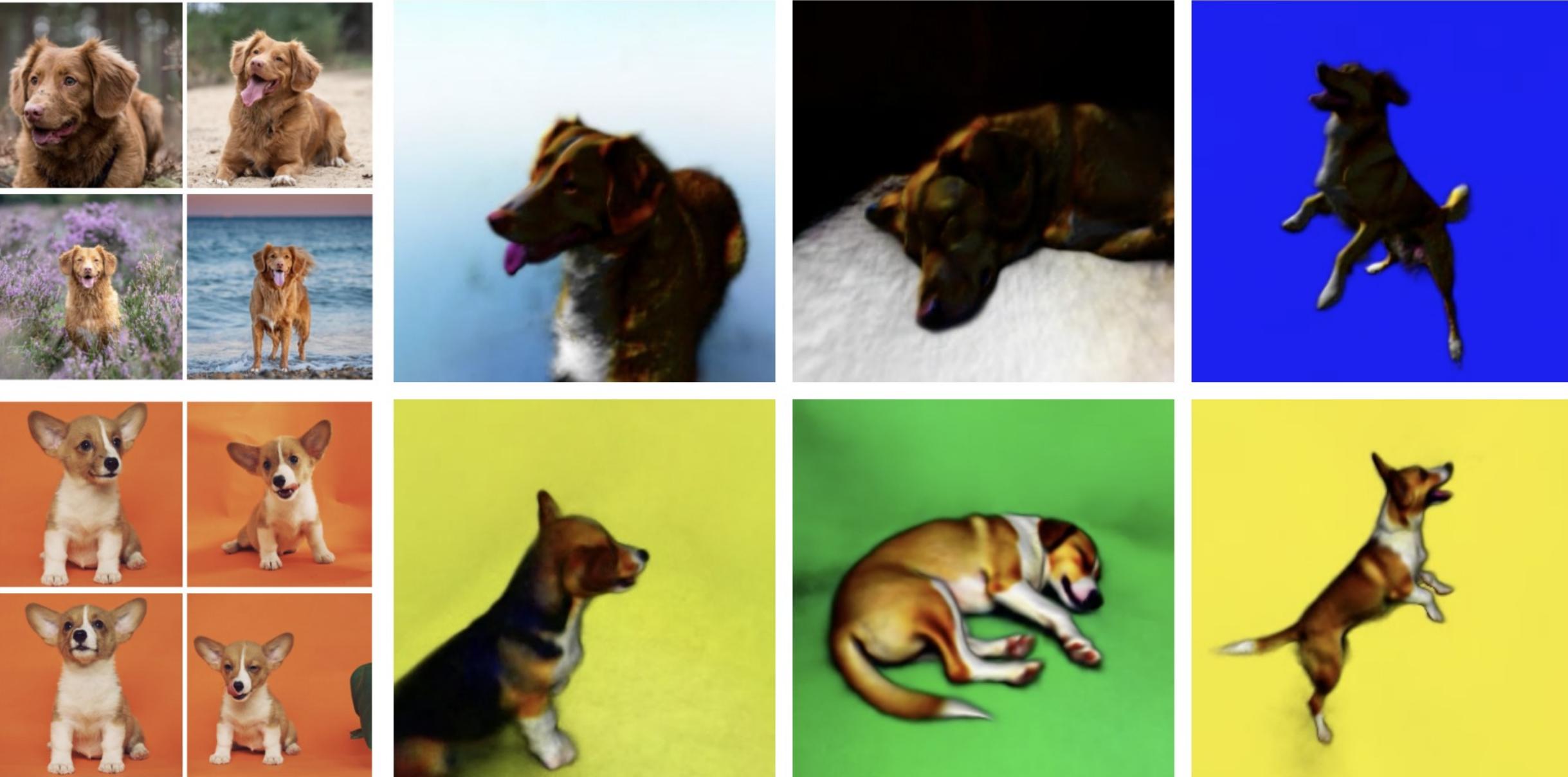


in a bucket

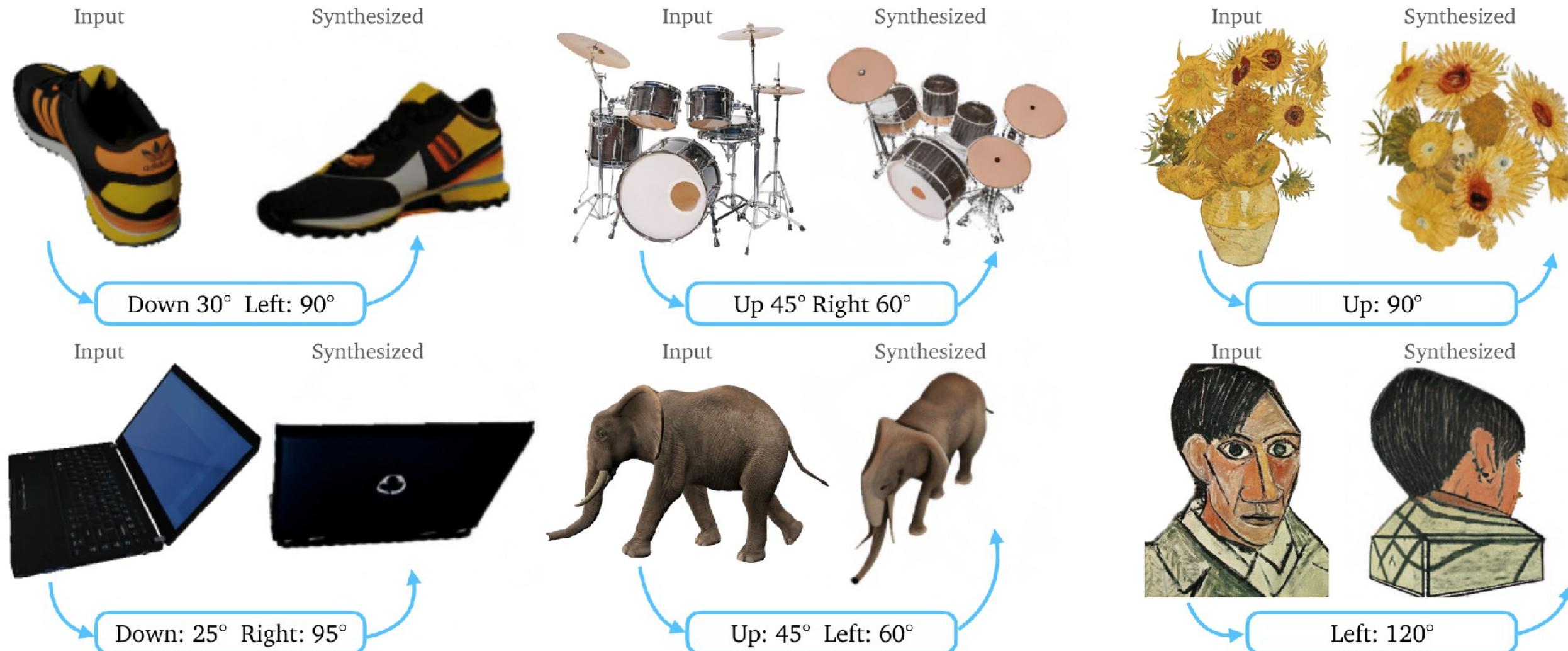


getting a haircut

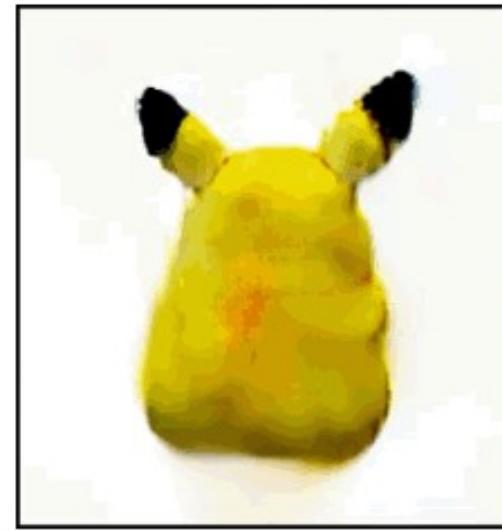
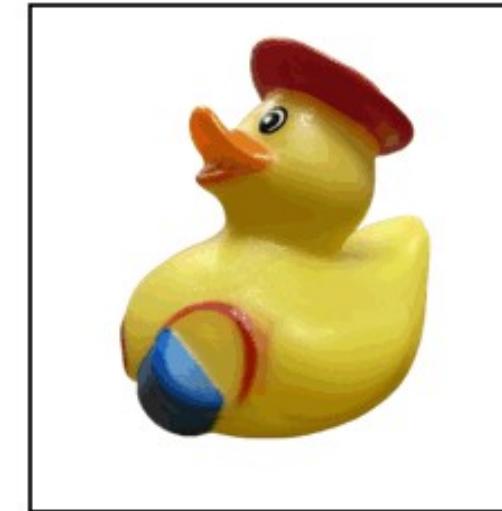
# Dreambooth 3D



# Zero-123: viewangle control



# Zero-123: examples



# Magic123



# Key takeaways

- Text-to-2D models can solve to text-to-3D
- Generation with SDS loss
- Model fine-tuning can address some of SDS issues