

A practitioner's guide to diffusion models

CS180/280A, 11/4

Aleksander Holynski



A raccoon wearing a tuxedo and top hat
with the moon in the background.





1000

What are diffusion models?

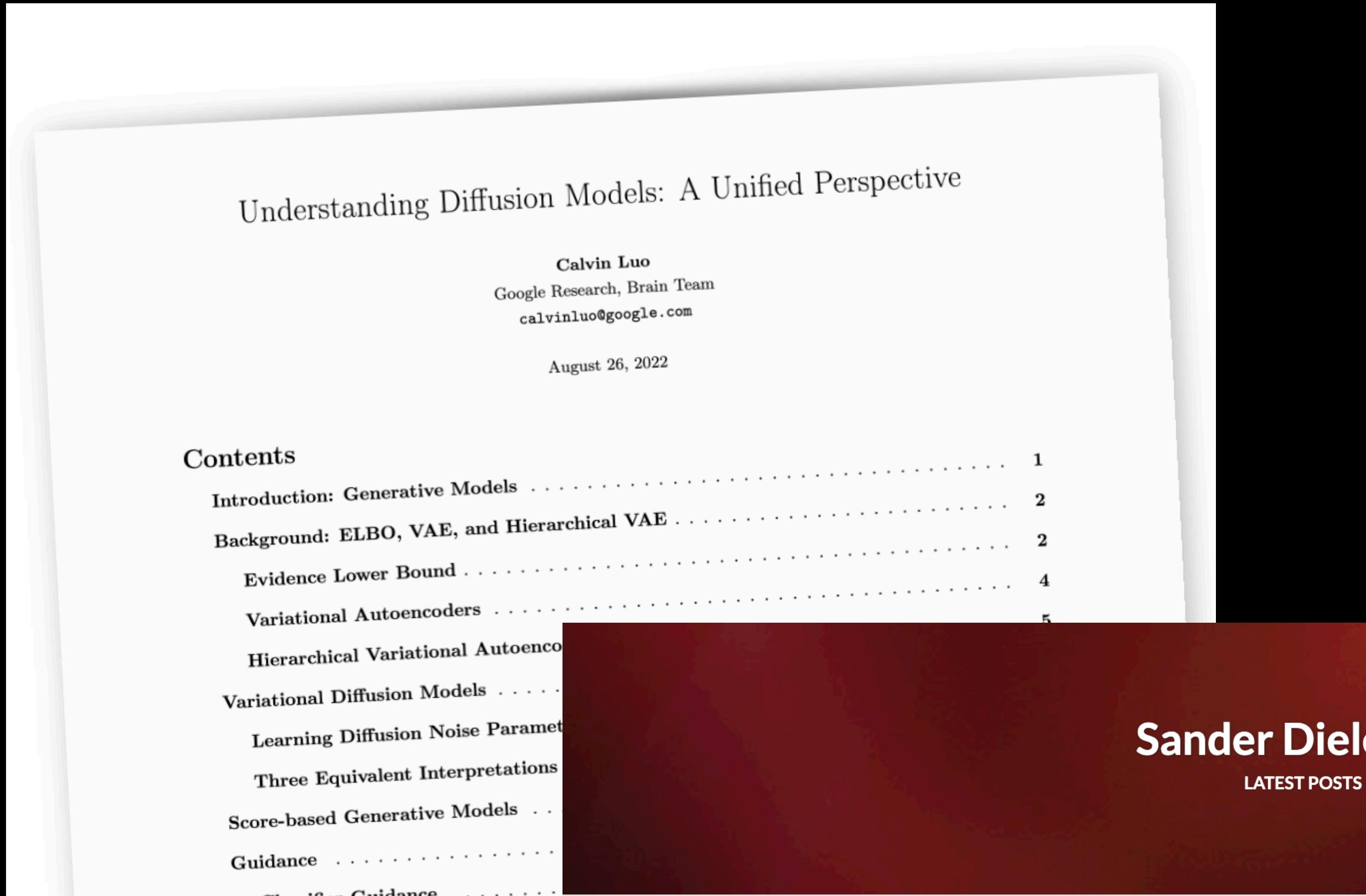
How do I train one?

How do I use them?

What's covered

- Intuition on diffusion models—how do they work?
- How to train one—(data, architecture, training processes)
- How to use them—(sampling, guidance, distillation, zero-shot editing)
- What can they be used for? (editing, image priors, controllability)

What's not covered: *most of the math*

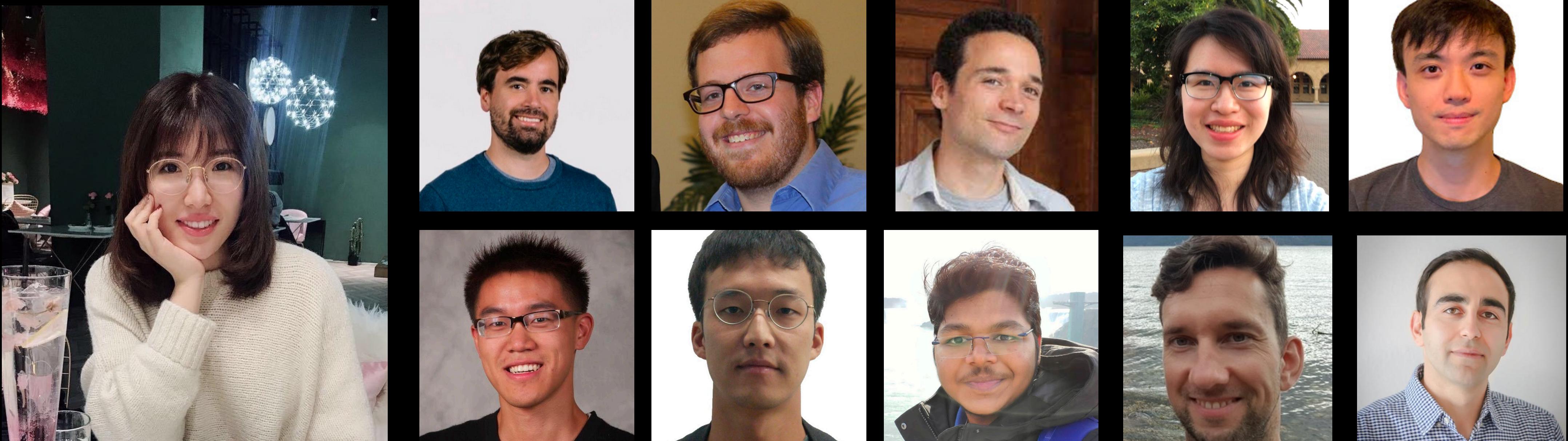


The image is a thumbnail for a blog post by Sander Dieleman. The title is "The paradox of diffusion distillation". Below the title is a short description: "Diffusion models split up the difficult task of generating data from a high-dimensional distribution into many denoising tasks, each of which is much easier. We train them to solve just one of these tasks at a time. To sample, we make many predictions in sequence. This iterative refinement is where their power comes from." At the bottom left is the date "FEBRUARY 28, 2024 BY SANDER DIELEMAN" and at the bottom right is "© READING TIME ~46 MINUTES".

The image shows the cover page of a tutorial titled "Tutorial on Diffusion Models for Imaging and Vision" by Stanley Chan. The date is March 28, 2024. The abstract states: "The astonishing growth of generative tools in recent years has empowered many exciting applications in text-to-image generation and text-to-video generation. The underlying principle behind these generative tools is the concept of *diffusion*, a particular sampling mechanism that has overcome some shortcomings that were deemed difficult in the previous approaches. The goal of this tutorial is to discuss the essential ideas underlying the diffusion models. The target audience of this tutorial includes undergraduate and graduate students who are interested in doing research on diffusion models or applying these models to solve other problems." The contents table includes sections on VAE and DDPM.

Contents	
1 The Basics: Variational Auto-Encoder (VAE)	2
1.1 VAE Setting	2
1.2 Evidence Lower Bound	4
1.3 Training VAE	7
1.4 Loss Function	9
1.5 Inference with VAE	9
2 Denoising Diffusion Probabilistic Model (DDPM)	10
2.1 DDPM Setting	10
2.2 Evidence Lower Bound	13
2.3 Training DDPM	14
2.4 Loss Function	15
2.5 Inference with DDPM	18
2.6 DDPM vs VAE	20
2.7 DDPM vs Score-based GAN	23
2.8 DDPM vs DDIM	25
2.9 DDPM vs DDPM	27
3 Summary and Future Work	30
3.1 Summary	30
3.2 Future Work	33

Acknowledgements



Some slides were borrowed from
Denoising Diffusion-based Generative Modeling CVPR2022 tutorial
(<https://cvpr2022-tutorial-diffusion-models.github.io/>)

Part 1:

What are diffusion models?

Denoising Diffusion Models

Emerging as powerful generative models, outperforming GANs



["Diffusion Models Beat GANs on Image Synthesis"](#)
Dhariwal & Nichol, OpenAI, 2021



["Cascaded Diffusion Models for High Fidelity Image Generation"](#)
Ho et al., Google, 2021

Deep generative models

Learning to generate data

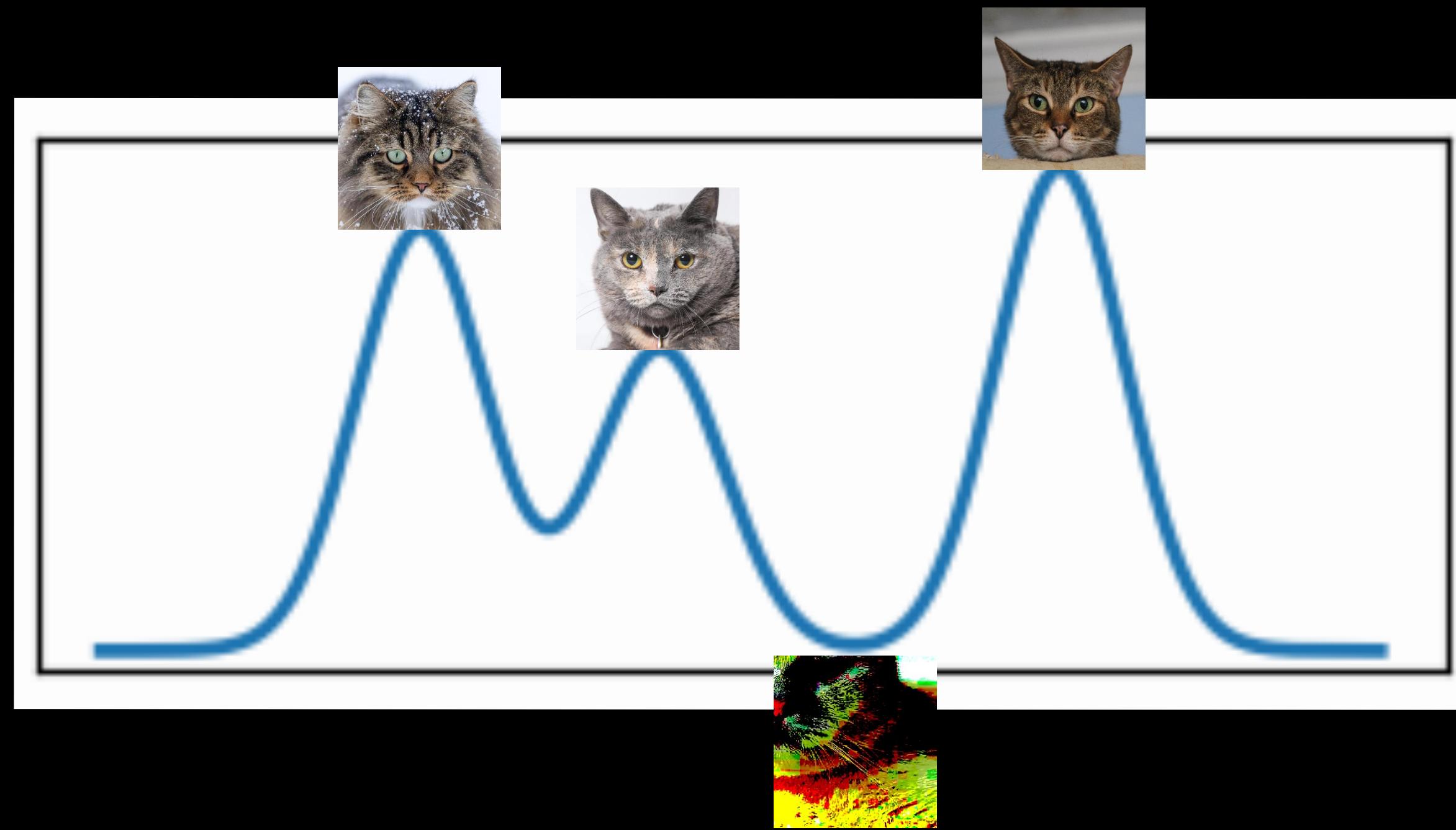


Deep generative models

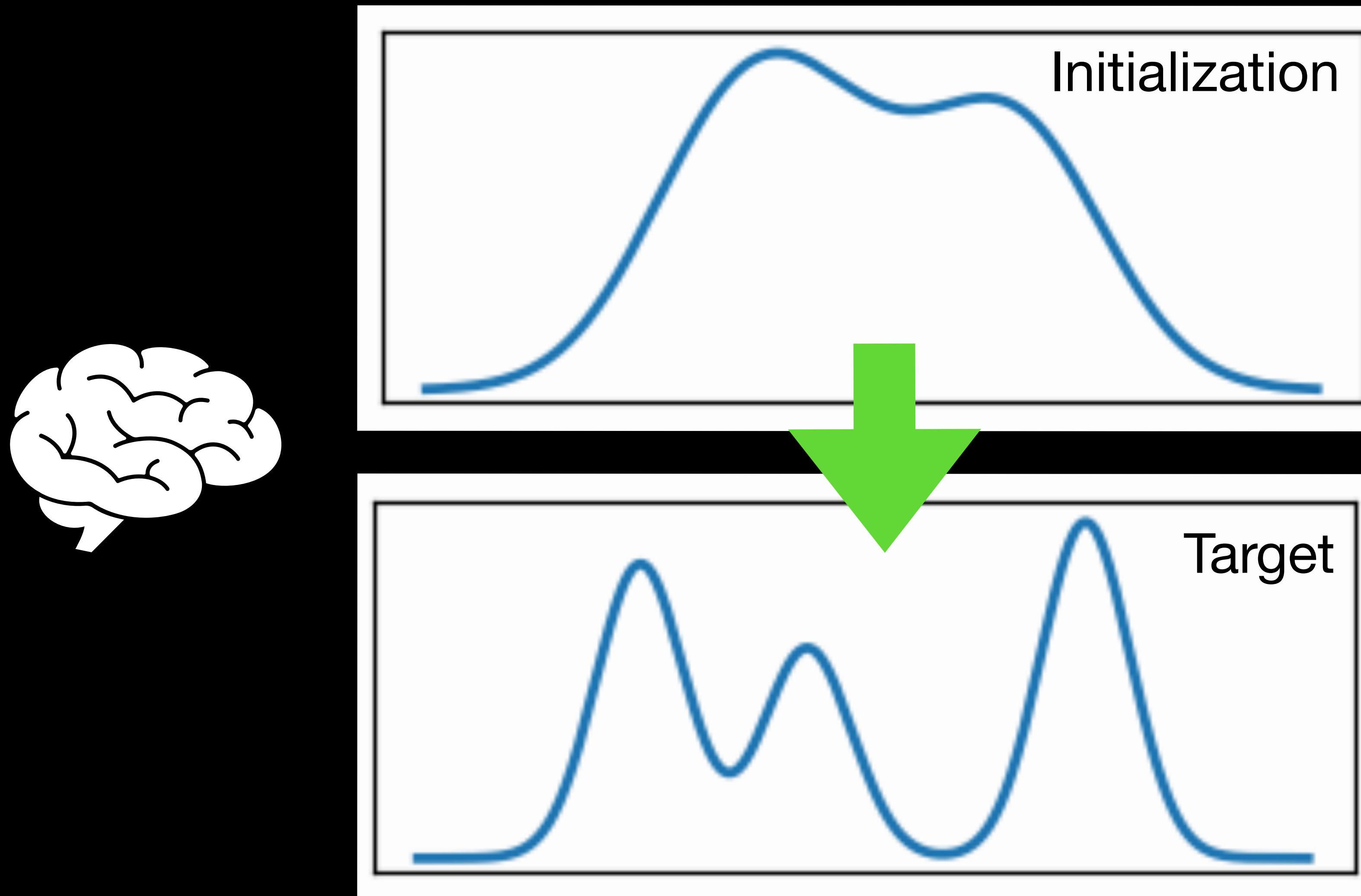
Learning to generate data



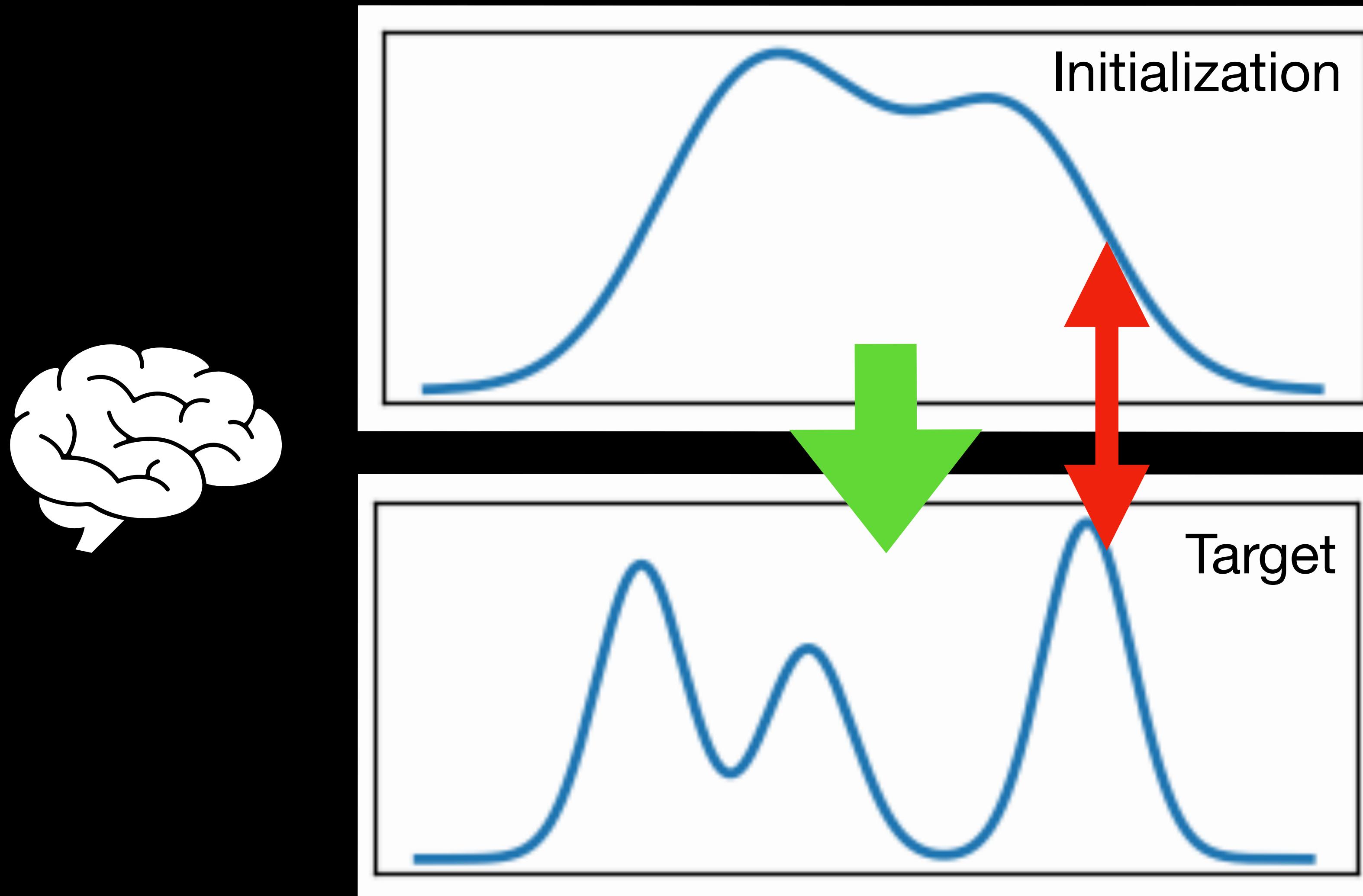
Neural Network



Training a generative model



Training a generative model



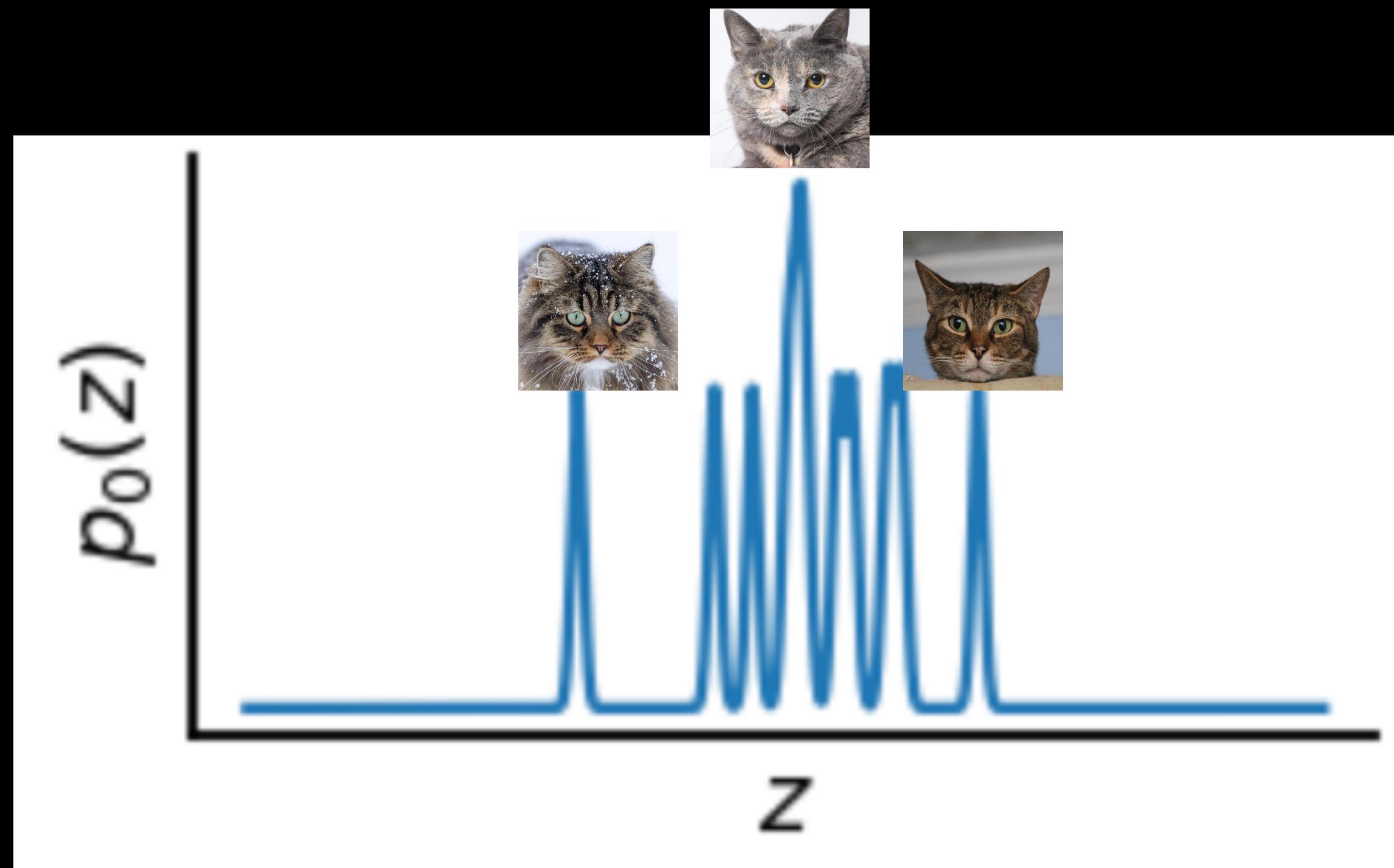
Training a generative model

“Give me image #57838905”



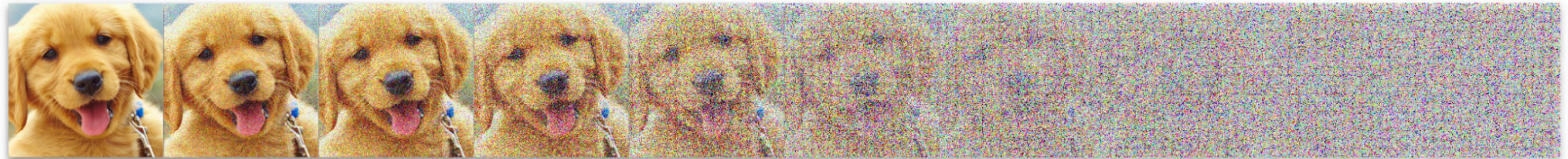
Compute a loss

Training a generative model

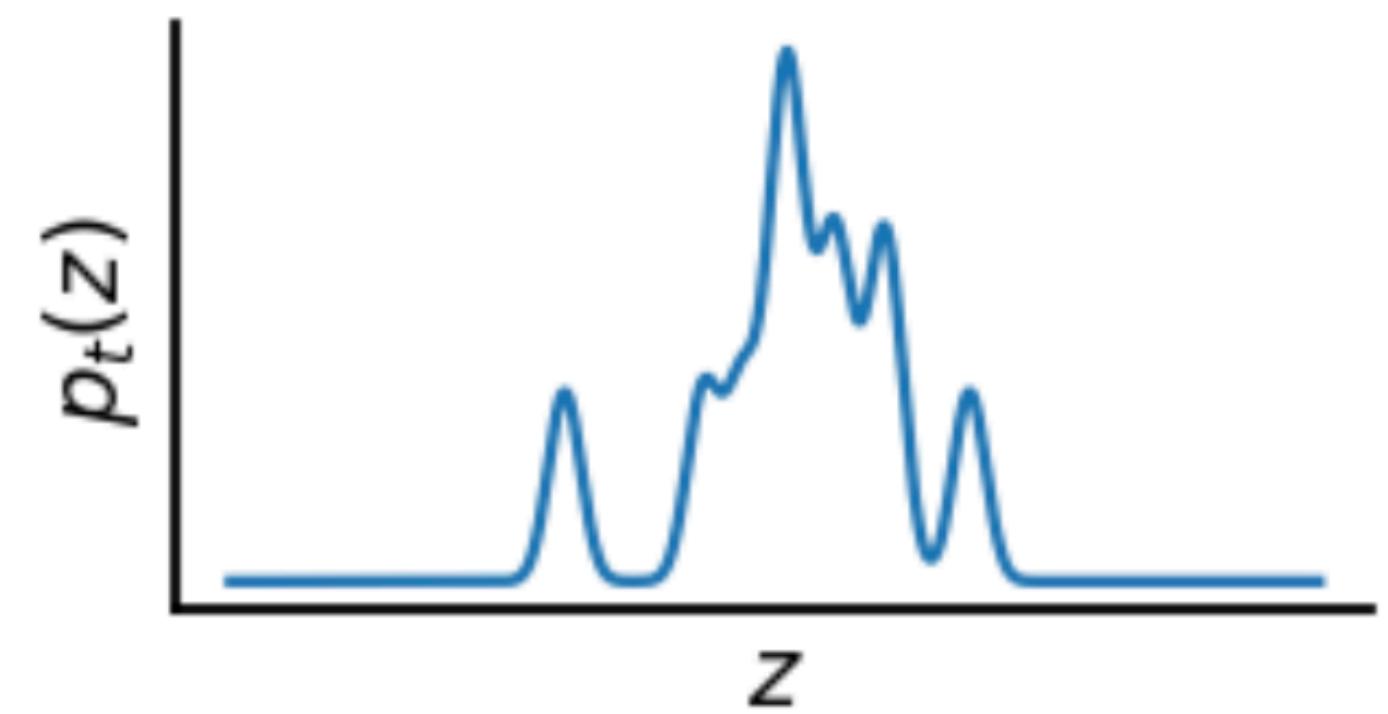
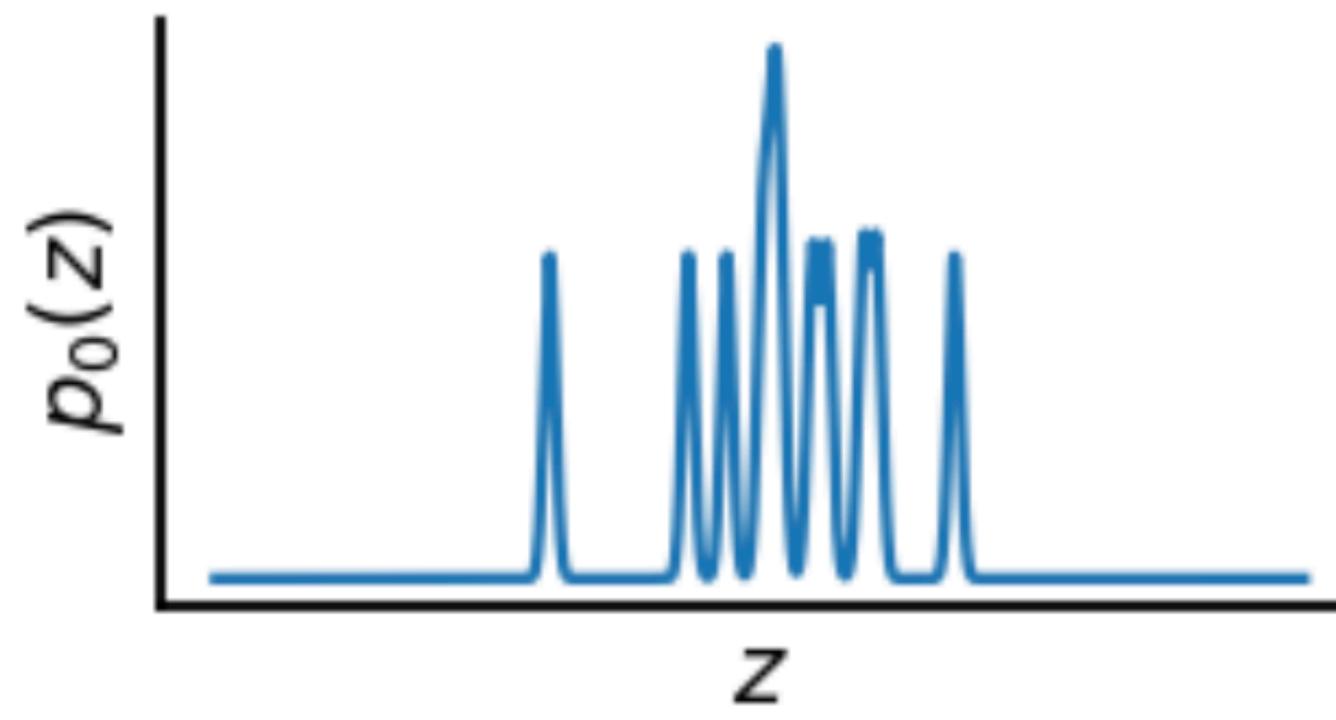


Part 2:

Intuition on diffusion



$p_0(z_0), \dots, p_t(z_t), \dots, p_T(z_T)$



Creating is hard!



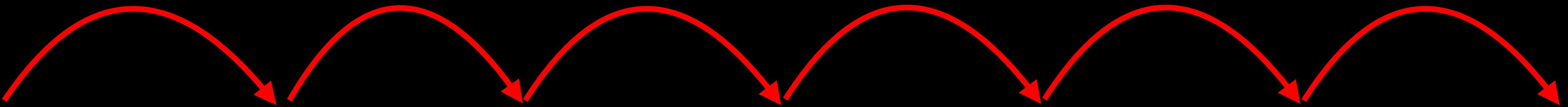
Creating is hard!

Destroying is easy!

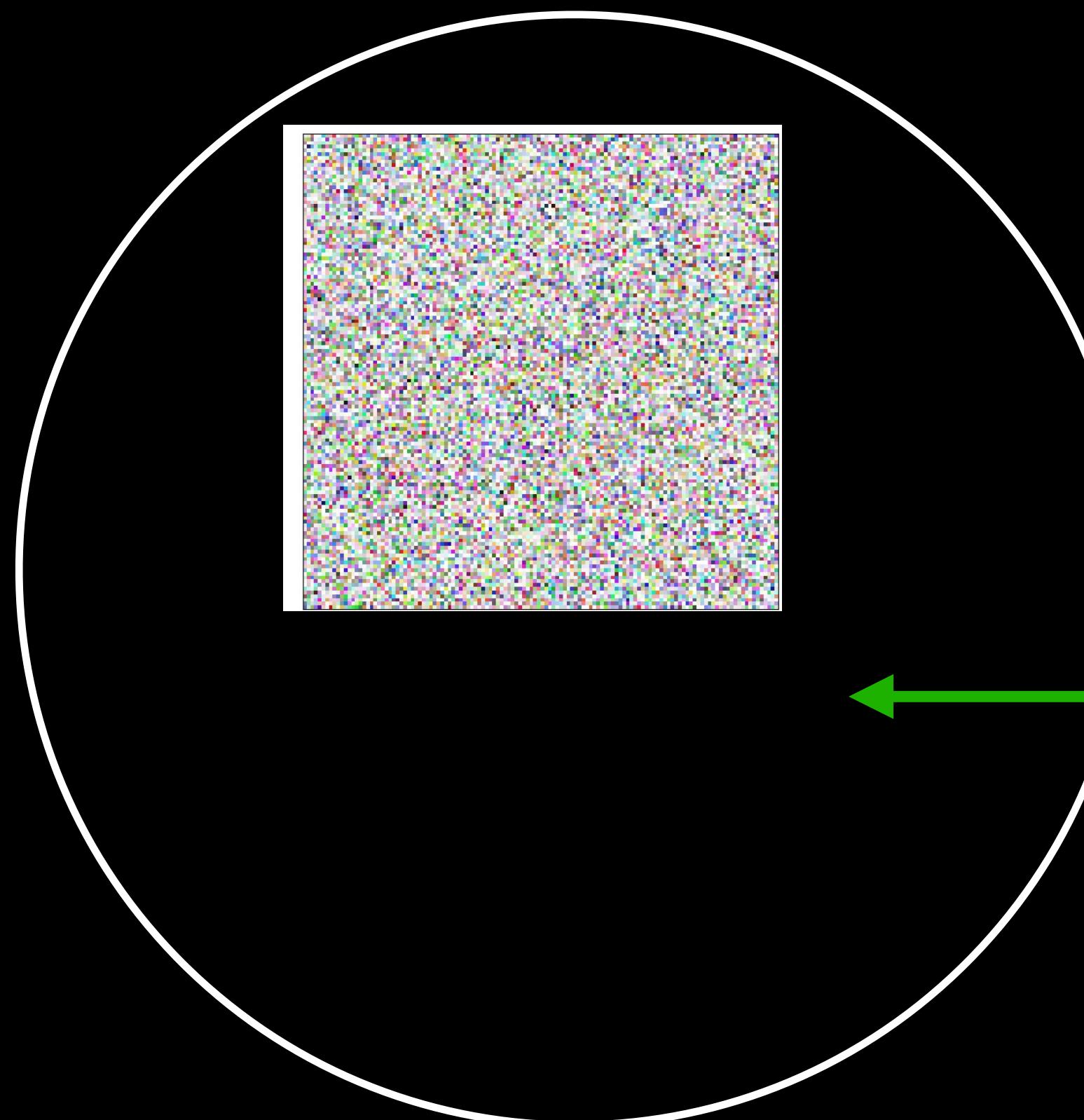




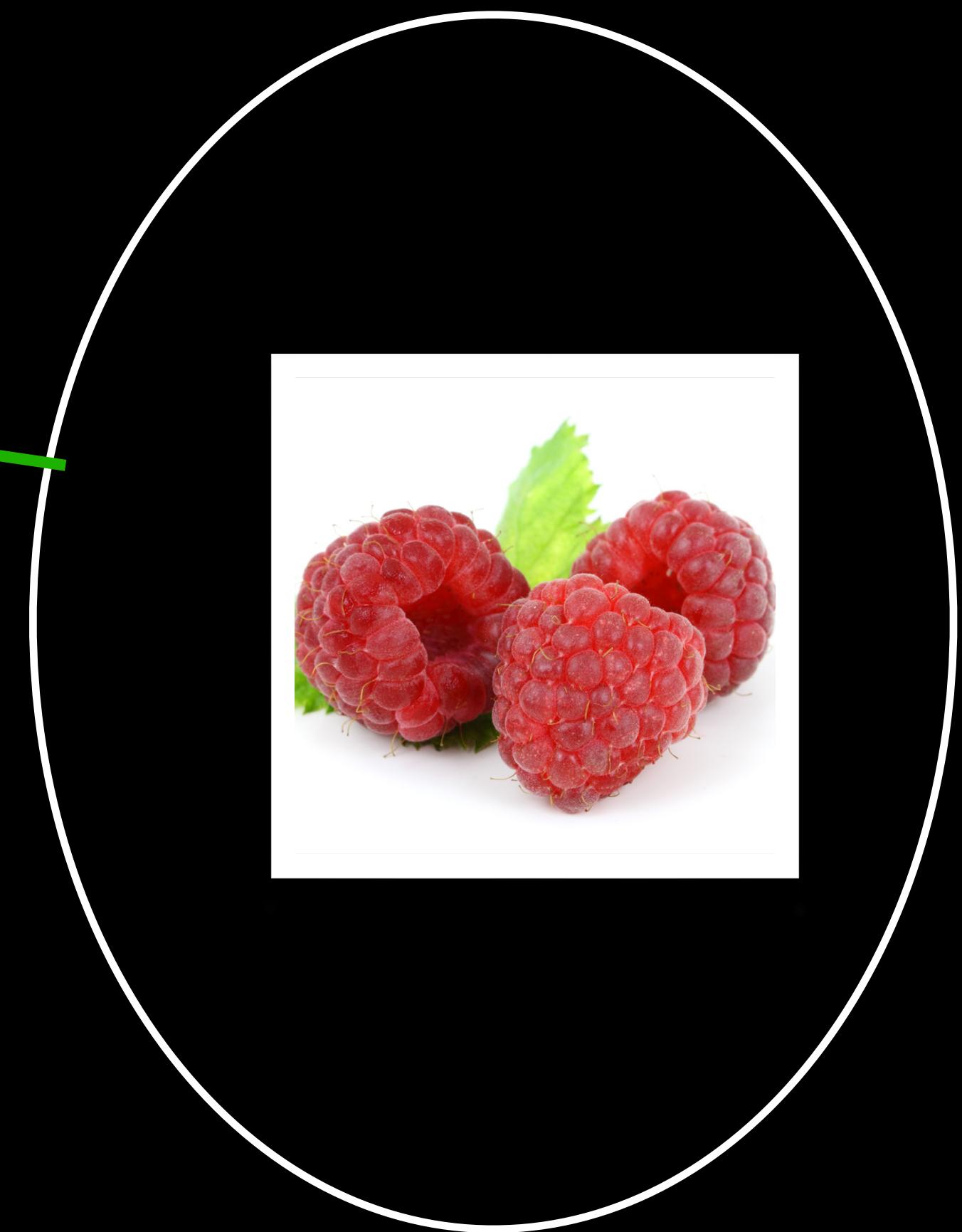
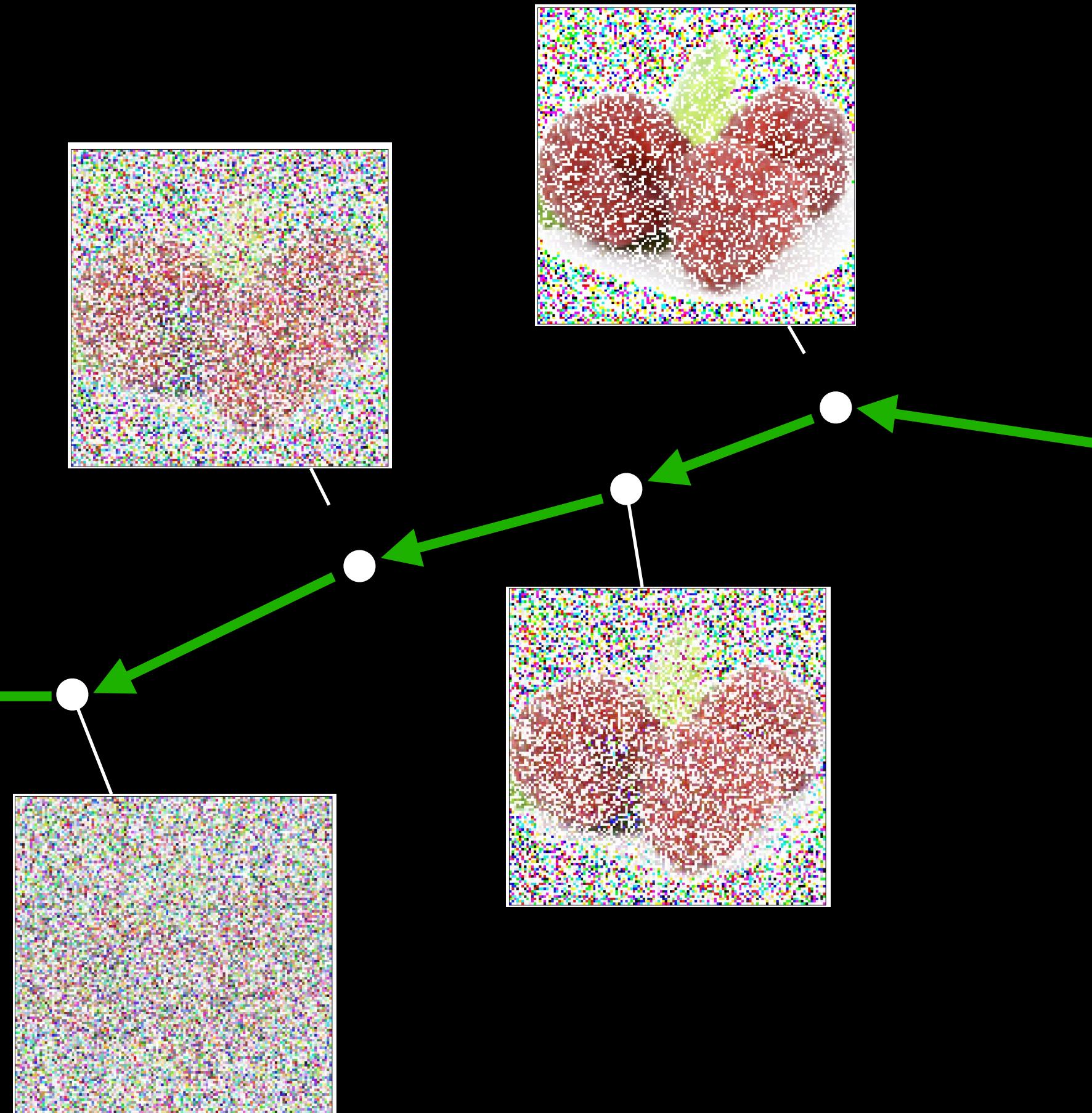
Forward destruction trajectory



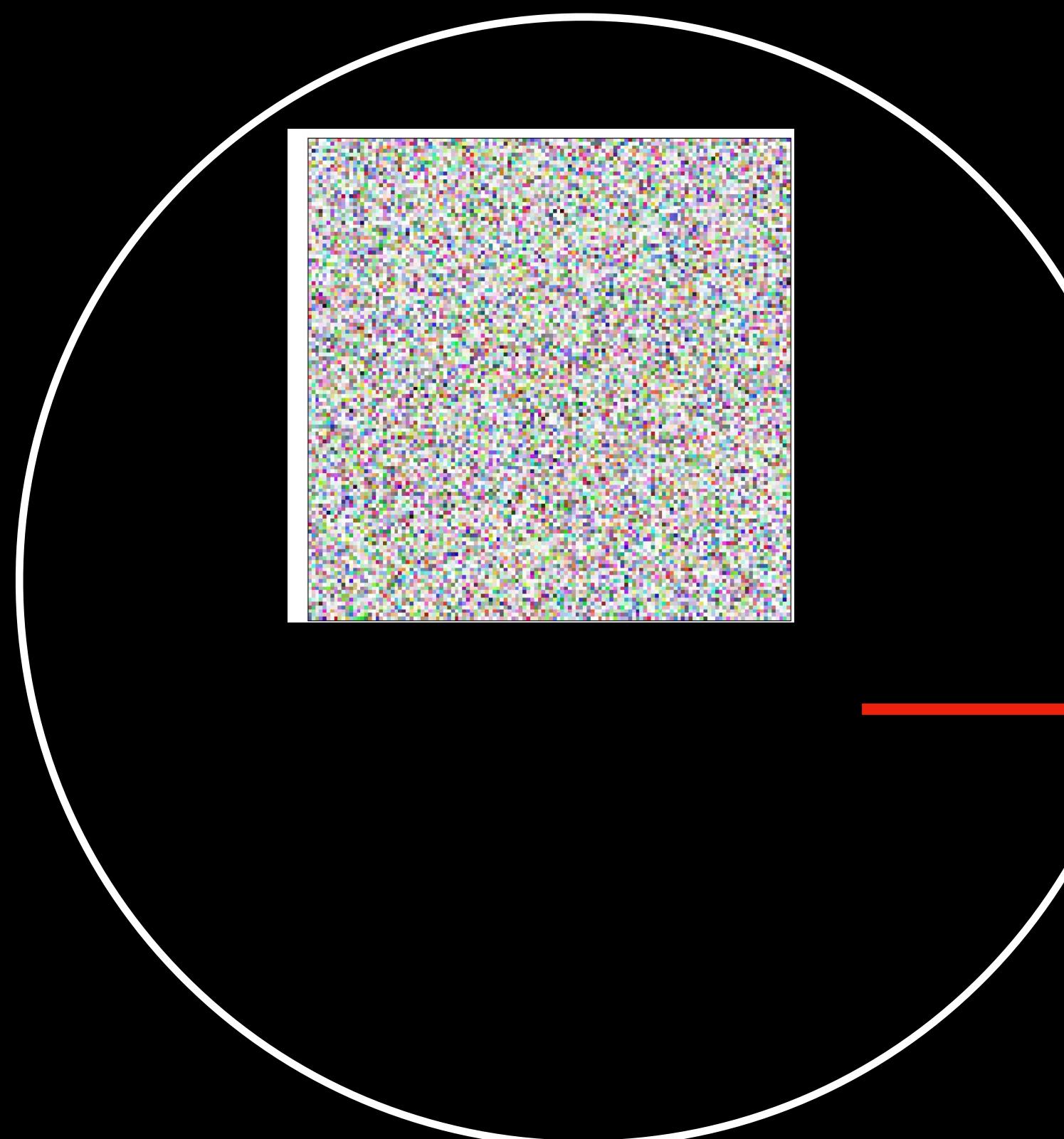
Reverse destruction trajectory



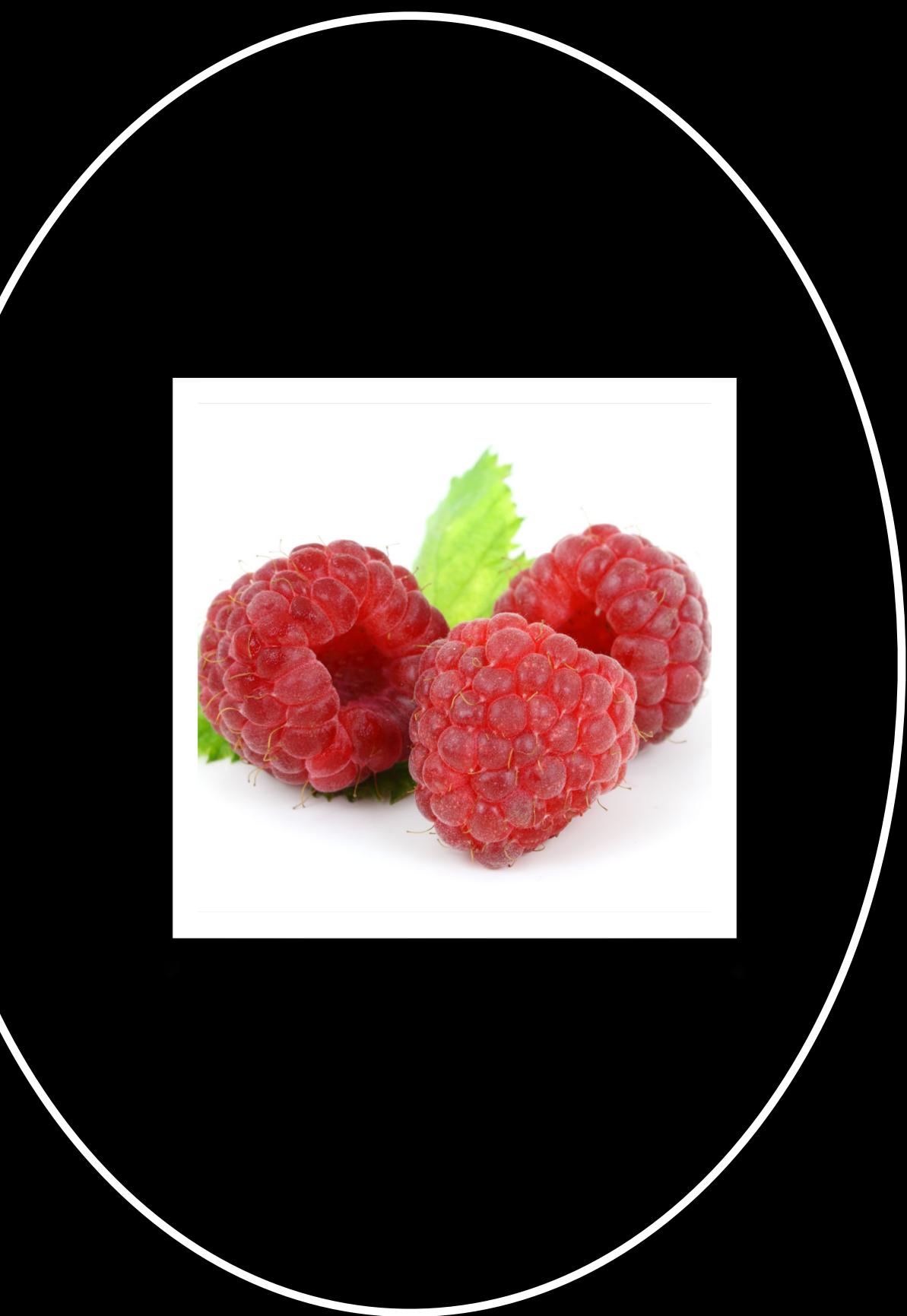
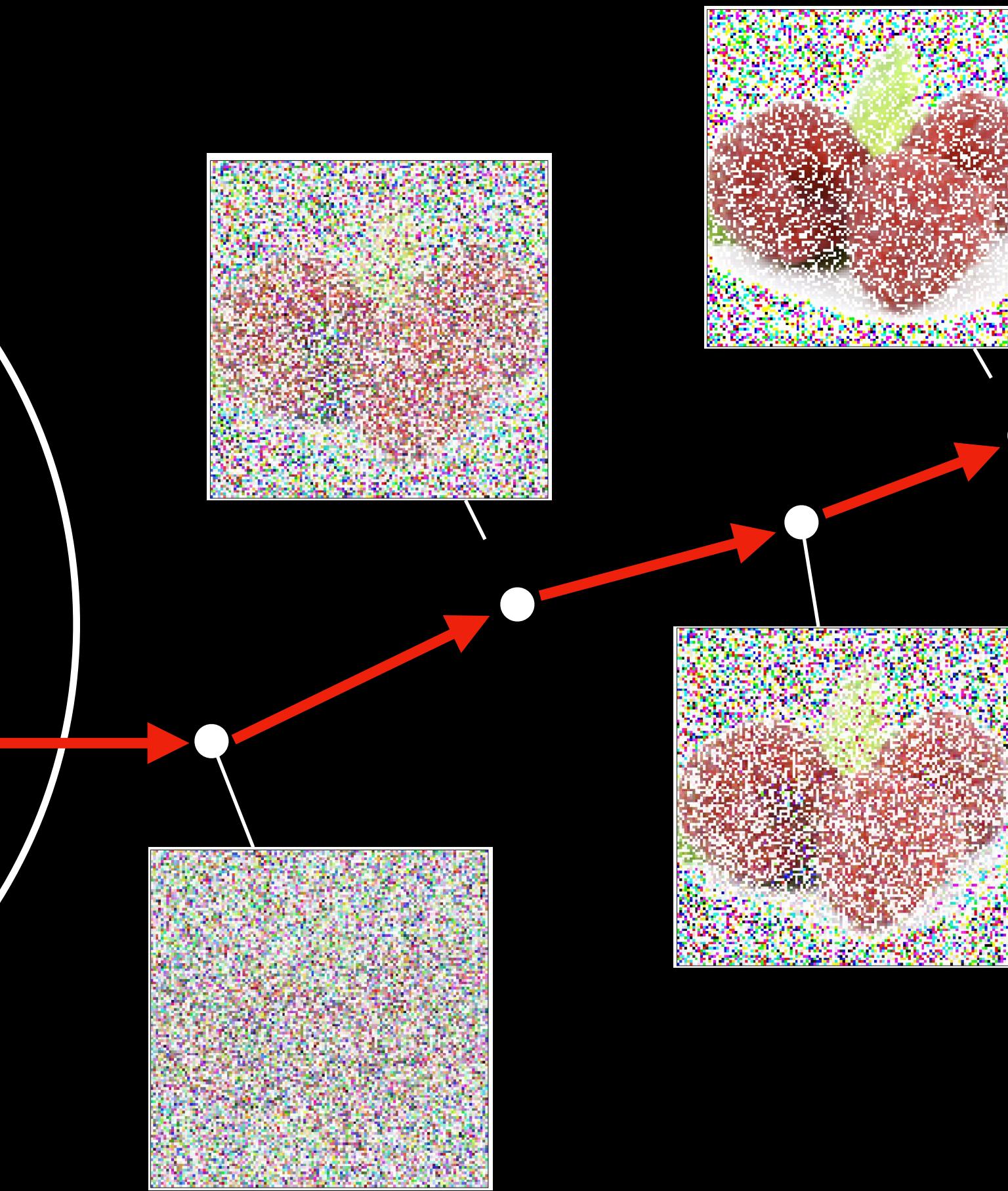
Pure Noise



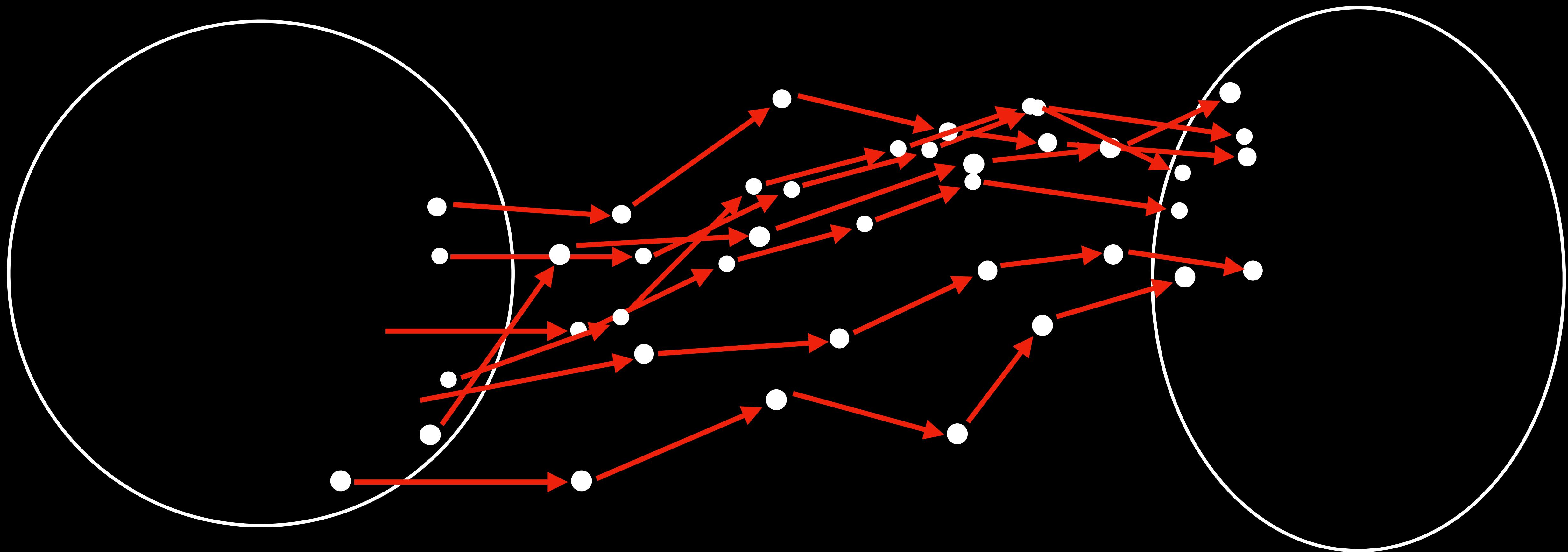
Real Images



Pure Noise

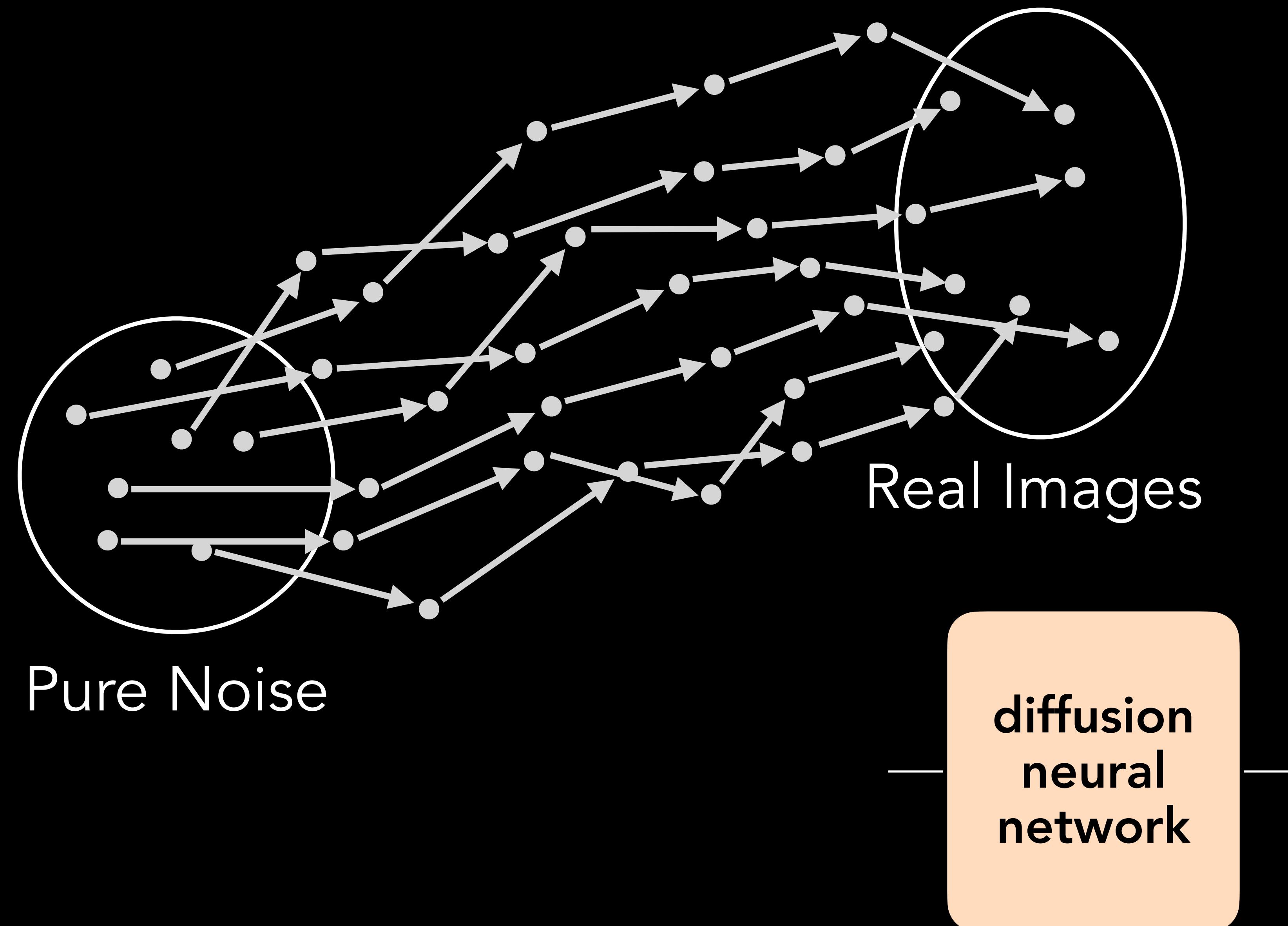


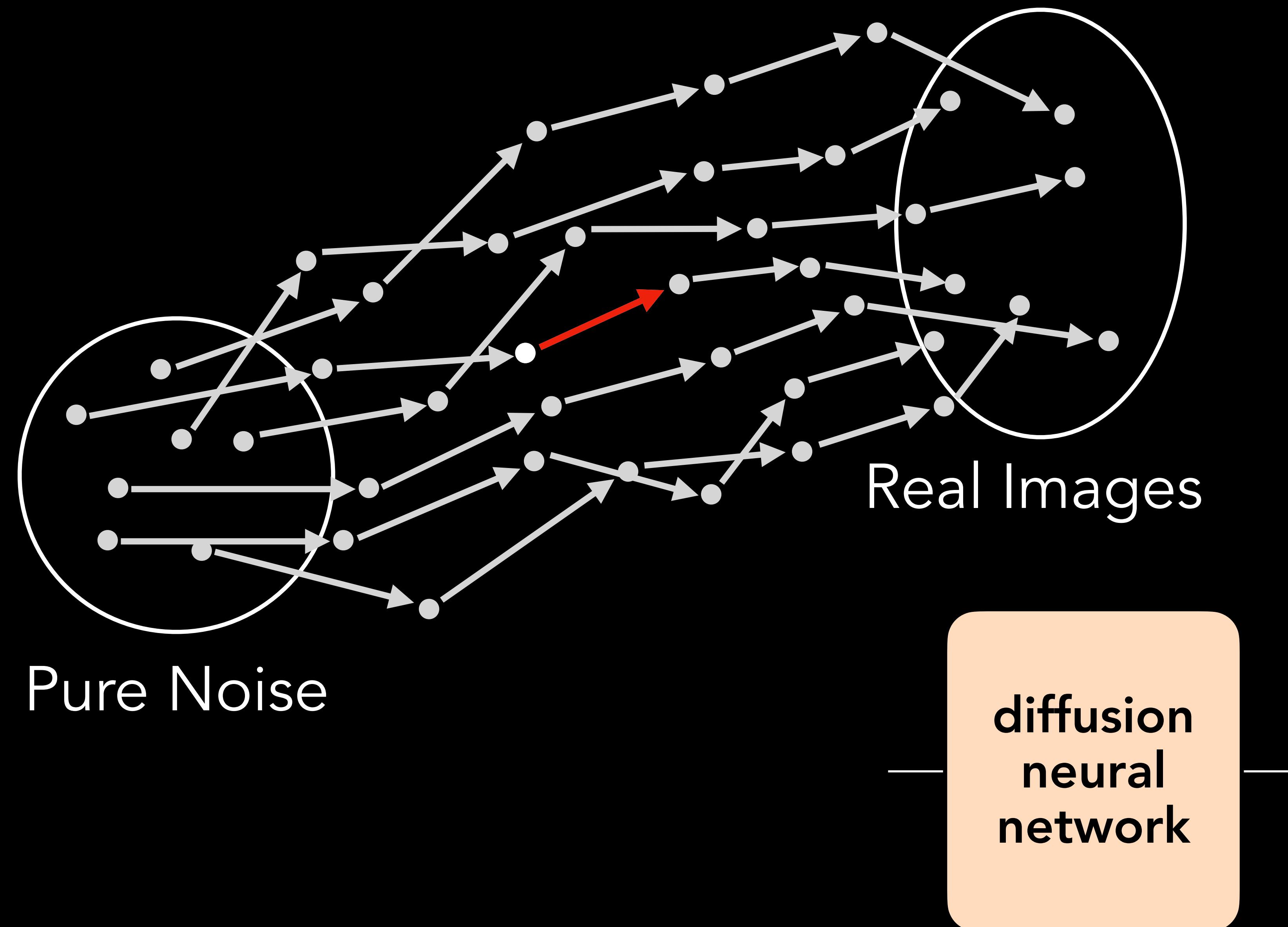
Real Images

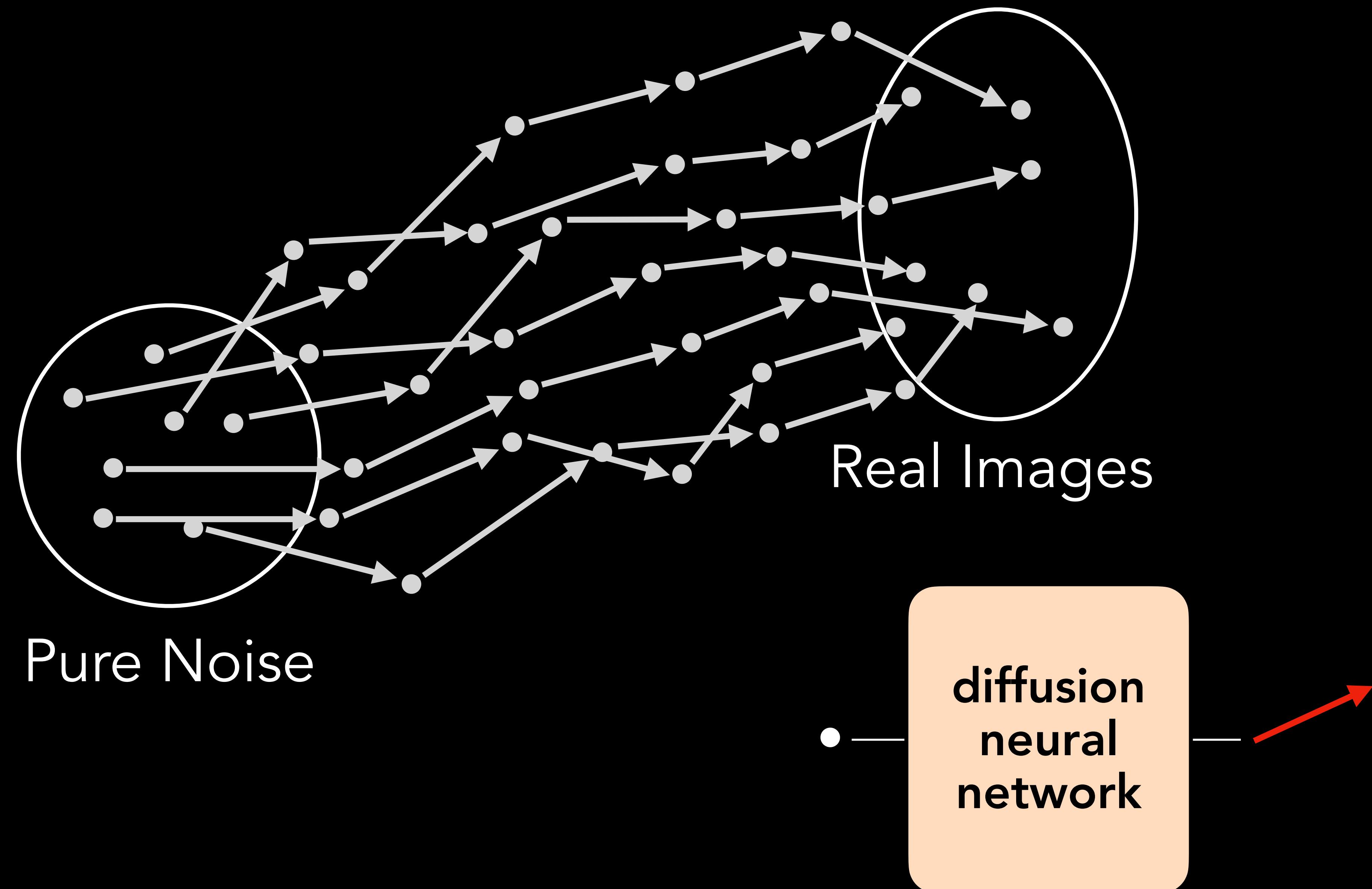


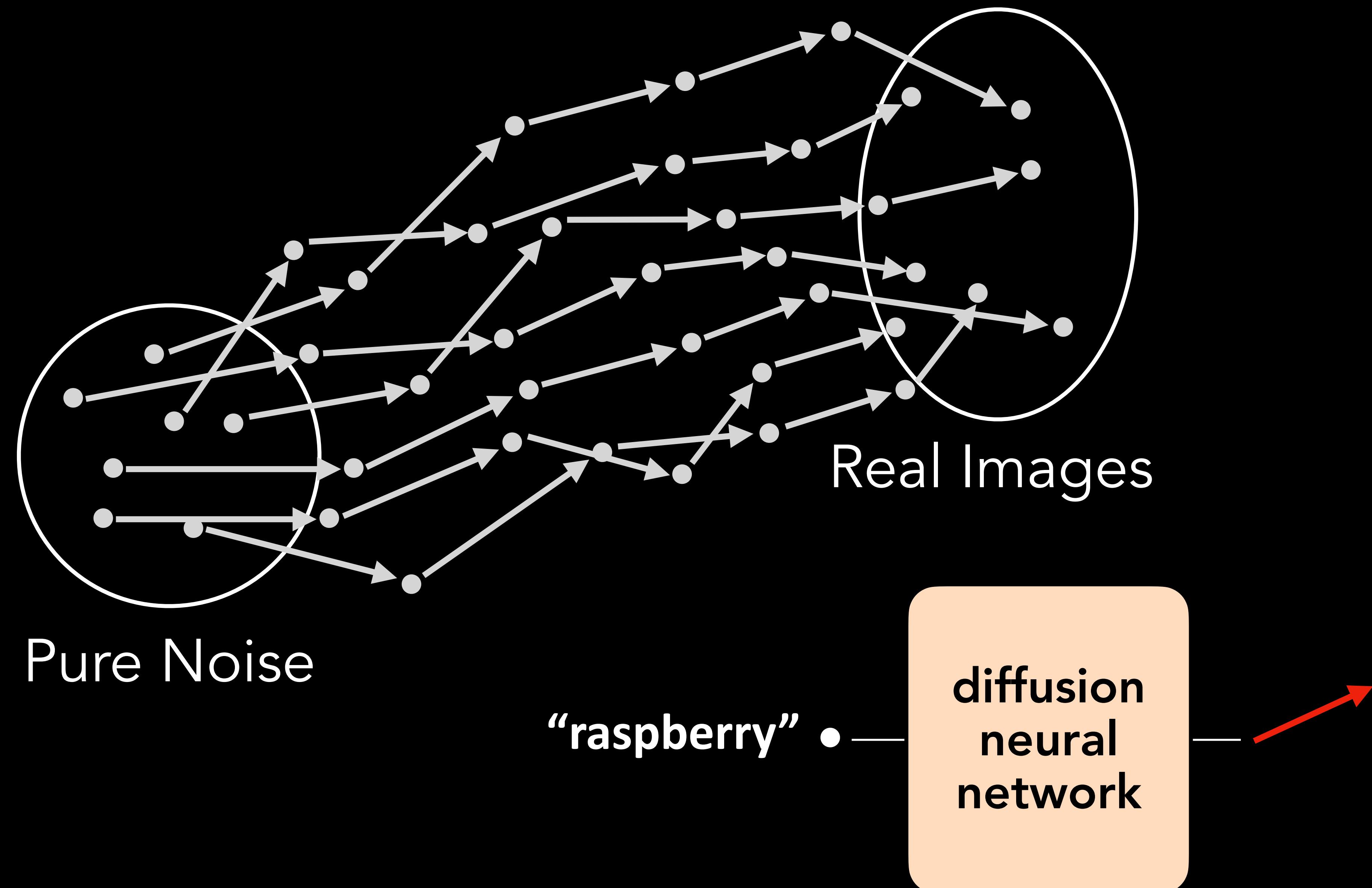
Pure Noise

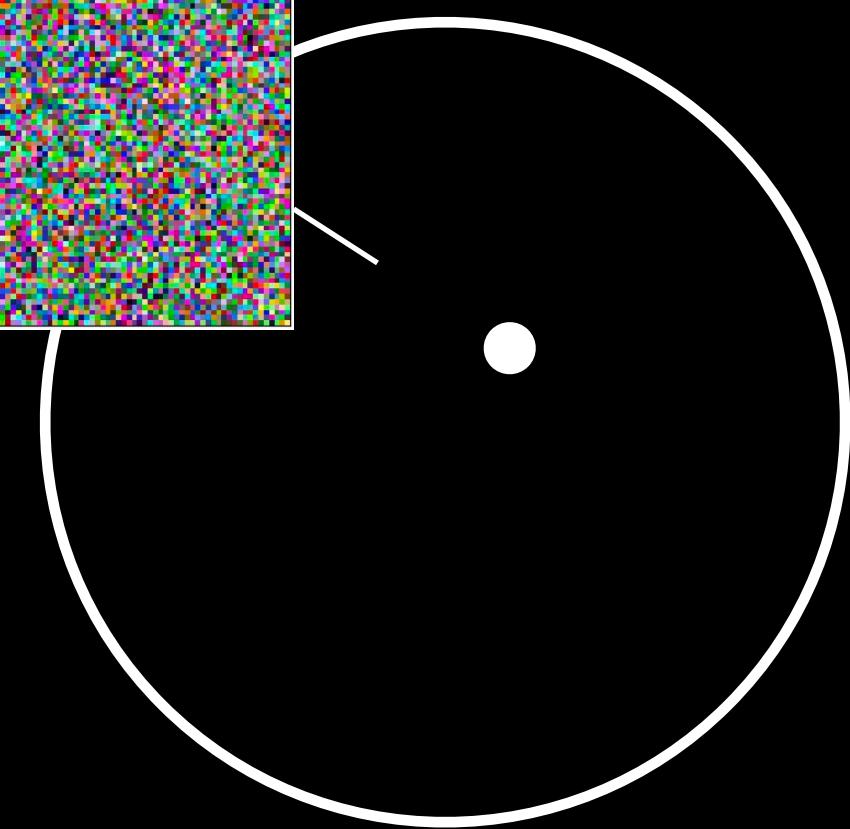
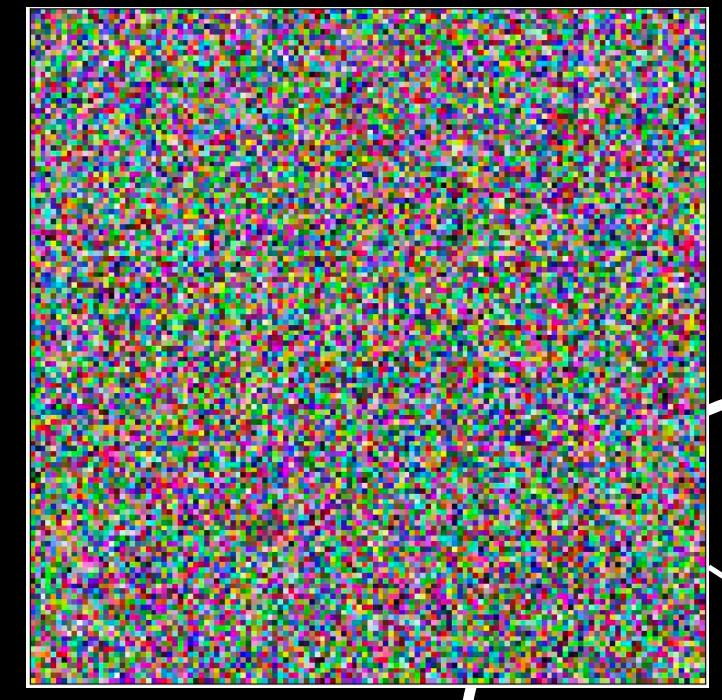
Real Images



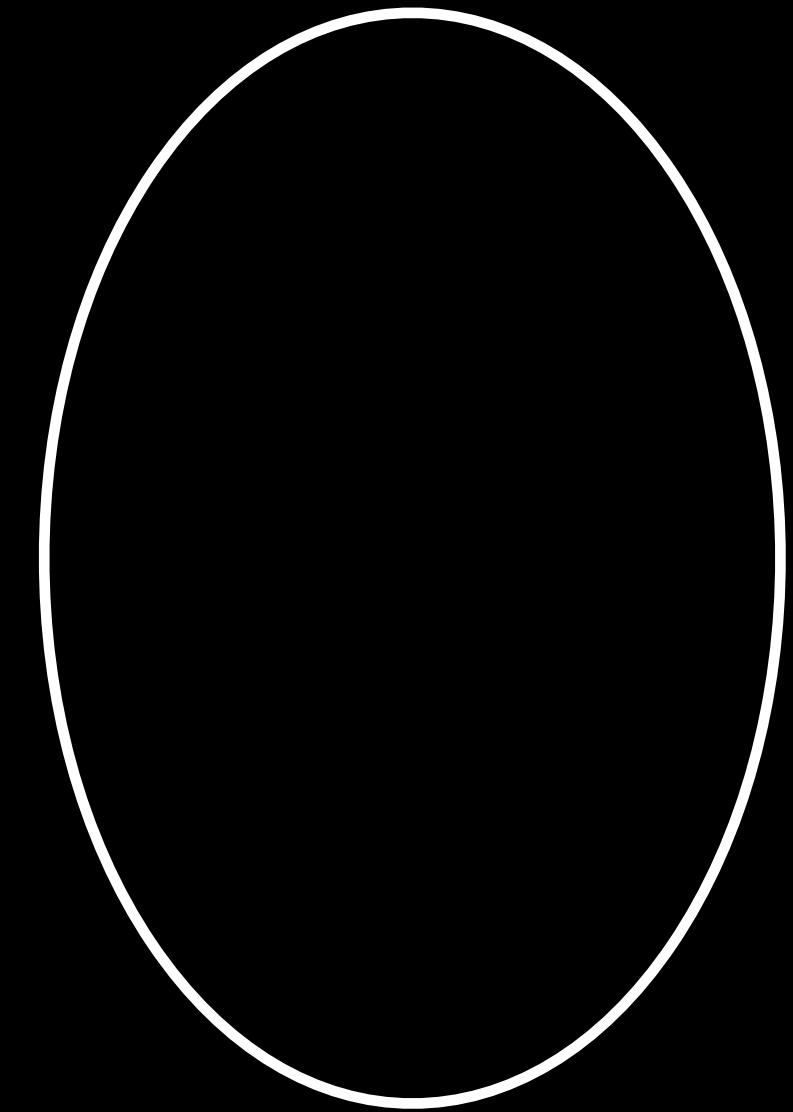








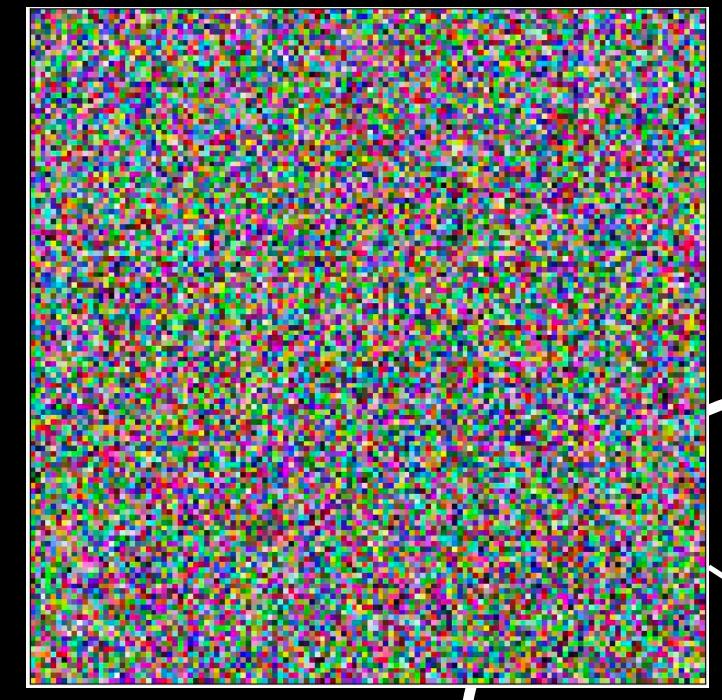
Pure Noise



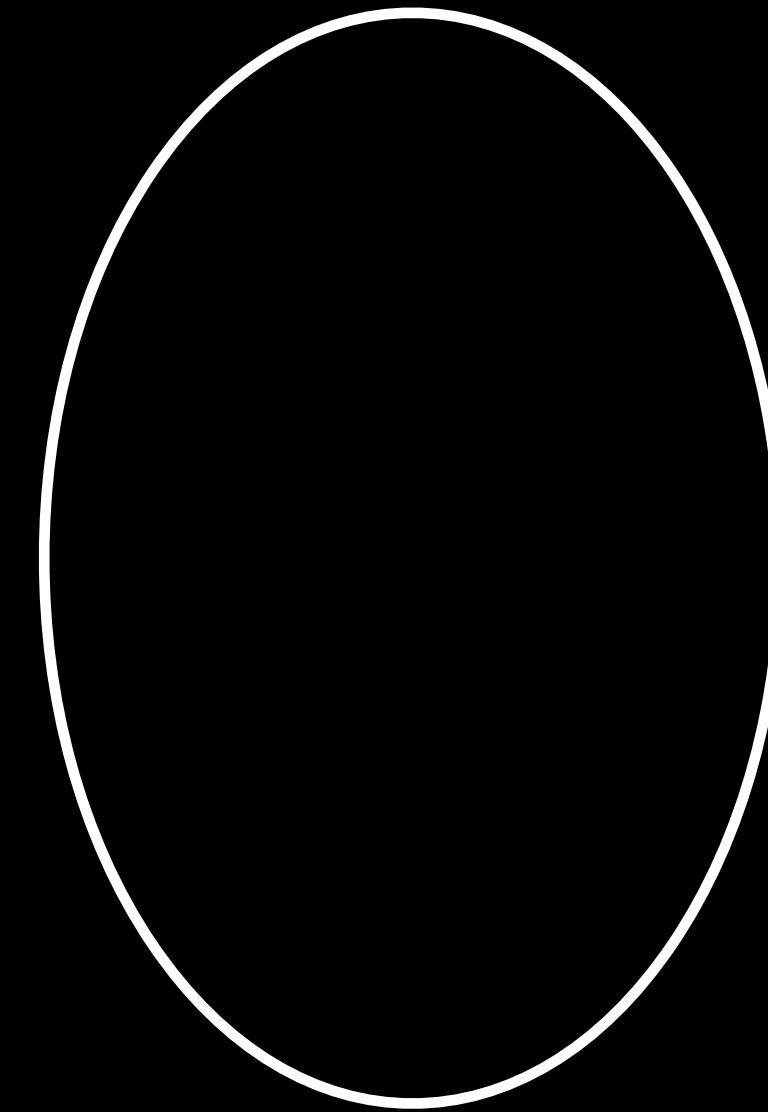
Real Images

“raspberry”

**diffusion
neural
network**



Pure Noise

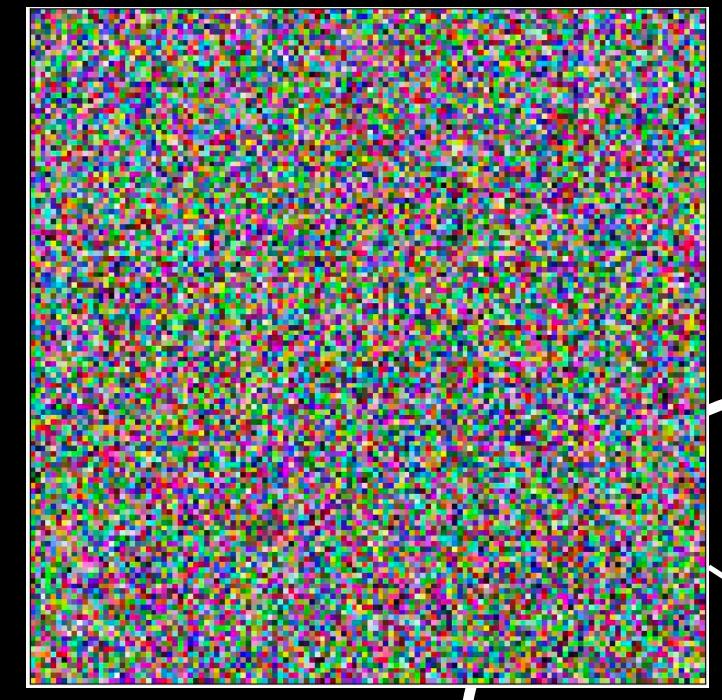


Real Images

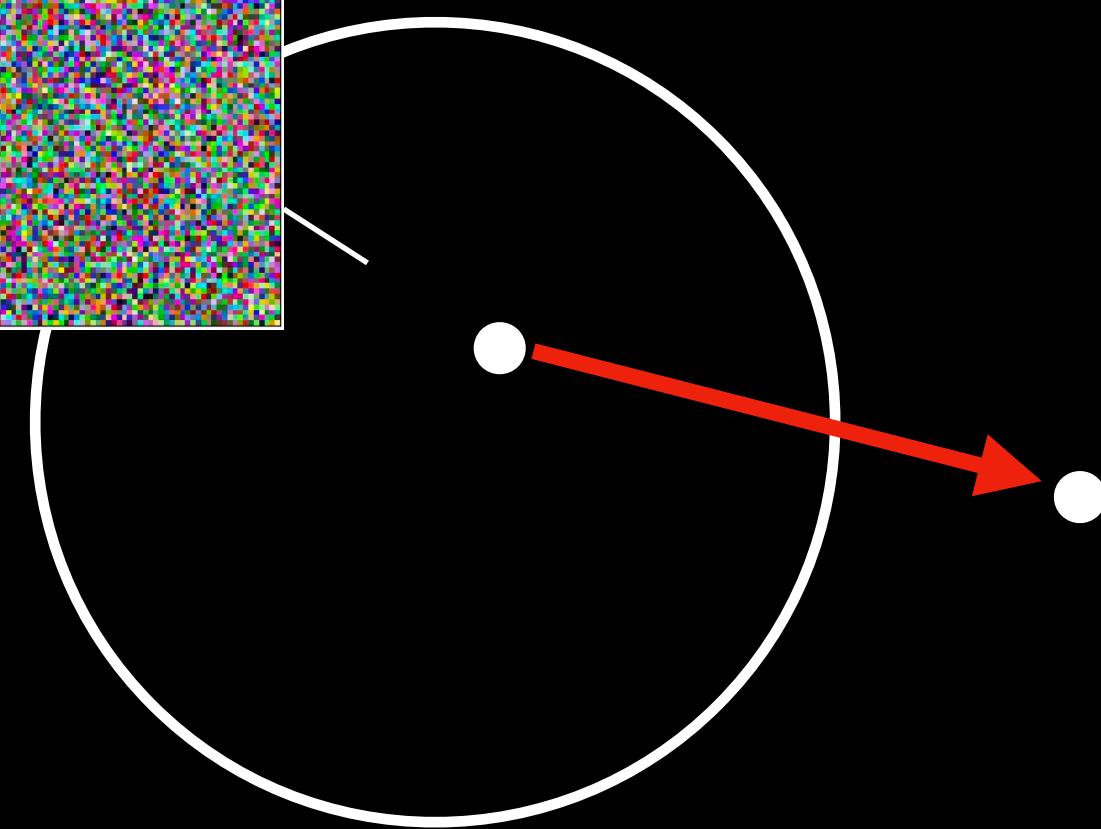
“raspberry” • —

diffusion
neural
network

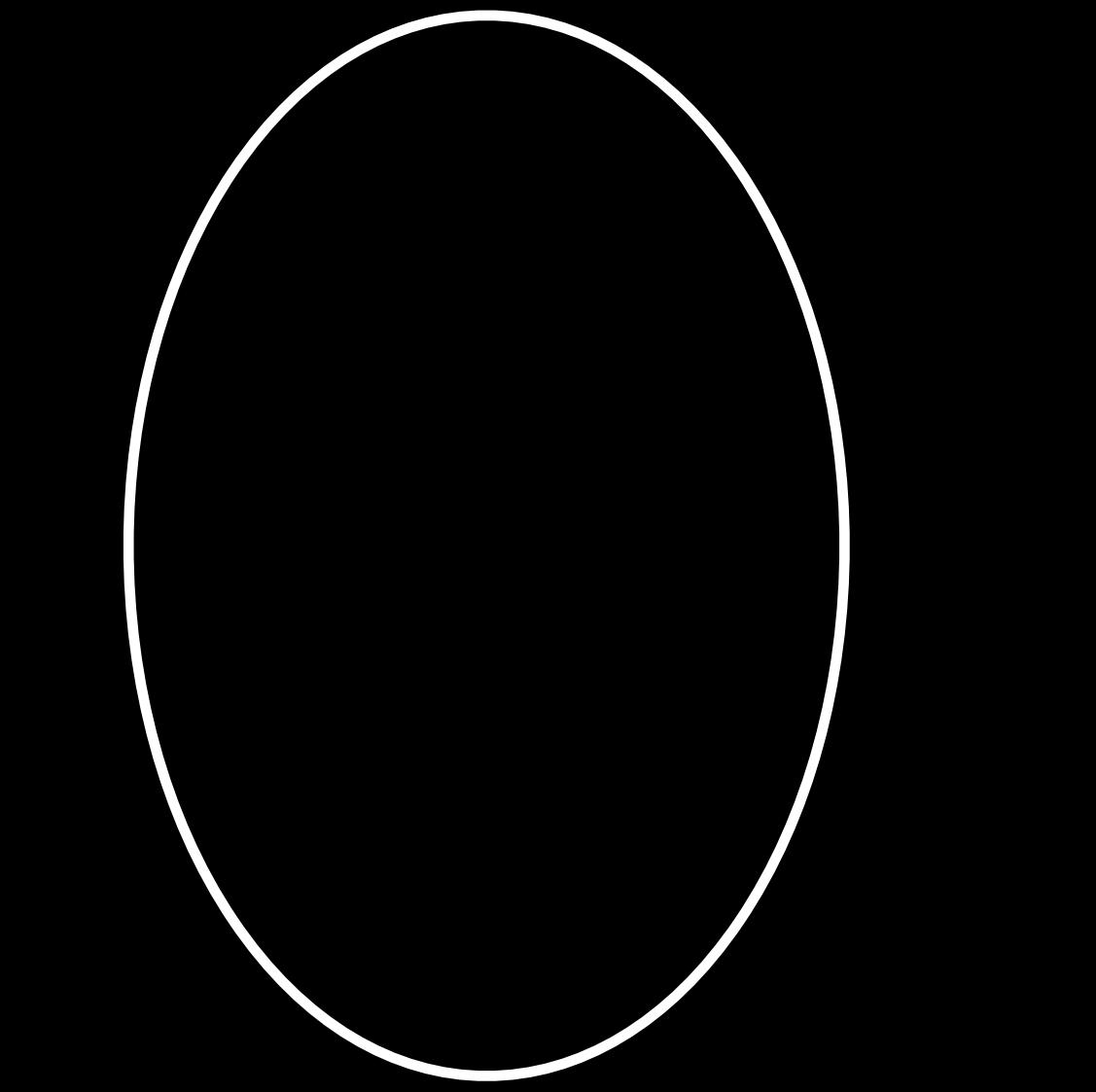




Pure Noise

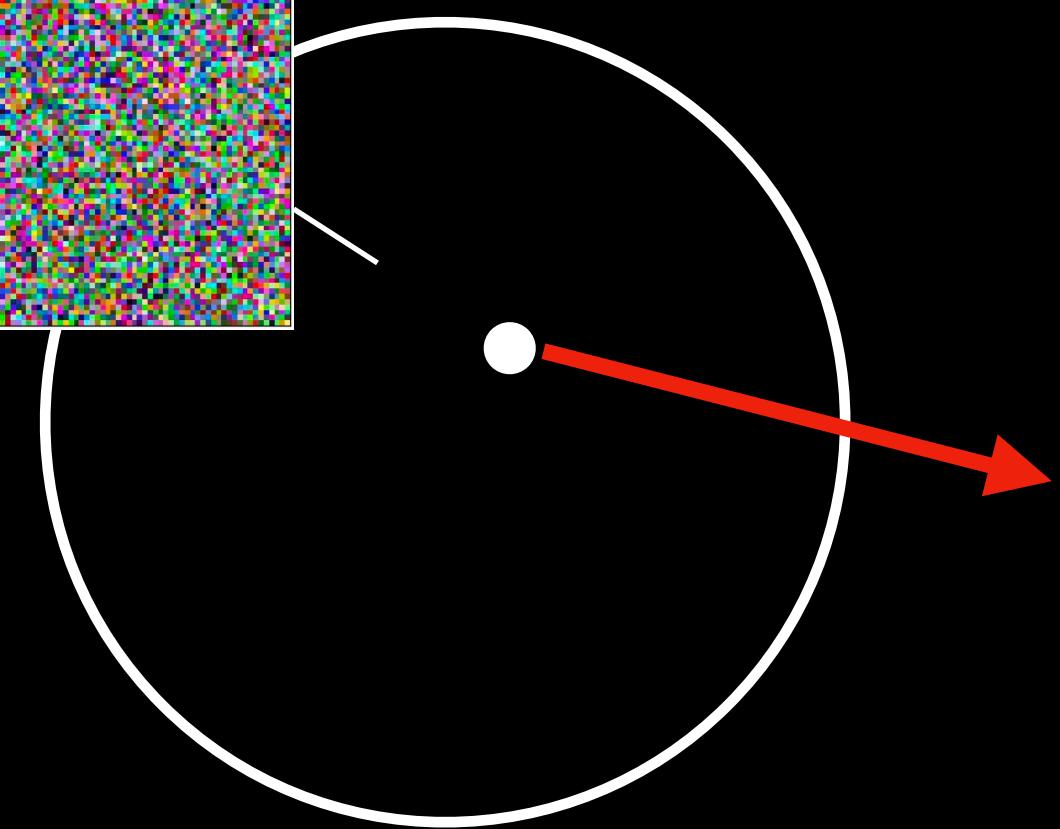
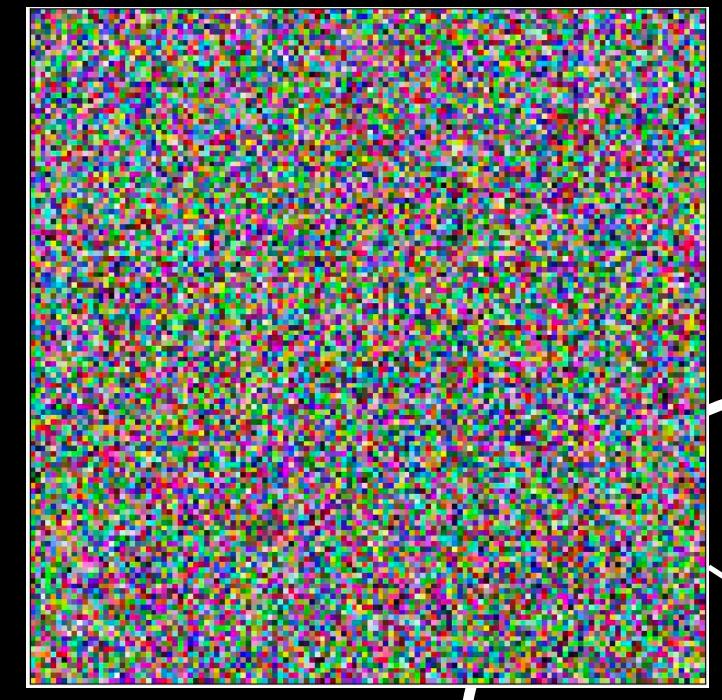


“raspberry”

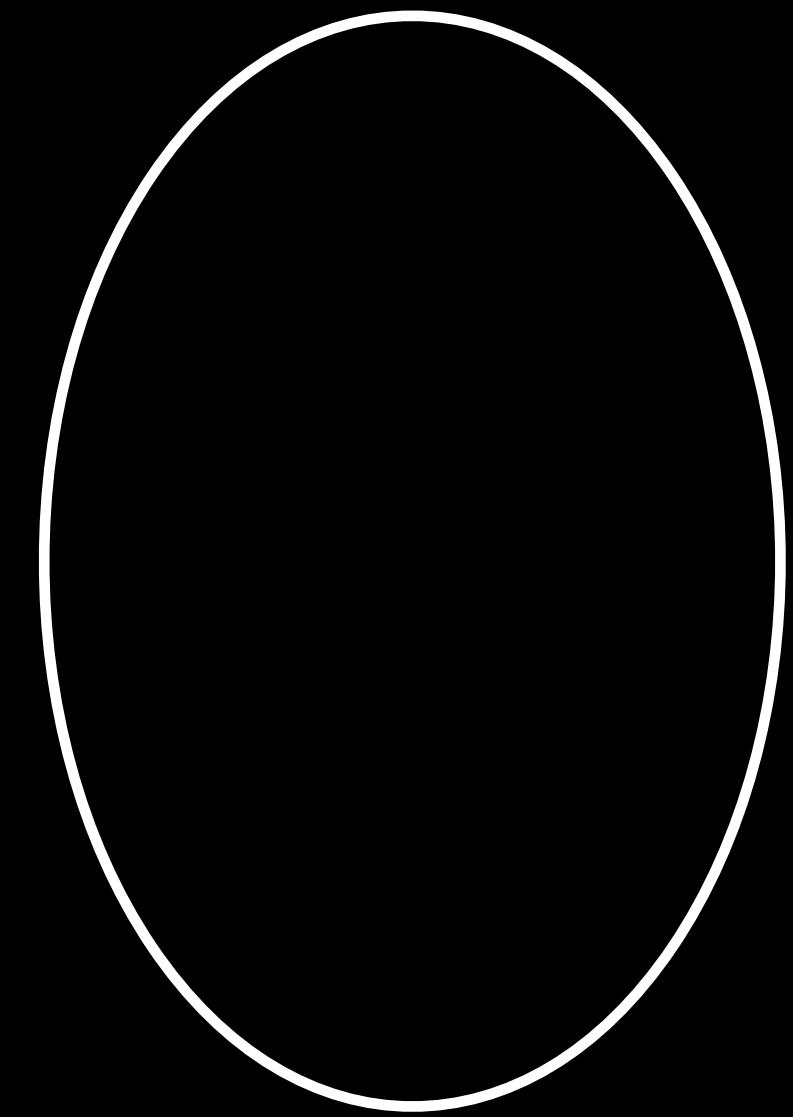


Real Images

**diffusion
neural
network**



Pure Noise

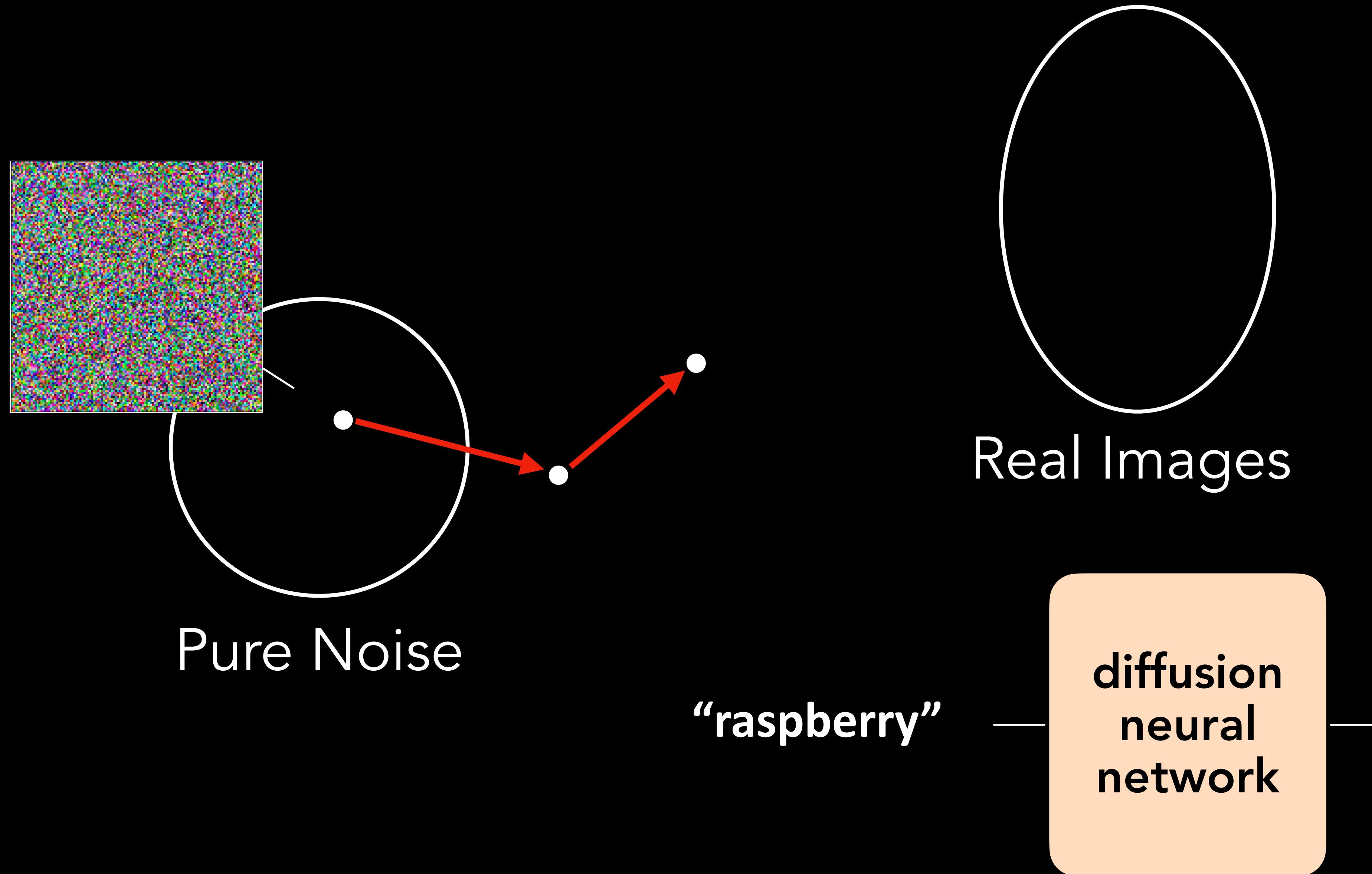


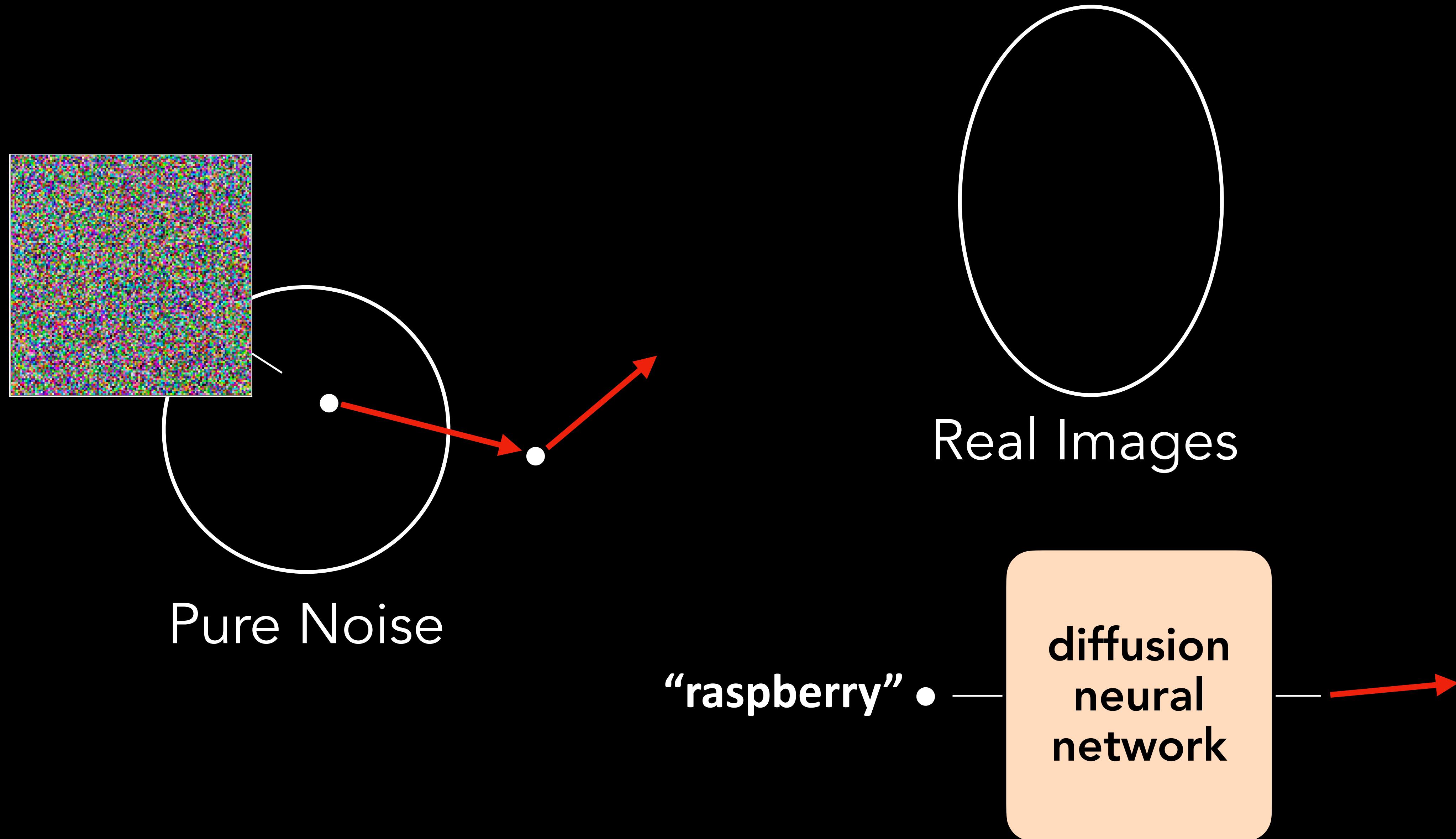
Real Images

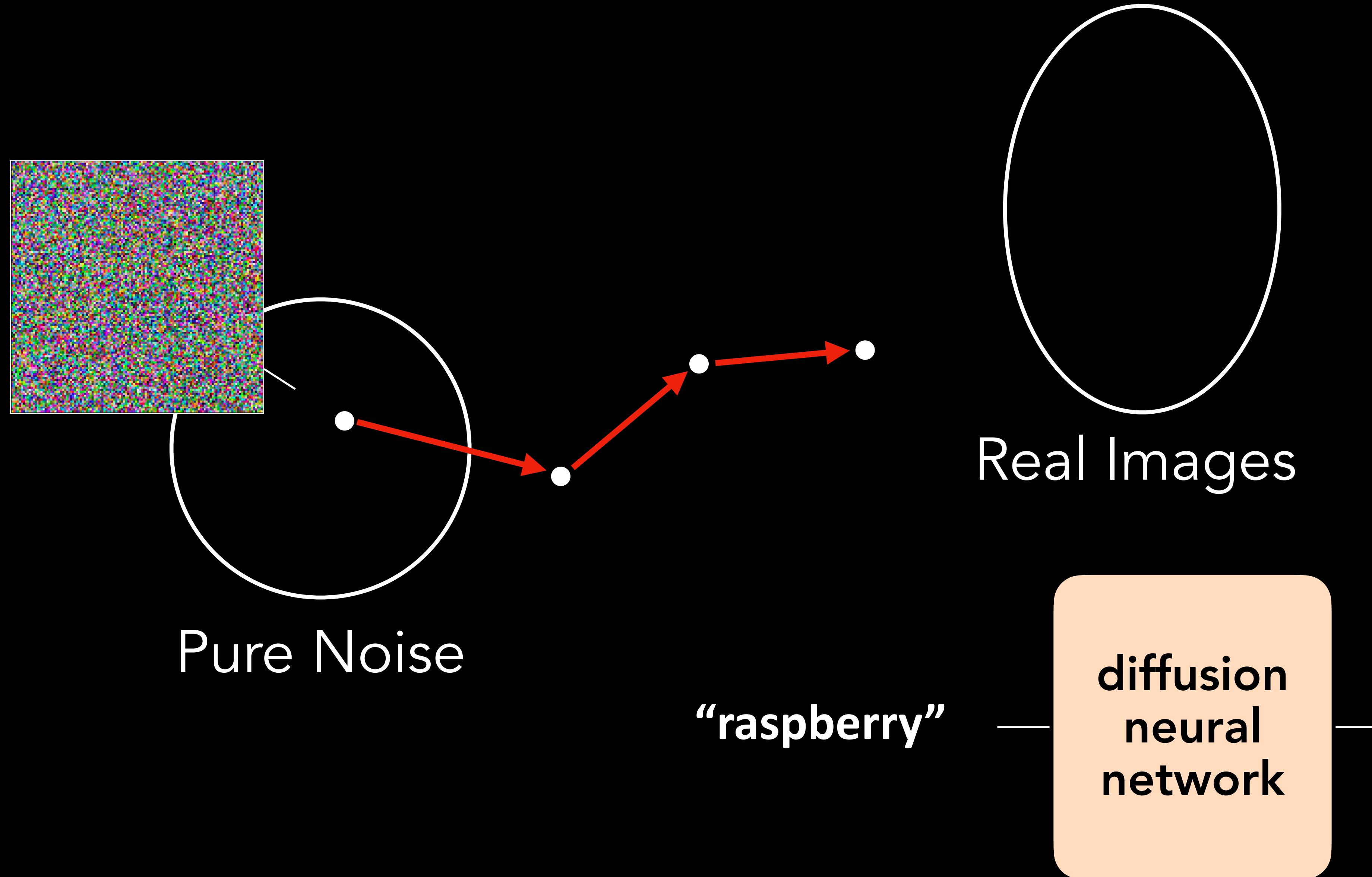
“raspberry”

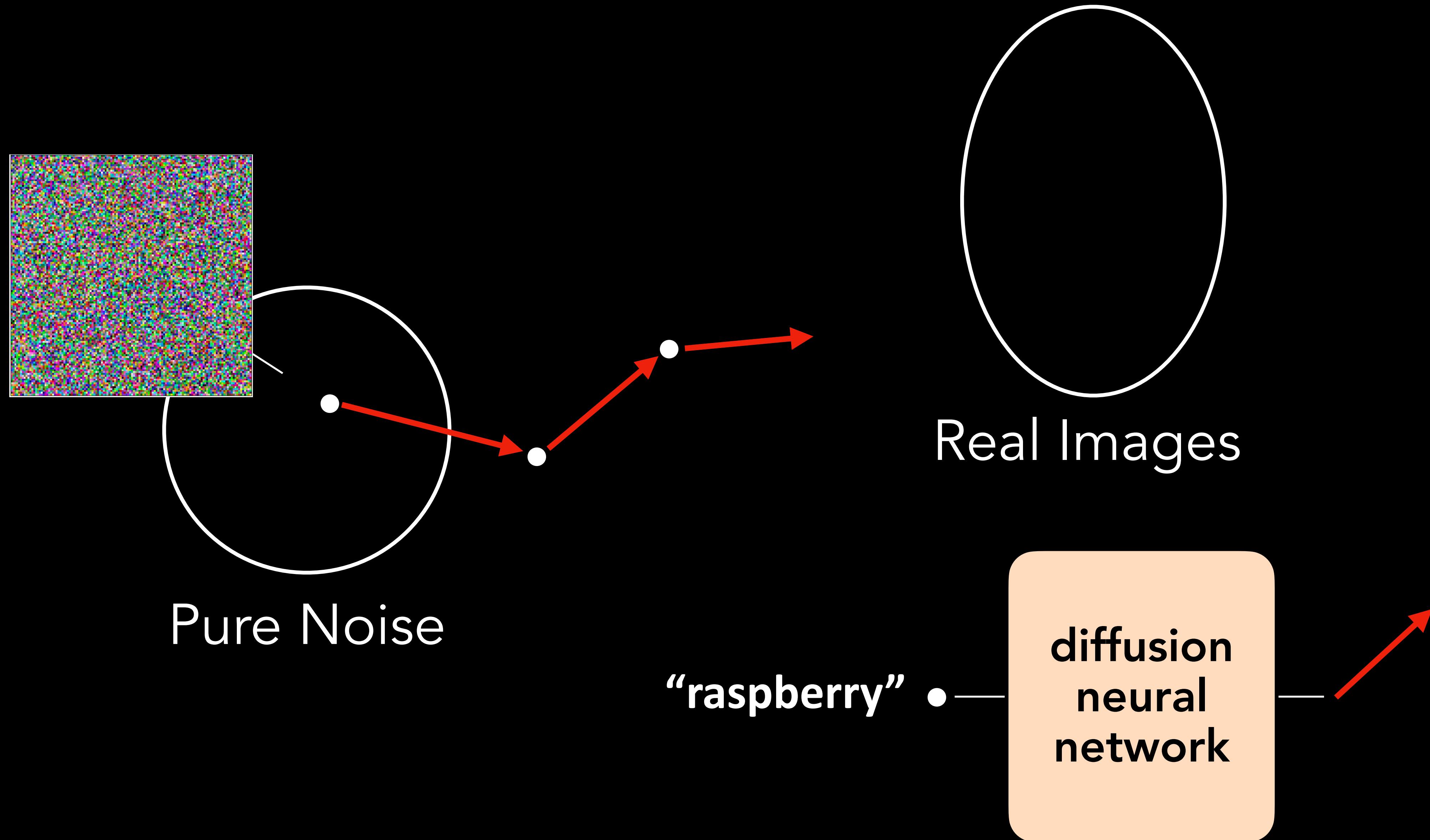
diffusion
neural
network

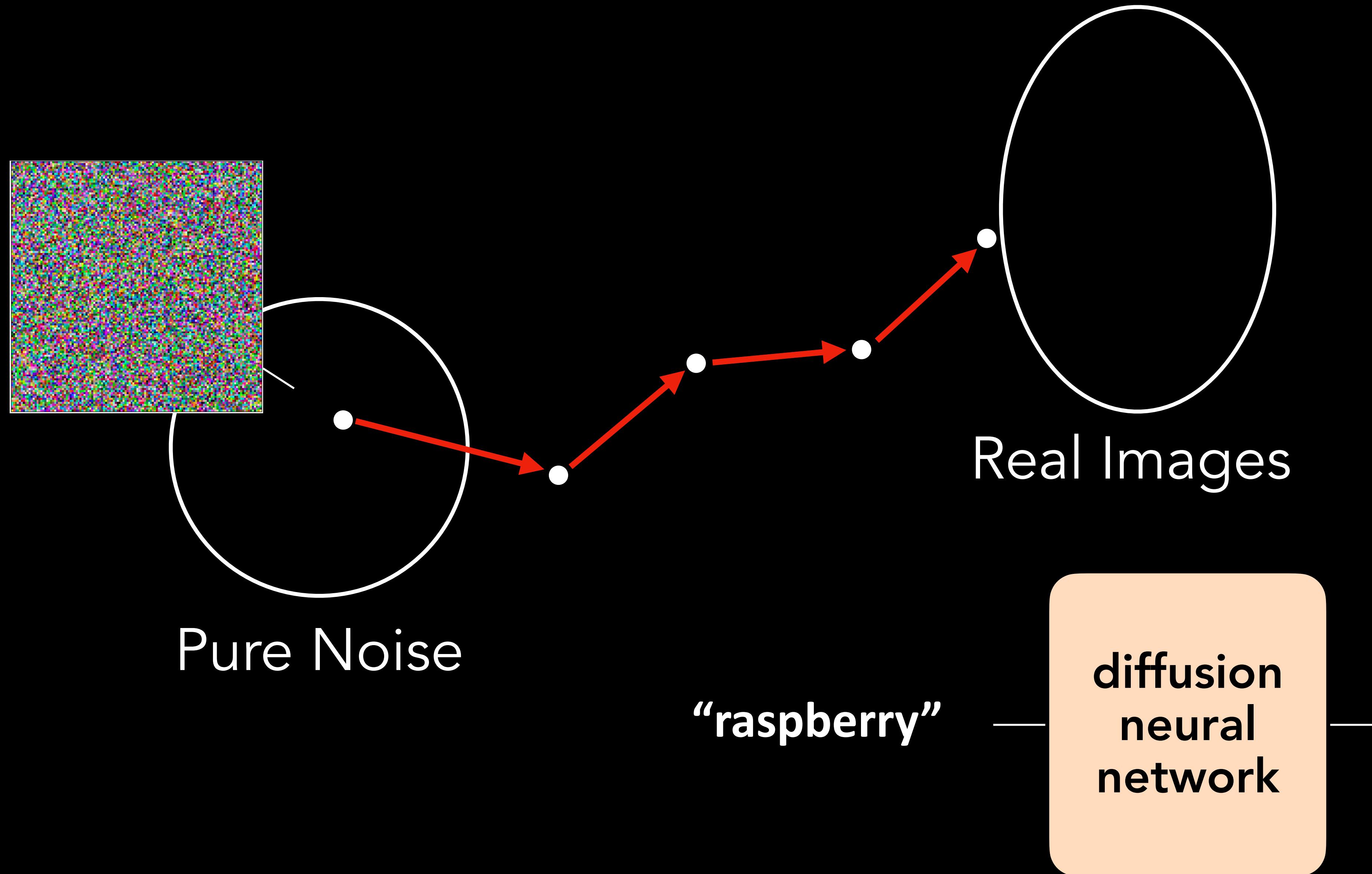


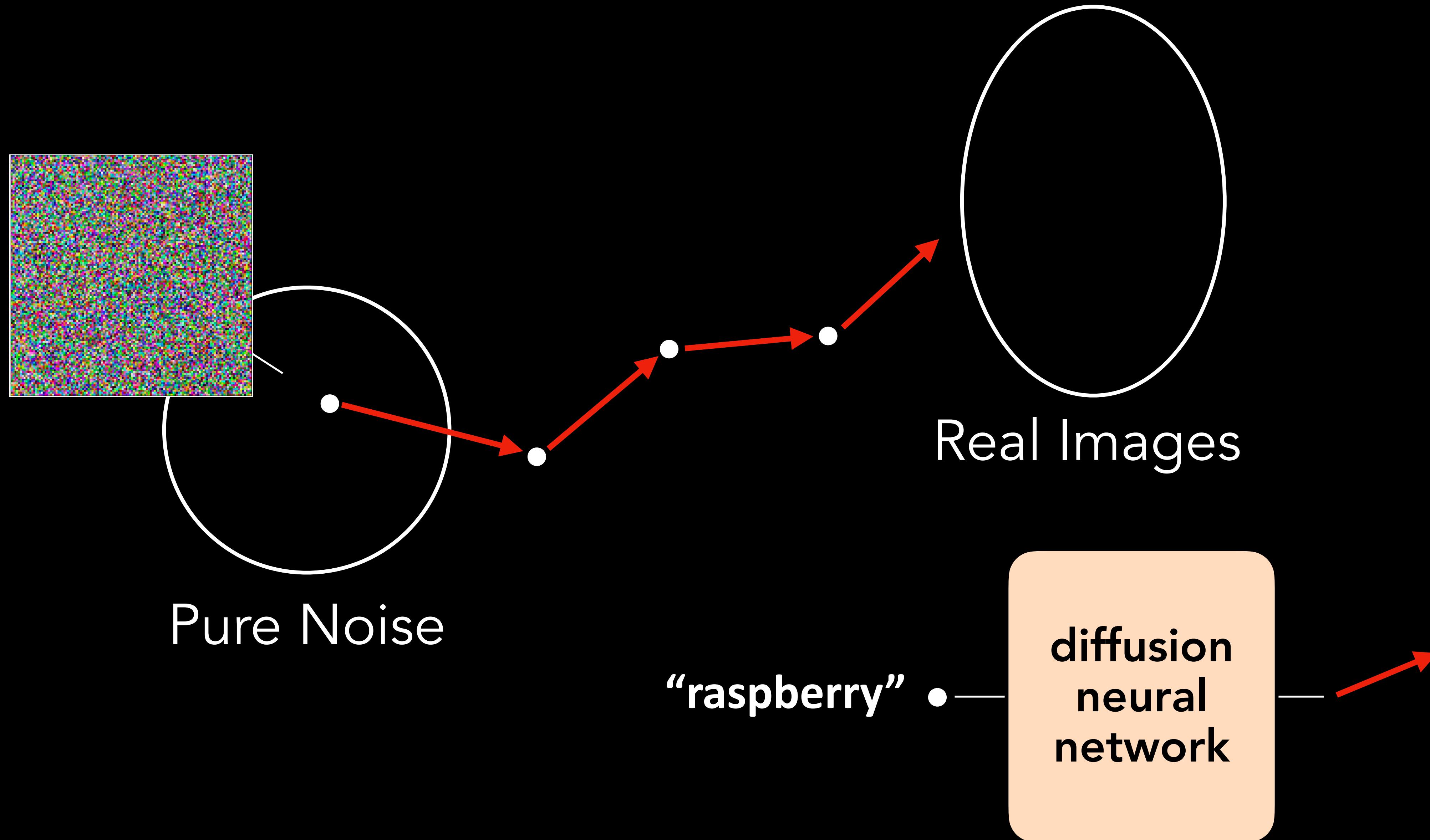


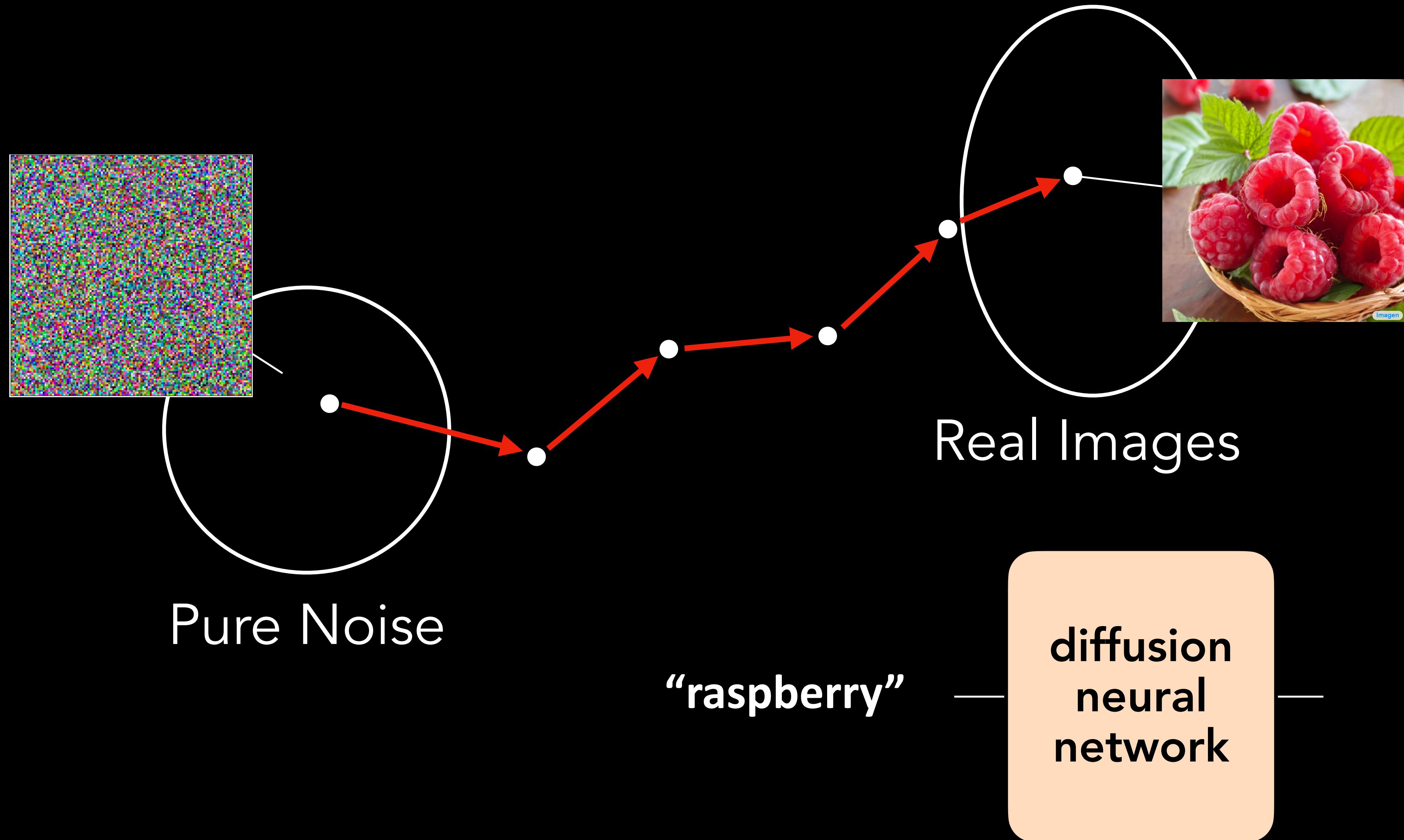




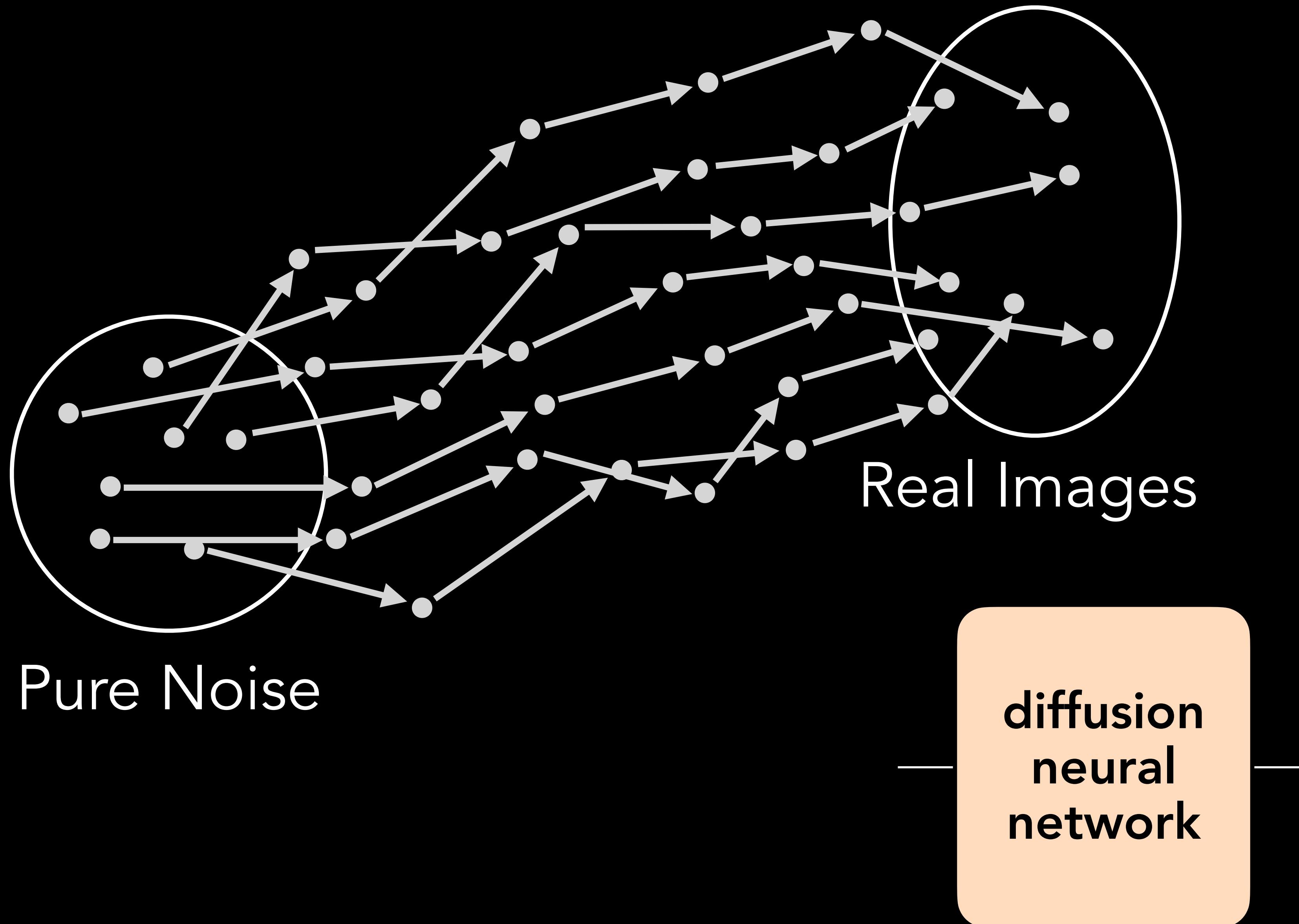


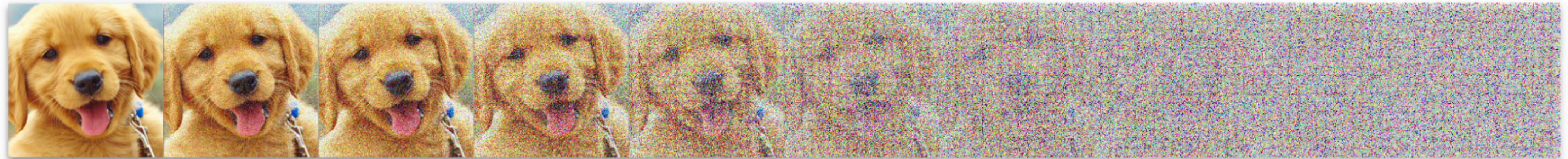




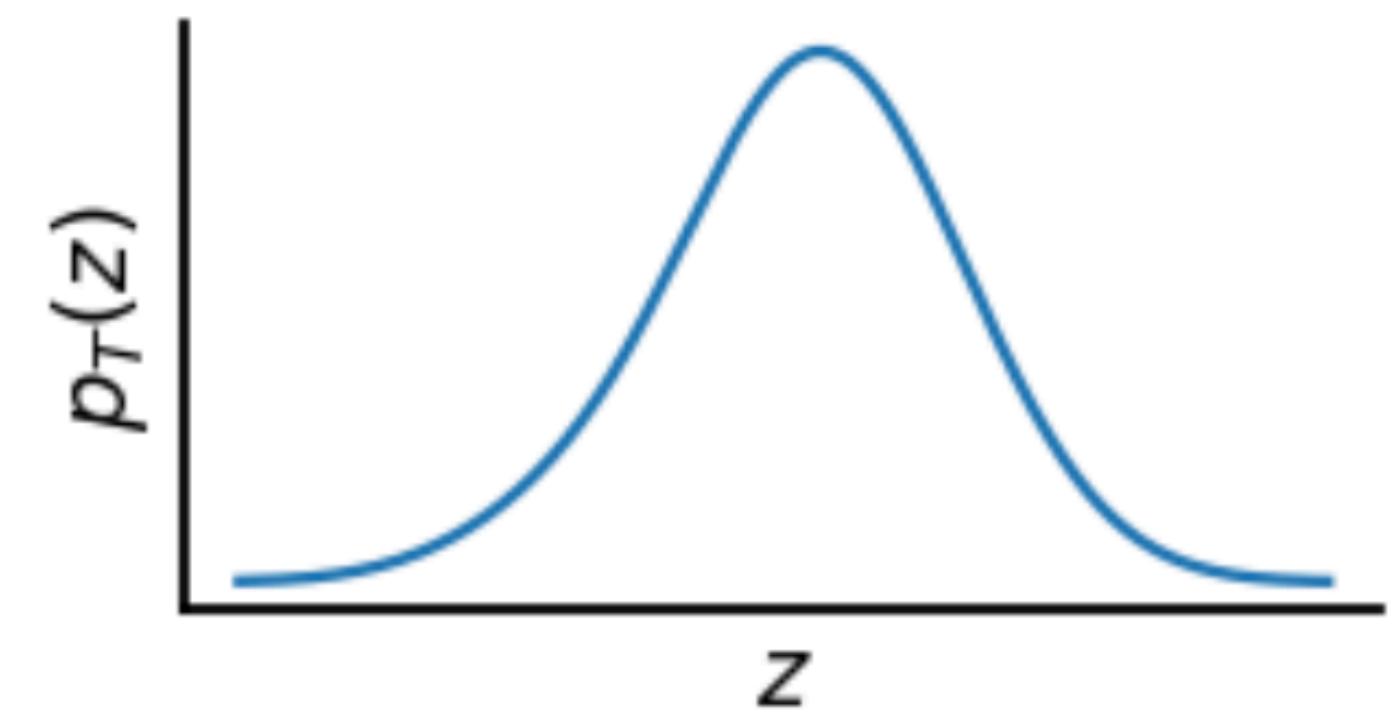
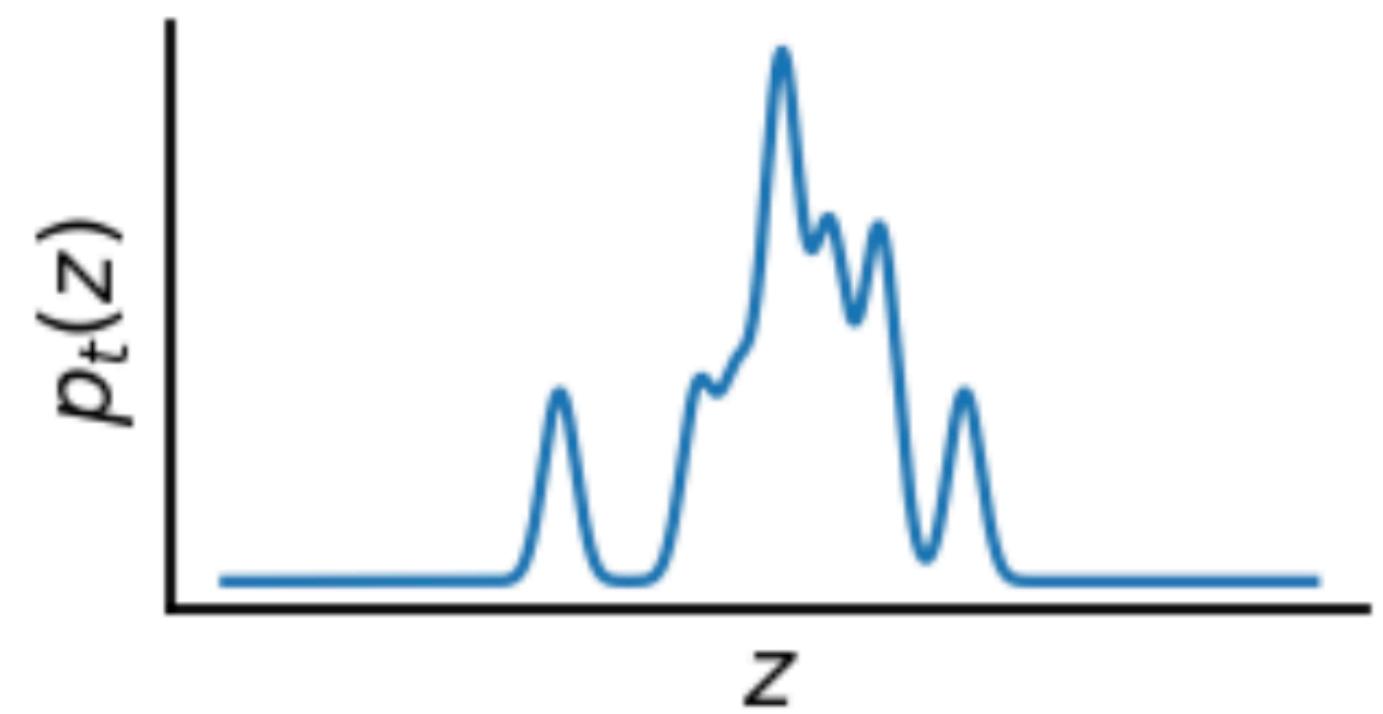
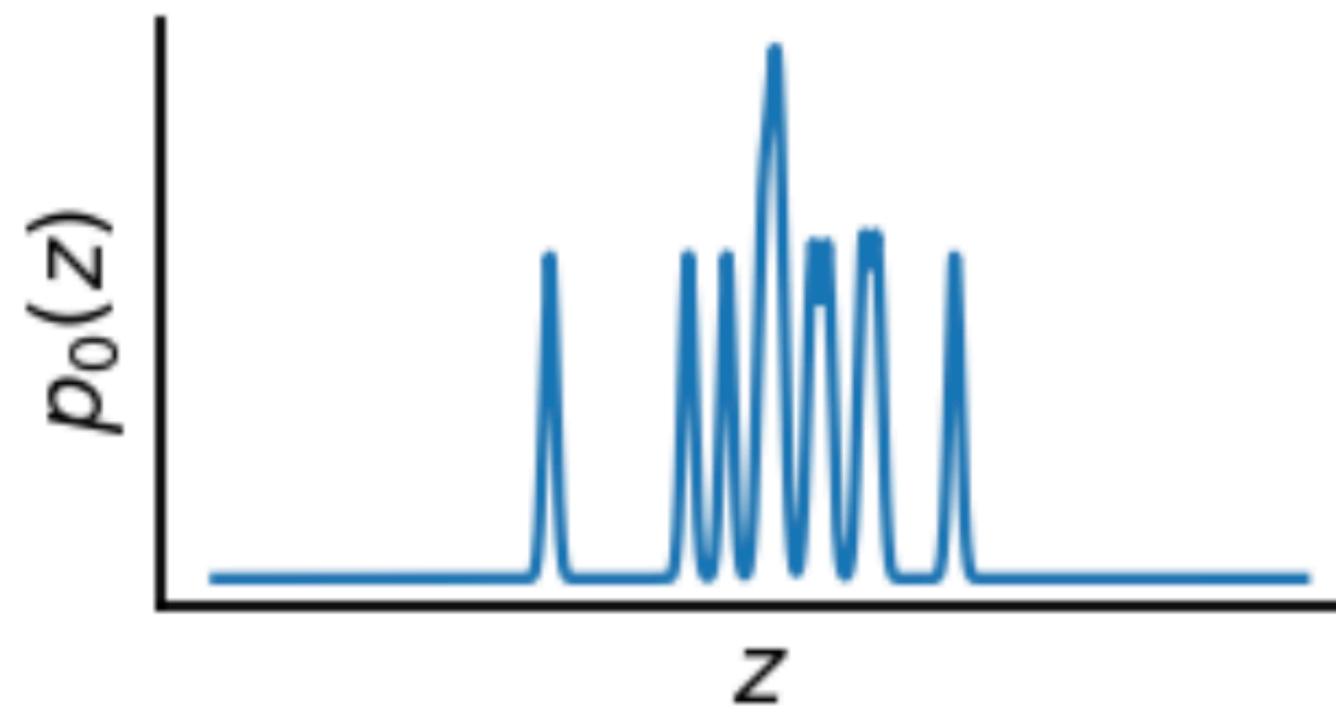


Learns more than just the training images!

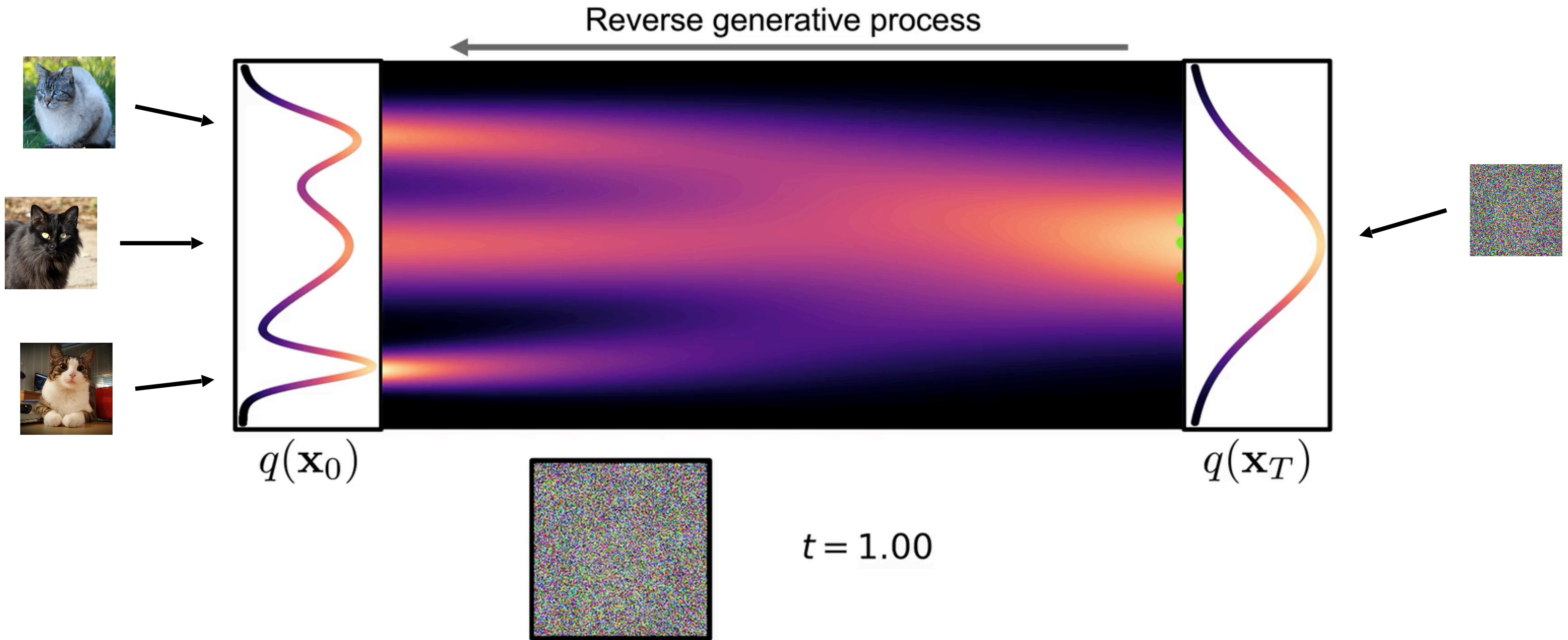




$p_0(z_0), \dots, p_t(z_t), \dots, p_T(z_T)$

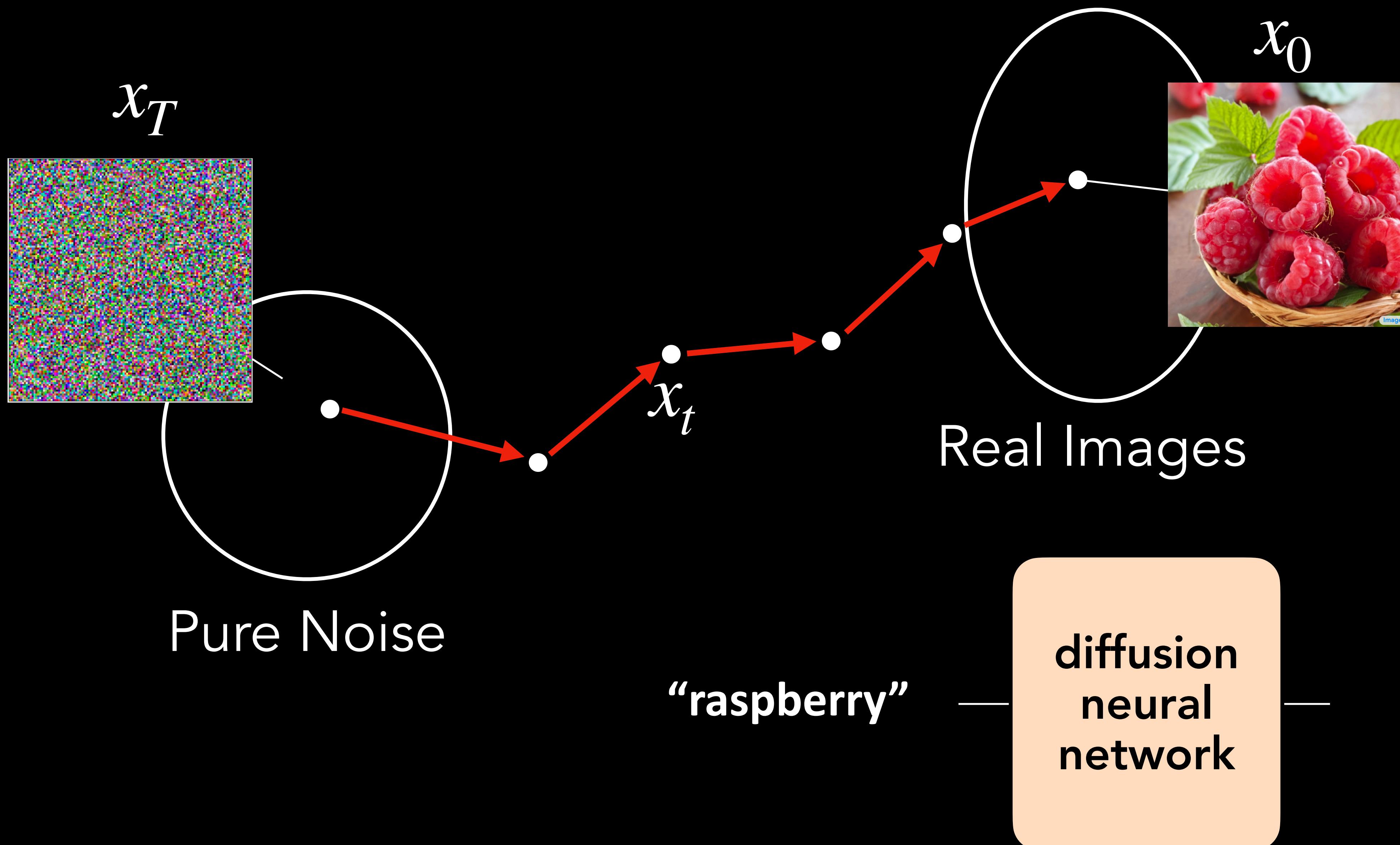


The Generative Reverse Stochastic Differential Equation



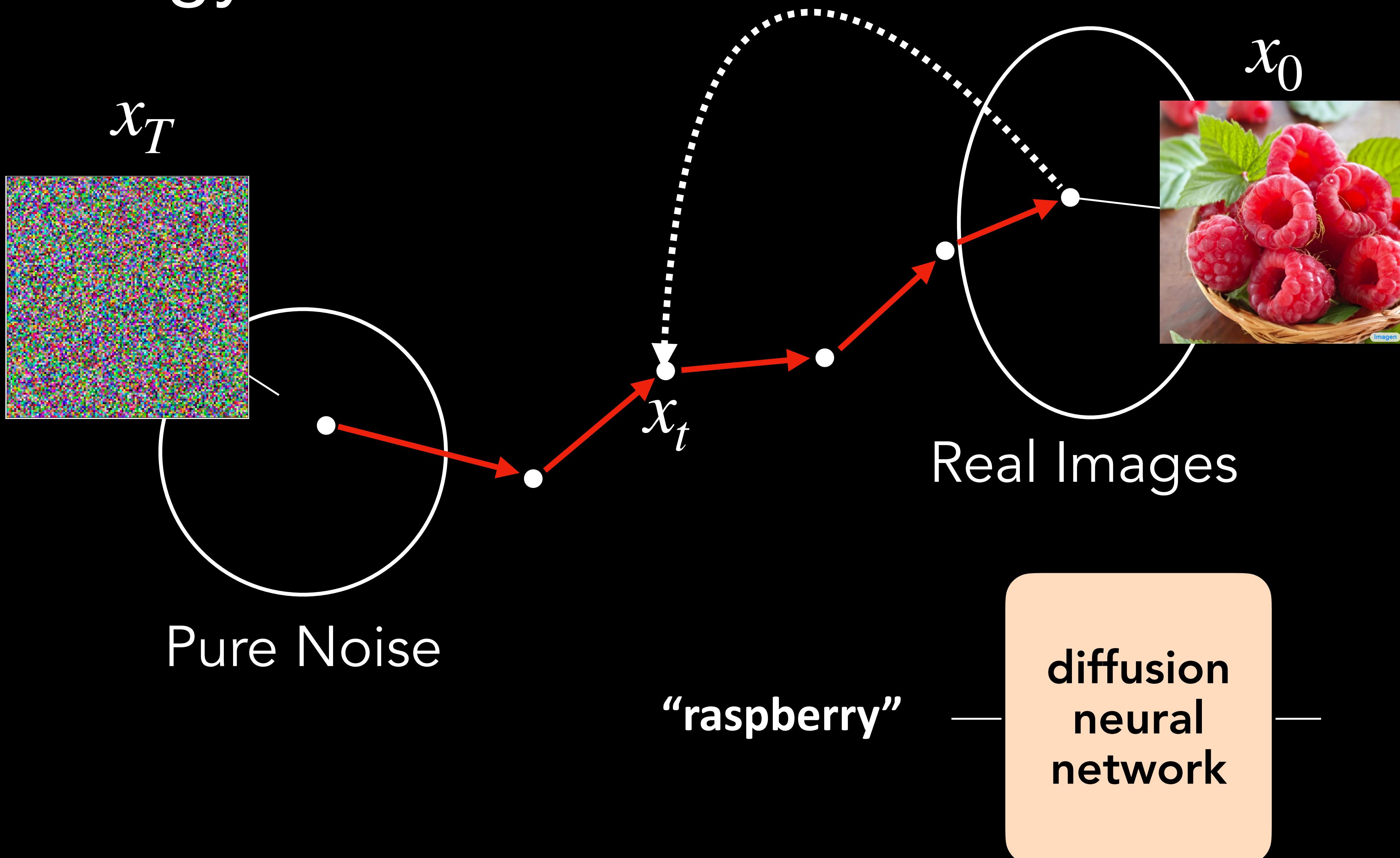
Part 3: How to train a diffusion model

Terminology

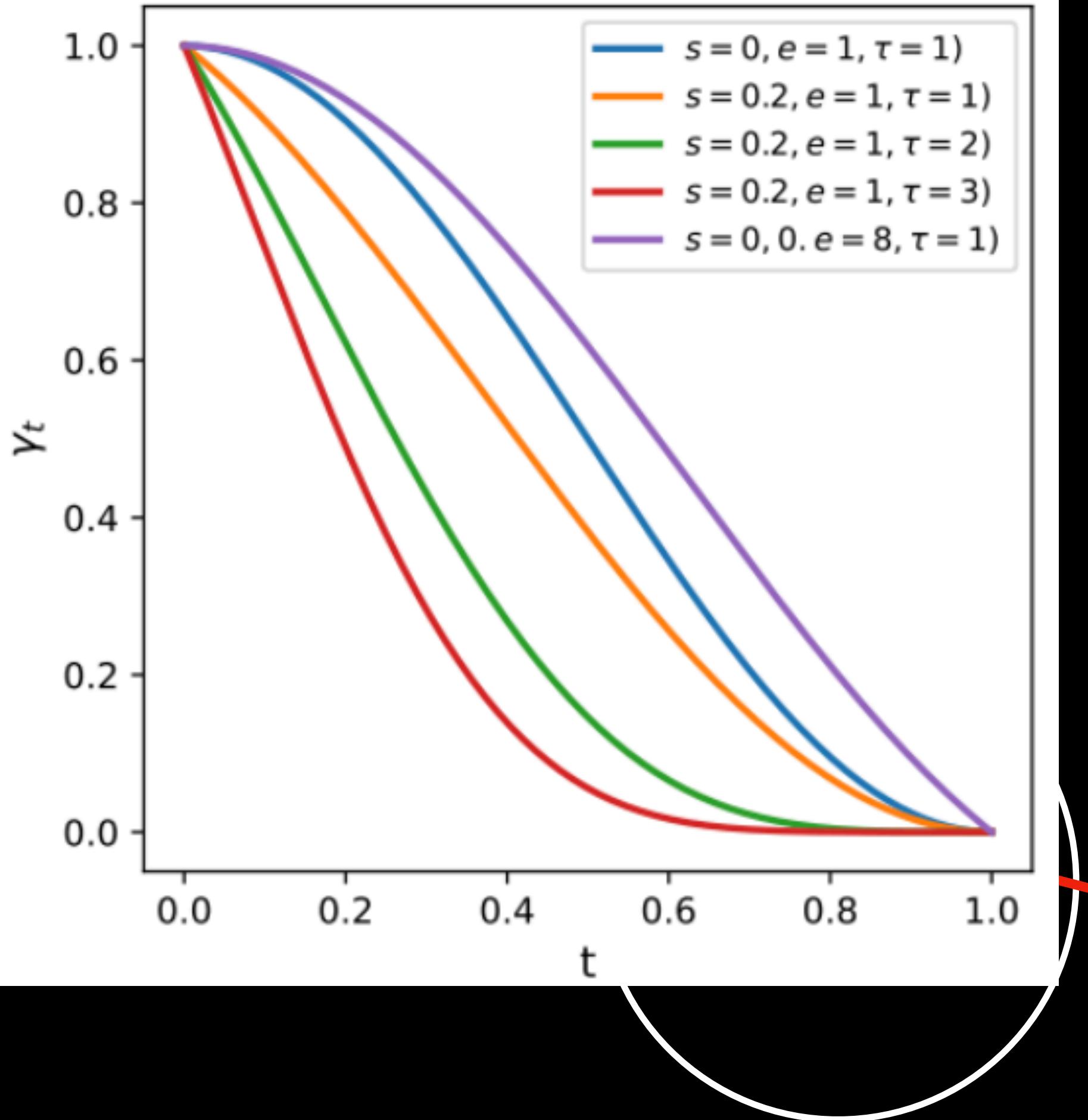


Terminology

$$\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$$

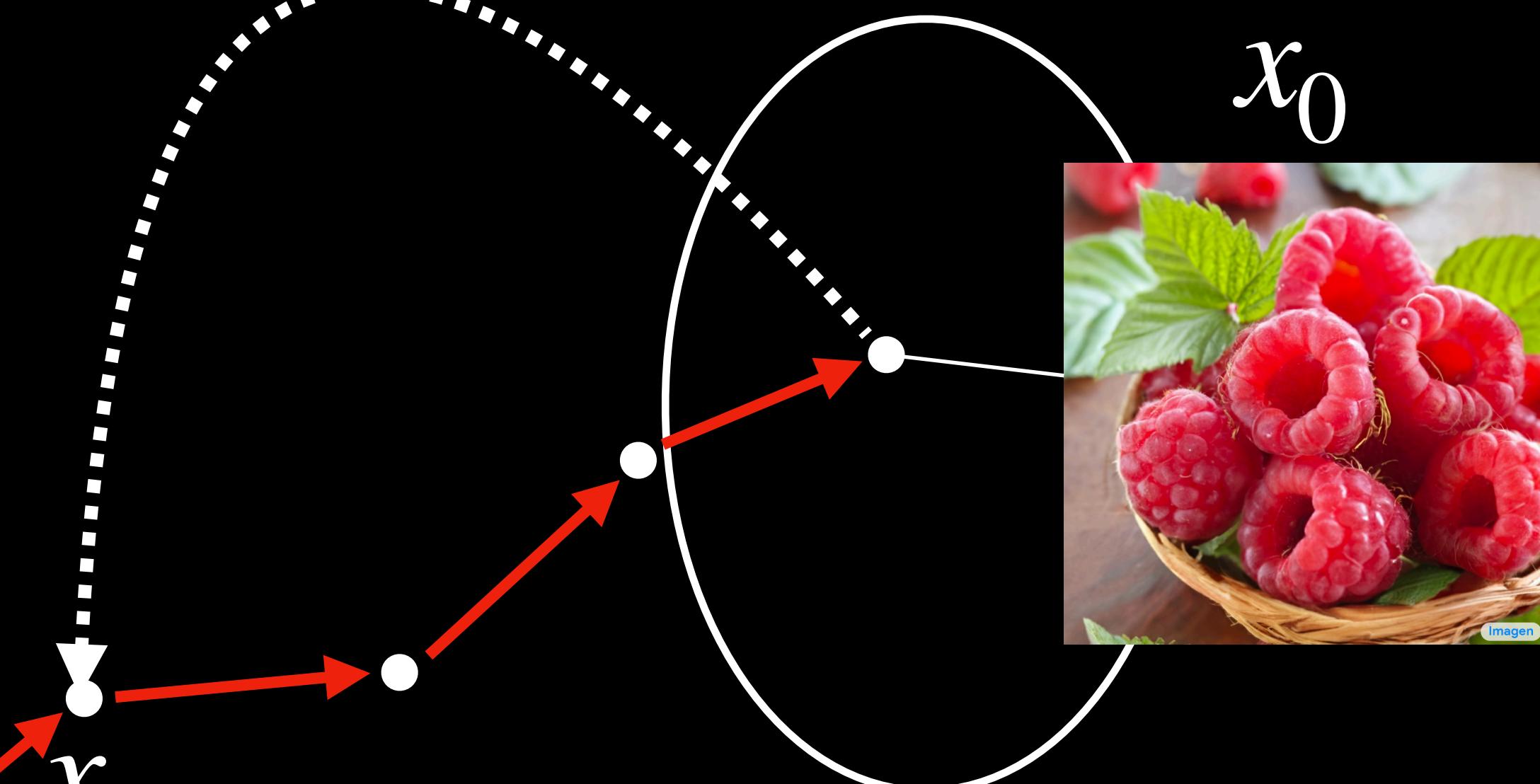


Cosine-based



Pure Noise

$$\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$$



"raspberry"

diffusion
neural
network

*training
algorithm*

repeat

$x_0 :=$ sample from dataset

$t :=$ sample from uniform distribution $[1:T]$

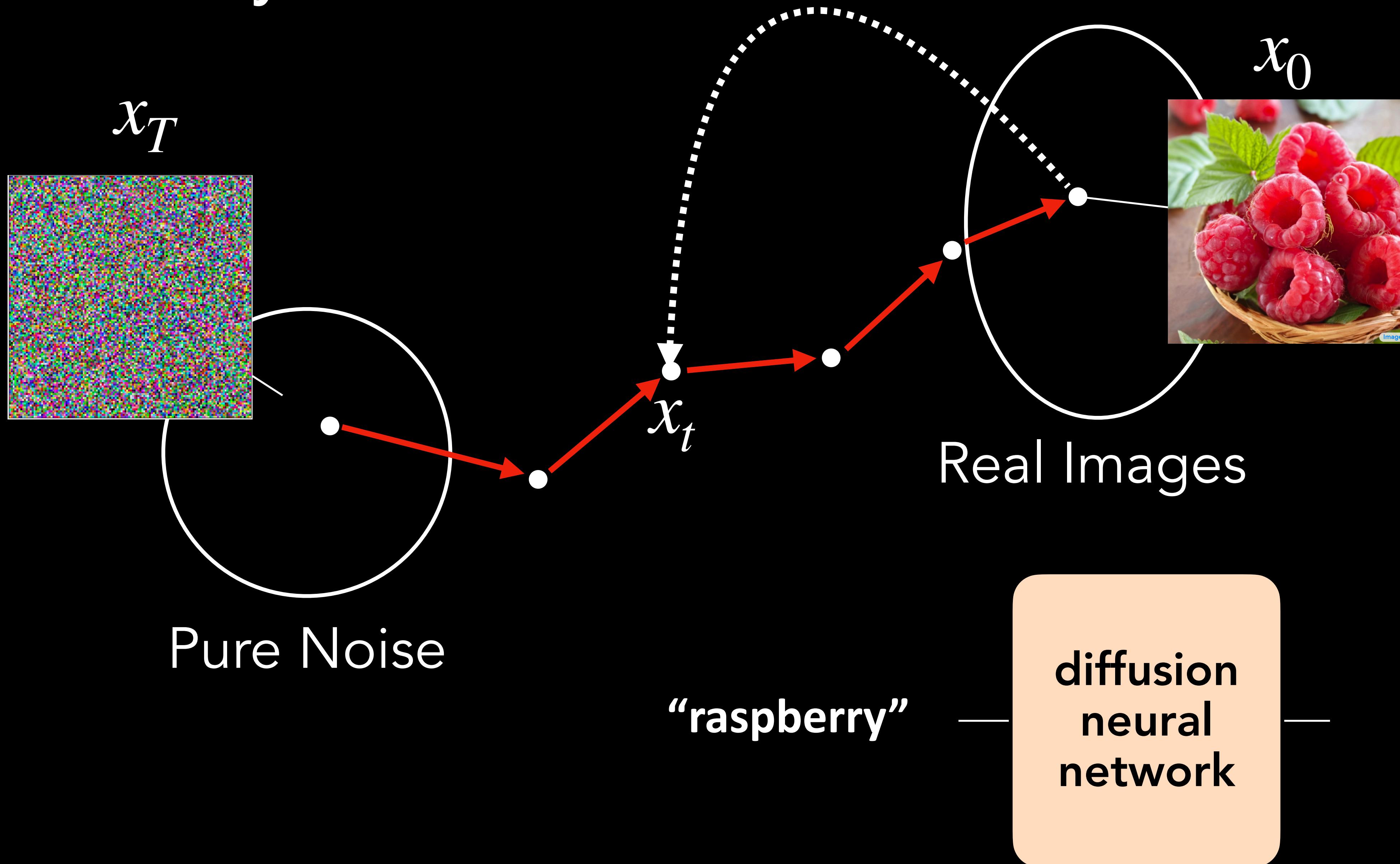
$\text{epsilon} :=$ sample from standardized normal

take gradient descent step

until converged

$$\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|$$

More visually...



What's important?

Architecture

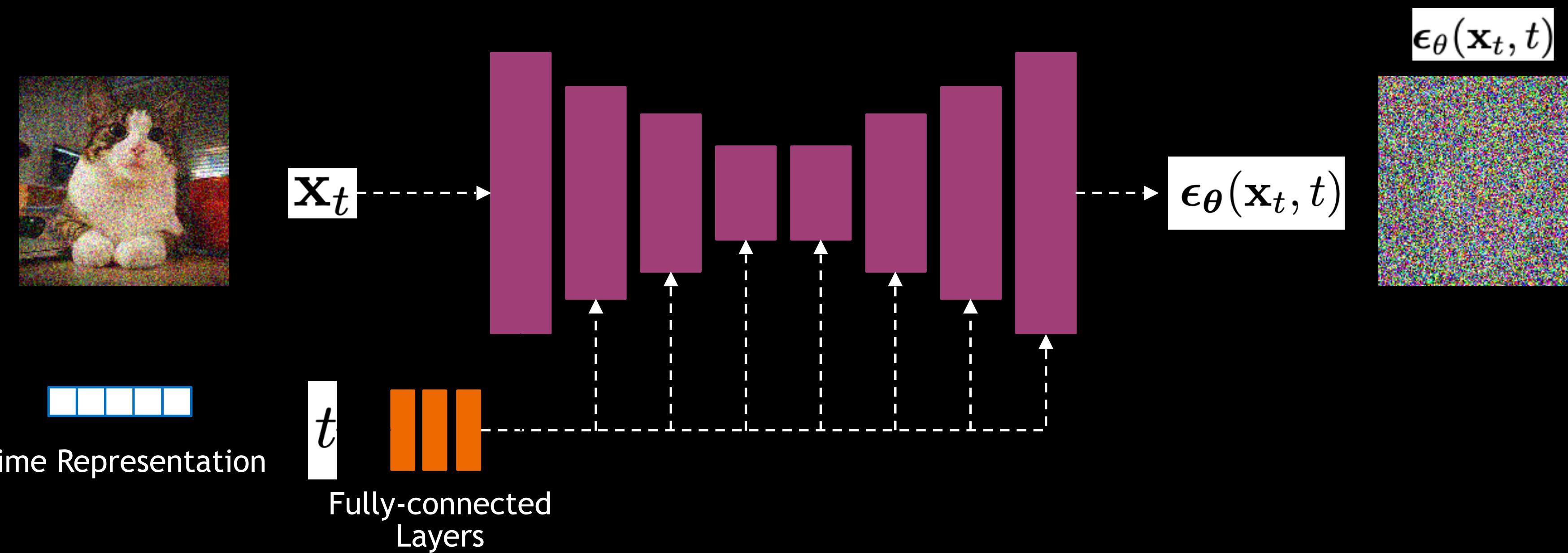
Capacity

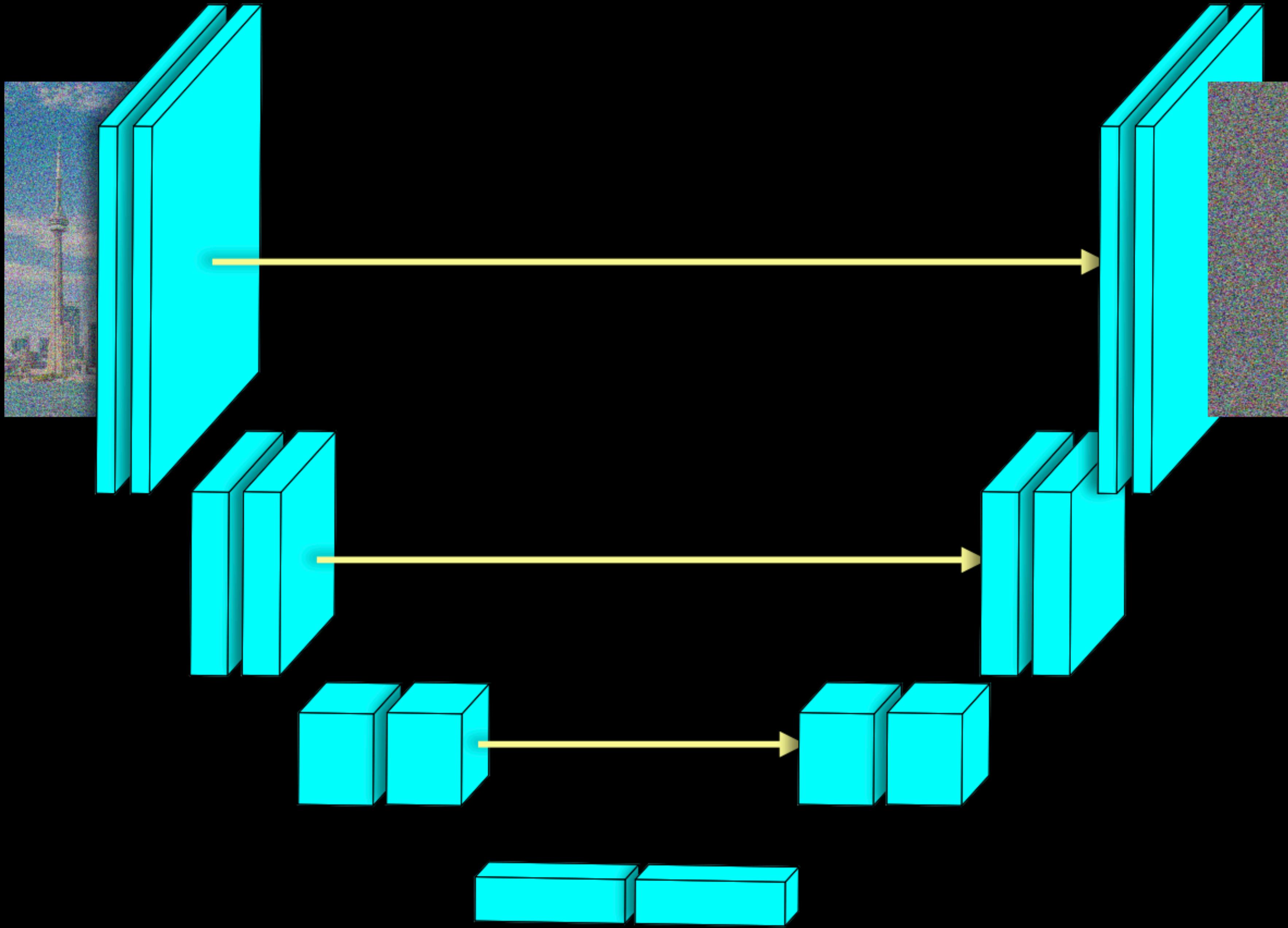
Data

Noise schedule

Conditioning

Architecture

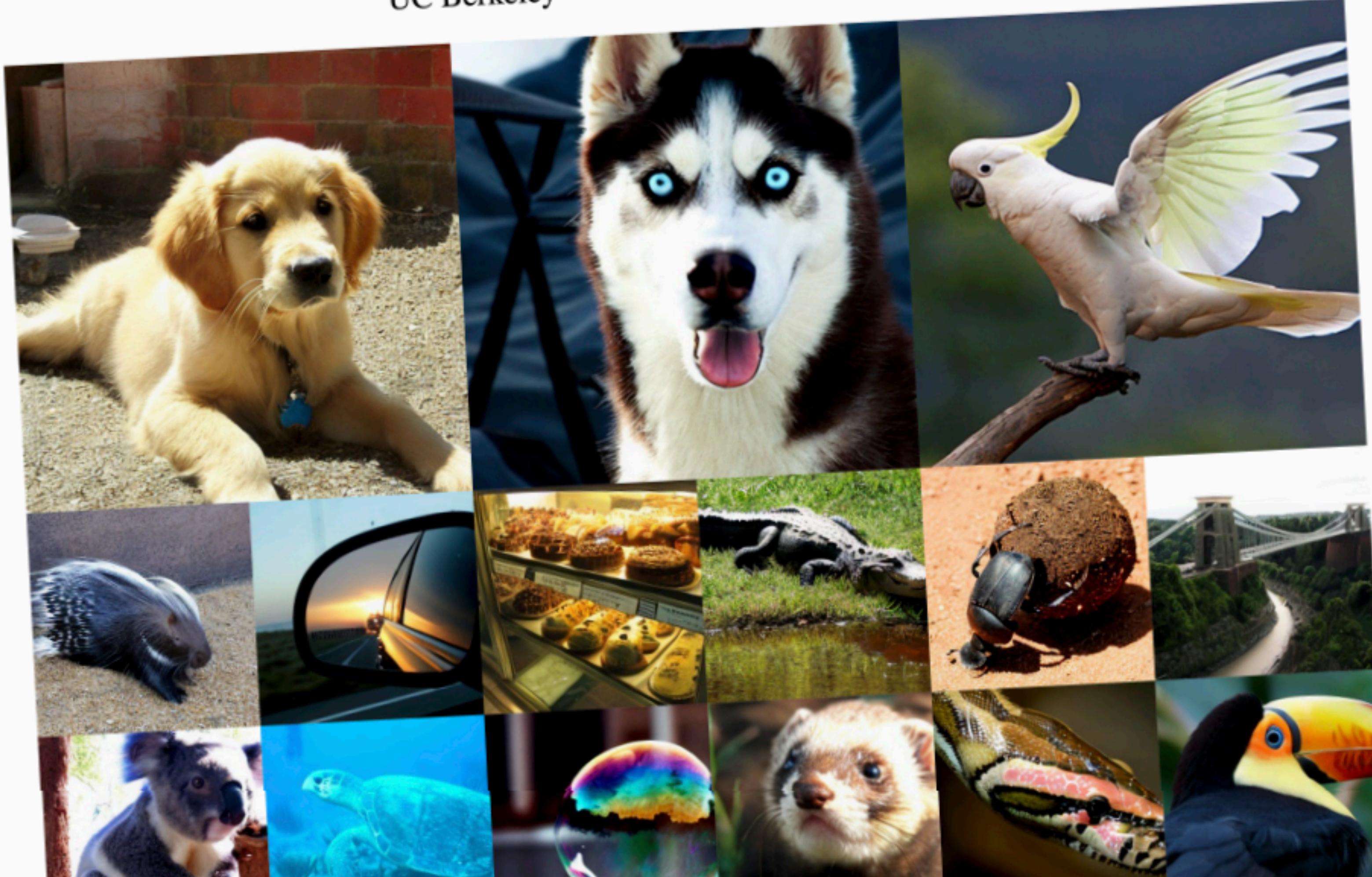




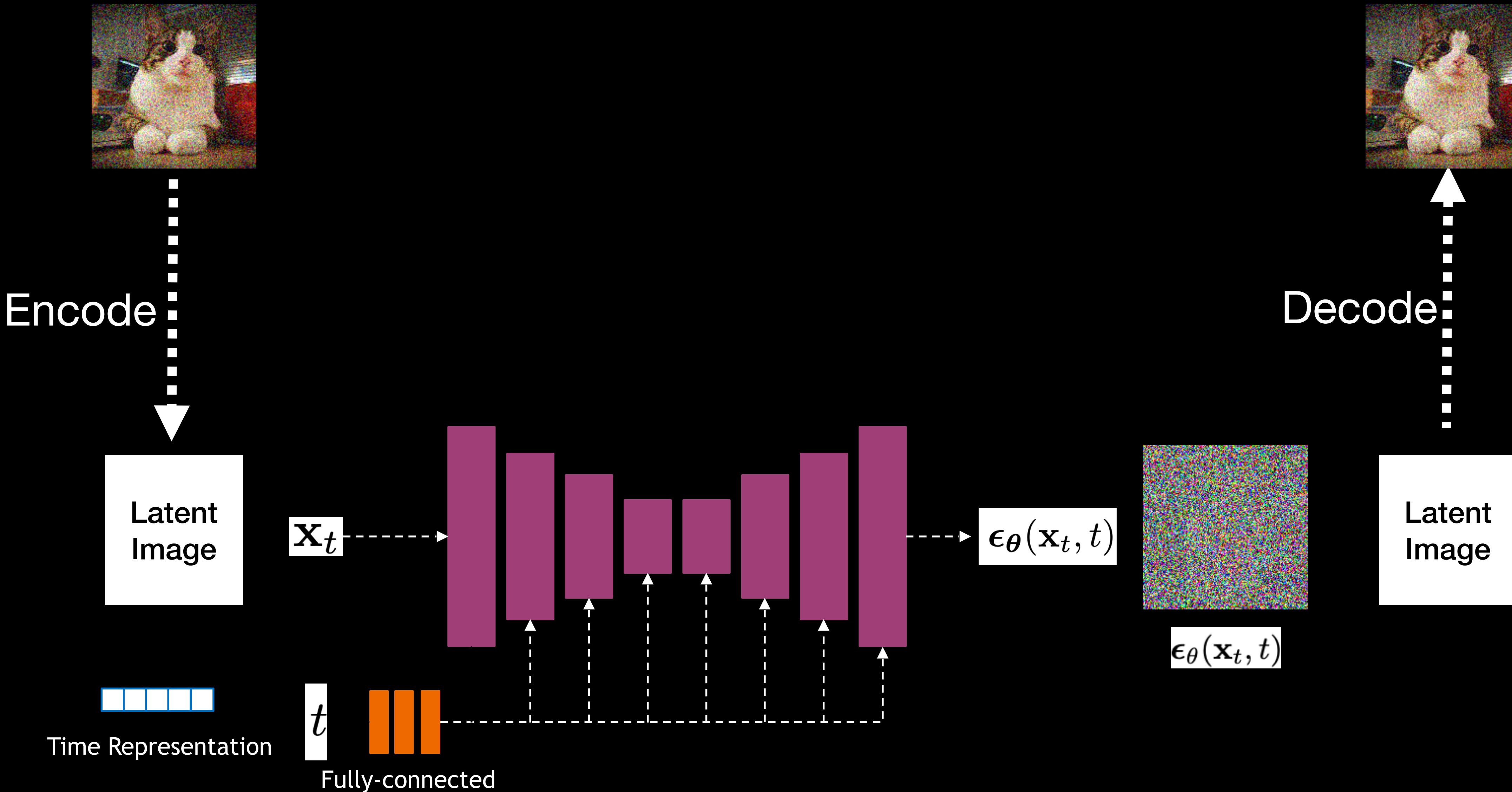
Scalable Diffusion Models with Transformers

William Peebles*
UC Berkeley

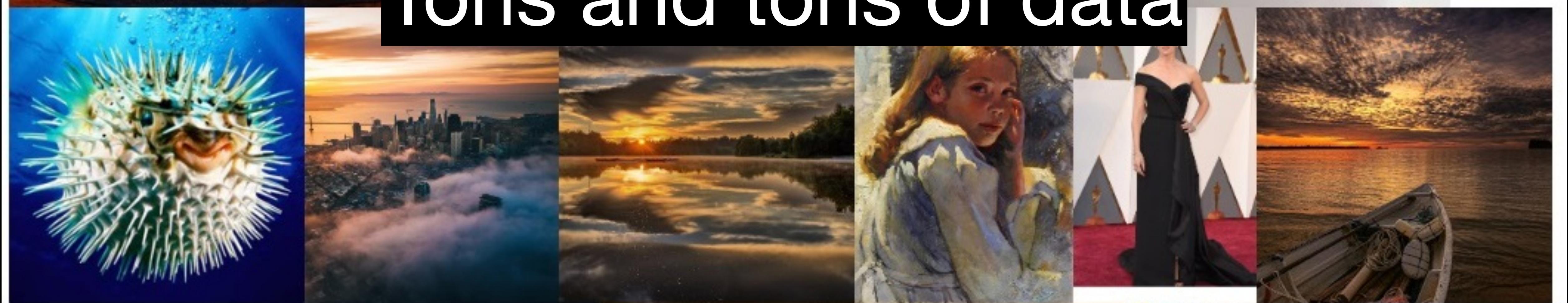
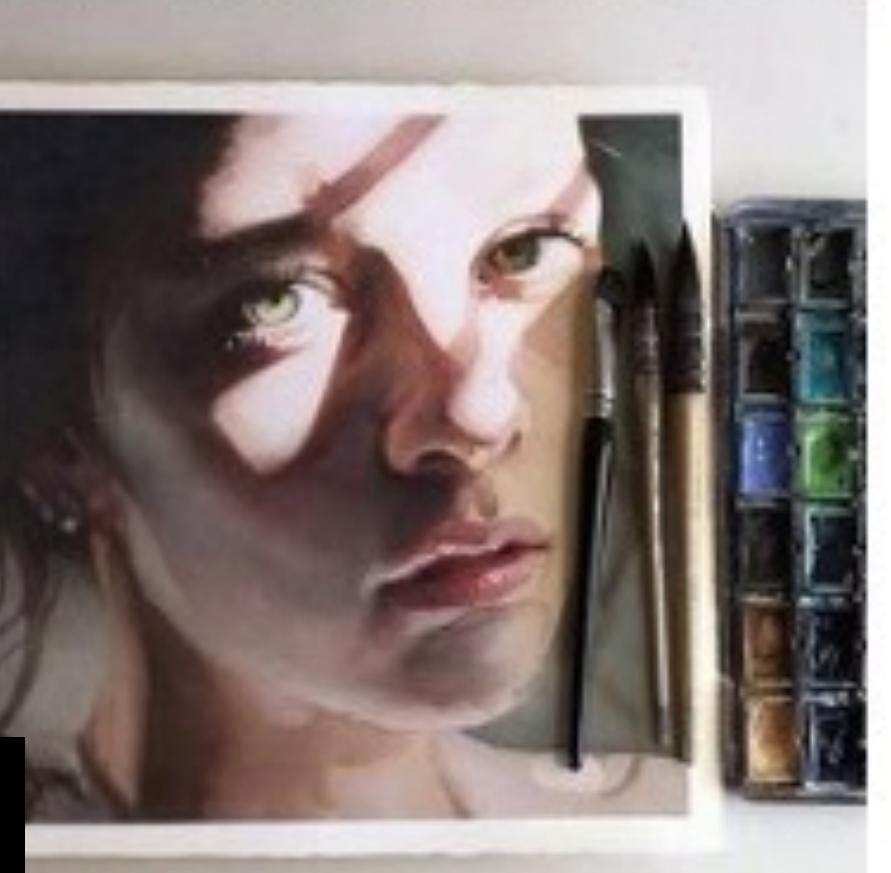
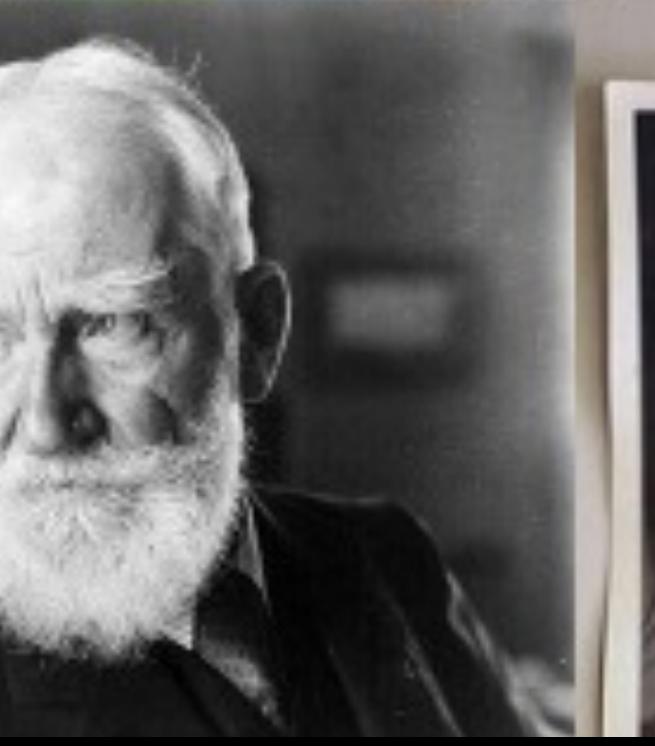
Saining Xie
New York University



Latent models



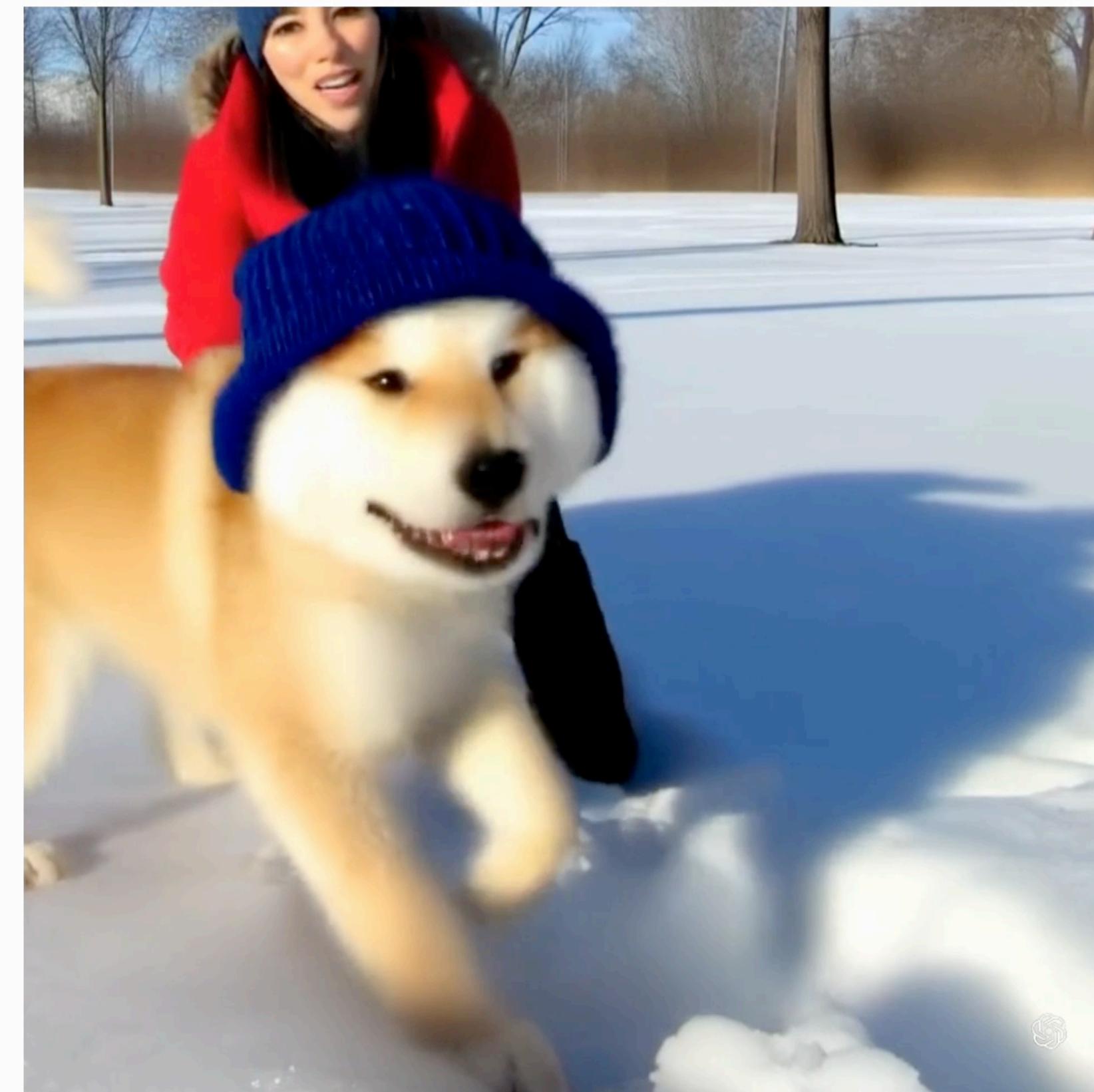
Tons and tons of data



Bigger is better



Base compute

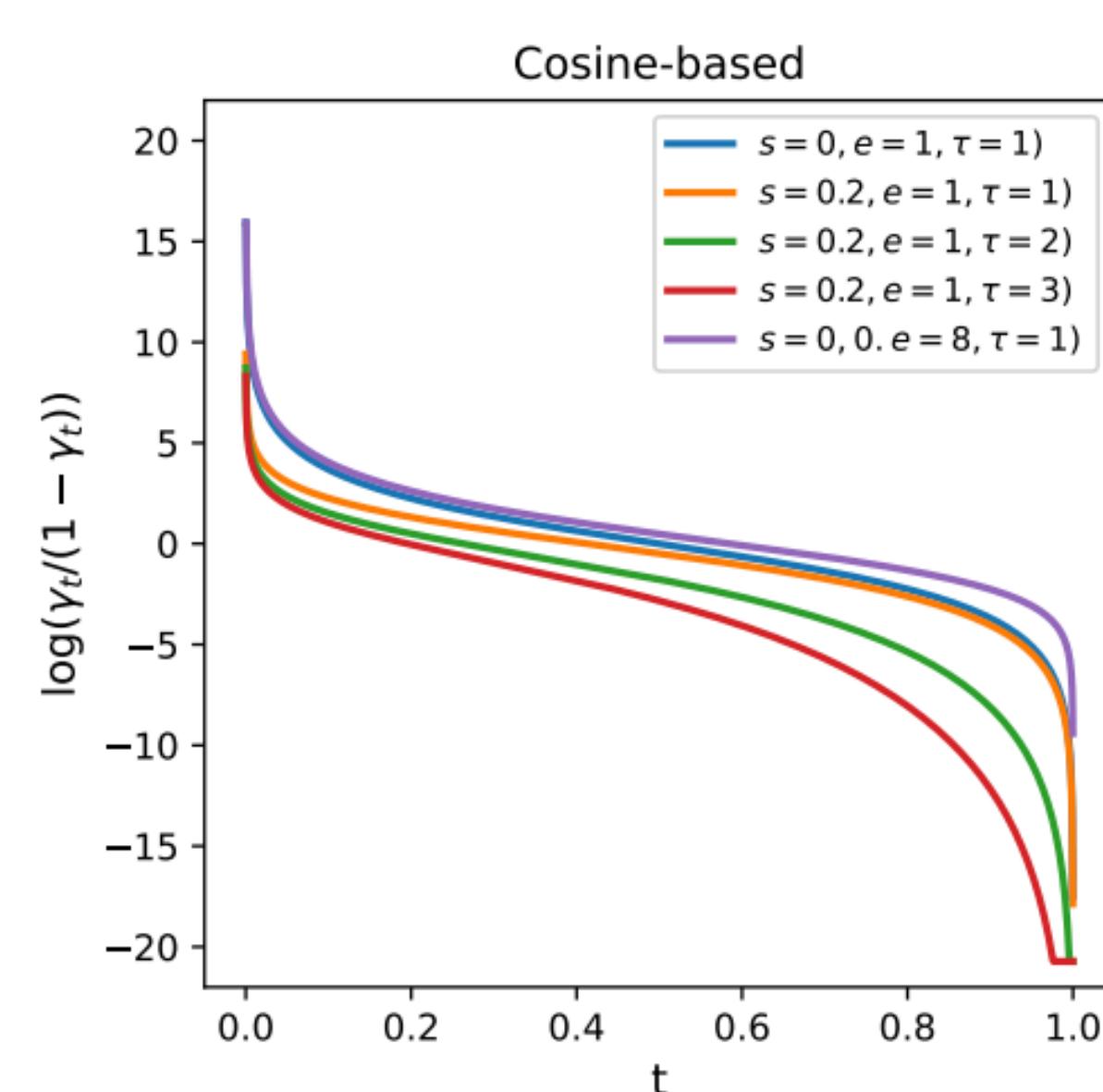
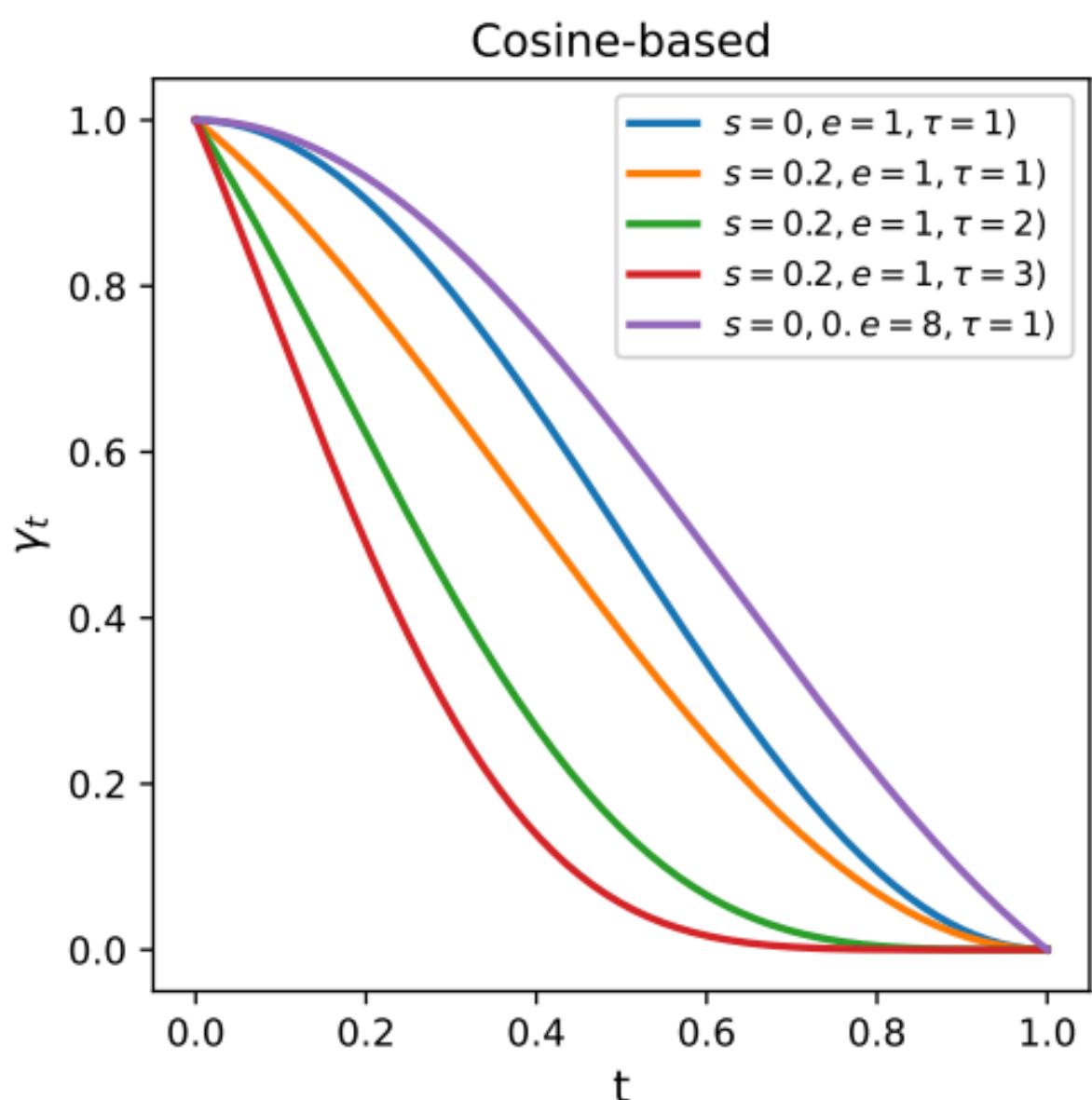


4x compute

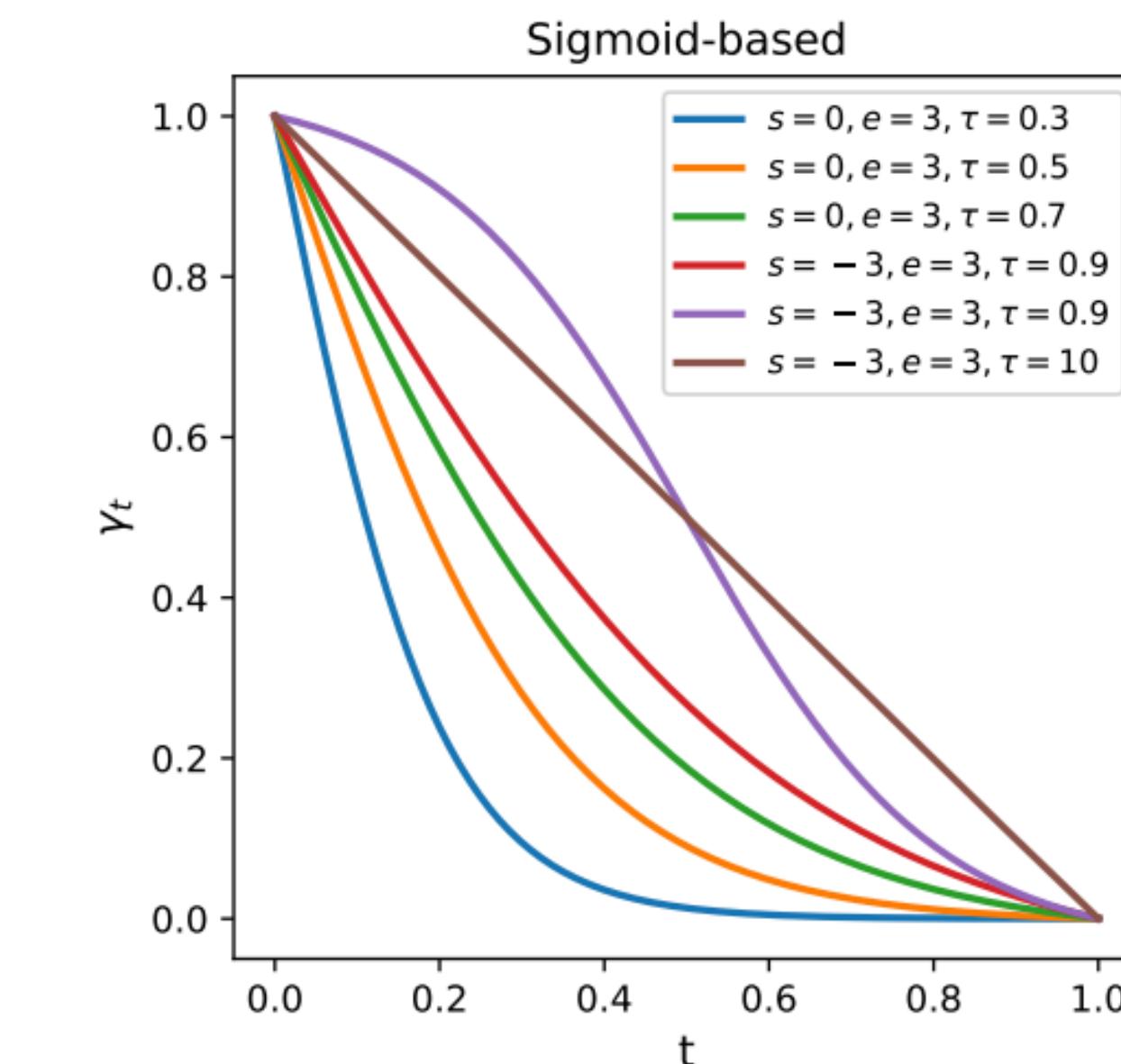


32x compute

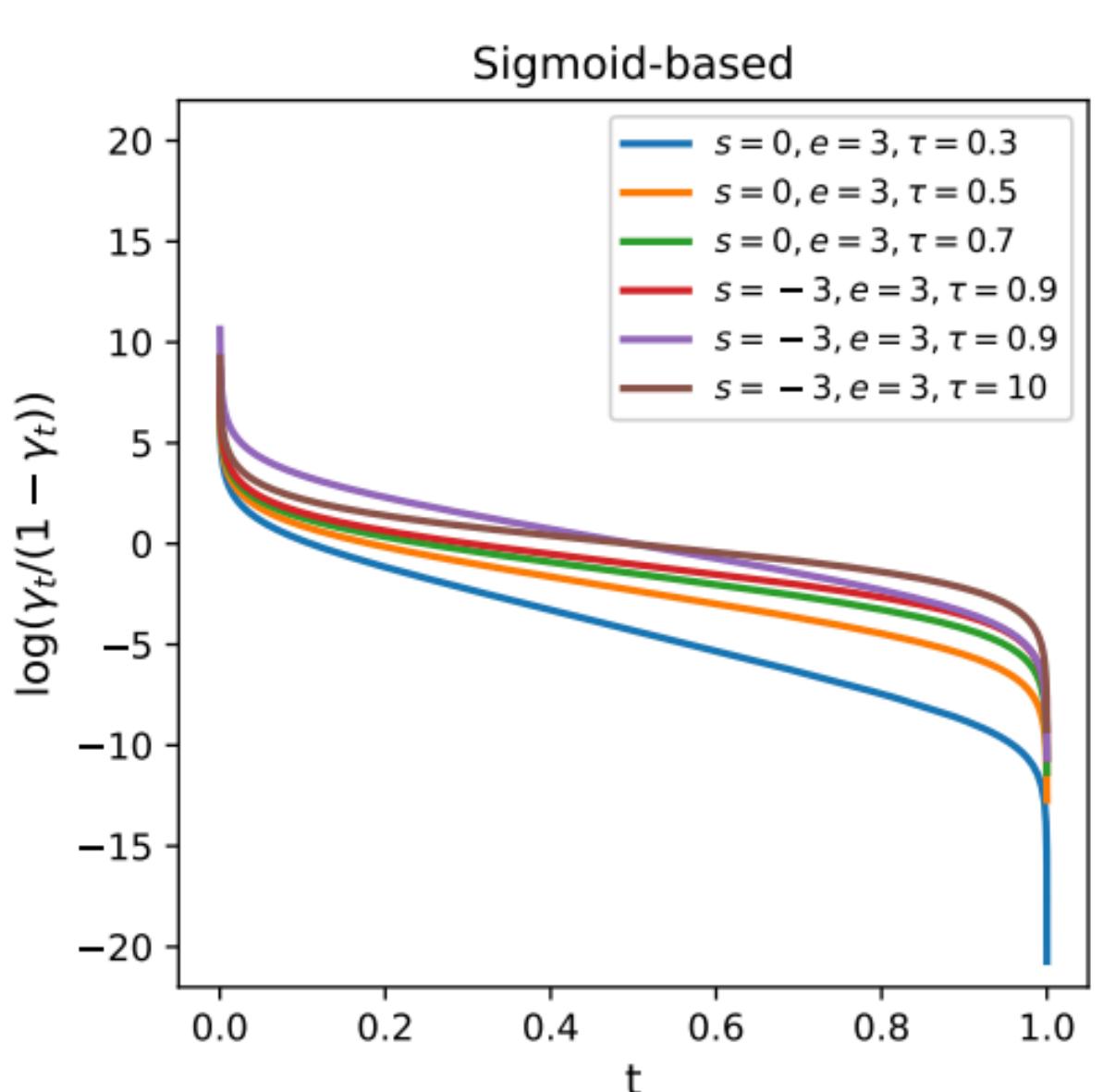
Noise schedules!



(a) Cosine



(c) Sigmoid



(d) Sigmoid (logSNR)

Noise schedules!



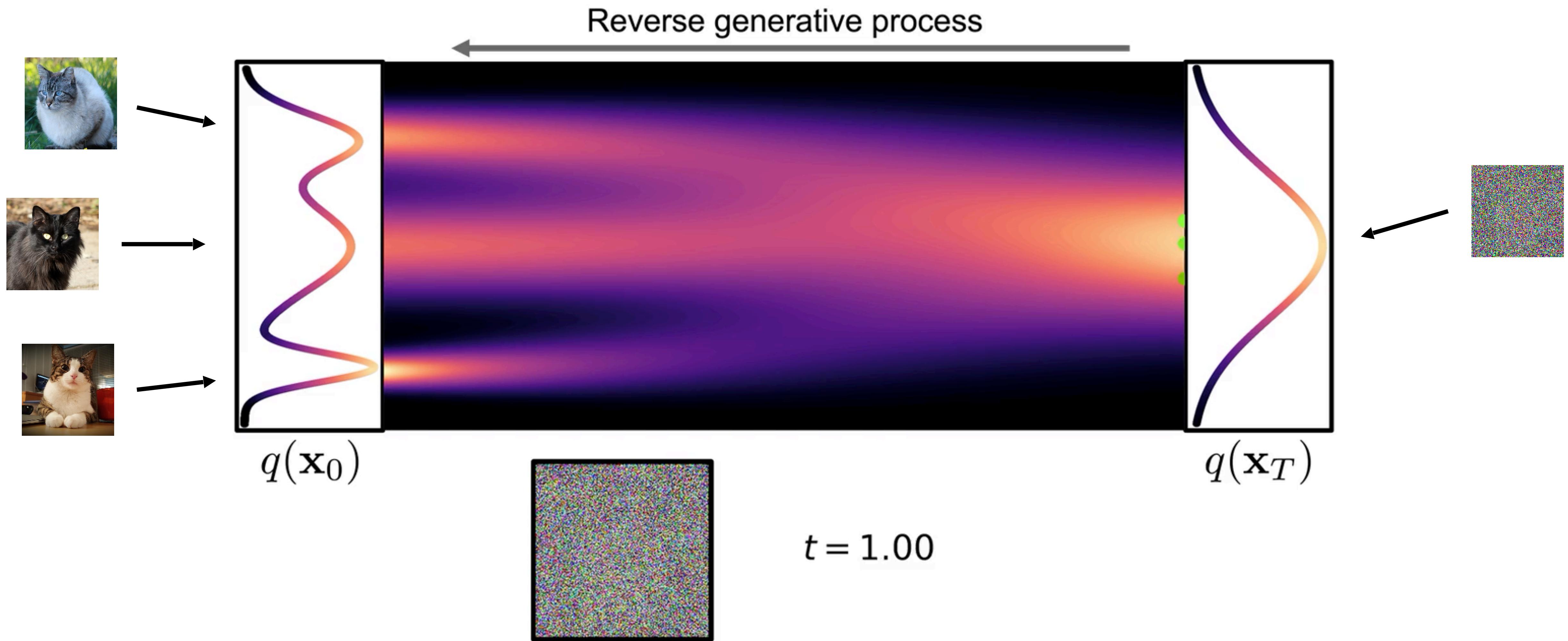
Figure 2: Noised images ($\mathbf{x}_t = \sqrt{\gamma}\mathbf{x}_0 + \sqrt{1-\gamma}\boldsymbol{\epsilon}$) with the same noise level ($\gamma = 0.7$). We see that higher resolution natural images tend to exhibit higher degree of redundancy in (nearby) pixels, therefore less information is destroyed with the same level of independent noise.

Training checklist

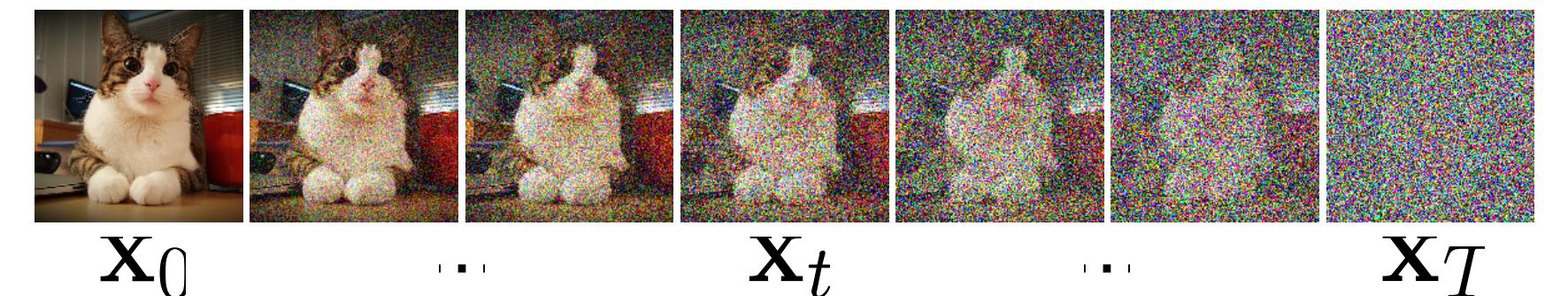
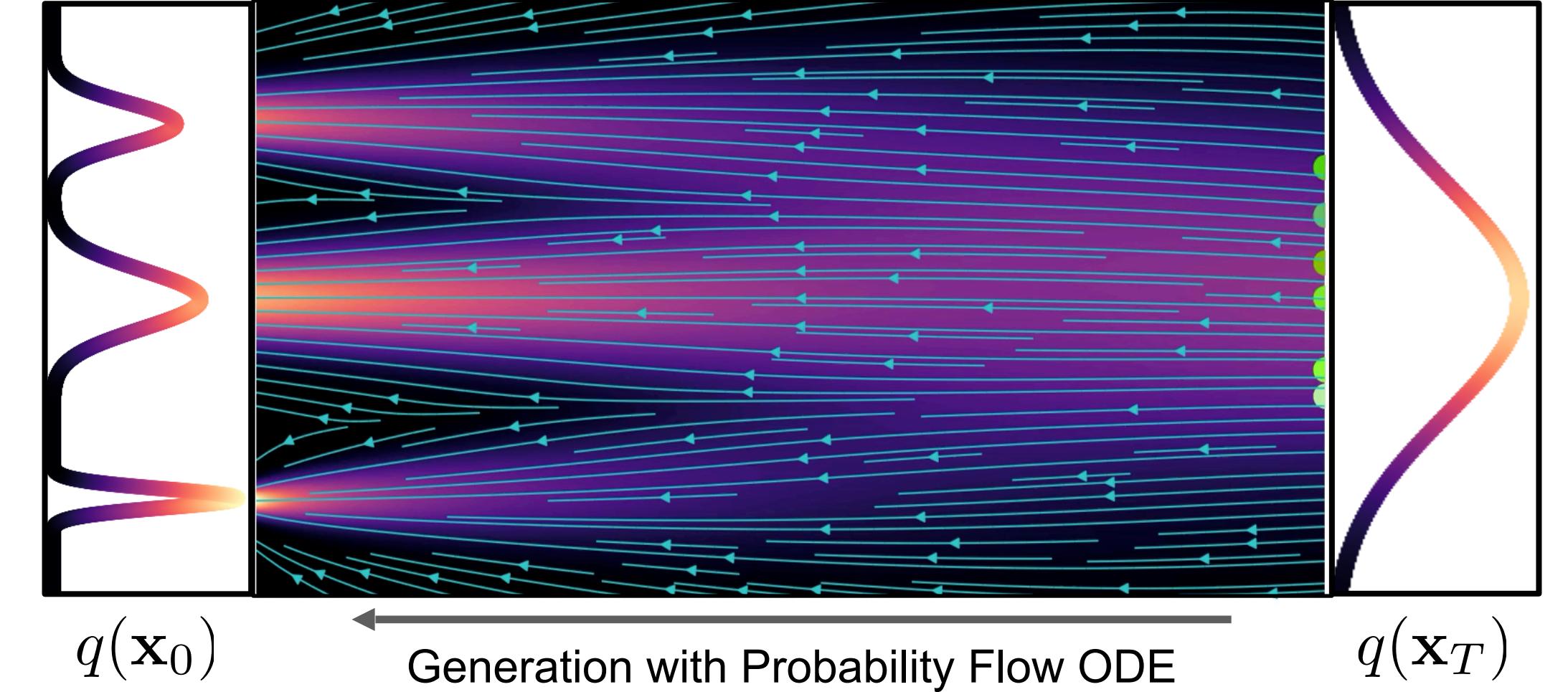
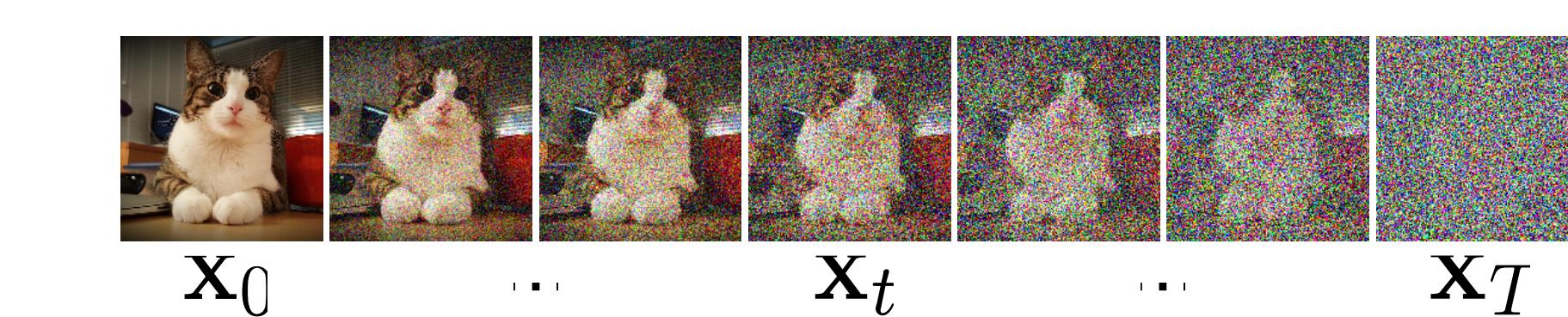
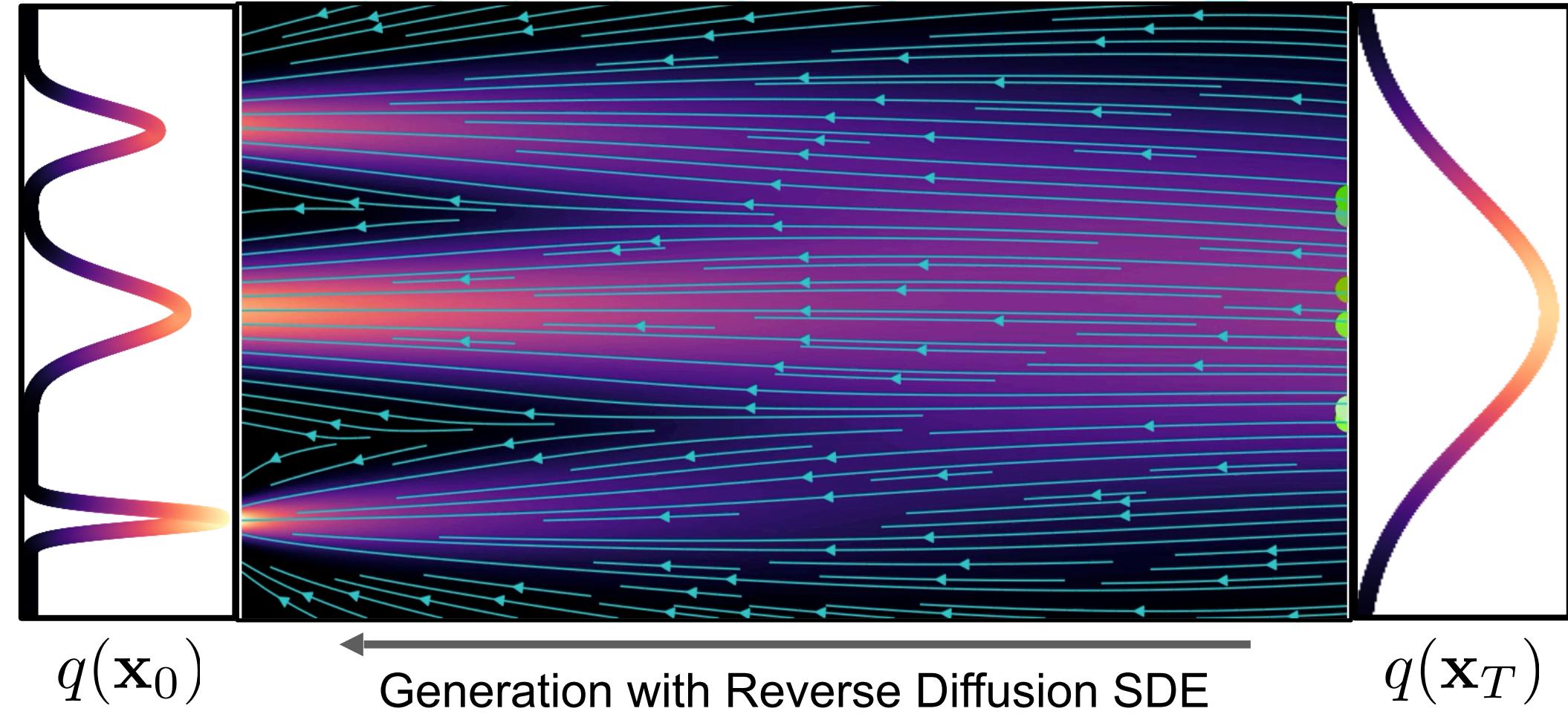
- Decide on a data distribution to model
 - Download or create a dataset
- Come up with a noise schedule
- Pick an architecture + conditioning signal
- Latent / not latent?
- Get access to a ton of compute
- Start training!

Part 4: How do you use a diffusion model?

The visualization of sampling we saw before



Two ways of sampling—stochastic and deterministic (SDE vs ODE)



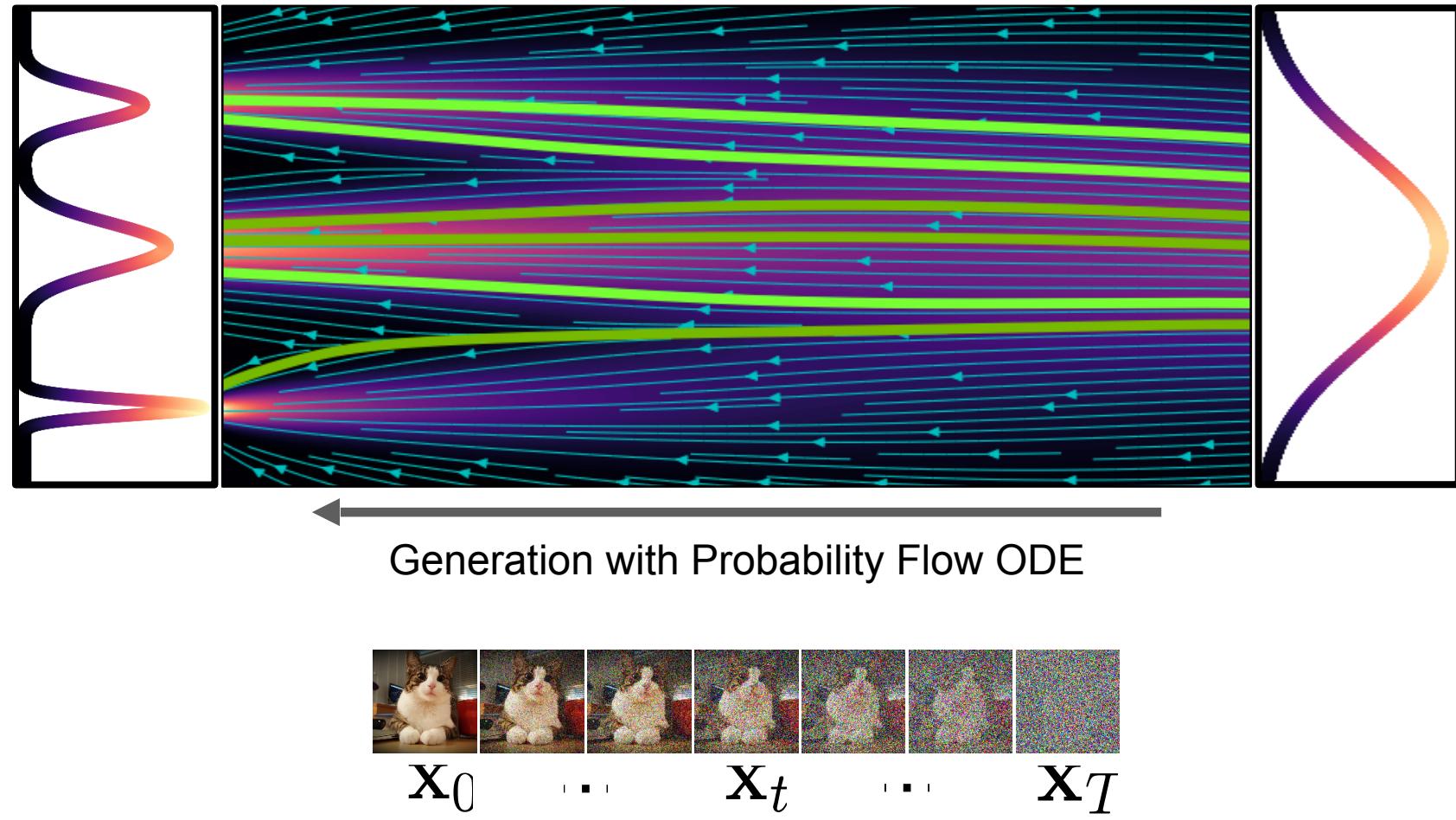
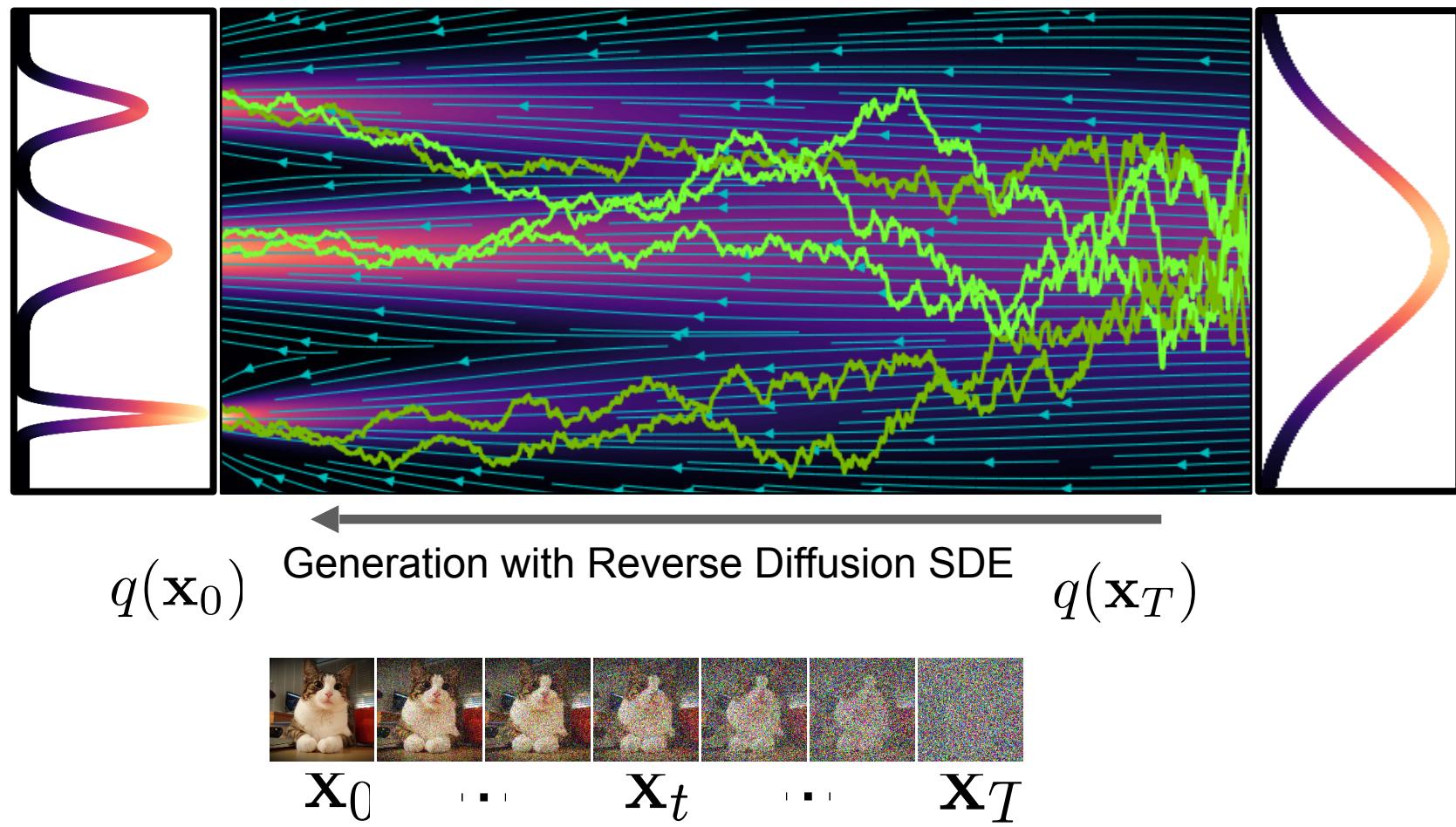
- **Generative Reverse Diffusion SDE (stochastic):**

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t) [\mathbf{x}_t + 2\mathbf{s}_{\theta}(\mathbf{x}_t, t)] dt + \sqrt{\beta(t)} d\bar{\omega}_t$$

- **Generative Probability Flow ODE (deterministic):**

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t) [\mathbf{x}_t + \mathbf{s}_{\theta}(\mathbf{x}_t, t)] dt$$

Pros and cons



SDE vs. ODE Sampling: Pro's and Con's

- **Pro:** Continuous noise injection can help to compensate errors during diffusion process (Langevin sampling actively pushes towards correct distribution).
- **Con:** Often slower, because the stochastic terms themselves require fine discretization during solve.

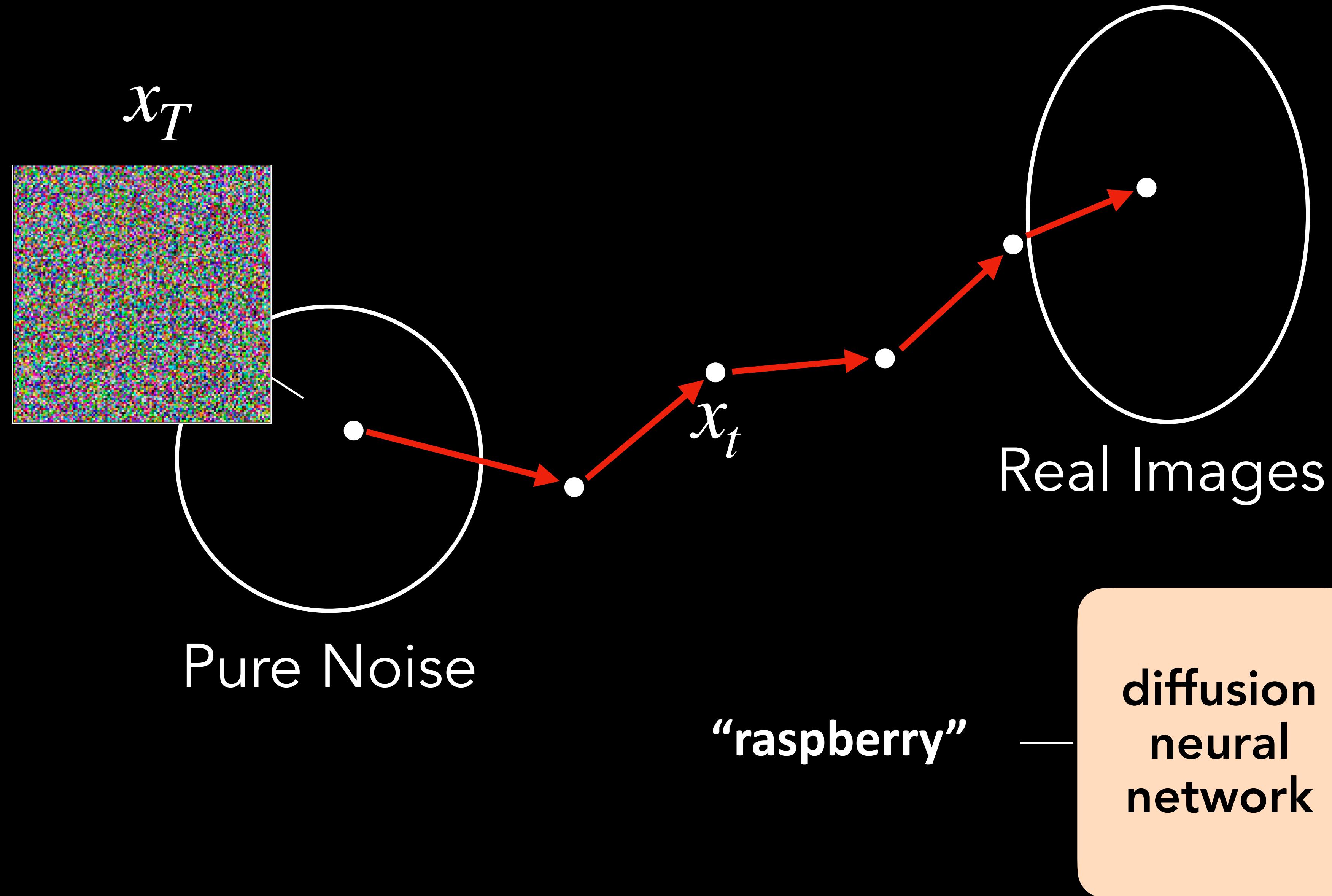
- **Pro:** Can leverage fast ODE solvers. Best when targeting very fast sampling.
- **Con:** No “stochastic” error correction, often slightly lower performance than stochastic sampling.

One common issue:

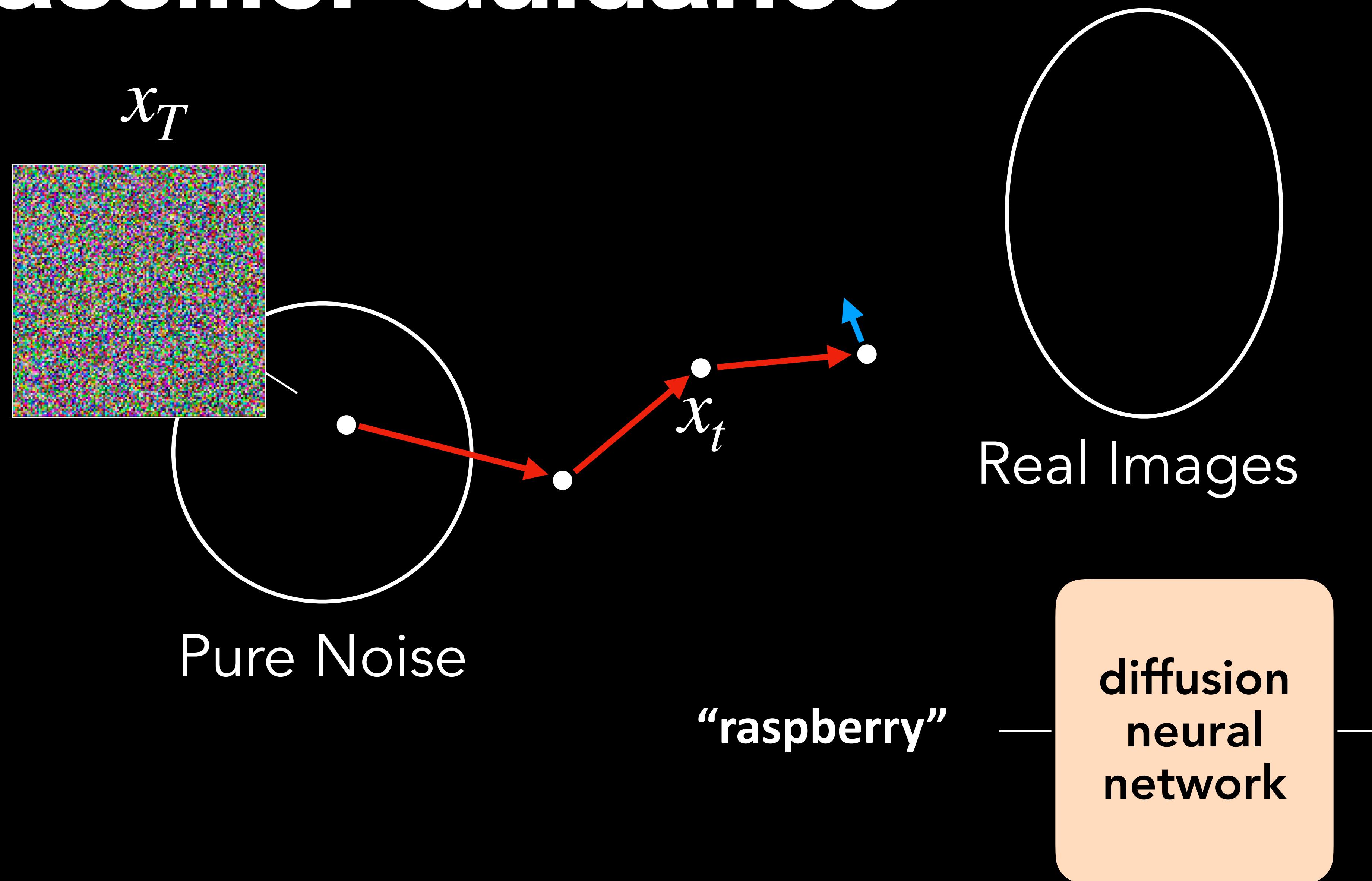
Trained models may ignore conditioning signals!

A solution:

Guidance!



Classifier Guidance



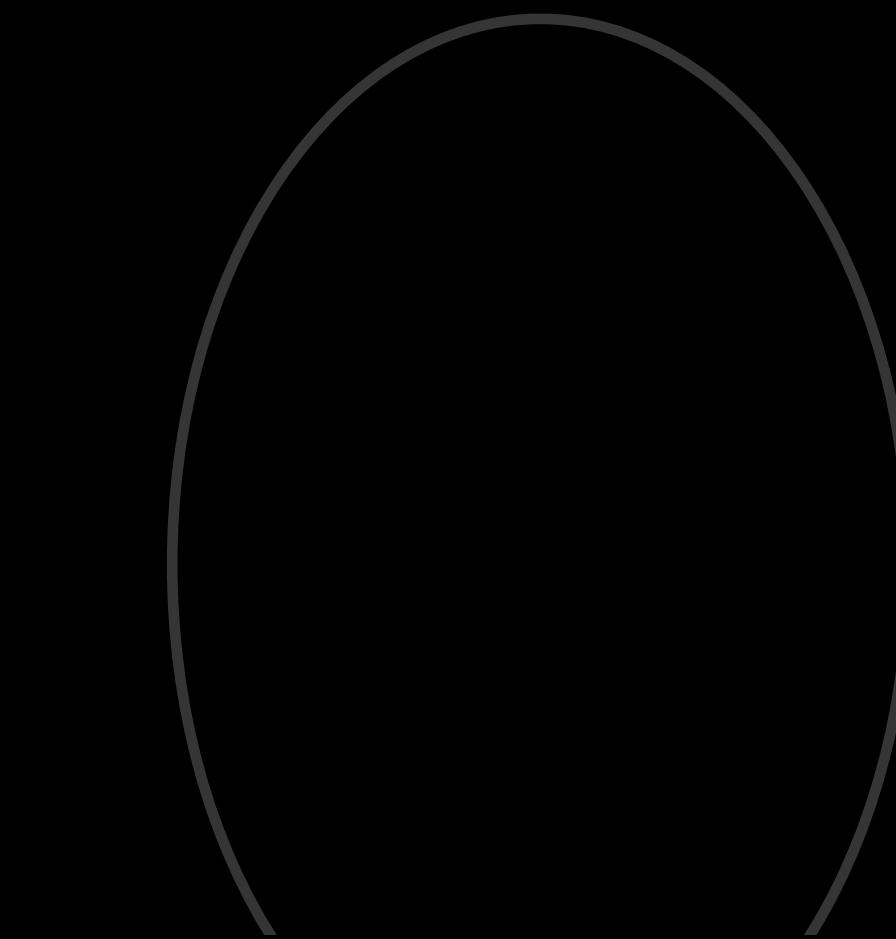
Classifier Guidance

x_T



Pure Noise

“raspberry”

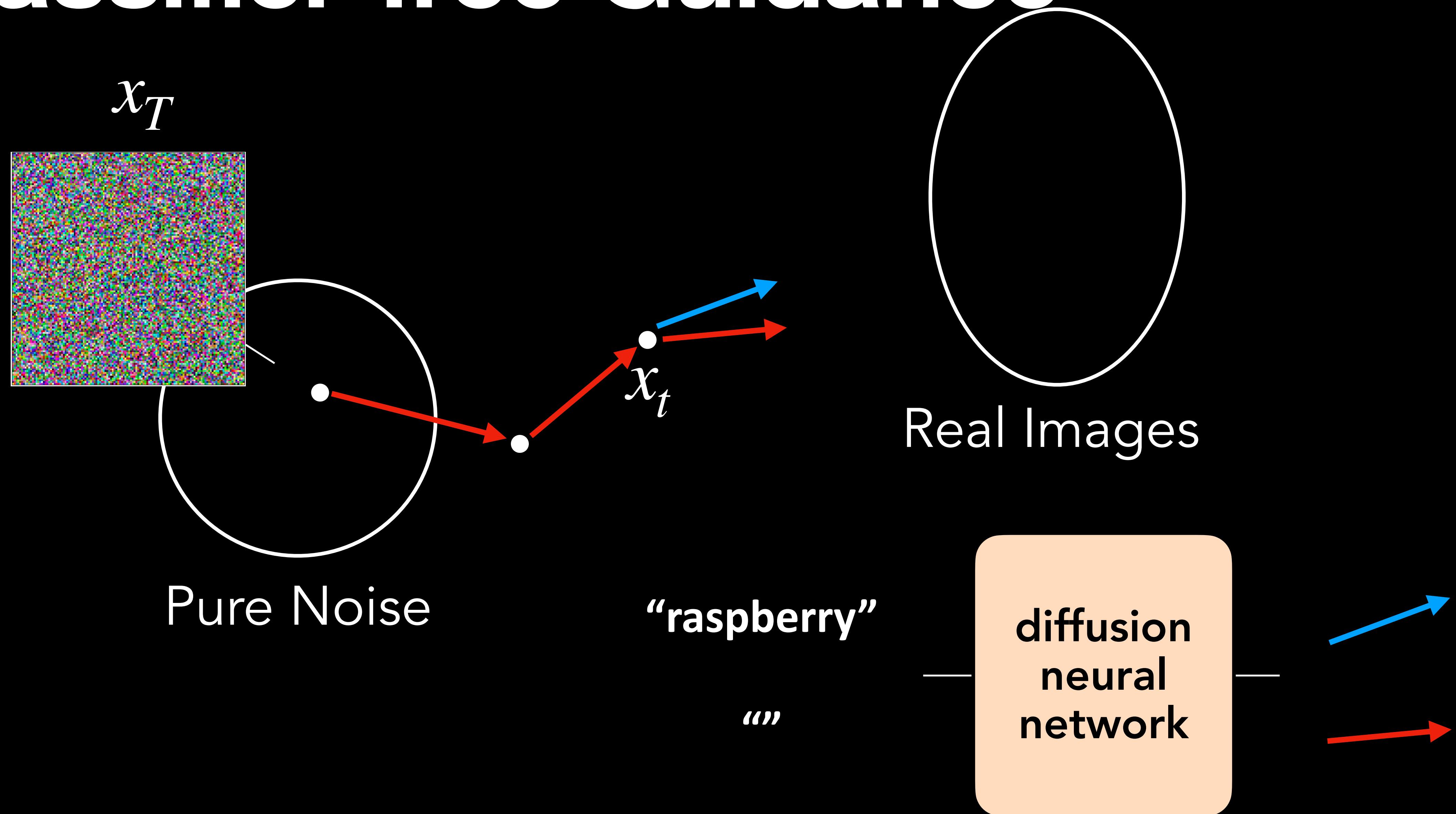


Real Images

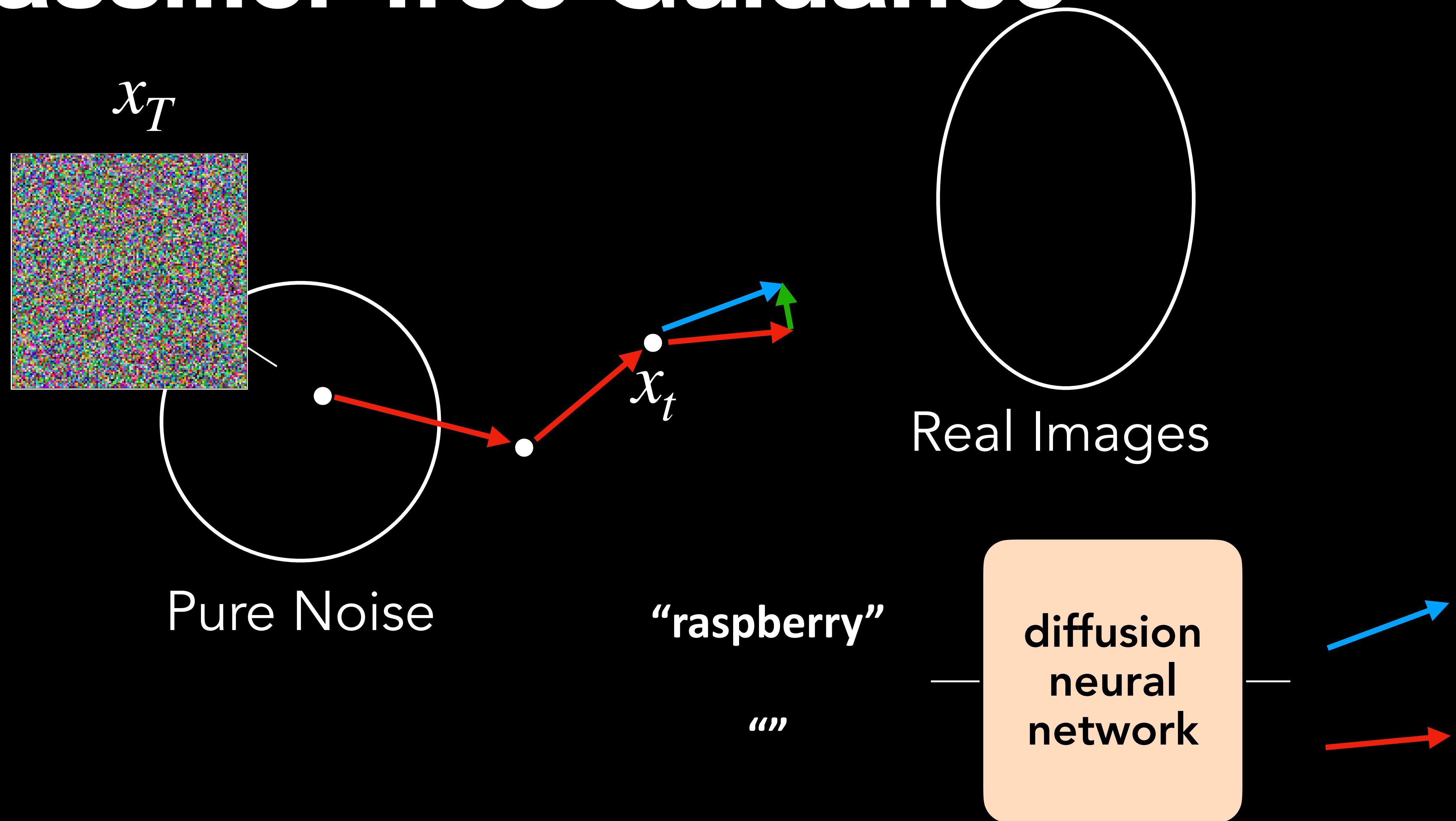
diffusion
neural
network

Need to train a classifier for each conditioning input!

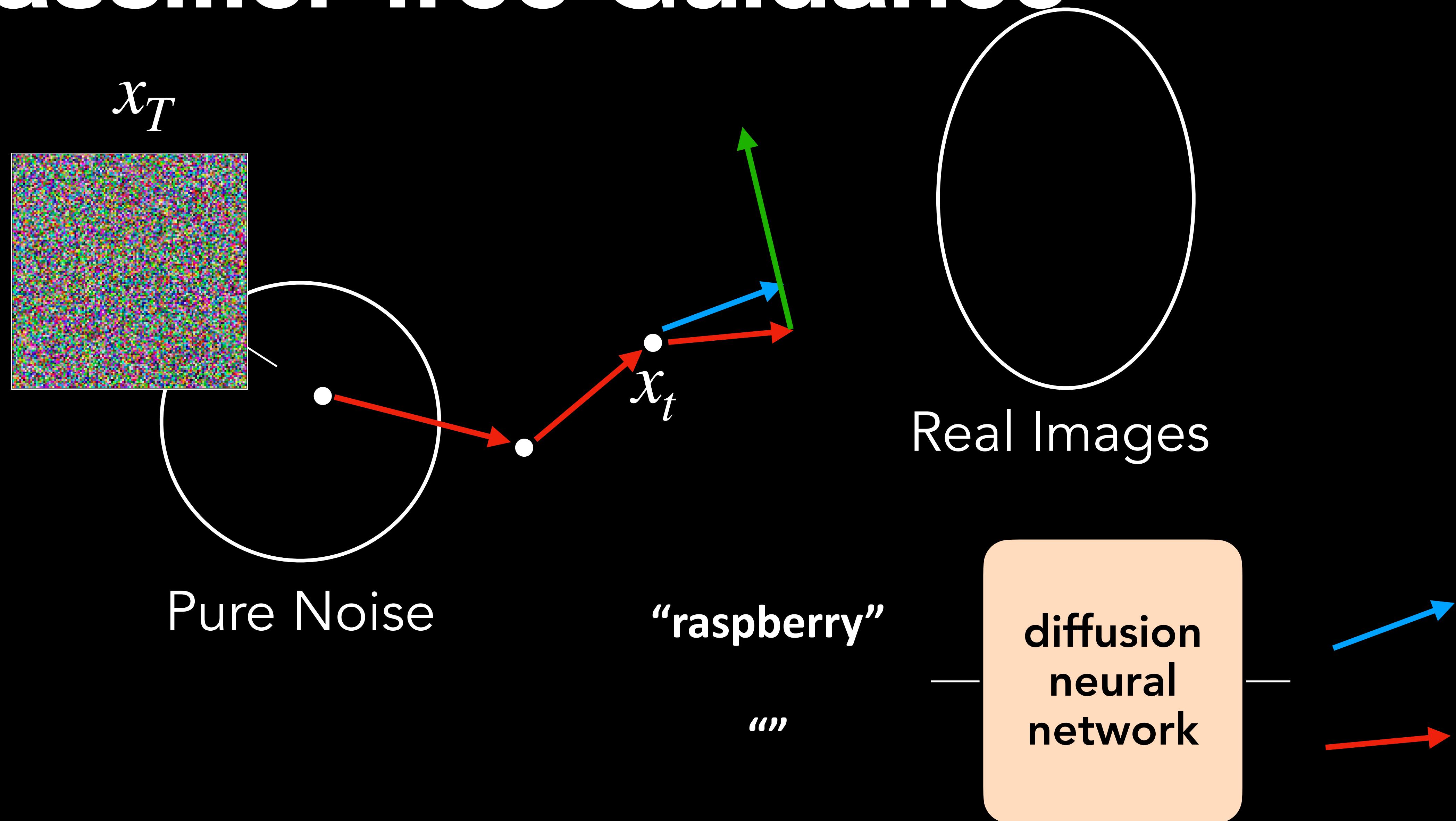
Classifier-free Guidance



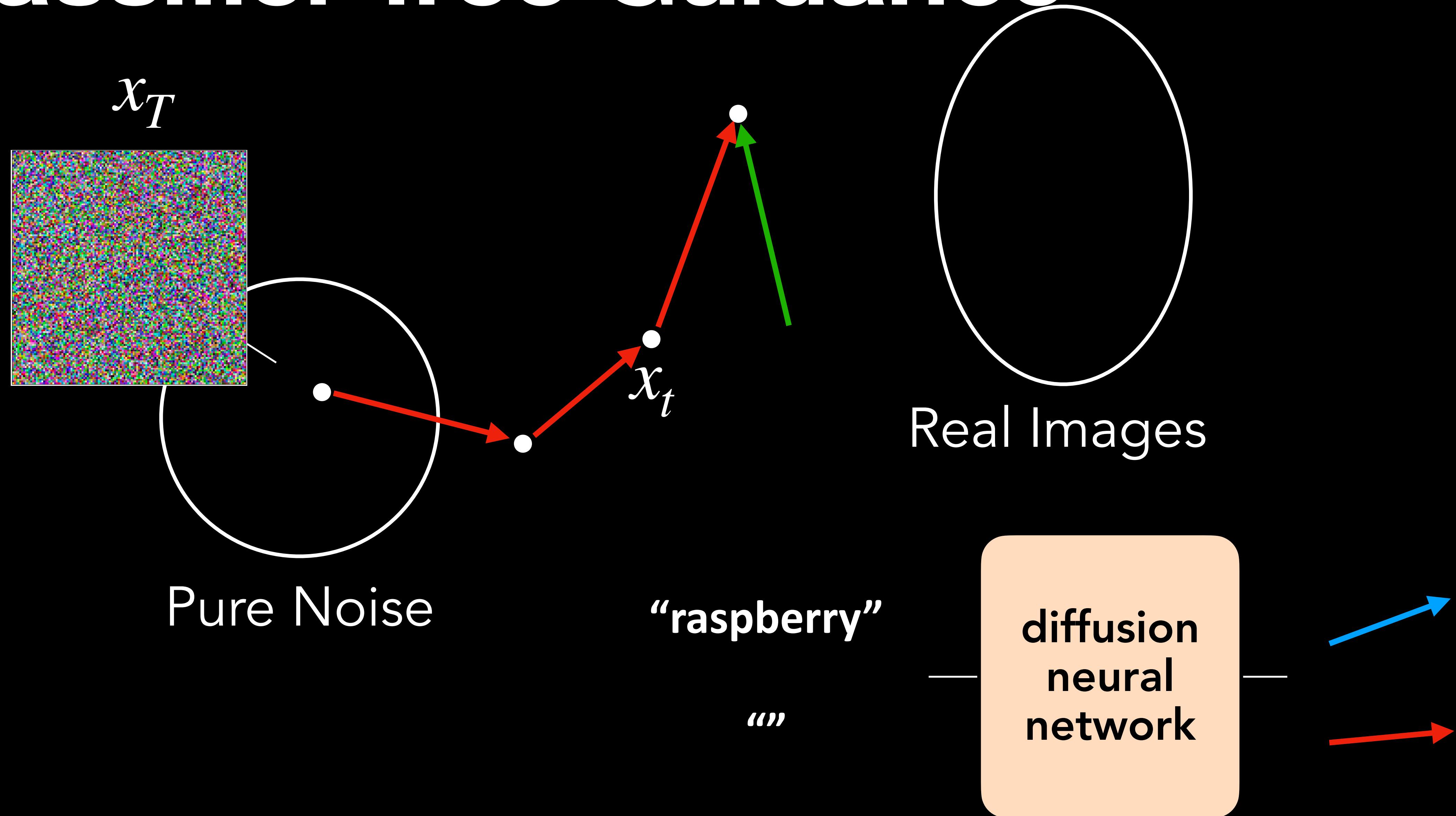
Classifier-free Guidance



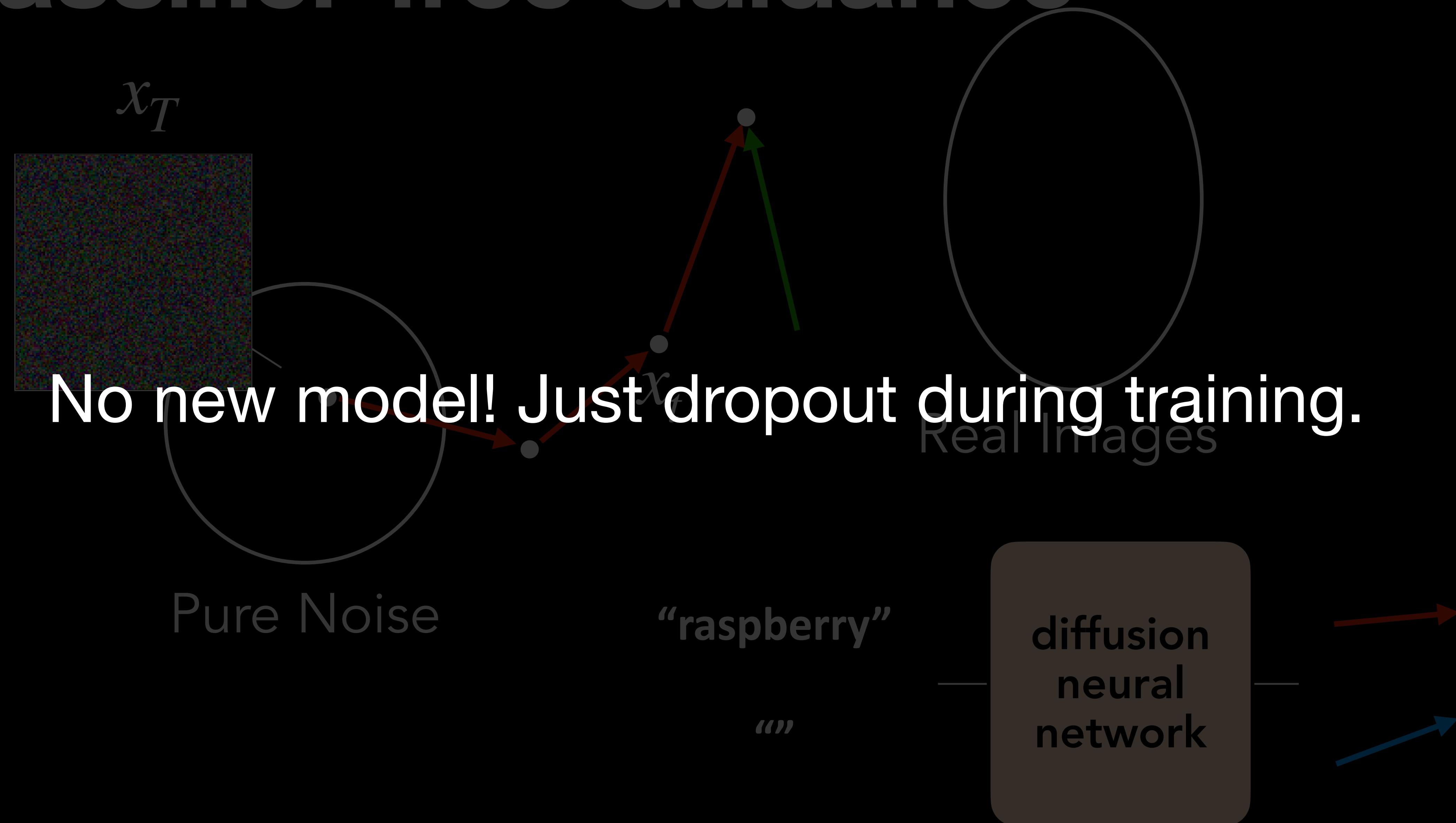
Classifier-free Guidance



Classifier-free Guidance

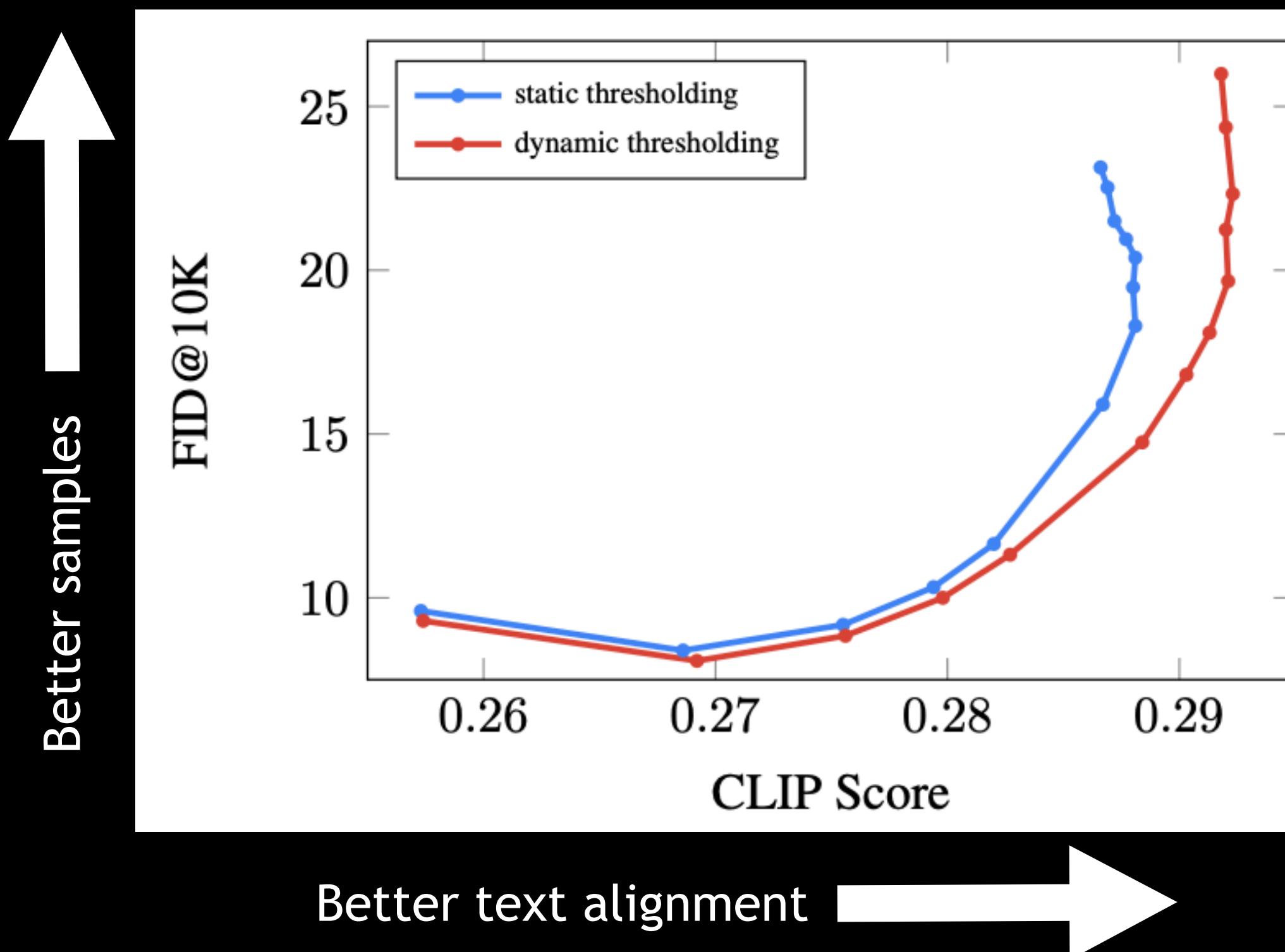


Classifier-free Guidance



Classifier-free guidance

- Large classifier-free guidance weights → better text alignment, worse image fidelity



Another issue:

Sampling is slow!

Often need 10-1000+ steps to generate high quality samples

How to make sampling faster?

One solution: distill diffusion models into models using just 4-8 sampling steps!

One example:

Progressive distillation for fast sampling of diffusion models, Salimans & Ho, ICLR 2022

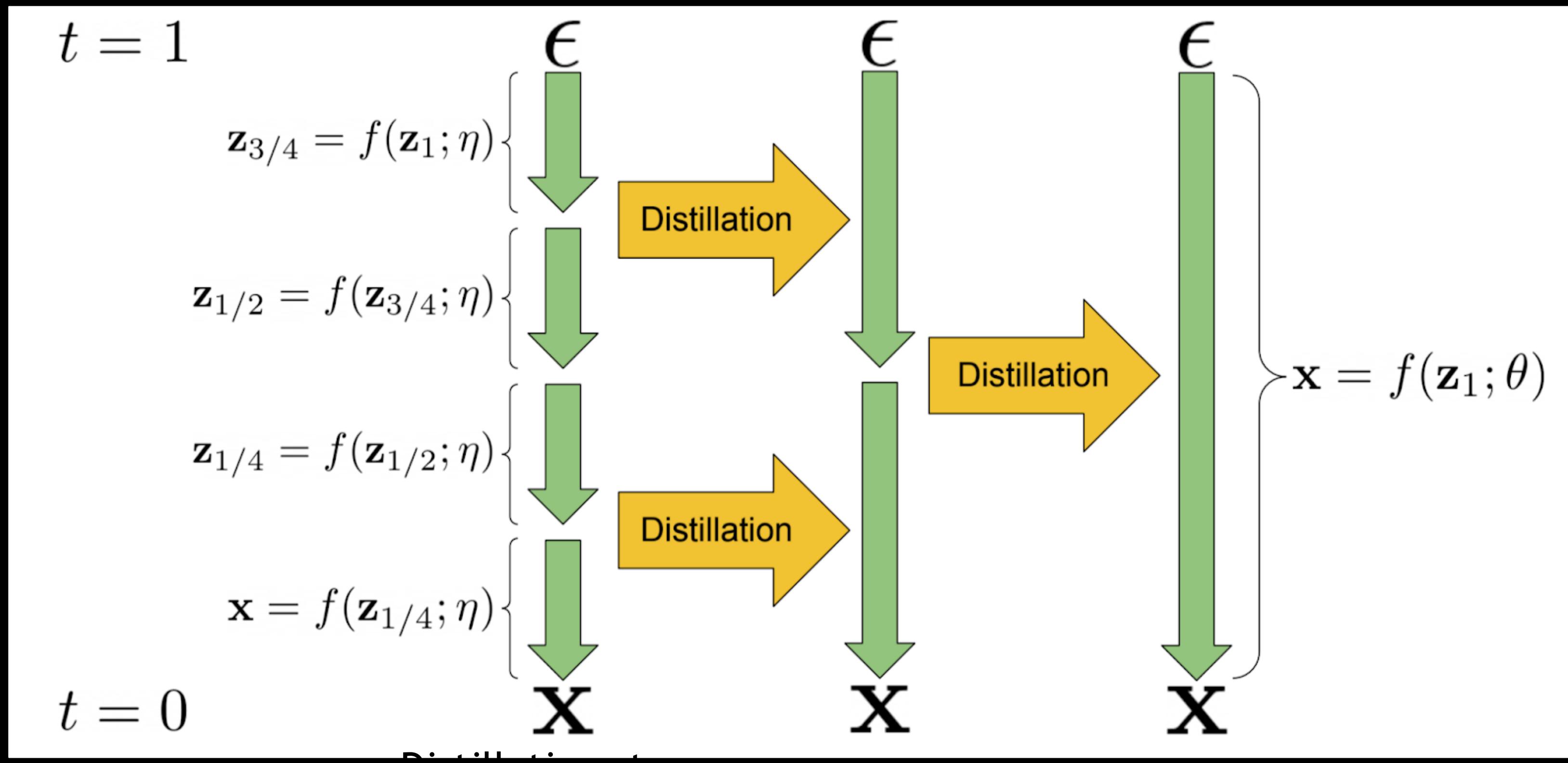
On Distillation of Guided Diffusion Models, Meng et al., CVPR 2023

Progressive distillation

Distill a deterministic ODE sampler (i.e. DDIM sampler) to the same model architecture.

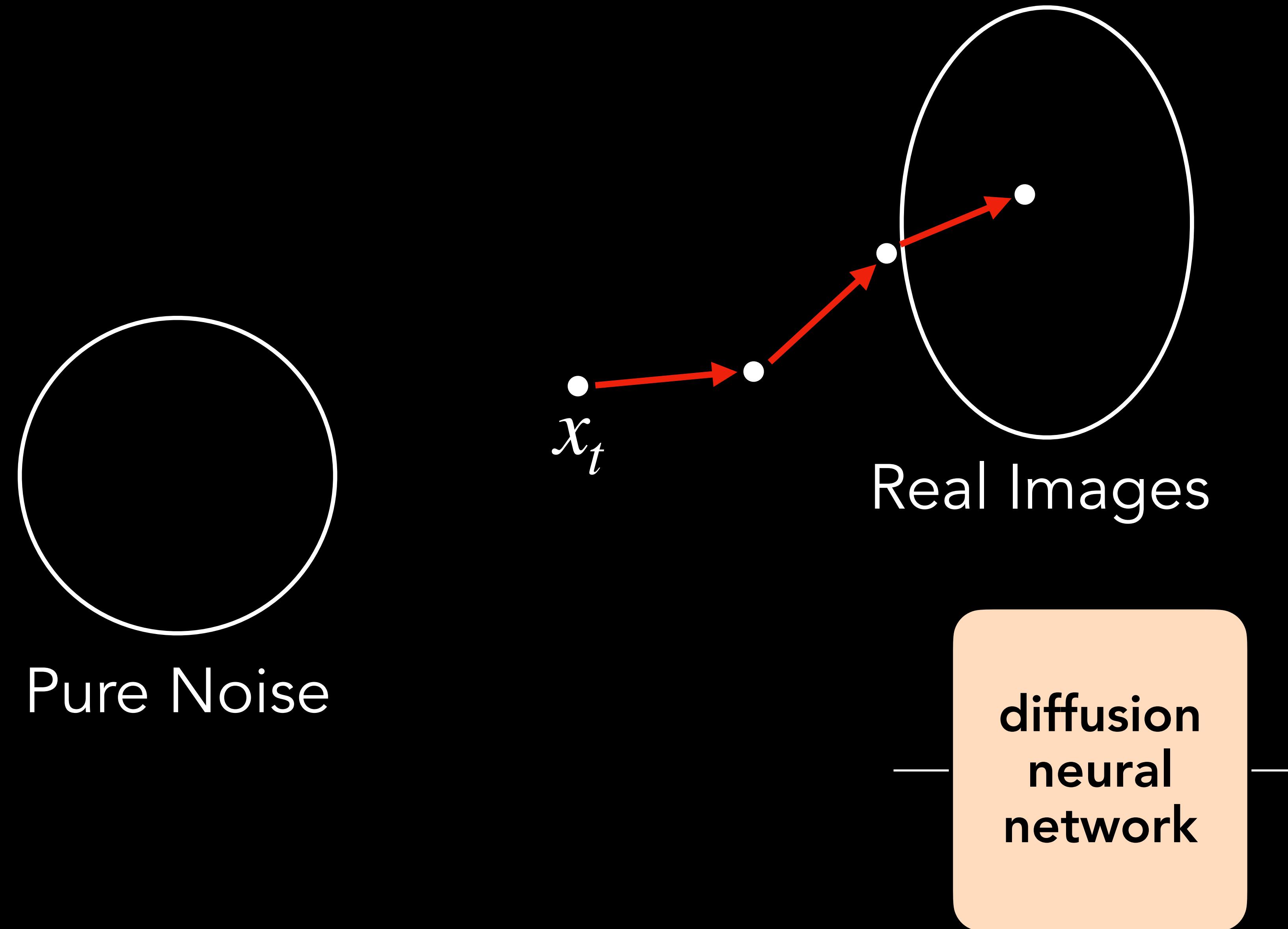
A “student” model is learned to distill two adjacent sampling steps of the “teacher” model to one sampling step.

At next stage, the “student” model from previous stage will serve as the new “teacher” model.



Some other sampling tricks

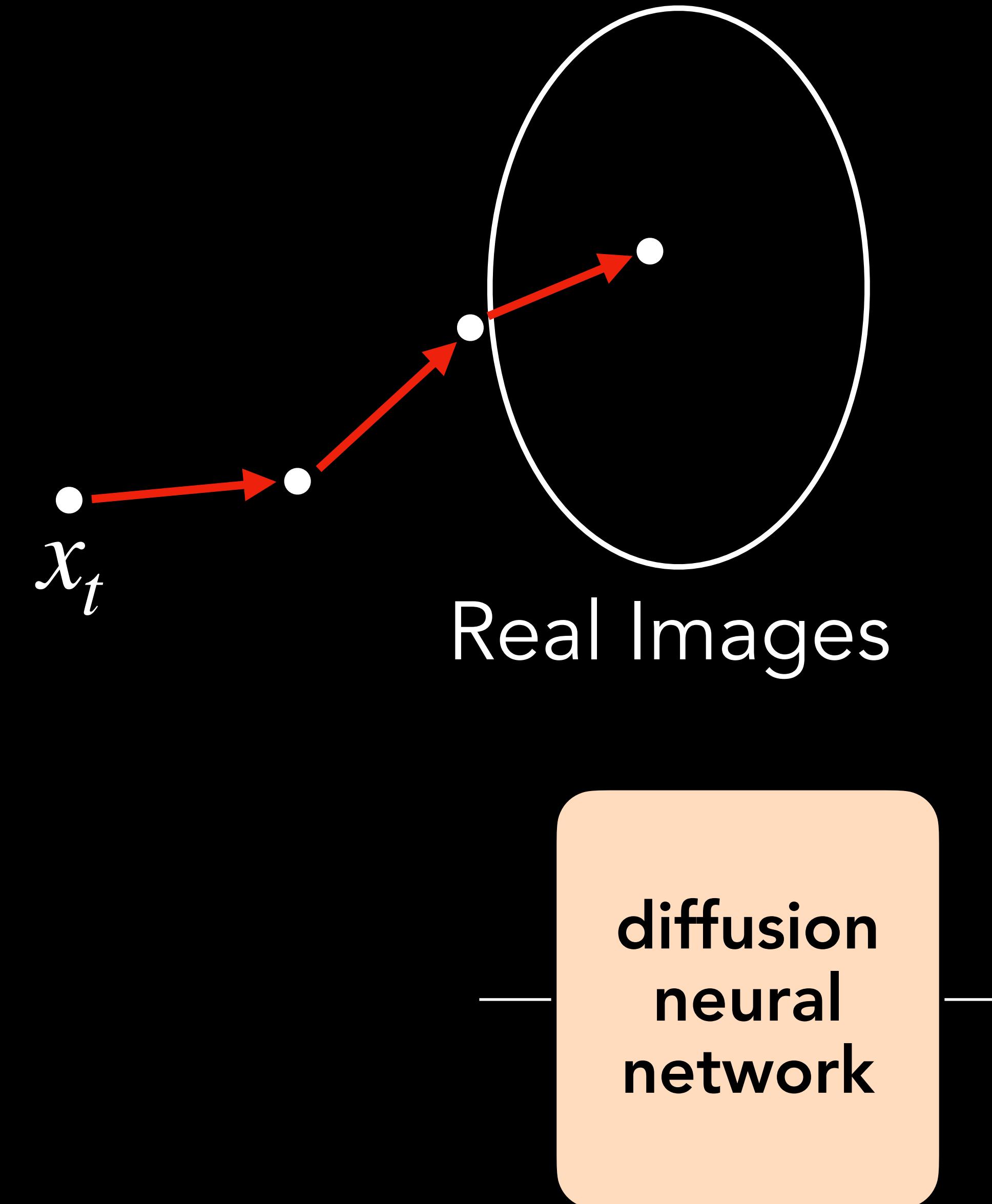
SDEdit: don't start from pure noise!



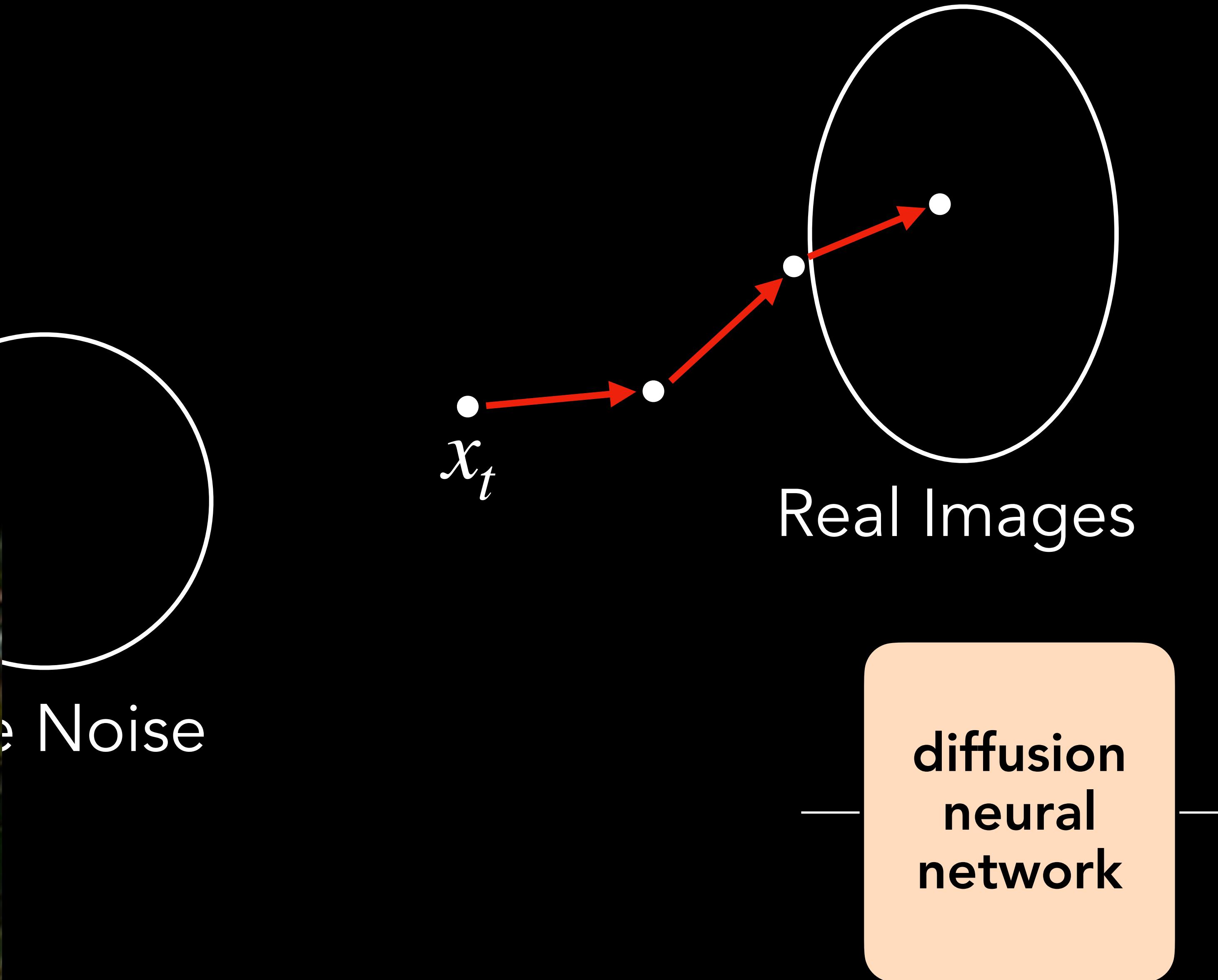
SDEdit: don't start from pure noise!



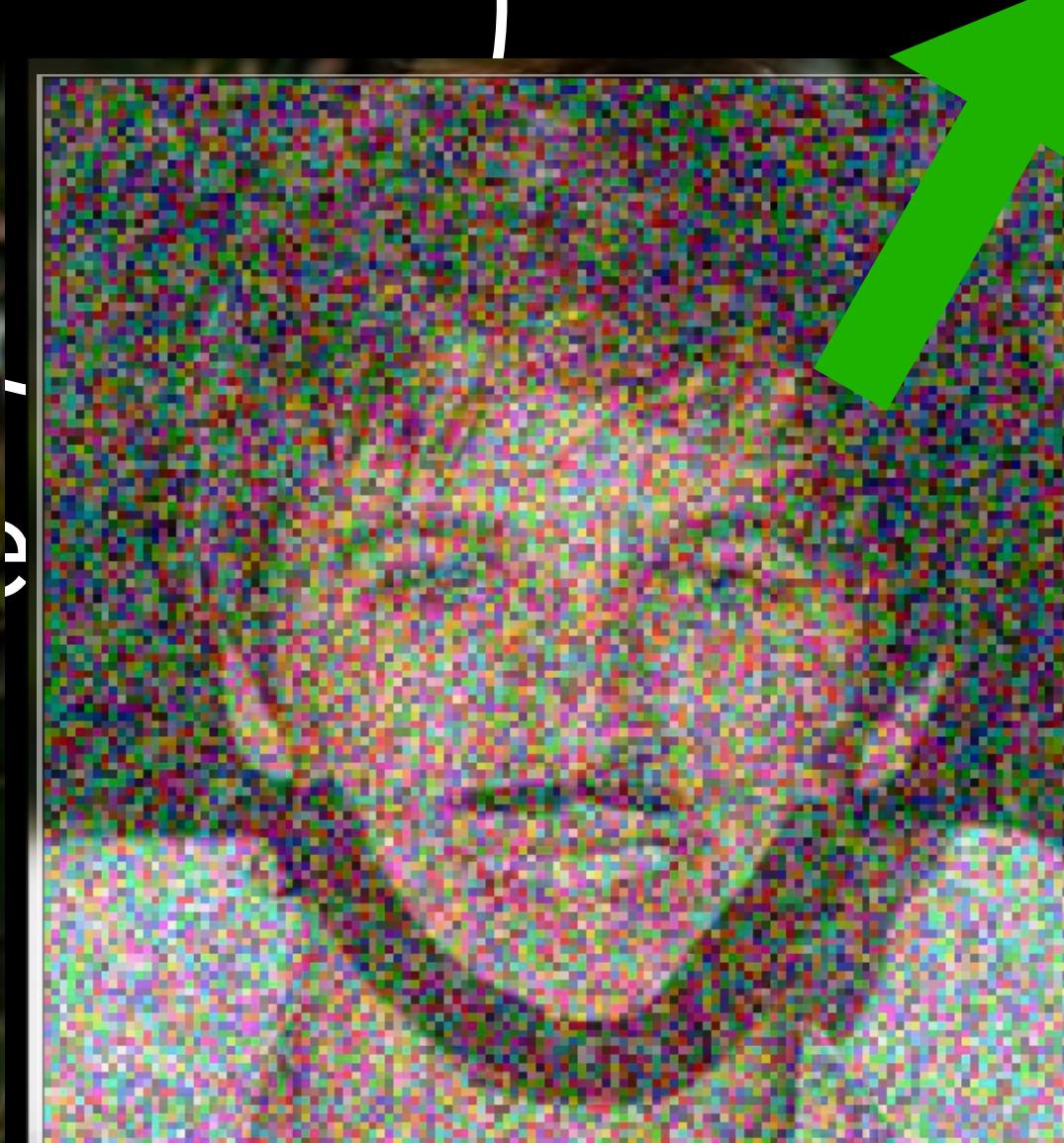
Pure Noise



SDEdit: don't start from pure noise!

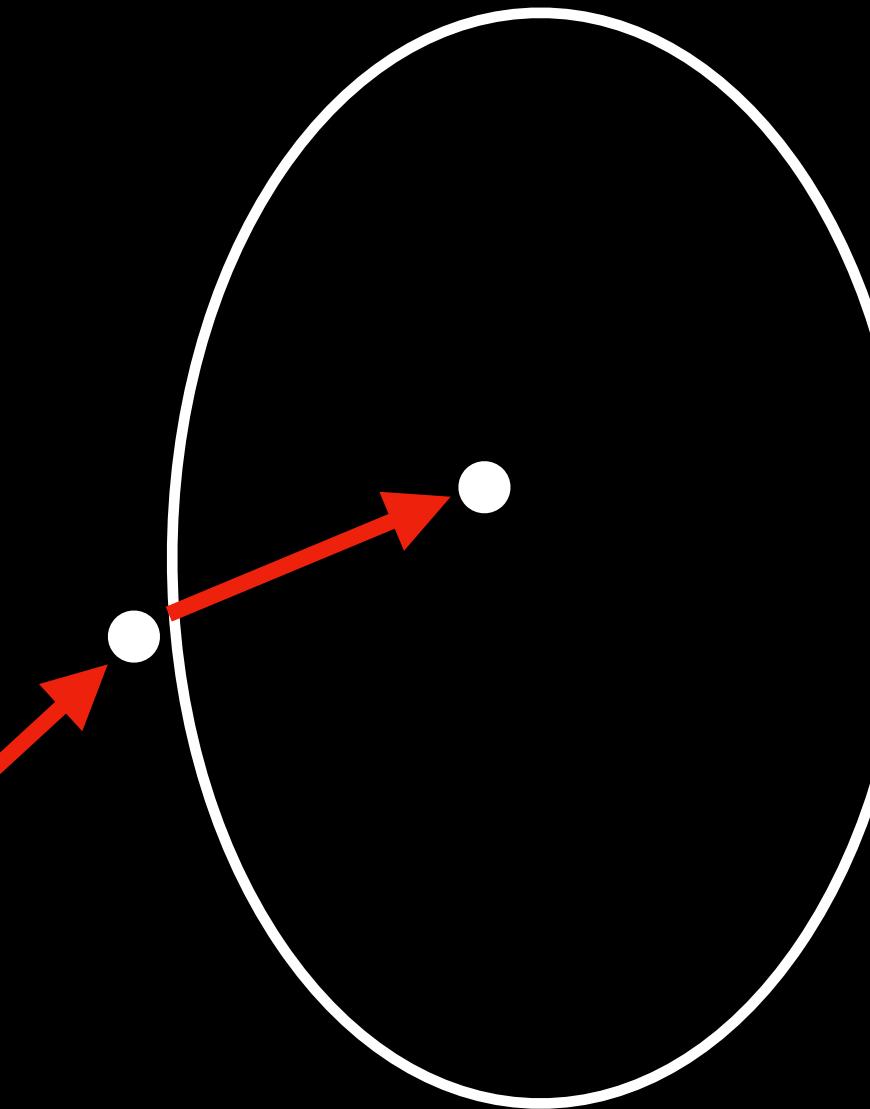


SDEdit: don't start from pure noise!



\dot{x}_t

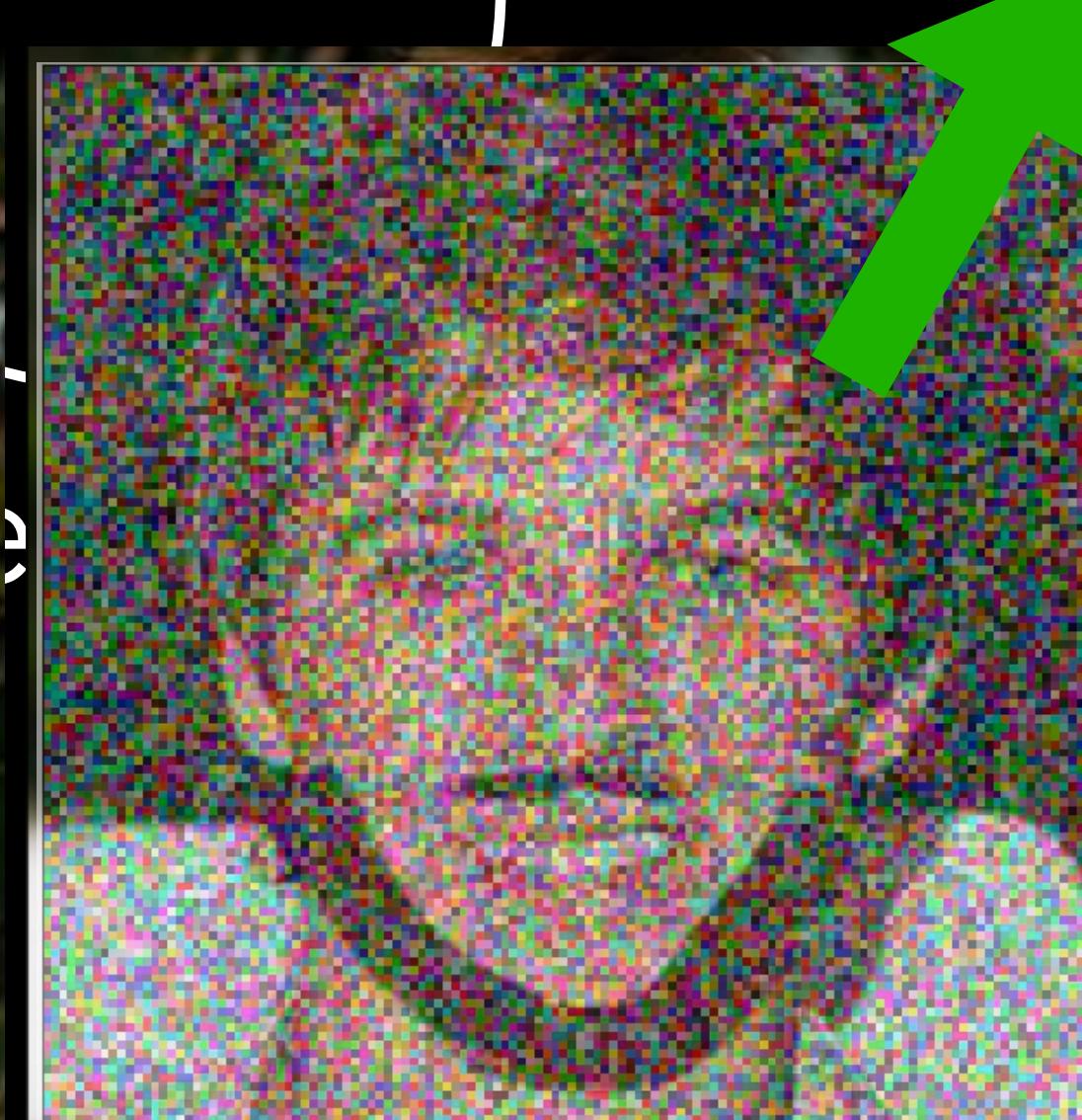
**“Guy with
beard”**



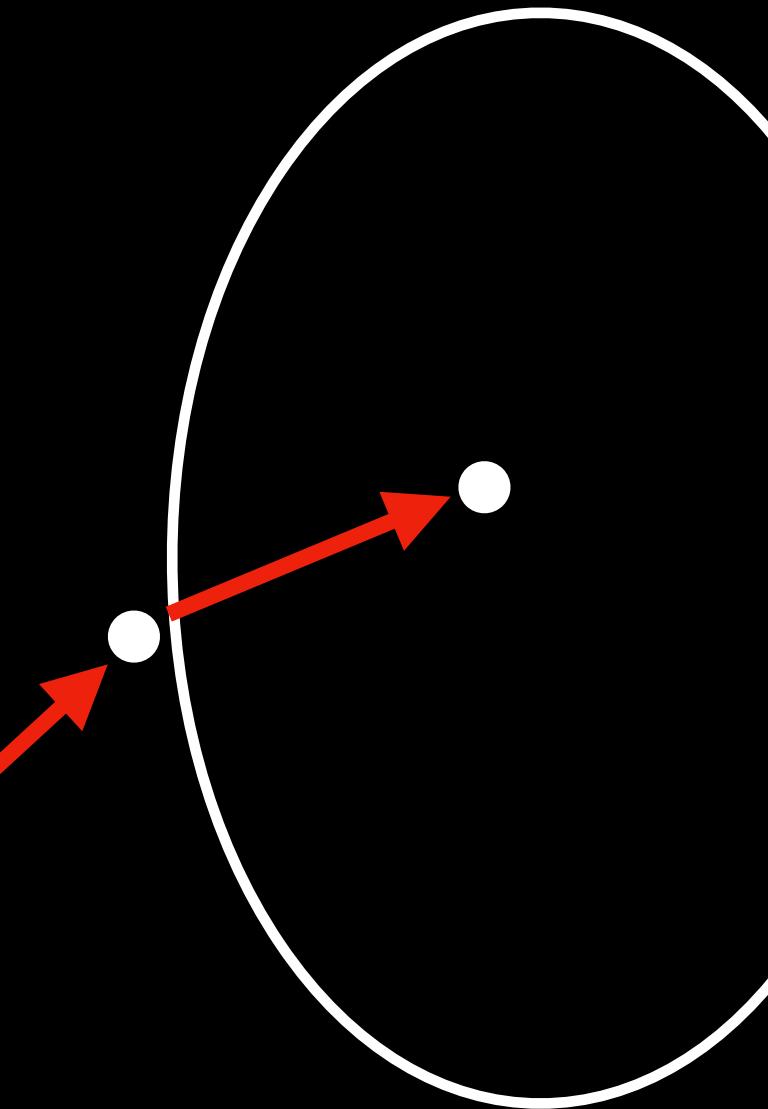
Real Images

**diffusion
neural
network**

SDEdit: don't start from pure noise!

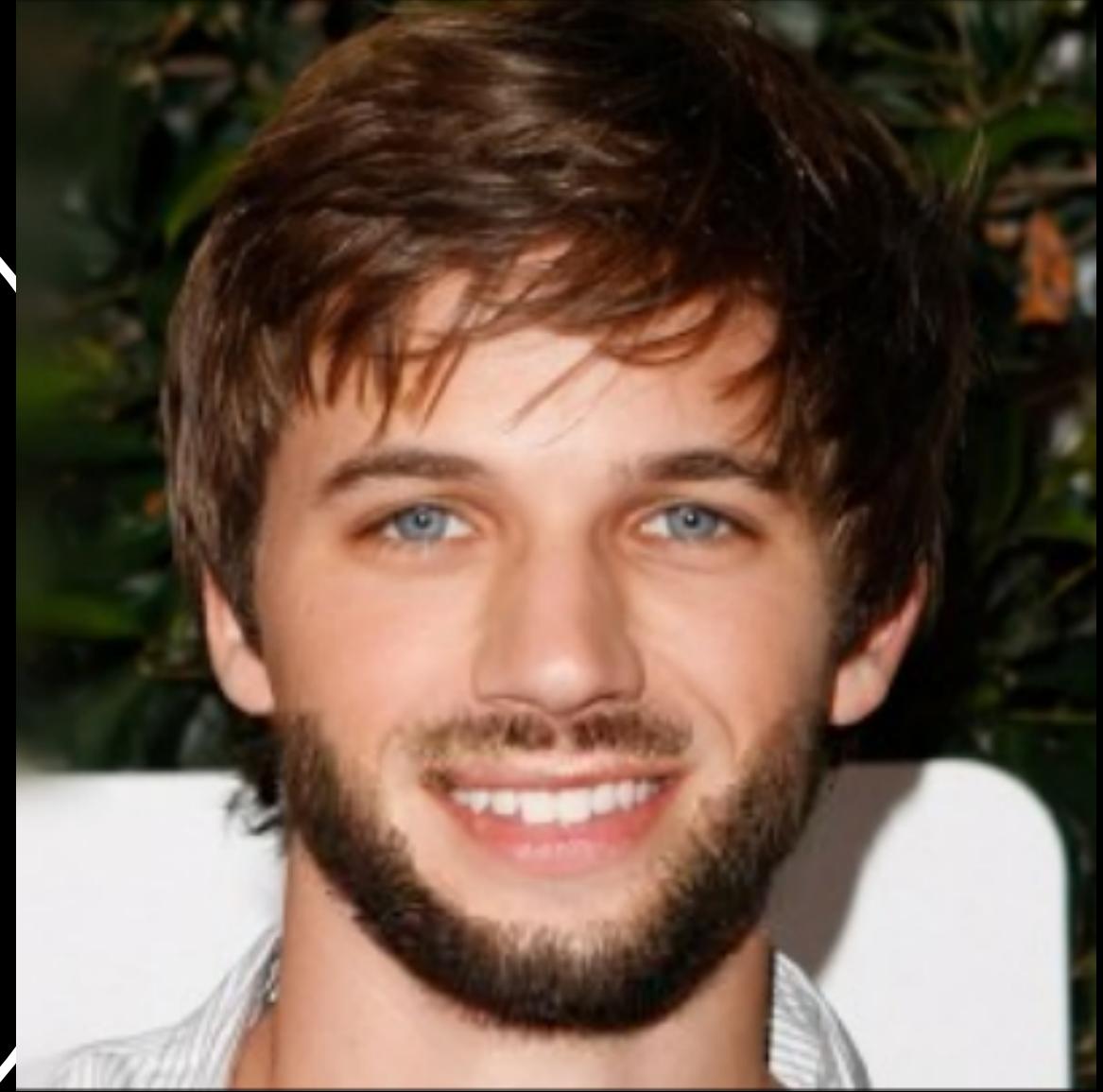
 x_t

“Guy with
beard”

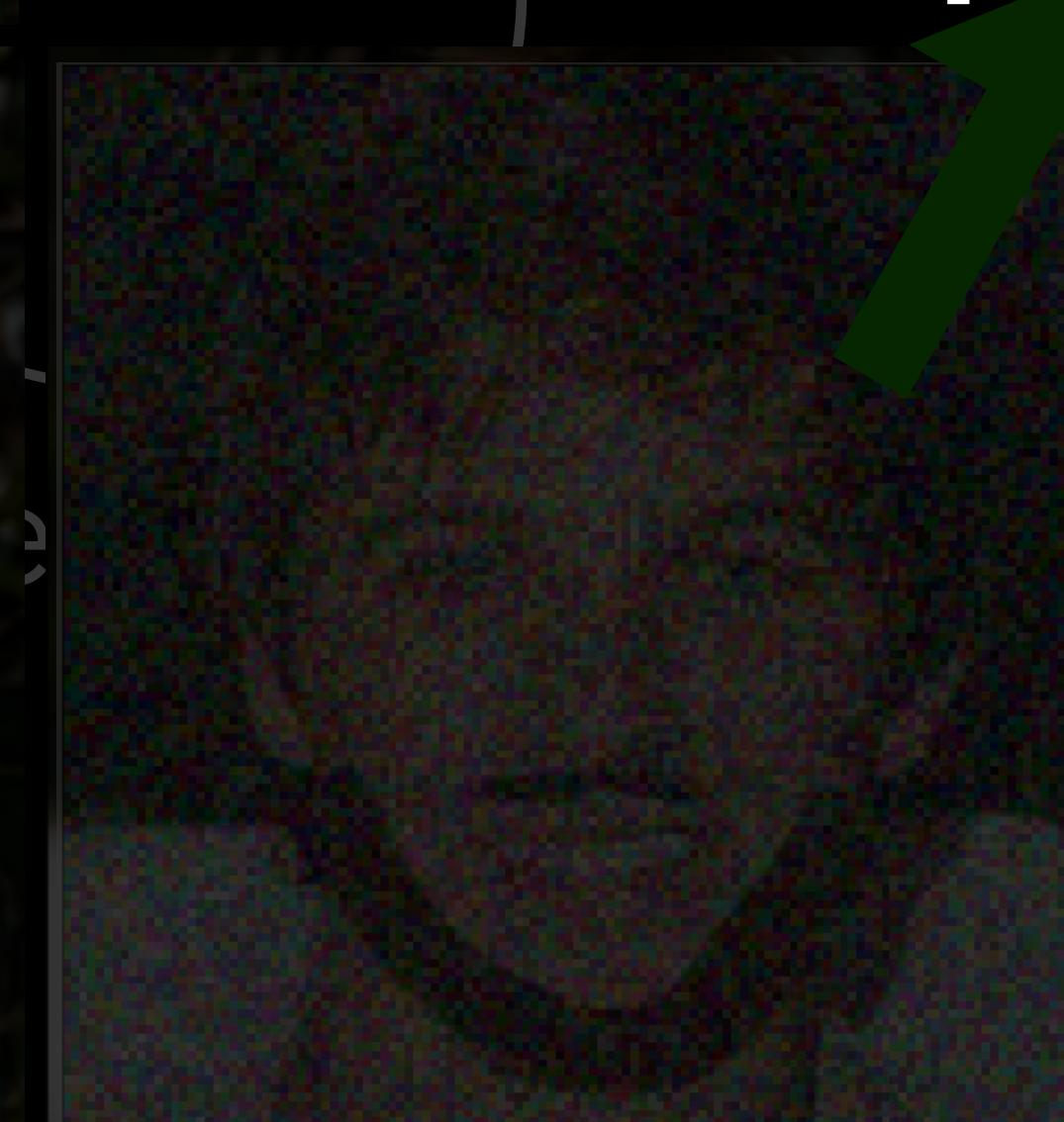
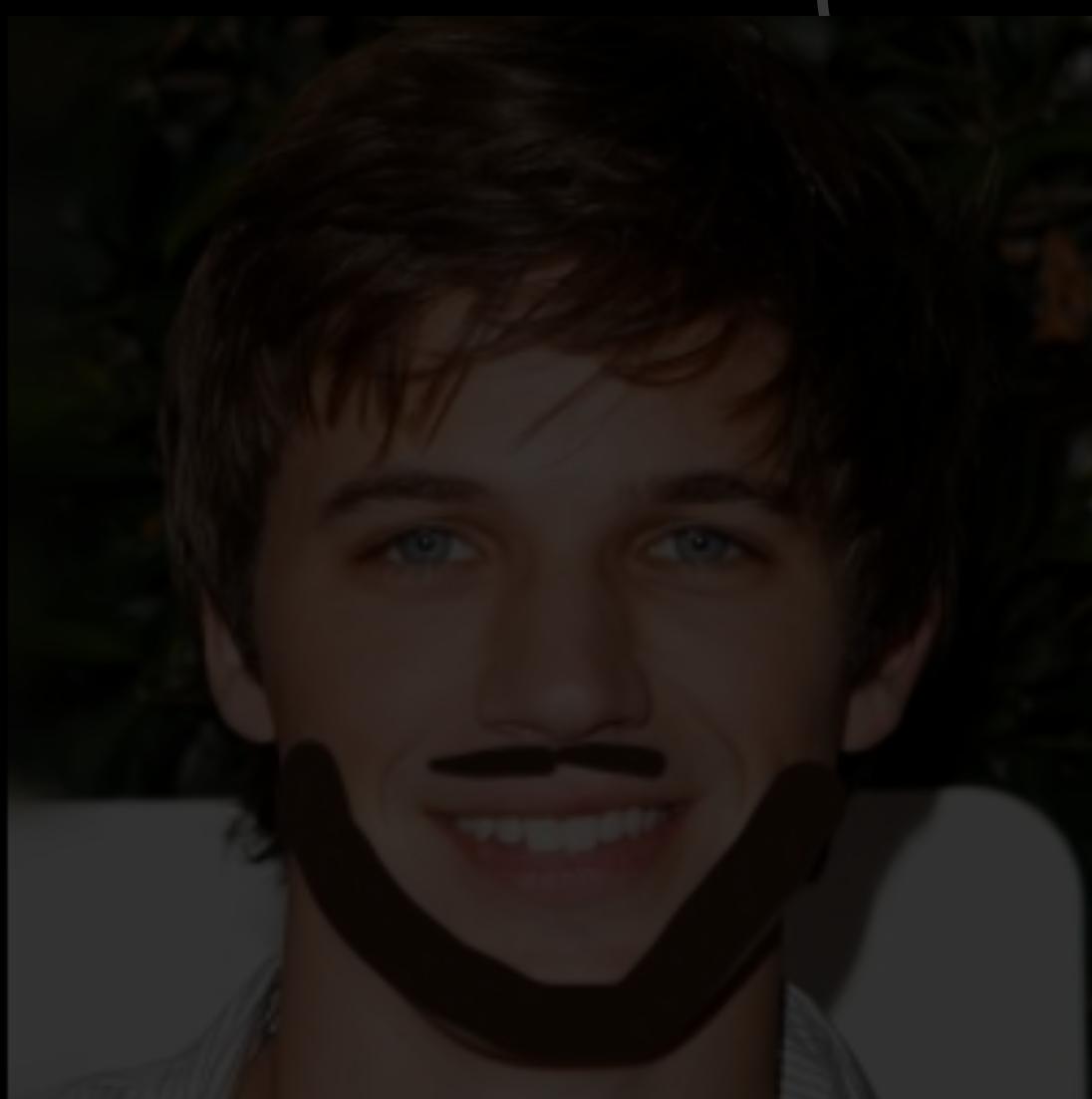
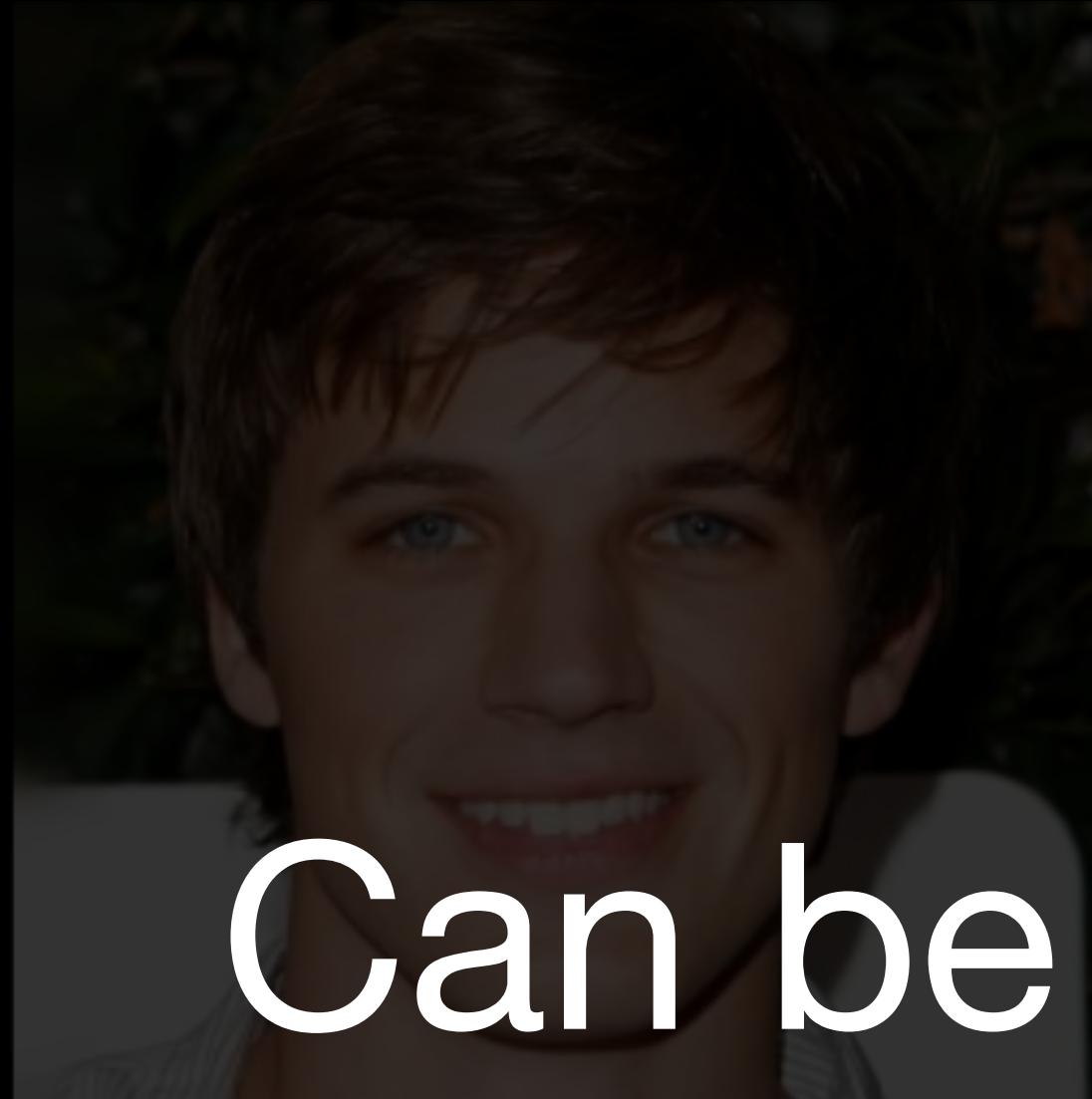


Real Images

diffusion
neural
network



SDEdit: don't start from pure noise!

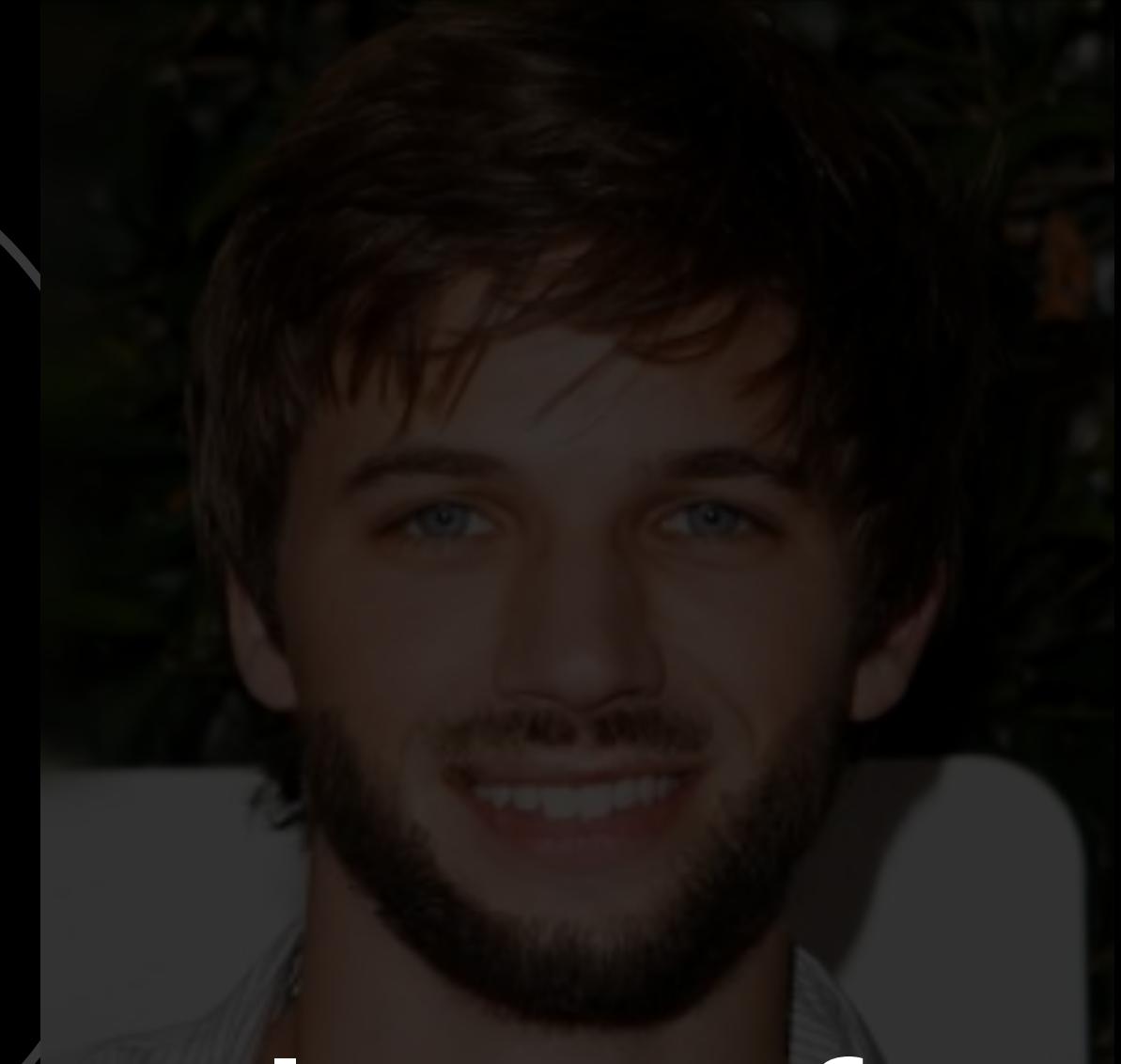


Can be hard to pick the right value of t

Real Images

“Guy with
beard”

diffusion
neural
network



Prompt-to-Prompt: Generating two similar images

Photo of a cat riding a bicycle



Photo of a cat riding a car



Prompt-to-Prompt: Generating two similar images

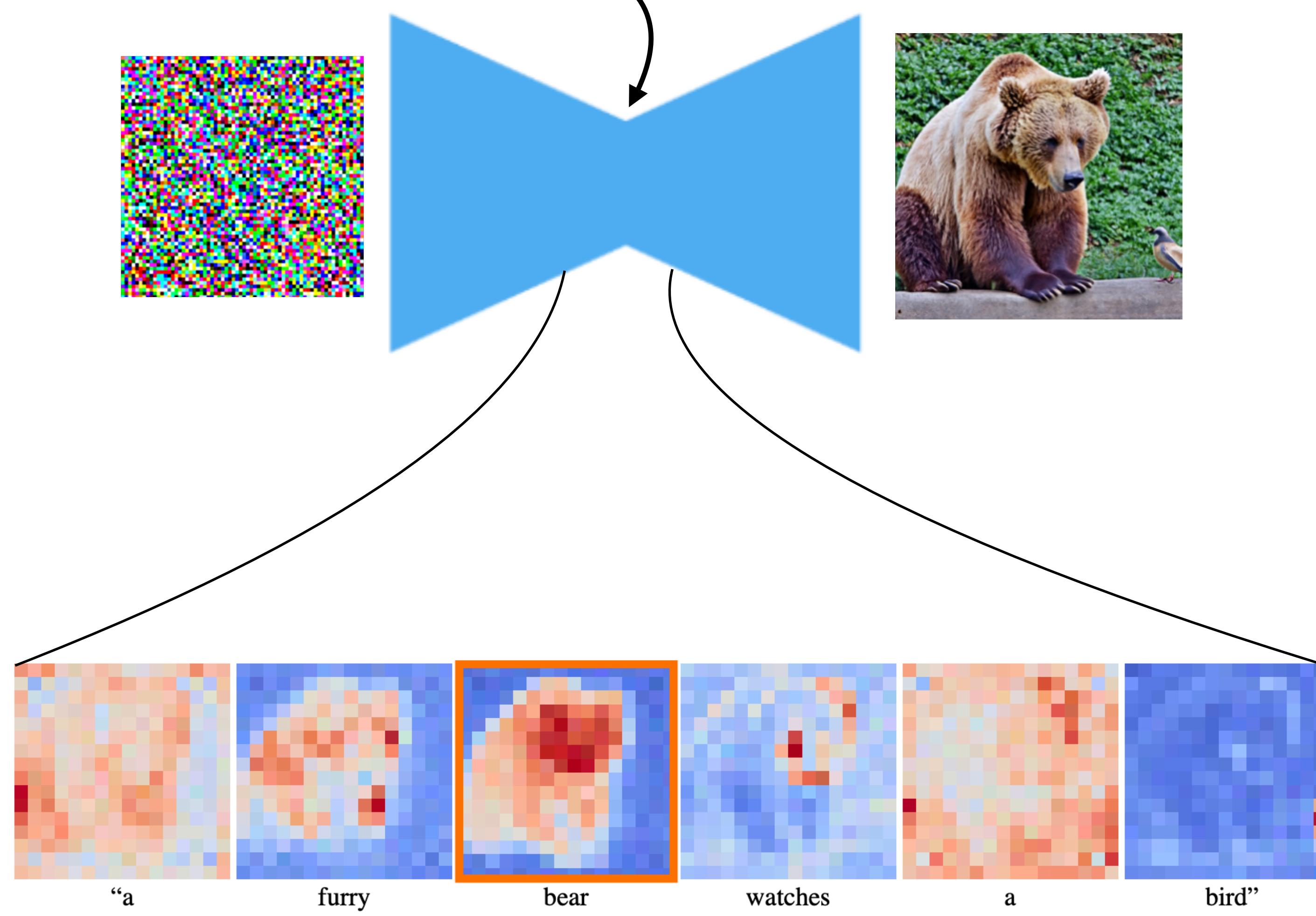
Photo of a cat riding a bicycle

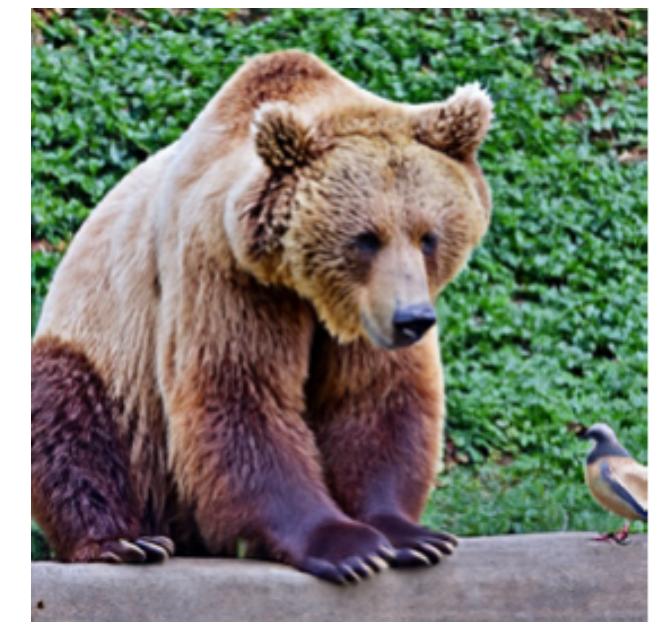


Photo of a cat riding a car

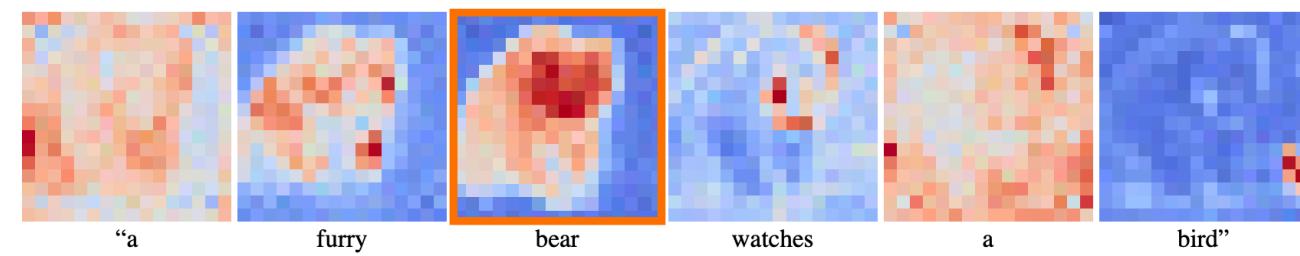
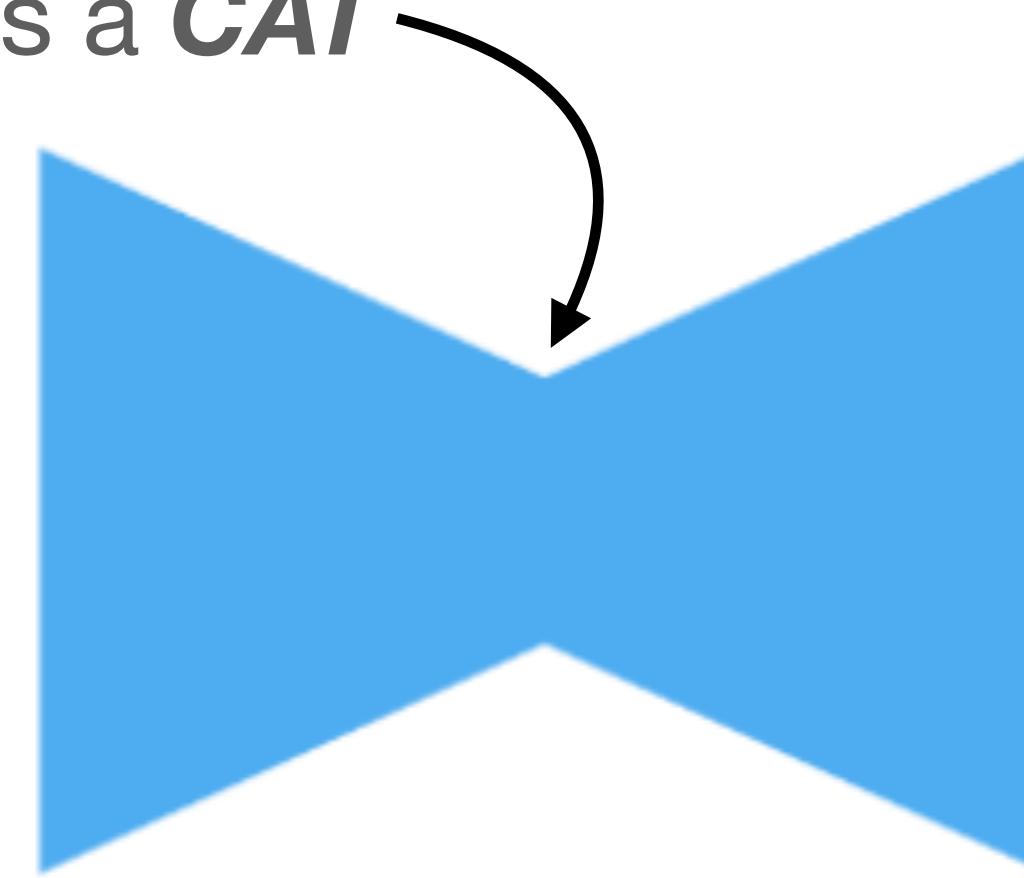
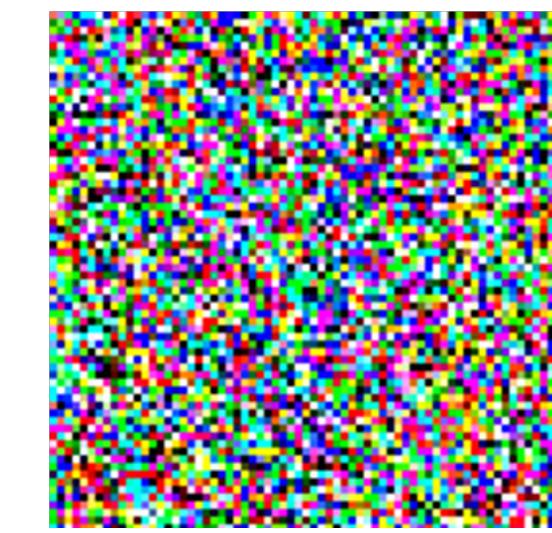


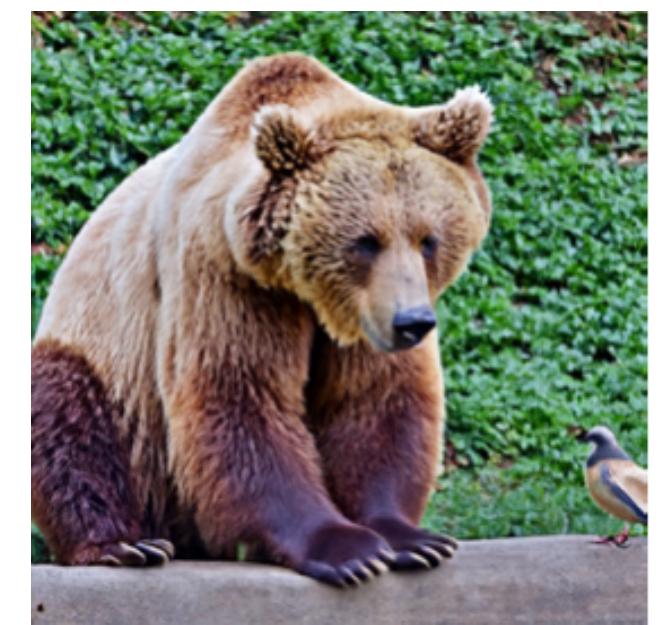
A furry bear watches a bird



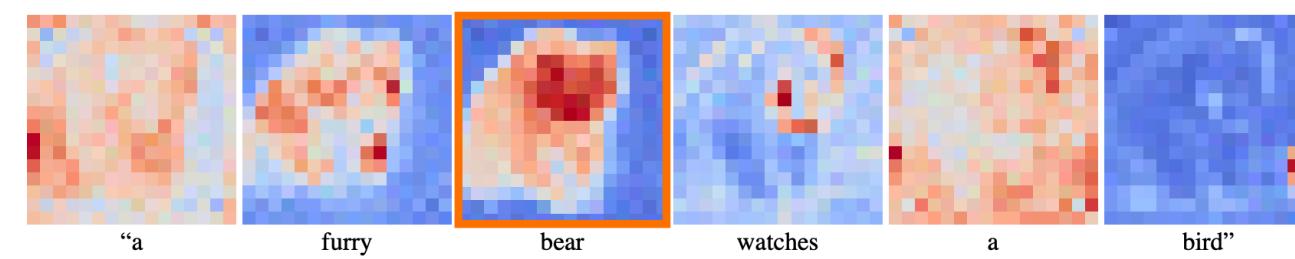
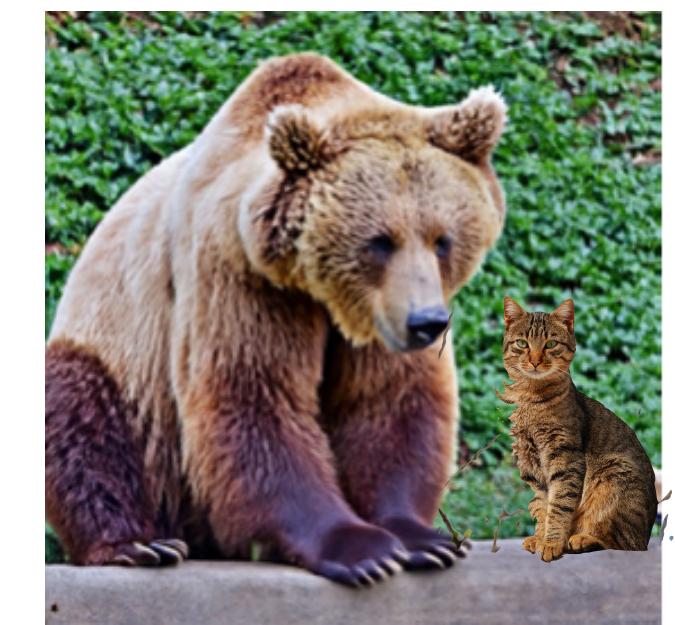
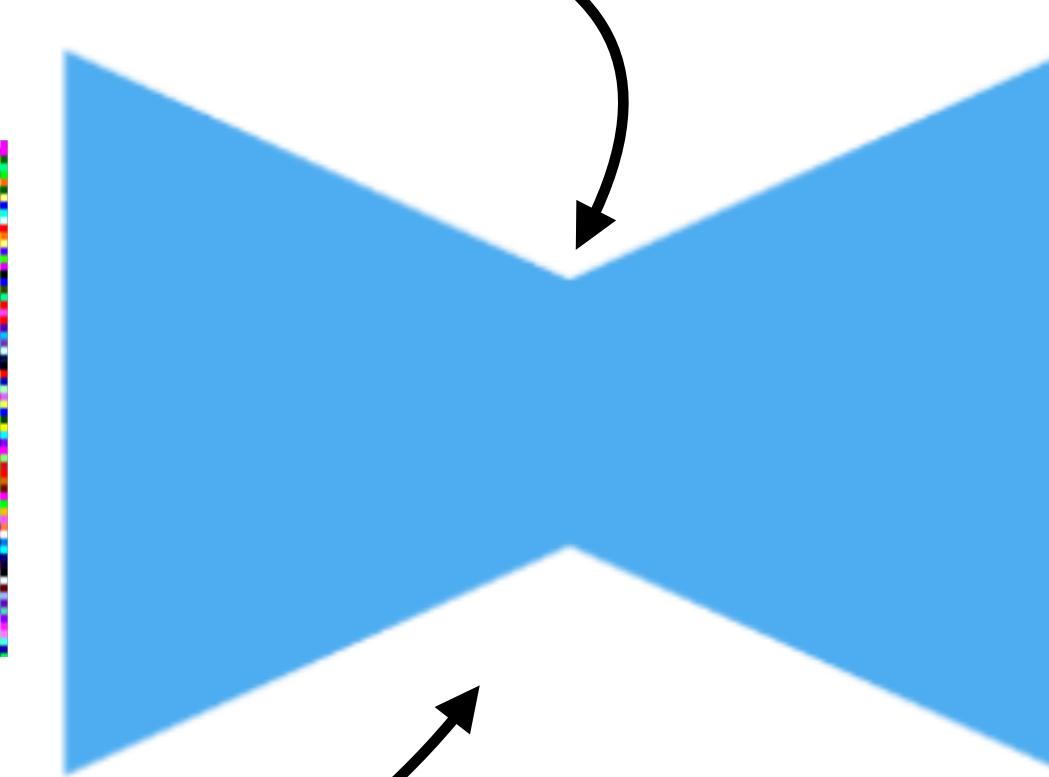


A furry bear watches a **CAT**

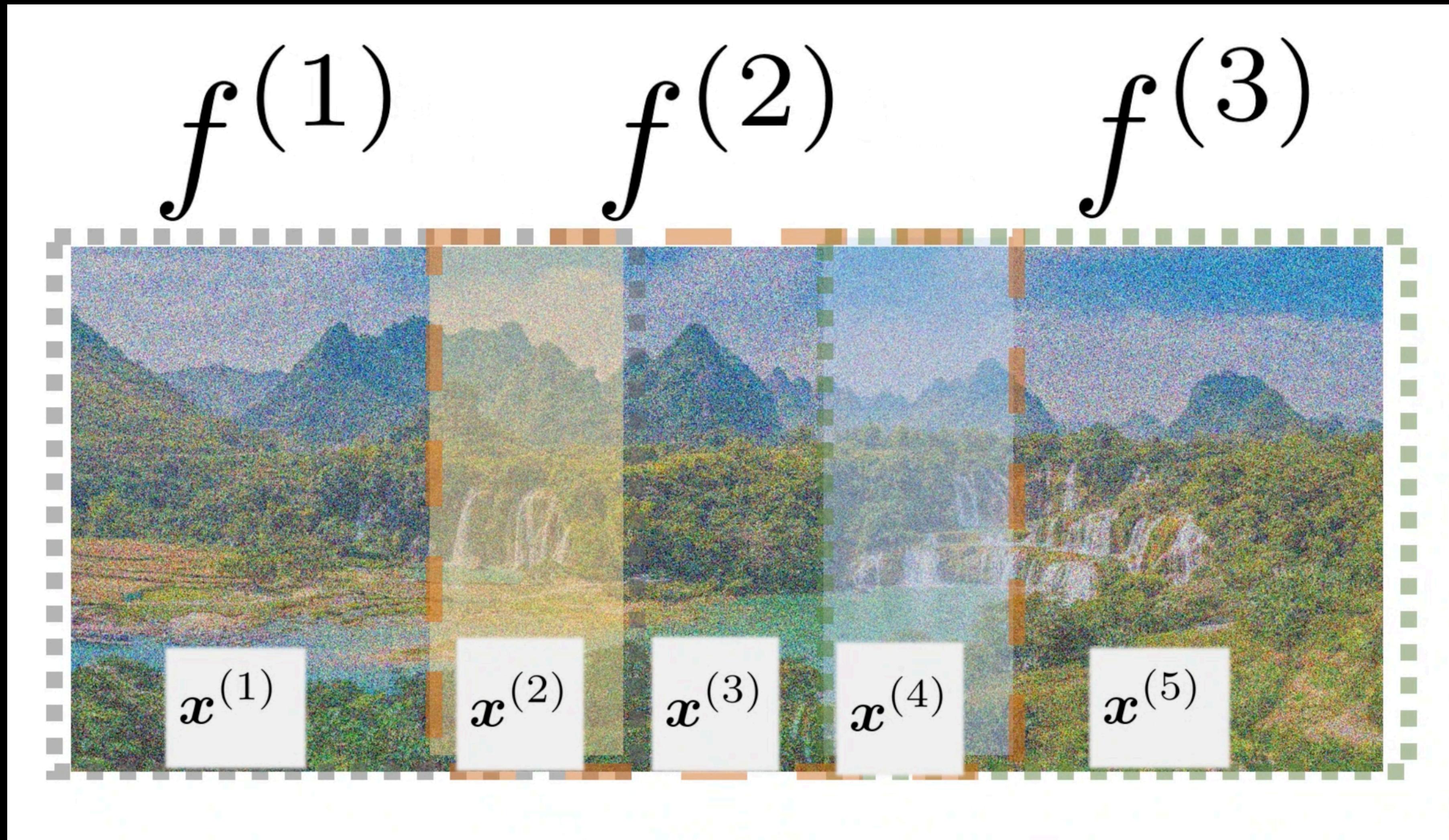




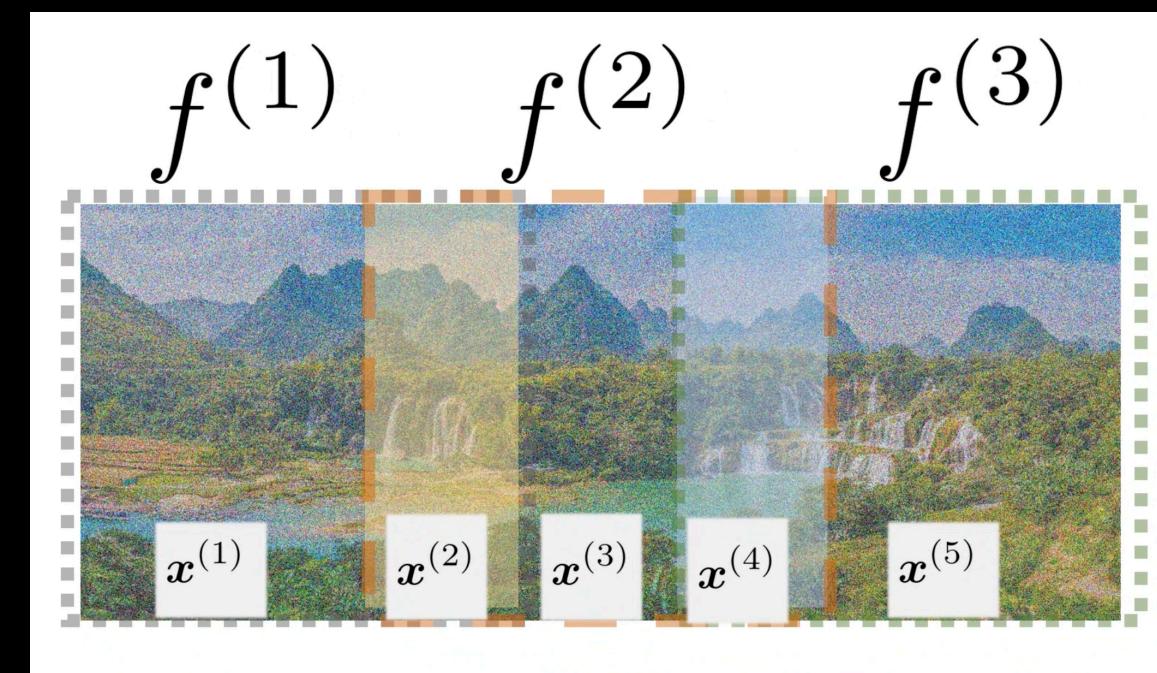
A furry bear watches a **CAT**



Mixing scores



Mixing scores



Can't do this with latent models!

