

Diffusion Models 2

Practicals

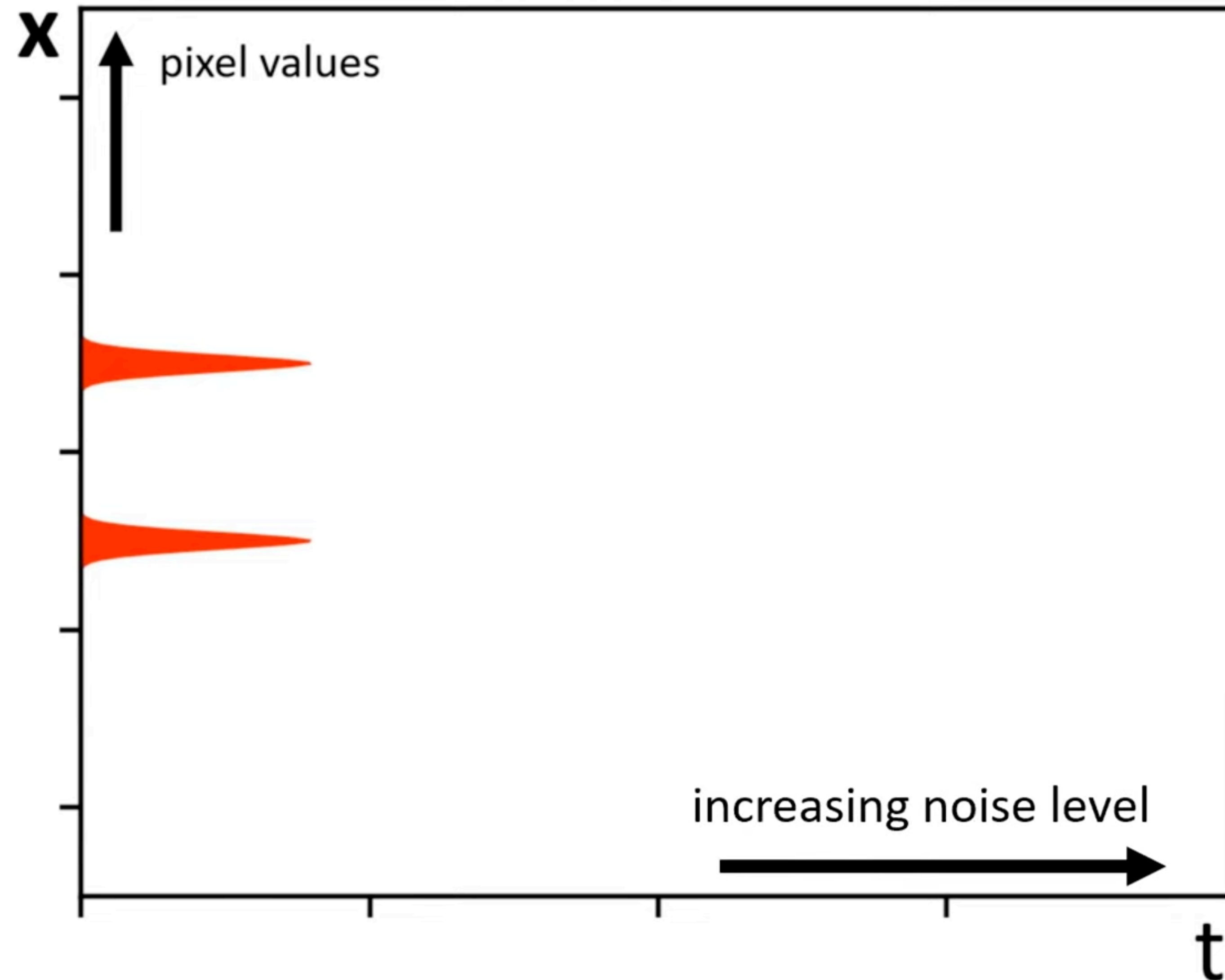
CS280

Songwei Ge & David McAllister

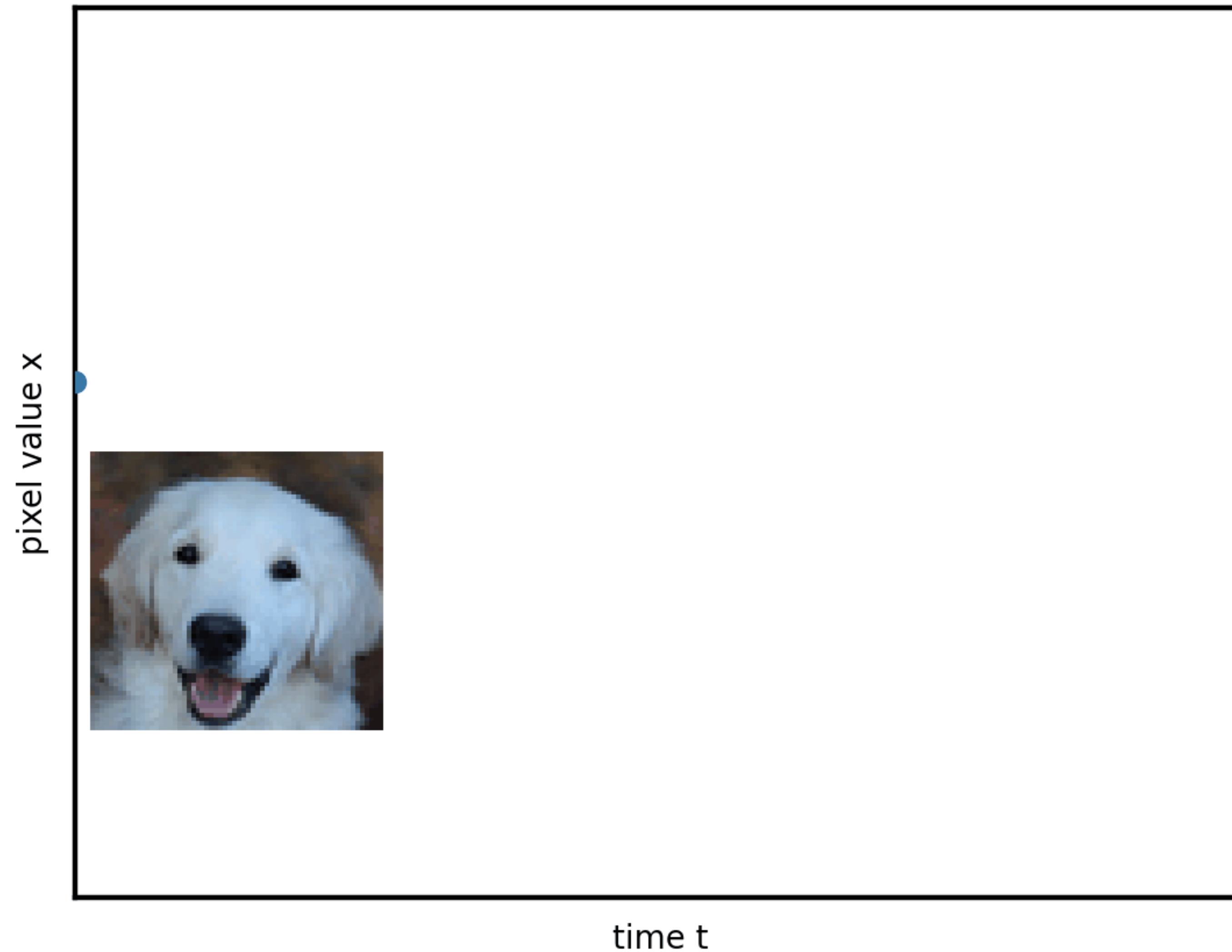
Outline

- Flow matching model samplers
- Inversion/Distillation
- Guidance
- Application

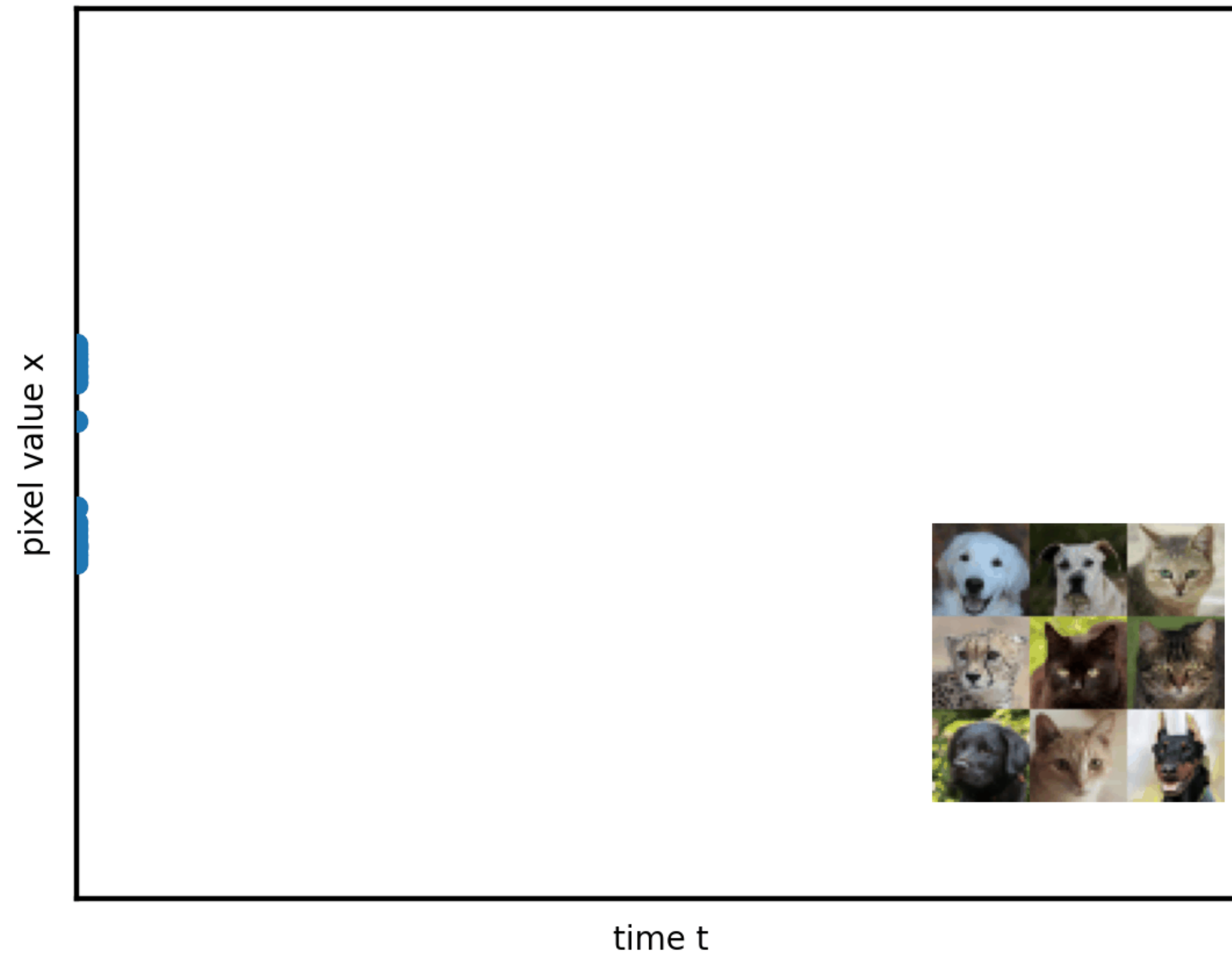
Revisit diffusion models with a 1D example



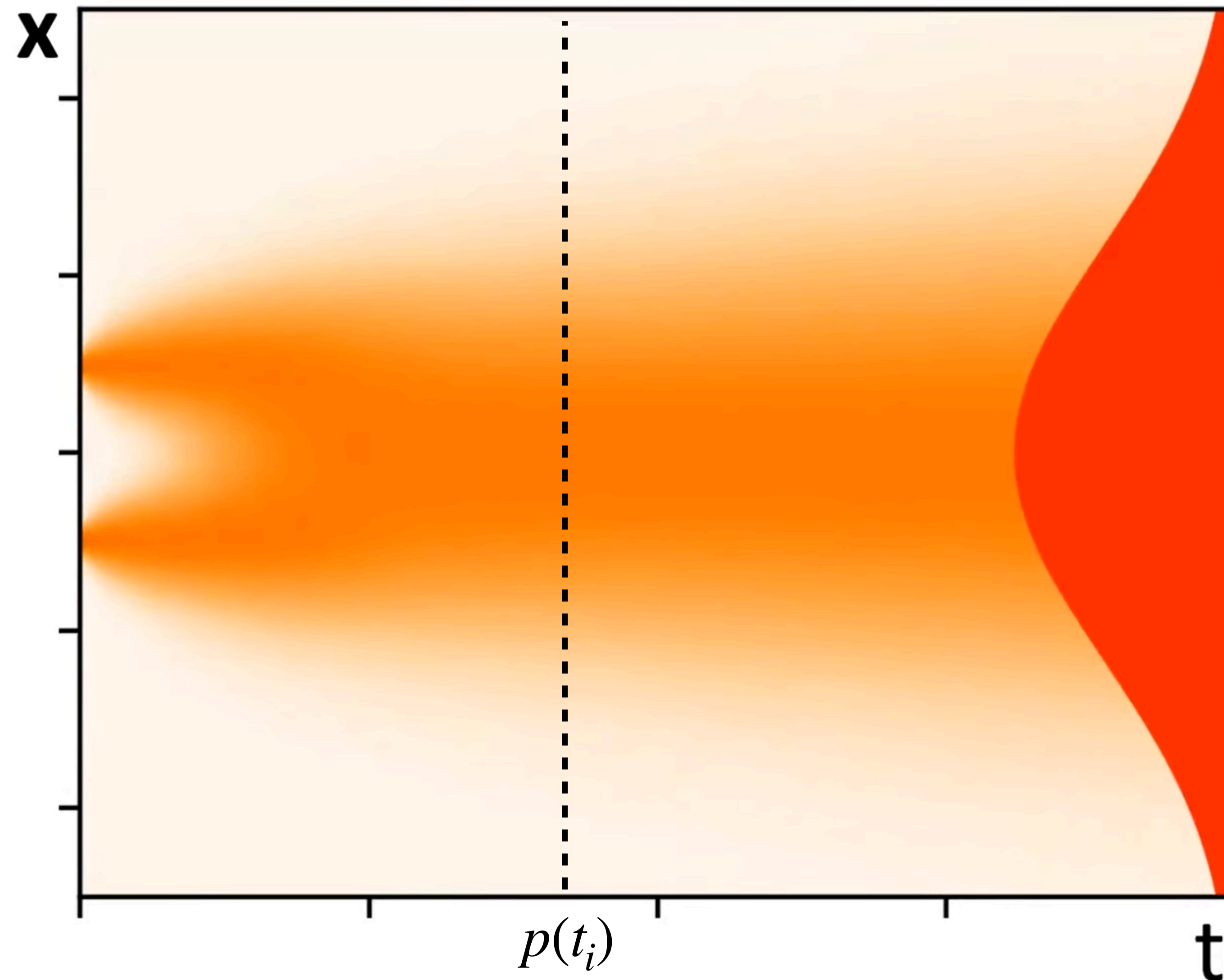
Revisit diffusion models with a 1D example



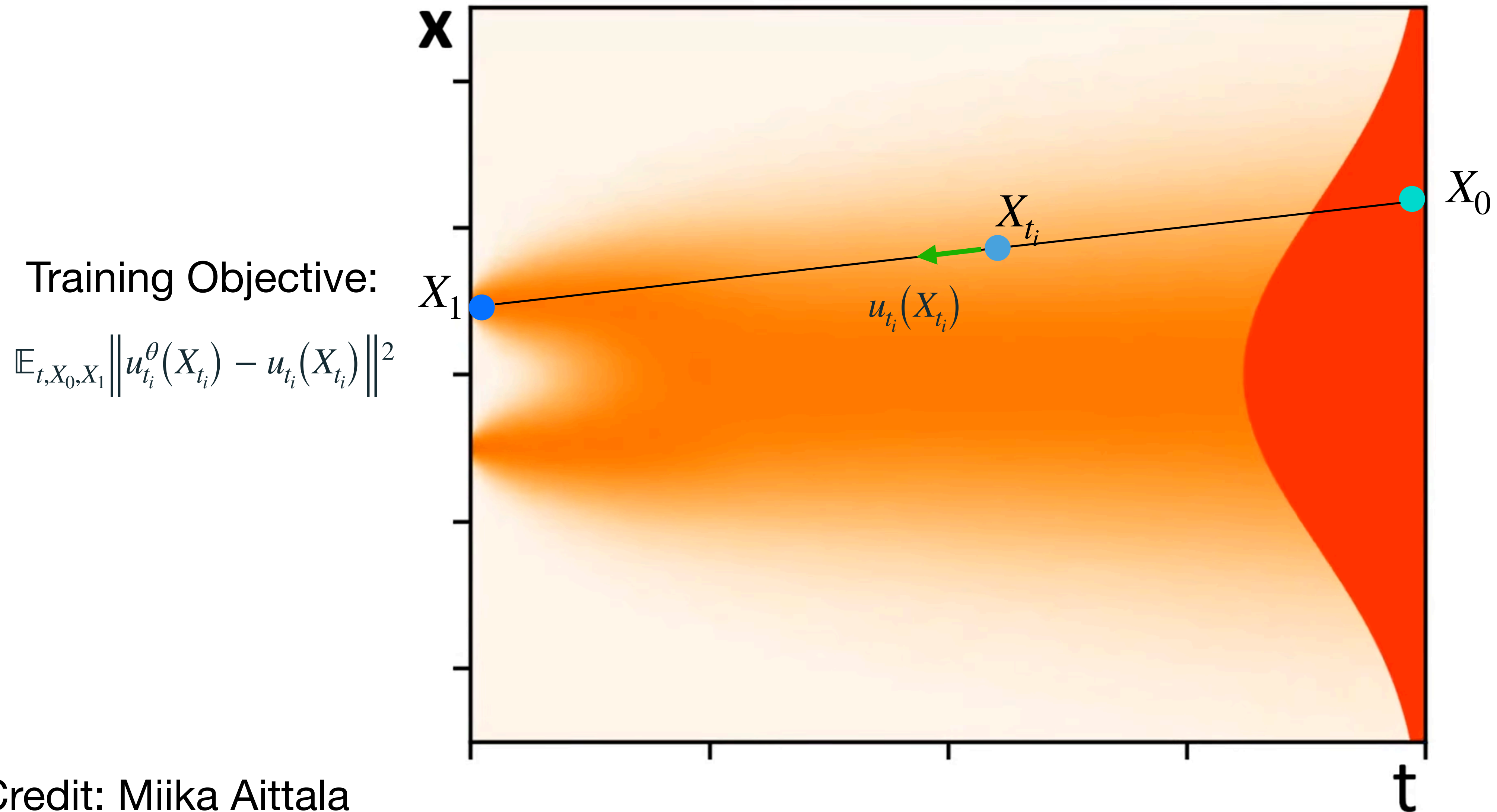
Revisit diffusion models with a 1D example



Revisit diffusion models with a 1D example



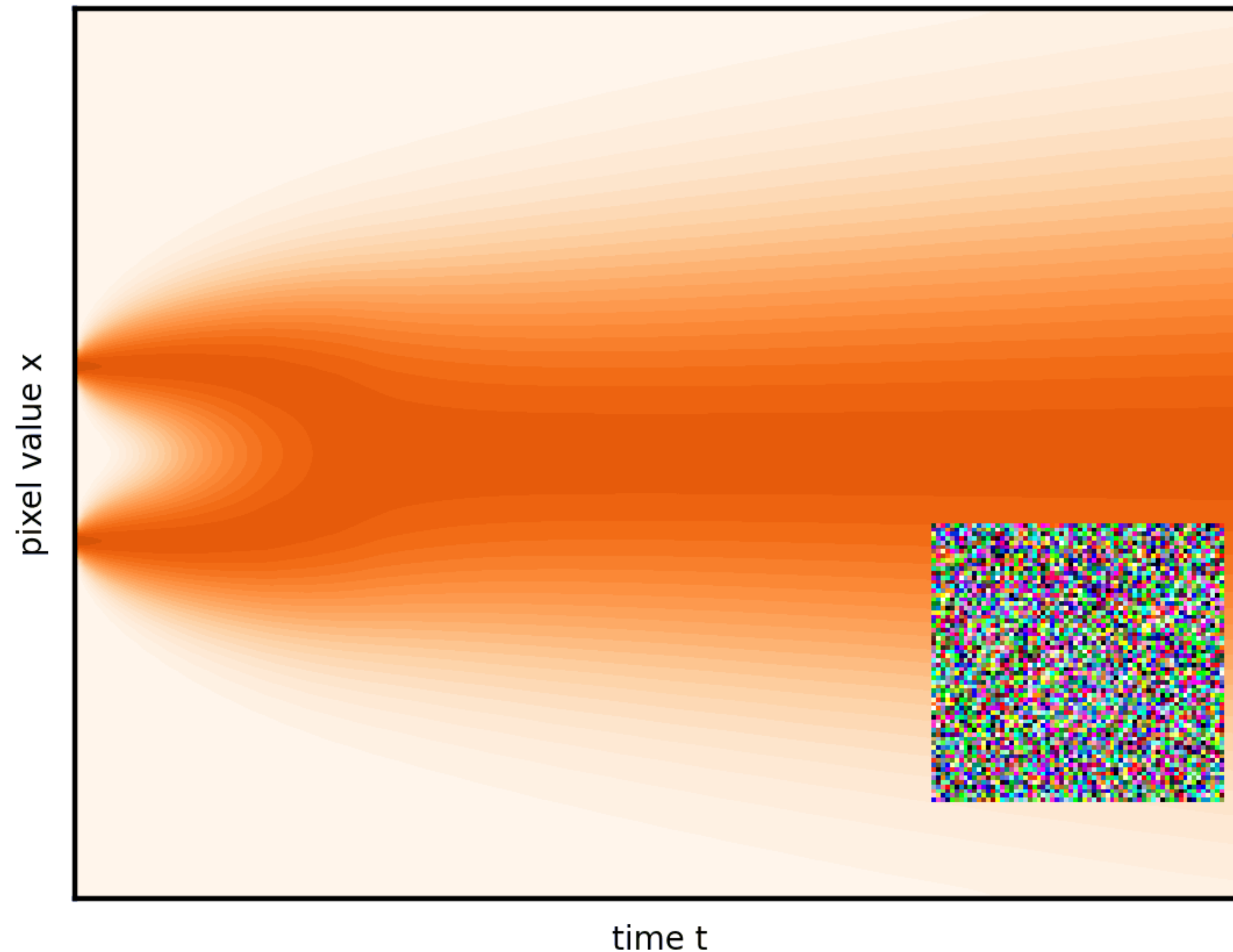
Flow matching model training



Sampling by solving the flow ODE

Flow ODE

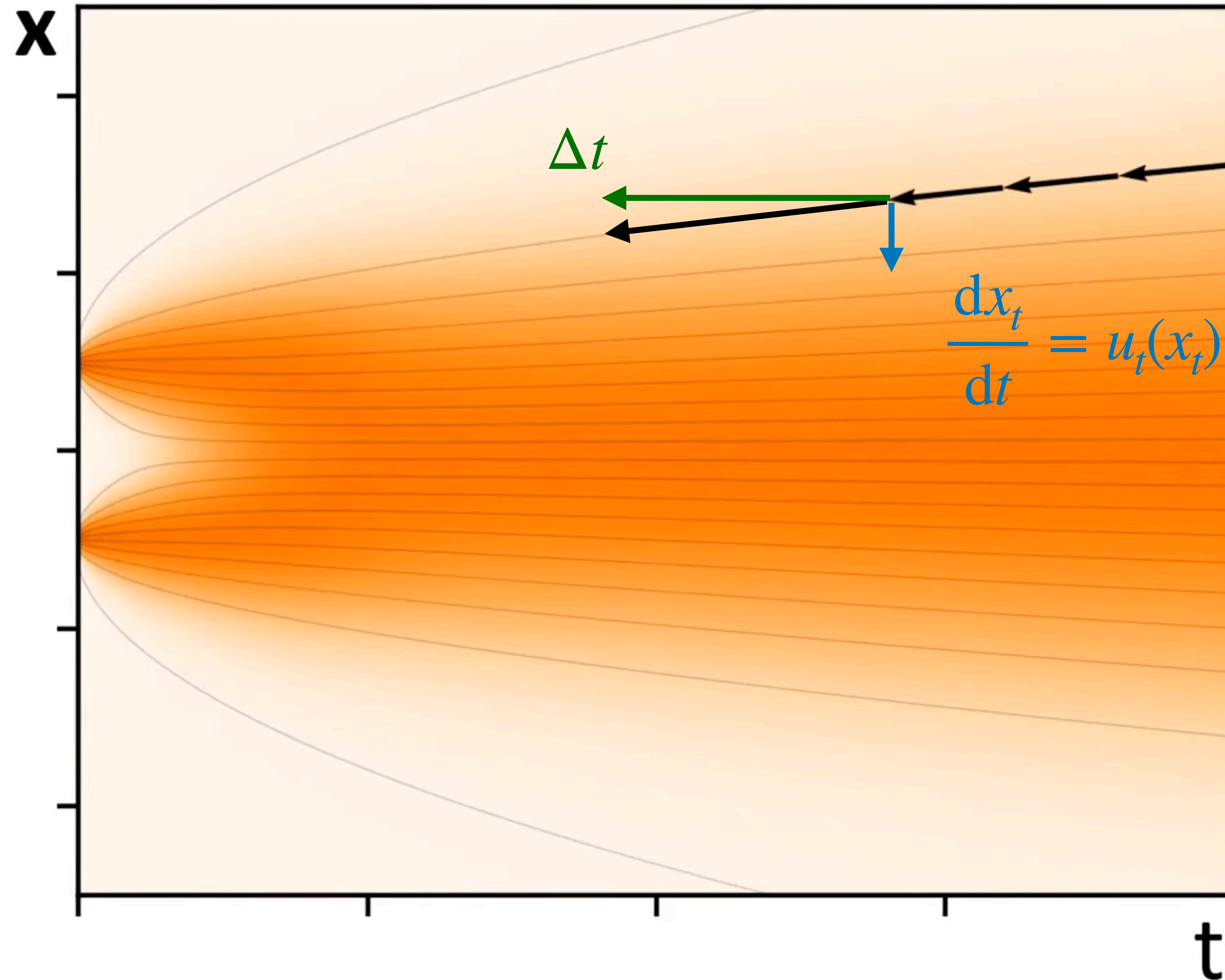
$$\frac{dx_t}{dt} = u_t^\theta(x_t)$$



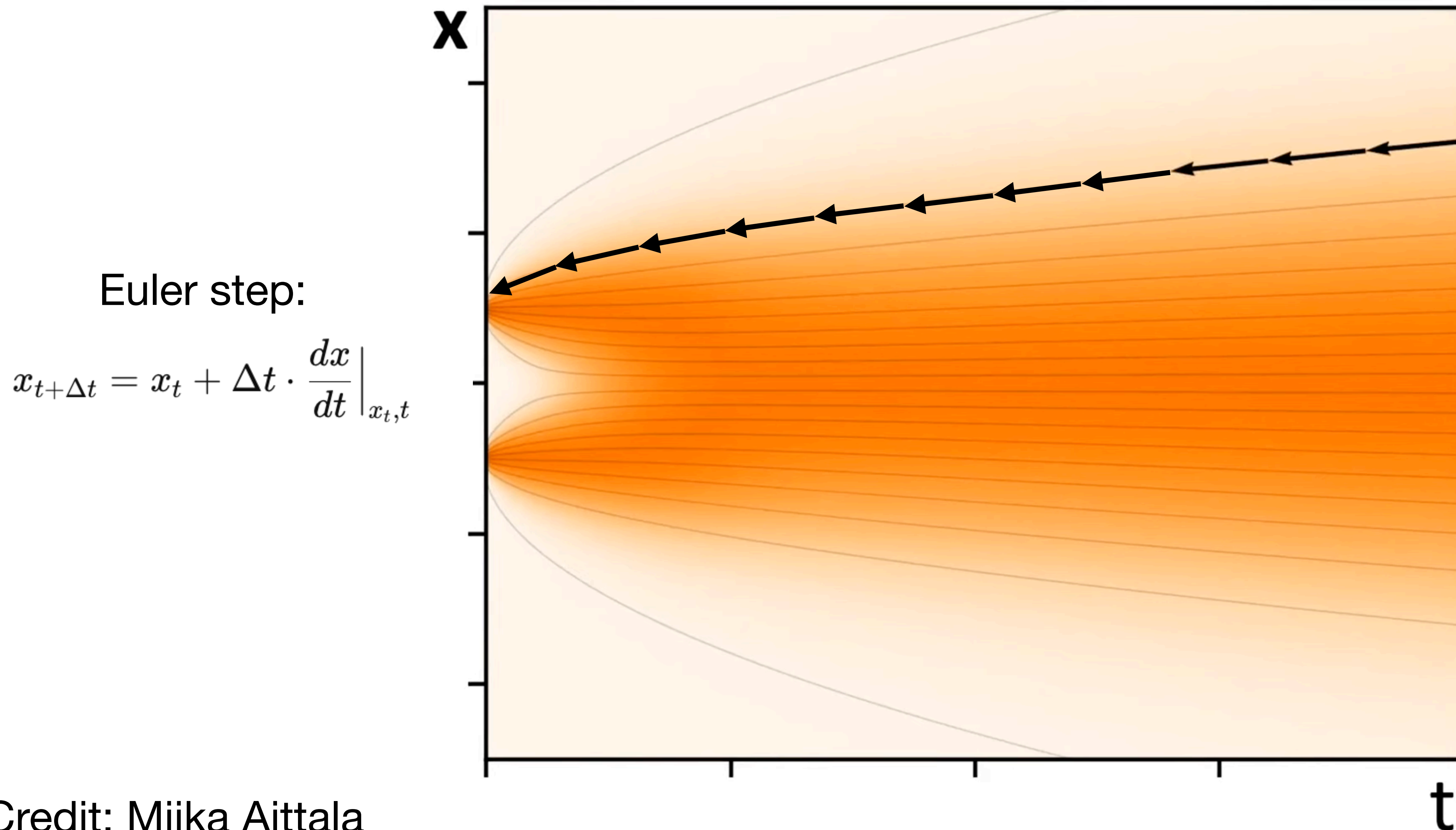
Solving the flow ODE with discretization

Euler step:

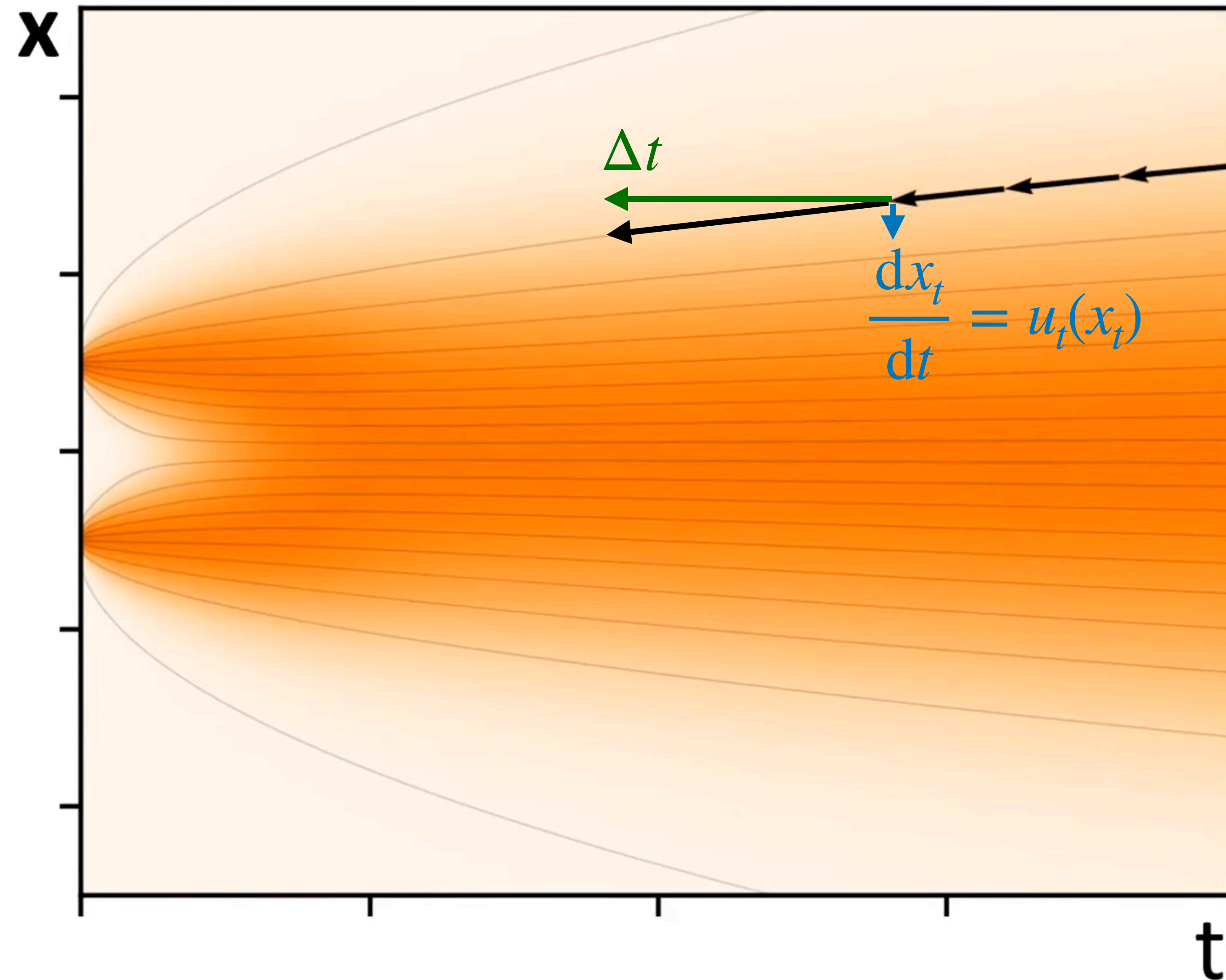
$$x_{t+\Delta t} = x_t + \Delta t \cdot \left. \frac{dx}{dt} \right|_{x_t, t}$$



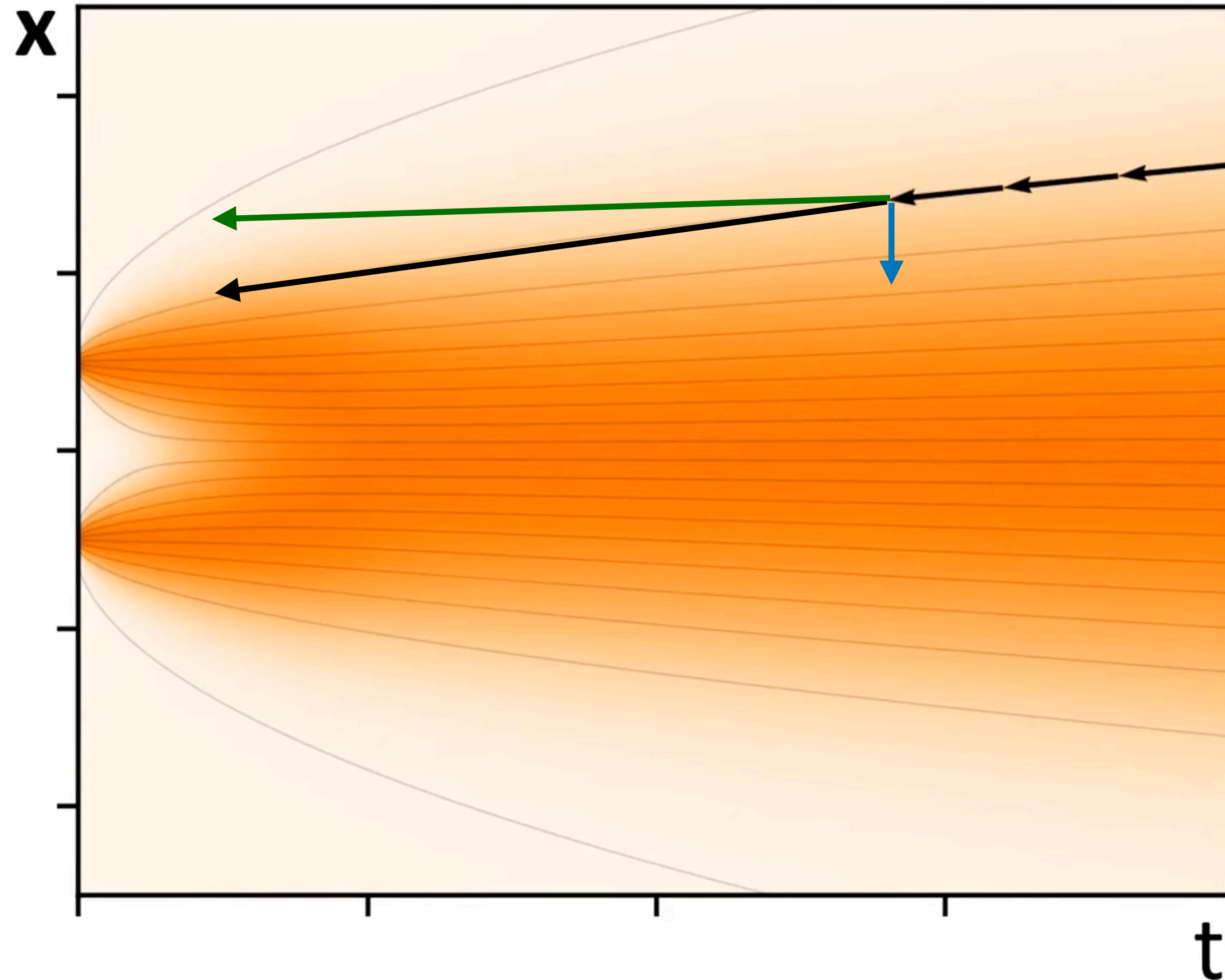
Solving the flow ODE with discretization



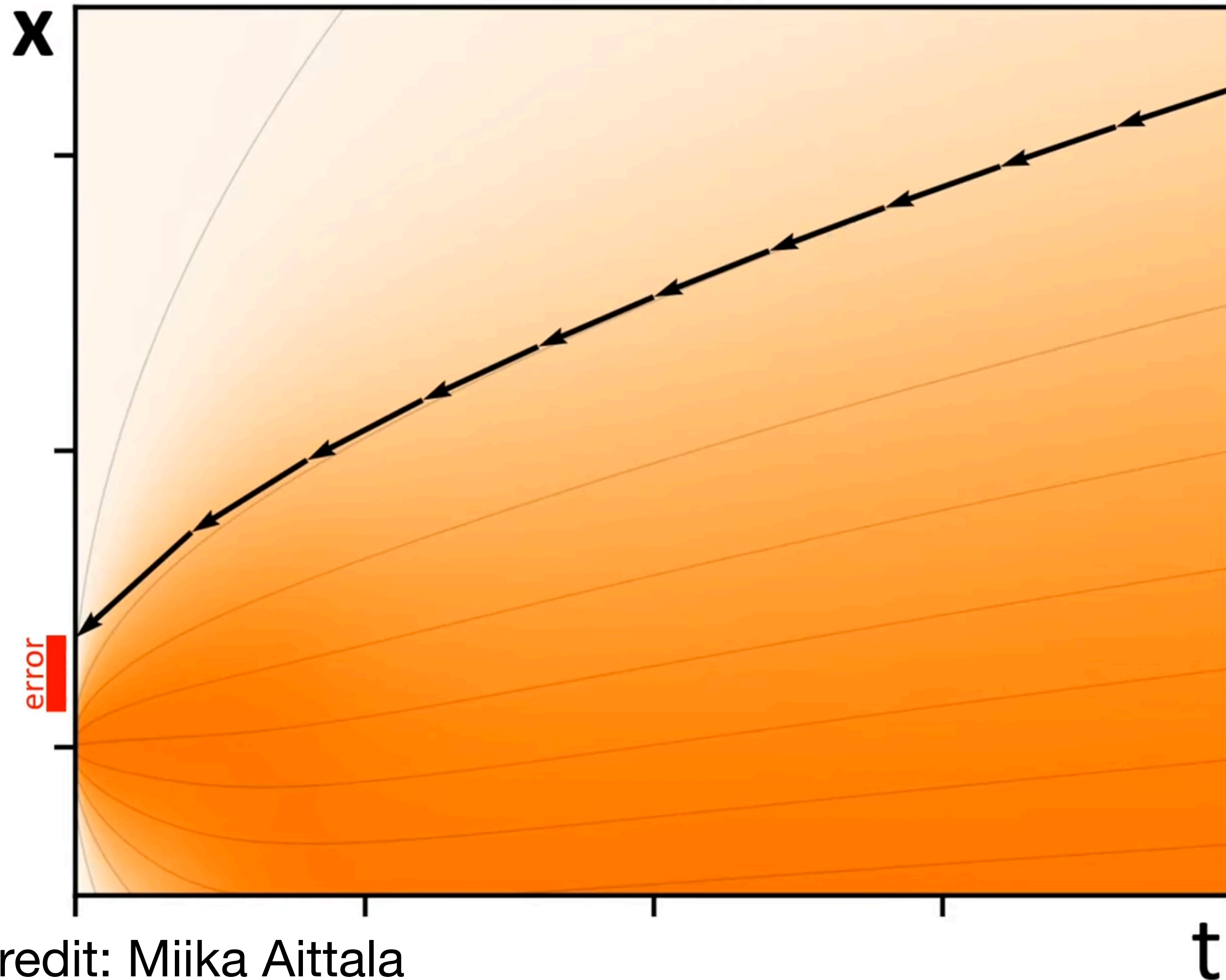
What can go wrong with this sampling process



What can go wrong with this sampling process

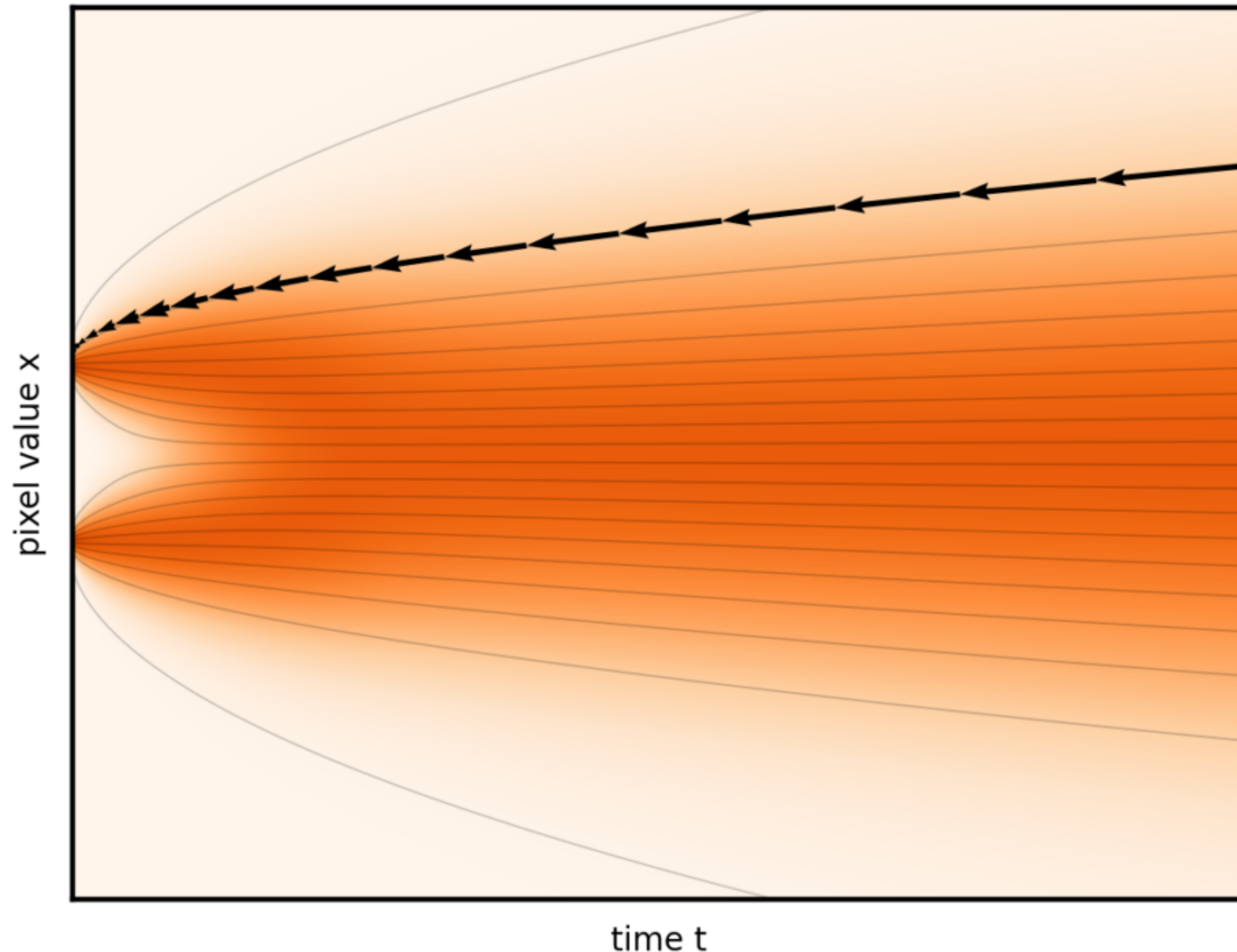


Error Sources when Solving the flow ODE



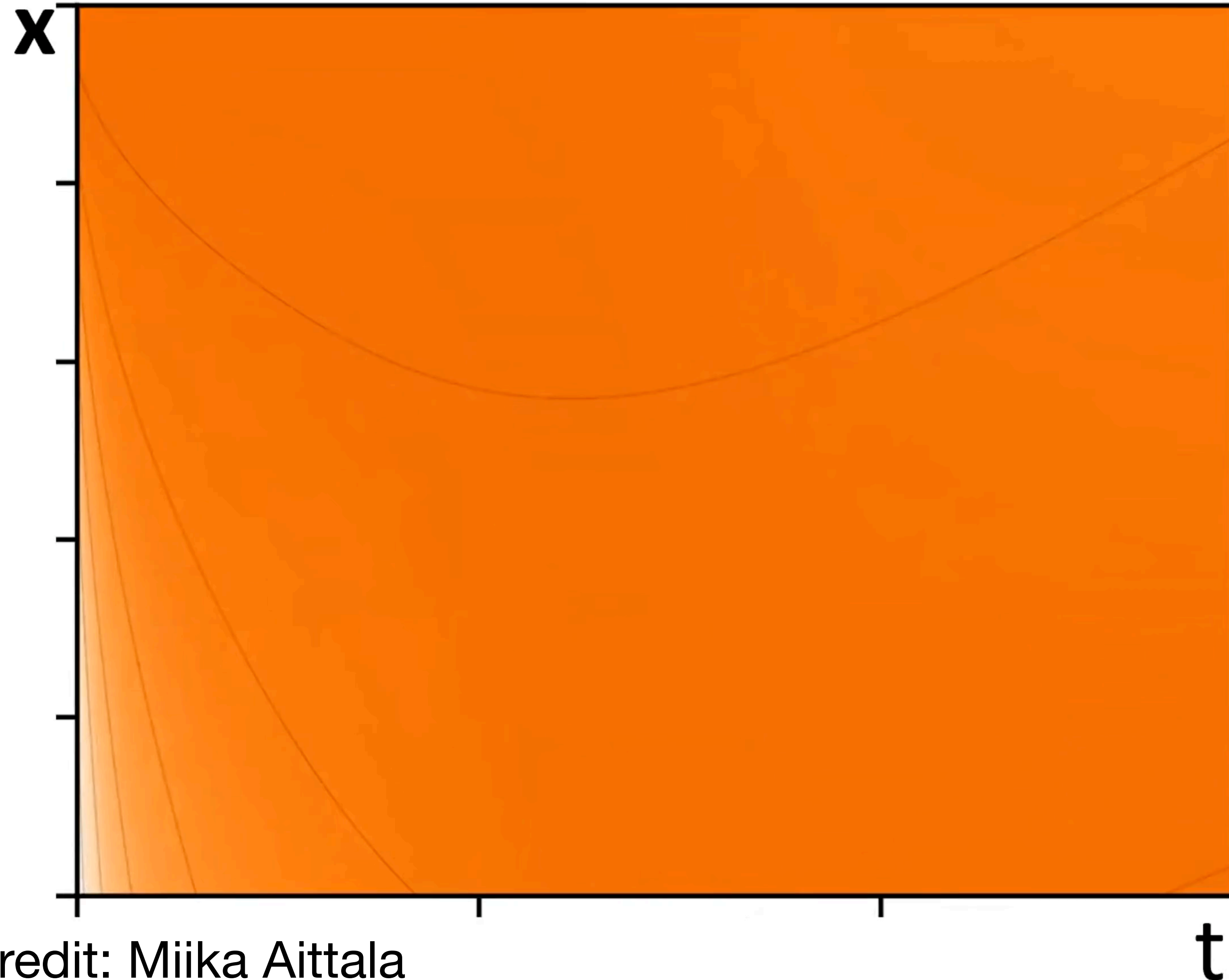
1. Truncation error: fails to approximate ideal trajectory by finite steps.
1. Naive solution: sampling with more steps.

Smart Time Step Schedule



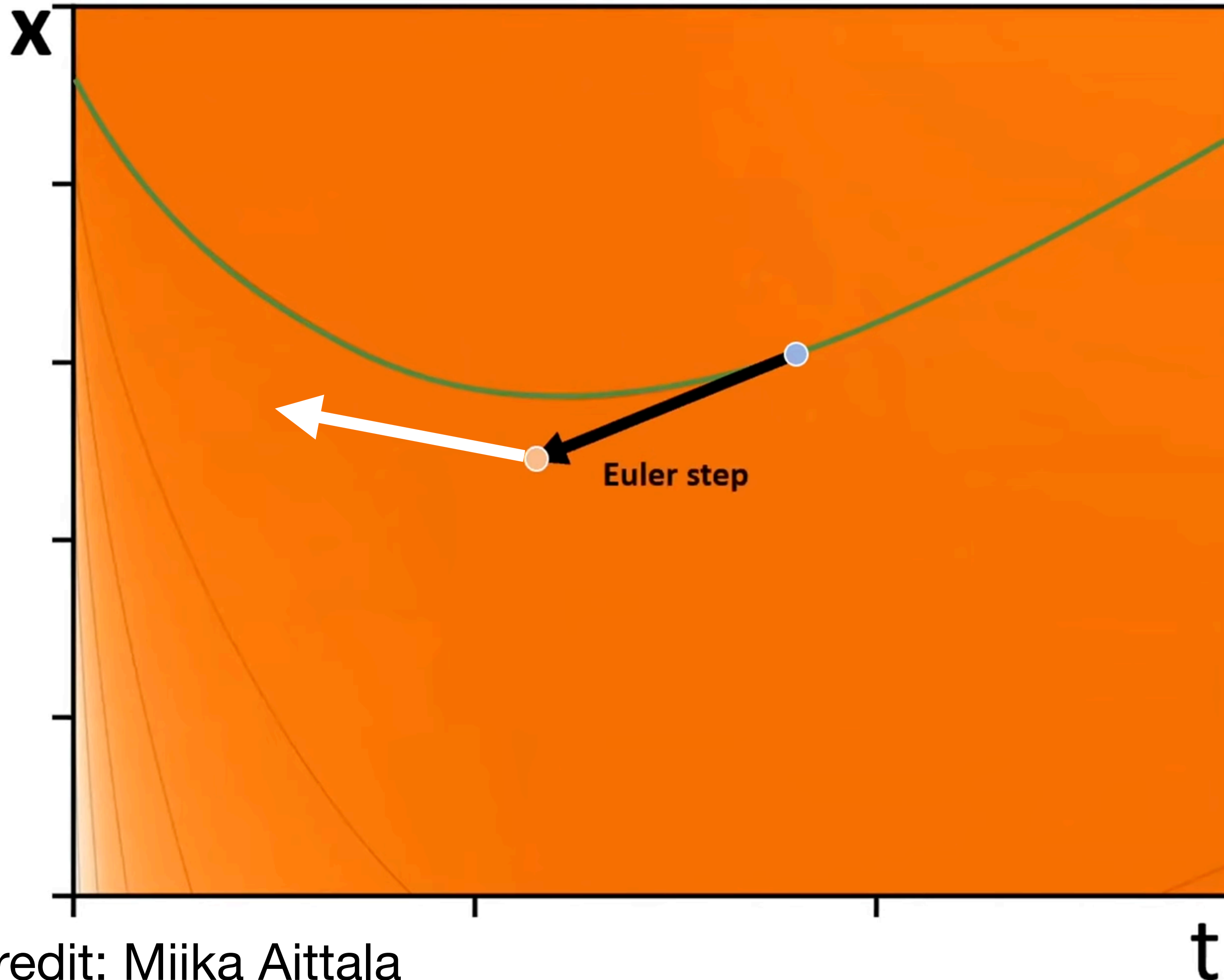
1. Truncation error: fails to approximate ideal trajectory by finite steps.
 1. Naive solution: sampling with more steps.
 2. Time steps are long at high noise levels and short at low noise levels

Advanced ODE Solvers



1. Truncation error: fails to approximate ideal trajectory by finite steps.
 1. Naive solution: sampling with more steps.
 2. Time steps are long at high noise levels and short at low noise levels
 3. Higher-order ODE solver

Higher-order solvers



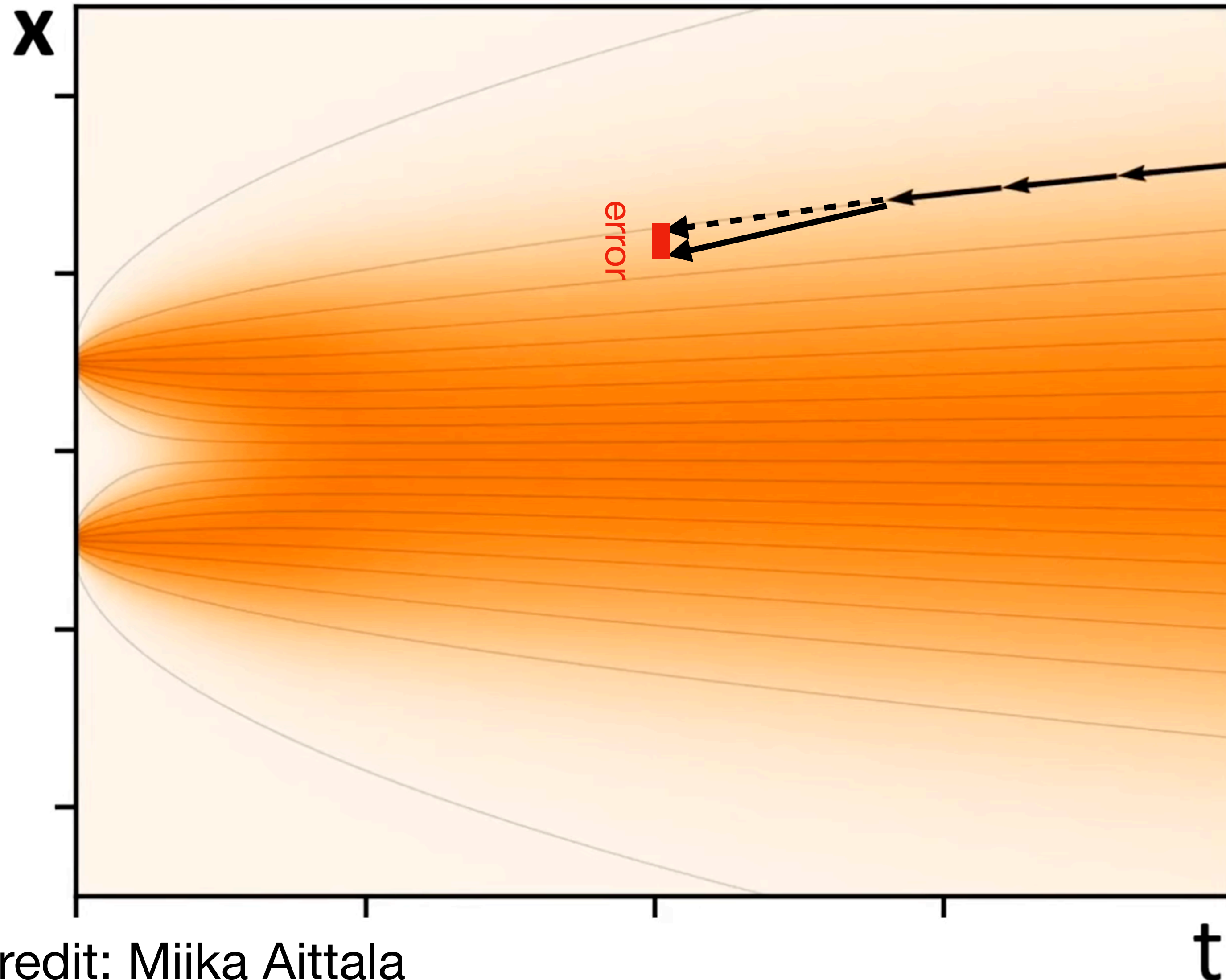
- Clever sub-steps \rightarrow higher accuracy though more cost
- 2nd order Heun method as an example

Higher-order solvers



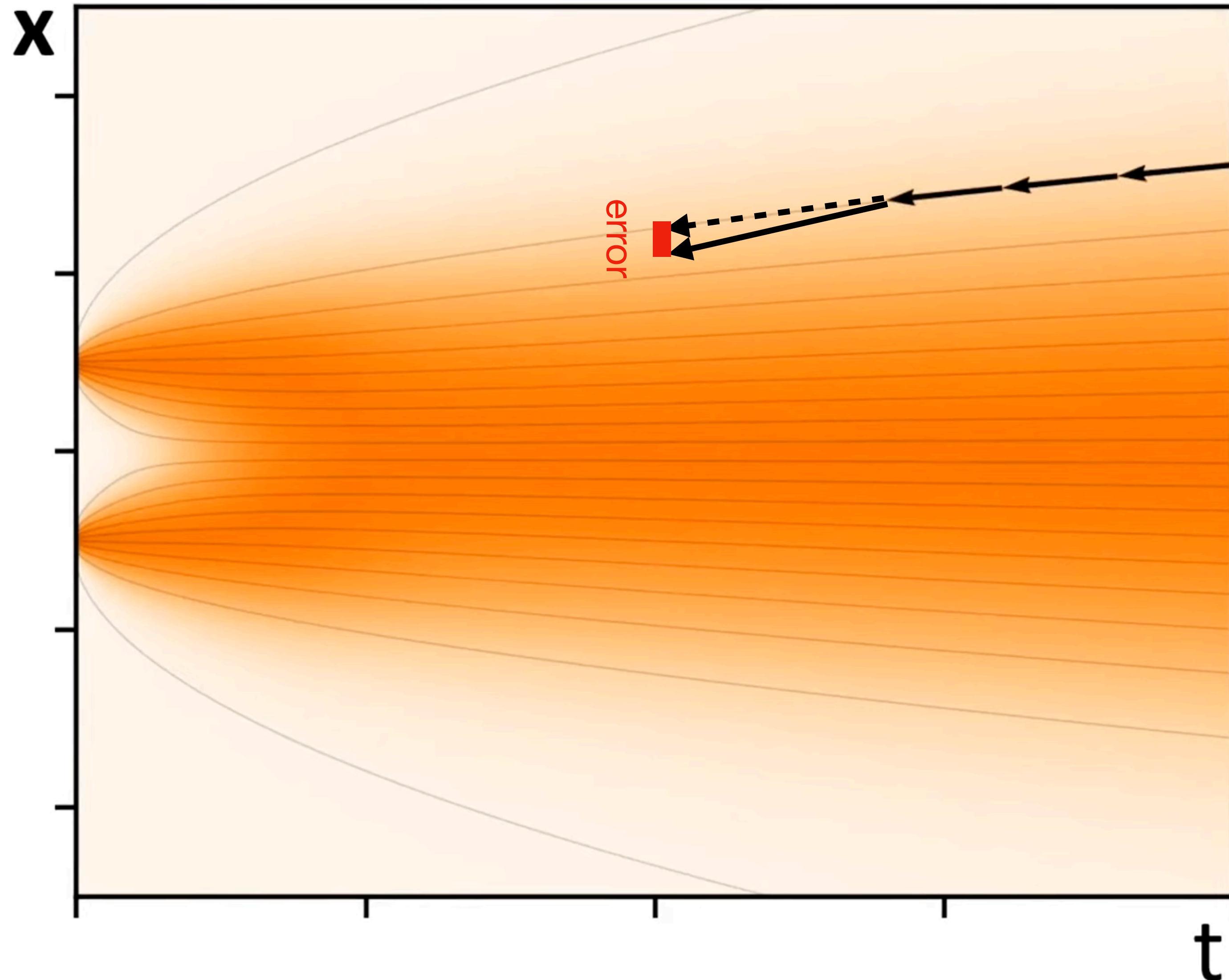
- Clever sub-steps -> higher accuracy though more cost
- 2nd order Heun method as an example

Error Sources when Solving the flow ODE



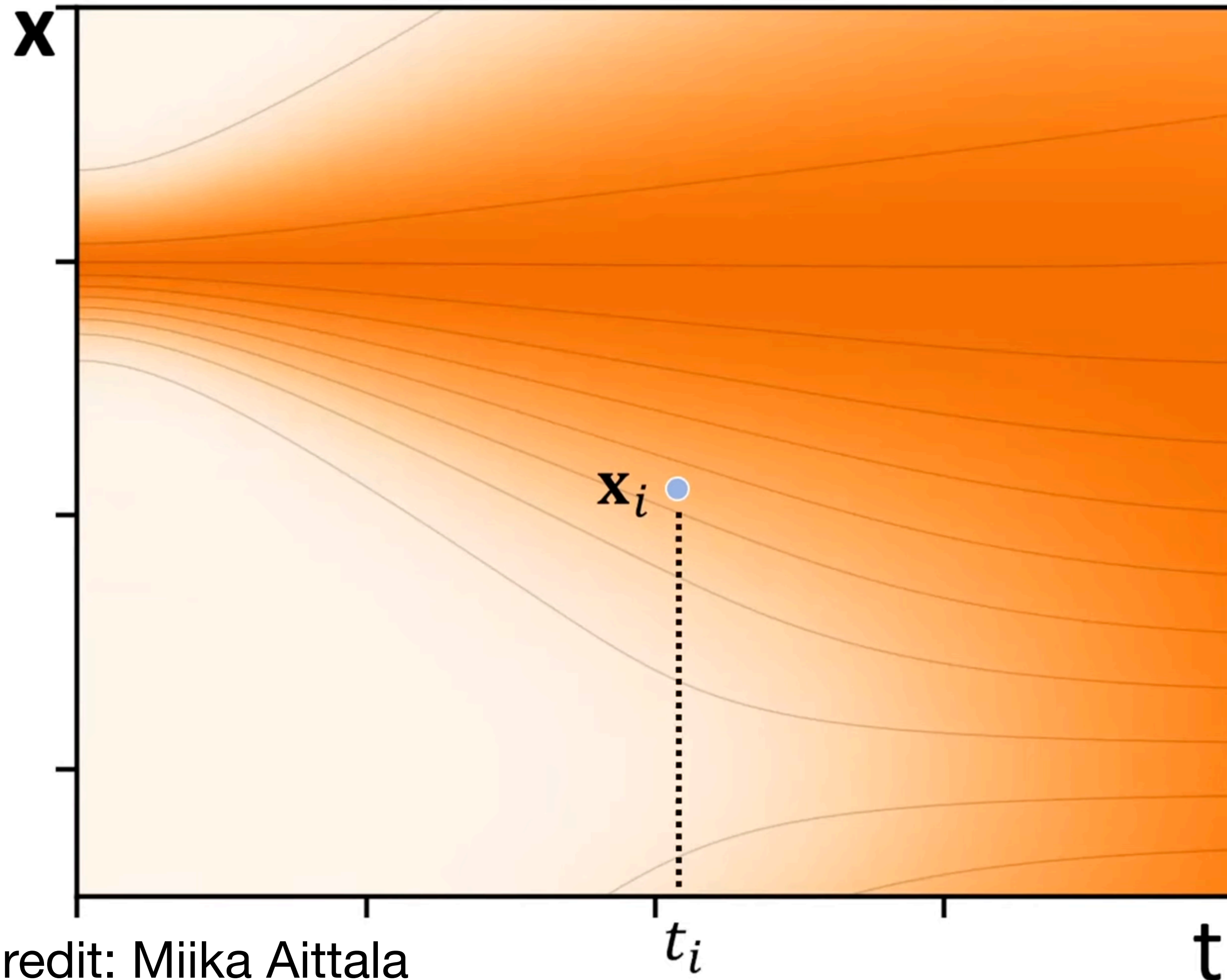
1. Truncation error: fails to approximate ideal trajectory by finite steps.
2. Model fails to approximate the marginal flow.

Stochastic Sampler



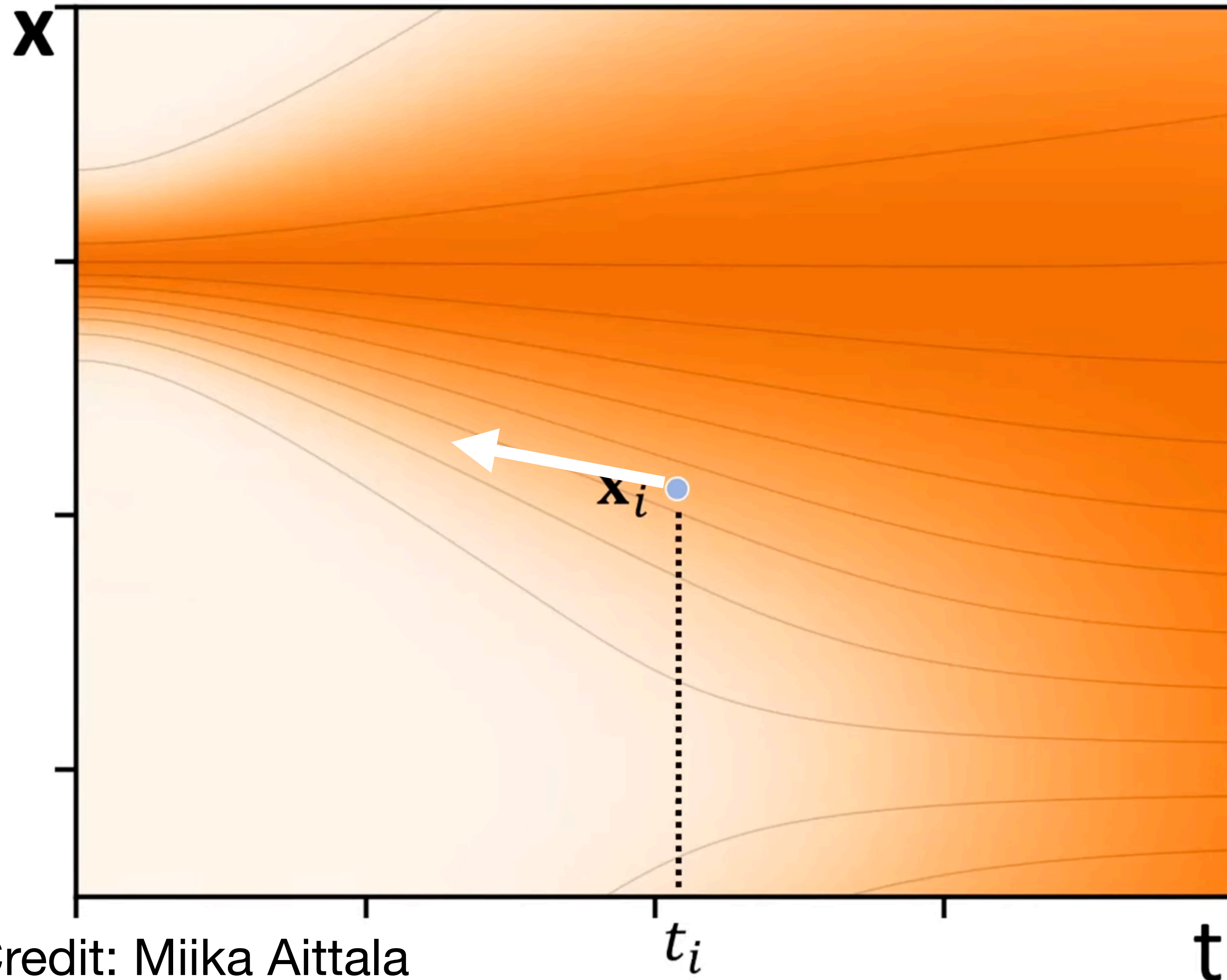
1. Truncation error: fails to approximate ideal trajectory by finite steps.
2. Model fails to approximate the marginal flow.
 1. Stochastic sampler (SDE) injects fresh noise throughout the evolution in addition to reducing the noise.

Stochastic Sampler: why does it improve quality



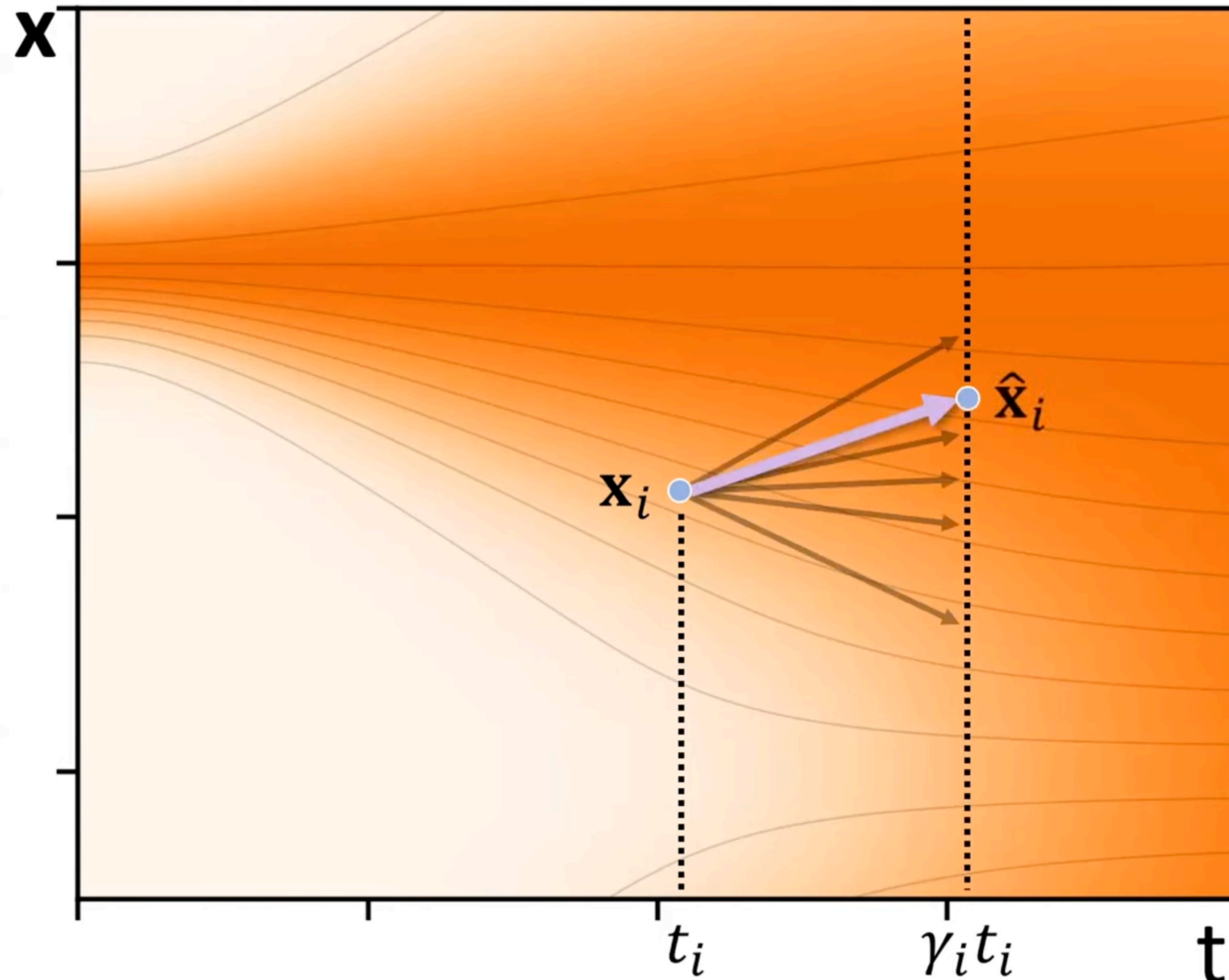
1. Truncation error: fails to approximate ideal trajectory by finite steps.
2. Model fails to approximate the marginal flow.
 1. Stochastic sampler (SDE) injects fresh noise throughout the evolution in addition to reducing the noise.

Stochastic Sampler: why does it improve quality



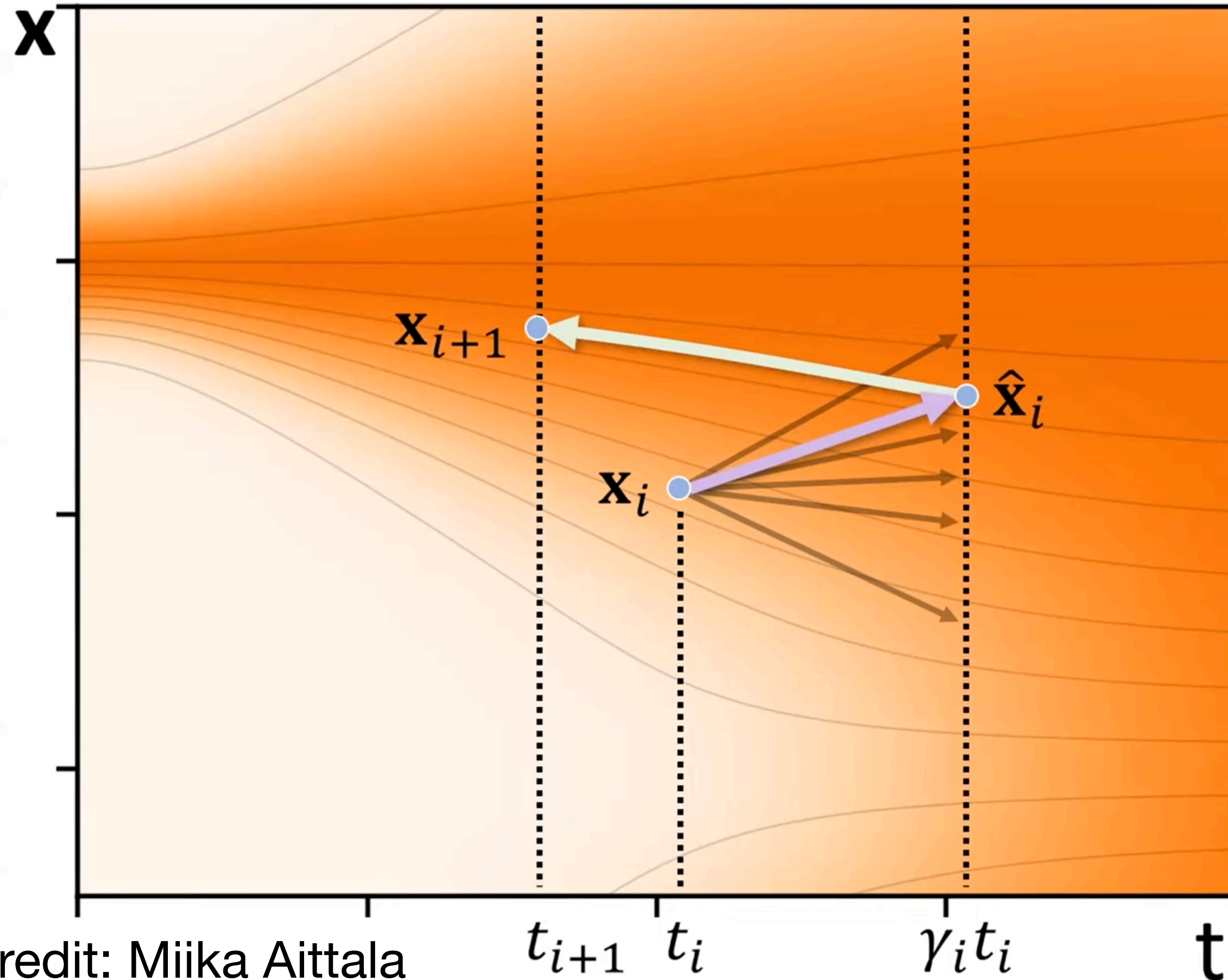
1. Truncation error: fails to approximate ideal trajectory by finite steps.
2. Model fails to approximate the marginal flow.
 1. Stochastic sampler (SDE) injects fresh noise throughout the evolution in addition to reducing the noise.

Stochastic Sampler: why does it improve quality



1. Truncation error: fails to approximate ideal trajectory by finite steps.
2. Model fails to approximate the marginal flow.
 1. Stochastic sampler (SDE) injects fresh noise throughout the evolution in addition to reducing the noise.

Stochastic Sampler: why does it improve quality

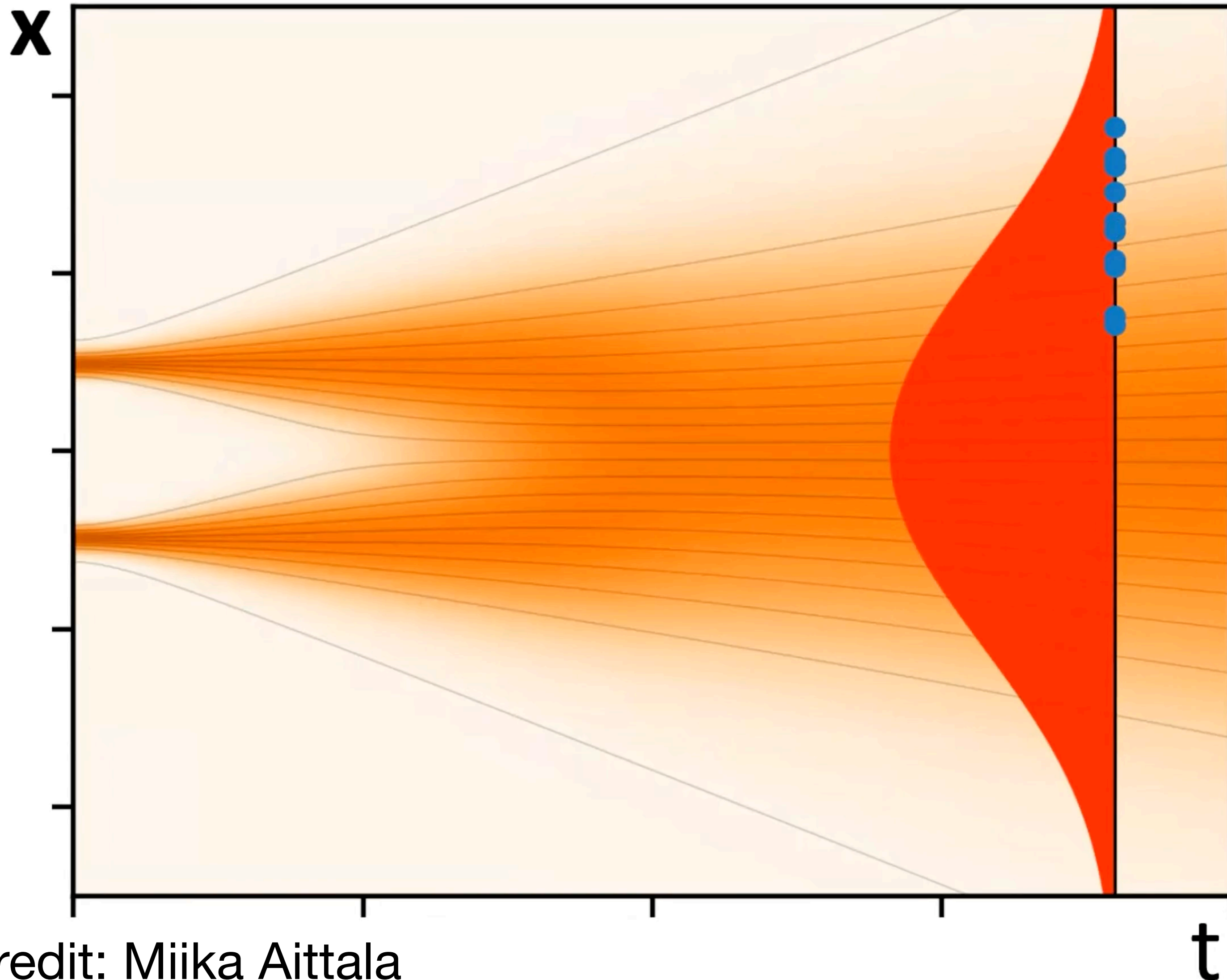


1. Truncation error: fails to approximate ideal trajectory by finite steps.

2. Model fails to approximate the marginal flow.

1. Stochastic sampler (SDE) injects fresh noise throughout the evolution in addition to reducing the noise.

Stochastic Sampler Helps Explore the Distribution

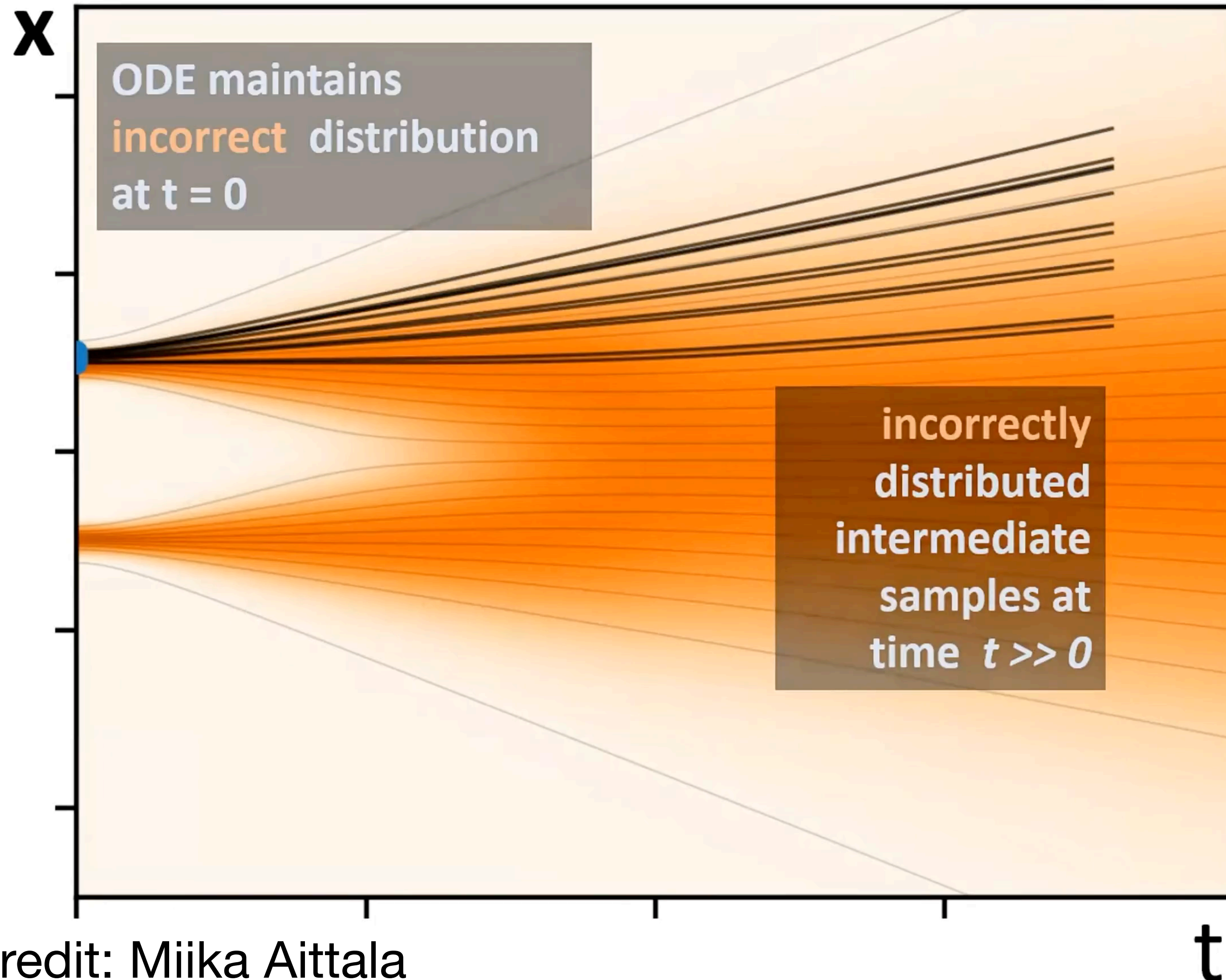


1. Truncation error: fails to approximate ideal trajectory by finite steps.

2. Model fails to approximate the marginal flow.

1. Stochastic sampler (SDE) injects fresh noise throughout the evolution in addition to reducing the noise.

Stochastic Sampler Helps Explore the Distribution



1. Truncation error: fails to approximate ideal trajectory by finite steps.
2. Model fails to approximate the marginal flow.
 1. Stochastic sampler (SDE) injects fresh noise throughout the evolution in addition to reducing the noise.

Is stochasticity always helpful?

CIFAR-10: **no**

Imagenet: **yes**

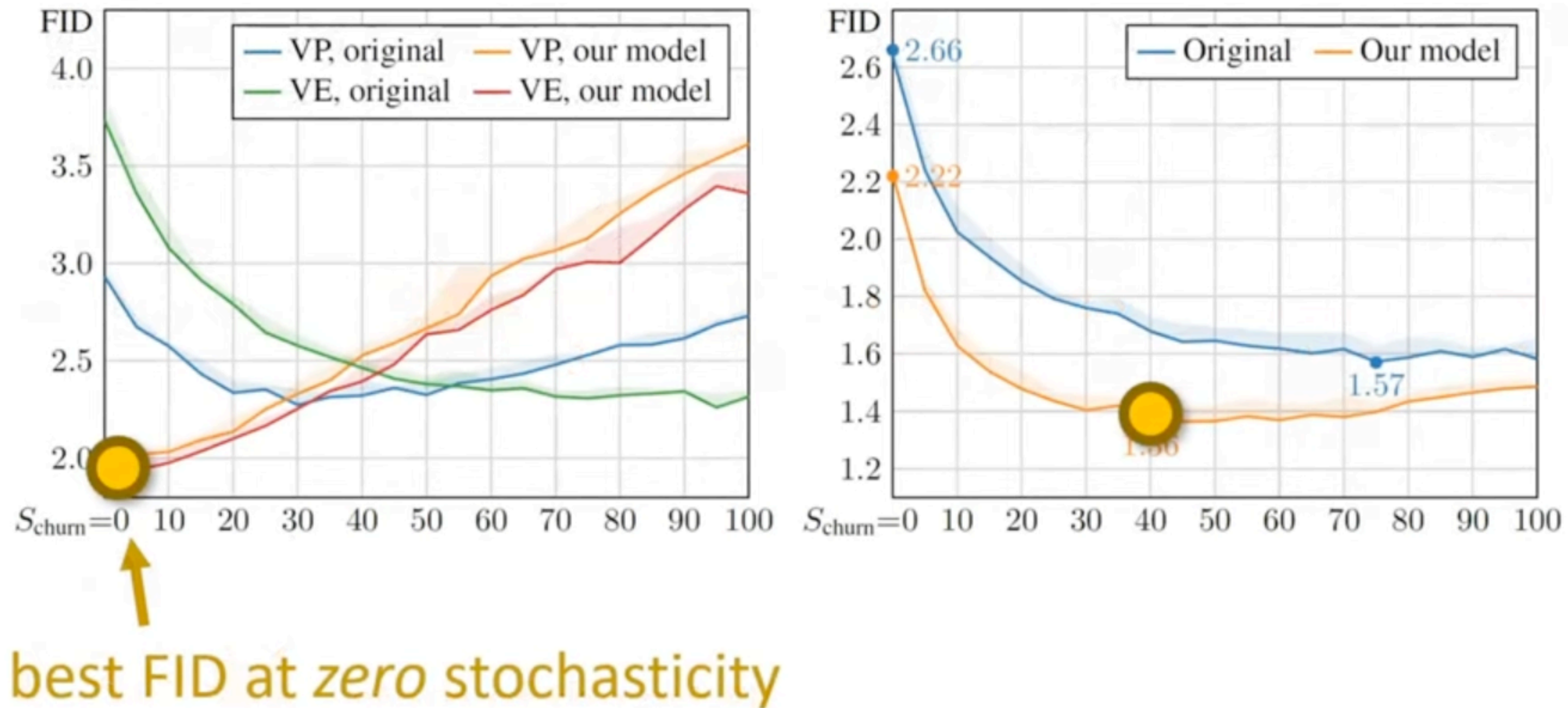


Image Generation by Solving the Flow ODE

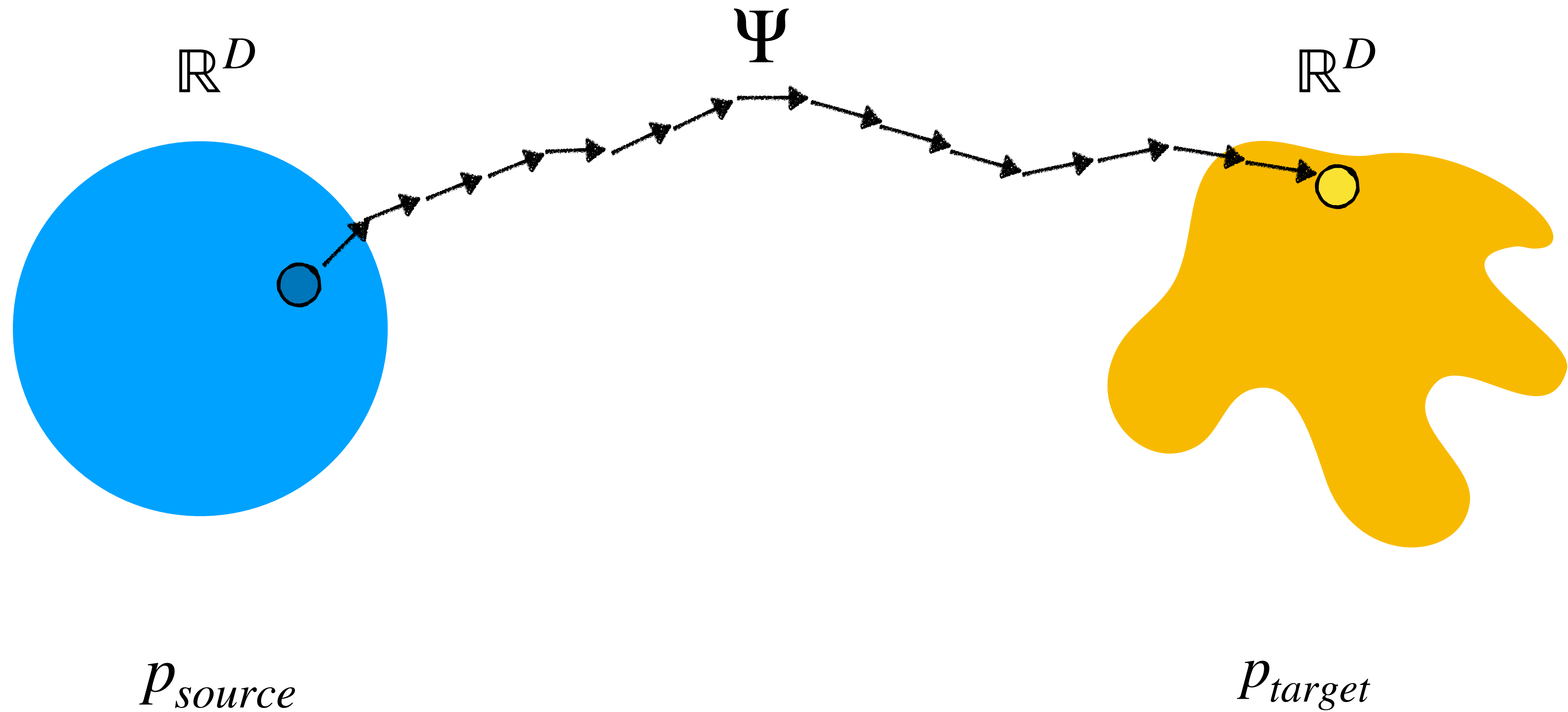
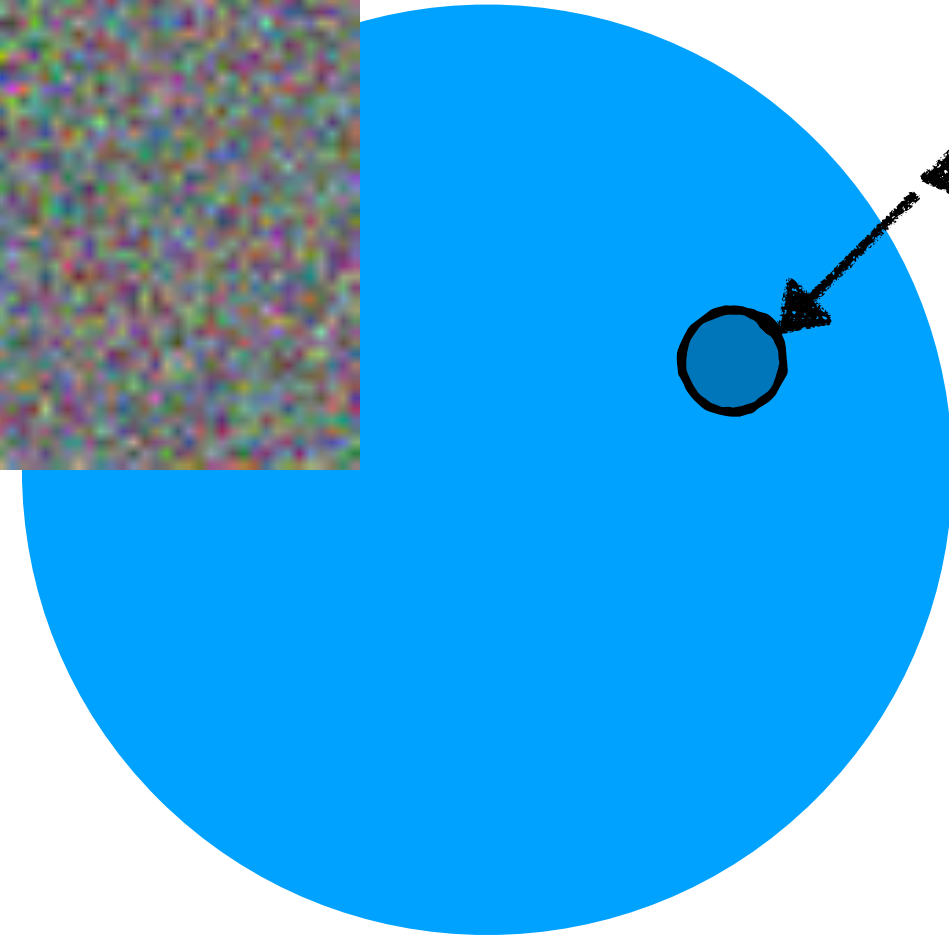


Image Editing with Diffusion Inversion

$c = \text{"A man wearing a suit..."}$

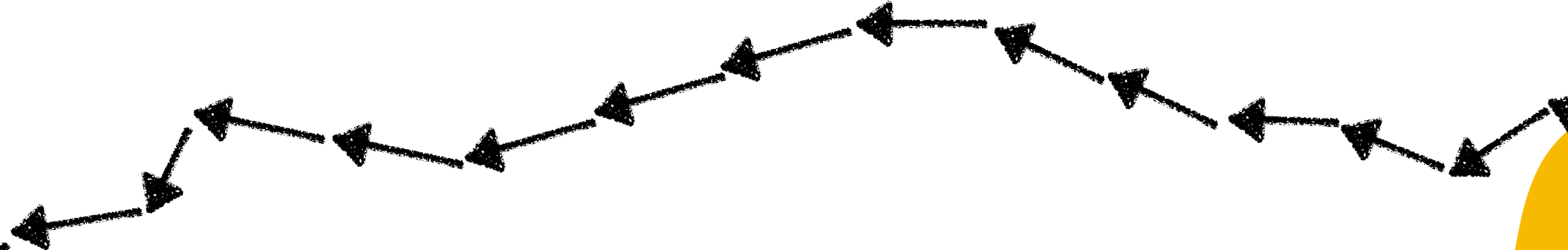


\mathbb{R}^D

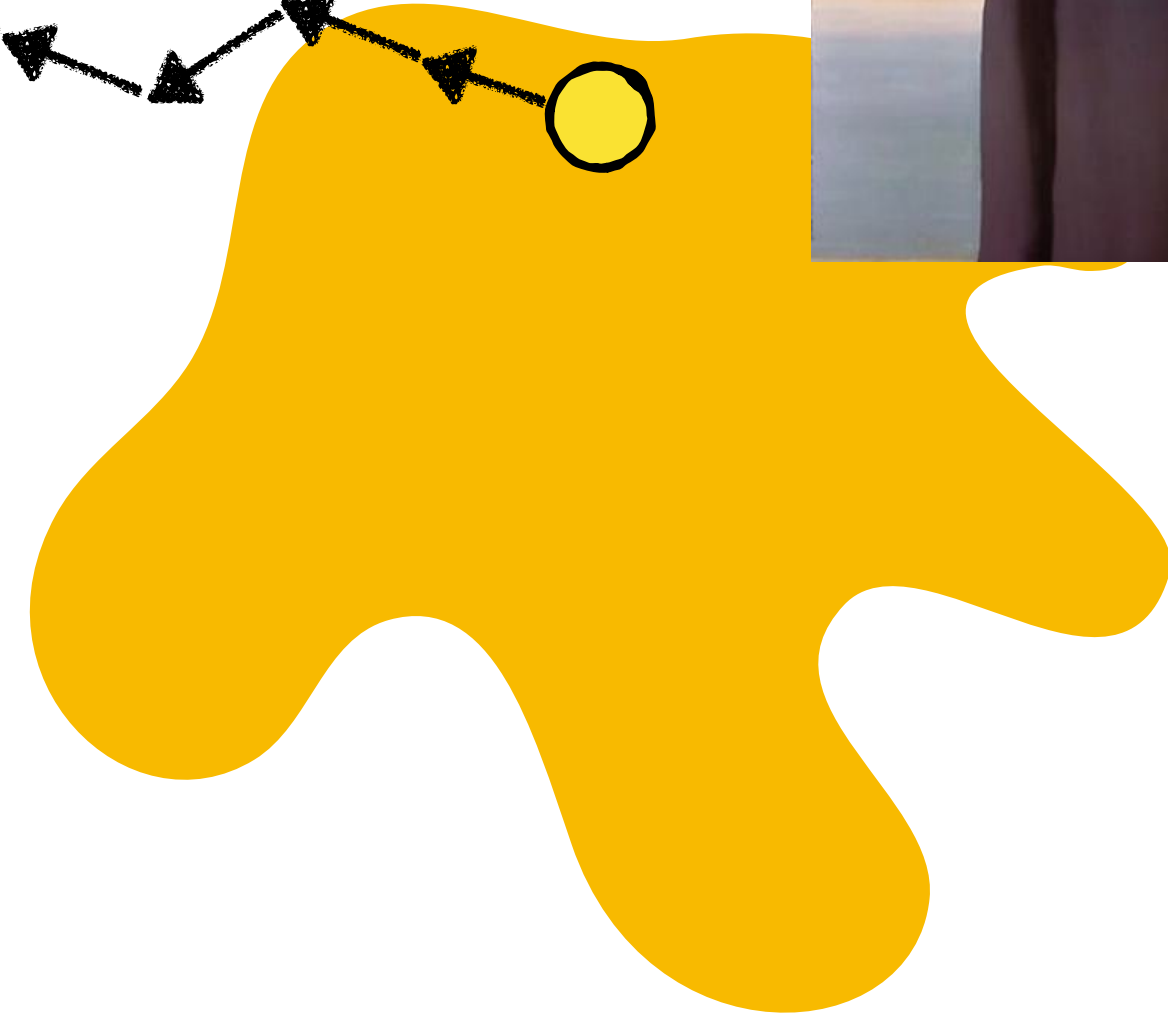


P_{source}

$\Psi^{-1}(x, c)$



\mathbb{R}^D



P_{target}

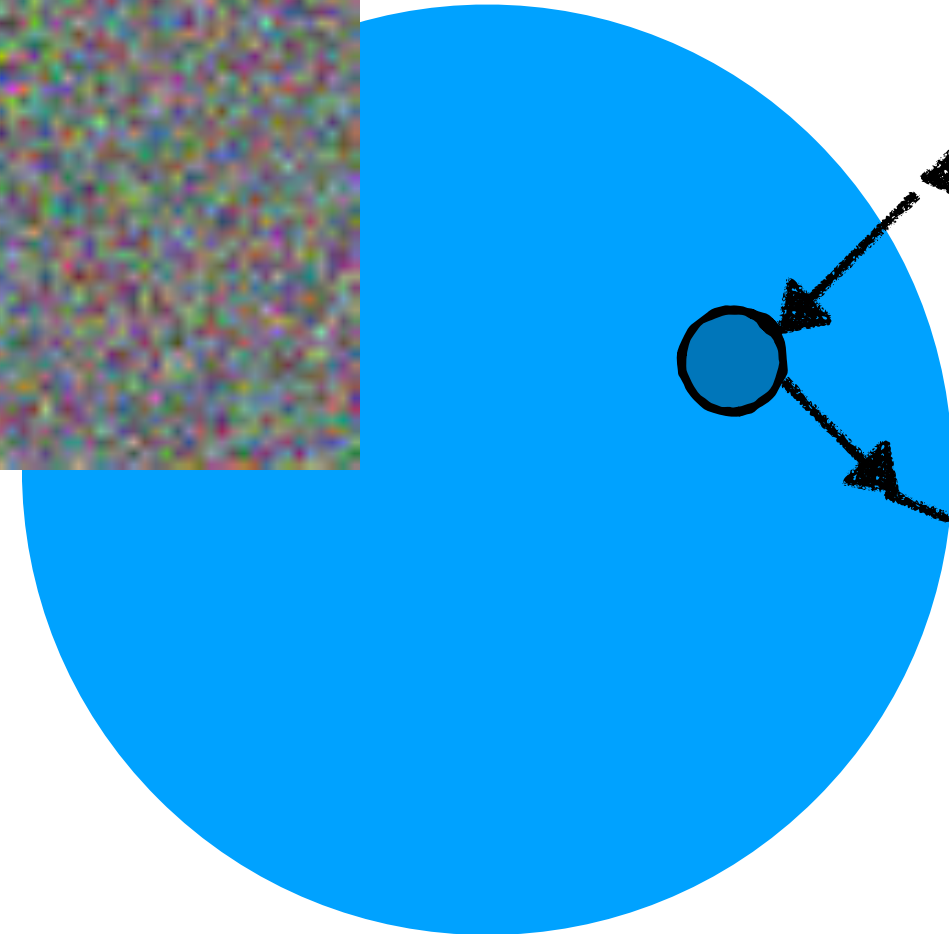


Image Editing with Diffusion Inversion

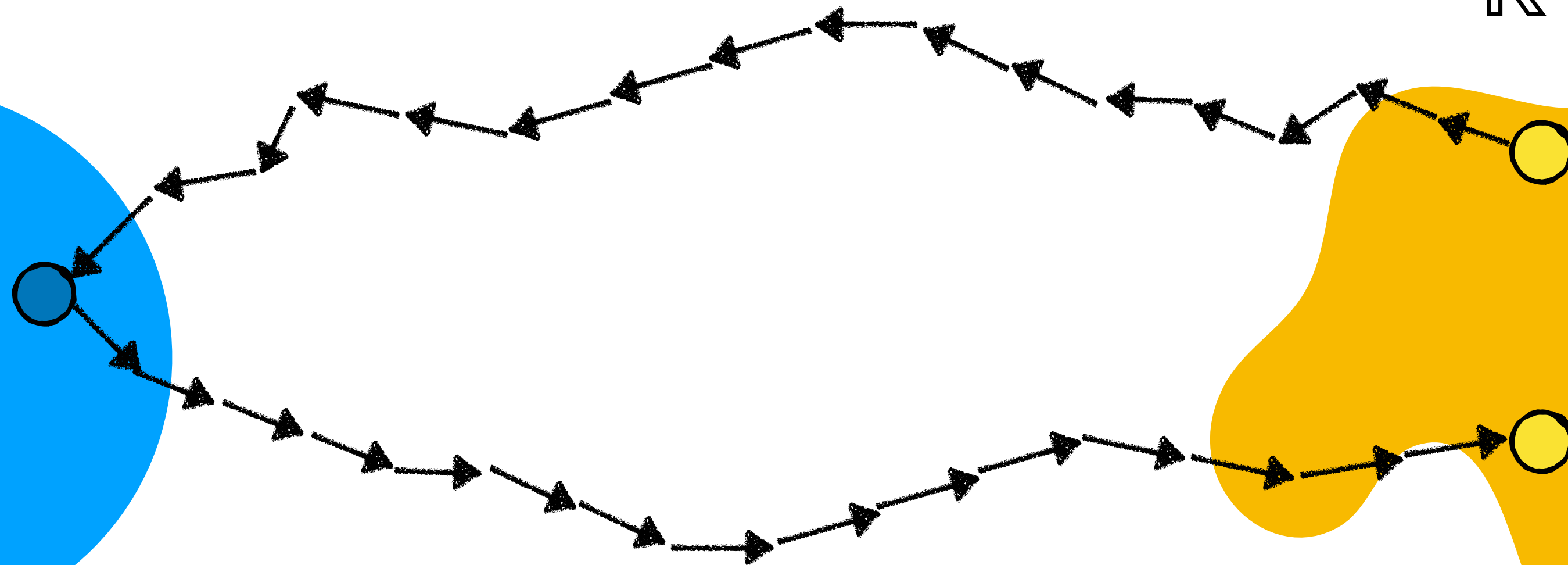
$c = \text{"A man wearing a suit..."}$



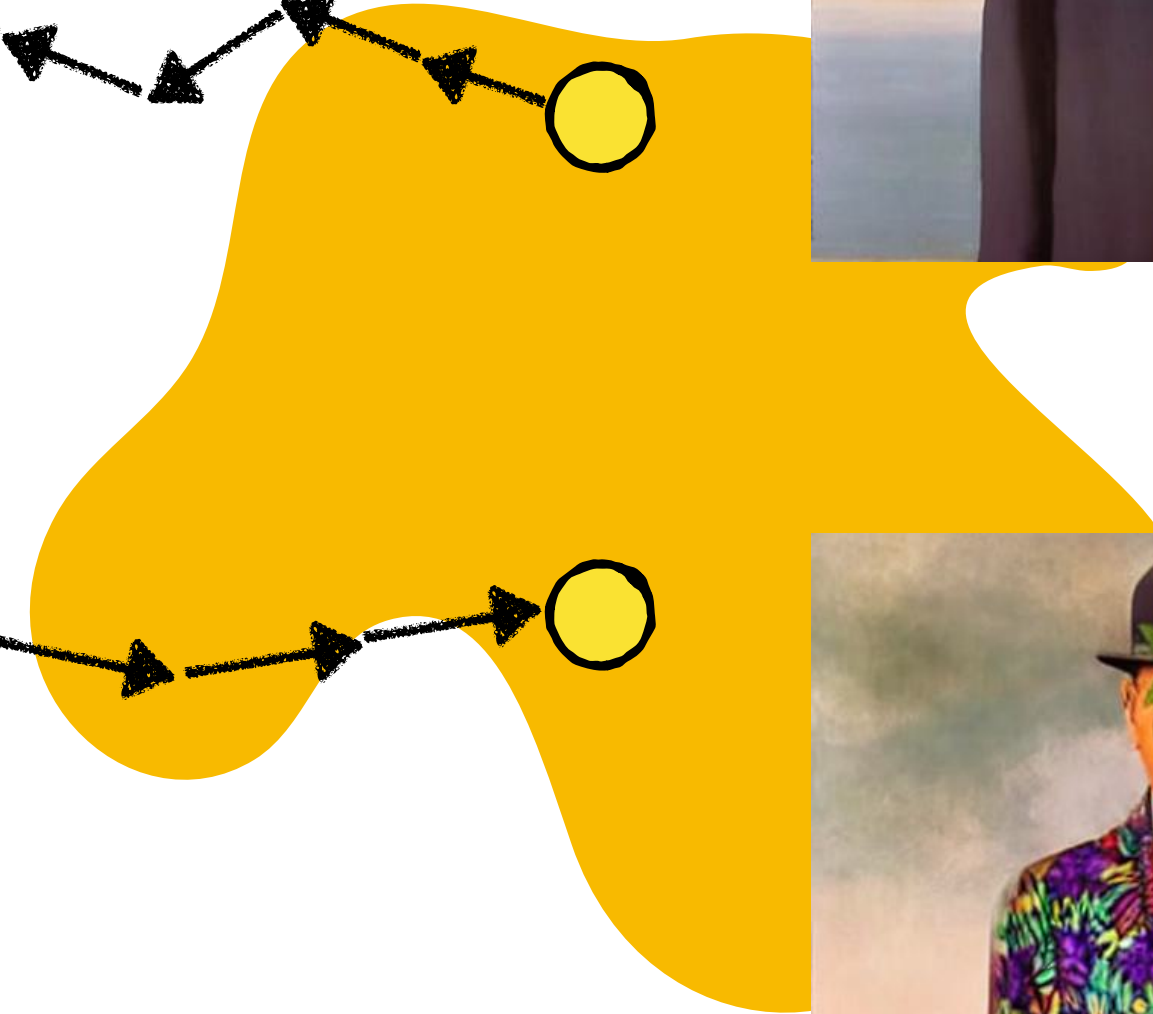
\mathbb{R}^D



$\Psi^{-1}(x, c)$



\mathbb{R}^D



$\Psi(x, c')$

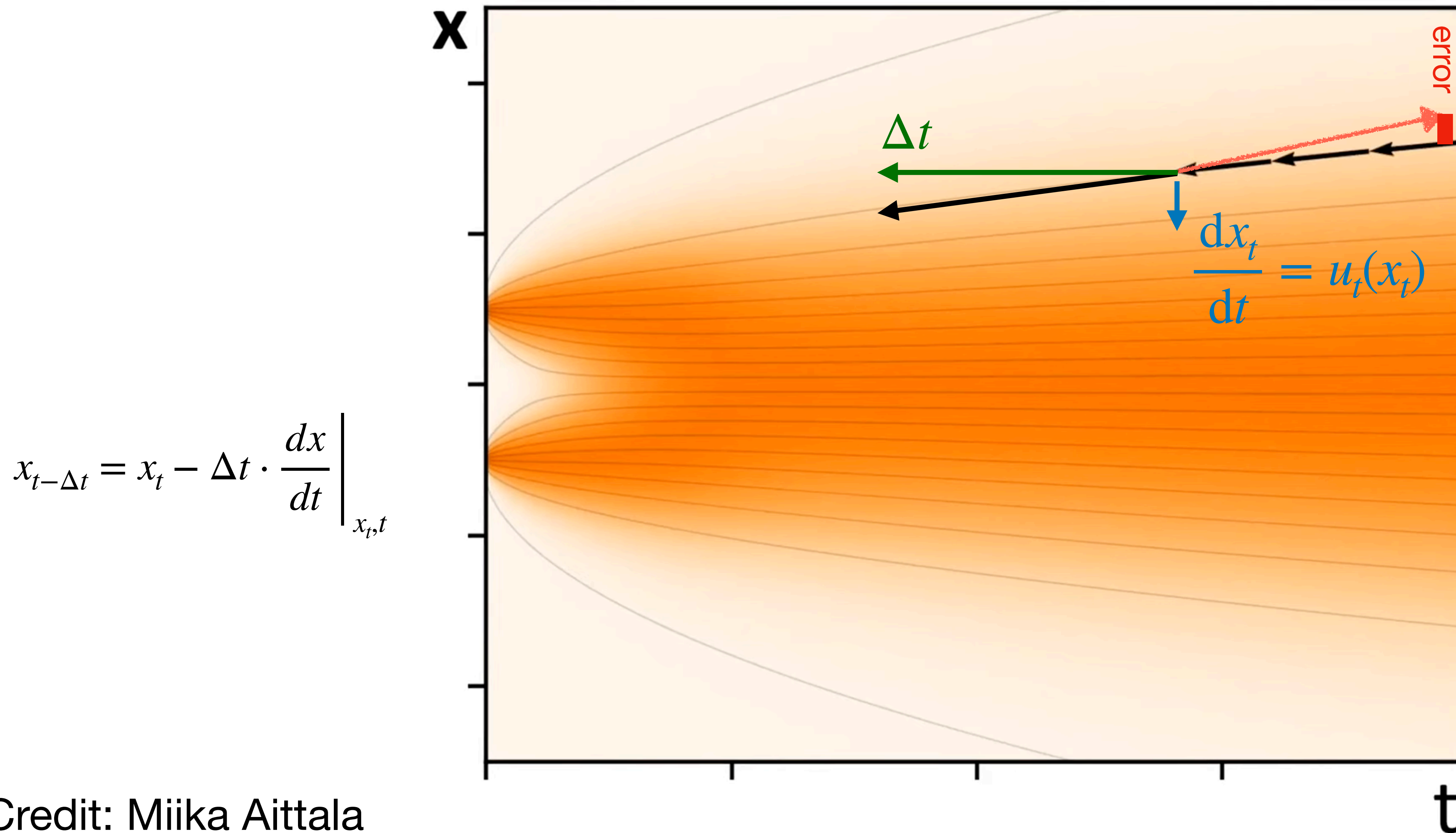
$c = \text{"A man wearing a Hawaiian shirt..."}$



P_{source}

P_{target}

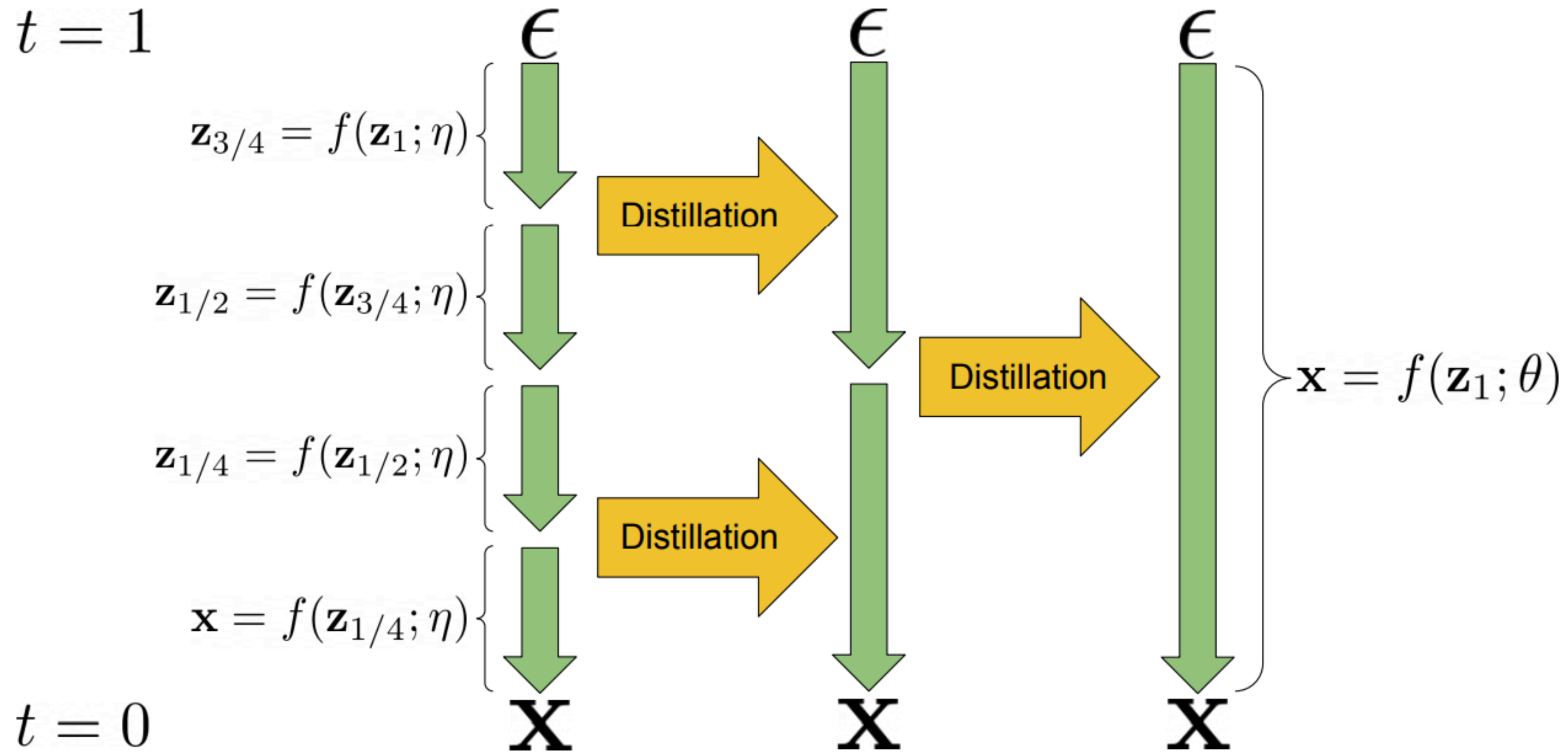
Inverting a diffusion model is not easy



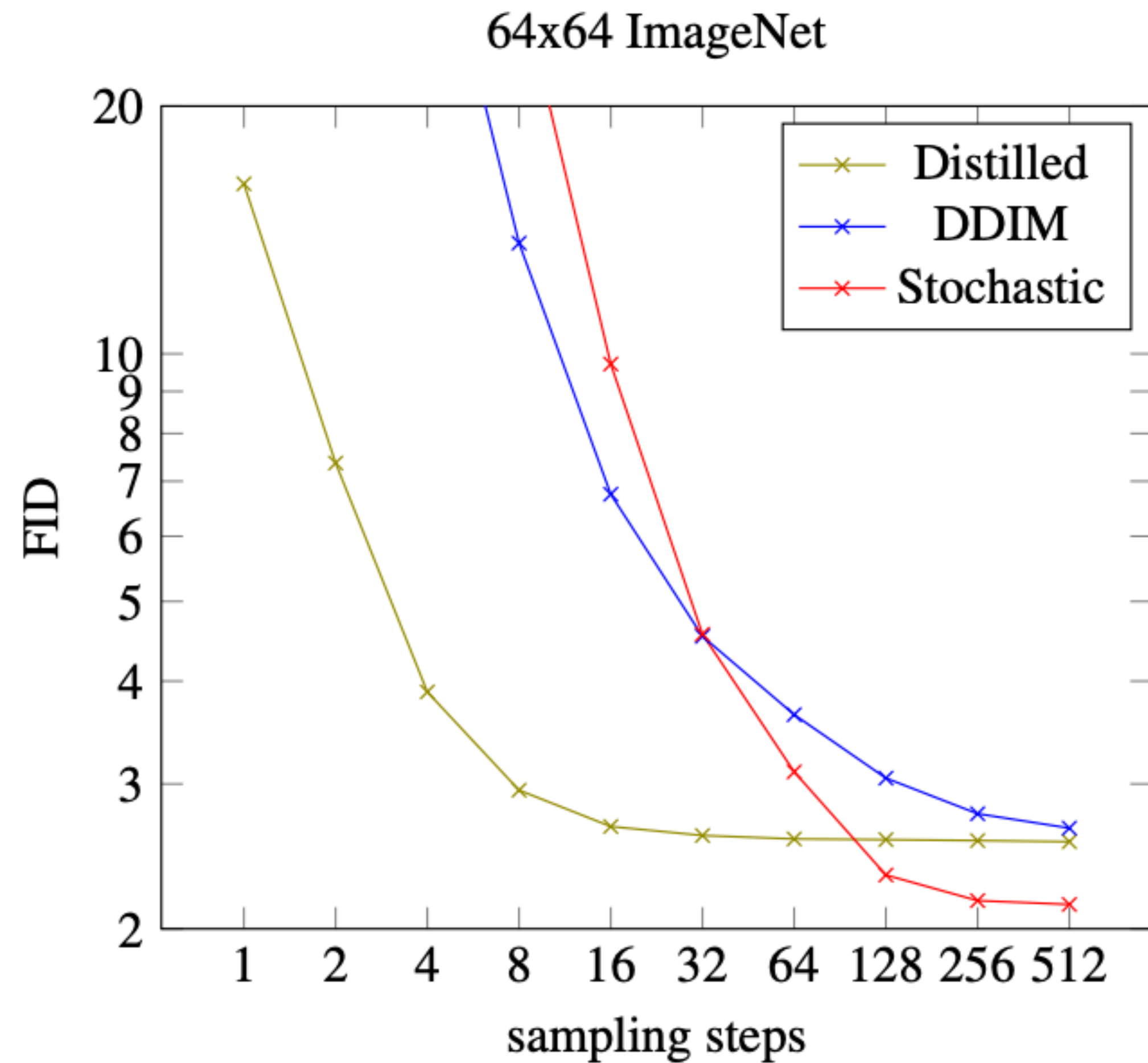
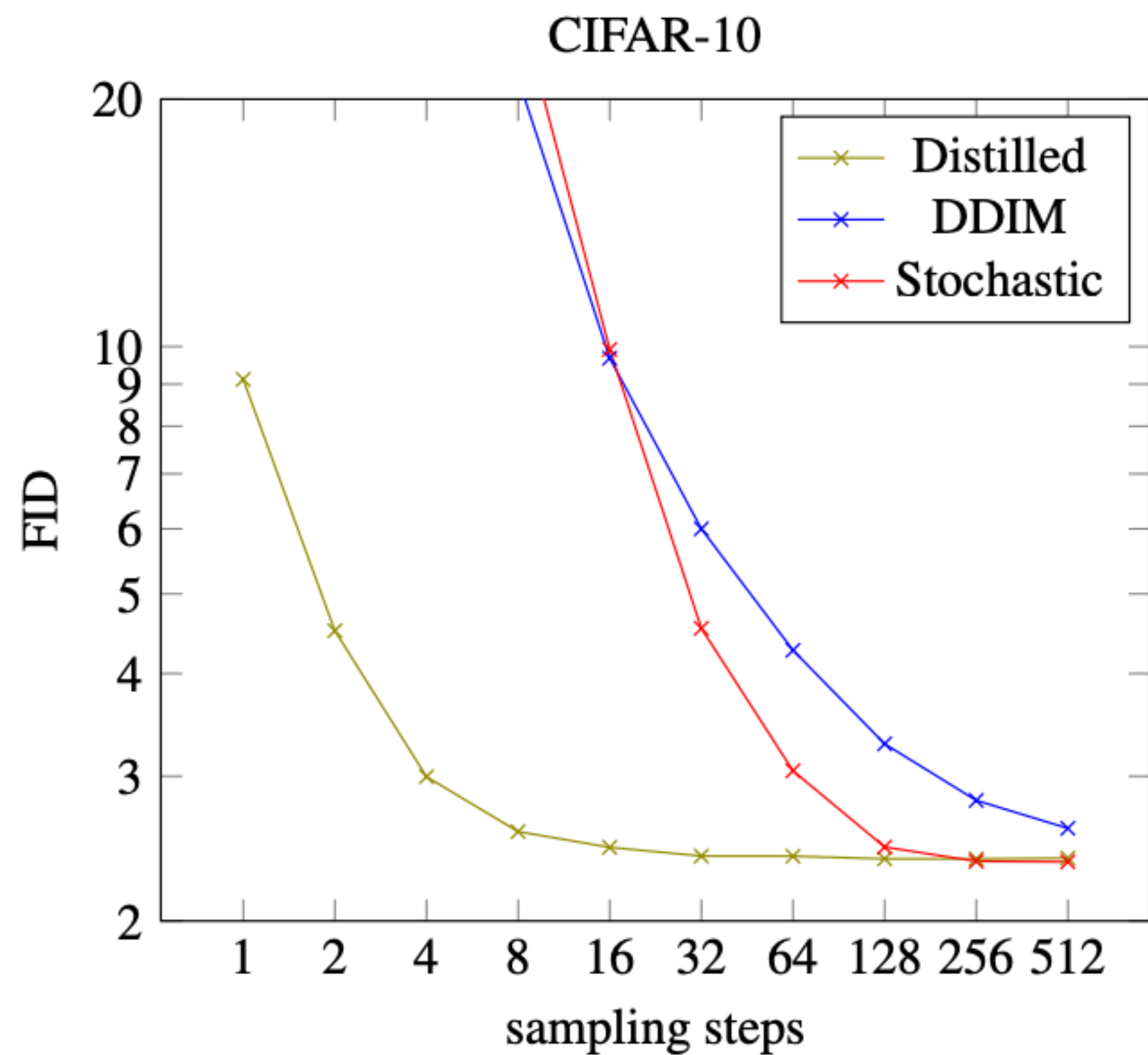
Sampling is slow!

Generating a 3s HD video could take 15 minutes!

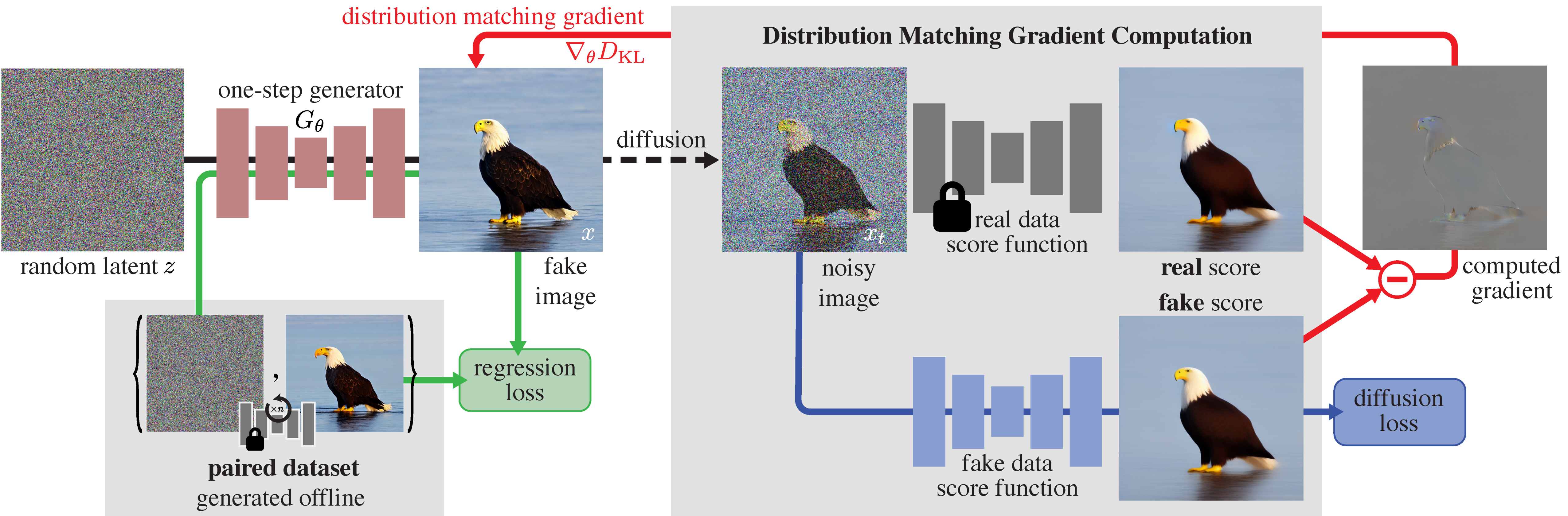
Progressive Distillation



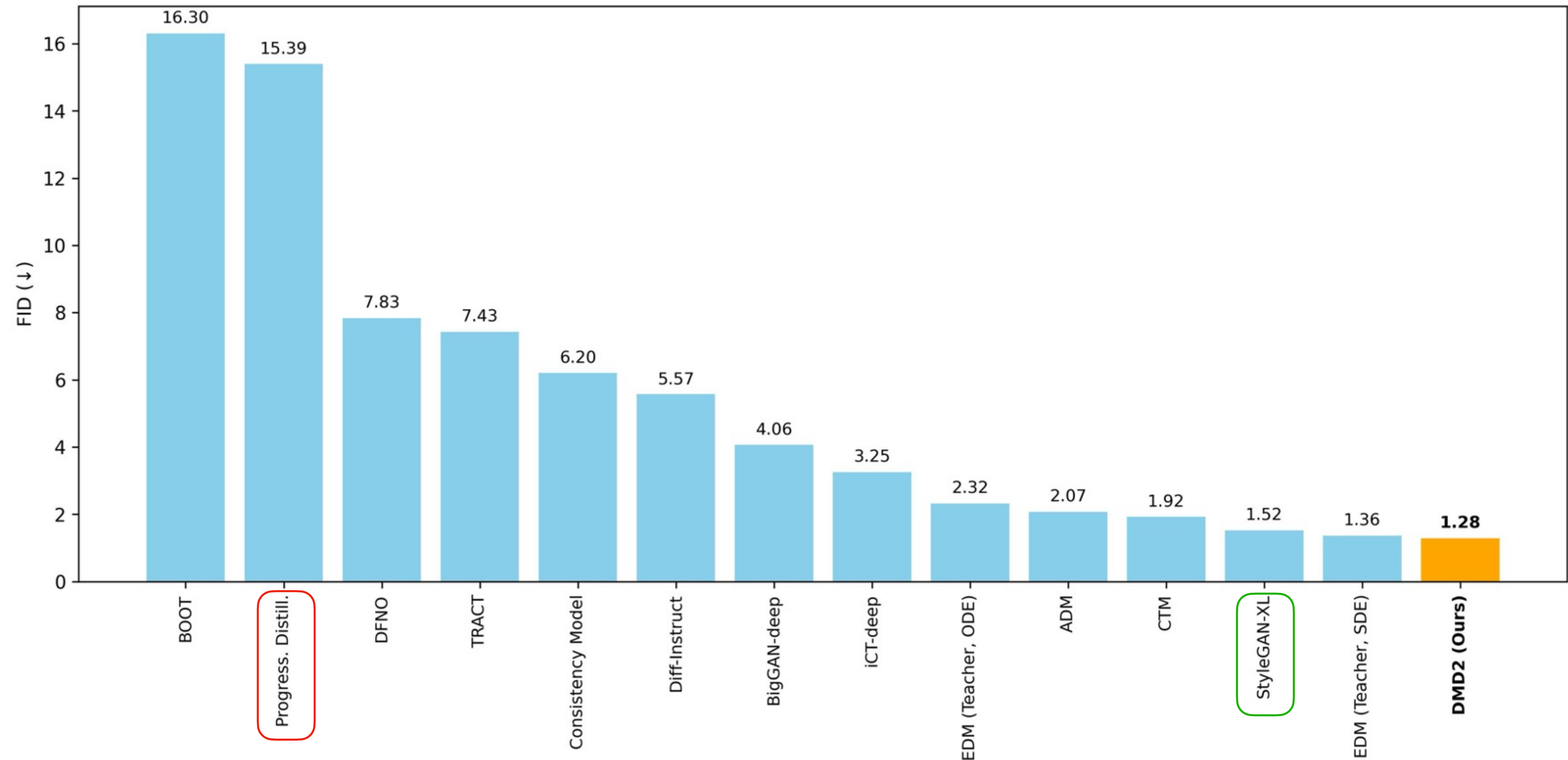
Progressive Distillation



Distribution Matching Distillation



Distribution Matching Distillation





“beret of raspberries”

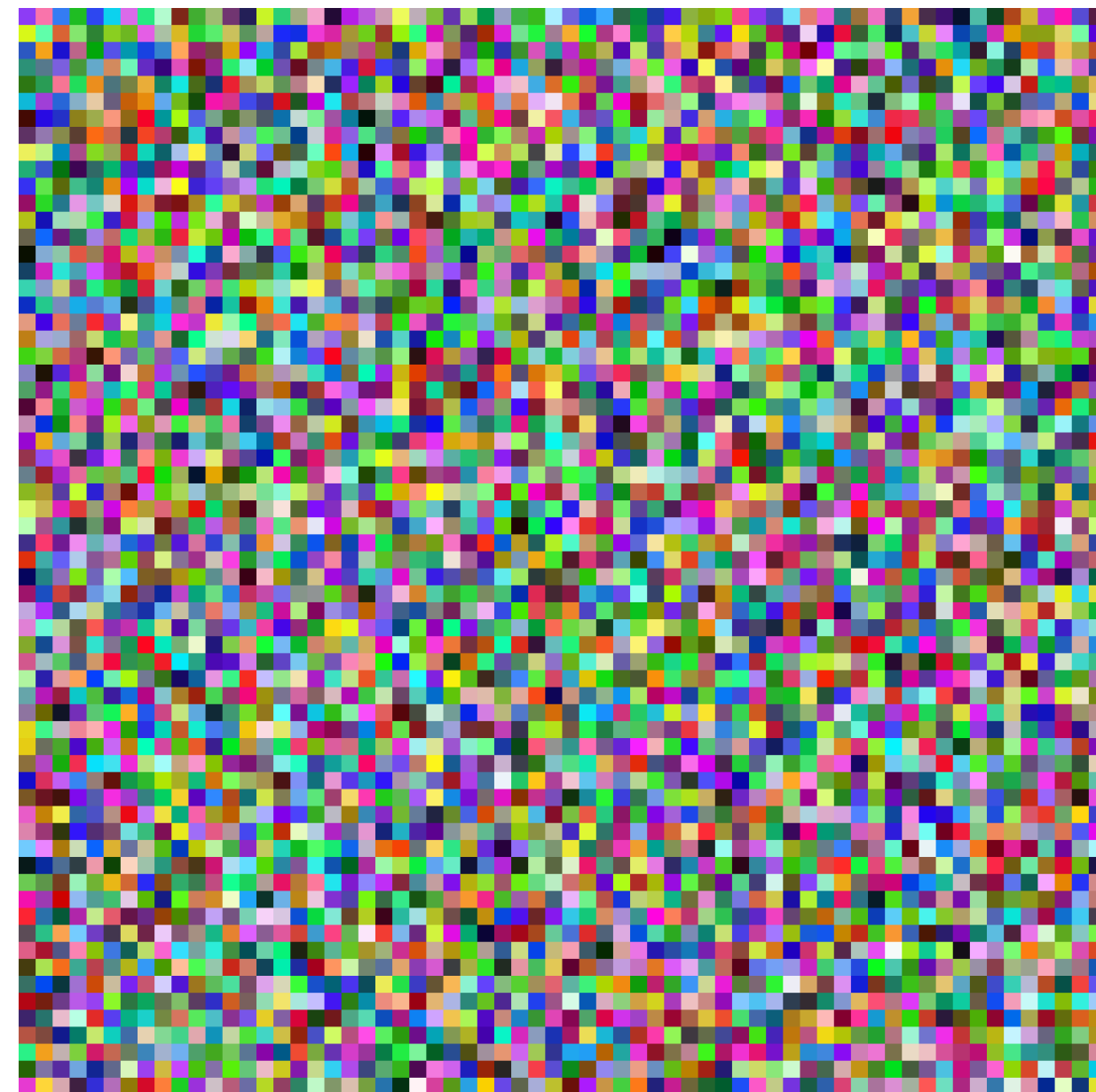
Conditioning & Guidance



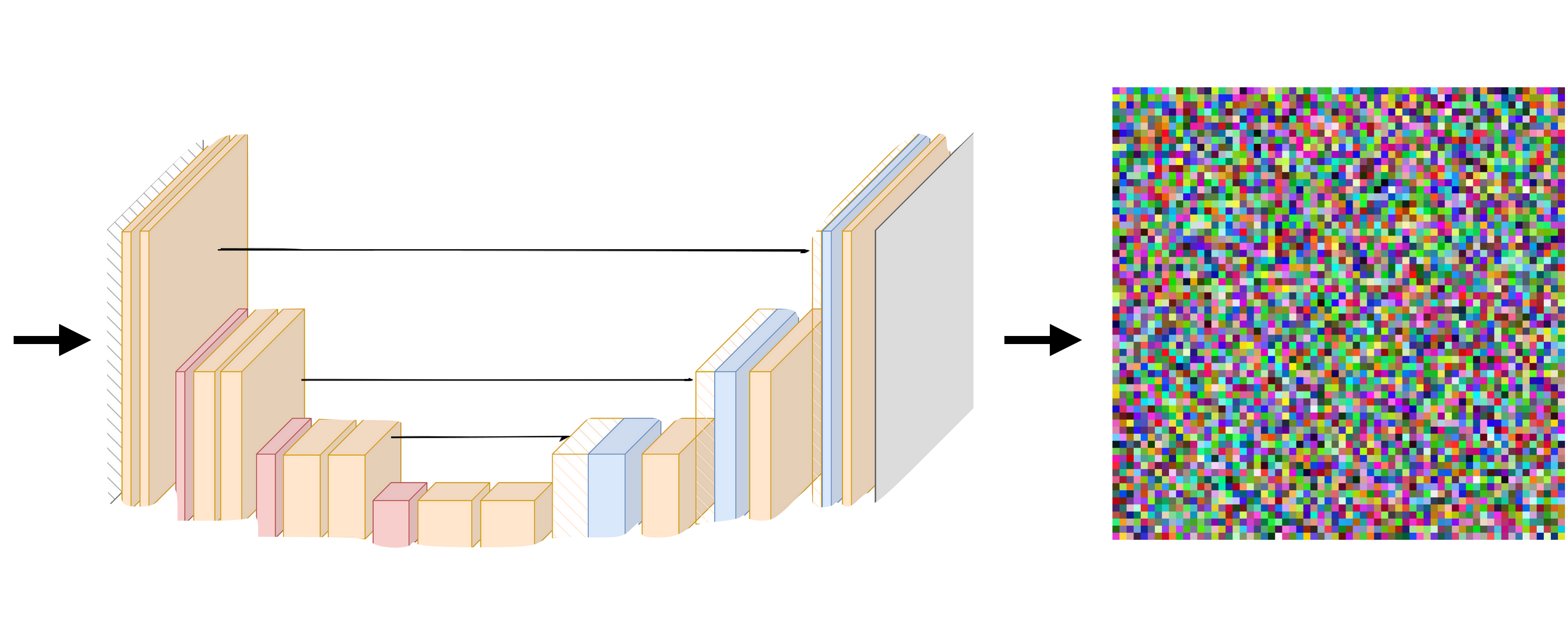
raspberry beret

Diffusion Guidance

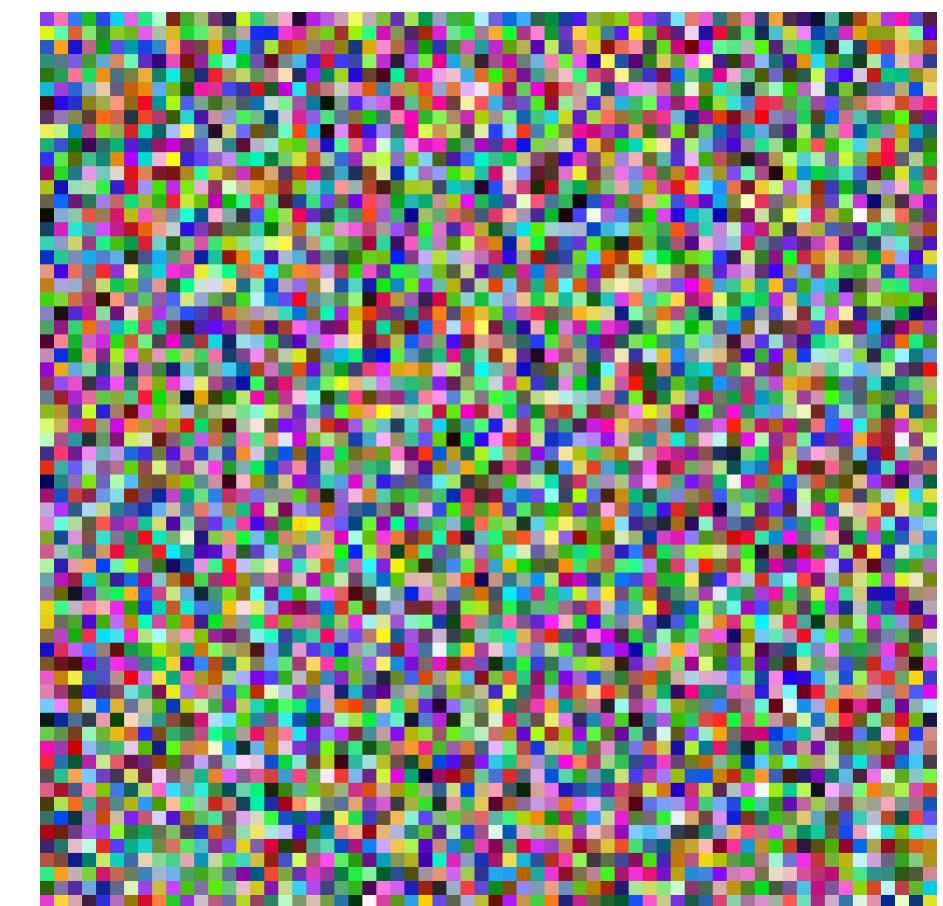
Motivation



\mathbf{X}_t Noise Image



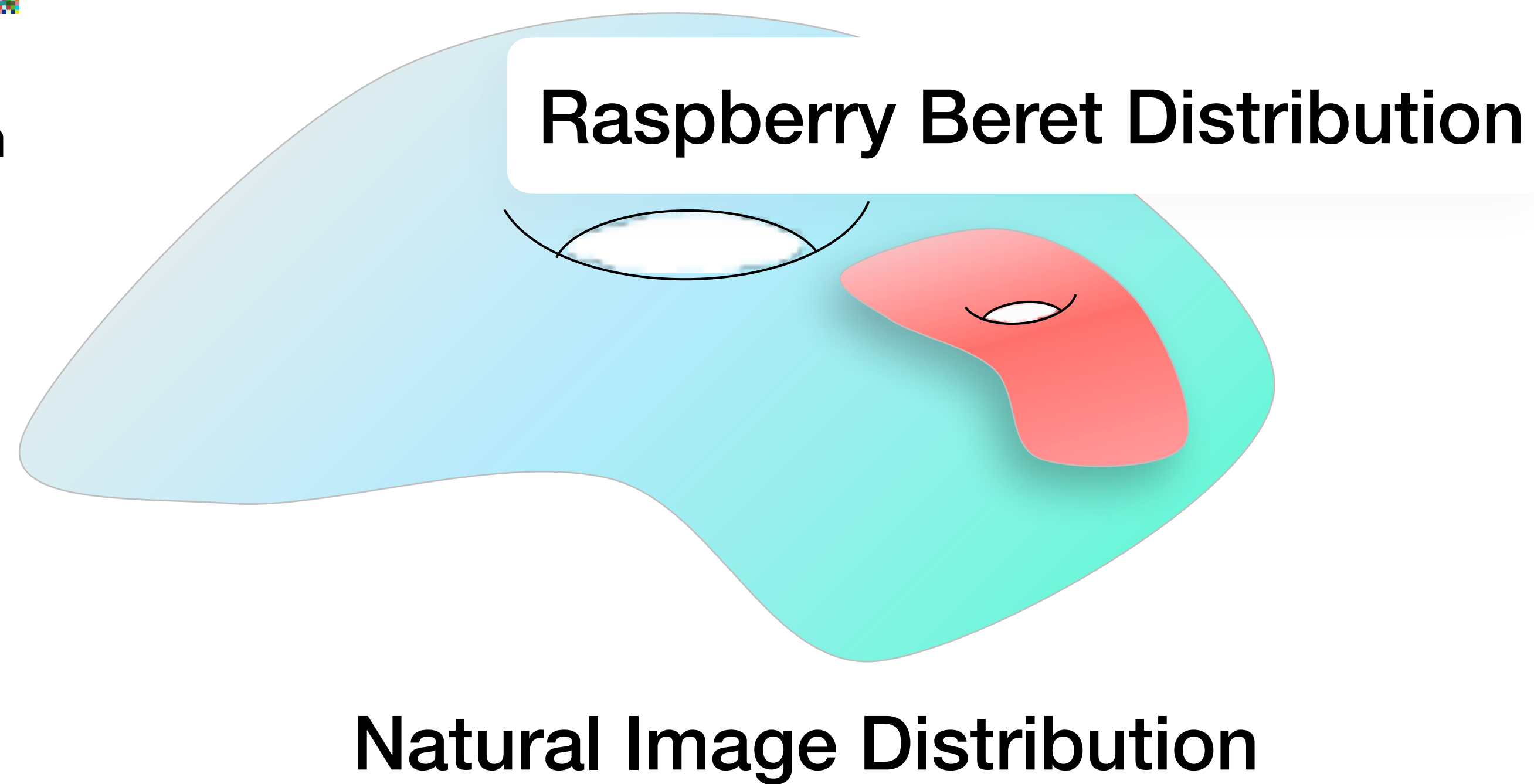
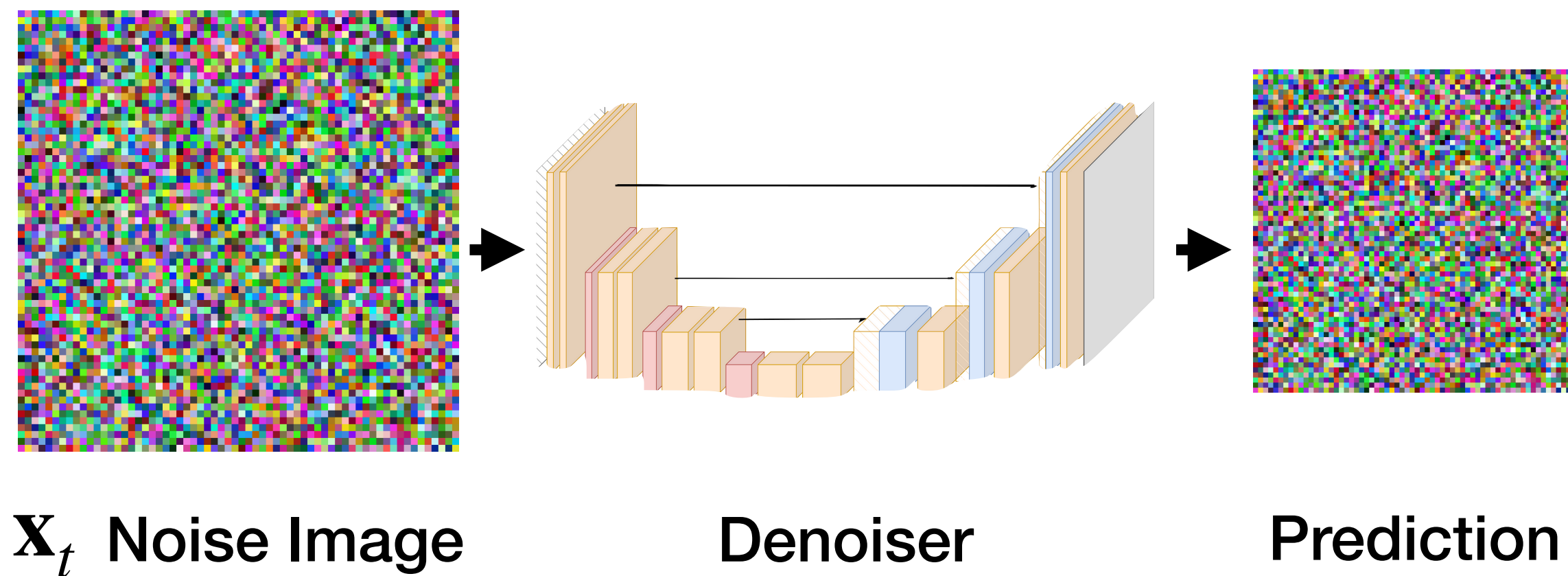
Denoiser



Prediction
(Flow, Epsilon etc.)

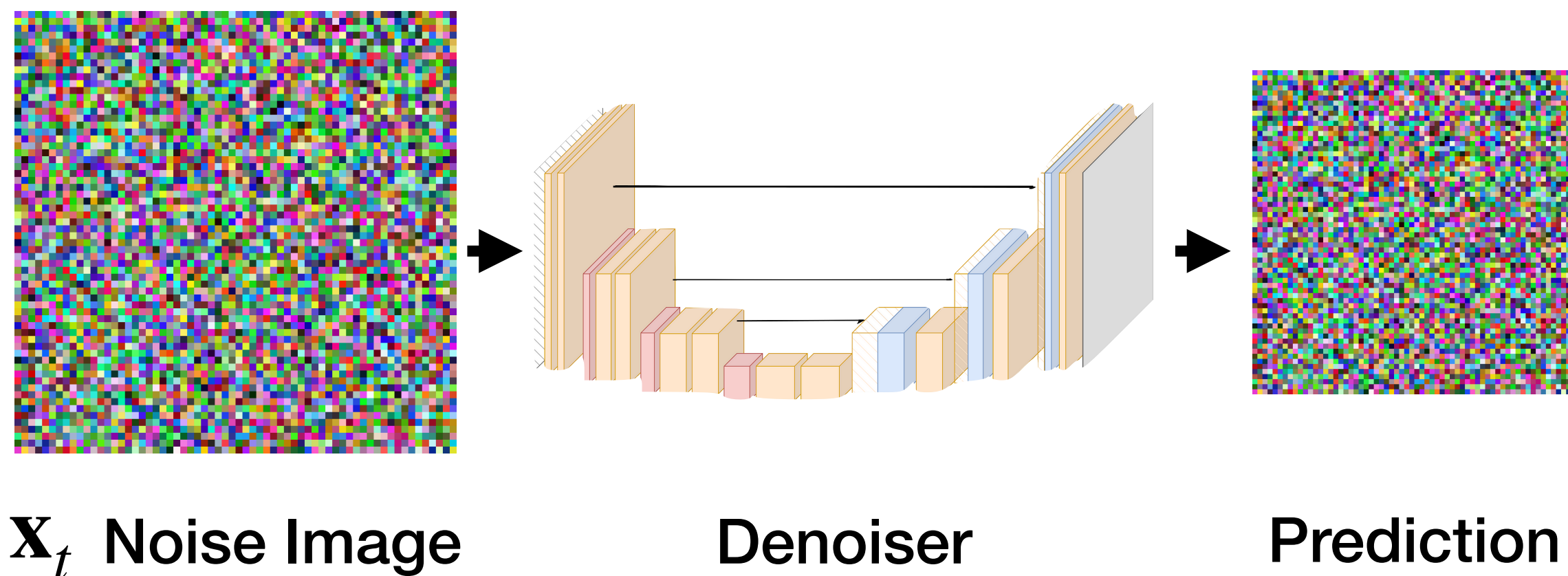
Diffusion Guidance

Motivation

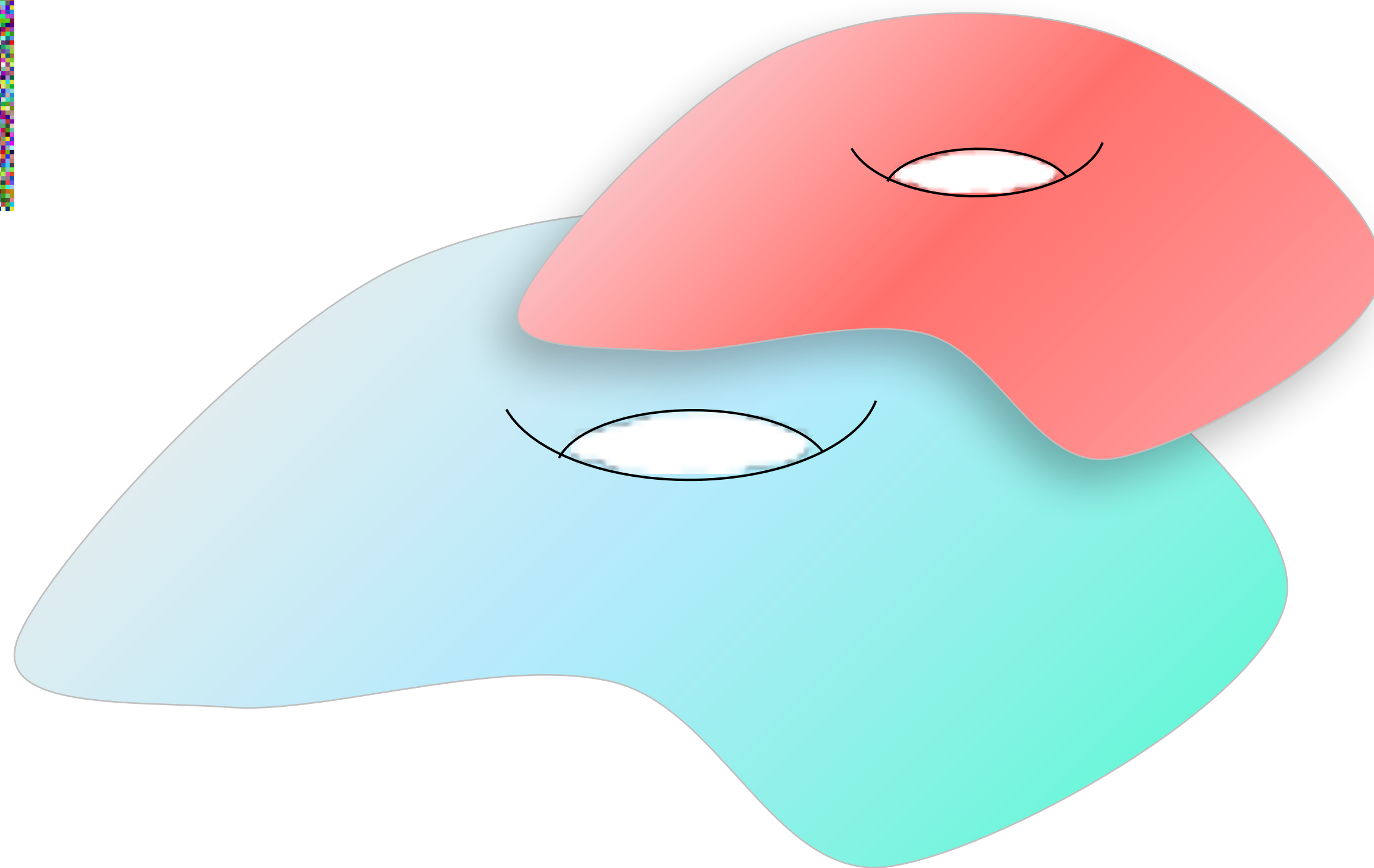


Diffusion Guidance

Motivation



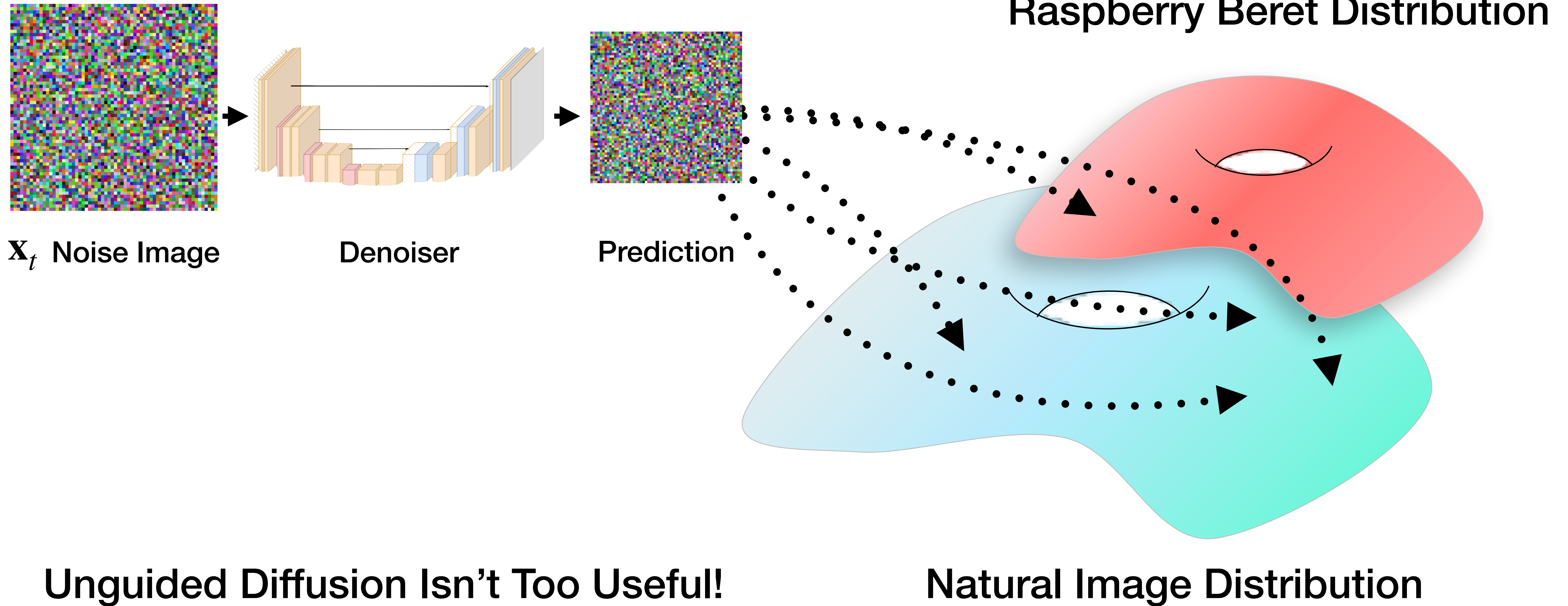
Raspberry Beret Distribution



Natural Image Distribution

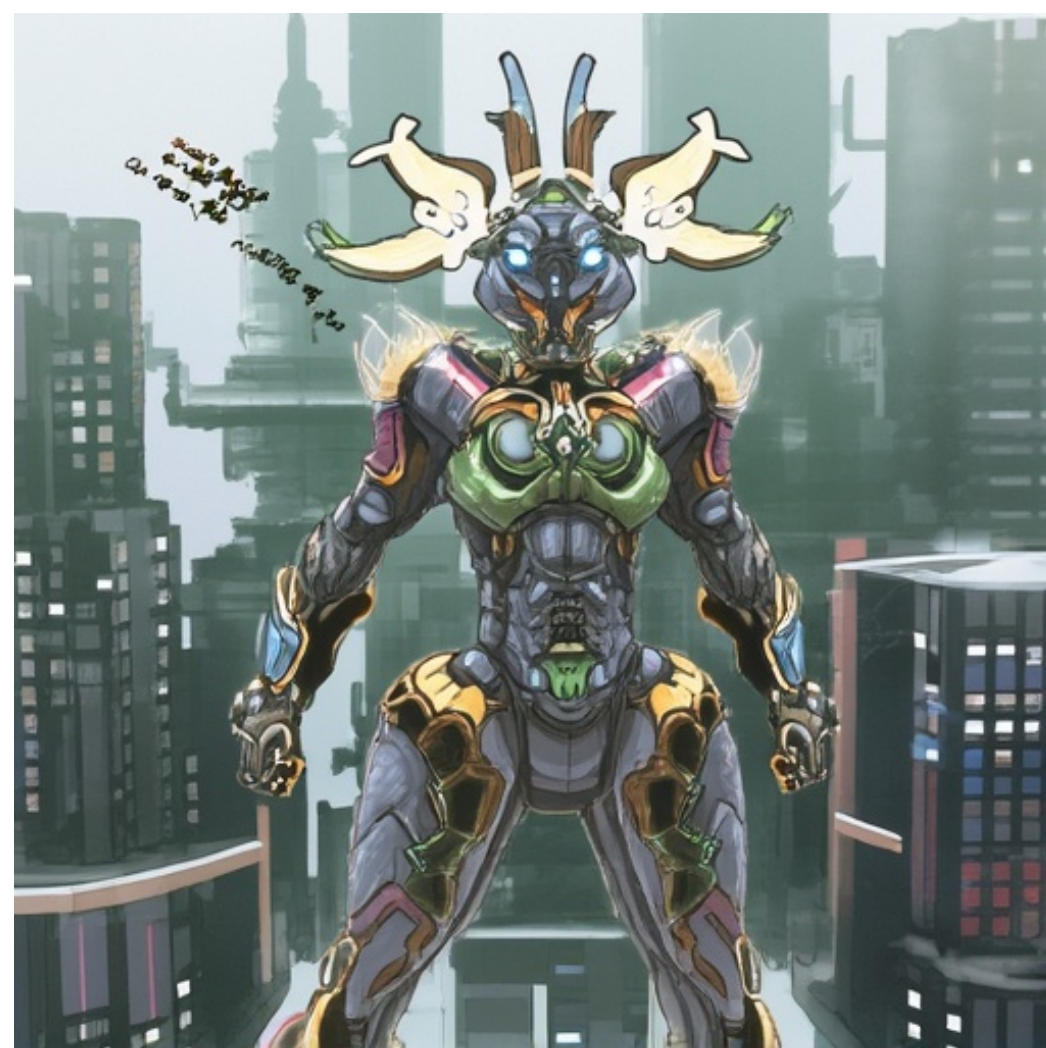
Diffusion Guidance

Motivation



Flux Pro Unguided Samples

Imitates Distribution of Internet Training Data



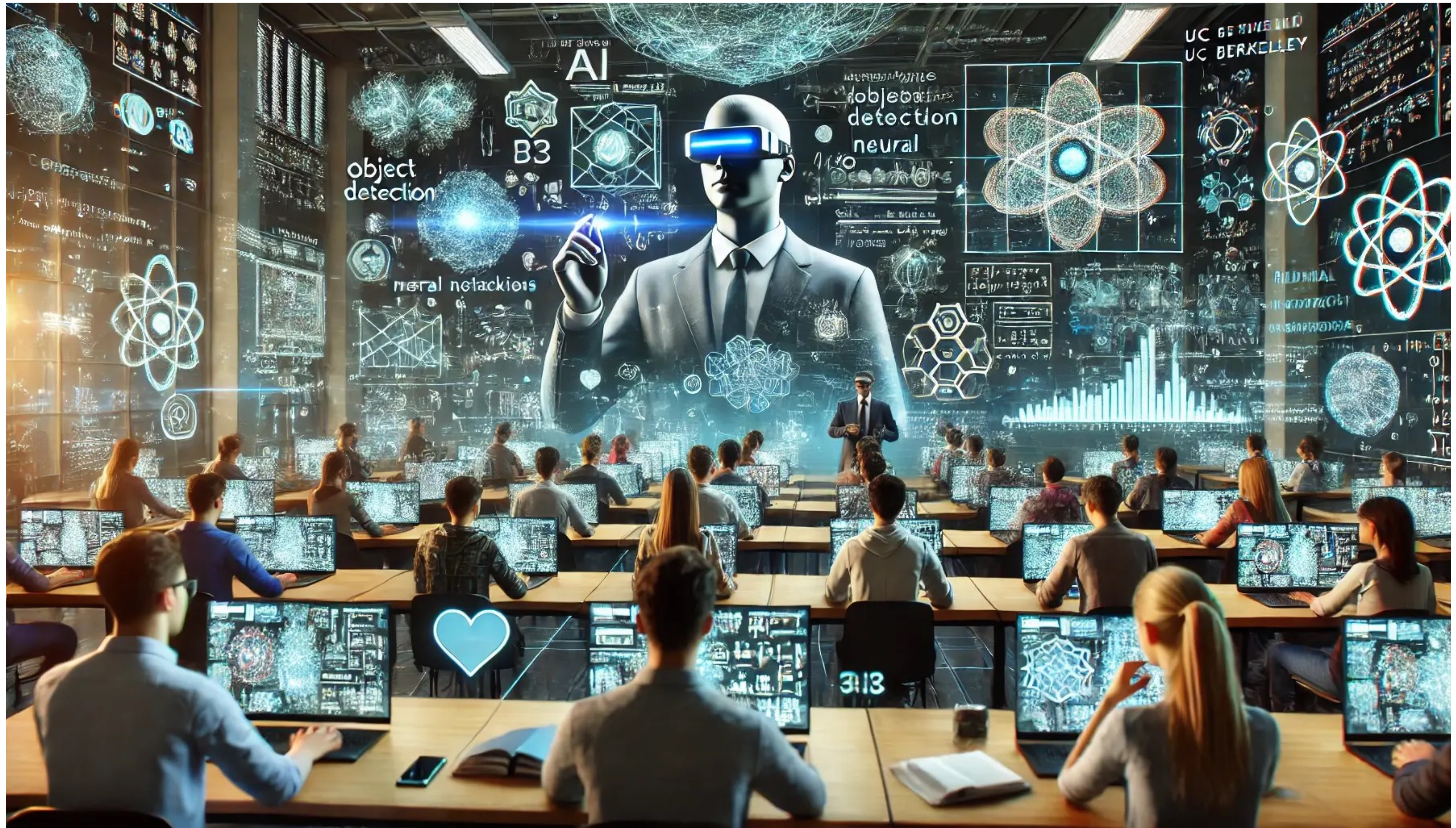
Flux Pro Guided Samples

“beret of raspberries”

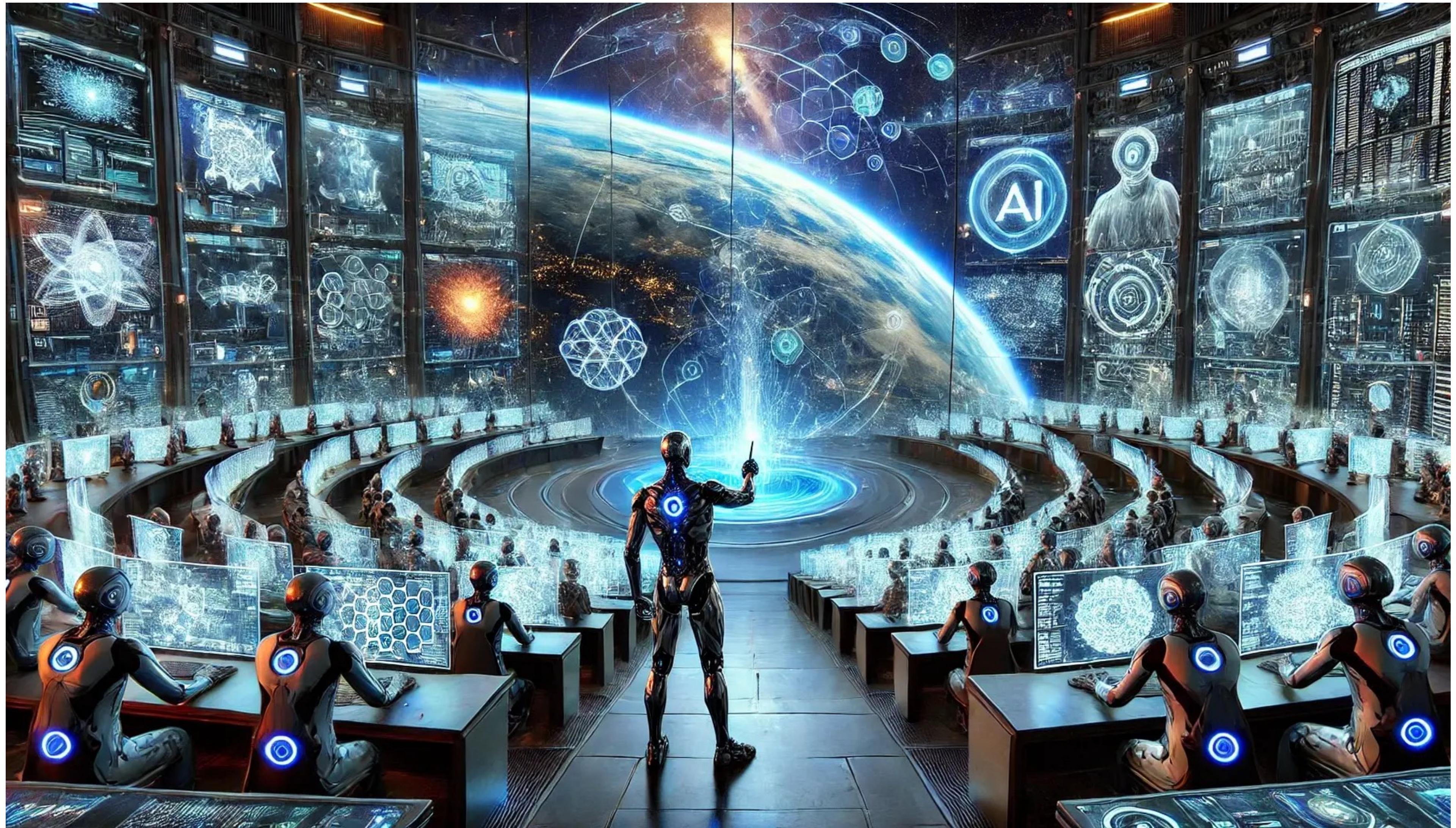




Generate a photo of a Berkeley computer vision class



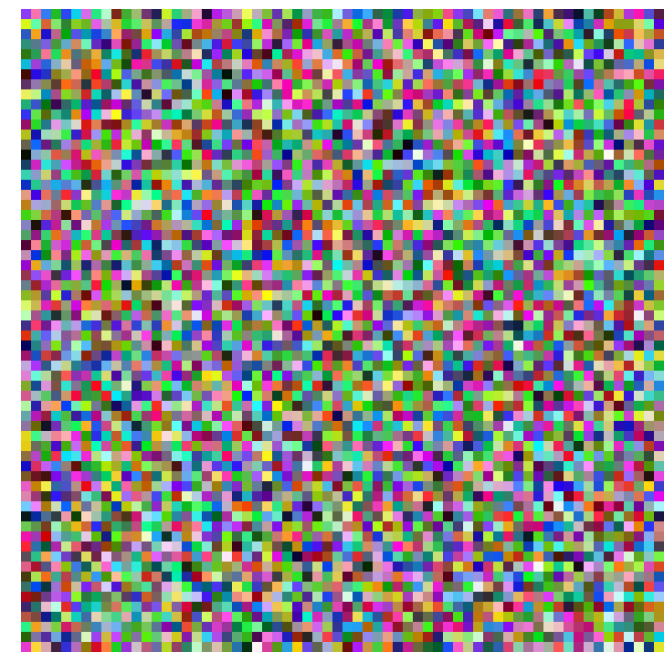
Make it more epic



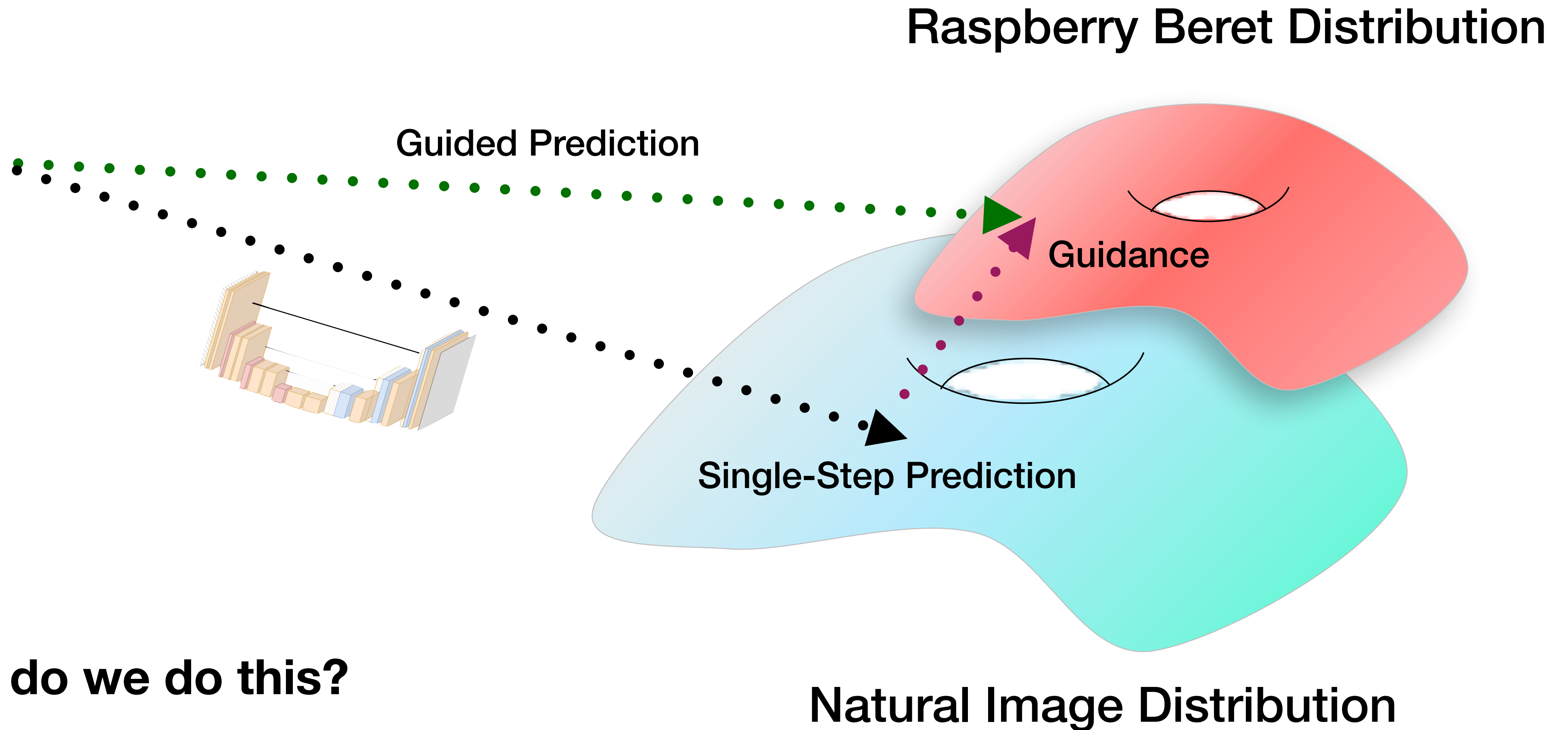
Make it THE MOST EPIC computer vision class that you can ever think of

Diffusion Guidance

Push Toward a Conditional Mode



\mathbf{X}_t Noise Image



Two Approaches

Original

Classifier Guidance

Guide with a
pretrained classifier.

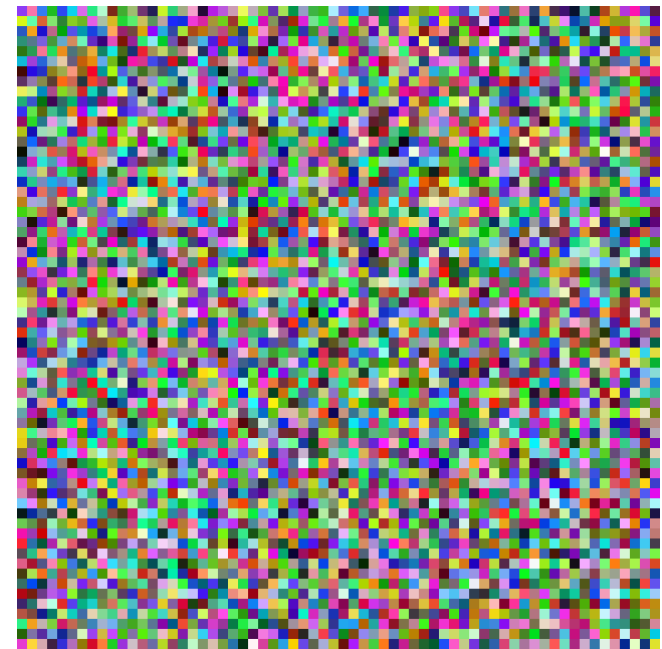
Current

Classifier-Free Guidance

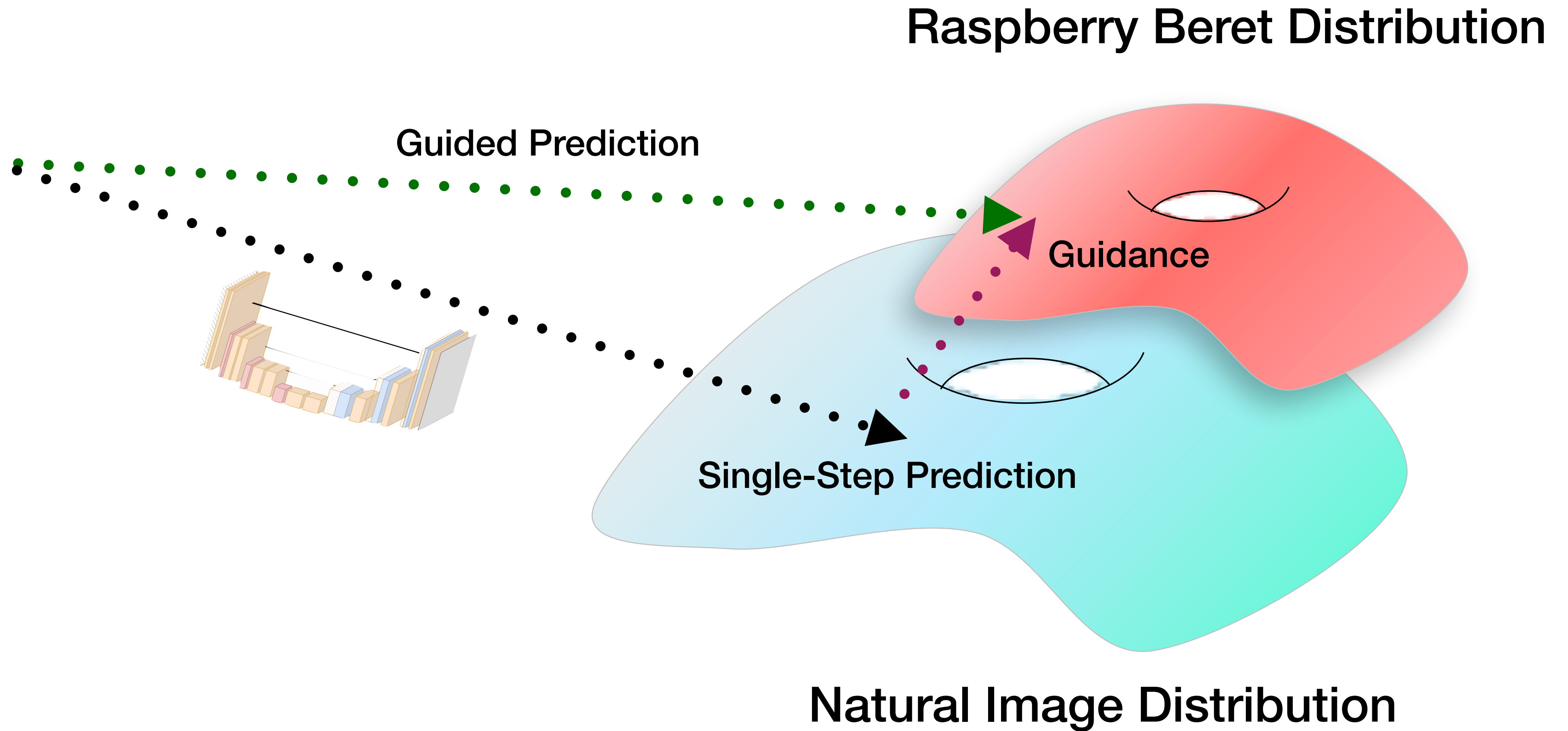
Guide a diffusion
model with itself.

Diffusion Guidance

Push Toward a Conditional Mode

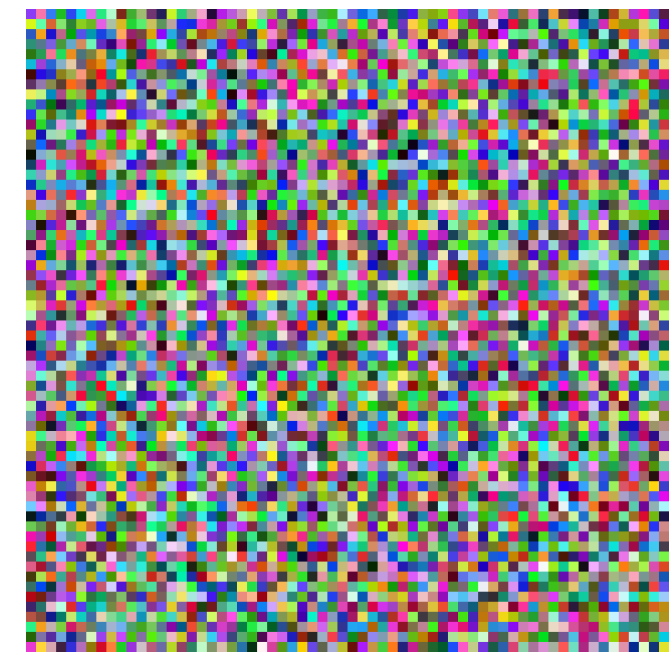


\mathbf{X}_t Noise Image

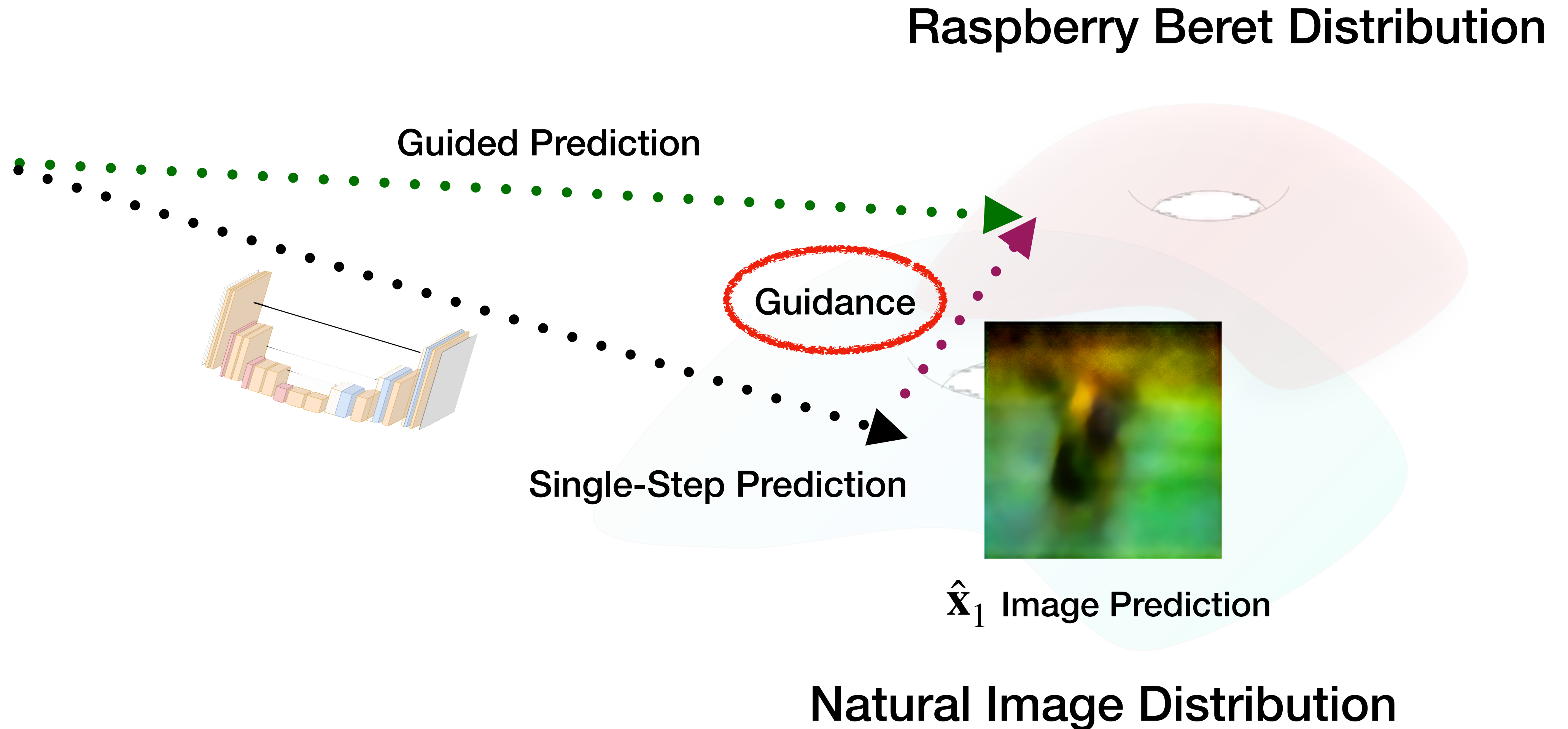


Classifier Guidance

Using a Pretrained Classifier

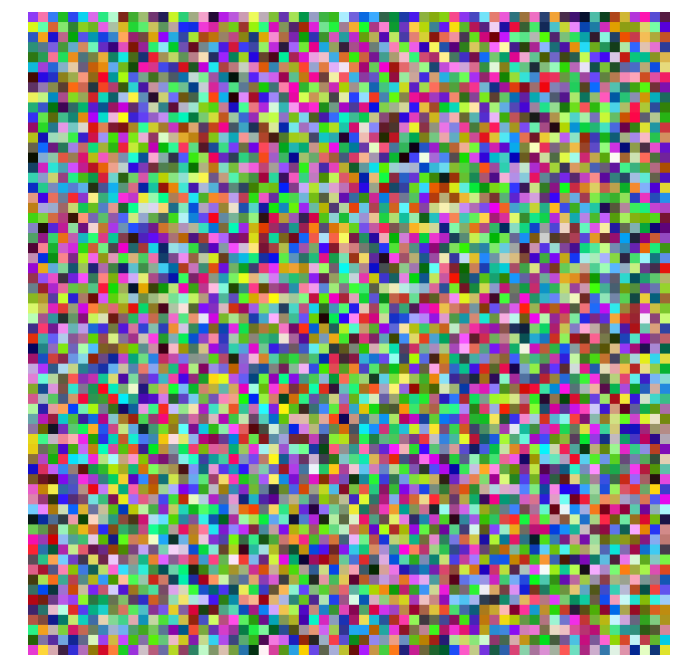


\mathbf{X}_t Noise Image



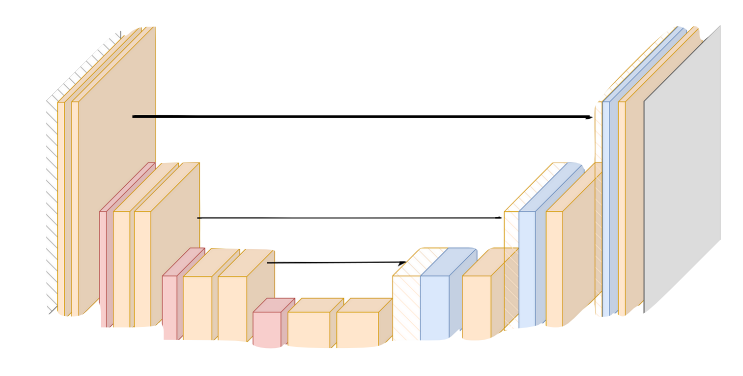
Classifier Guidance

Using a Pretrained Classifier



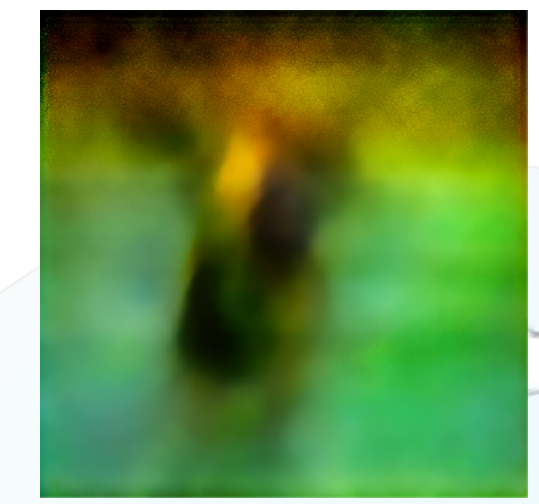
\mathbf{X}_t Noise Image

...

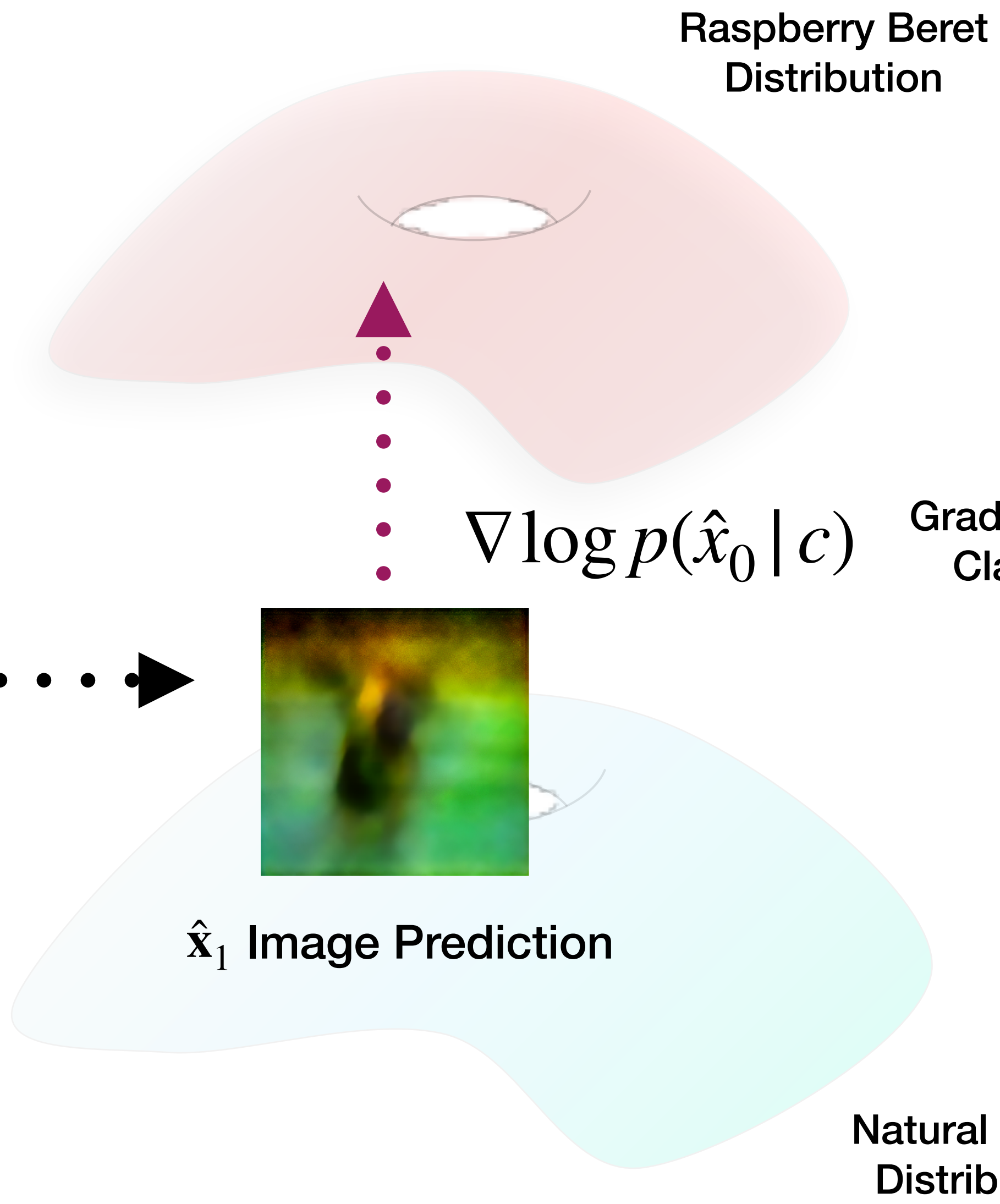


Diffusion Model

...



$\hat{\mathbf{x}}_1$ Image Prediction



Raspberry Beret Distribution

$$\nabla \log p(\hat{x}_0 | c)$$

Gradient from Classifier

Natural Image Distribution

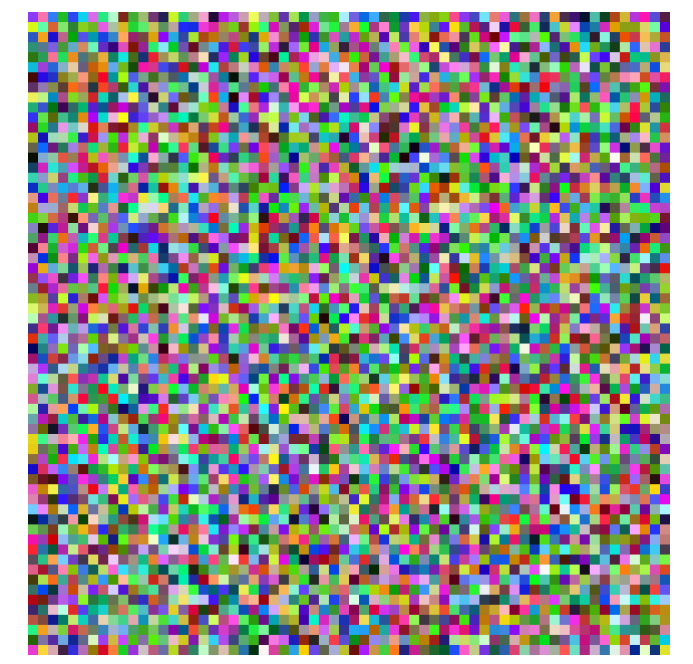
For Example...

			
mite	container ship	motor scooter	leopard
<ul style="list-style-type: none"> mite black widow cockroach tick starfish 	<ul style="list-style-type: none"> container ship lifeboat amphibian fireboat drilling platform 	<ul style="list-style-type: none"> motor scooter go-kart moped bumper car golfcart 	<ul style="list-style-type: none"> leopard jaguar cheetah snow leopard Egyptian cat
			
grille	mushroom	cherry	Madagascar cat
<ul style="list-style-type: none"> convertible grille pickup beach wagon fire engine 	<ul style="list-style-type: none"> agaric mushroom jelly fungus gill fungus dead-man's-fingers 	<ul style="list-style-type: none"> dalmatian grape elderberry ffordshire bullterrier currant 	<ul style="list-style-type: none"> squirrel monkey spider monkey titi indri howler monkey

IMAGENET

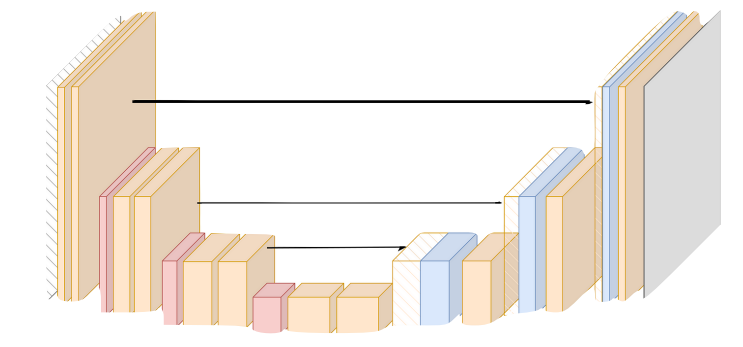
Classifier Guidance

Using a Pretrained Classifier



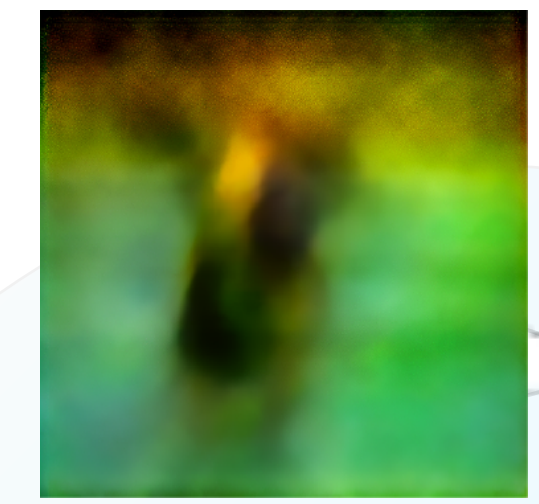
\mathbf{X}_t Noise Image

...



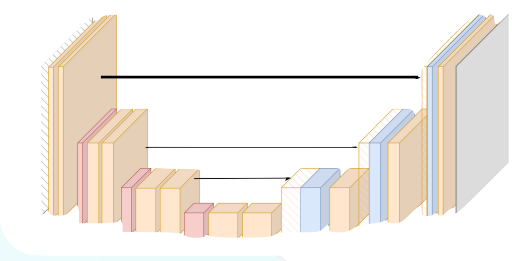
Diffusion Model

...



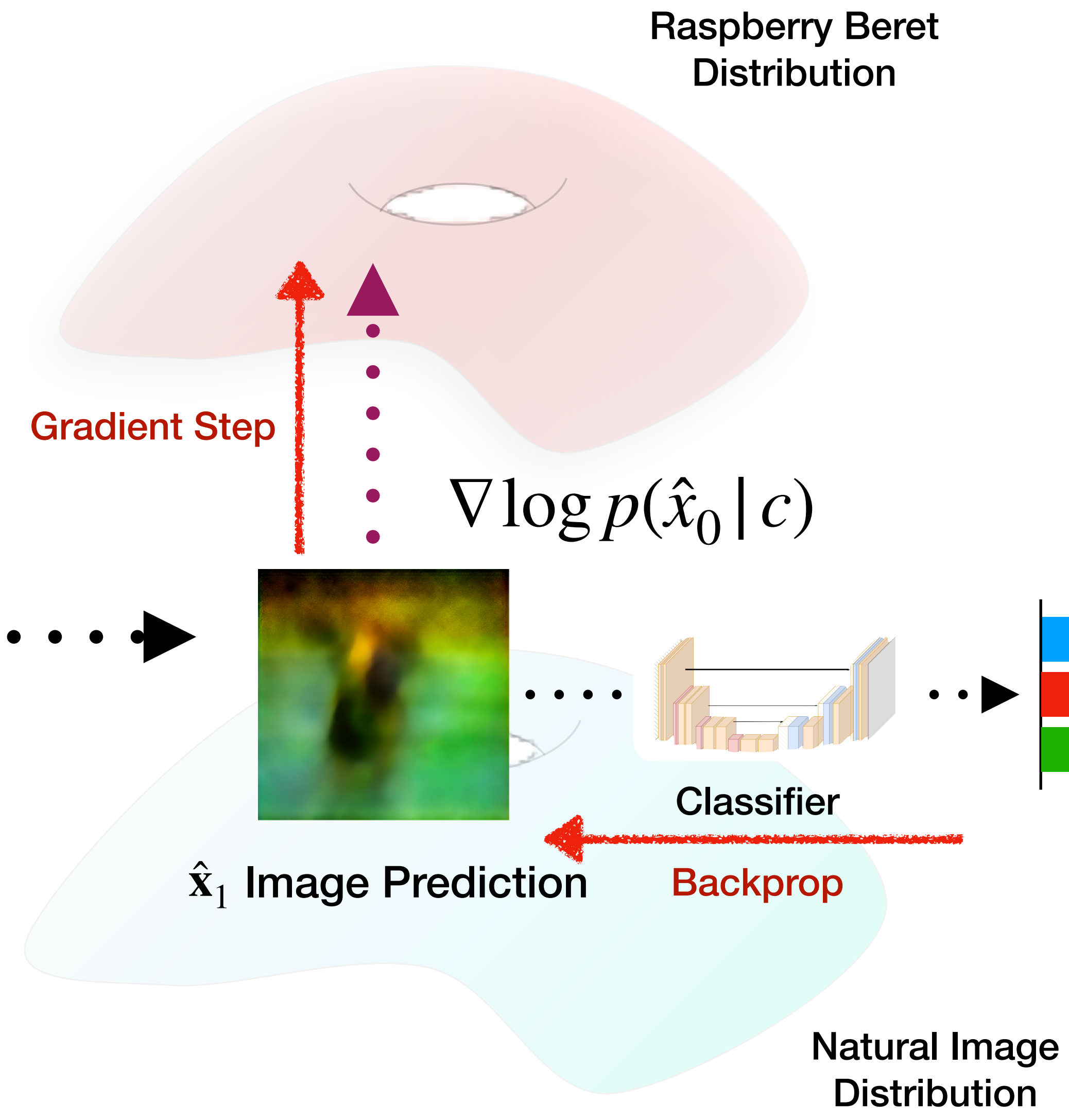
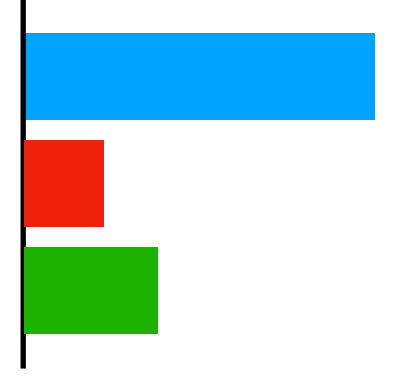
$\hat{\mathbf{x}}_1$ Image Prediction

...



Classifier

...



Raspberry Beret Distribution

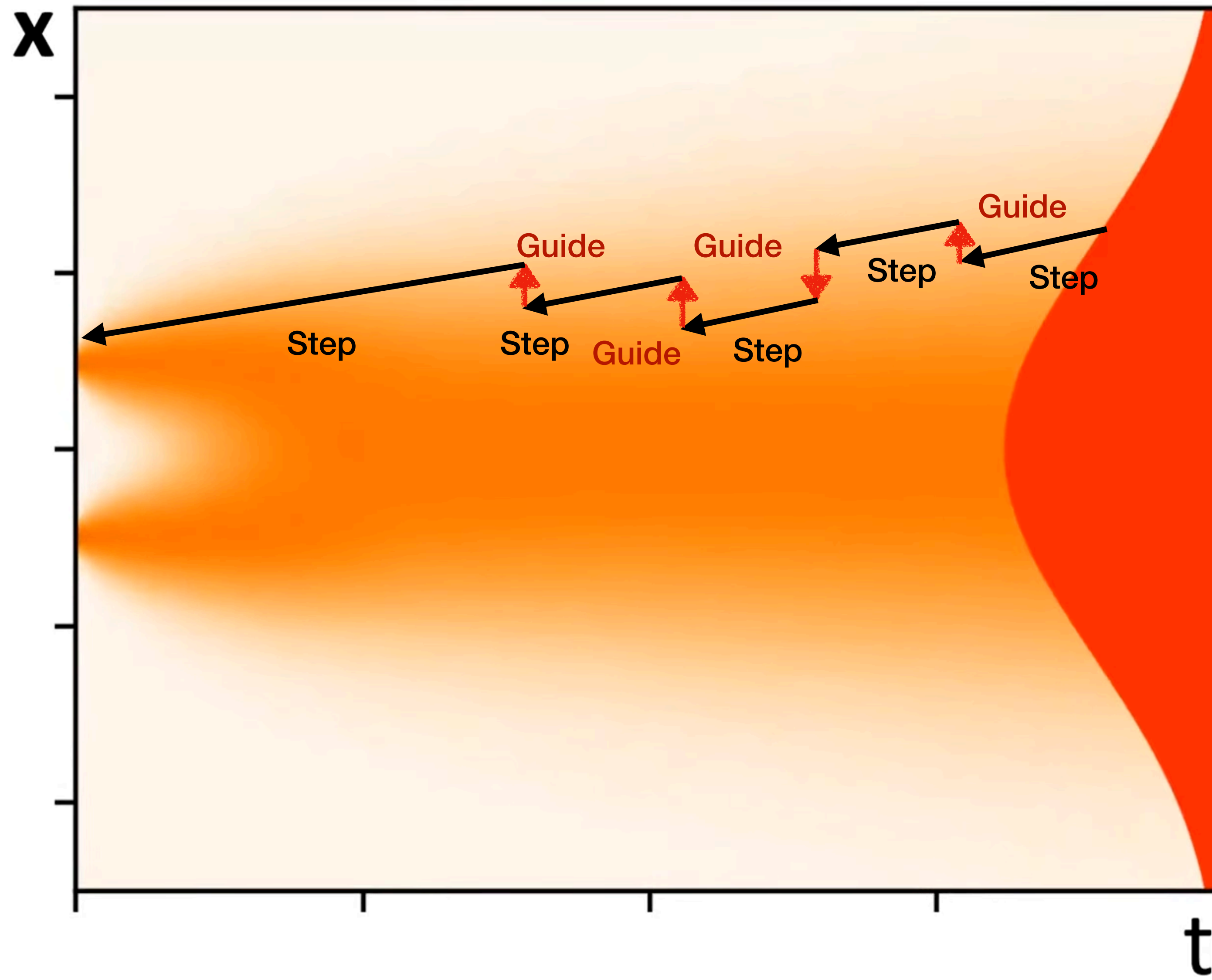
Gradient Step

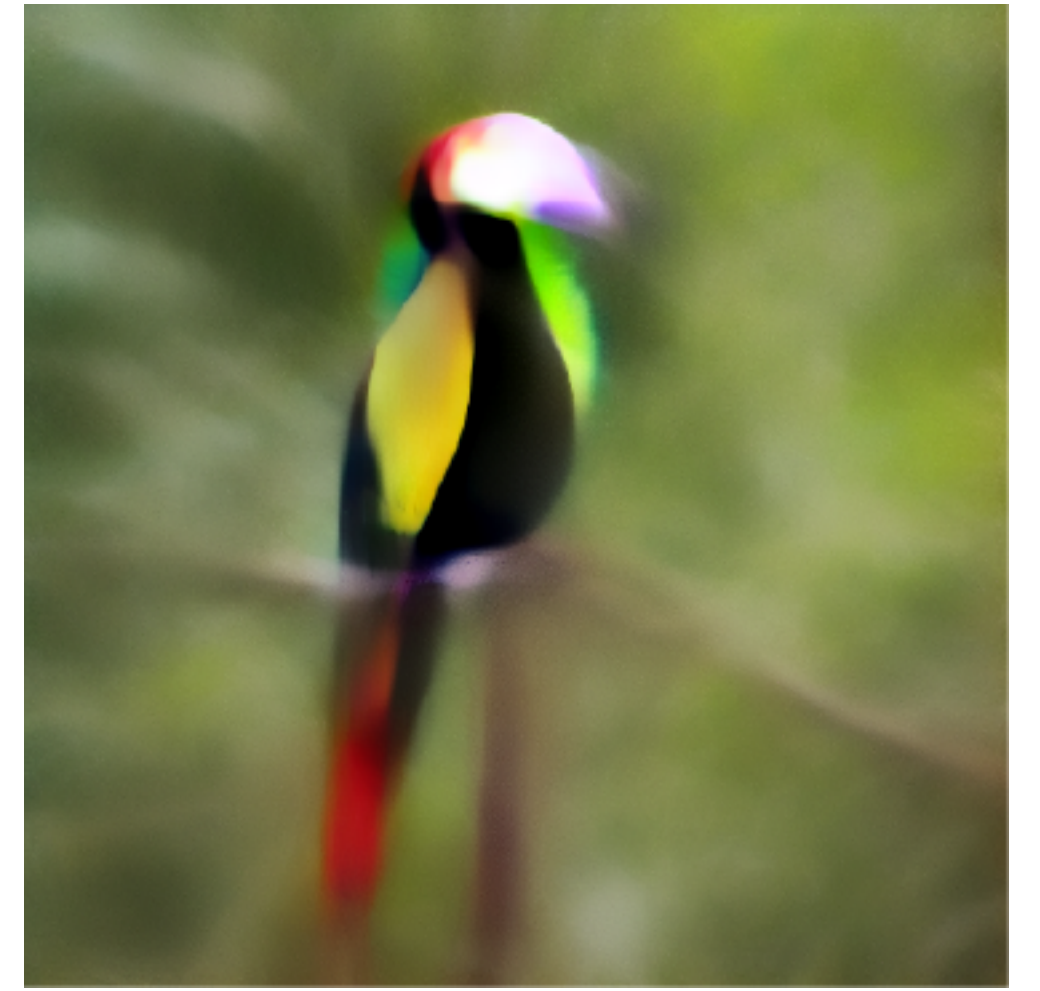
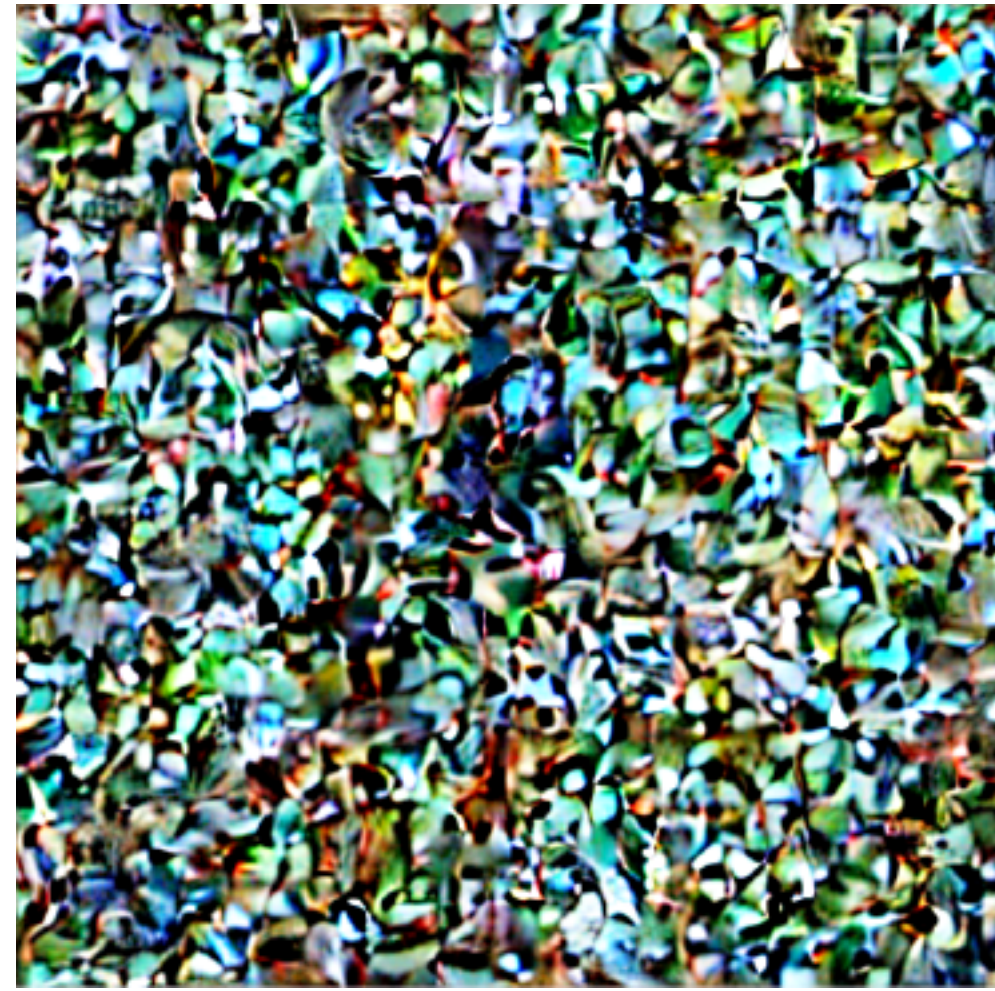
$$\nabla \log p(\hat{x}_0 | c)$$

Backprop

Natural Image Distribution

Sampling ODE View





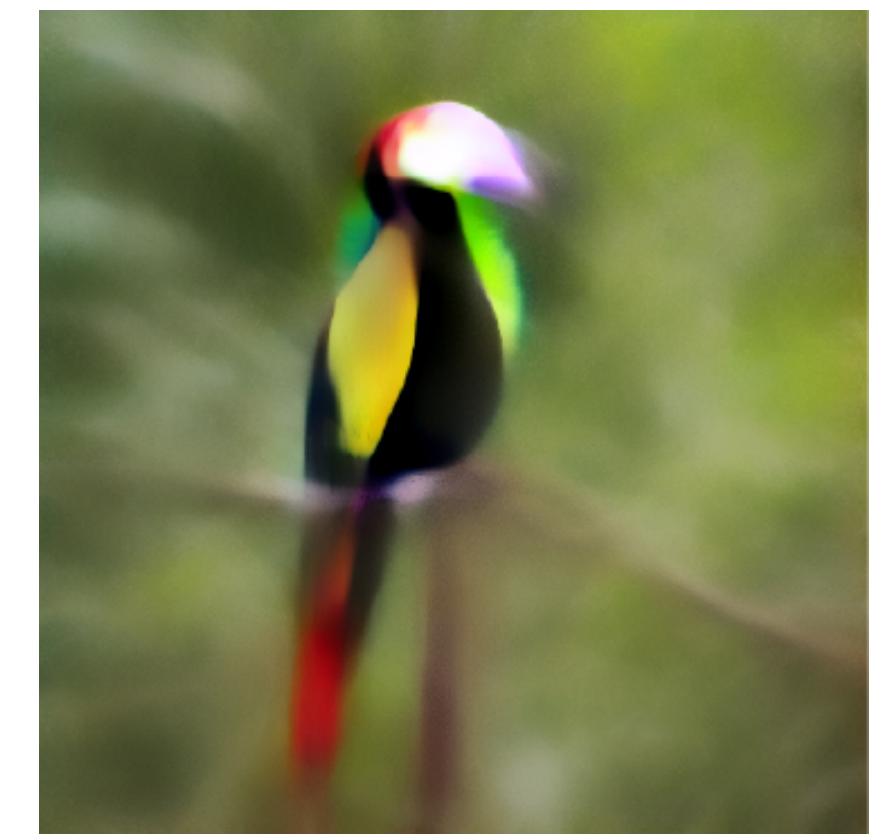
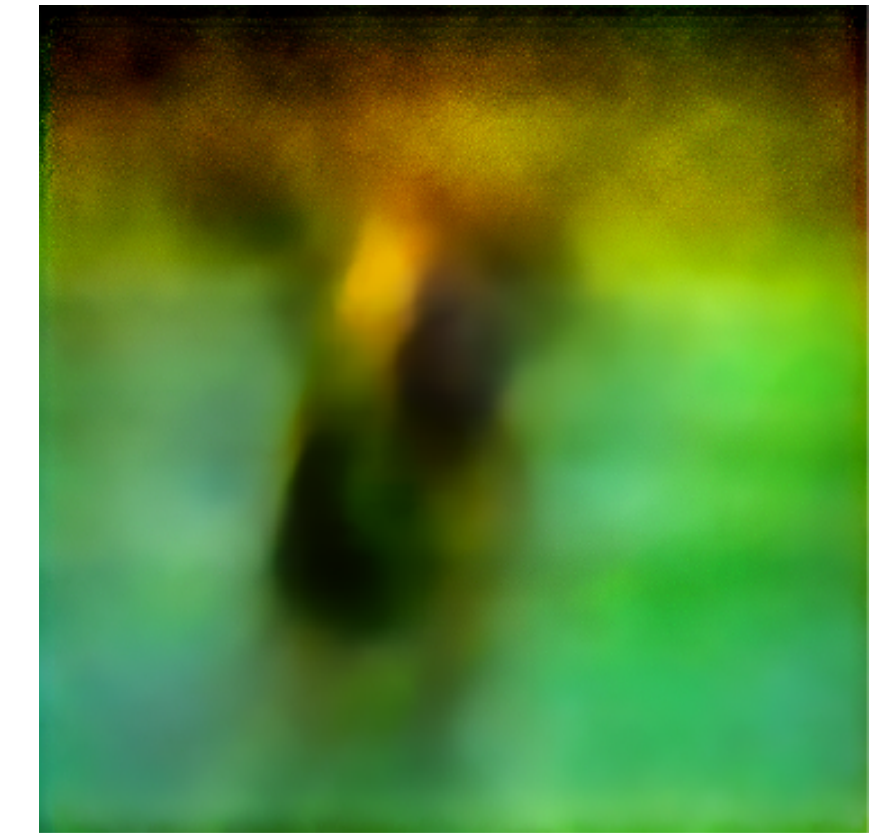
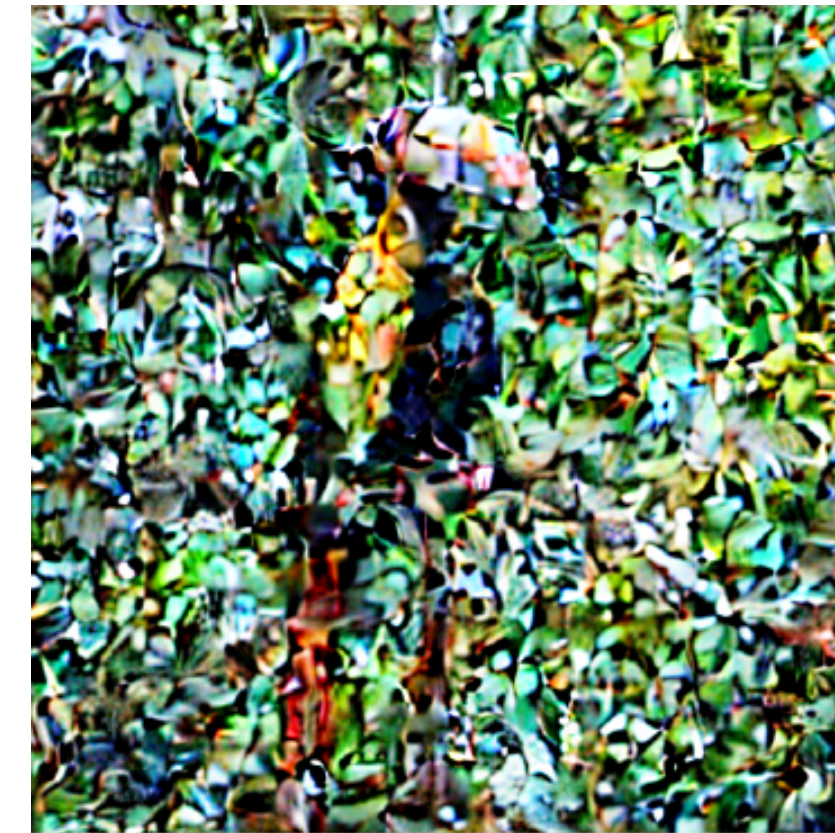
Classify These Images!

Classifier Guidance Pathologies

Intermediate Diffusion Steps are O.O.D. for Classifier



ImageNet Training Data



X_t Noise Image

\hat{x}_1 Image Prediction

Two Approaches

Original

Classifier Guidance

Guide with a
pretrained classifier.

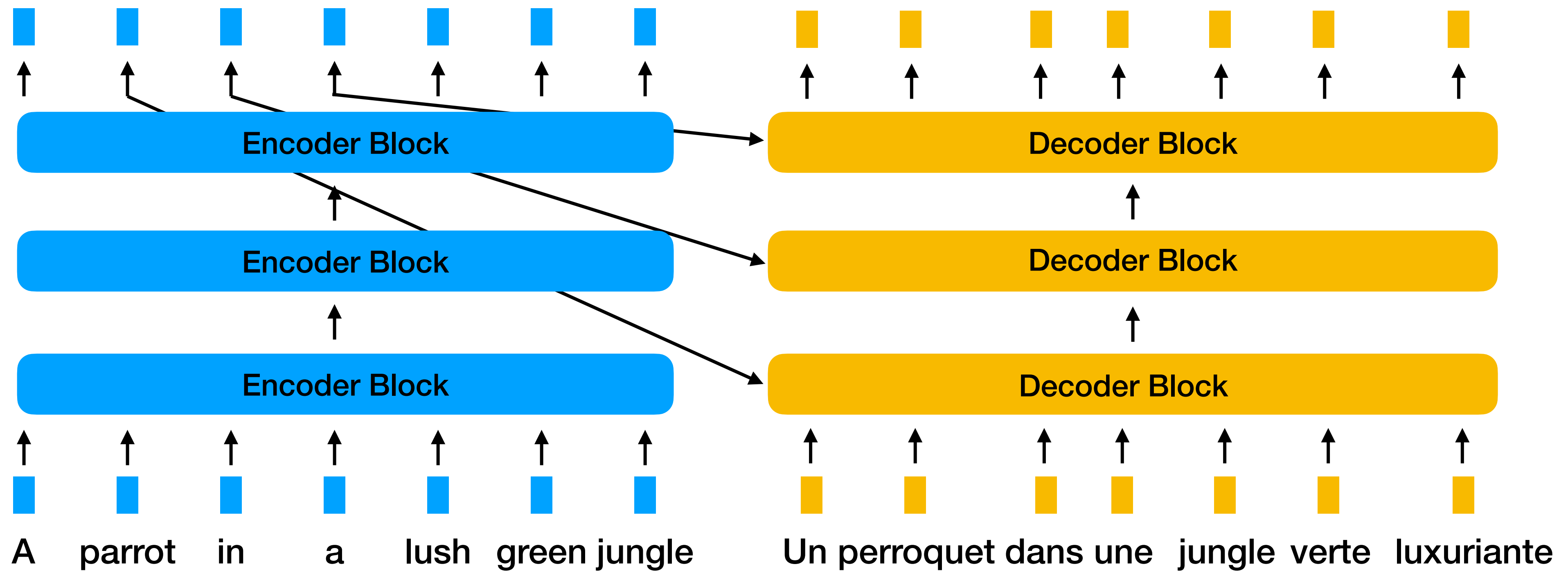
Current

Classifier-Free Guidance

Guide a diffusion
model with itself.

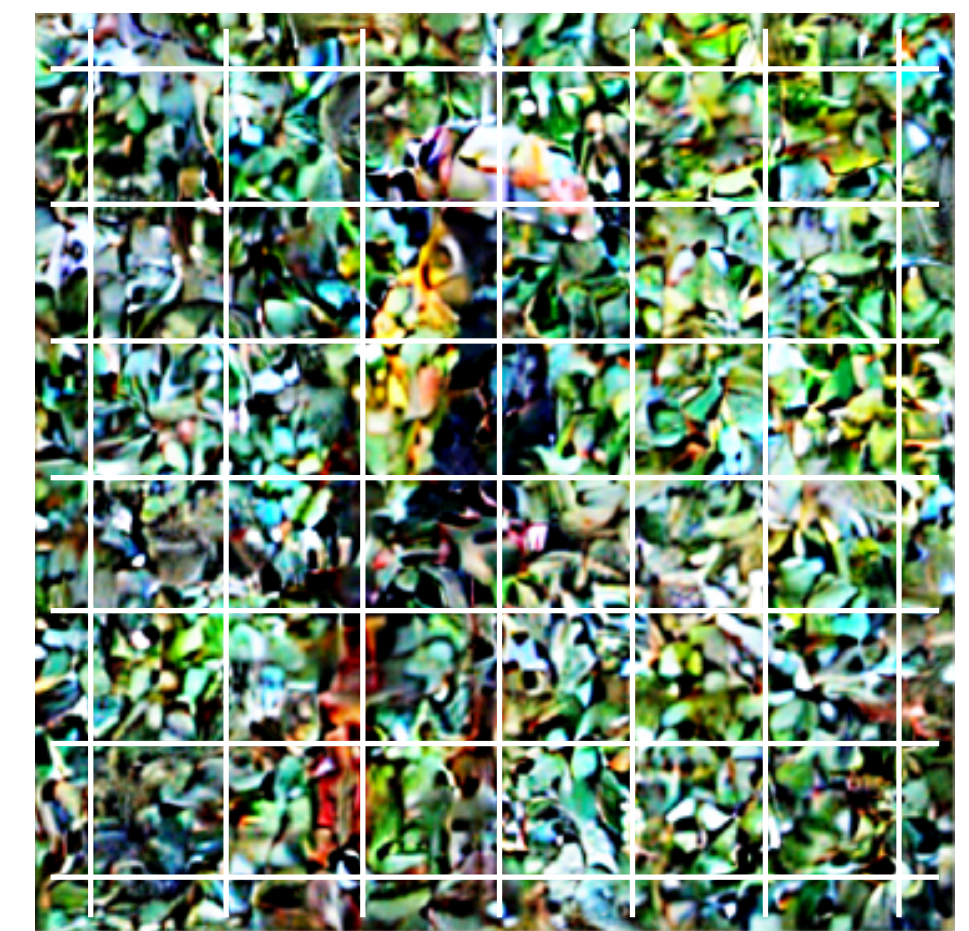
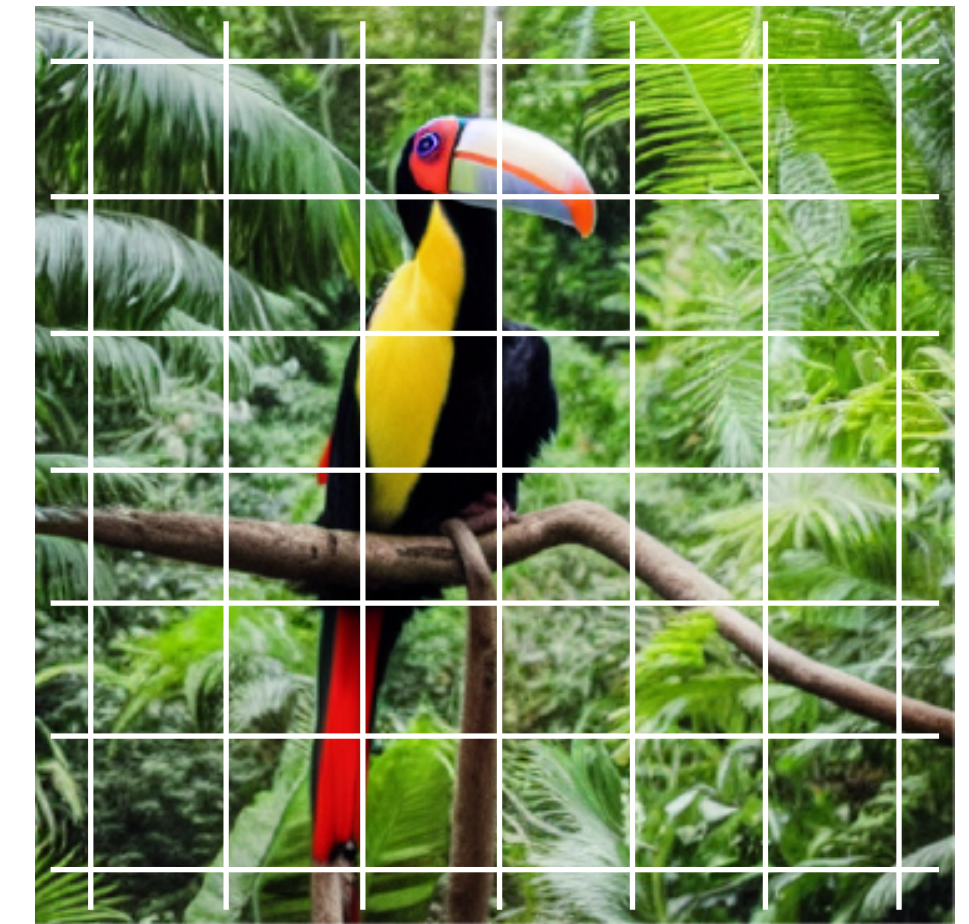
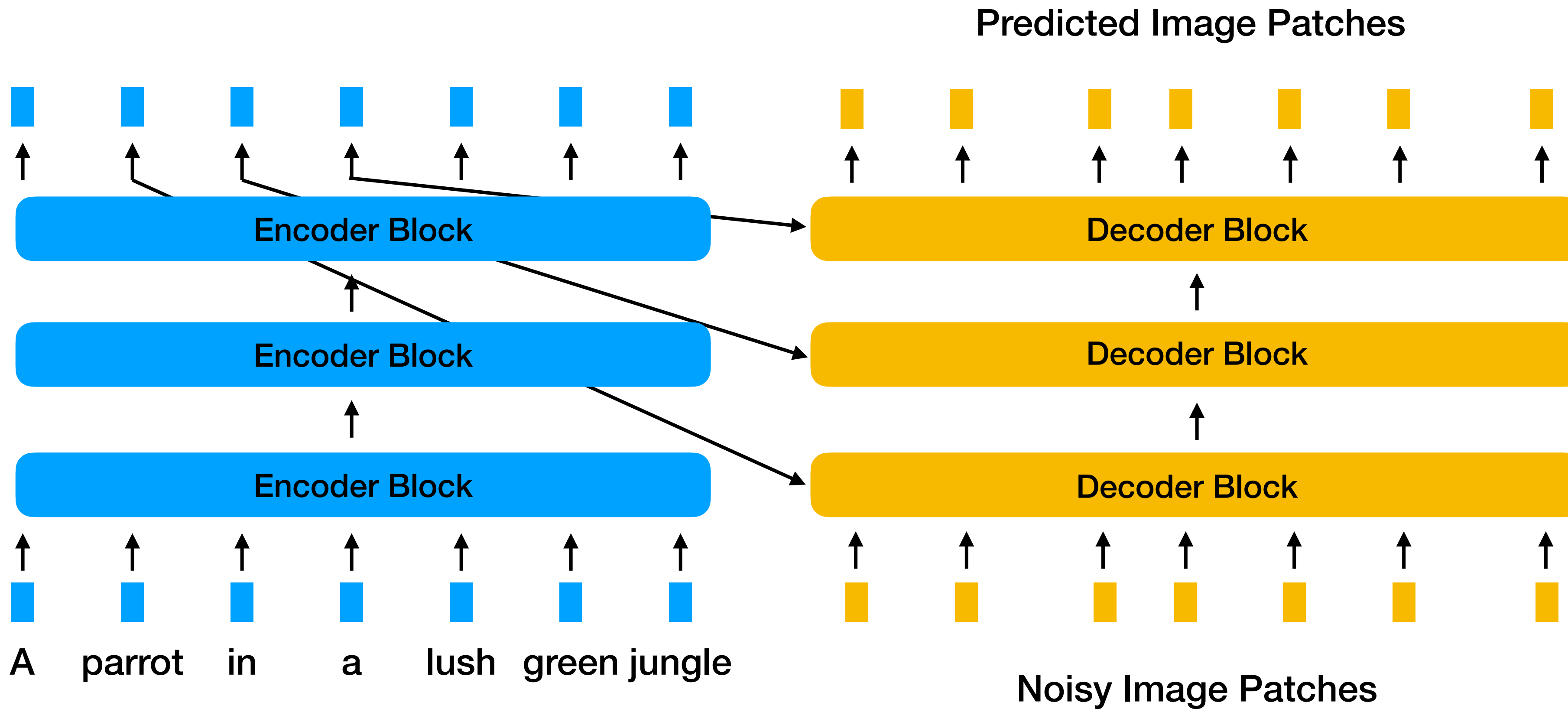
Encoder-Decoder Architecture

Attention is All You Need



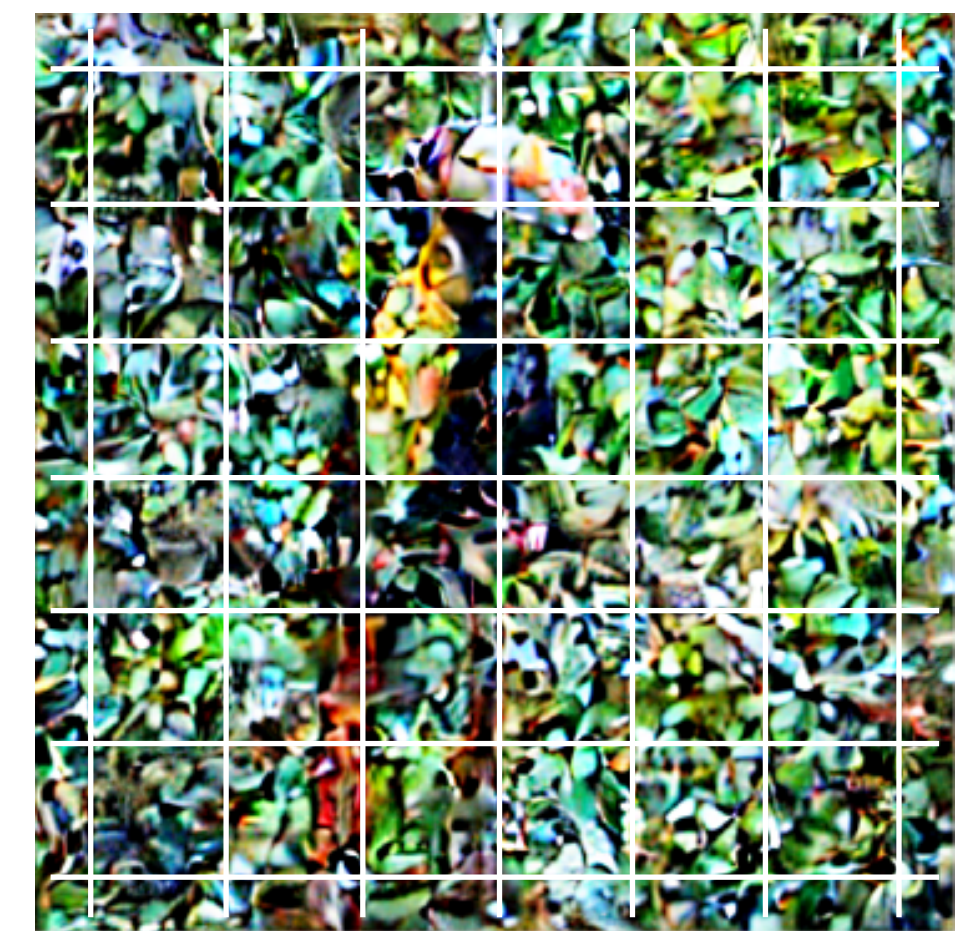
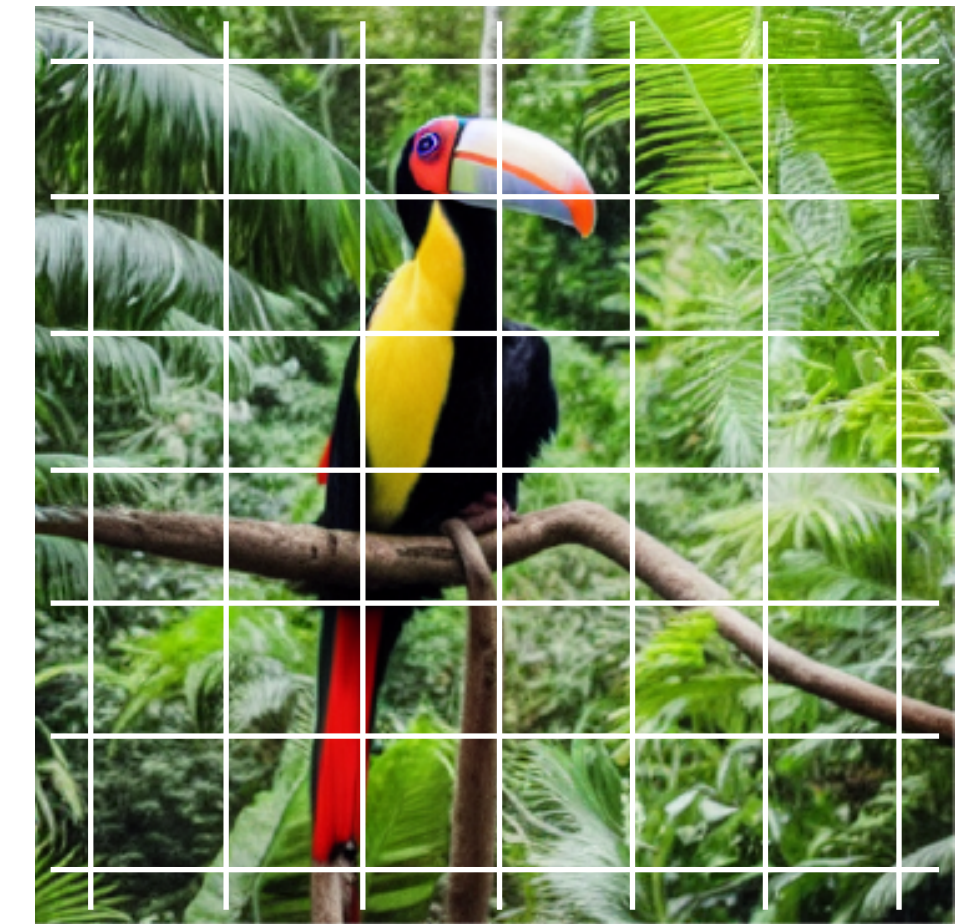
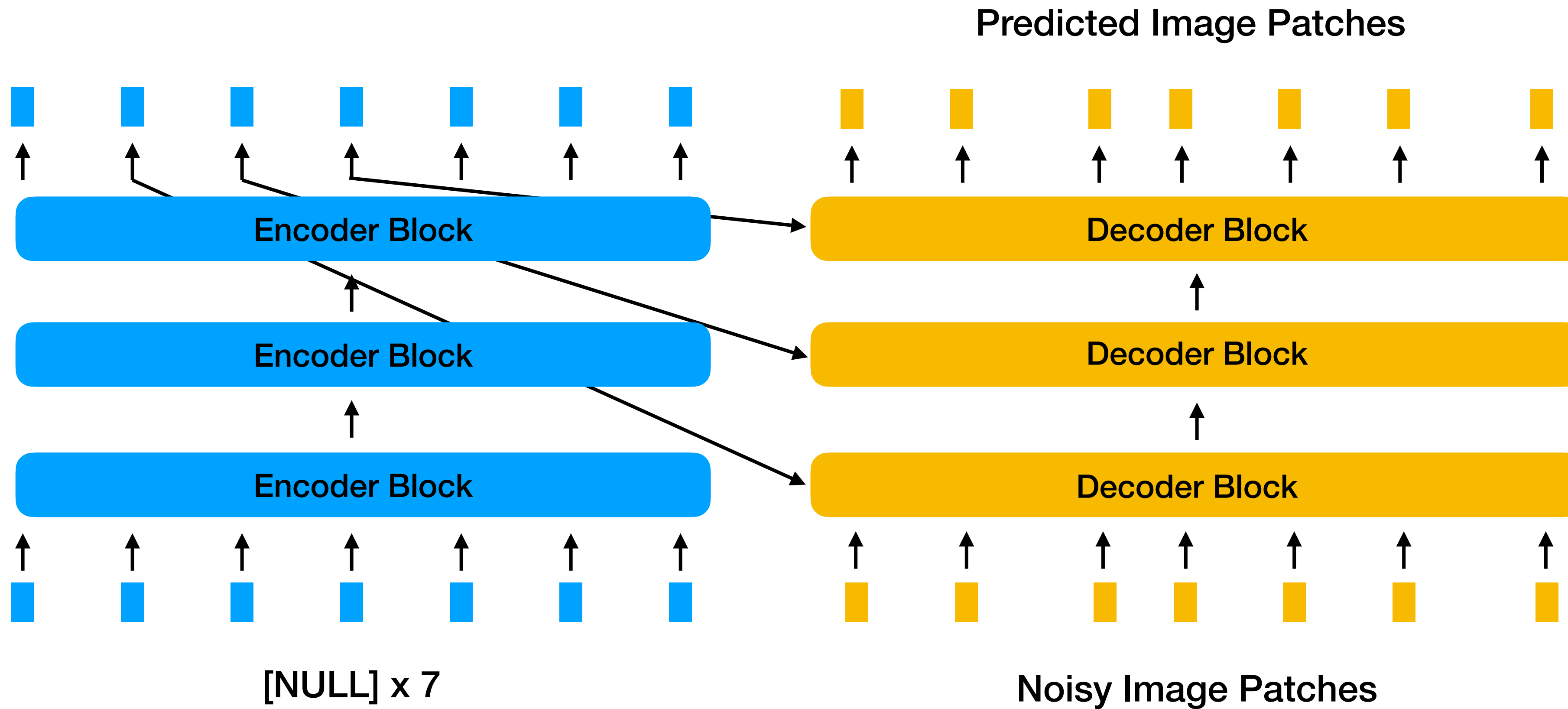
Diffusion Transformer Architecture

DiT, PixArt Alpha, MMDiT, etc.



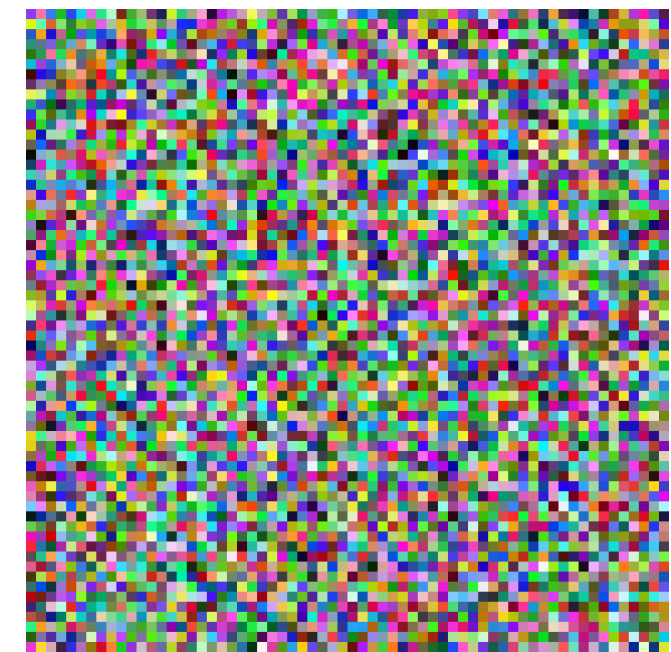
Dropout Conditioning

Maybe 30% of Training Samples

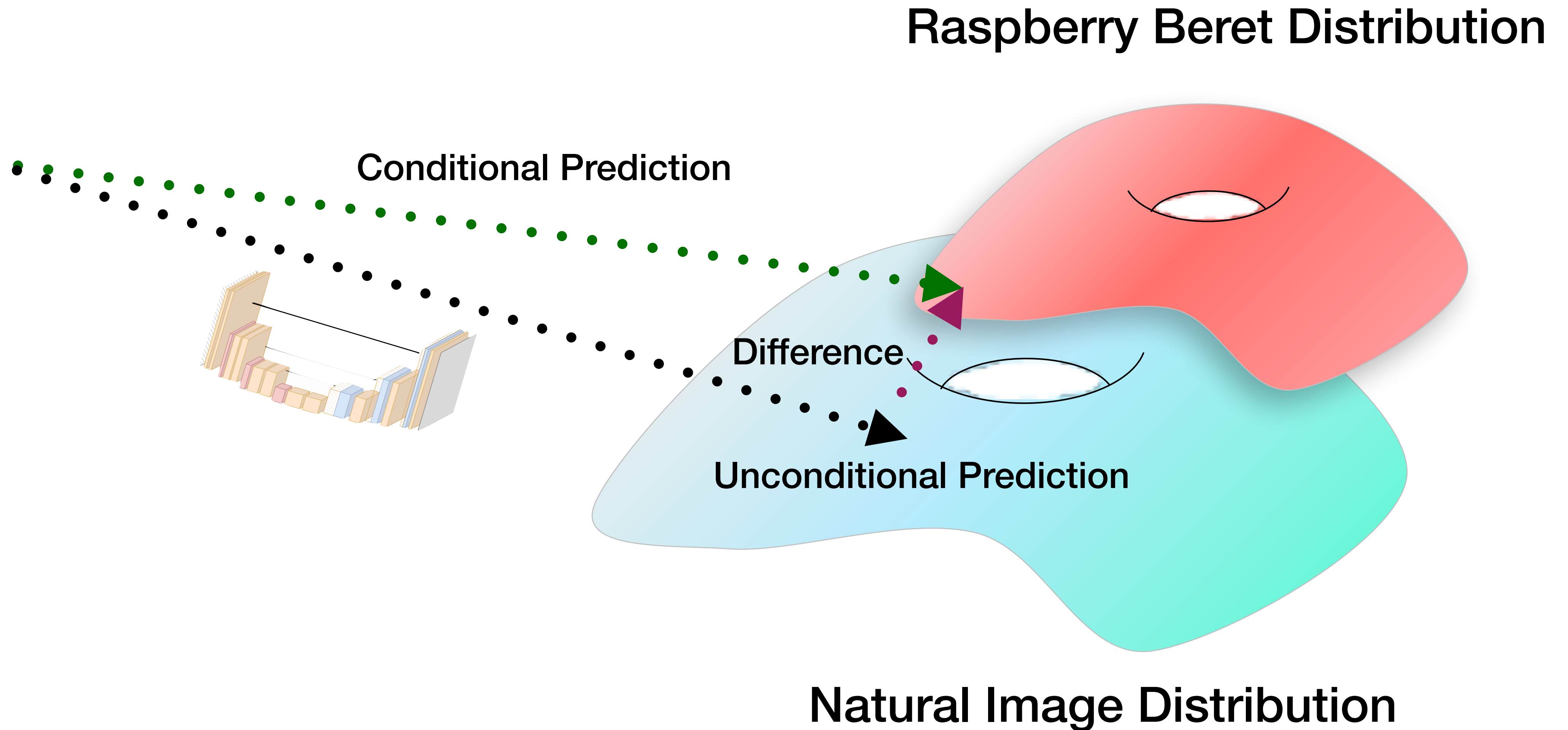


Conditional Diffusion

Model Knows Unconditional, Text-Conditioned Distributions

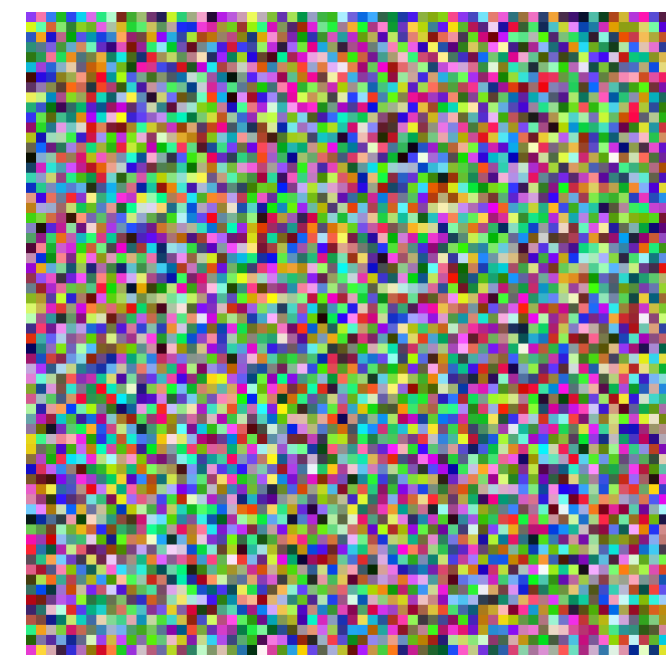


\mathbf{X}_t Noise Image

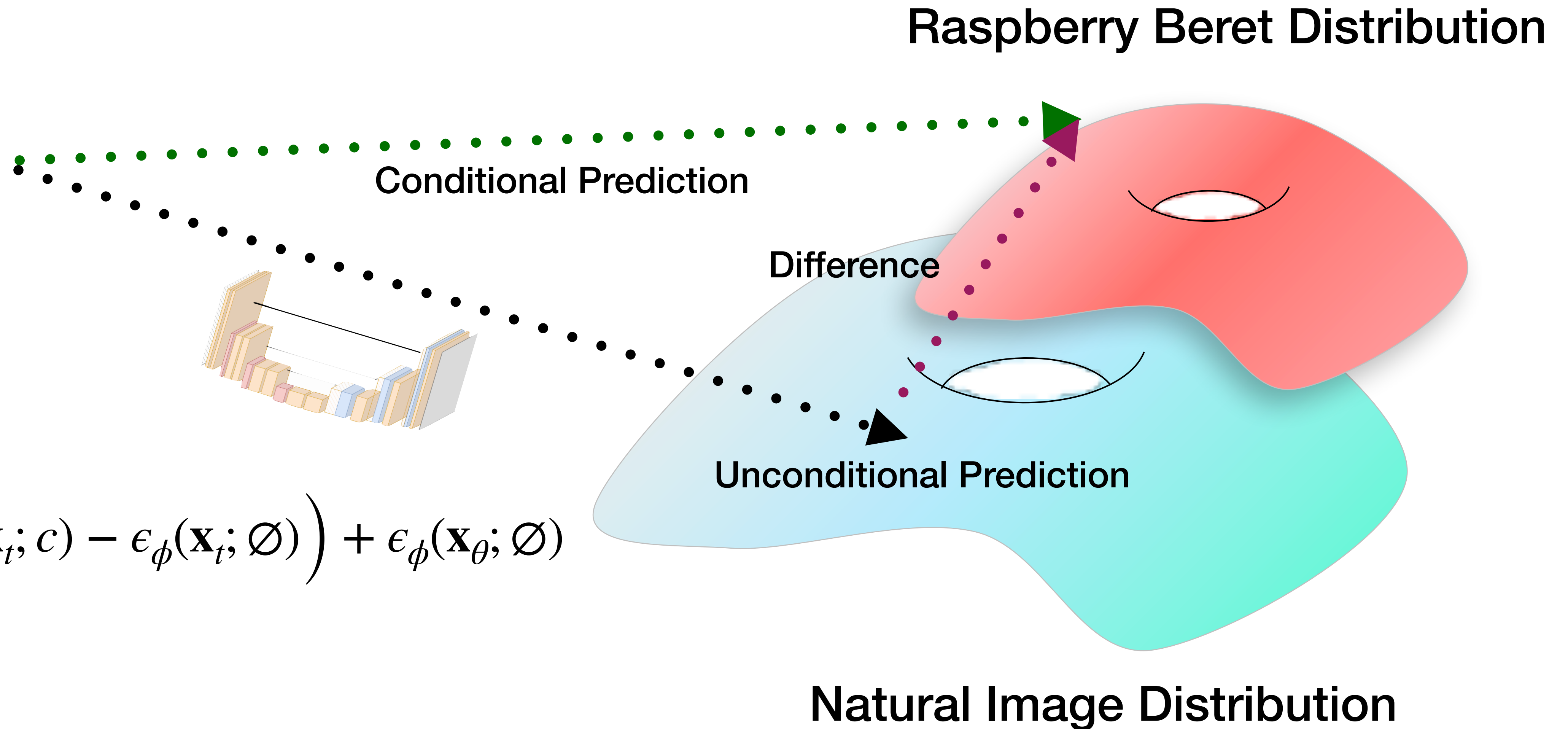


Classifier-Free Guidance

Model Knows Unconditional, Text-Conditioned Distributions



\mathbf{X}_t Noise Image



Raspberry Beret Distribution

Conditional Prediction

Difference

Unconditional Prediction

Natural Image Distribution

$$\epsilon_{CFG} = s \cdot \left(\epsilon_{\phi}(\mathbf{x}_t; c) - \epsilon_{\phi}(\mathbf{x}_t; \emptyset) \right) + \epsilon_{\phi}(\mathbf{x}_{\theta}; \emptyset)$$

ϵ_{ϕ} Predicted Noise

\emptyset Null Prompt

c Target Prompt

s CFG Scale

Classifier-Free Guidance

TLDR: Amplify Delta Between Conditional, Unconditional Predictions

$$\epsilon_{\phi, CFG} = s \cdot \left(\epsilon_{\phi}(\mathbf{x}_t; c) - \epsilon_{\phi}(\mathbf{x}_t; \emptyset) \right) + \epsilon_{\phi}(\mathbf{x}_\theta; \emptyset)$$

ϵ_{ϕ} Predicted Noise

\emptyset Null Prompt

c Target Prompt

s CFG Scale

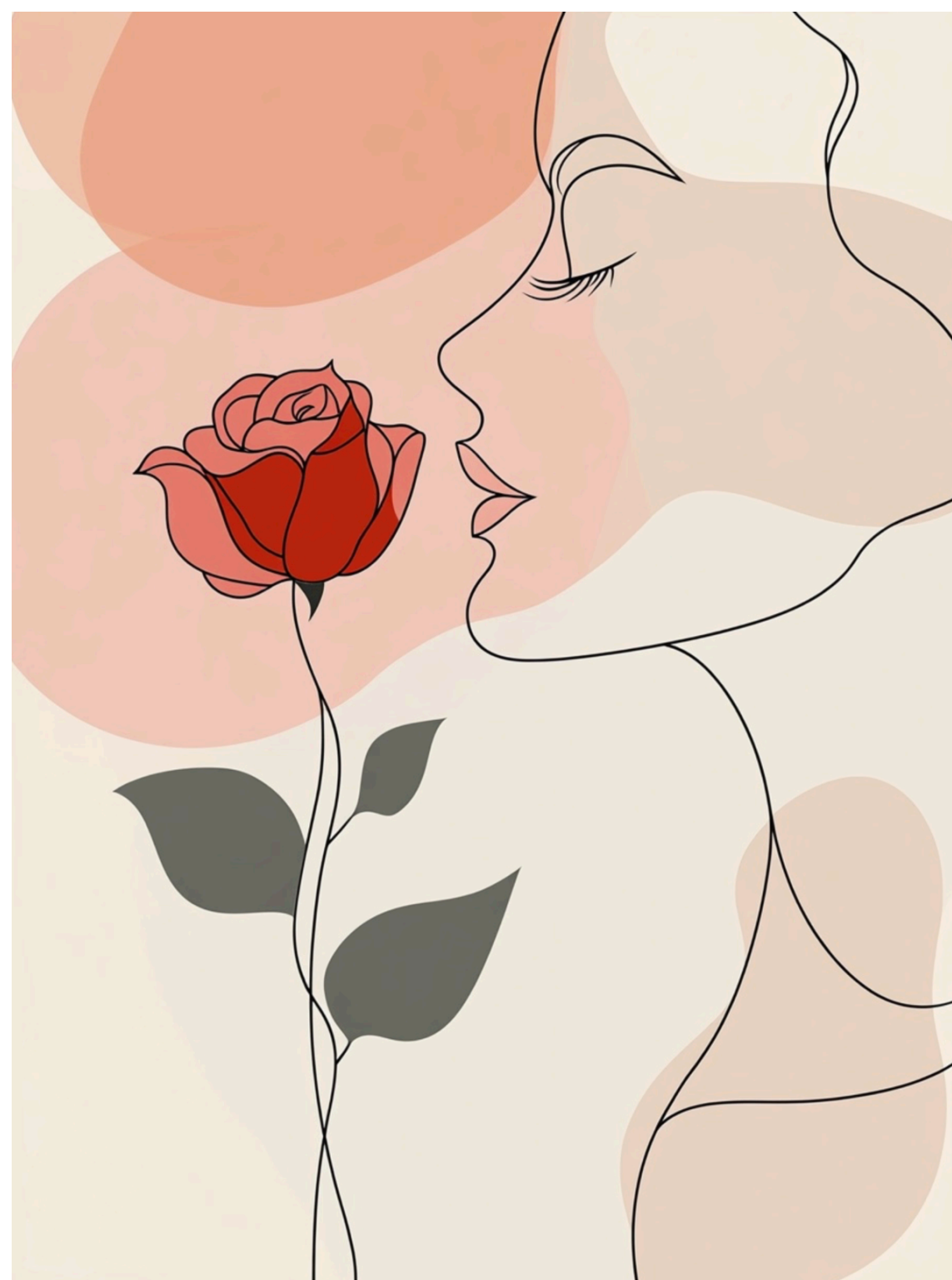
Classifier-Free Guidance

Drives Generations Toward Conditional Mode



$s = 1$ (Conditional Generation)

$s = 3$ (Conditional Generation)



Text-Conditioned Diffusion Samples
(Midjourney)



Text-Conditioned Diffusion Sample

(Veo 2)

$$p(x | c)$$

Guidance Lets Us Sample a Conditional Distribution

$$p(x | c)$$

Guidance Lets Us Sample a Conditional Distribution

What if we're creative about the conditioning?

AI Creations - OmniHuman



Condition on the First frame and Audio, Generate the Video



Generate a playable world
set in a futuristic city

Condition on the First Frame and Actions, Generate an Interactive Environment

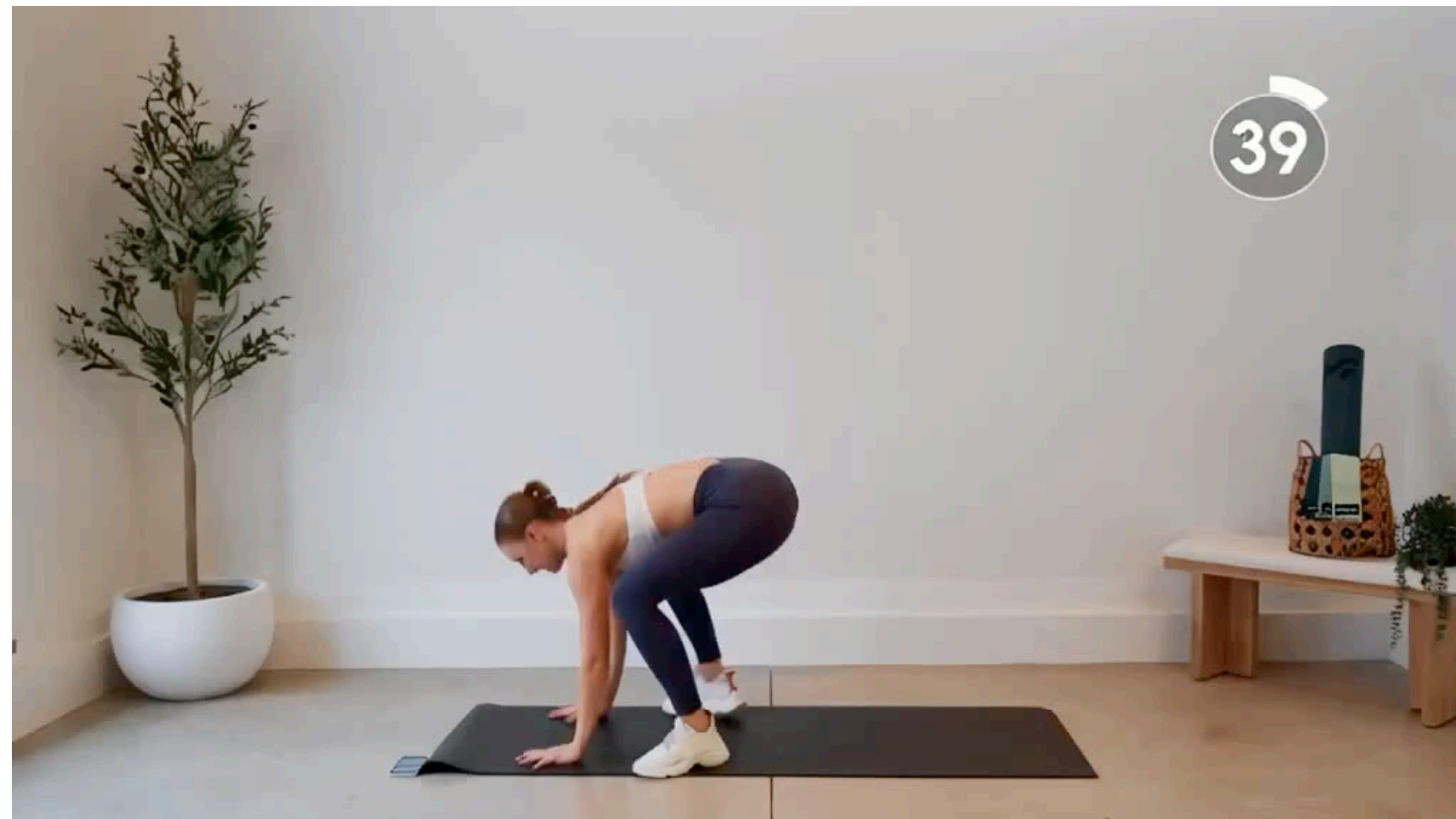


Condition on Video Pixels, Generate Depth Map

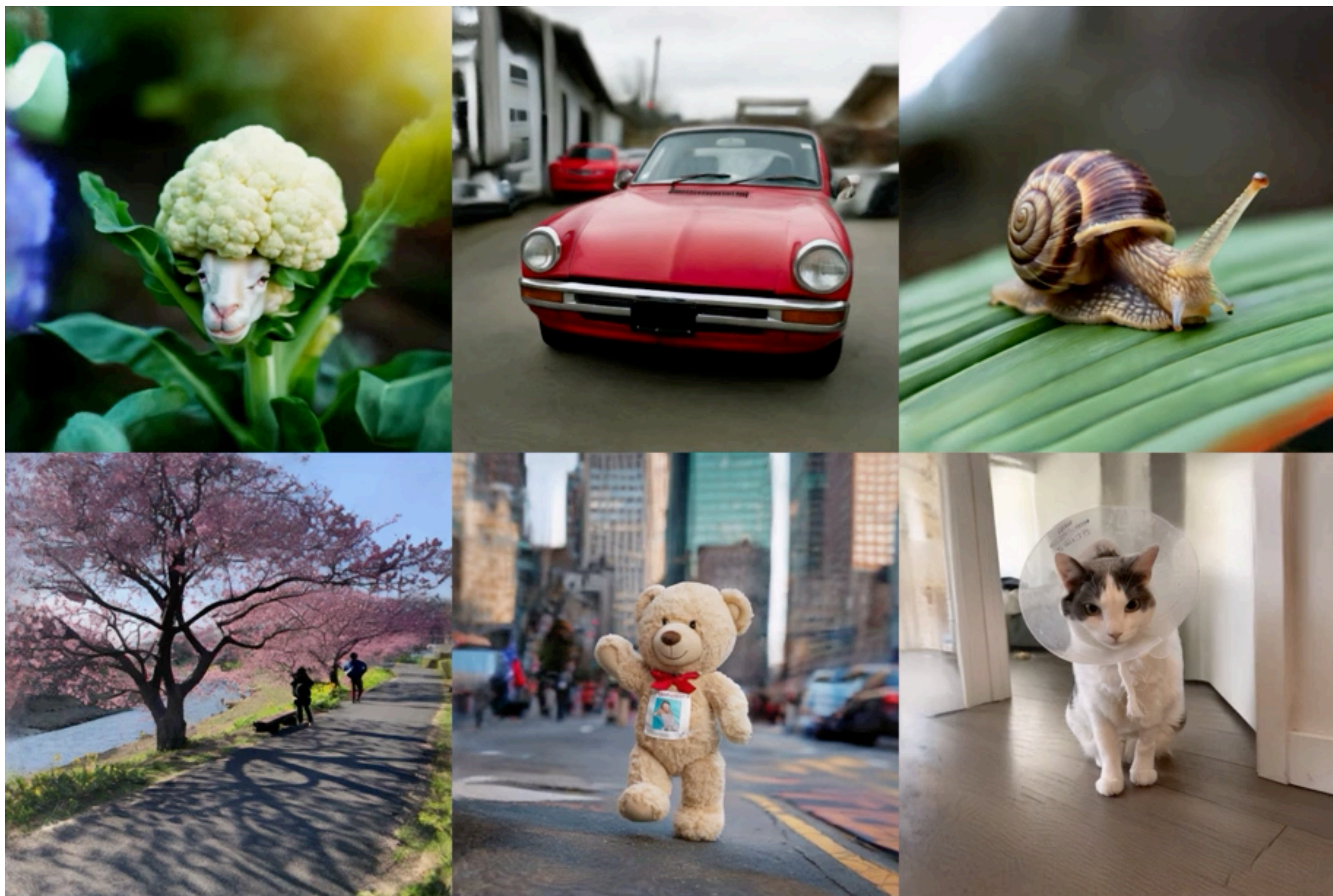
Conditioning



Samples



Condition on Start and End Frame, Generate Intermediate Frames



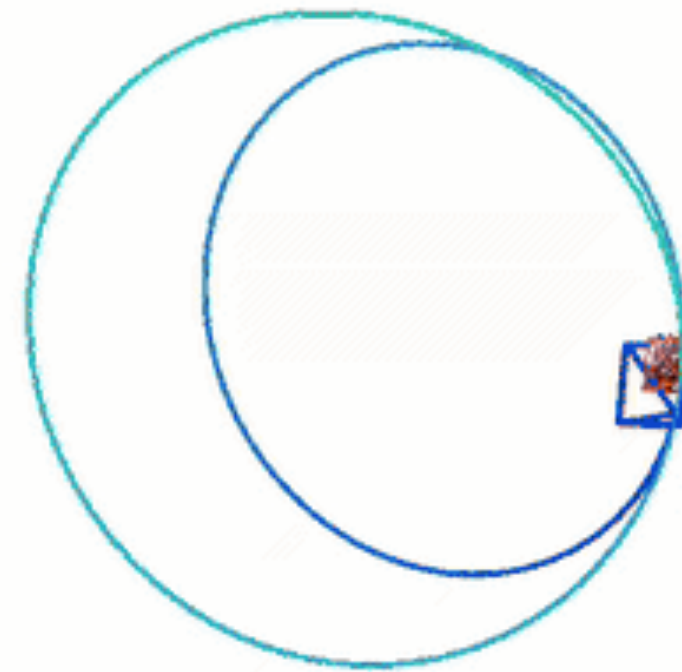
Condition on Image + Camera, Generate Novel Views

"CAT3D: Create Anything in 3D with Multi-View Diffusion Models" Gao et al. (2024)

INPUT VIEW



CAMERA CONTROL



Condition on Image + Camera, Generate Novel Views



Condition on Ego View, Generate Human Poses