



University  
of Basel

Department of  
Biomedical Engineering

# A gentle introduction to **Diffusion Models for Medical Image Analysis**

30.01.2024

Paul Friedrich, MSc

Julia Wolleb, PhD

*Center for medical Image Analysis and Navigation (CIAN), University of Basel*

# A Quick Recap on Generative Modelling

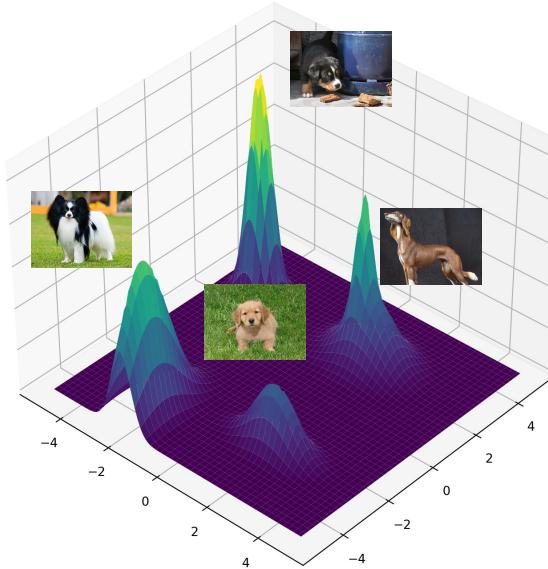


...



Data samples  
(e.g. Stanford Dogs)

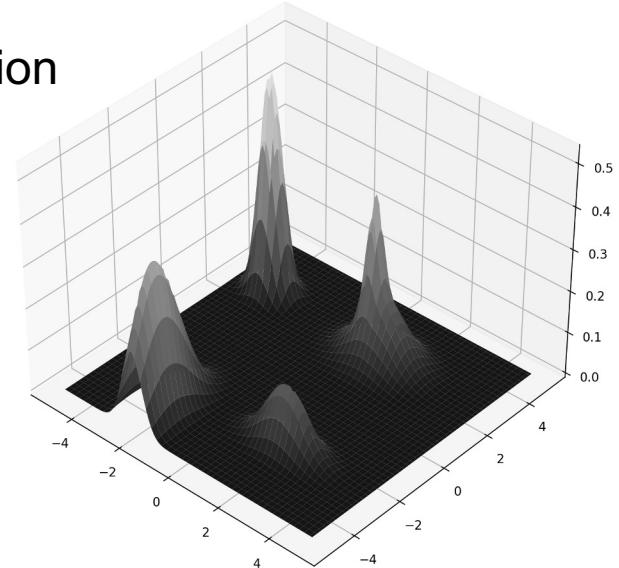
# A Quick Recap on Generative Modelling



Data distribution  
(unknown)

$\approx$

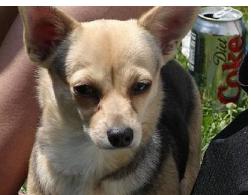
Model distribution



Datapoints are *i.i.d.* samples  
of this underlying data distribution

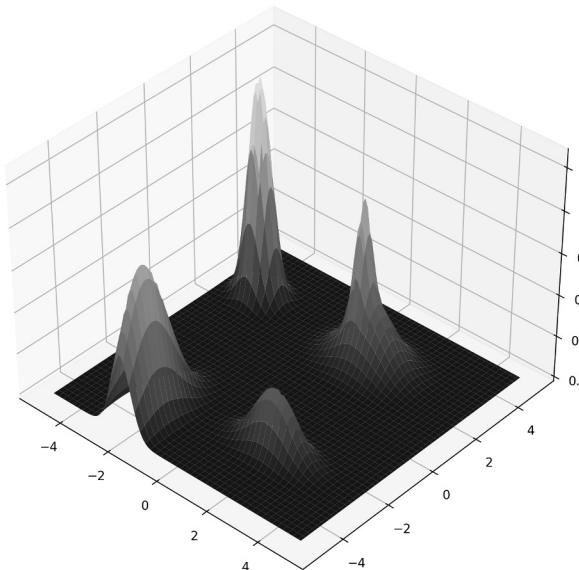
We can define a parameterized distribution  
that we tune to be close to the data distribution

# A Quick Recap on Generative Modelling



Novel data points

Sampling



Model distribution  
= **Generative Model**

Probability evaluation



High  
probability



Low  
probability

# The Landscape of Deep Generative Models

Denoising Diffusion Models

Generative Adversarial Networks

Variational Autoencoders

Normalizing Flows

Energy-based Models

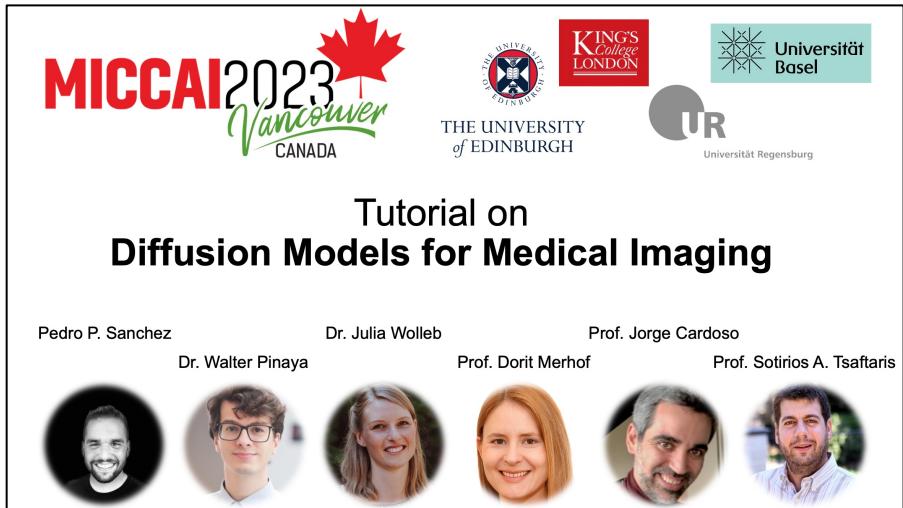
Autoregressive Models

Restricted Boltzmann Machines

# Agenda

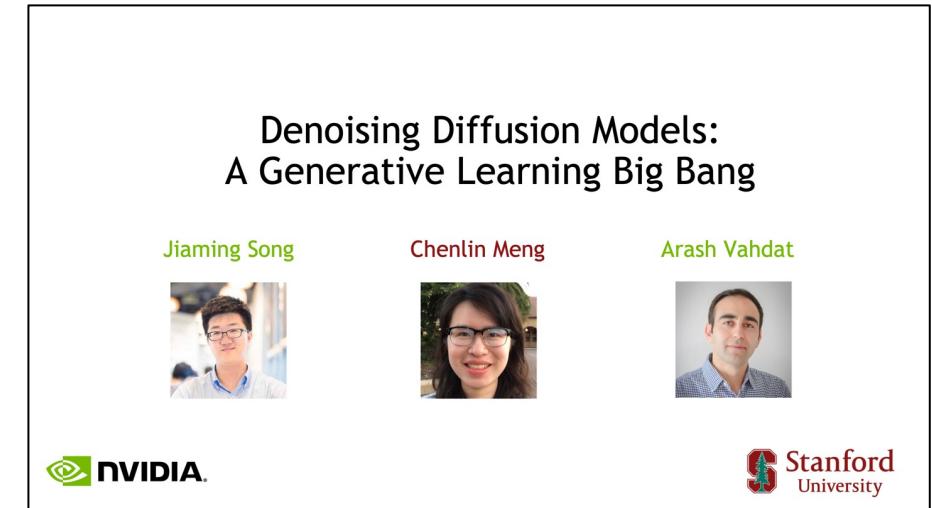
- ❑ Introduction to Diffusion Models [~30 min]
  - ❑ Physical Intuition & General Concepts
  - ❑ Denoising Diffusion Probabilistic Models
  - ❑ A Score-Based View on Diffusion Models
- ❑ Advanced Topics [~30 min]
  - ❑ Sampling Strategies
  - ❑ Inference-time Conditioning
  - ❑ Training-time Conditioning
- ❑ Applications in Medical Imaging [~30 min]
  - ❑ Synthesis
  - ❑ Inpainting
  - ❑ Segmentation
  - ❑ Anomaly Detection
  - ❑ Reconstruction
  - ❑ Registration

# Acknowledgement



The slide features the MICCAI 2023 logo (Vancouver, Canada) at the top left. To the right are logos for The University of Edinburgh, King's College London, Universität Basel, and Universität Regensburg. Below these, the title "Tutorial on Diffusion Models for Medical Imaging" is centered. Under the title, six names are listed with their corresponding circular profile pictures: Pedro P. Sanchez, Dr. Walter Pinaya, Dr. Julia Wolleb, Prof. Jorge Cardoso, Prof. Dorit Merhof, and Prof. Sotirios A. Tsaftaris.

MICCAI 2023 - Tutorial



The slide features the title "Denoising Diffusion Models: A Generative Learning Big Bang" in large bold letters. Below the title are three names: Jiaming Song, Chenlin Meng, and Arash Vahdat, each with a small circular profile picture. At the bottom left is the NVIDIA logo, and at the bottom right is the Stanford University logo.

CVPR 2023 - Tutorial

**Some slides are adopted/inspired by/copied from these very nice tutorials!**

# The Physical Intuition behind Diffusion Models

---

**Deep Unsupervised Learning using  
Nonequilibrium Thermodynamics**

---

**Jascha Sohl-Dickstein**

Stanford University

JASCHA@STANFORD.EDU

**Eric A. Weiss**

University of California, Berkeley

EWEISS@BERKELEY.EDU

**Niru Maheswaranathan**

Stanford University

NIRUM@STANFORD.EDU

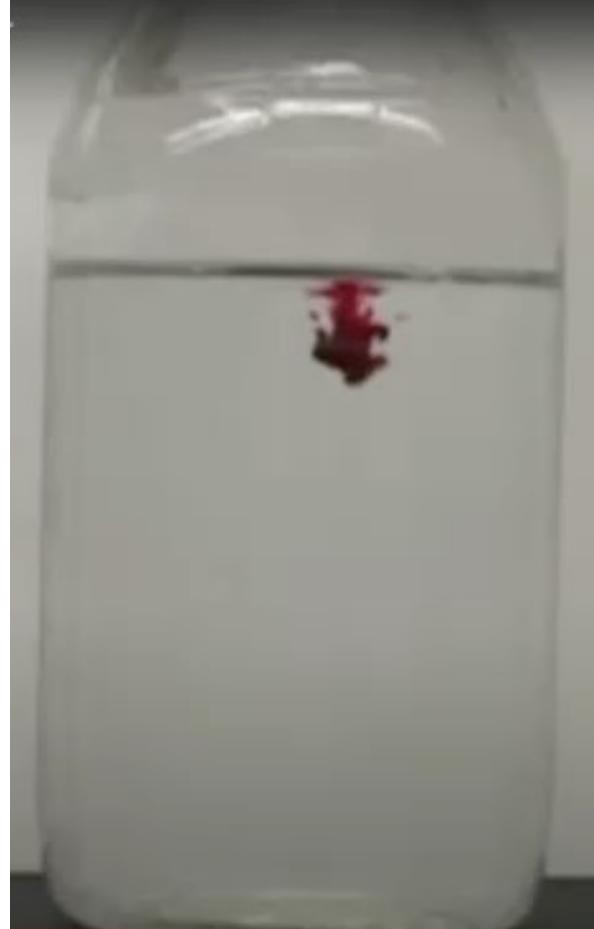
**Surya Ganguli**

Stanford University

SGANGULI@STANFORD.EDU

PMLR (2015)

# The Physical Intuition behind Diffusion Models (Macroscopic)



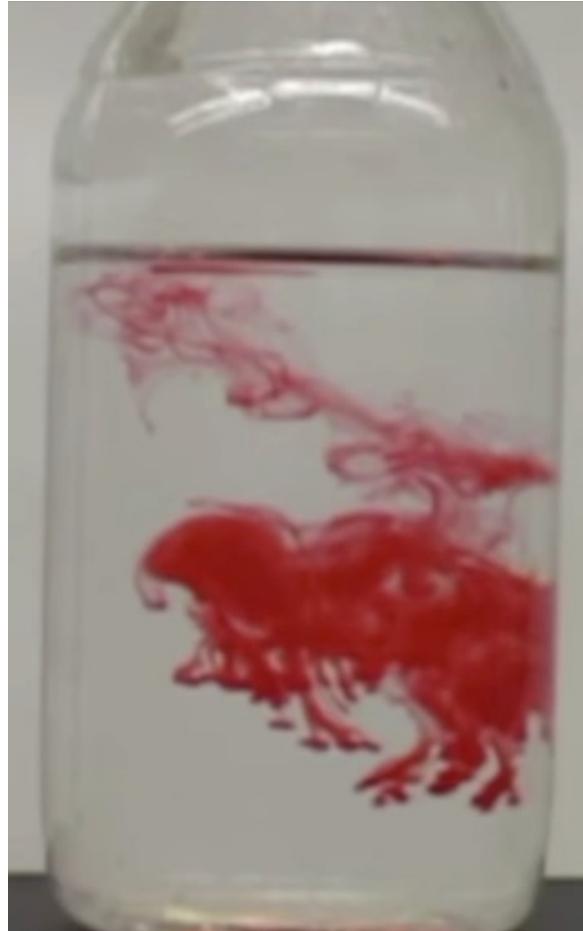
The dye density represents our probability density.

Goal: We want to *learn this probability density*

**Observing this diffusion process:**

- Original data distribution is perturbed over time
- Data distribution → Uniform distribution (mapping to a simple distribution)

# The Physical Intuition behind Diffusion Models (Macroscopic)



The dye density represents our probability density.

Goal: We want to *learn this probability density*

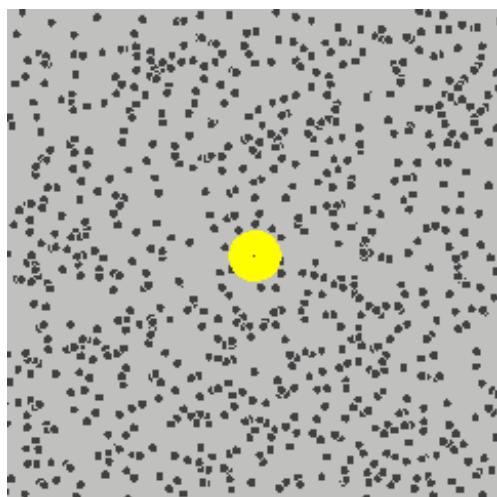
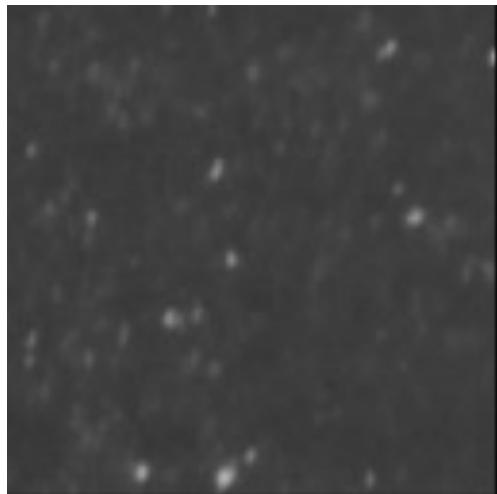
**Observing this diffusion process:**

- Original data distribution is perturbed over time
- Data distribution → Uniform distribution (mapping to a simple distribution)

**Can we learn to revert this process (run it backwards)?**

- Uniform distribution → Data distribution
- **Yes**, but we first need a way to model the system.

# The Physical Intuition Behind Diffusion Models (Microscopic)



We can try to model the diffusion process by **modelling the Brownian motion** of single particles.

- We can observe that the position updates follow **small Gaussians**
  - This holds true for the **forward as well as the reverse process** (for small enough  $\Delta t$ )
  - We can define a known diffusion process with a chain of Gaussian position updates
  - We try to learn the reverse process by estimating the mean and the covariance of the backward steps
- Same mechanism is used in the Diffusion Models we will see today!

# Denoising Diffusion Probabilistic Models

---

## Denoising Diffusion Probabilistic Models

---

Jonathan Ho  
UC Berkeley  
[jonathanho@berkeley.edu](mailto:jonathanho@berkeley.edu)

Ajay Jain  
UC Berkeley  
[ajayj@berkeley.edu](mailto:ajayj@berkeley.edu)

Pieter Abbeel  
UC Berkeley  
[pabbeel@cs.berkeley.edu](mailto:pabbeel@cs.berkeley.edu)

NeurIPS (2020)

---

## Improved Denoising Diffusion Probabilistic Models

---

Alex Nichol<sup>\*†</sup> Prafulla Dhariwal<sup>\*†</sup>

ICML (2021)

---

## Diffusion Models Beat GANs on Image Synthesis

---

Prafulla Dhariwal<sup>\*</sup>  
OpenAI  
[prafulla@openai.com](mailto:prafulla@openai.com)

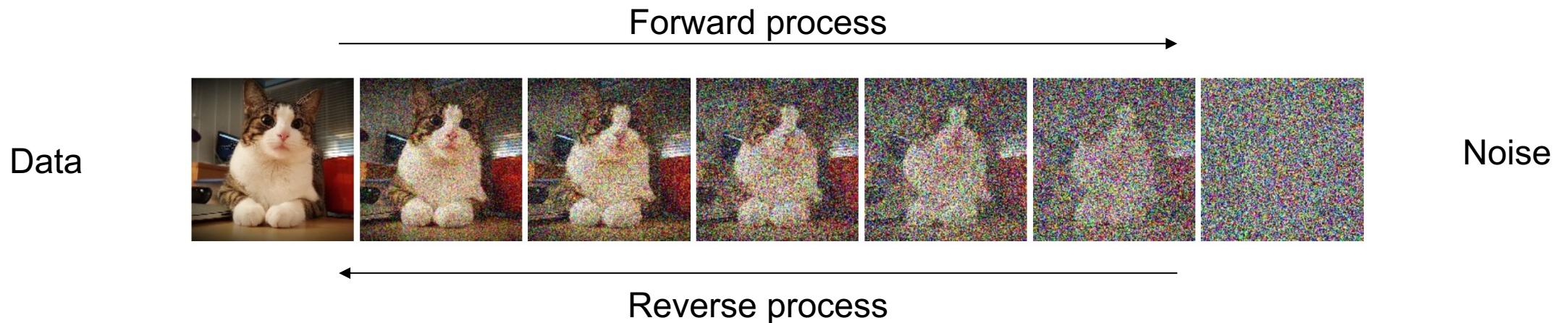
Alex Nichol<sup>\*</sup>  
OpenAI  
[alex@openai.com](mailto:alex@openai.com)

NeurIPS (2021)

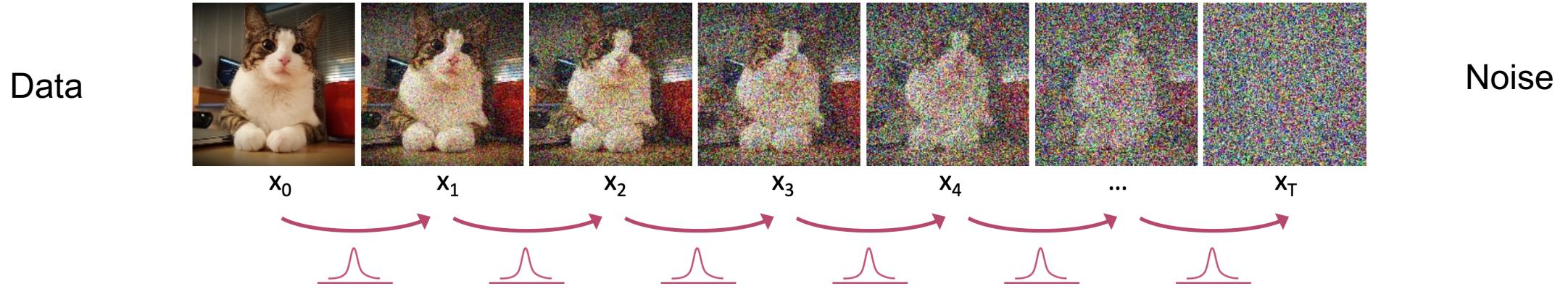
# How do Diffusion Models work?

Diffusion Models consist of two main components:

- A ***fixed forward diffusion process*** that gradually adds noise to the image
- A ***learned reverse diffusion process*** that gradually removes noise from the image



# The Forward Diffusion Process



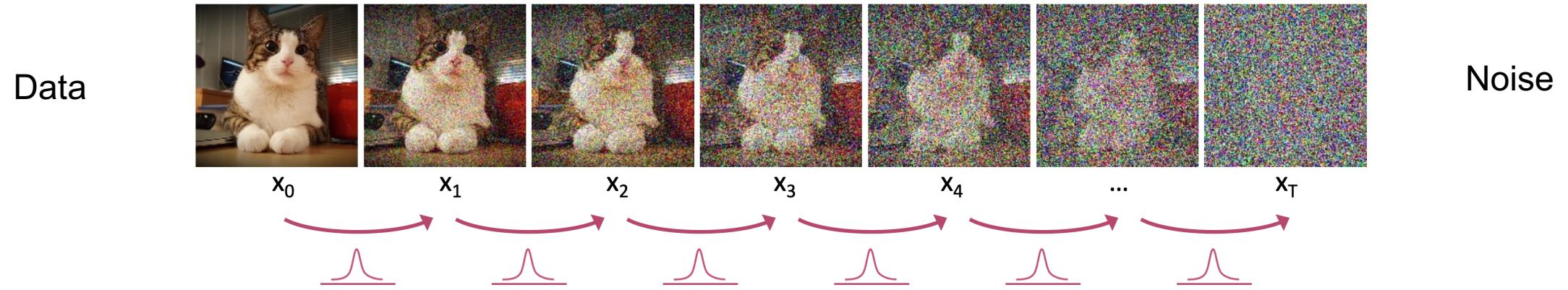
We model the **forward process** as a Markov chain:

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$$

with each transition being a parameterized Gaussian:

$$q(x_t|x_{t-1}) = N(\sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$$

# The Forward Diffusion Process



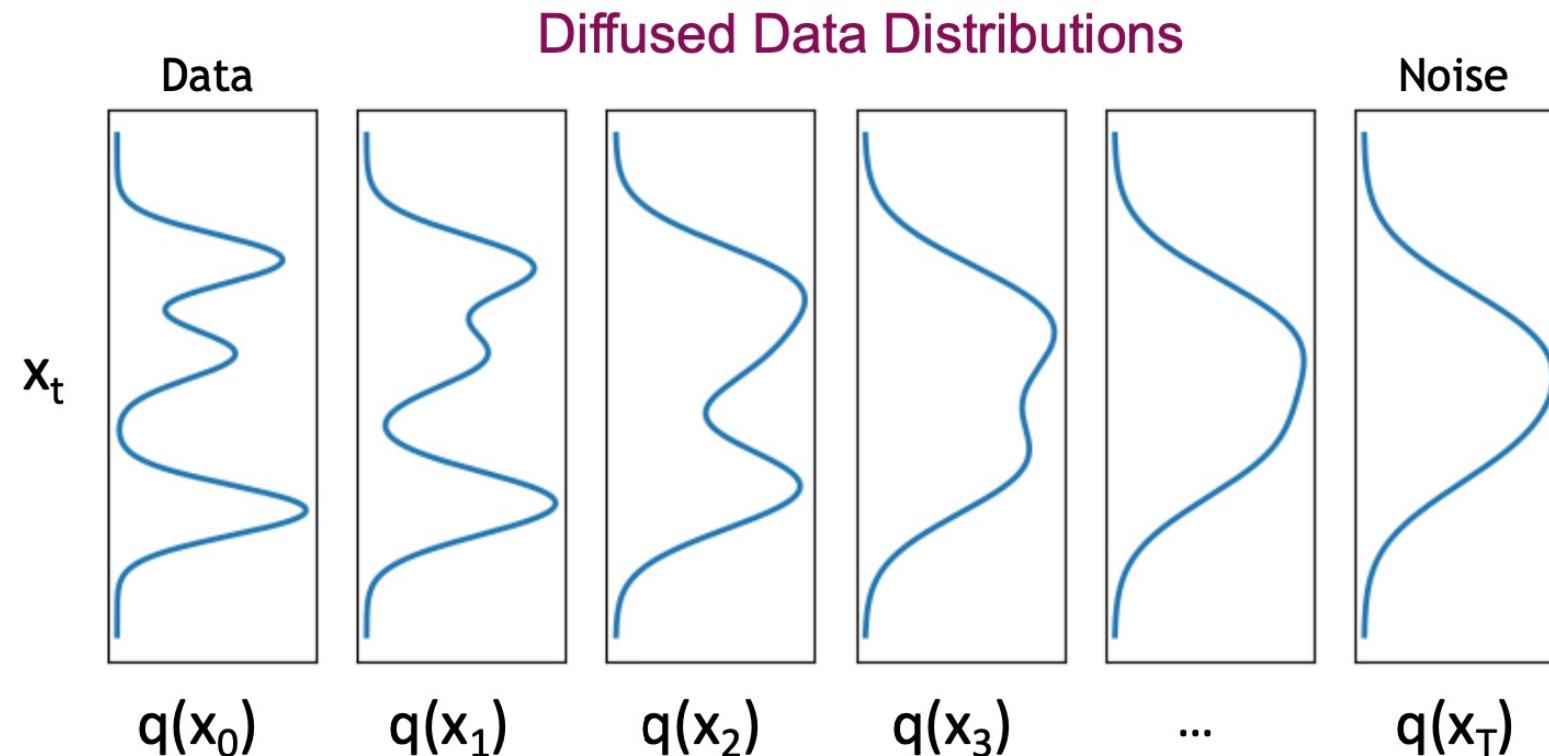
$$q(x_t | x_{t-1}) = N(\sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$$

We define a variance schedule for  $\beta_t$

$$q(x_T | x_0) \approx N(0, I)$$

We usually choose  $T \approx 1000$  (this remains a design choice)

# The Forward Diffusion Process



We defined a forward process that transforms our data distribution to noise.

# The Reverse Diffusion Process

Recall, the diffusion process is designed in a way that:

$$q(x_T) \approx N(0, I)$$

We could generate new samples by:

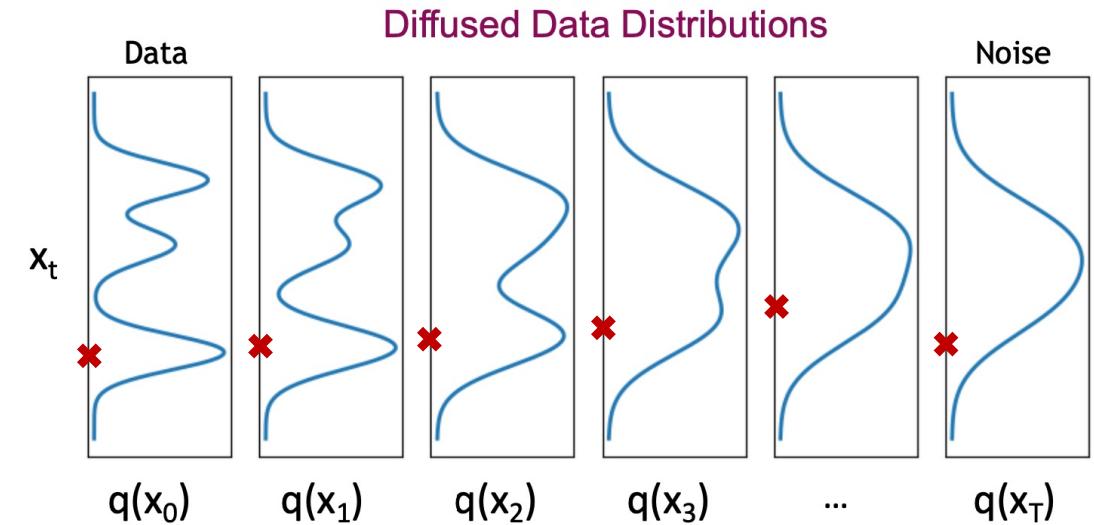
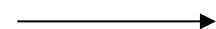
- Sampling  $x_T$ :

$$x_T \sim N(0, I)$$

- Iteratively sample  $x_{t-1}$  for  $T$  timesteps:

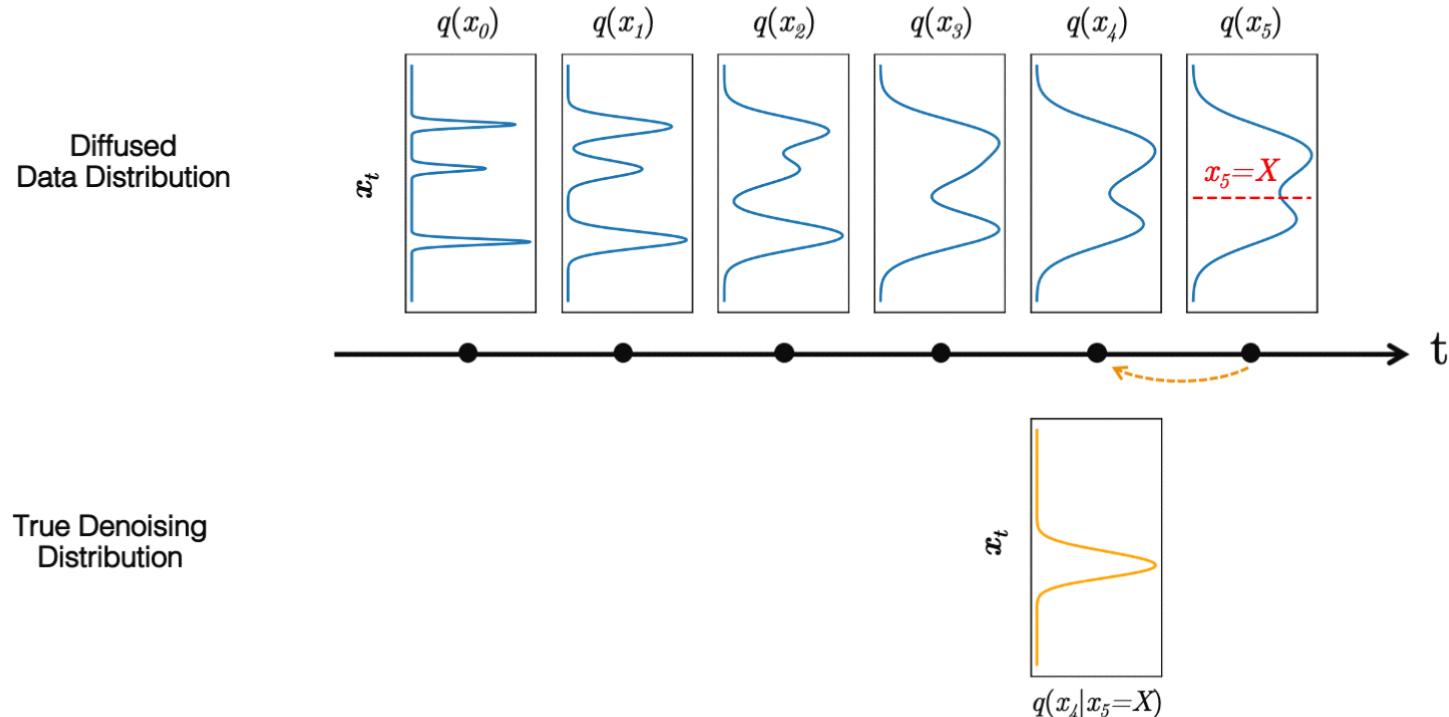
$$x_{t-1} \sim q(x_{t-1} | x_t)$$

True denoising distribution



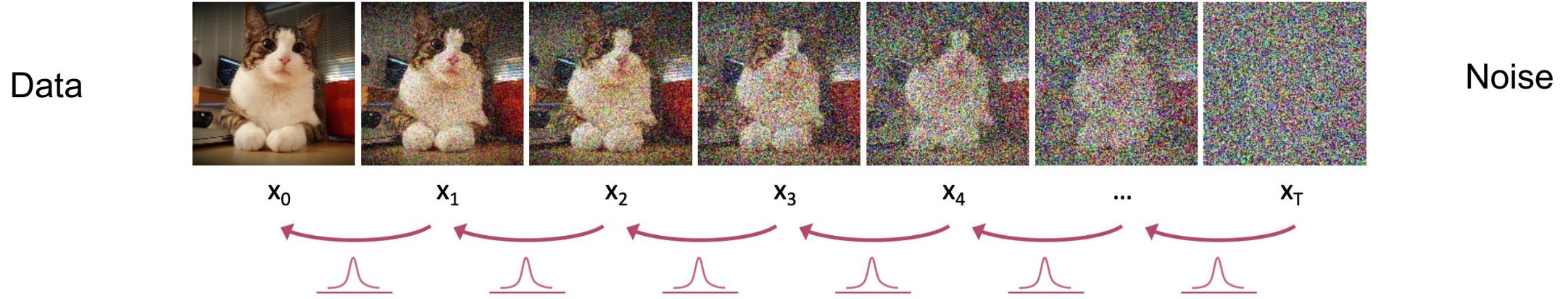
This distribution is unknown. Can we estimate it?

# The Reverse Diffusion Process



**YES!** We can approximate the true denoising distribution (as a normal distribution) for small steps.

# The Reverse Diffusion Process



We approximate the true denoising distribution  $q(x_{t-1}|x_t)$  as being normal distributed:

$$p_\theta(x_{t-1}|x_t) = N(\underbrace{\mu_\theta(x_t, t)}, \sigma_t^2 I)$$

Mean is estimated by a neural network

# How can we train such a model?

Keep in mind:  $p_\theta(x_{t-1}|x_t) = N(\mu_\theta(x_t, t), \sigma_t^2 I)$

Ho et al. (2020) found that we can parameterize  $\mu_\theta(x_t, t)$  as follows:

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right)$$

the noise to be removed

with  $\alpha_t := 1 - \beta_t$  and  $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$ .

---

## Algorithm 1 Training

---

- 1: **repeat**
- 2:  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3:  $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4:  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on  
$$\nabla_\theta \|\boldsymbol{\epsilon} - \epsilon_\theta(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$$
- 6: **until** converged

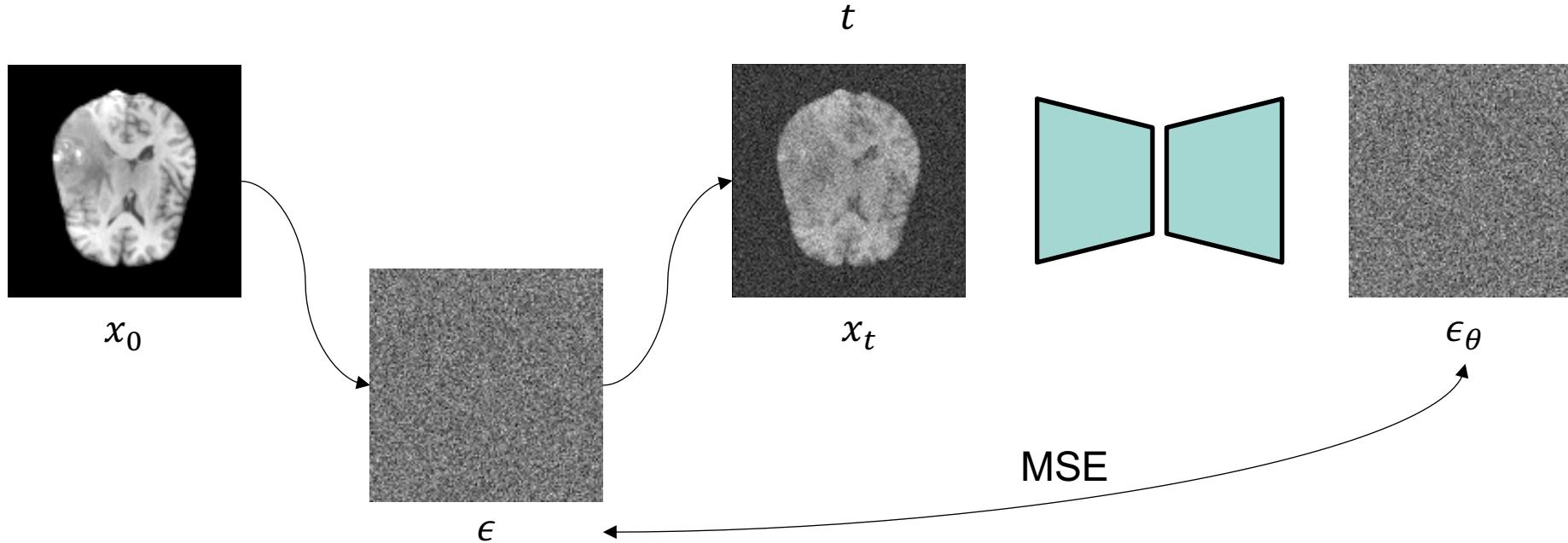
---

## Algorithm 2 Sampling

---

- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for**  $t = T, \dots, 1$  **do**
- 3:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$
- 4:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
- 5: **end for**
- 6: **return**  $\mathbf{x}_0$

# How can we train such a model?



## Algorithm 1 Training

```
1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
       $\nabla_\theta \|\epsilon - \epsilon_\theta(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \epsilon, t)\|^2$ 
6: until converged
```

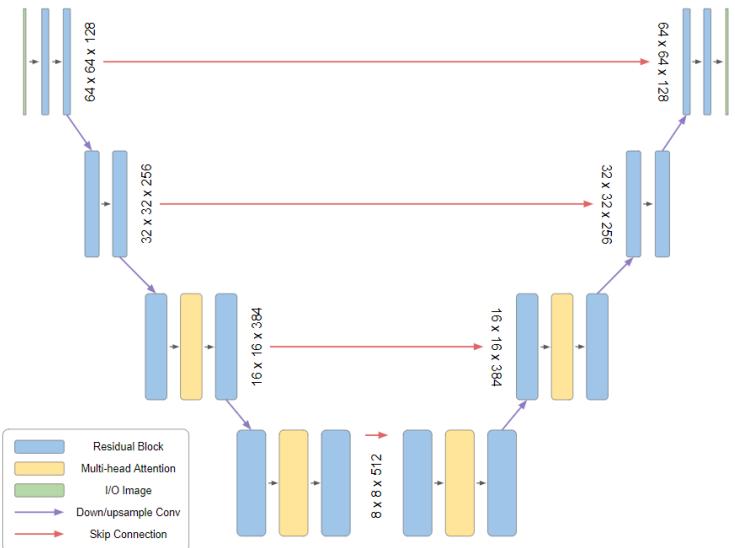
## Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

# Implementation Details

Keep in mind: We want to predict the noise to be removed from a corrupted image.

$\epsilon_\theta(x_t, t)$  is usually implemented as a U-Net:

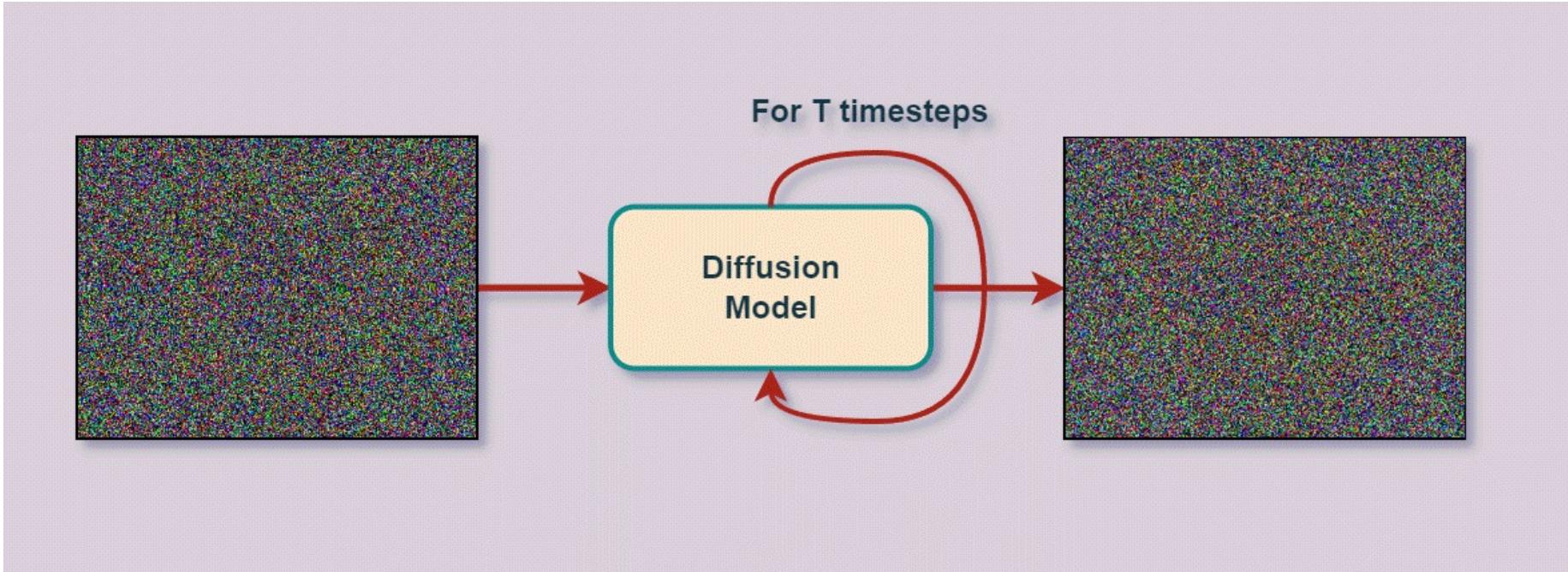


BUT, we could also use:

- Transformers
- VQ-VAEs
- ...

This remains a design choice.

# Generating Samples



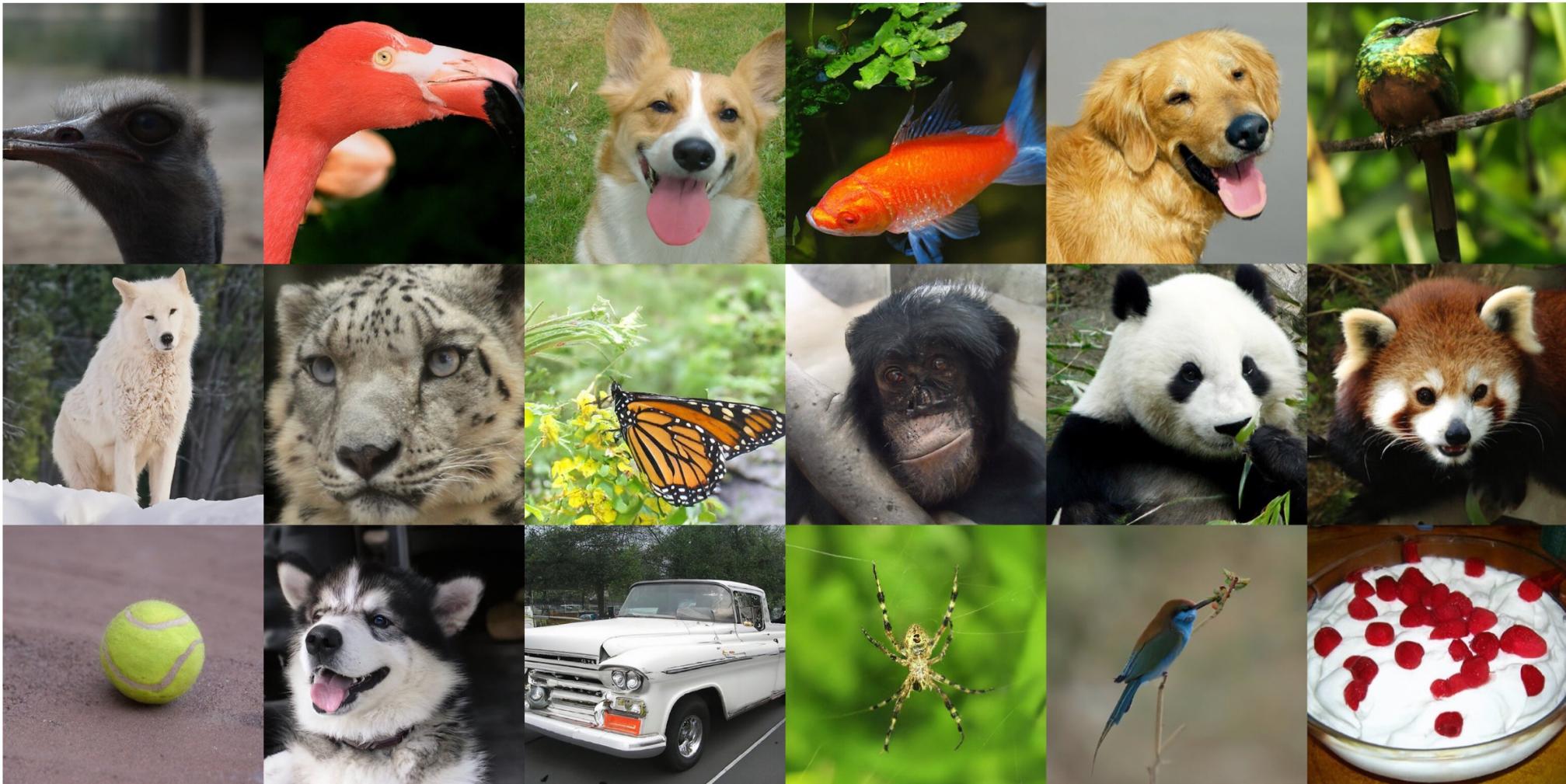
## Algorithm 1 Training

```
1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
       $\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$ 
6: until converged
```

## Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( \mathbf{x}_t - \frac{1 - \bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

# Sample Quality



Samples from model trained on ImageNet (512 x 512)

# Pros & Cons



Samples from model trained on  
CelebA-HQ (256 x 256)

## Pros:

- + High sample quality & diversity
- + Build on a strong theoretical foundation
- + Easy and stable to train (just a simple MSE Loss + just one network)

## Cons:

- Very slow (sampling usually requires multiple model evaluations)  
→ we will see some strategies to speed this process up

# Playing Around with Diffusion Models



# A Score-Based View on Diffusion Models

---

## Generative Modeling by Estimating Gradients of the Data Distribution

---

**Yang Song**  
Stanford University  
[yangsong@cs.stanford.edu](mailto:yangsong@cs.stanford.edu)

**Stefano Ermon**  
Stanford University  
[ermon@cs.stanford.edu](mailto:ermon@cs.stanford.edu)

NeurIPS (2019)

---

## SCORE-BASED GENERATIVE MODELING THROUGH STOCHASTIC DIFFERENTIAL EQUATIONS

---

**Yang Song\***  
Stanford University  
[yangsong@cs.stanford.edu](mailto:yangsong@cs.stanford.edu)

**Abhishek Kumar**  
Google Brain  
[abhishek@google.com](mailto:abhishek@google.com)

**Jascha Sohl-Dickstein**  
Google Brain  
[jaschasd@google.com](mailto:jaschasd@google.com)

**Stefano Ermon**  
Stanford University  
[ermon@cs.stanford.edu](mailto:ermon@cs.stanford.edu)

**Diederik P. Kingma**  
Google Brain  
[durk@google.com](mailto:durk@google.com)

**Ben Poole**  
Google Brain  
[pooleb@google.com](mailto:pooleb@google.com)

ICLR (2021)

---

## Maximum Likelihood Training of Score-Based Diffusion Models

---

**Yang Song\***  
Computer Science Department  
Stanford University  
[yangsong@cs.stanford.edu](mailto:yangsong@cs.stanford.edu)

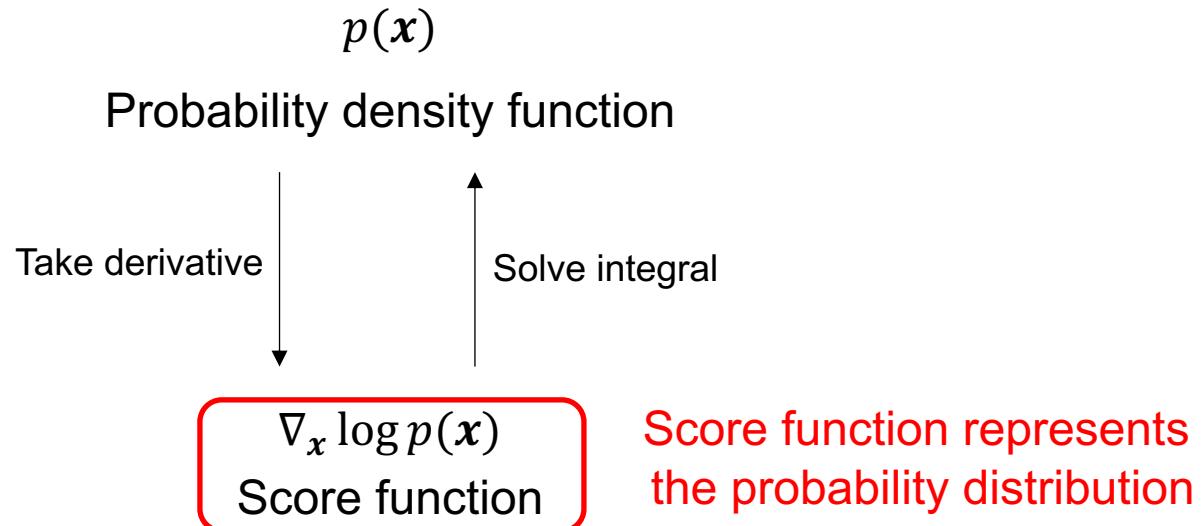
**Iain Murray**  
School of Informatics  
University of Edinburgh  
[i.murray@ed.ac.uk](mailto:i.murray@ed.ac.uk)

**Conor Durkan\***  
School of Informatics  
University of Edinburgh  
[conor.durkan@ed.ac.uk](mailto:conor.durkan@ed.ac.uk)

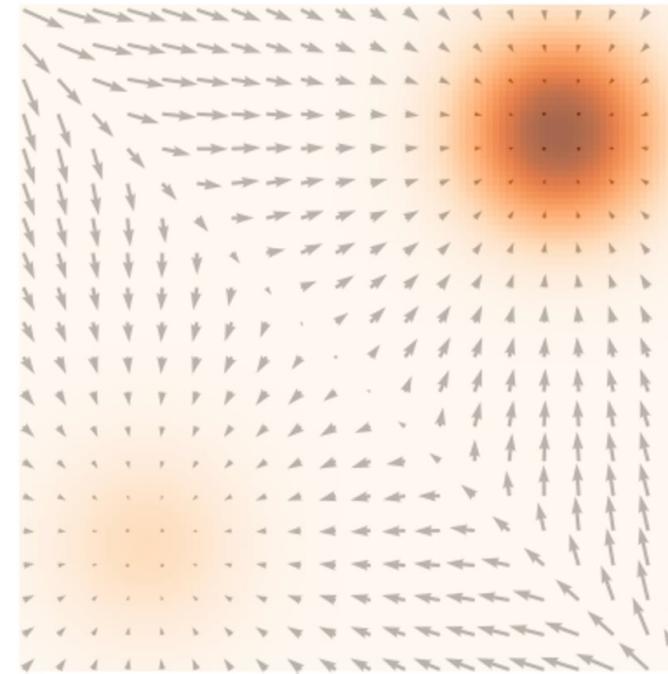
**Stefano Ermon**  
Computer Science Department  
Stanford University  
[ermon@cs.stanford.edu](mailto:ermon@cs.stanford.edu)

NeurIPS (2021)

# What do we mean when talking about a score function?

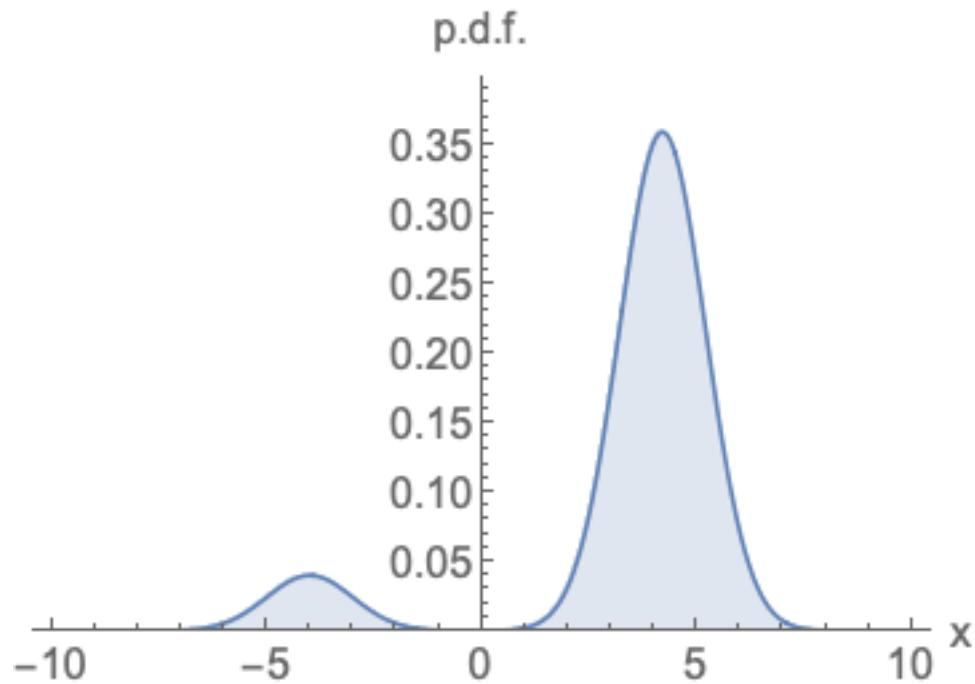


The score function **preserves all information** of the density function, but is much **easier to handle** → **why?**

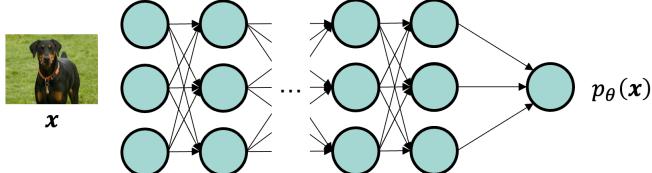


Probability density function (color coded)  
Score function (vector field)

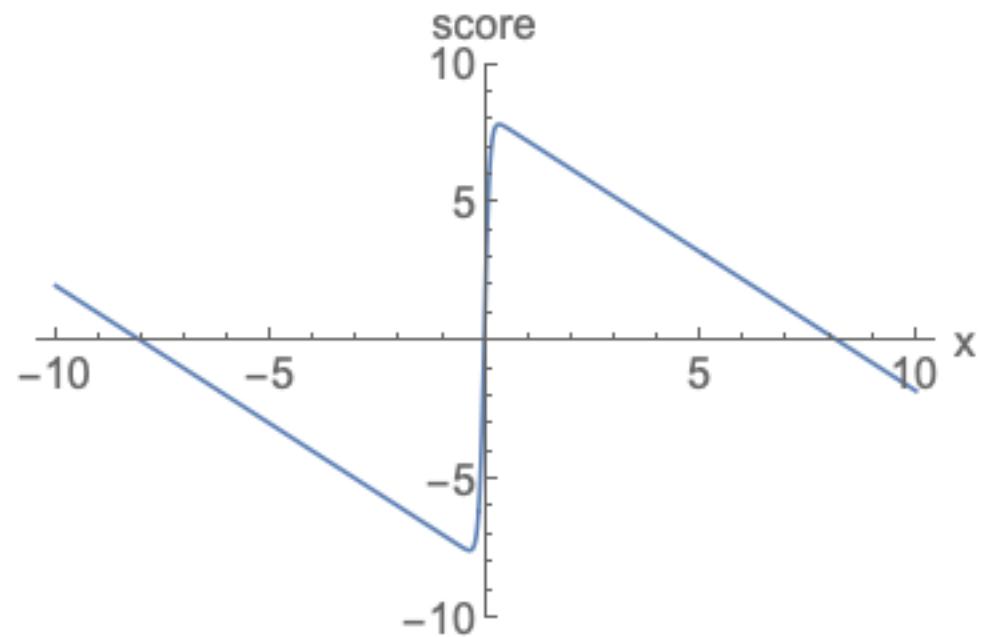
# Score functions bypass the normalizing constant



Probability density function



$$\frac{e^{f_\theta(x)}}{\cancel{Z_\theta}} = p_\theta(x)$$



Score function

$$\begin{aligned}\nabla_x \log p_\theta(x) &= \nabla_x f_\theta(x) - \cancel{\nabla_x \log Z_\theta} \\ &= \nabla_x f_\theta(x) - 0\end{aligned}$$

We always need to ensure normalization.

Score function doesn't rely on normalizing constant.

# Can we estimate such a score function from data?

We know it's possible to **train a properly normalized statistical model to estimate the data density function using maximum likelihood** → can we do something similar to estimate a score function/model?

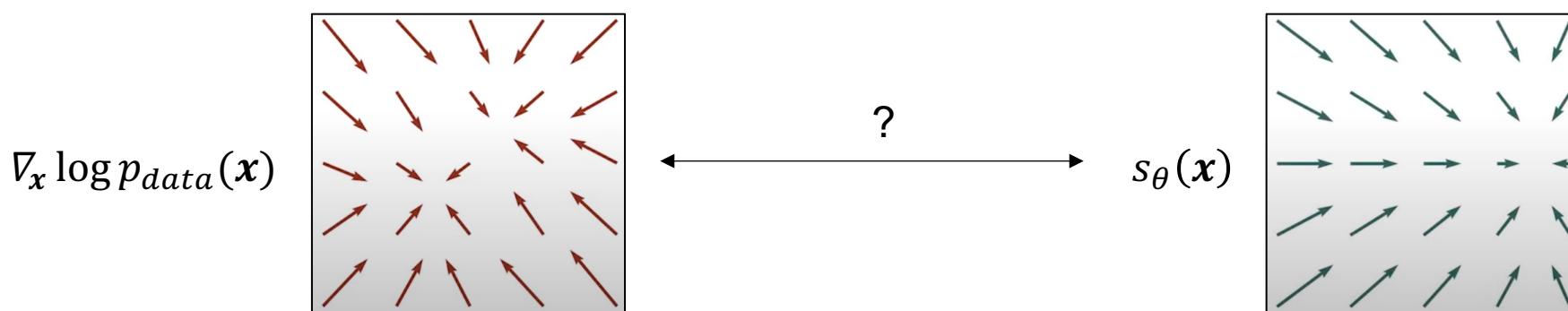
**Given (Data):**  $\{x_1, x_2, \dots, x_N\} \sim p_{data}(x)$

**Goal (Score function):**  $\nabla_x \log p_{data}(x)$

We can't compute this as we don't know  $p_{data}(x)$

**Score model:**  $s_\theta(x): \mathbb{R}^d \rightarrow \mathbb{R}^d \approx \nabla_x \log p_{data}(x)$

**Objective:**  $\mathbb{E}_{p_{data}(x)}[\|\nabla_x \log p_{data}(x) - s_\theta(x)\|_2^2]$  (Fisher divergence to compare the vector fields)



# Score matching

There exists a so-called **score matching objective**, that is similar to the Fisher divergence (up to a constant):

$$\mathbb{E}_{p_{data}(x)} \left[ \frac{1}{2} \|s_\theta(x)\|_2^2 + \text{trace}(\nabla_x s_\theta(x)) \right] \longrightarrow \text{doesn't rely on the ground truth score}$$

As a constant doesn't matter for optimization, this score matching objective defines the same optimum as the Fisher divergence and can **effectively be estimated by the empirical mean over the training data set**:

$$\approx \frac{1}{N} \sum_{i=1}^N \left[ \frac{1}{2} \|s_\theta(x_i)\|_2^2 + \text{trace}(\nabla_x s_\theta(x_i)) \right]$$

Journal of Machine Learning Research 6 (2005) 695–709

Submitted 11/04; Revised 3/05; Published 4/05

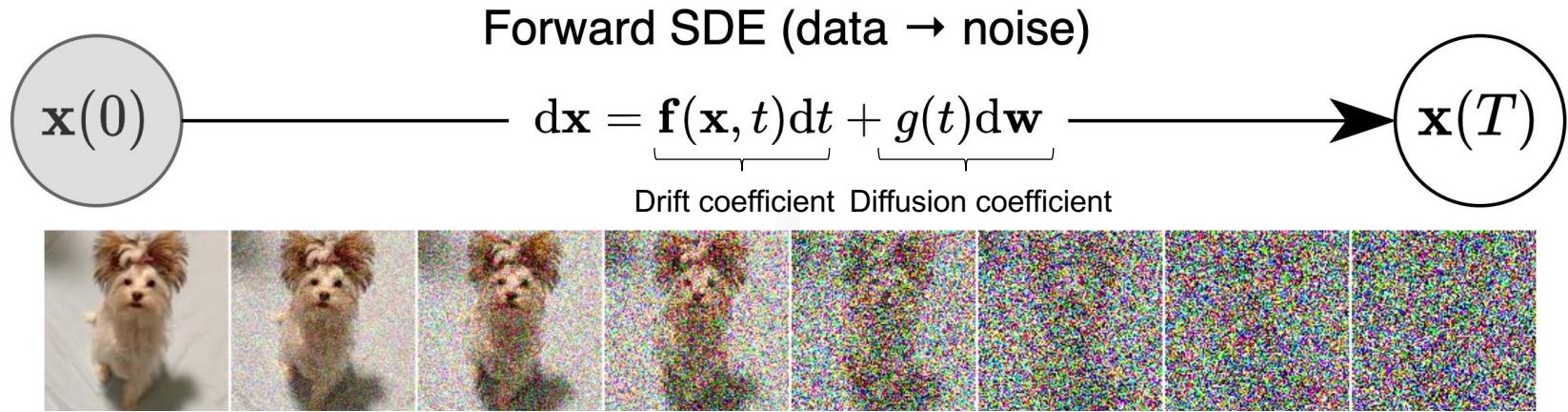
<https://andrewcharlesjones.github.io/journal/21-score-matching.html>

Estimation of Non-Normalized Statistical Models  
by Score Matching

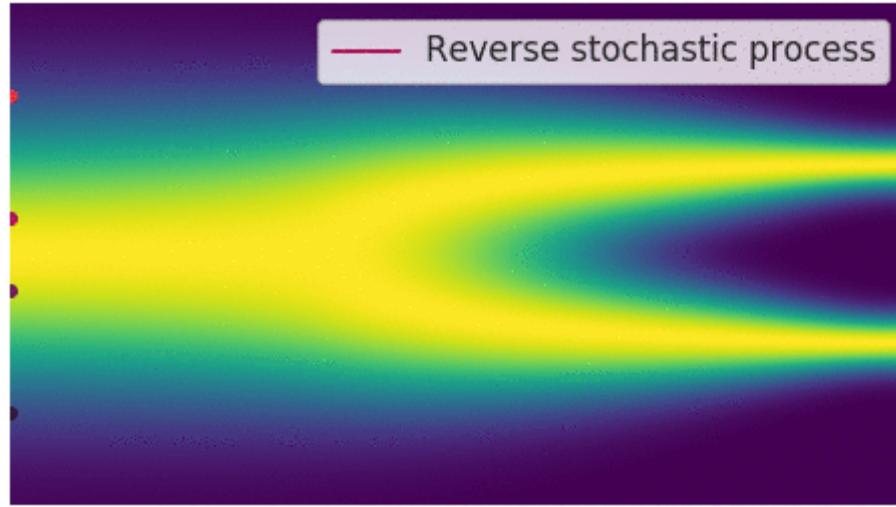
Aapo Hyvärinen  
*Helsinki Institute for Information Technology (HIIT)  
Department of Computer Science  
FIN-00014 University of Helsinki, Finland*

AAPO.HYVARINEN@HELSINKI.FI

# Modelling the Diffusion Process with SDEs



# Modelling the Diffusion Process with SDEs



$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g^2(t) \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + g(t) d\mathbf{w}$$

If we can **model this score function**, we can **solve the reverse SDE** using Euler, Milstein or Runge-Kutta method.

<https://yang-song.net/blog/2021/score/> (More information on this very nice blog post)

# Agenda

- ❑ Introduction to Diffusion Models [~30 min]

- ❑ Physical Intuition & General Concepts
- ❑ Denoising Diffusion Probabilistic Models
- ❑ A Score-Based View on Diffusion Models

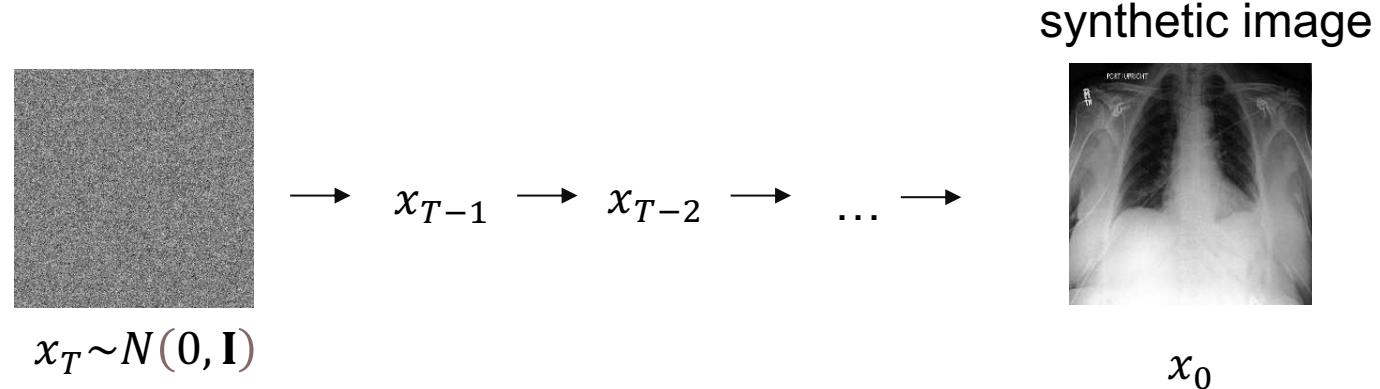
- ❑ Advanced Topics [~30 min]

- ❑ Sampling Strategies
- ❑ Inference-time Conditioning
- ❑ Training-time Conditioning

- ❑ Applications in Medical Imaging [~30 min]

- ❑ Synthesis
- ❑ Inpainting
- ❑ Segmentation
- ❑ Anomaly Detection
- ❑ Reconstruction
- ❑ Registration

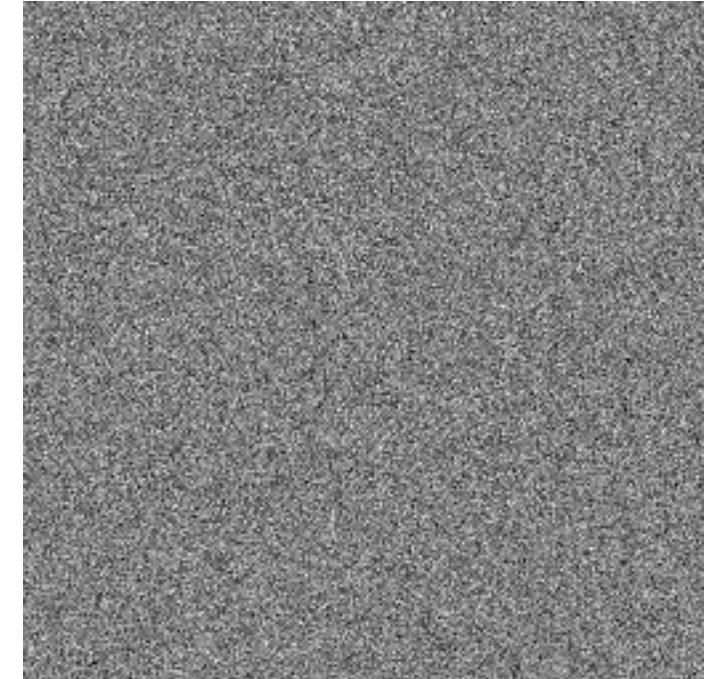
# Fake Image Generation



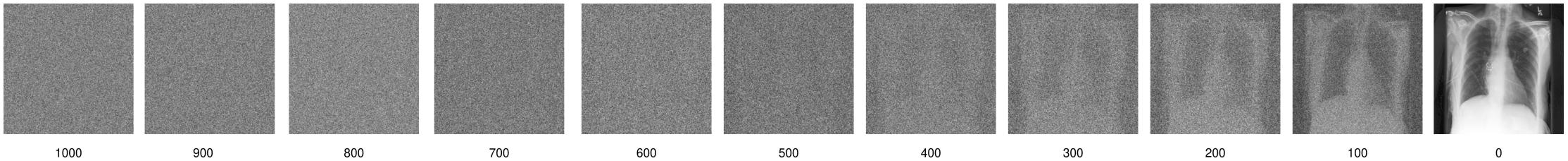
**U-Net**

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right) + \sigma_t \epsilon, \quad \text{with } \epsilon \sim \mathcal{N}(0, \mathbf{I}).$$

Random component



# Schedulers: How to Accelerate Sampling?



Published as a conference paper at ICLR 2021

## DENOISING DIFFUSION IMPLICIT MODELS

Jiaming Song, Chenlin Meng & Stefano Ermon  
Stanford University  
`{tsong, chenlin, ermon}@cs.stanford.edu`

### ABSTRACT

Denoising diffusion probabilistic models (DDPMs) have achieved high quality image generation without adversarial training, yet they require simulating a Markov chain for many steps in order to produce a sample. To accelerate sampling, we present denoising diffusion implicit models (DDIMs), a more efficient class of iterative implicit probabilistic models with the same training procedure as DDPMs. In DDPMs, the generative process is defined as the reverse of a particular Markovian diffusion process. We generalize DDPMs via a class of non-Markovian diffusion processes that lead to the same training objective. These non-Markovian

"Denoising diffusion probabilistic models (DDPMs) have achieved high quality image generation, yet they require simulating a Markov chain for many steps in order to produce a sample."



We need to make the generation process faster.

# From DDPMs to DDIMs

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}} \left( \underbrace{\frac{\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_\theta^{(t)}(\mathbf{x}_t)}{\sqrt{\alpha_t}}}_{\text{“predicted } \mathbf{x}_0\text{”}} \right) + \underbrace{\sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta^{(t)}(\mathbf{x}_t)}_{\text{“direction pointing to } \mathbf{x}_t\text{”}} + \underbrace{\sigma_t \epsilon_t}_{\text{random noise}}$$

DDPM sampling scheme

$$\sigma_t = \sqrt{(1 - \alpha_{t-1})/(1 - \alpha_t)} \sqrt{1 - \alpha_t/\alpha_{t-1}}$$

DDIM sampling scheme

$$\sigma_t = 0$$

We remove the random component

The training process stays the same.

# An Excursion into ODEs

- The connection to ordinary differential equations (ODEs) can be seen when we rewrite the DDIM denoising step as

$$\frac{x_{t-1}}{\sqrt{\bar{\alpha}_{t-1}}} = \frac{x_t}{\sqrt{\bar{\alpha}_t}} + \left( \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{\bar{\alpha}_{t-1}}} - \sqrt{\frac{1 - \bar{\alpha}_t}{\bar{\alpha}_t}} \right) \epsilon_\theta(x_t, t).$$

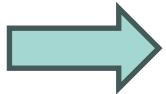
prediction

previous  
value

step size

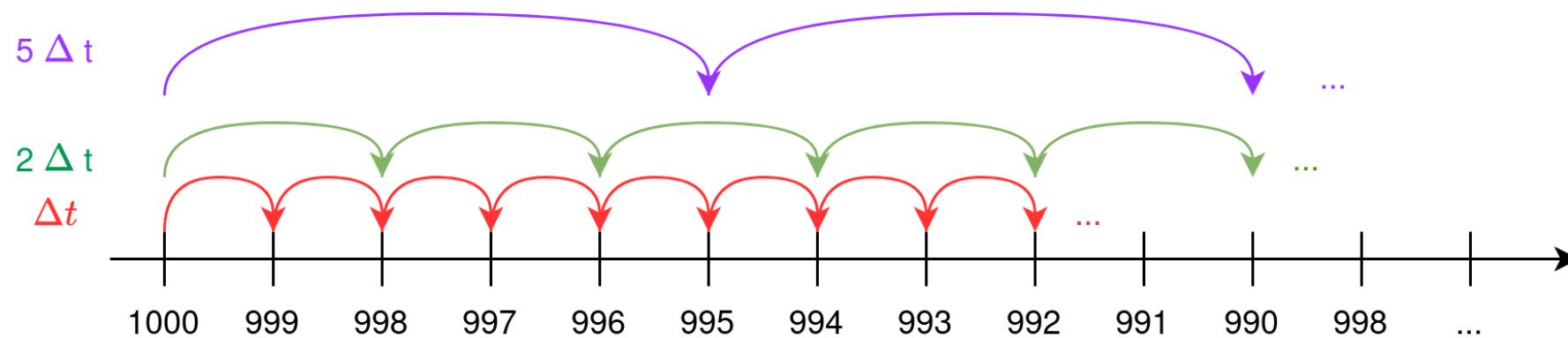
derivative

- This can be interpreted as the Euler approximation of an ODE.
- DDIM is a **probability flow** ODE from a SDE [1].
- We can speed up the generation process by choosing a larger step size.

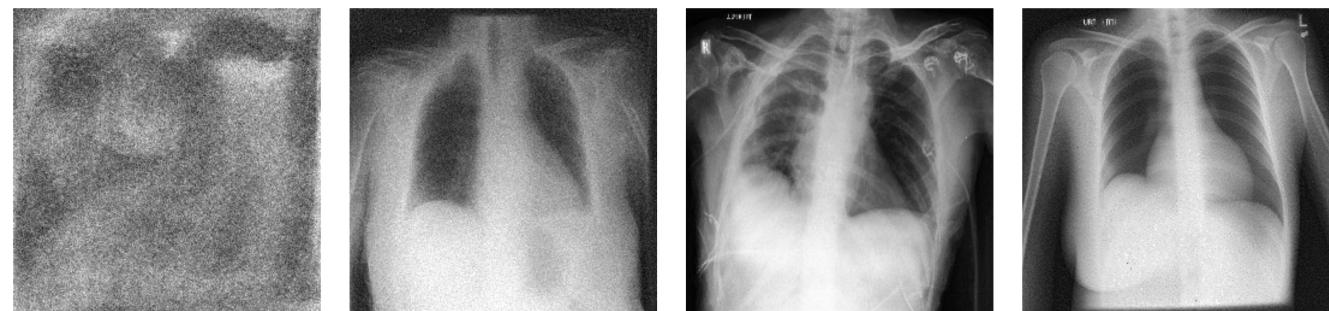


Faster, but less accurate

# DDIM Accelerated Sampling



- By skipping  $k$  steps, we have a step size of  $k\Delta t$ .
- Sampling is  $k$  times faster.
- We trade image quality for speed.



Total amount of steps

# Various Schedulers...

## Elucidating the Design Space of Diffusion-Based Generative Models

### PSEUDO NUMERICAL METHODS FOR DIFFUSION MODELS ON MANIFOLDS

Luping Liu, Yi Ren  
Zhejiang University  
{luping.liu, yiren.zj@zju.edu.cn}

Published as a conference paper at ICLR 2022

### PROGRESSIVE DISTILLATION FOR FAST SAMPLING OF DIFFUSION MODELS

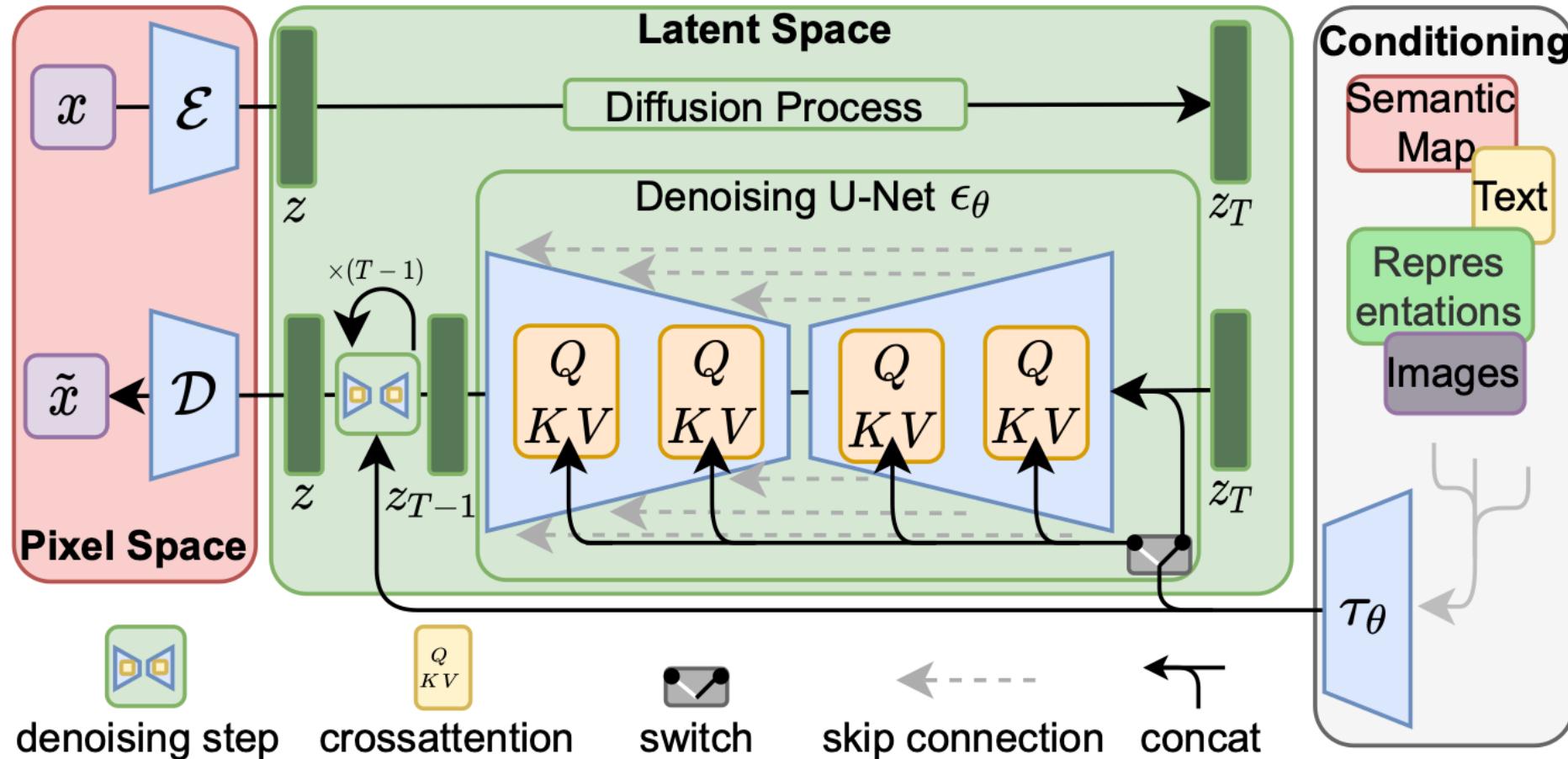
Tim Salimans & Jonathan Ho  
Google Research, Brain team  
{salimans,jonathanho}@google.com

#### ABSTRACT

Diffusion models have recently shown great promise for generative modeling, outperforming GANs on perceptual quality and autoregressive models at density estimation. A remaining downside is their slow sampling time: generating high quality samples takes many hundreds or thousands of model evaluations. Here we make two contributions to help eliminate this downside: First, we present new parameterizations of diffusion models that provide increased stability when using few sampling steps. Second, we present a method to distill a trained deterministic diffusion sampler, using many steps, into a new diffusion model that takes half as many sampling steps. We then keep progressively applying this distillation process.

- Choosing a different solver for the given ODE can improve speed and image quality.
- Other numerical approaches such as **Heun's Method** or **Runge Kutta** solvers can be explored.
- Knowledge distillation techniques can be used for fast sampling.

# Diffusion Models in the Latent Space



# Agenda

- ❑ Introduction to Diffusion Models [~30 min]

- ❑ Physical Intuition & General Concepts
- ❑ Denoising Diffusion Probabilistic Models
- ❑ A Score-Based View on Diffusion Models

- ❑ Advanced Topics [~30 min]

- ❑ Sampling Strategies
- ❑ Inference-time Conditioning
- ❑ Training-time Conditioning

- ❑ Applications in Medical Imaging [~30 min]

- ❑ Synthesis
- ❑ Inpainting
- ❑ Segmentation
- ❑ Anomaly Detection
- ❑ Reconstruction
- ❑ Registration

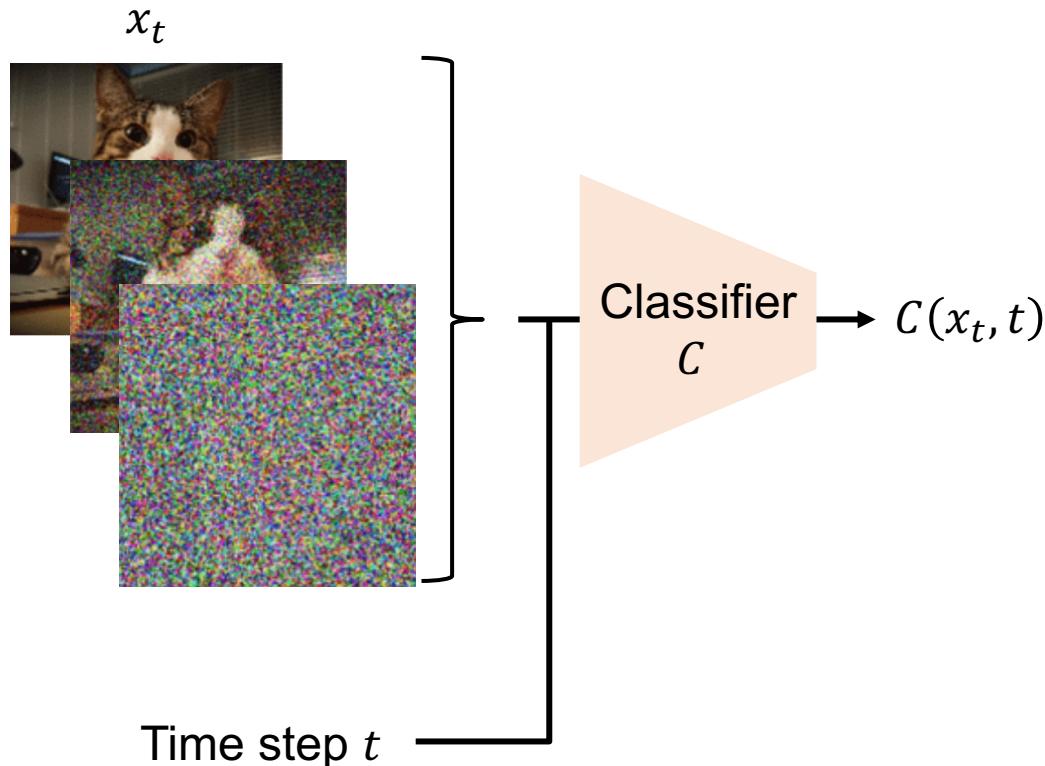
"mouse"

Diffusion Model

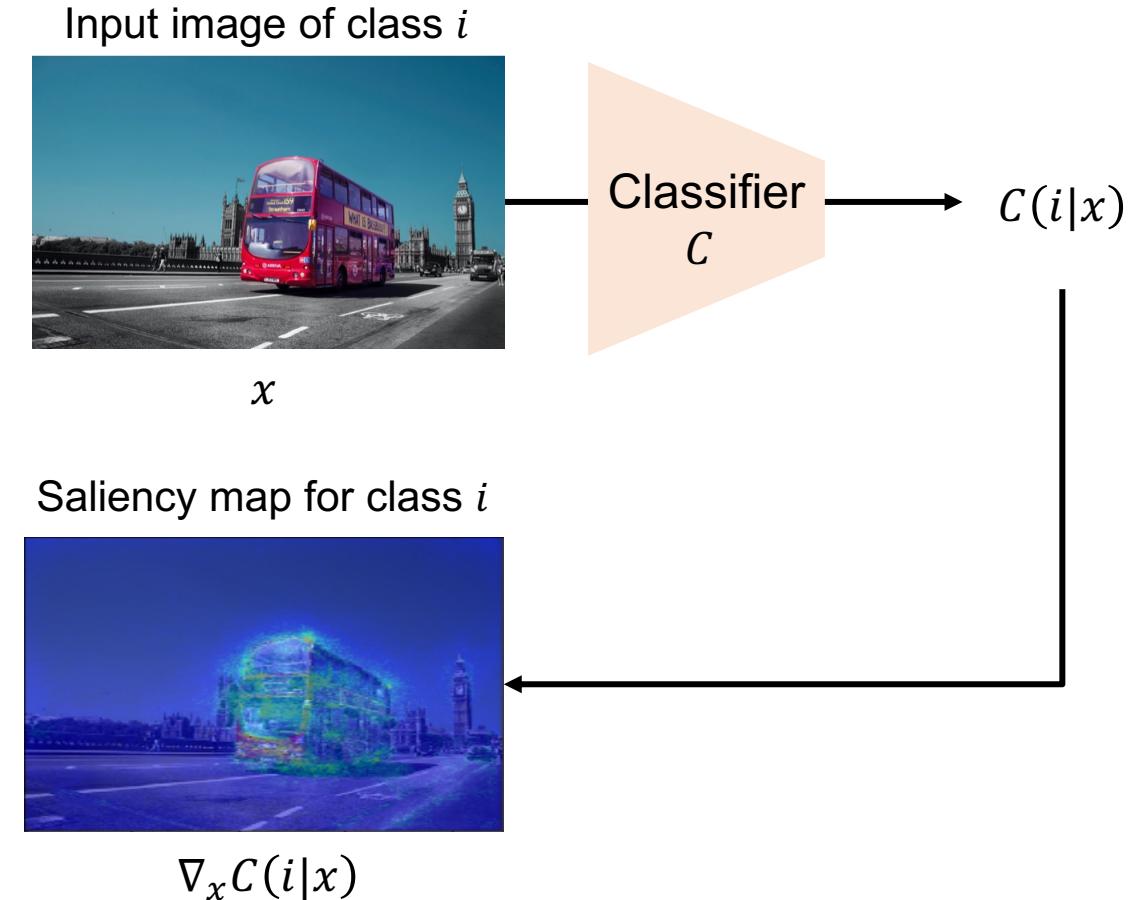


# Example: Classifier Guidance

We want a class-conditional diffusion model.

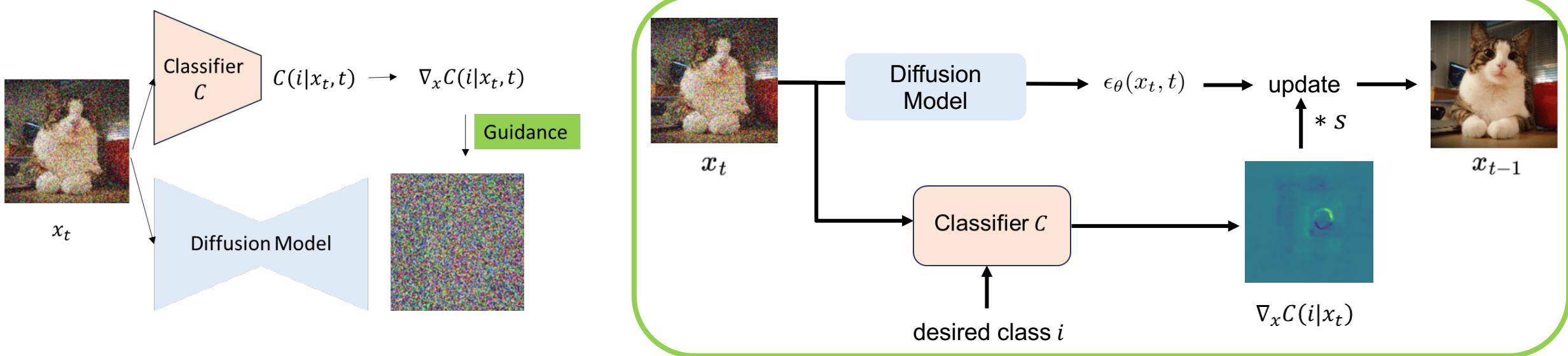


We consider the gradient with respect to the input pixels.



# Classifier Guidance

We use the gradient to guide the generation process towards a desired class.



Gradient guidance is not restricted to classification models. Other models (e.g., regression, segmentation, ...) work just in the same way.

# Classifier Guidance



goldfish

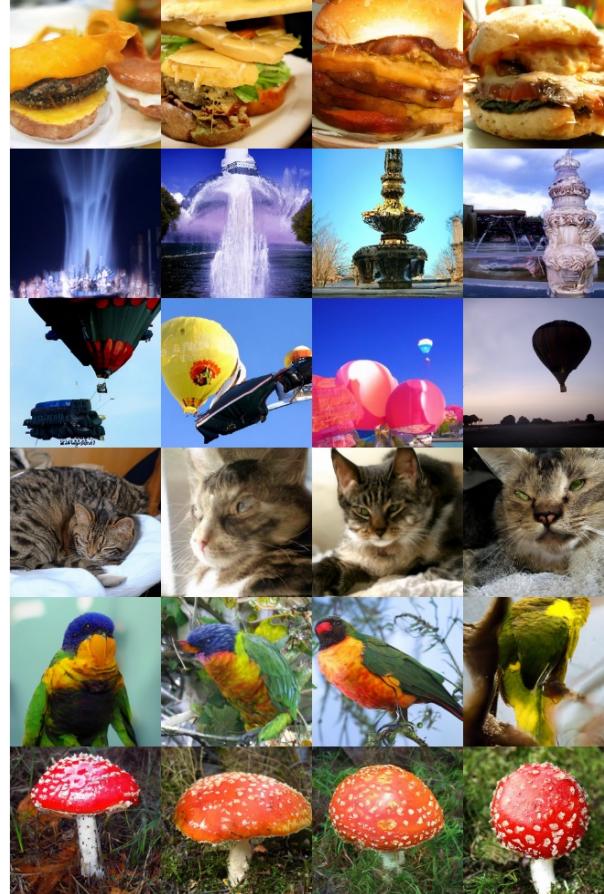
arctic fox

butterfly

African elephant

flamingo

tennis ball



cheeseburger

fountain

balloon

tabby cat

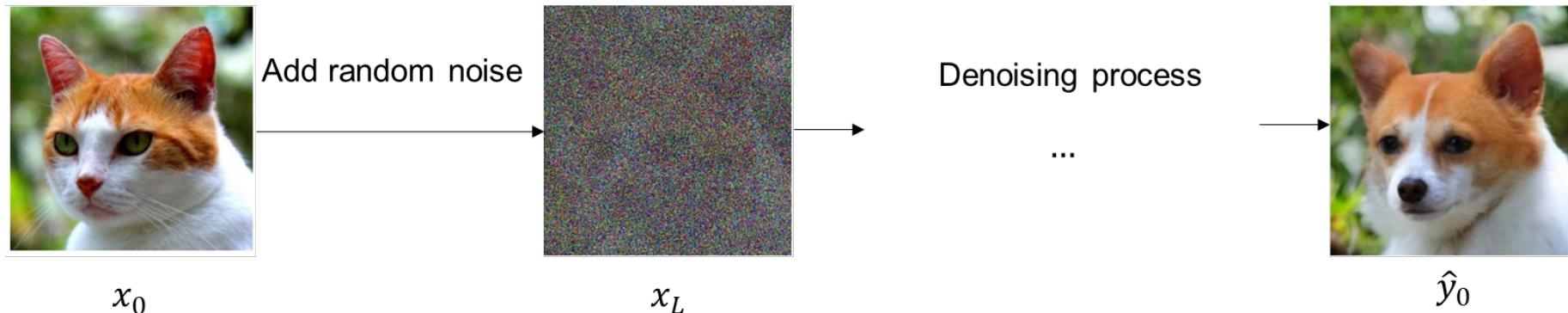
lorikeet

agaric

Check out this very nice tutorial: <https://sander.ai/2022/05/26/guidance.html>

# Image-to-image translation

We might want to translate an image to another...



- We add  $L$  steps of noise to an input image  $x_0$ .
- The smaller  $L$ , the less the image can be changed.
- The higher  $L$ , the more information is destroyed.



We need to find a way to keep the information of  $x_0$ .



We consider Denoising Diffusion Implicit Models (DDIMs).

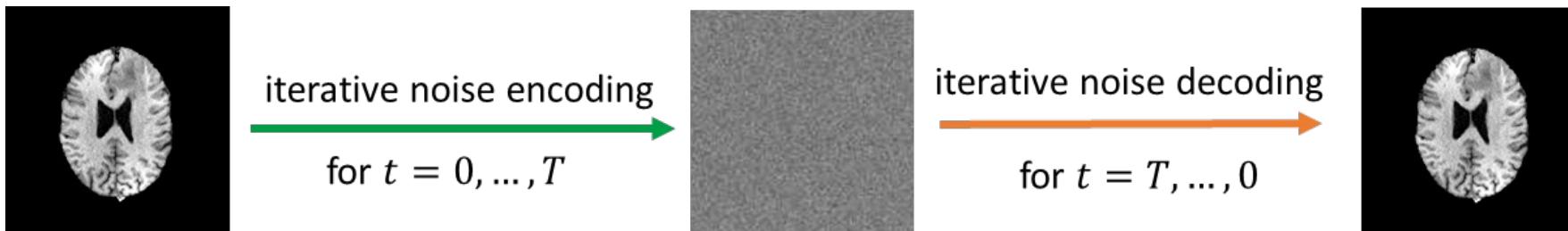
# DDIM Inversion

- Under the DDIM sampling scheme, we remove the random component.
- The connection to ordinary differential equations (ODEs) can be seen when we rewrite the denoising step as

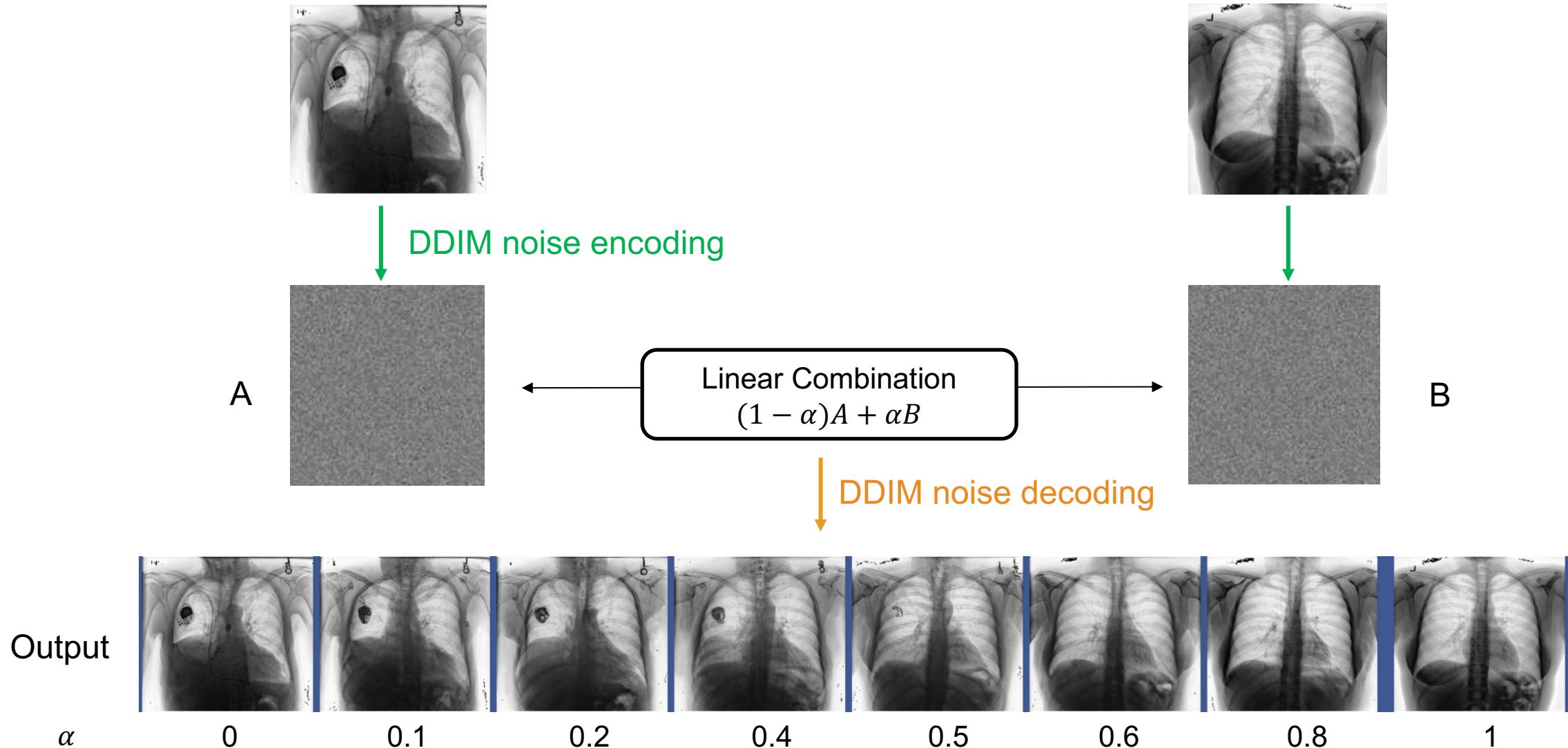
$$\frac{x_{t-1}}{\sqrt{\bar{\alpha}_{t-1}}} = \frac{x_t}{\sqrt{\bar{\alpha}_t}} + \left( \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{\bar{\alpha}_{t-1}}} - \sqrt{\frac{1 - \bar{\alpha}_t}{\bar{\alpha}_t}} \right) \epsilon_\theta(x_t, t). \quad \text{noise decoding}$$

- This can be interpreted as the Euler approximation of an ODE.
- Given infinitely small steps  $t$ , the reversed ODE can then be solved with

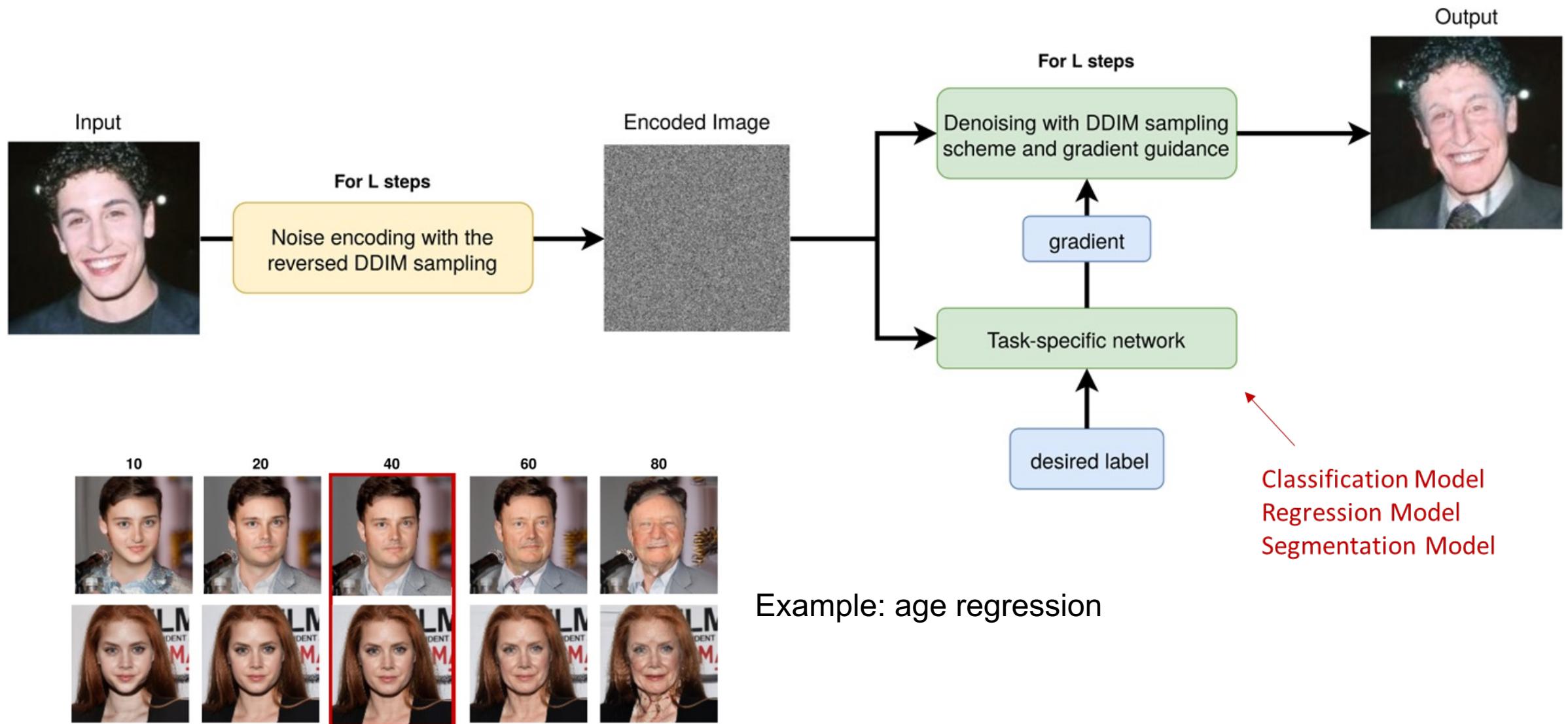
$$\frac{x_{t+1}}{\sqrt{\bar{\alpha}_{t+1}}} = \frac{x_t}{\sqrt{\bar{\alpha}_t}} + \left( \sqrt{\frac{1 - \bar{\alpha}_{t+1}}{\bar{\alpha}_{t+1}}} - \sqrt{\frac{1 - \bar{\alpha}_t}{\bar{\alpha}_t}} \right) \epsilon_\theta(x_t, t). \quad \text{noise encoding}$$



# Image Interpolation



# DDIM Inversion & Gradient Guidance



# Agenda

- ❑ Introduction to Diffusion Models [~30 min]

- ❑ Physical Intuition & General Concepts
- ❑ Denoising Diffusion Probabilistic Models
- ❑ A Score-Based View on Diffusion Models

- ❑ Advanced Topics [~30 min]

- ❑ Sampling Strategies
- ❑ Inference-time Conditioning
- ❑ Training-time Conditioning

- ❑ Applications in Medical Imaging [~30 min]

- ❑ Synthesis
- ❑ Inpainting
- ❑ Segmentation
- ❑ Anomaly Detection
- ❑ Reconstruction
- ❑ Registration

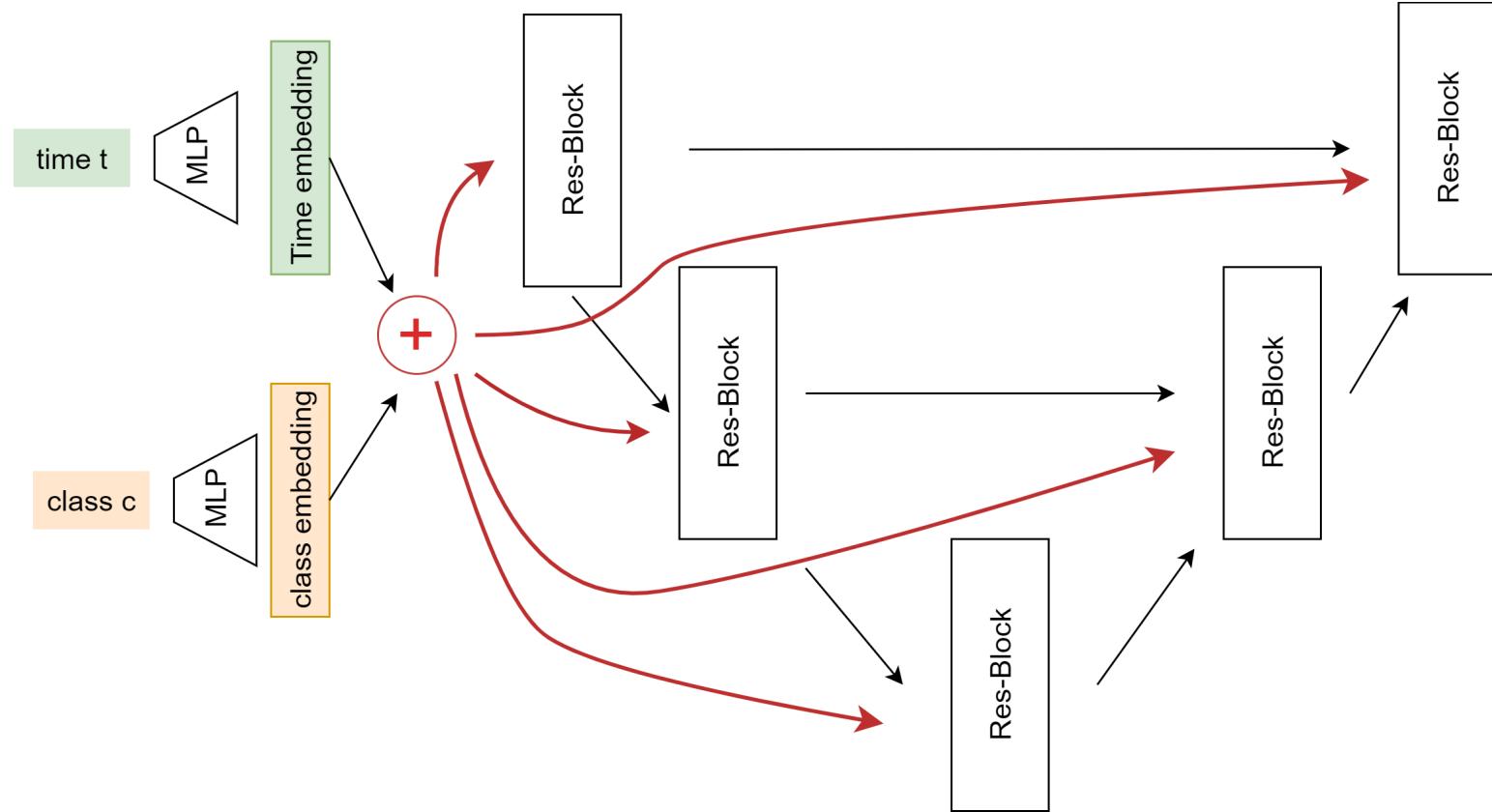
"mouse"

Diffusion Model

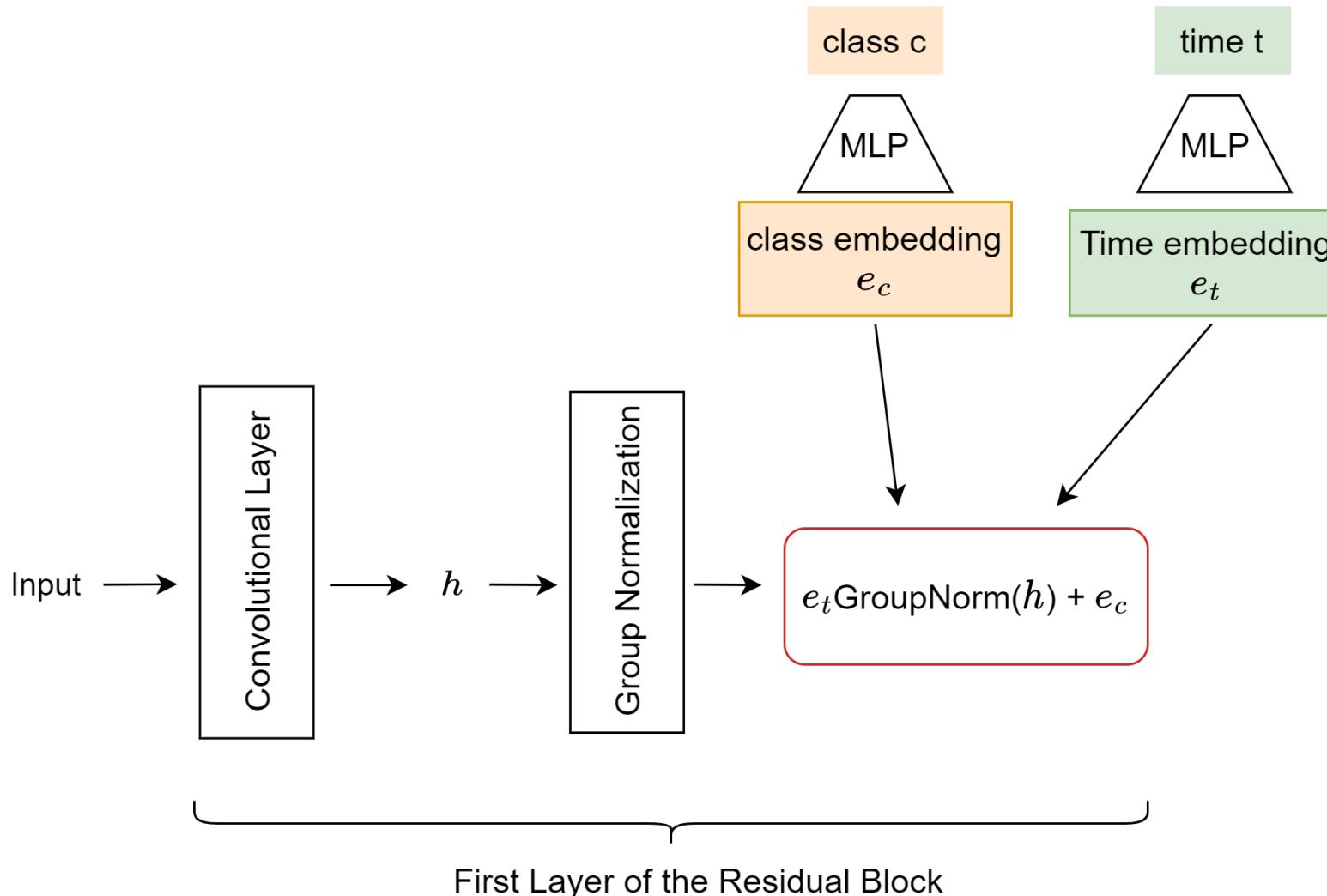


# Scalar Conditioning via Spatial Addition

- We train a class-conditional diffusion model by including a class label  $c$ .
- We compute a class embedding, and pass it to the residual blocks by **spatial addition**.

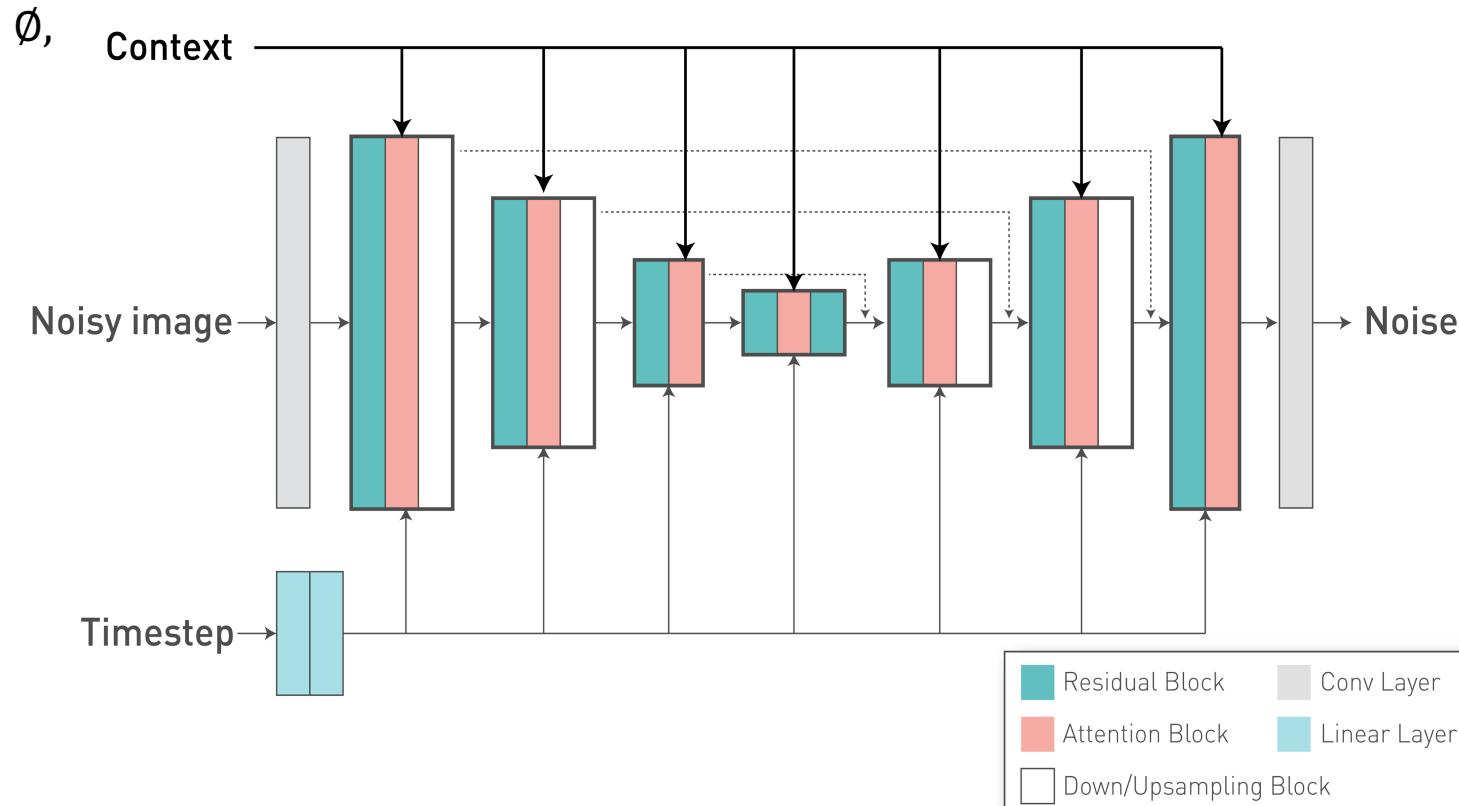


# Scalar Conditioning via Adaptive Group Normalization



- Similar to StyleGAN, we add time and class information using a group normalization layer.
- This happens in all residual blocks of the U-Net.

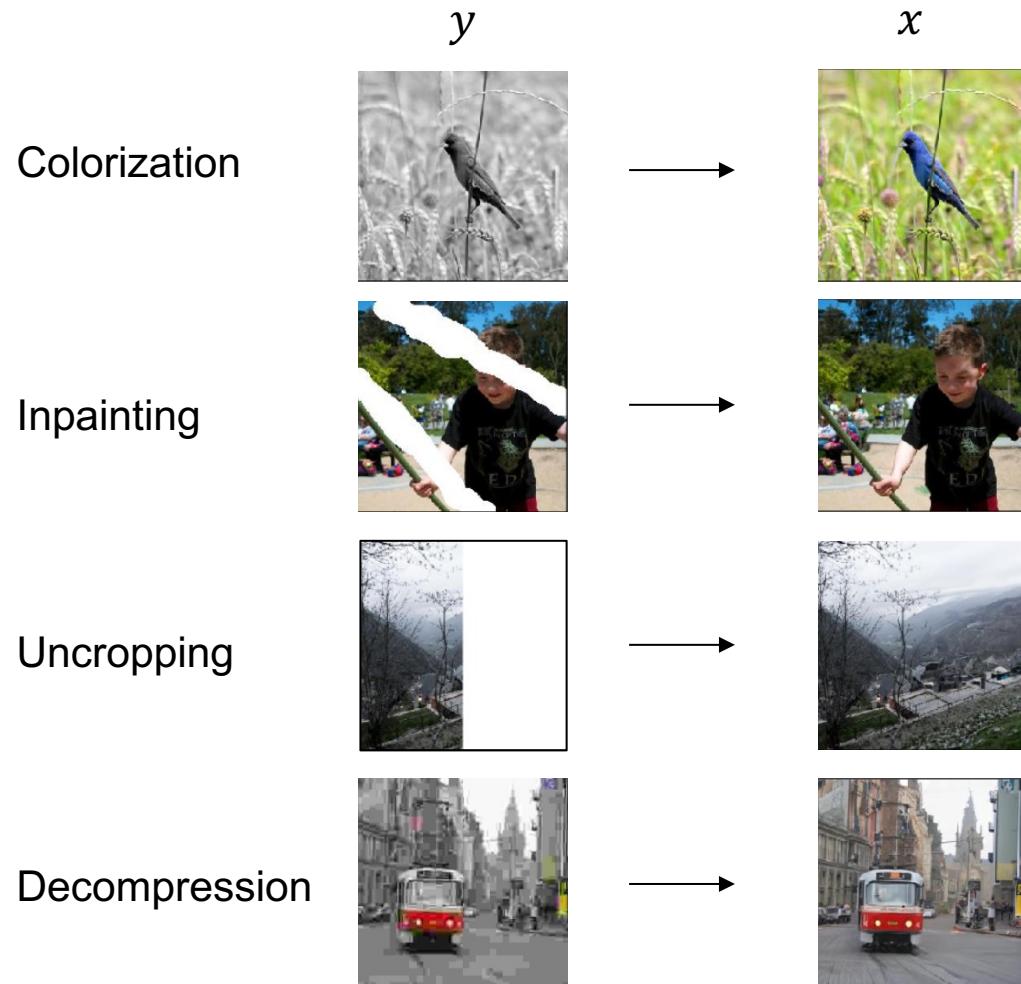
# Classifier-free Guidance



$$\widetilde{\epsilon}_{\theta}(x_t|y) = \epsilon_{\theta}(x_t|\emptyset) + w[\epsilon_{\theta}(x_t|y) - \epsilon_{\theta}(x_t|\emptyset)]$$

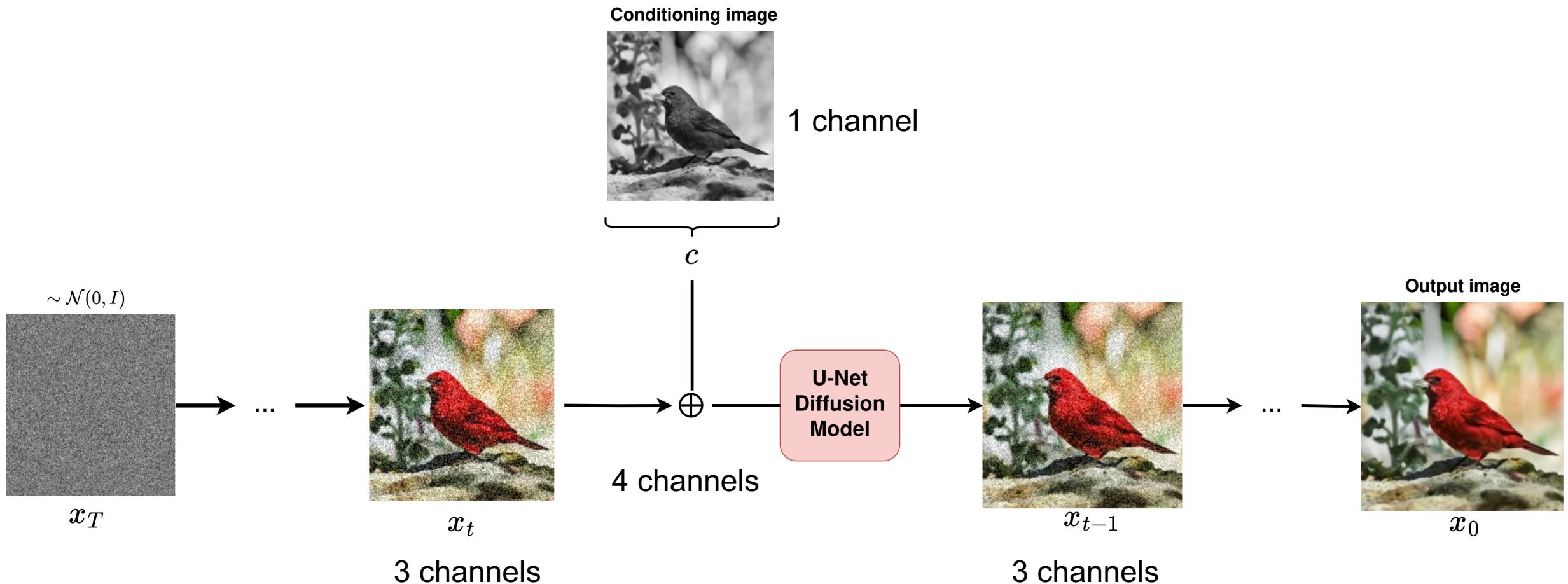
Check out this very nice tutorial: <https://sander.ai/2022/05/26/guidance.html>

# Image Conditioning through Concatenation

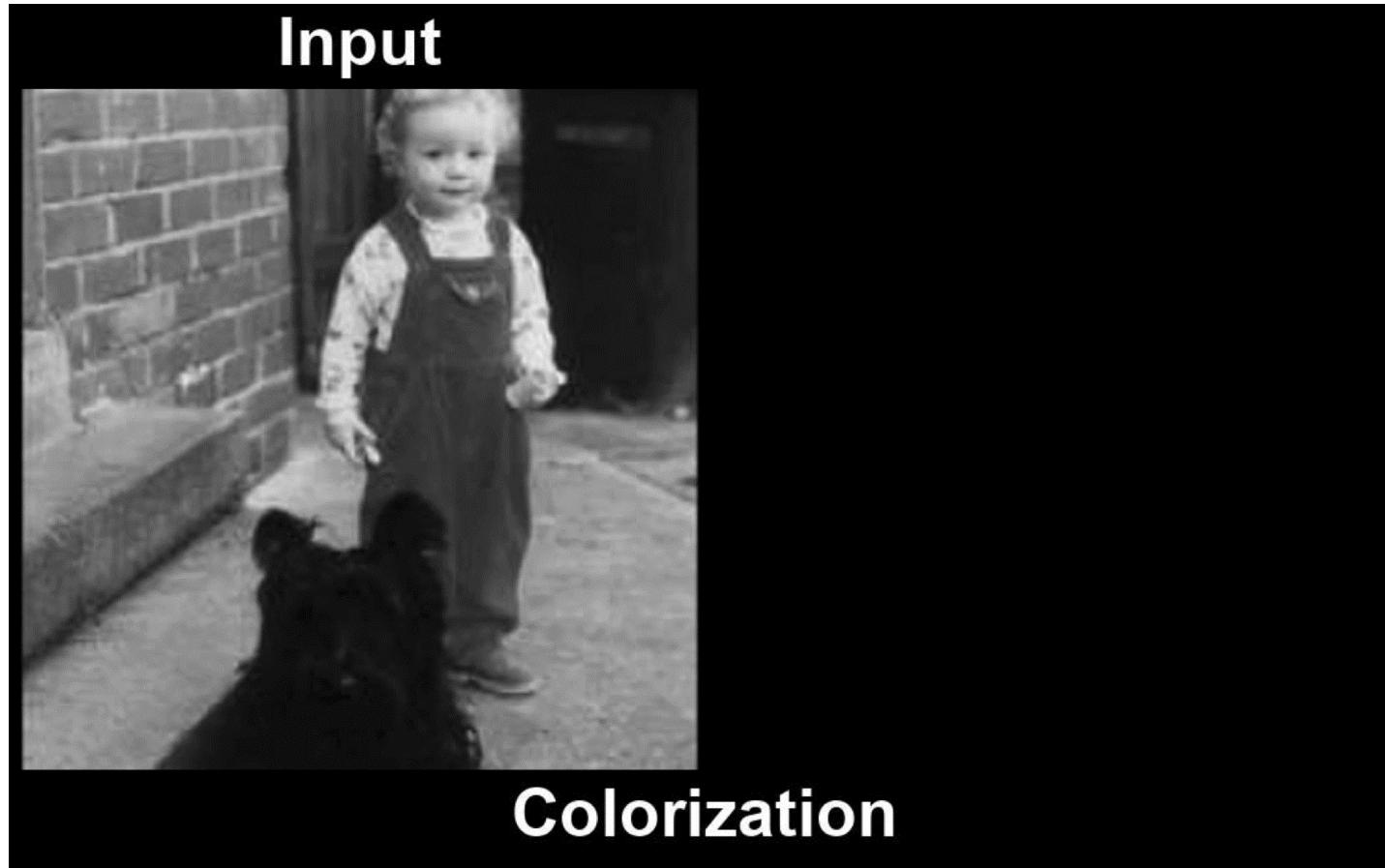


- For image generation of a fake image  $x$ , we can use a conditioning image  $y$ .
- This requires **paired** training.
- During training and sampling, we add information of the conditioning image  $x$  through **channel-wise concatenation**.

# Image Conditioning through Concatenation



# Palette: Image-to-Image Diffusion Models



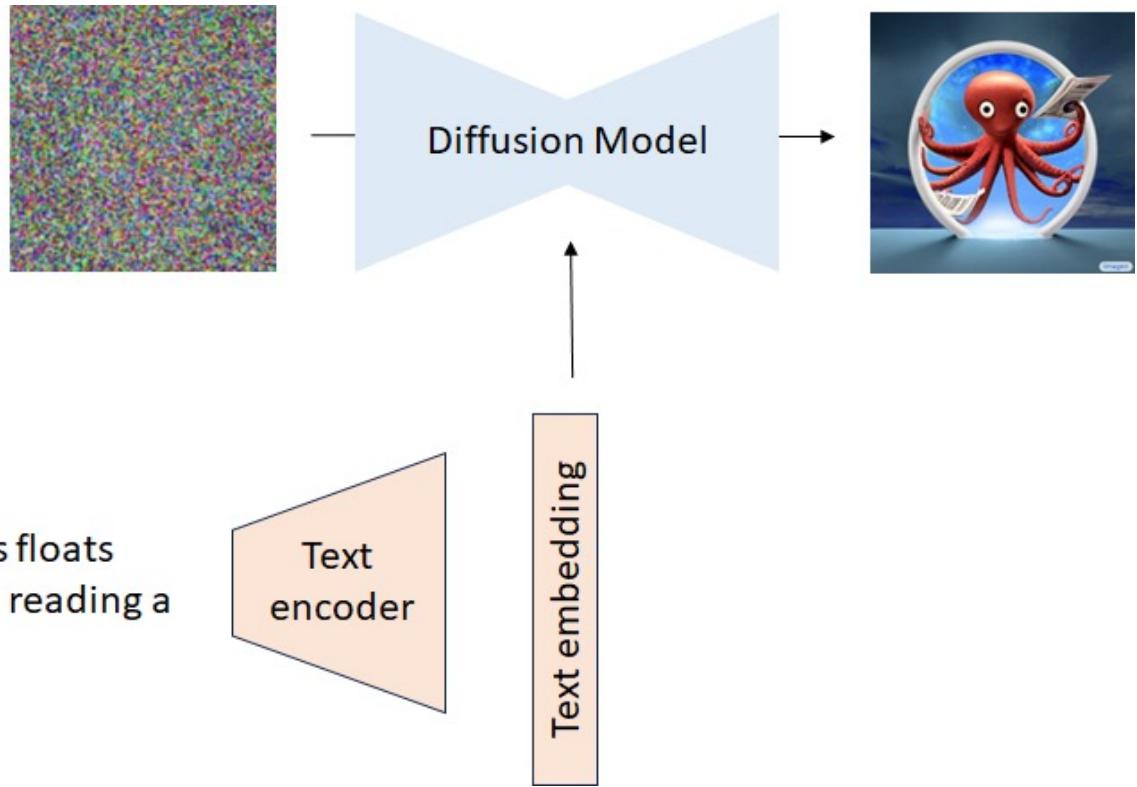
# Text Conditioning



"A small cactus wearing a straw hat and neon sunglasses in the Sahara desert."

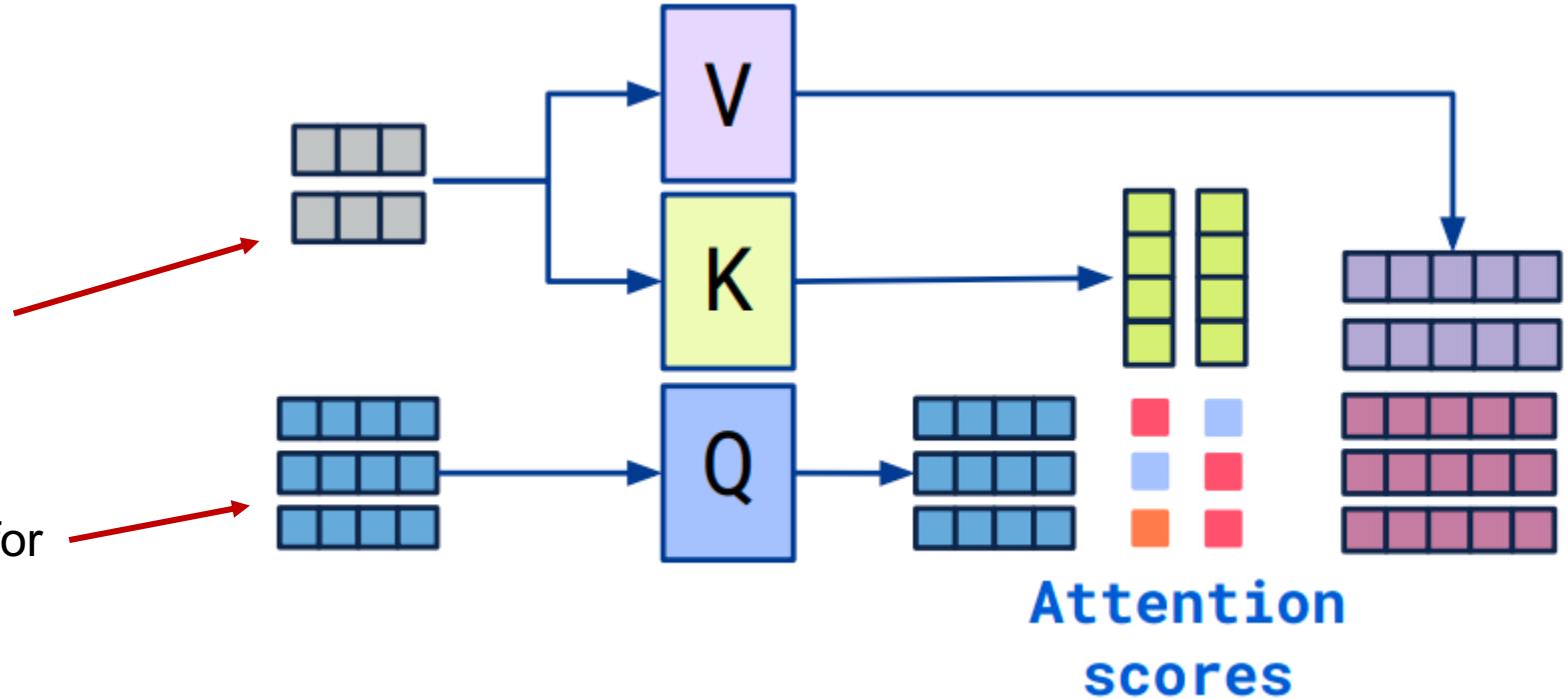
CLIP  
Dall-E  
Stable Diffusion  
Imagen  
...

An alien octopus floats through a portal reading a newspaper.

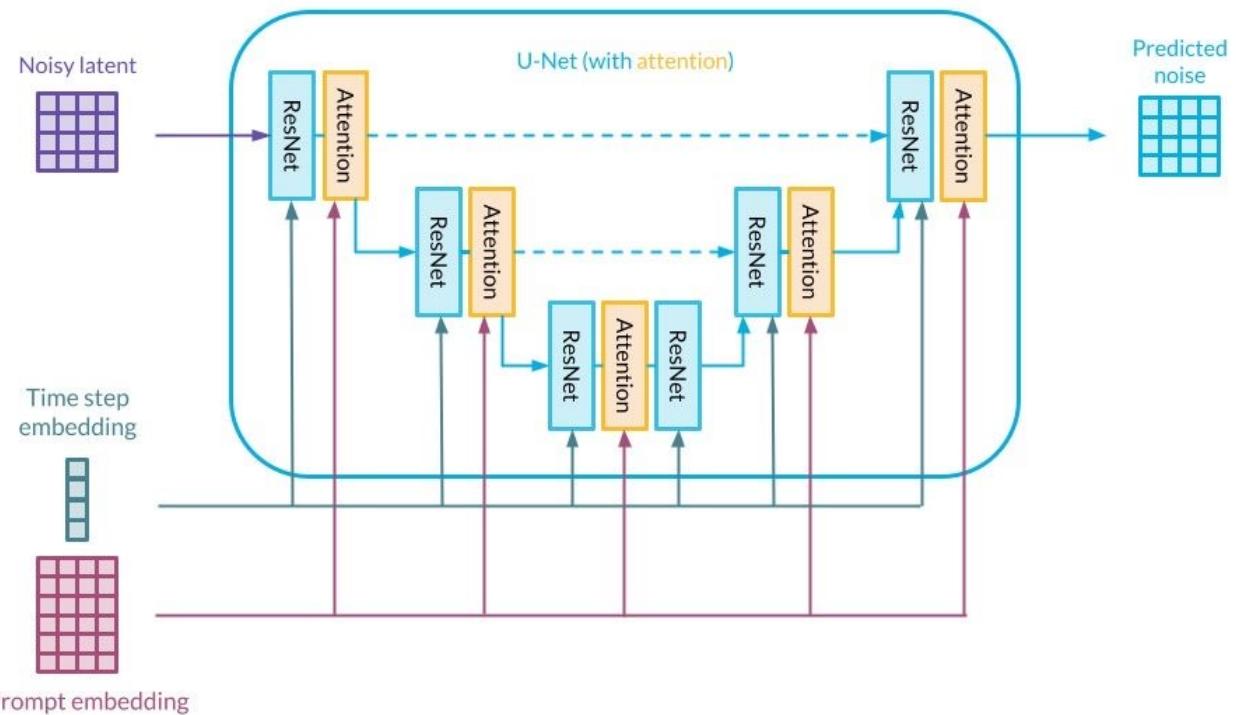
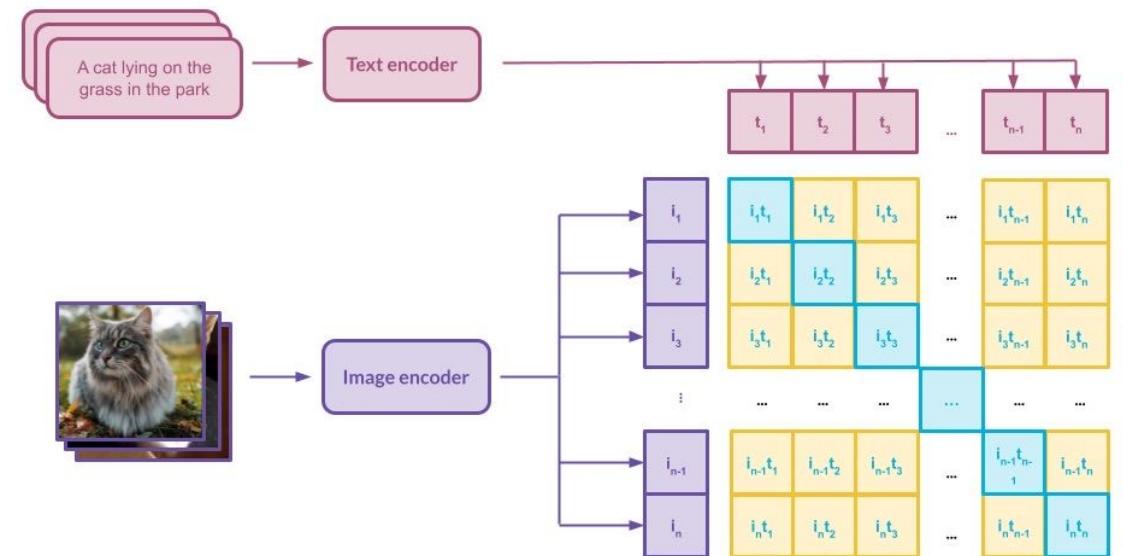


# Architecture - Conditioning

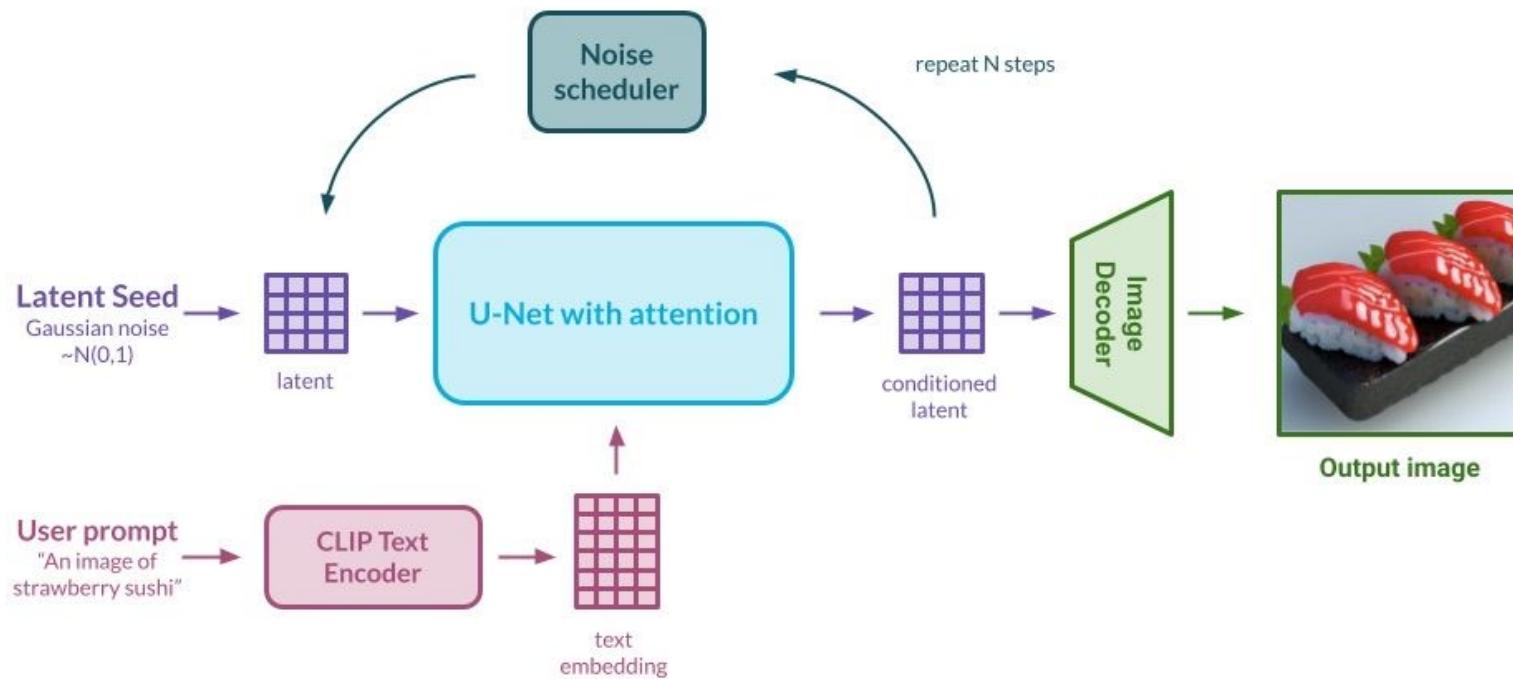
- Transformers **Cross-Attention**
- We use the text embedding to generate the key / value pair.
- We use the image embedding for the query.



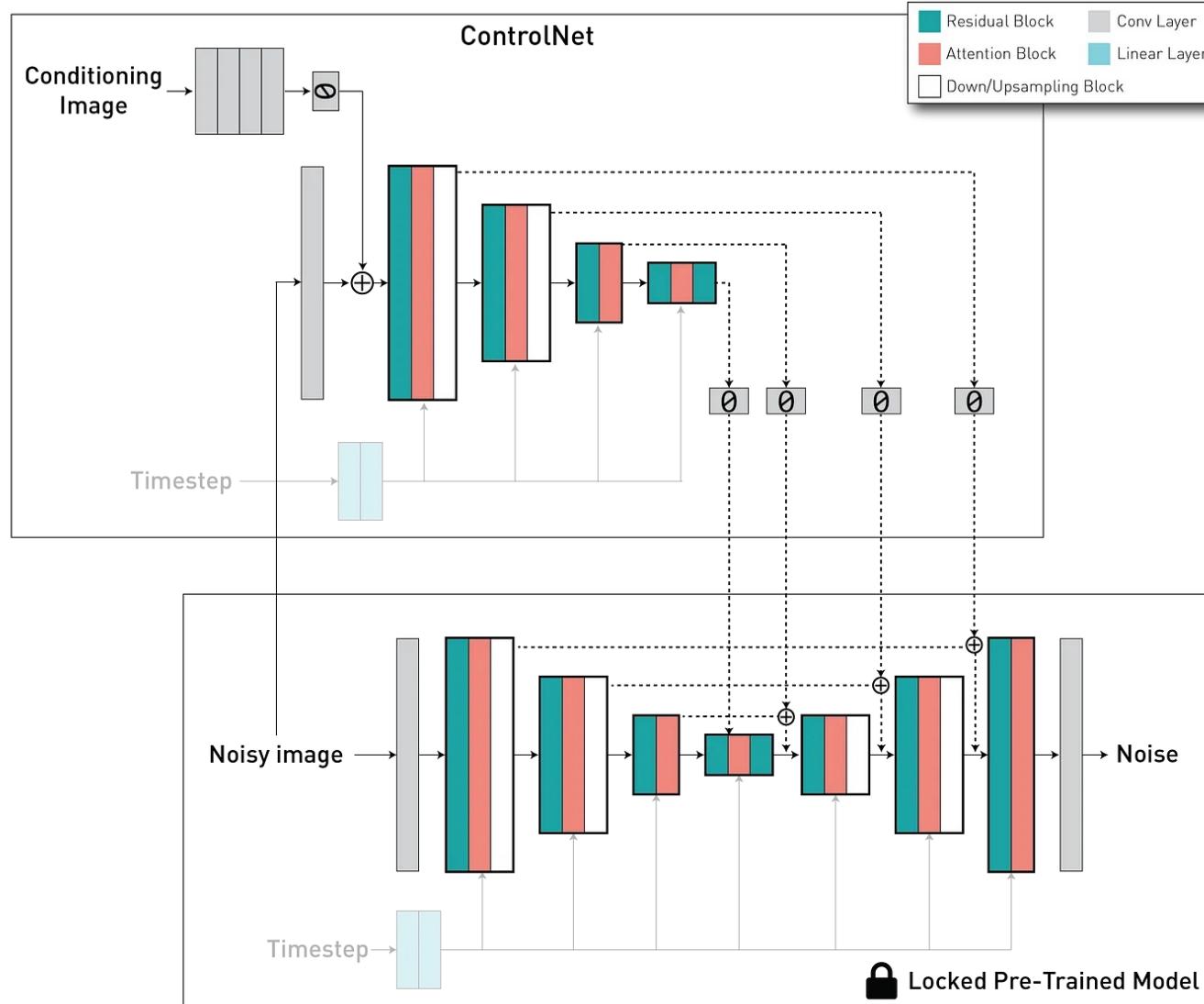
# Text Conditioning



# Text Conditioning

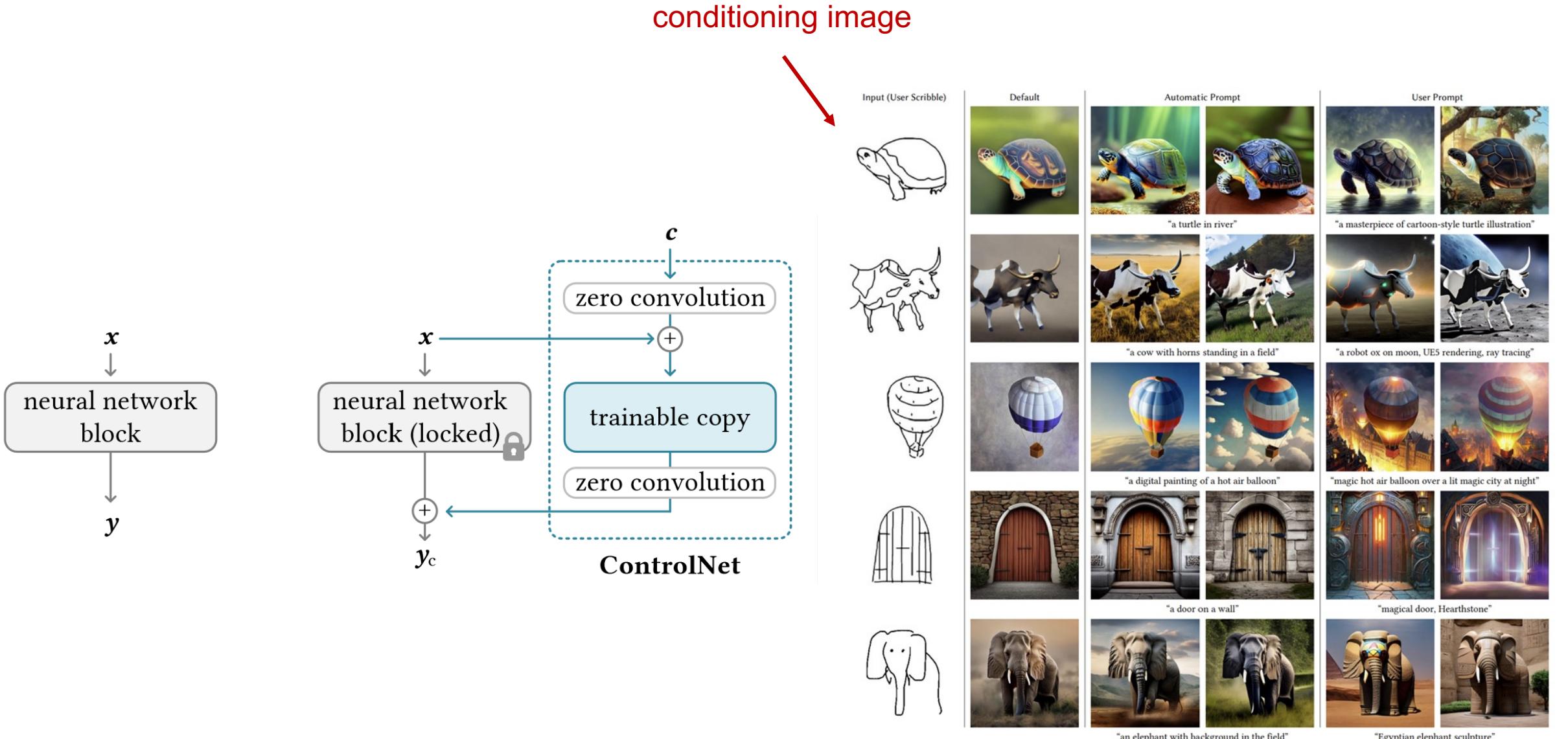


# ControlNet

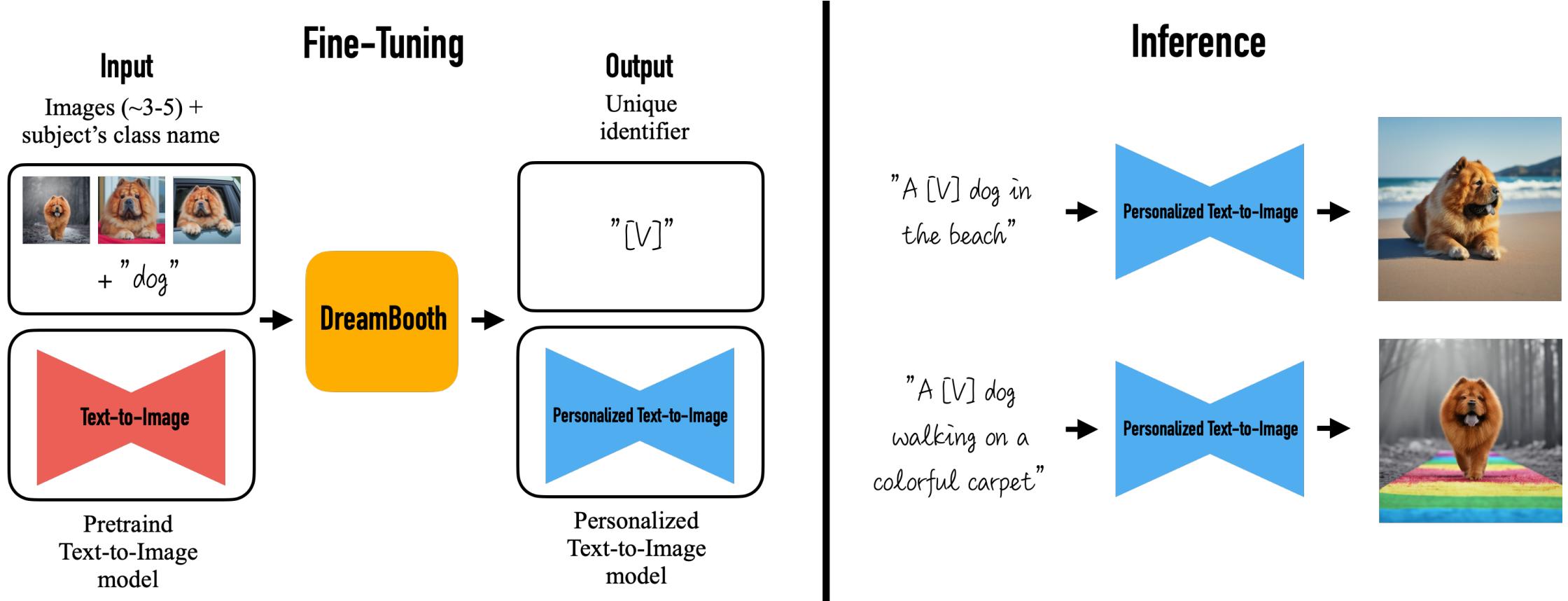


- We pretrain a diffusion model with text prompts.
- We freeze this model.
- We fine-tune a copy conditioned on  $c$ .
- We pass information to the frozen model through skip connections.

# ControlNet



# DreamBooth



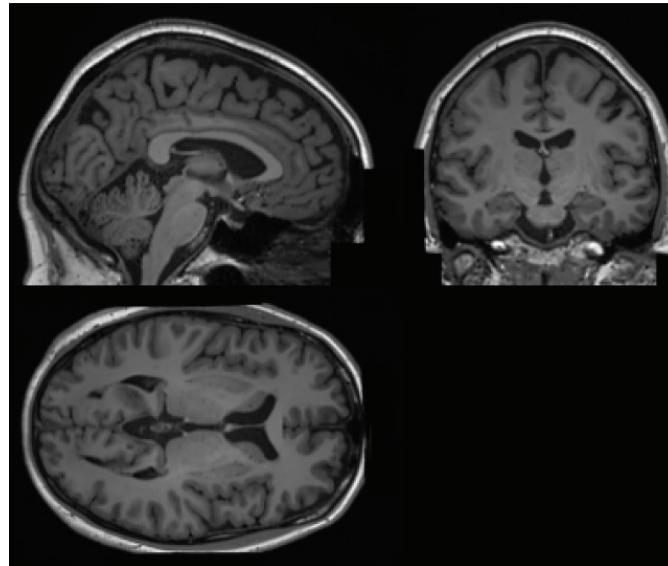
# Agenda

- ❑ Introduction to Diffusion Models [~30 min]
  - ❑ Physical Intuition & General Concepts
  - ❑ Denoising Diffusion Probabilistic Models
  - ❑ A Score-Based View on Diffusion Models
- ❑ Advanced Topics [~30 min]
  - ❑ Sampling Strategies
  - ❑ Inference-time Conditioning
  - ❑ Training-time Conditioning
- ❑ Applications in Medical Imaging [~30 min]
  - ❑ Synthesis
  - ❑ Inpainting
  - ❑ Segmentation
  - ❑ Anomaly Detection
  - ❑ Reconstruction
  - ❑ Registration

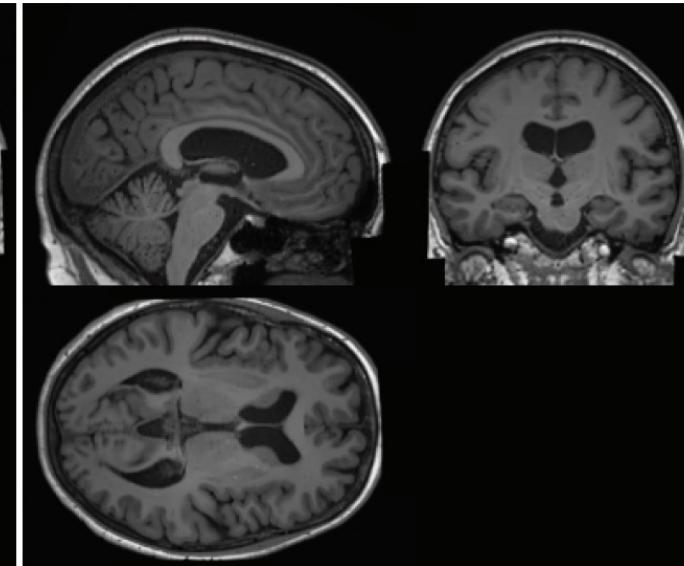
# Image Synthesis

Examples from the Community

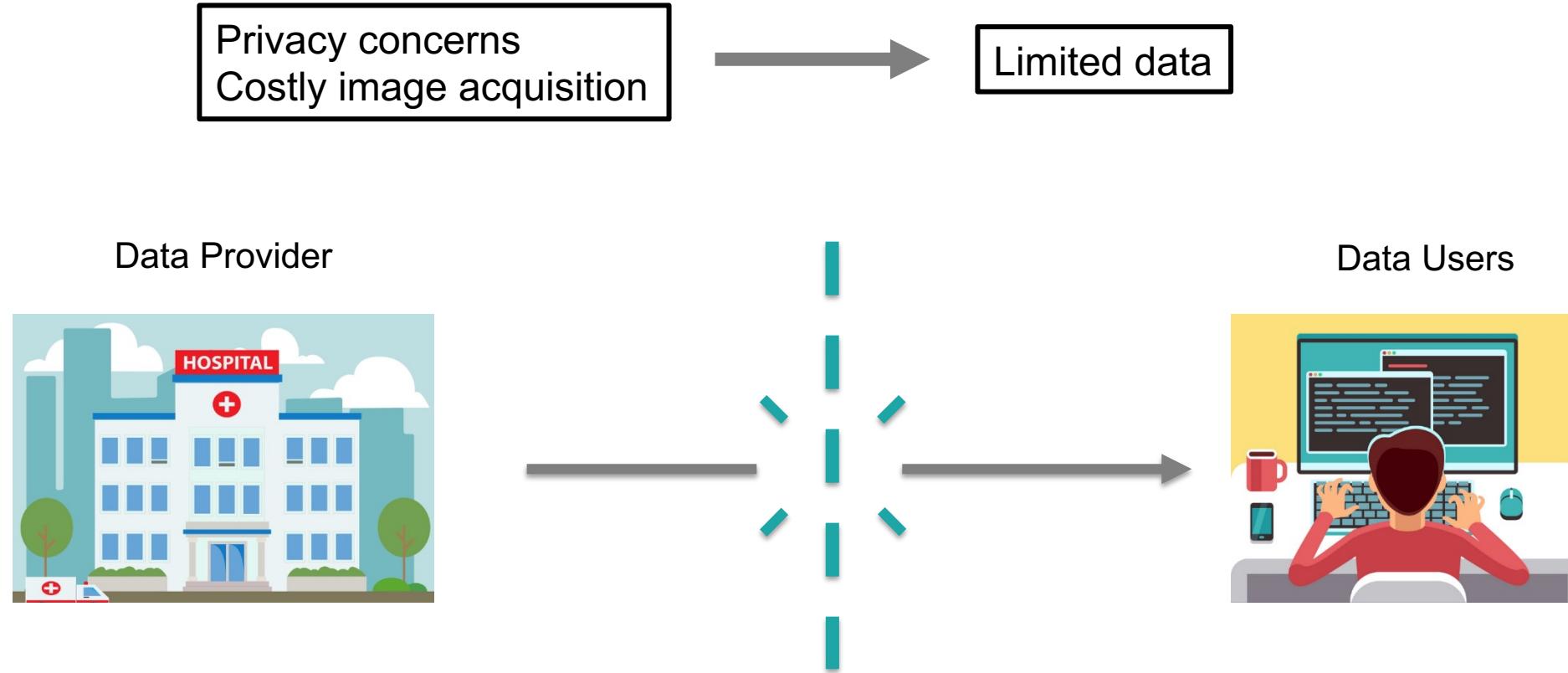
Real



Synthetic



# Why? Medical Images are Rare!



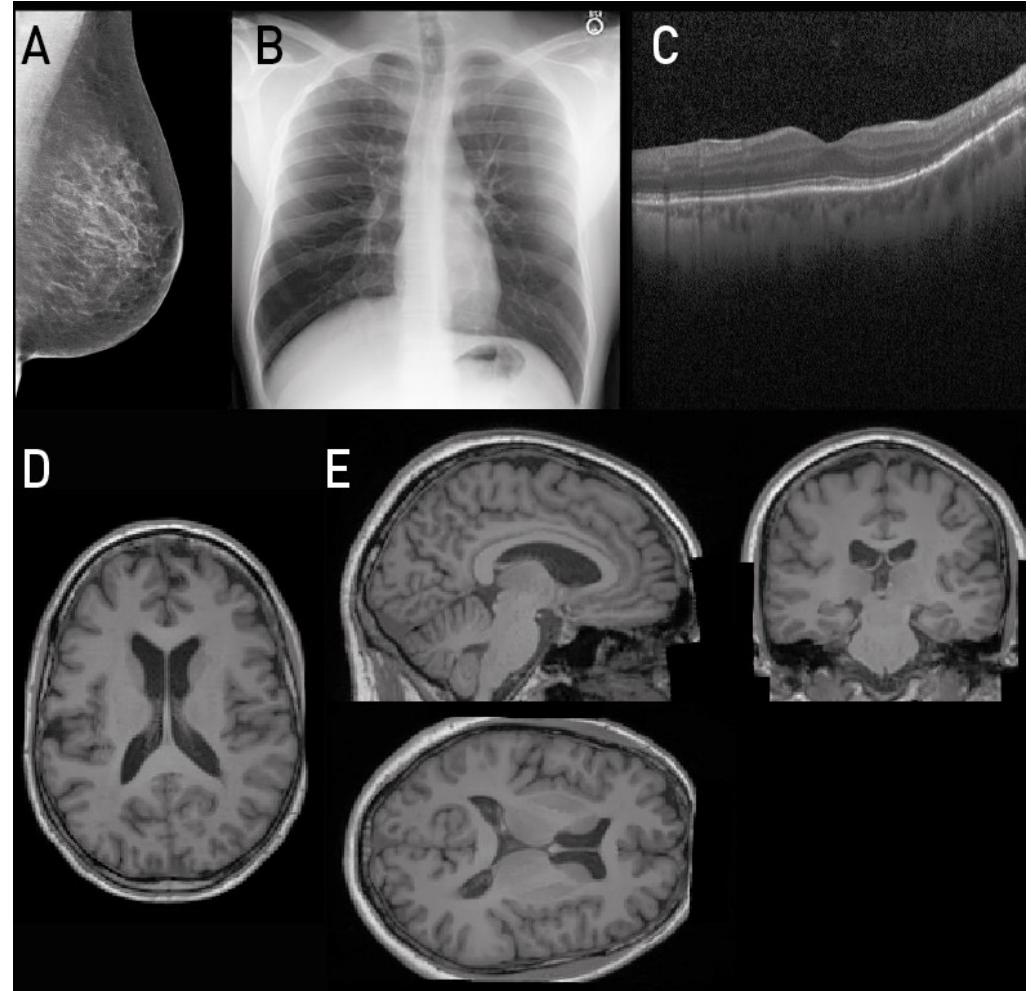
# Use of Synthetic Data

What can we use generated images for?

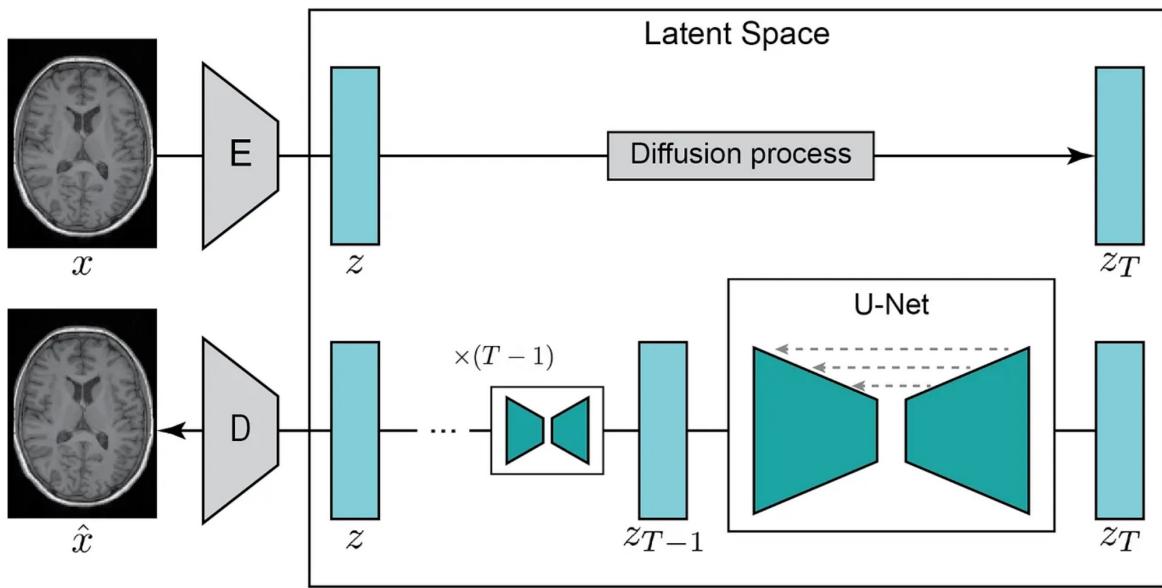
- Full private training
- Data augmentation
- Testing edge cases

Evaluation Criteria:

- Realism
- Diversity
- Privacy



# Generating High-Resolution 3D Brain Data

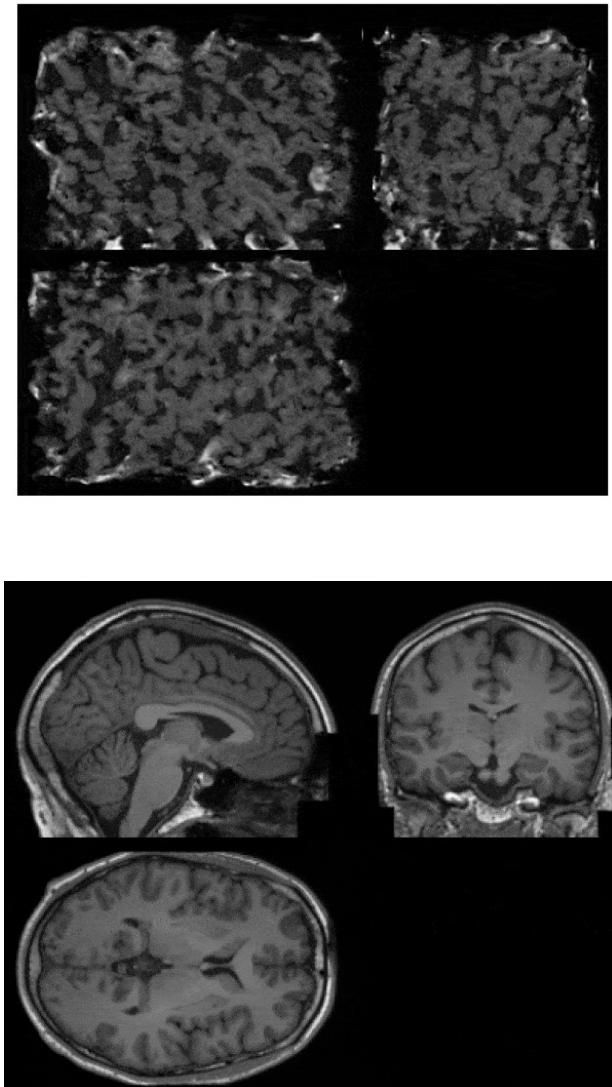


**Latent Diffusion Model** trained on data from UK Biobank (N = 31,740)

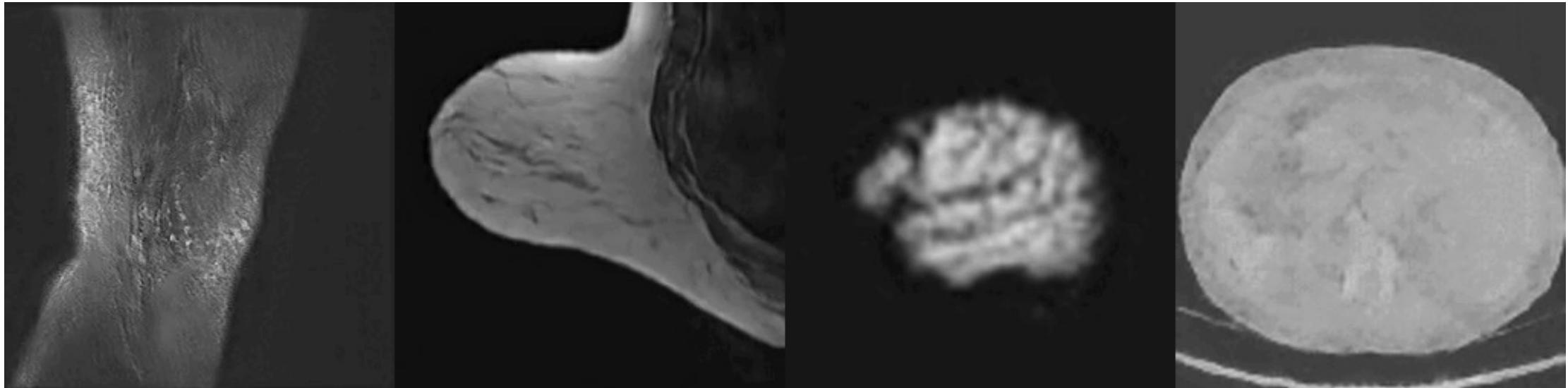
- Allows for the generation of T1-weighted brain MR images with a resolution of 160 x 224 x 160

Can be conditioned on covariates like:

- Age
- Gender
- Ventricular and Brain volumes



# High-Resolution 3D Medical Image Generation with LDMs



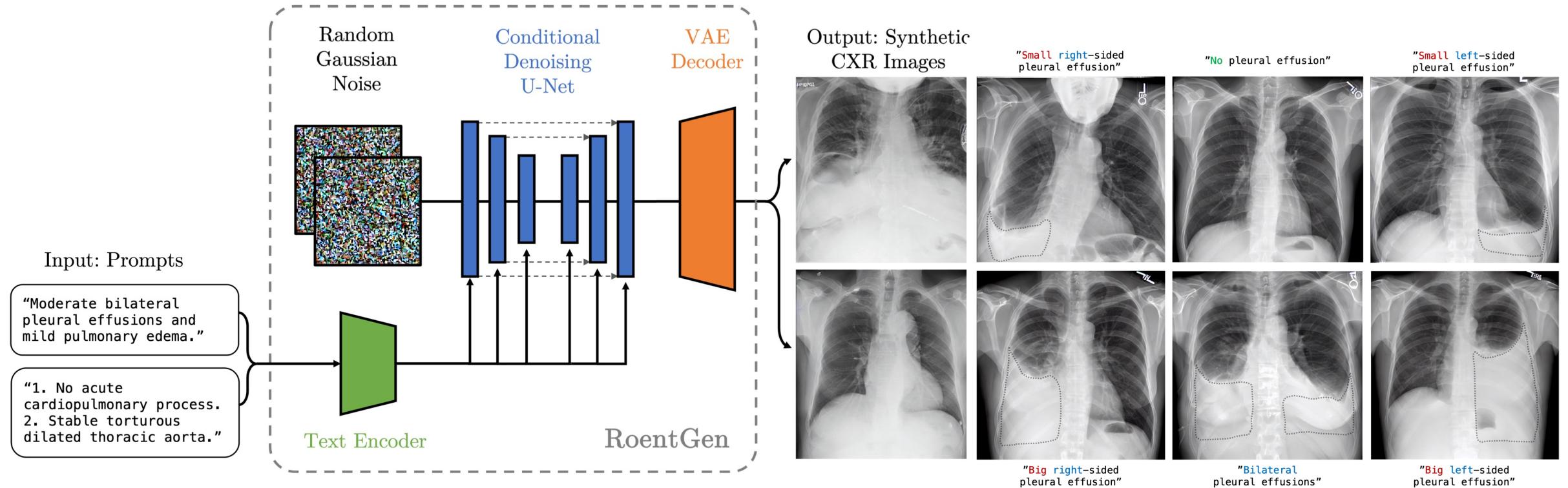
DUKE  
(256 x 256 x 32)

MRNet  
(256 x 256 x 32)

ADNI  
(64 x 64 x 64)

LIDC-IDRI  
(128 x 128 x 128)

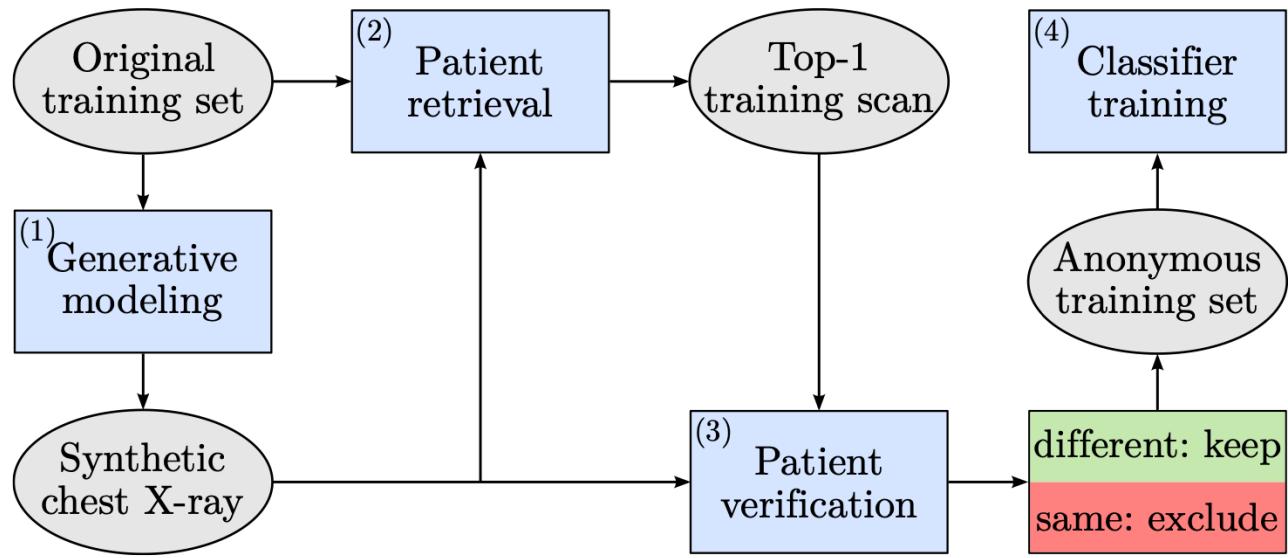
# Fine-Tuning Stable Diffusion



**Fine Tuned Stable Diffusion** pipeline for text-conditioned chest X-Ray generation

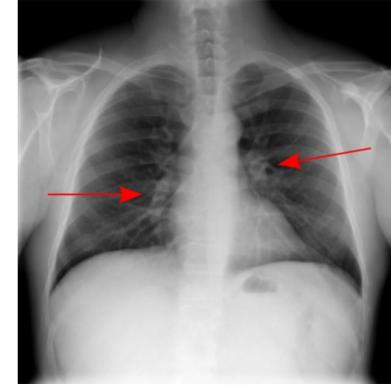
- U-Net and Text Encoder are jointly fine-tuned on the MIMIC-CXR dataset

# Generation of Anonymous Chest Radiographs

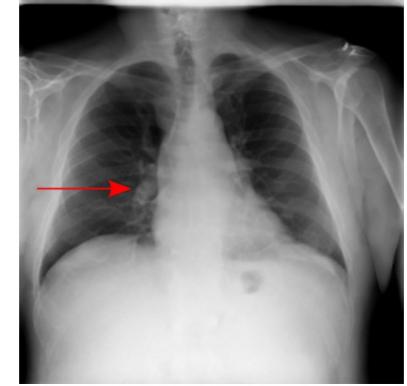


**Privacy enhanced sampling strategy for conditional chest X-Ray generation**

- Generated images are not per se private, as generative models may memorize training examples
- 1) Generate Image 2) Find most similar one from training data 3) Verify if patient is the same → include/exclude



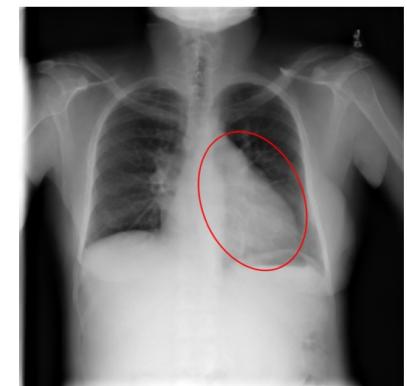
Infiltration



Nodule

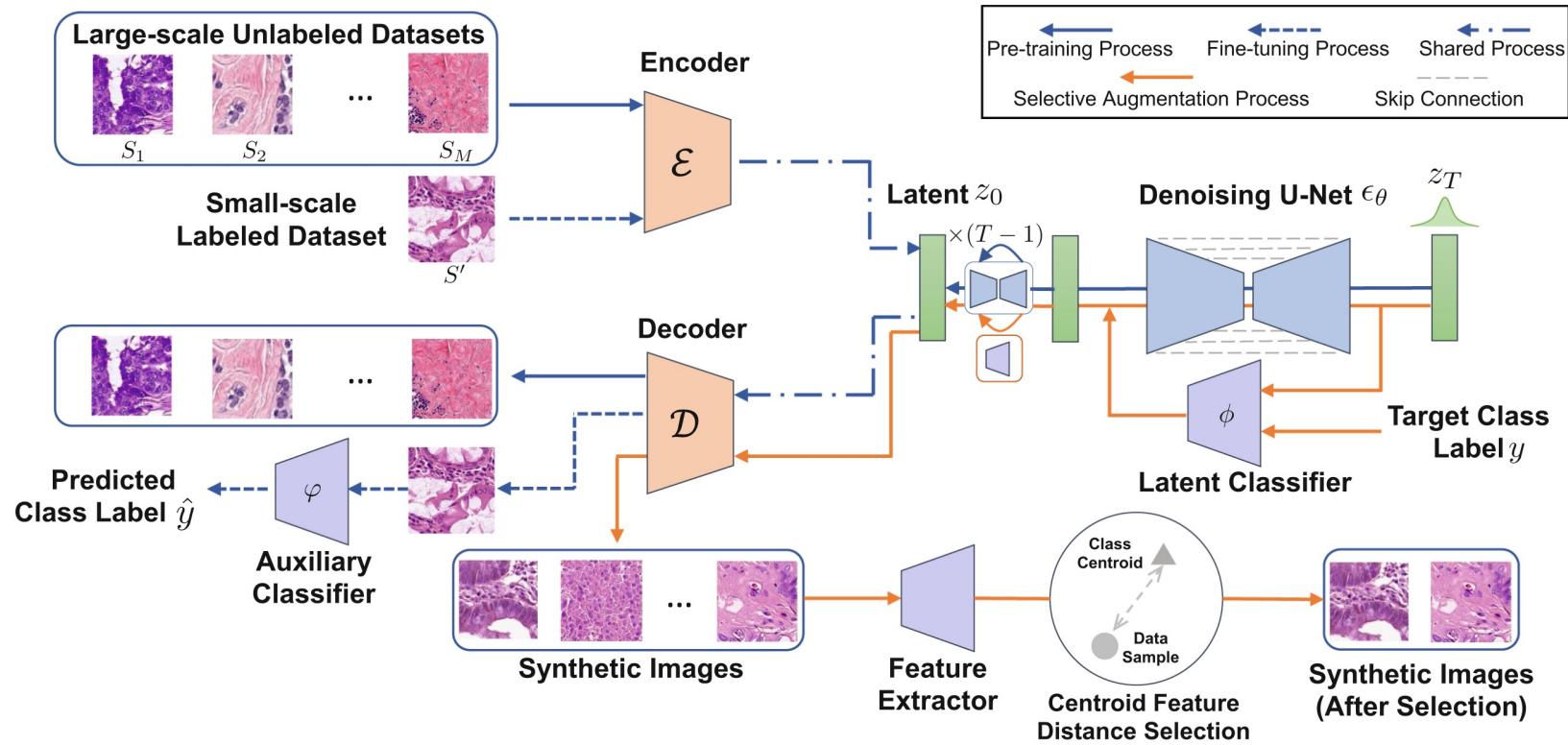


Mass



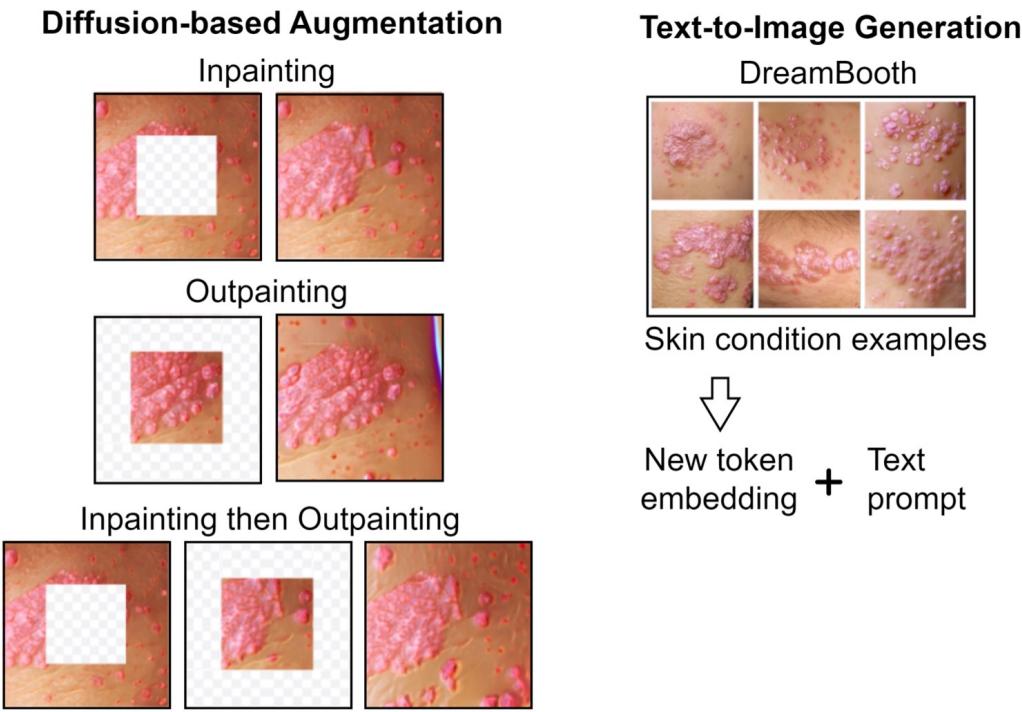
Cardiomegaly

# Synthetic Data Augmentation



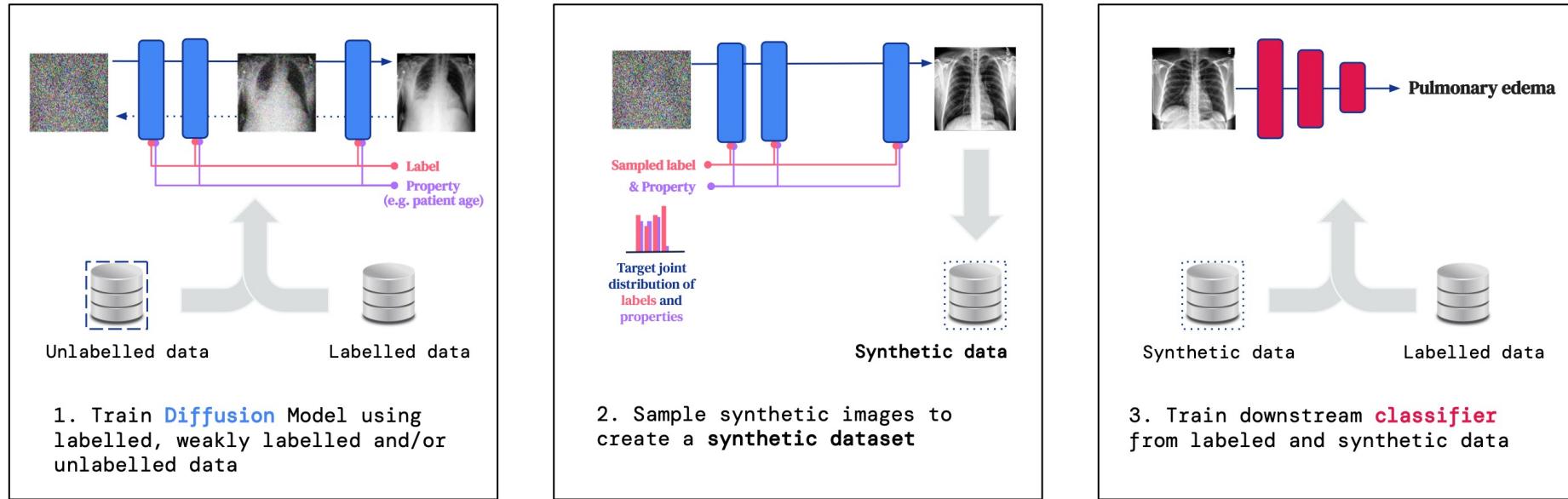
- 1) Pre-training an unconditional Latent Diffusion Model on a large unlabeled dataset
- 2) Conditional fine-tuning on an unseen labeled (small) dataset through a latent classifier
- 3) Selection of the highly-confident synthetic samples based on feature similarity with real data

# Synthetic Image Augmentation – Performance Gain



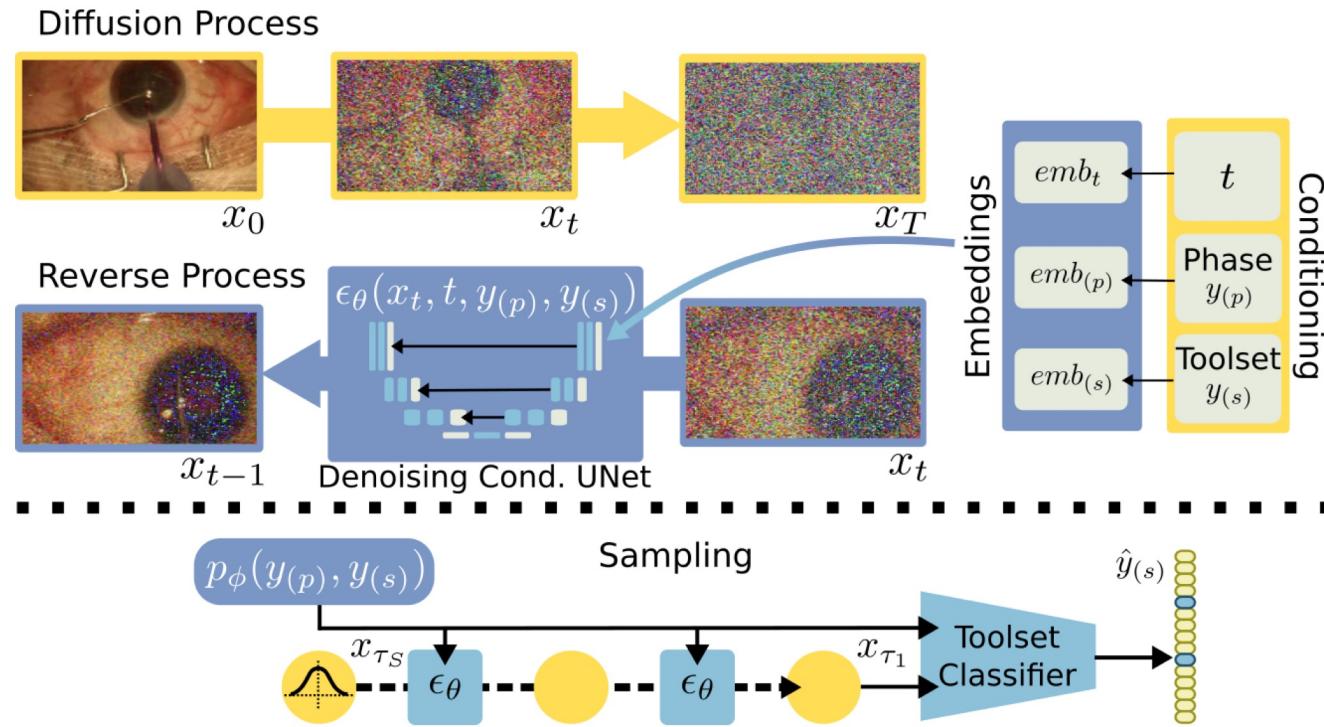
- Diffusion models can scalably generate images of skin disease and training models with that augmented data **improves performance in data-limited settings**
- Performance gains (in this setting) saturate at a **synthetic-to-real ratio of 10:1** and is substantially smaller than the gains obtained by adding real images → **collection of diverse real data remains more important**

# Synthetic Data for Distribution Shifts – Improving Fairness in AI



- Learning realistic augmentations from data is possible in a label-efficient manner using diffusion models
- Unlabeled data can be used to capture the data distribution of different conditions and subgroups → **steer the distribution of synthetic examples according to specific requirements**
- Learned augmentations can surpass heuristic, manually implemented ones by making **models more robust and statistically fair in- and out- of-distribution**

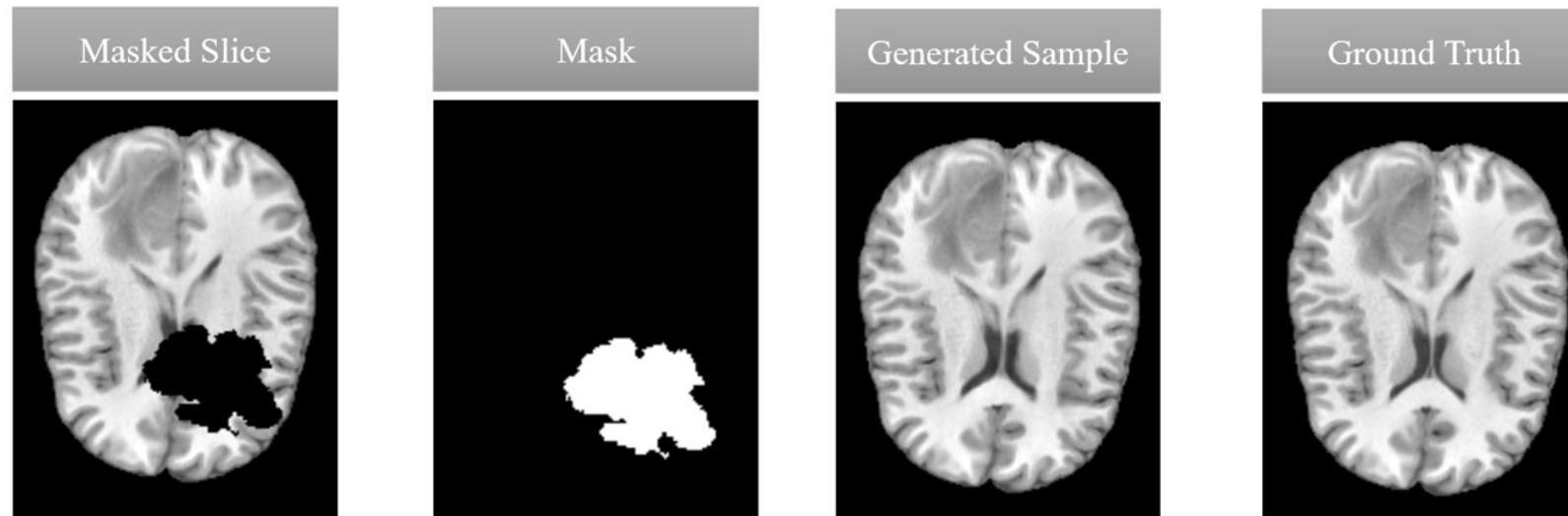
# Synthesising Rare Samples – Testing Edge Cases



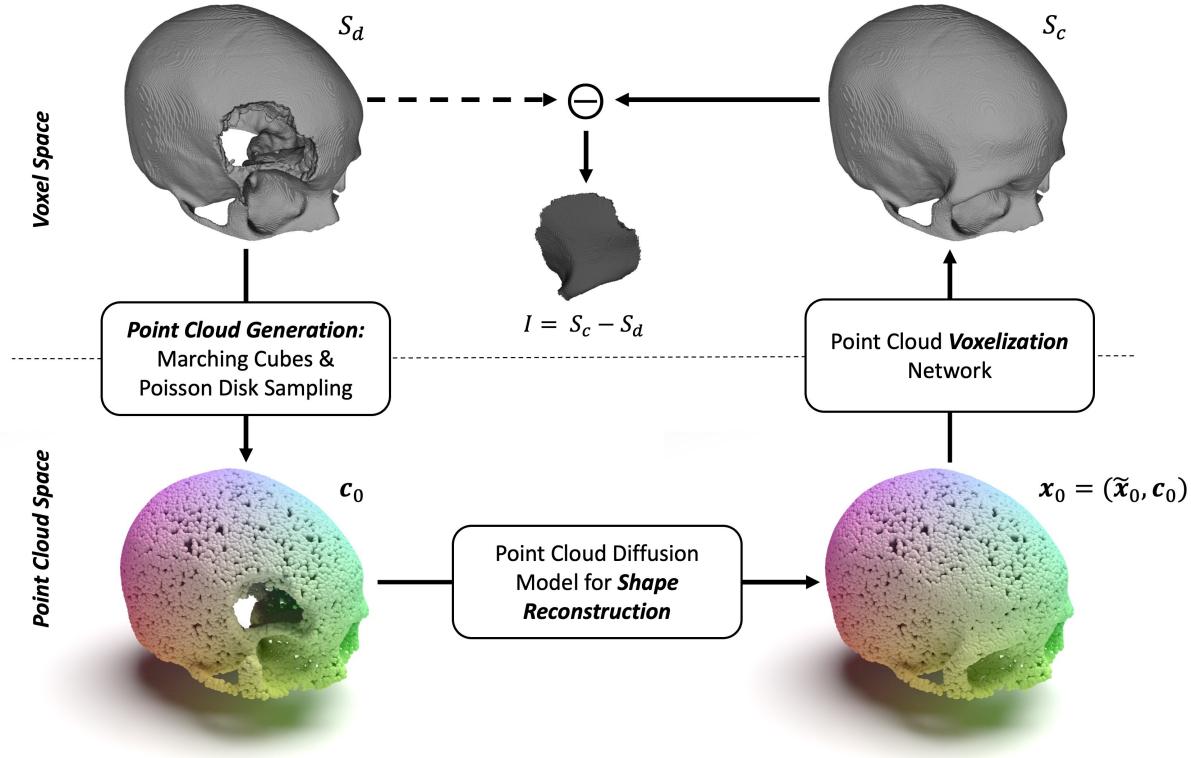
- Diffusion model is used to **synthesize rare cases** in instrument classification in cataract surgery
- Performance gain of > 10% for rare cases

# Inpainting

Examples from the Community

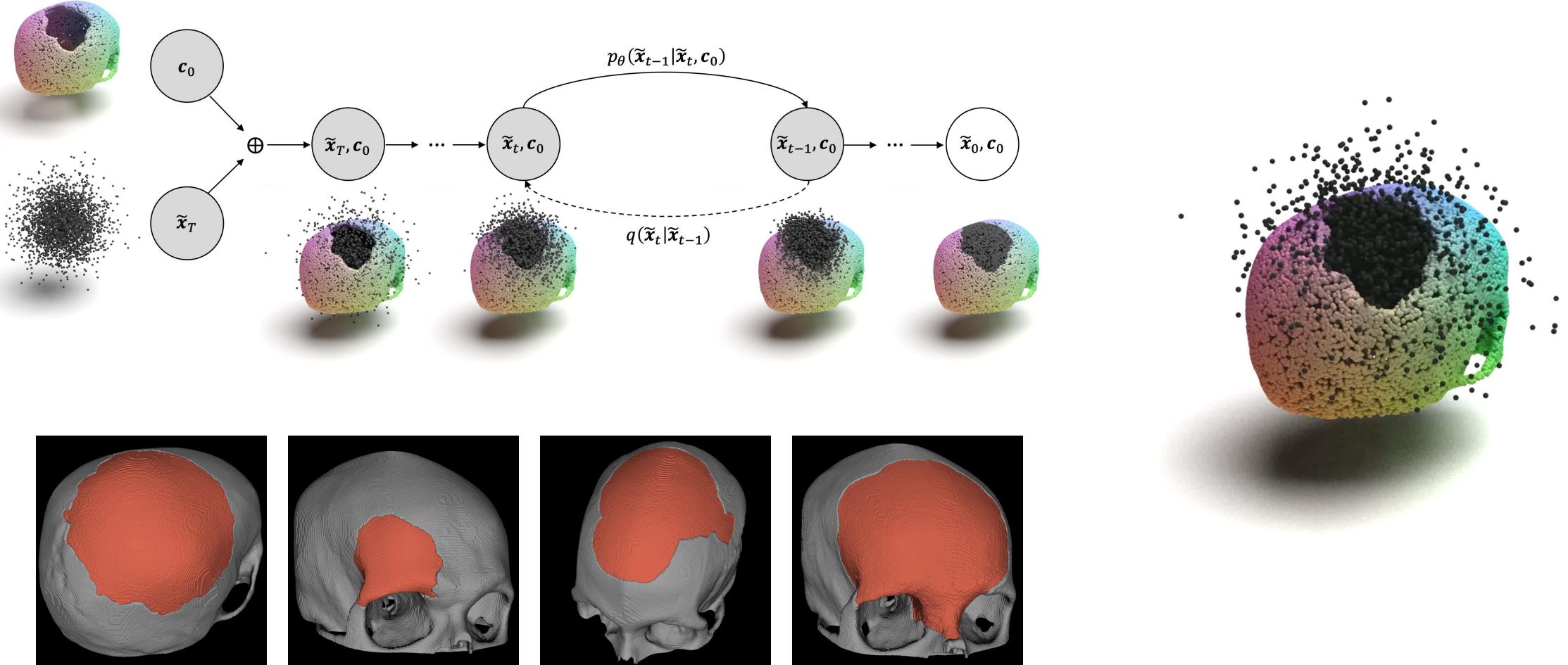


# Point Cloud Diffusion Models for Automatic Implant Generation



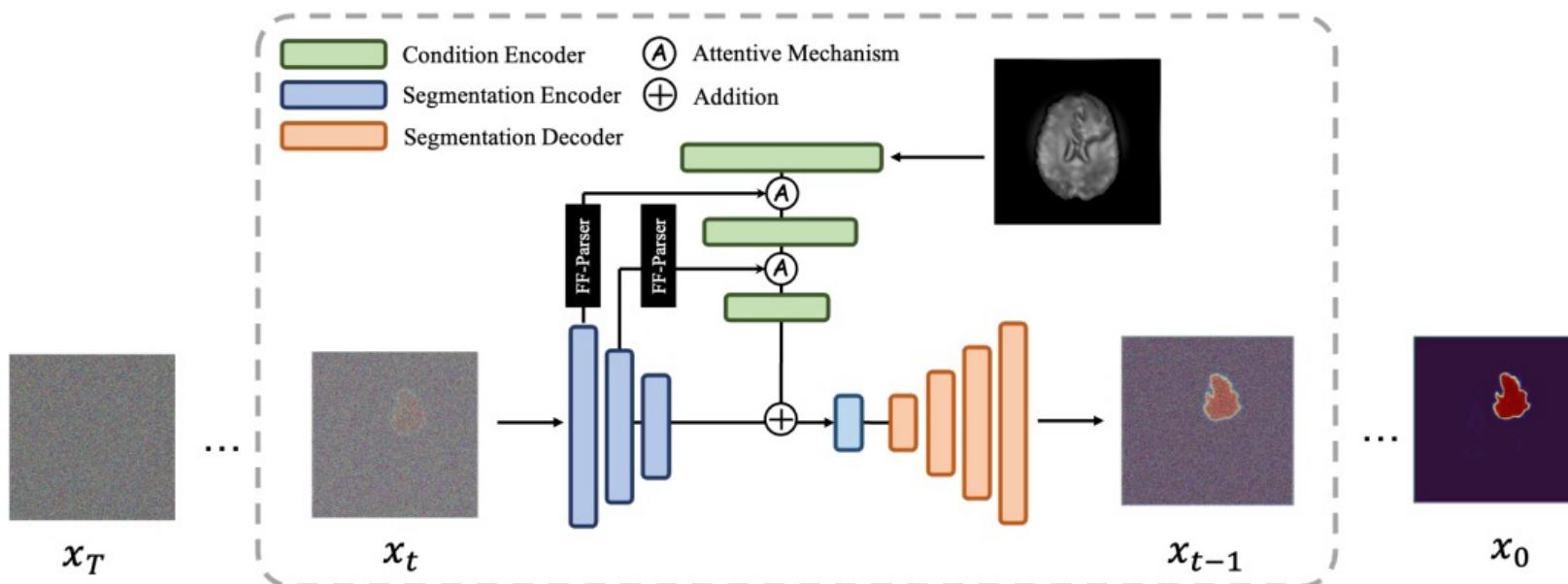
- Implant generation as shape completion task (basically 3D inpainting)
- Applying the diffusion model to high-resolution volumes ( $512 \times 512 \times 512$ ) is impossible due to limited GPU memory
- Diffusion models **can be applied to a wide variety of different data types** (like point clouds)

# Point Cloud Completion with Diffusion Models

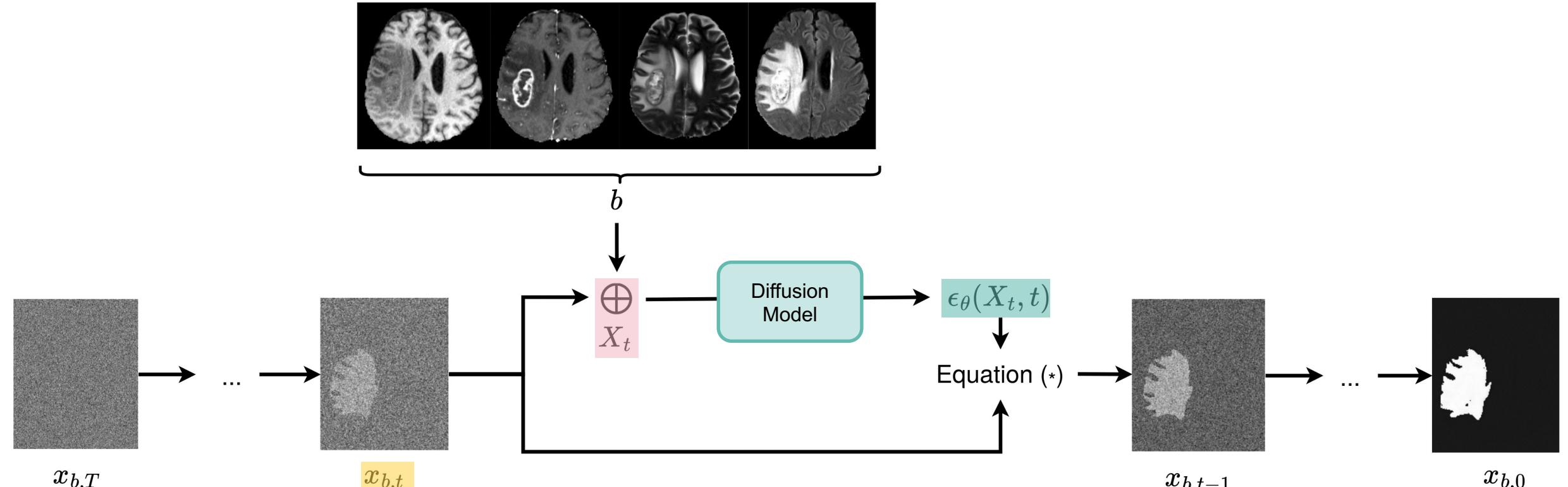


# Image Segmentation

Examples from the Community



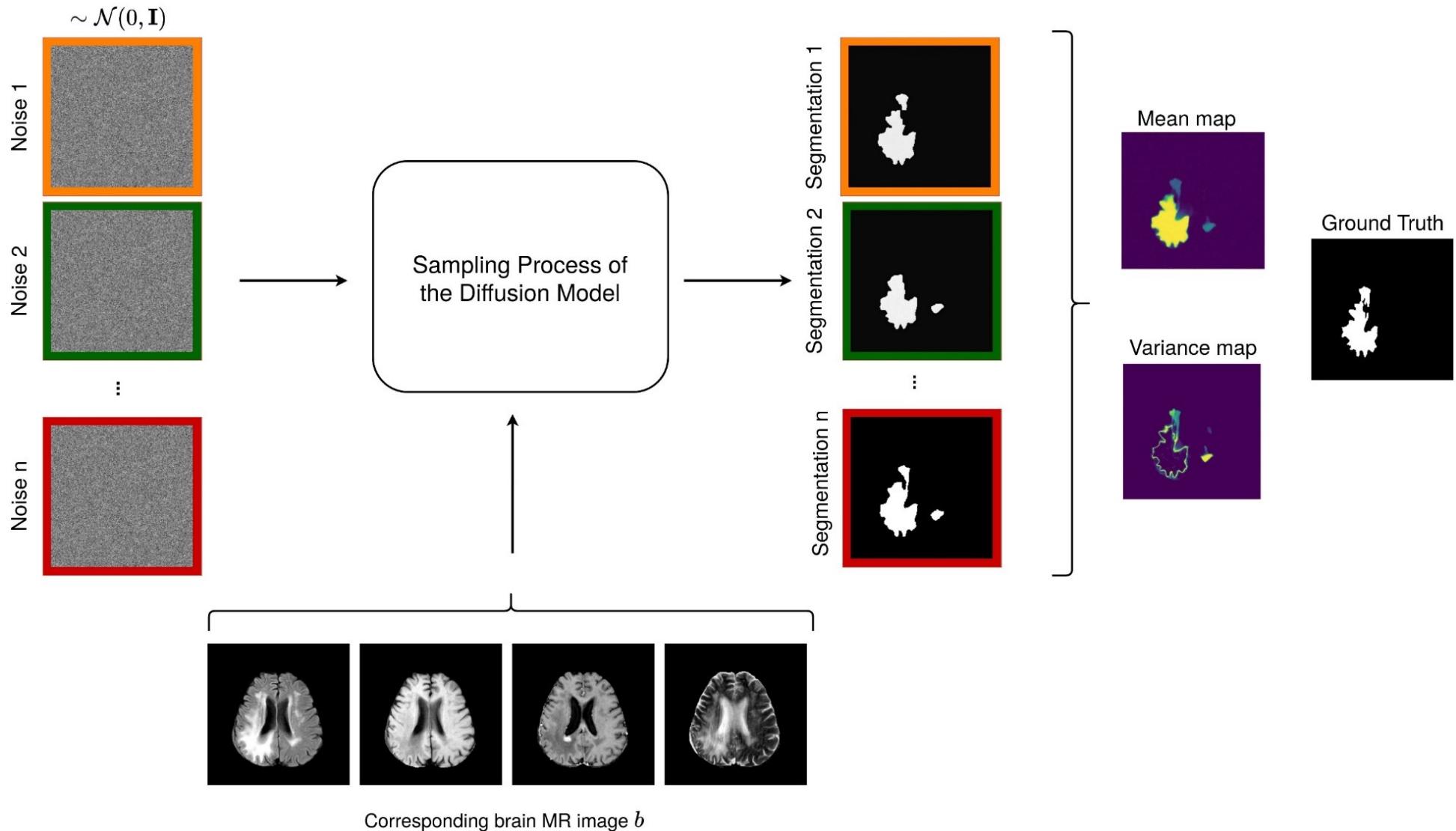
# Diffusion Models for Segmentation Mask Generation



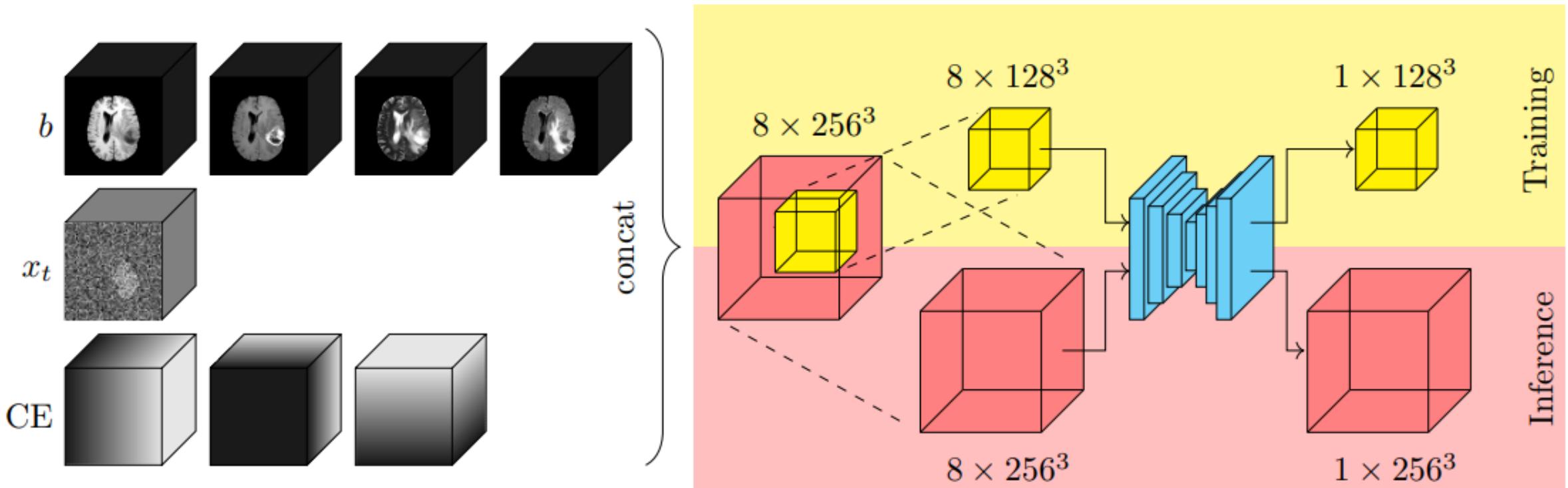
$$(*) \quad x_{b,t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x_{b,t} - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(X_t, t) \right) + \sigma_t \mathbf{z}, \quad \text{with } \mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$$

The anatomical information is added by concatenating the input images  $b$  to the noisy segmentation mask  $x_{b,t}$  in every step  $t$ .

# Generation of Segmentation Ensembles



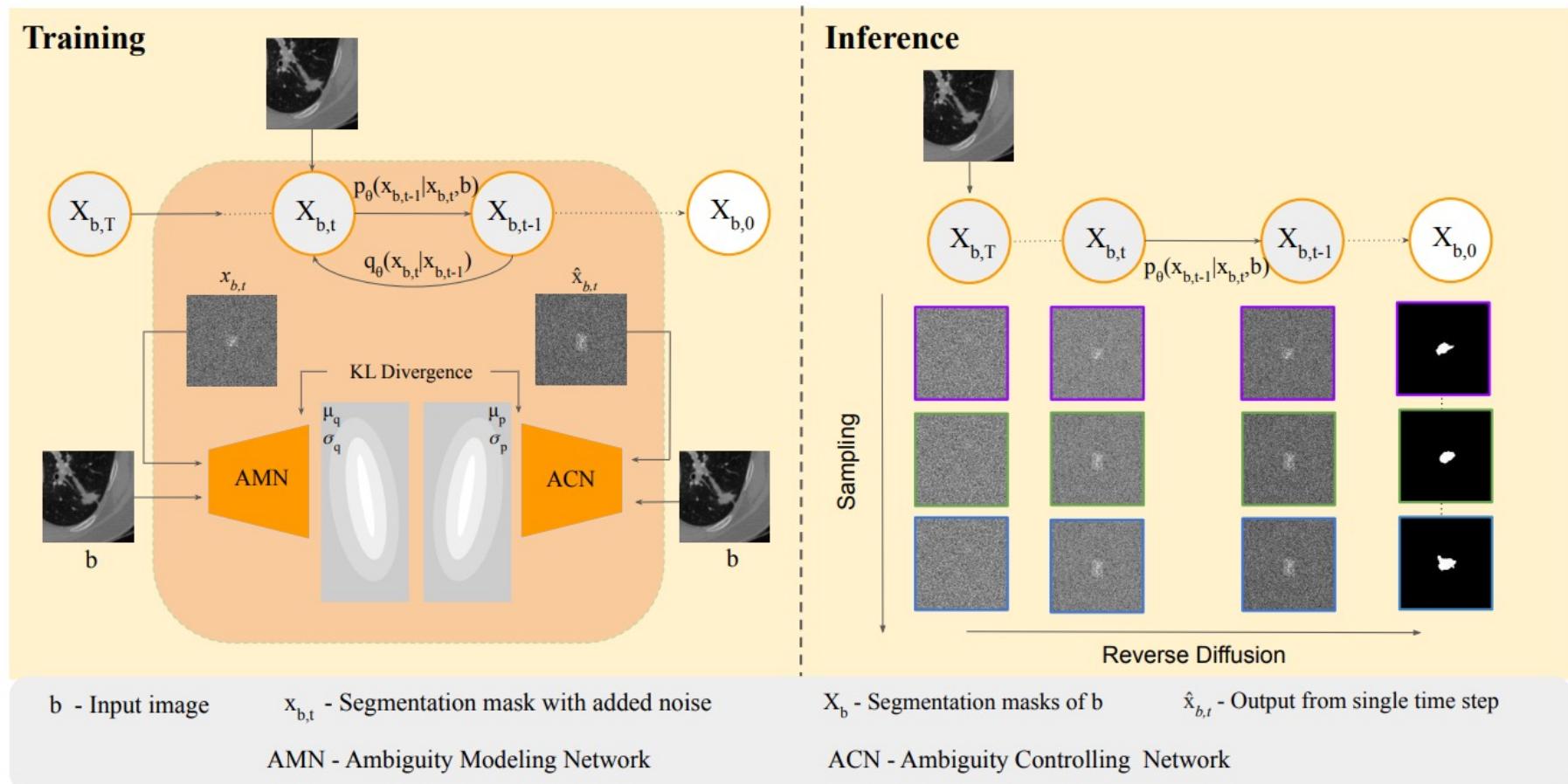
# 3D Segmentation with PatchDDM



- We add a position encoding in all 3 spatial dimensions.
- Training is on patches only, and saves memory and training time.
- Inference runs over the whole 3D volume.

# Ambiguous Segmentation

- Ambiguity Modelling Network (AMN) models the distribution of ground truth masks given an input image.
- Ambiguity Controlling Network (ACN) models the noisy output from the diffusion model conditioning on an input image.



# Segmentation with Diffusion Pre-training

## Diffusion

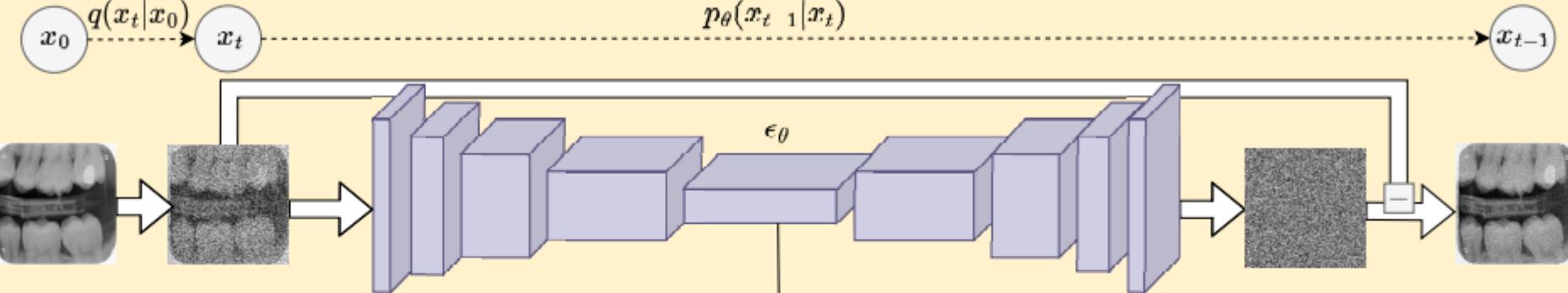
$$x_0 \in X_1$$

$$\epsilon \sim \mathcal{N}_{0,I}$$

$$t \sim \mathcal{U}_{1,T}$$

$$\nabla_{\theta} \|\epsilon_{\theta}(x_t, t) - \epsilon\|^2$$

## Pre-training



## Few label

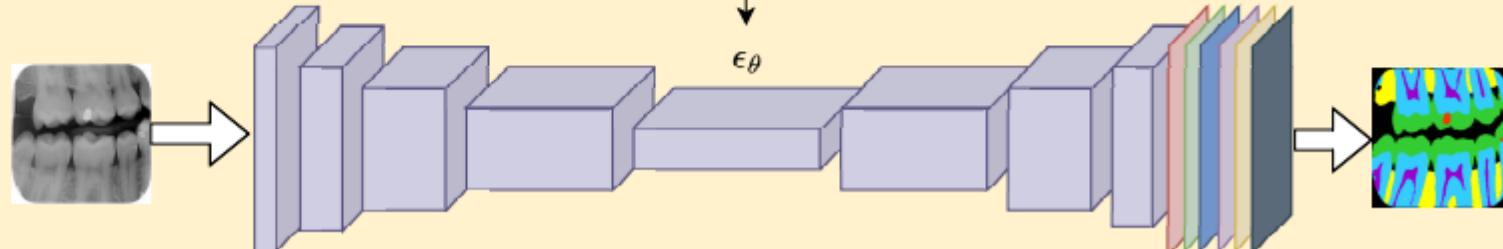
$$X_1 \cap X_2 = \emptyset$$

$$(x, y) \in X_2 \times Y$$

$$\hat{y} = \epsilon_{\theta}(x)$$

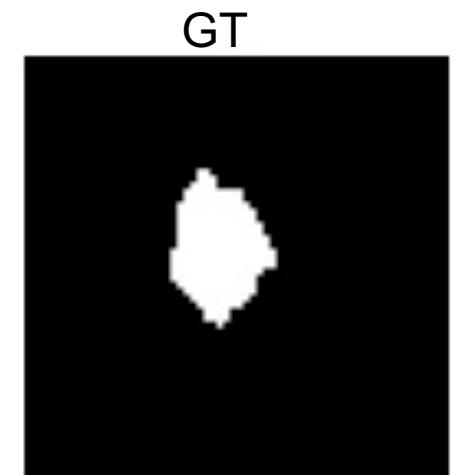
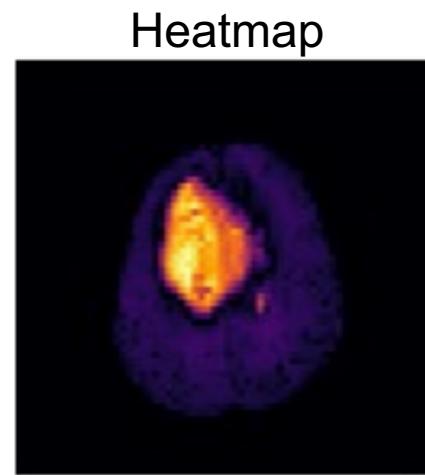
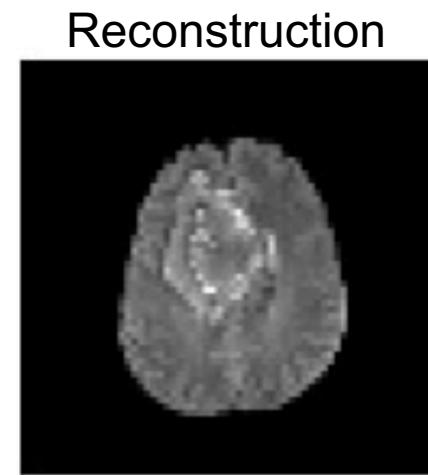
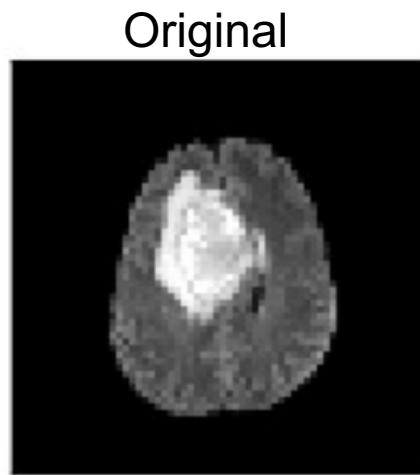
$$\nabla_{\theta} Loss(\hat{y}, y)$$

## Finetuning



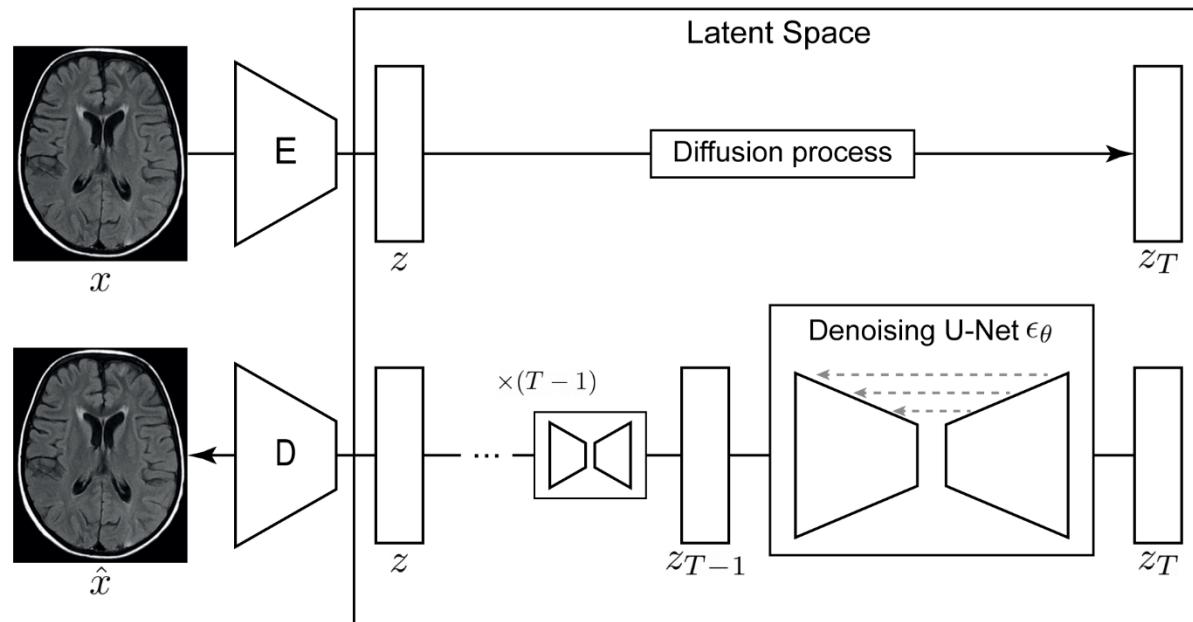
# Anomaly Detection

Examples from the Community

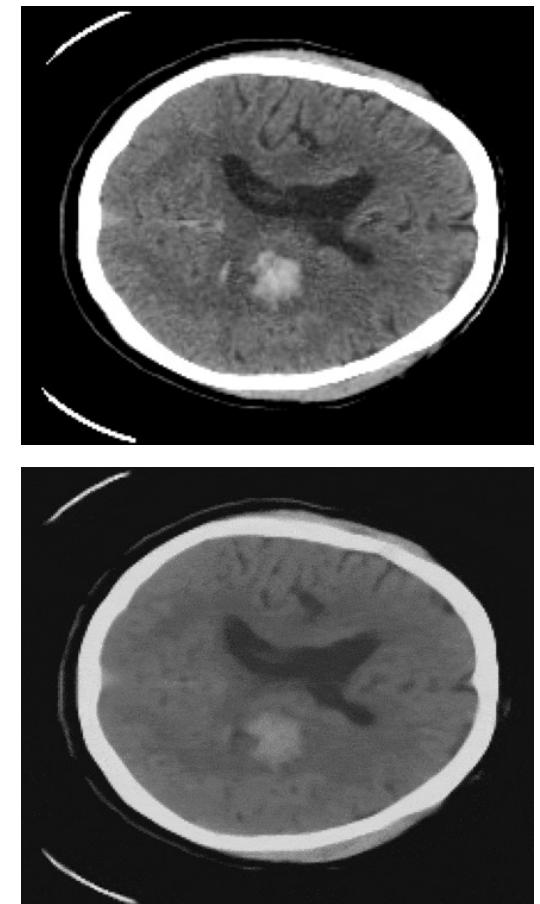


# Unsupervised Anomaly Segmentation

Latent Diffusion Model (LDM) learns the distribution of healthy brain data  
Compression (VQ-VAE) scales for high-resolution images



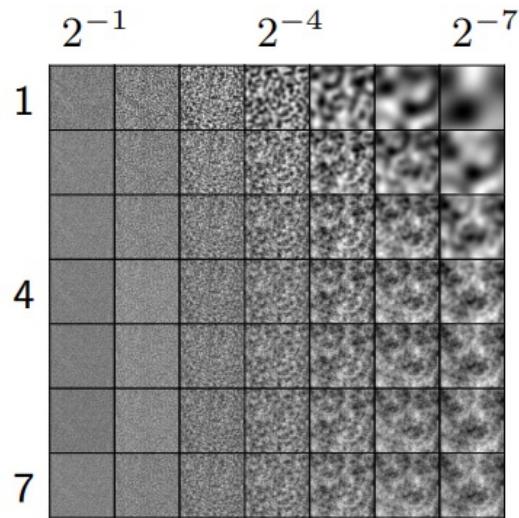
LDM identify regions with a low likelihood of being part of the healthy dataset



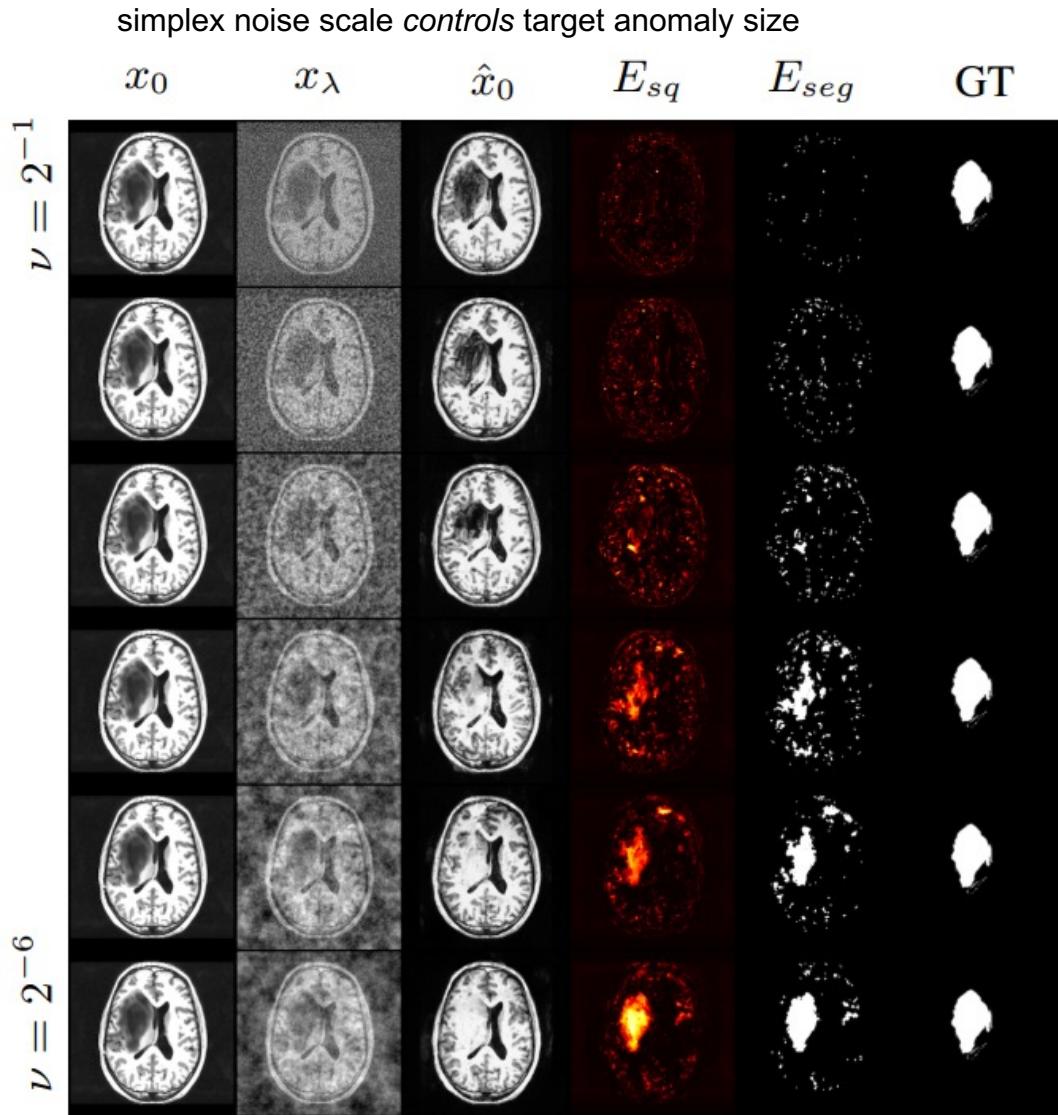
Reverse/denoising process is used to **inpaint** these regions and “heal” the possible anomalies

# Anomaly Detection with Simplex Noise

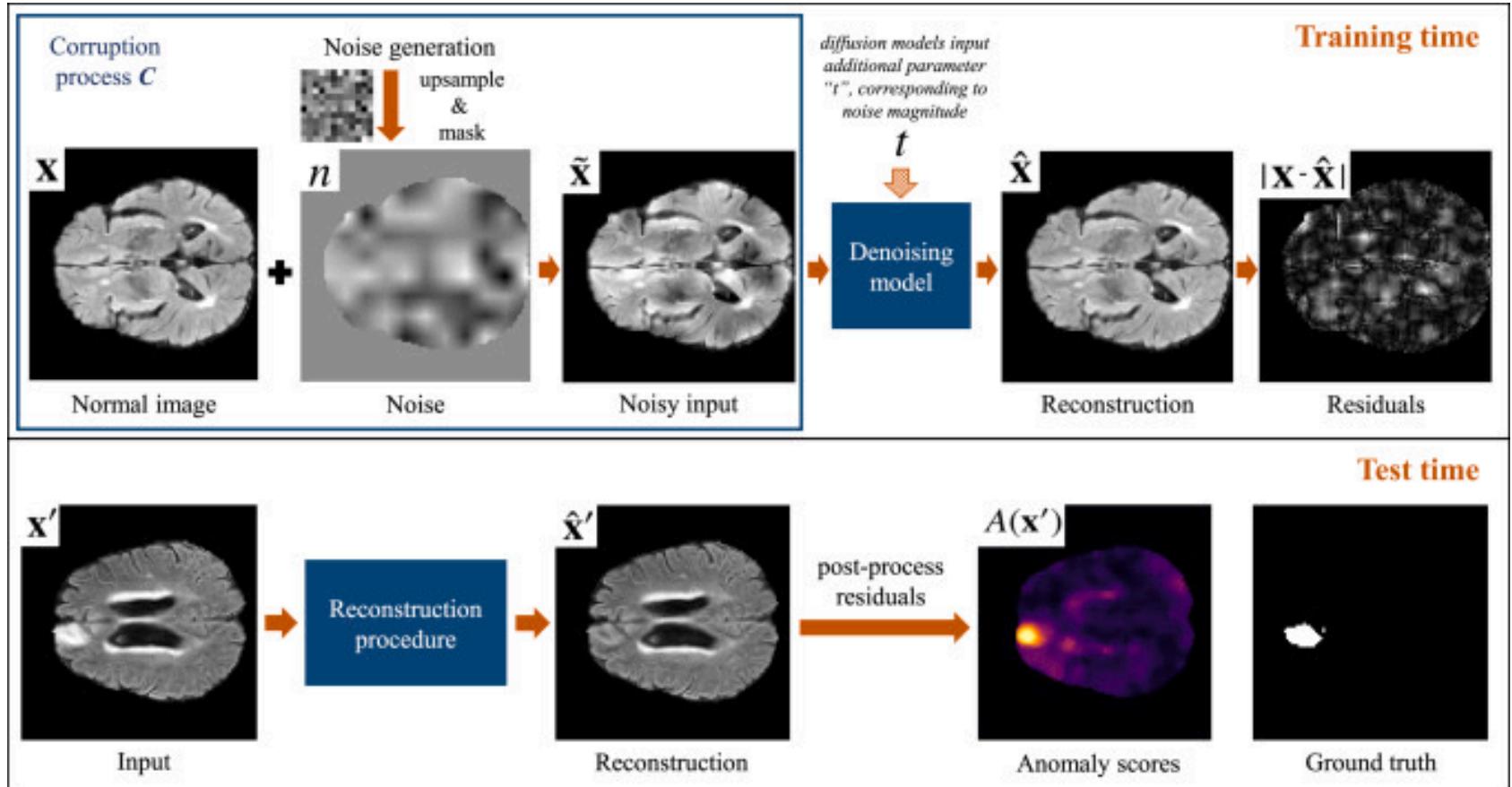
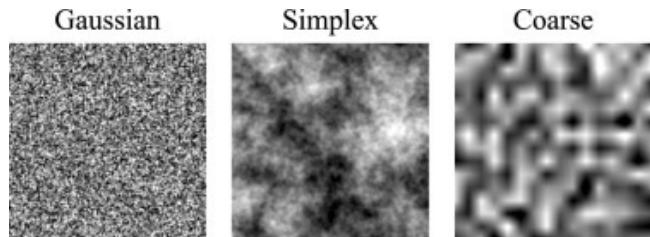
- Typical Gaussian noise is found to be insufficient for anomaly detection.
- Therefore, we explore the use of simplex noise for the corruption and sample generation of medical images.



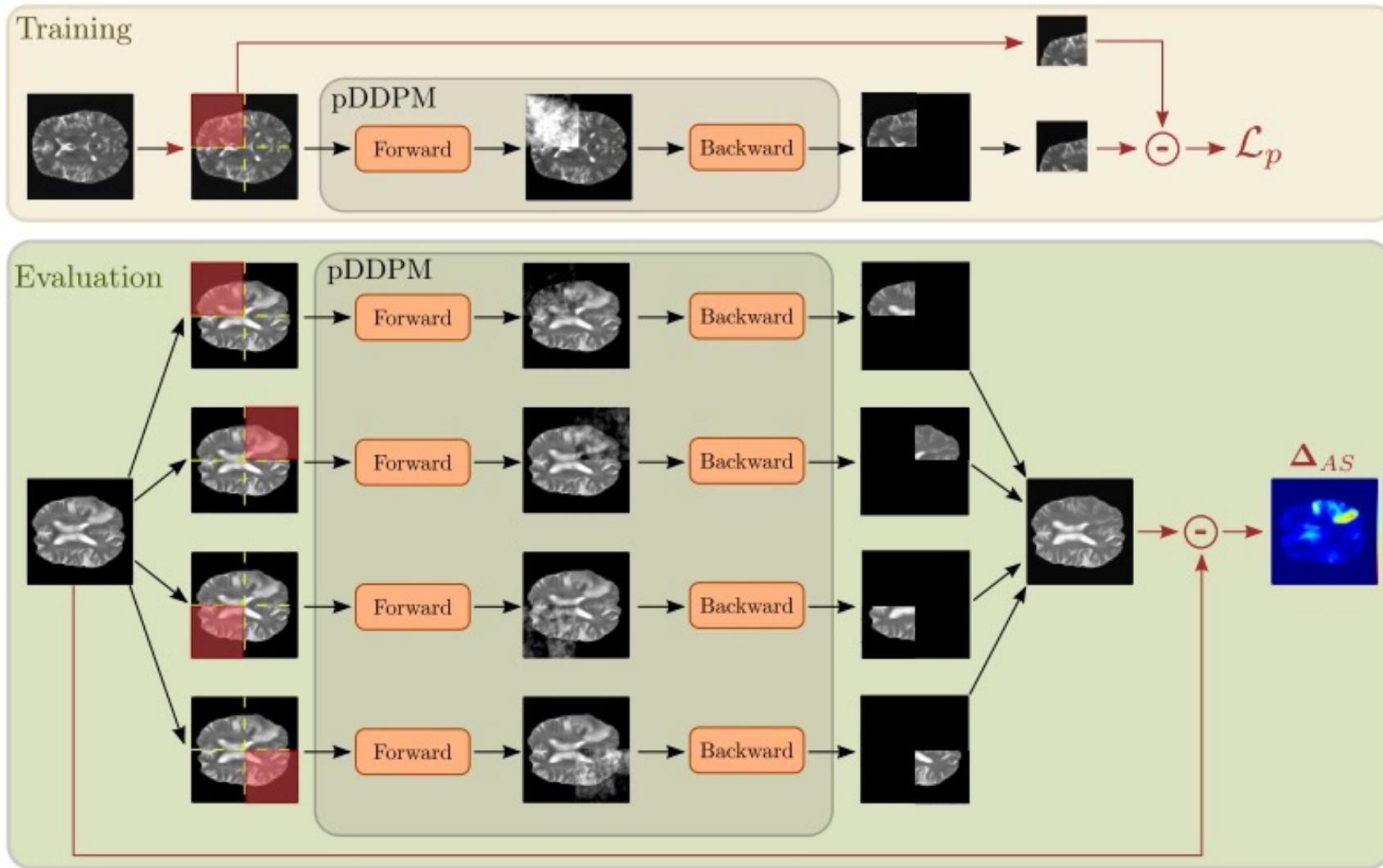
(a) Structures of simplex noise



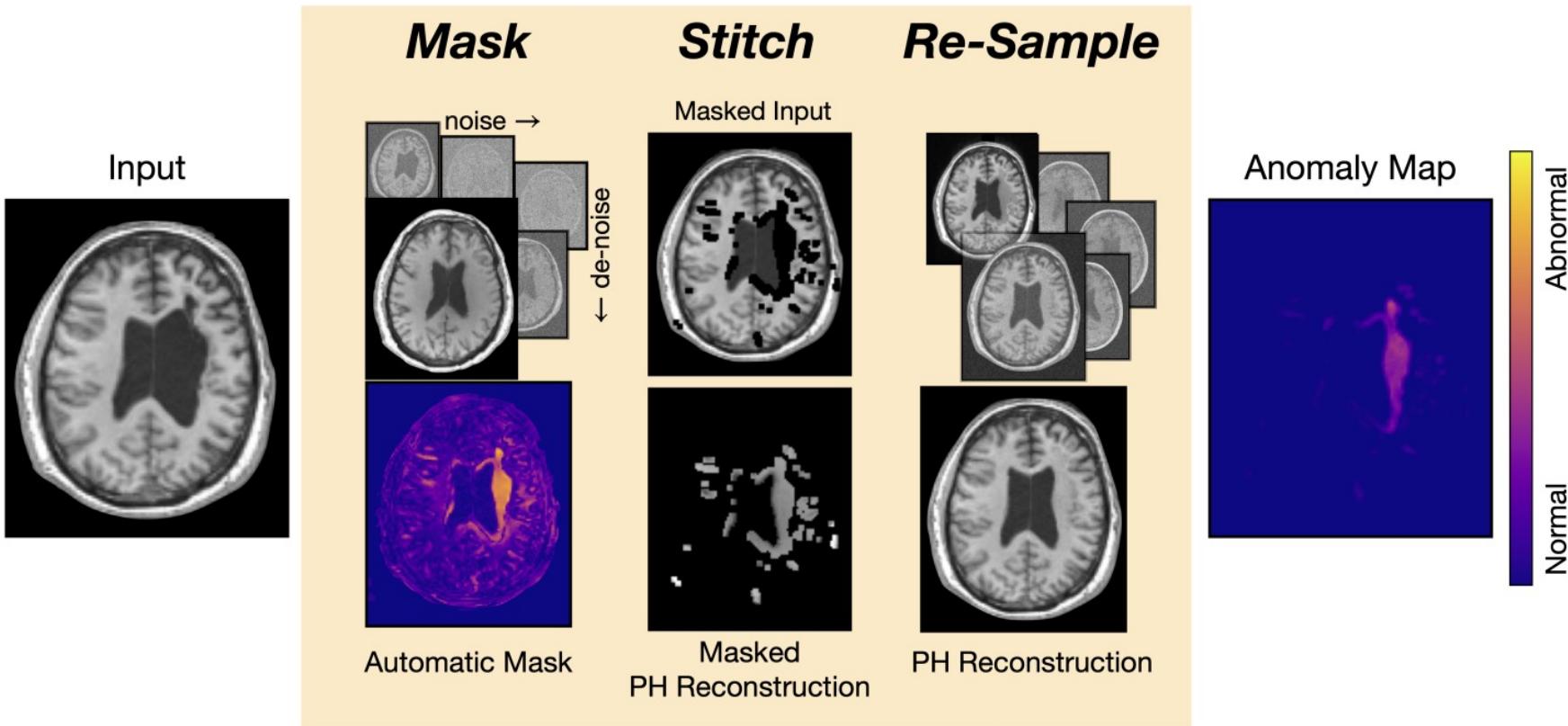
# Anomaly Detection with Coarse Noise



# Anomaly Detection from Patches

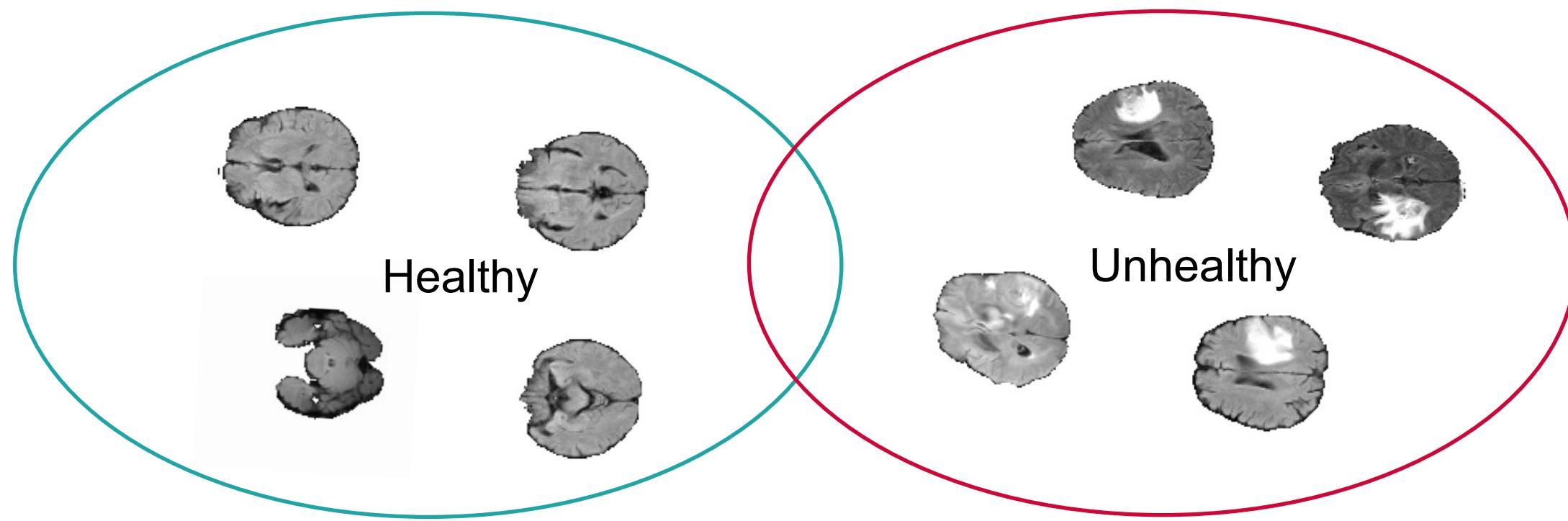


# Mask, Stitch, and Re-Sample

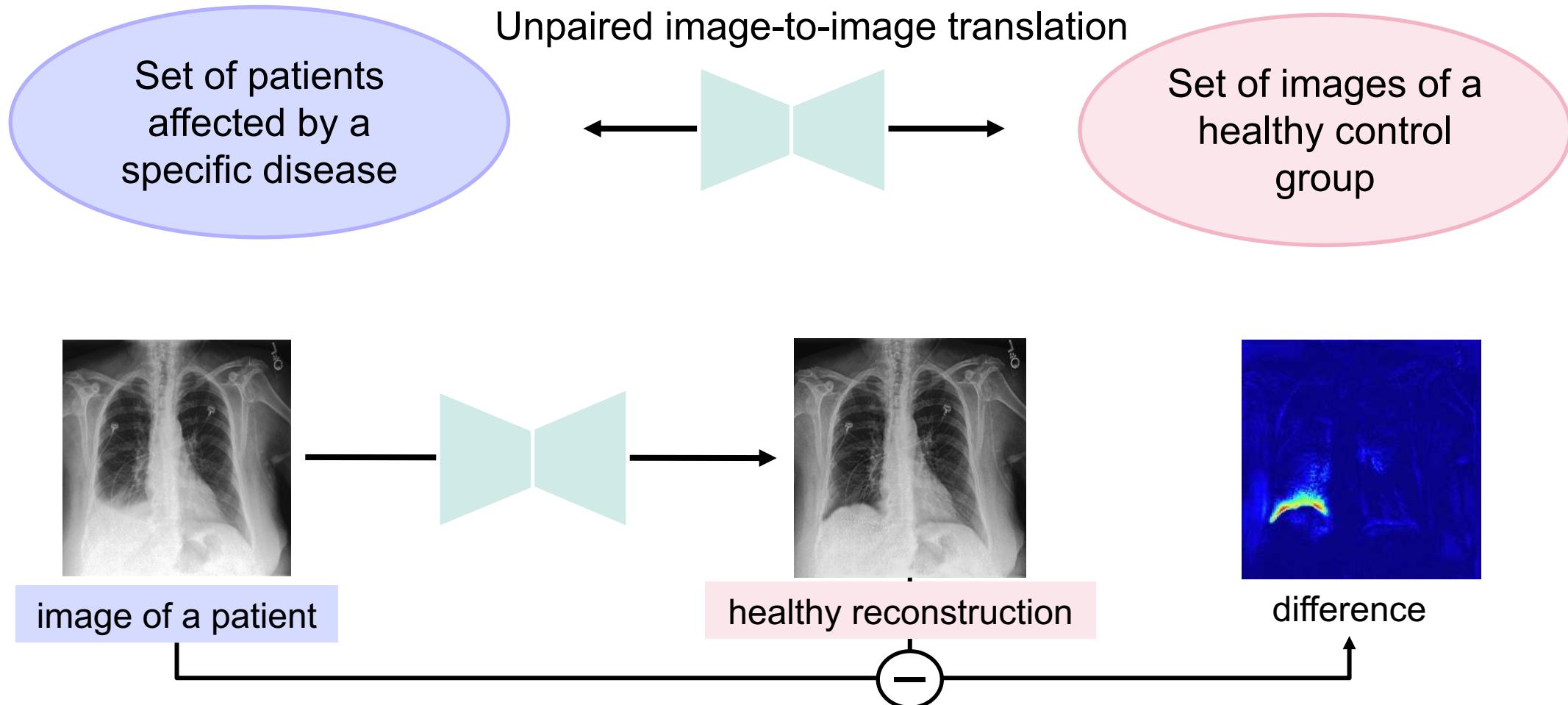


# Weakly Supervised Lesion Detection

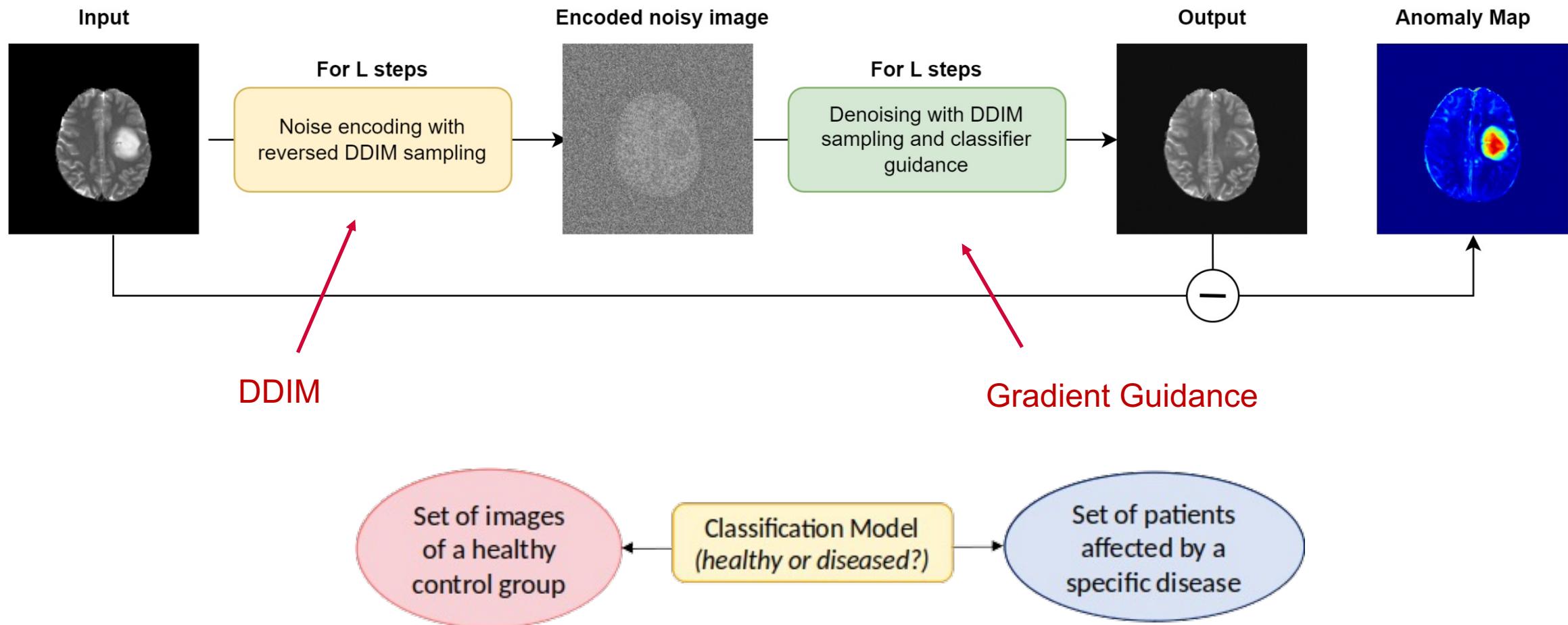
**Goal:** Pixel-wise anomaly detection using image-level labels only



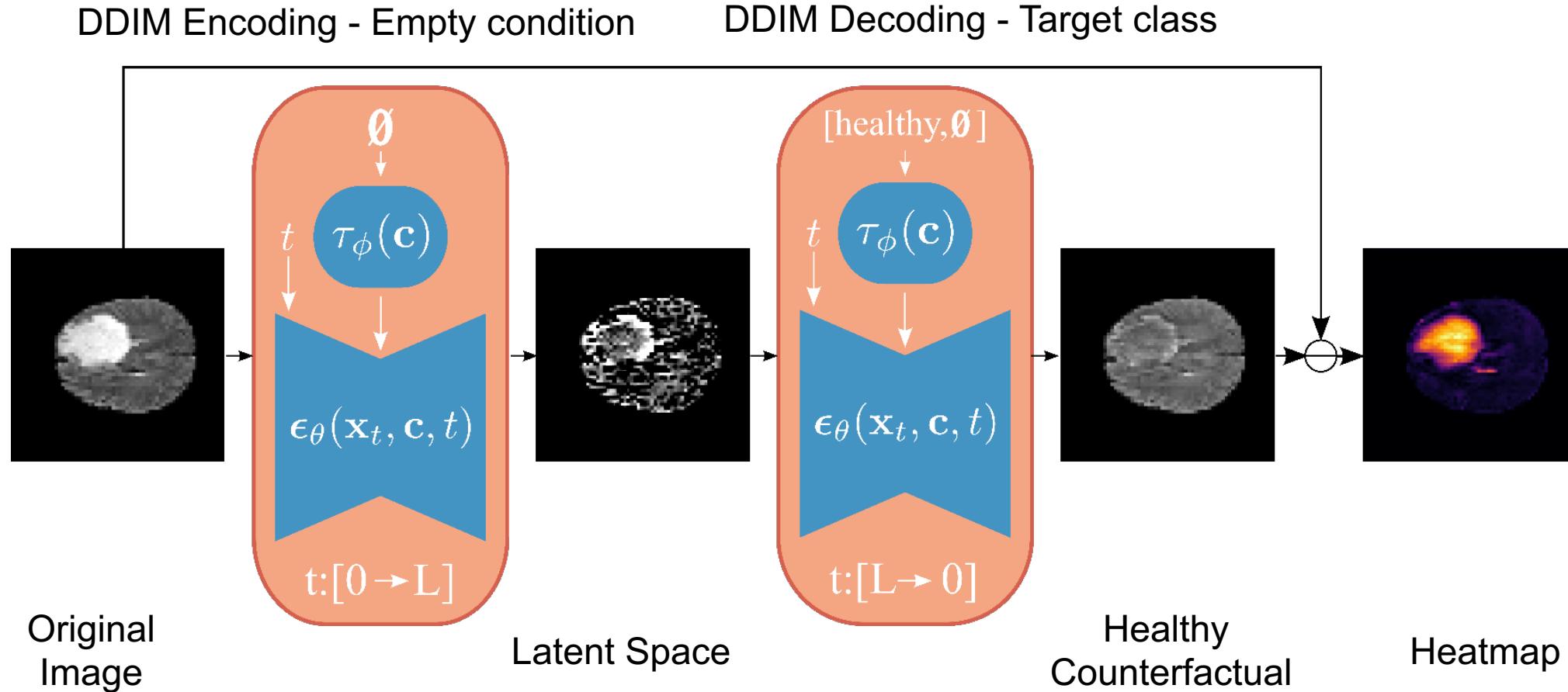
# Weakly Supervised Lesion Detection



# Weakly Supervised Lesion Detection

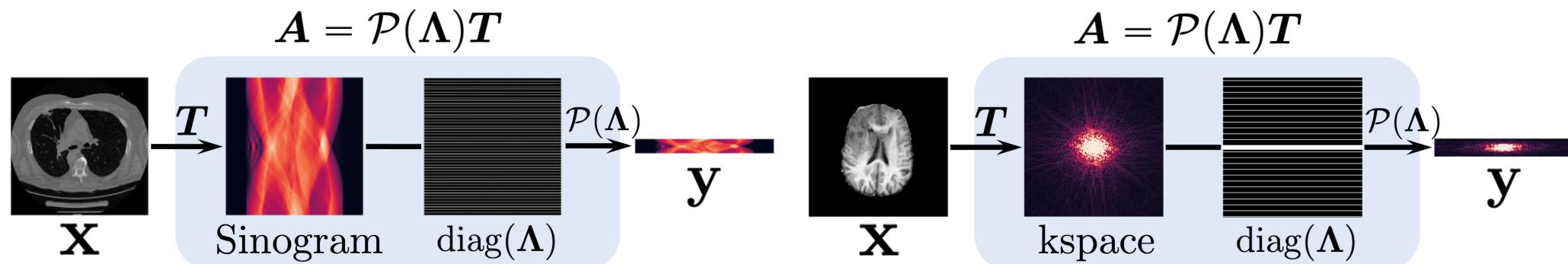


# Lesion Localization with Classifier-free Guidance

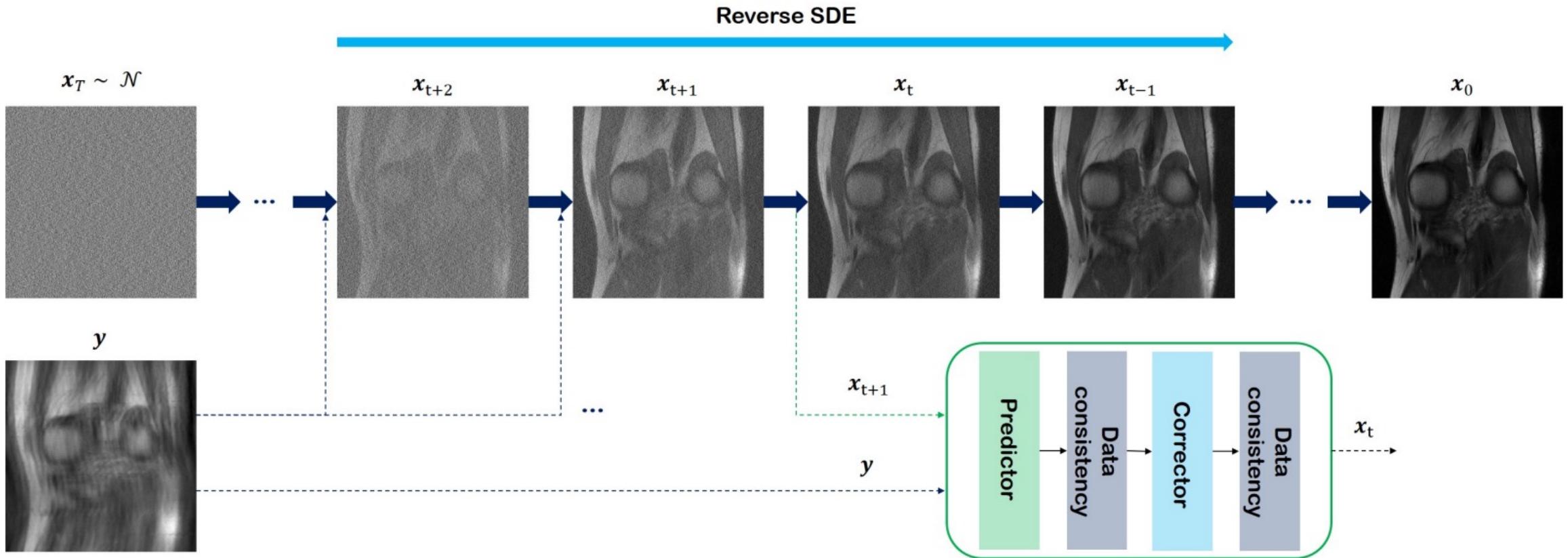


# Image Reconstruction

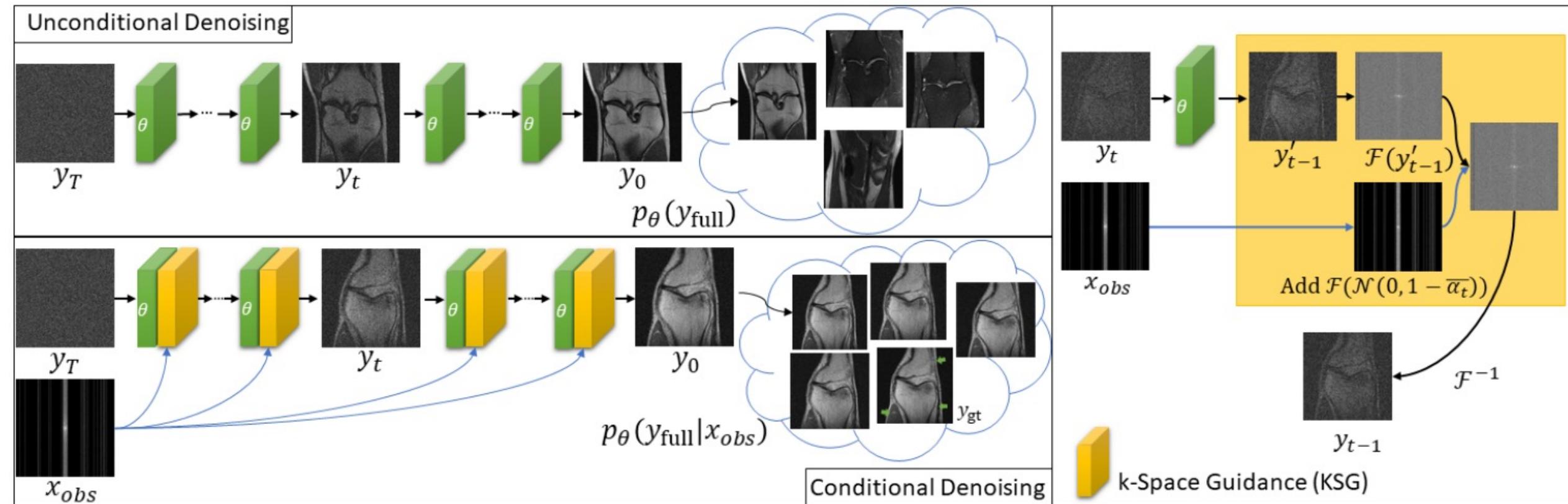
Examples from the Community



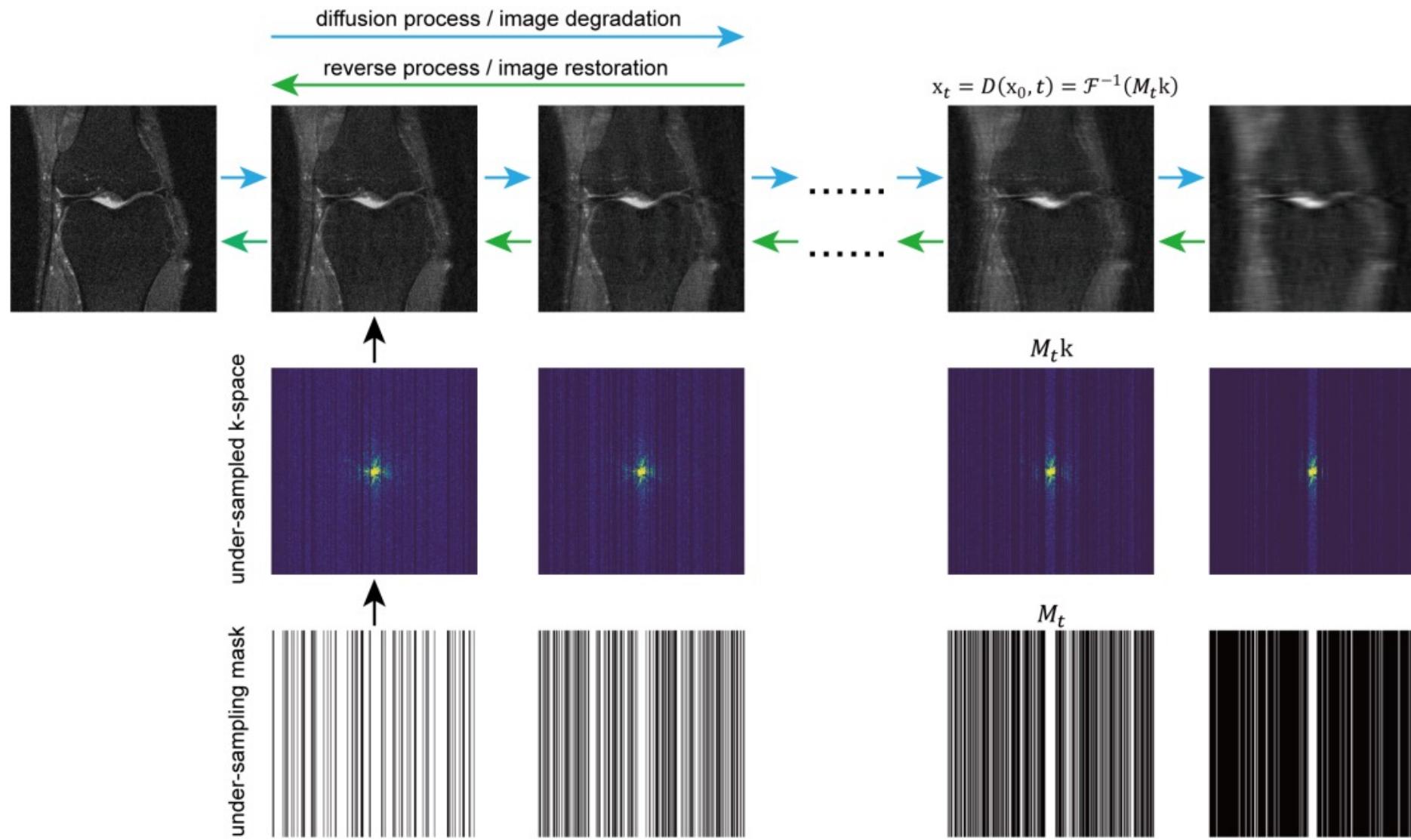
# Score-based Diffusion Models for Accelerated MRI



# Undersampled MR Reconstruction via Diffusion Model Sampling



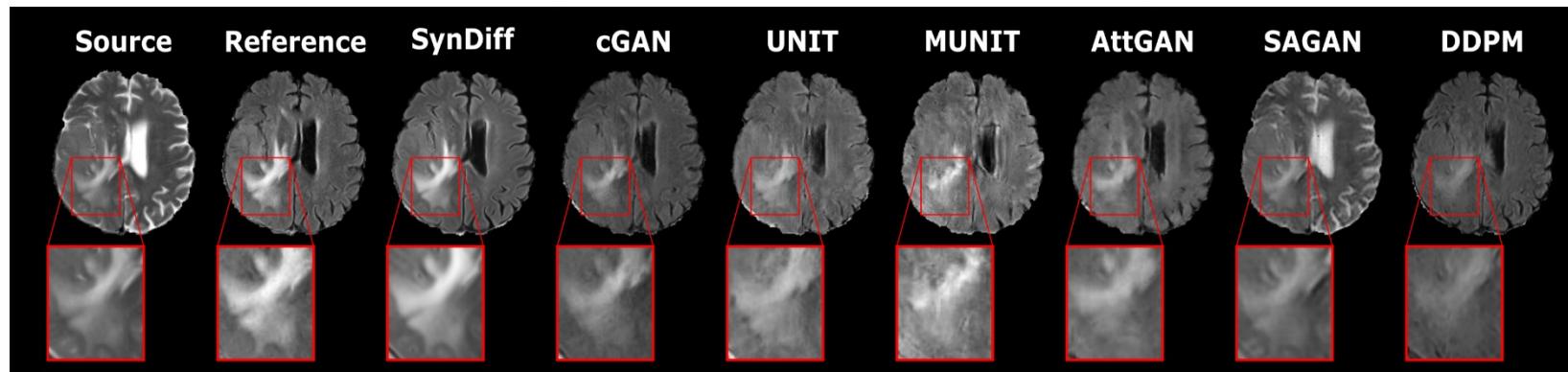
# K-space Cold Diffusion



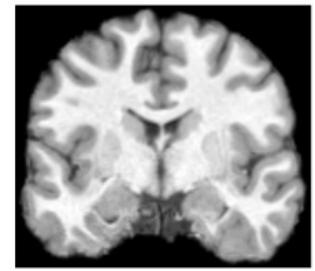
# Image-to-Image Translation

Examples from the Community

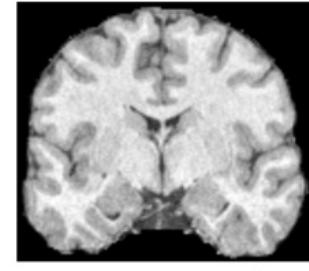
MRI Contrast Translation



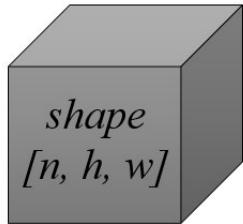
# Diffusion Models for Contrast Harmonization



Bad comparability



3T



for all slices  $b_i$  of  $B$  with  
 $i \in \{1, \dots, n\}$  and shape  $[h, w]$

Scan  $B$  of  $\mathcal{S}$

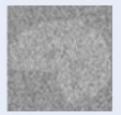


Diffusion Model

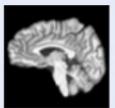
$$\sim \mathcal{N}(0, \mathbf{I})$$



$x_{b_i, T}$



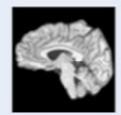
$x_{b_i, t}$



$b_i$

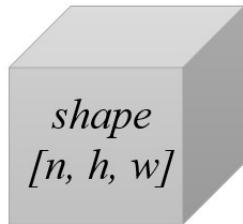


$x_{b_i, t-1}$



$x_{b_i, 0}$

1.5T



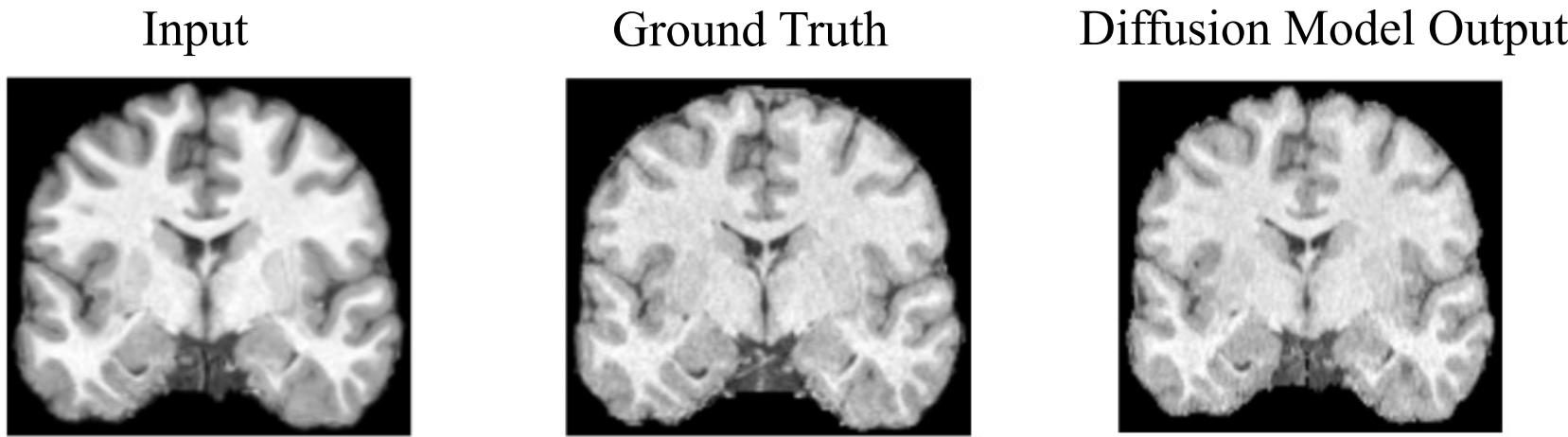
Scan  $B_{transformed}$   
in contrast of  $\mathcal{T}$



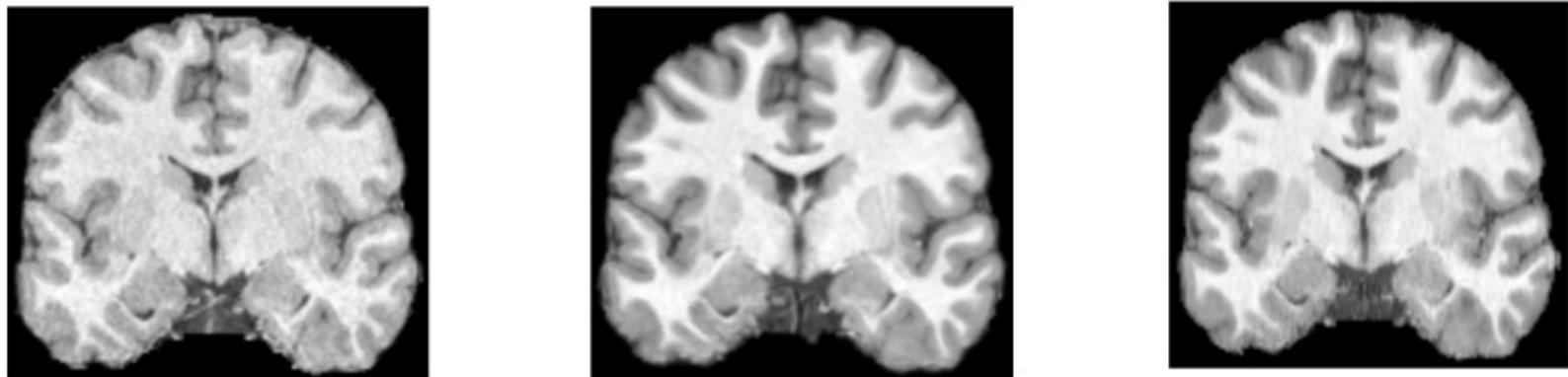
stack all  
slices  $x_{b_i, 0}$  for  
 $i \in \{1, \dots, n\}$

# Contrast Harmonization Results

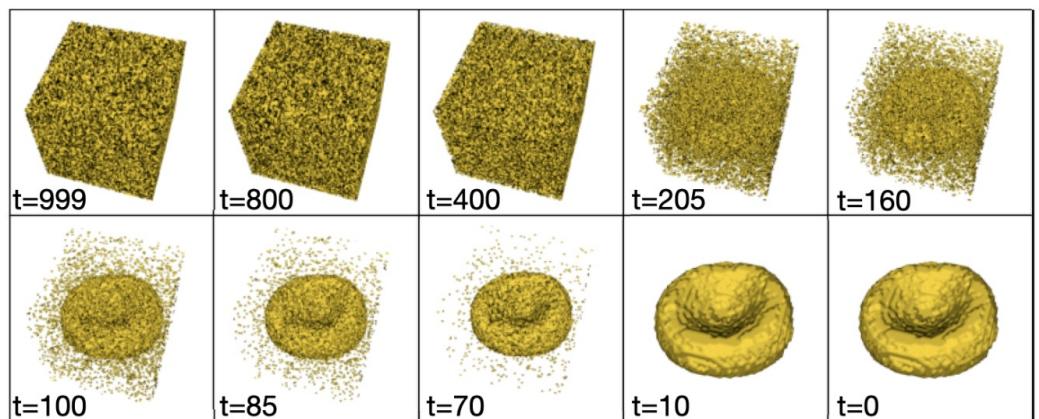
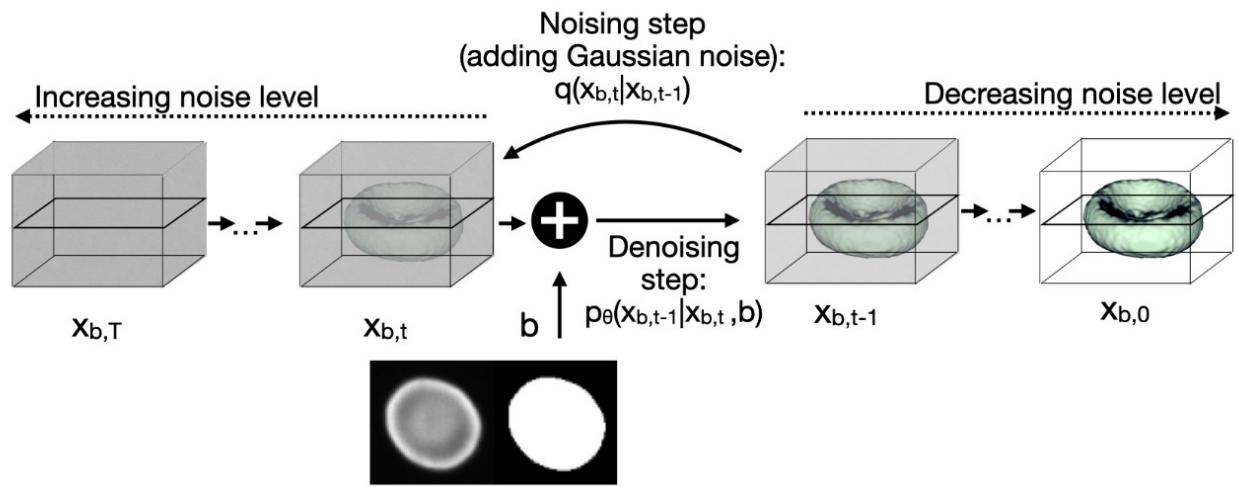
3 T to 1.5 T



1.5 T to 3 T



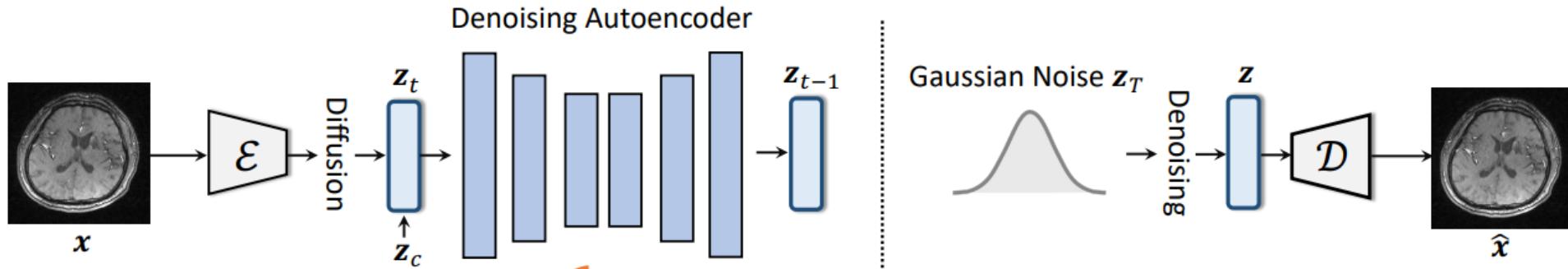
# 3D Shapes from 2D Microscopy



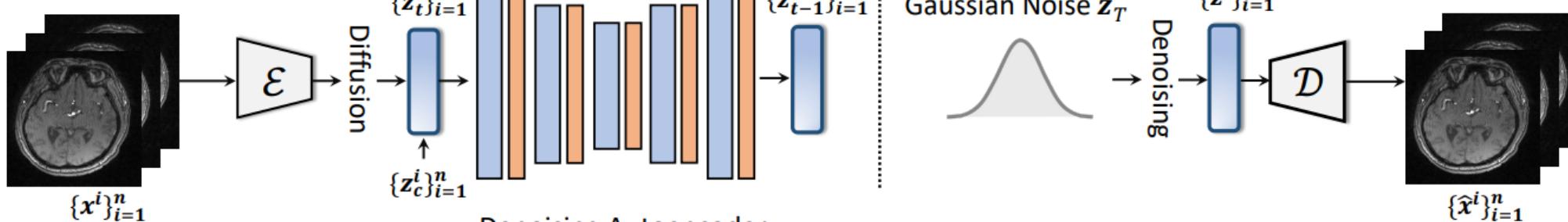
Class	2D input	3D groundtruth	3D predictions					
Spherocytes								
Stromocytes								
Discocytes								
Erythrocytes								
Cell cluster								
Multiblob								
Keratocyte								
Kinocyte								
Acanthocyte								

# 3D with 2D model

## Slice-wise Model



Insert & Quick Fine-tuning

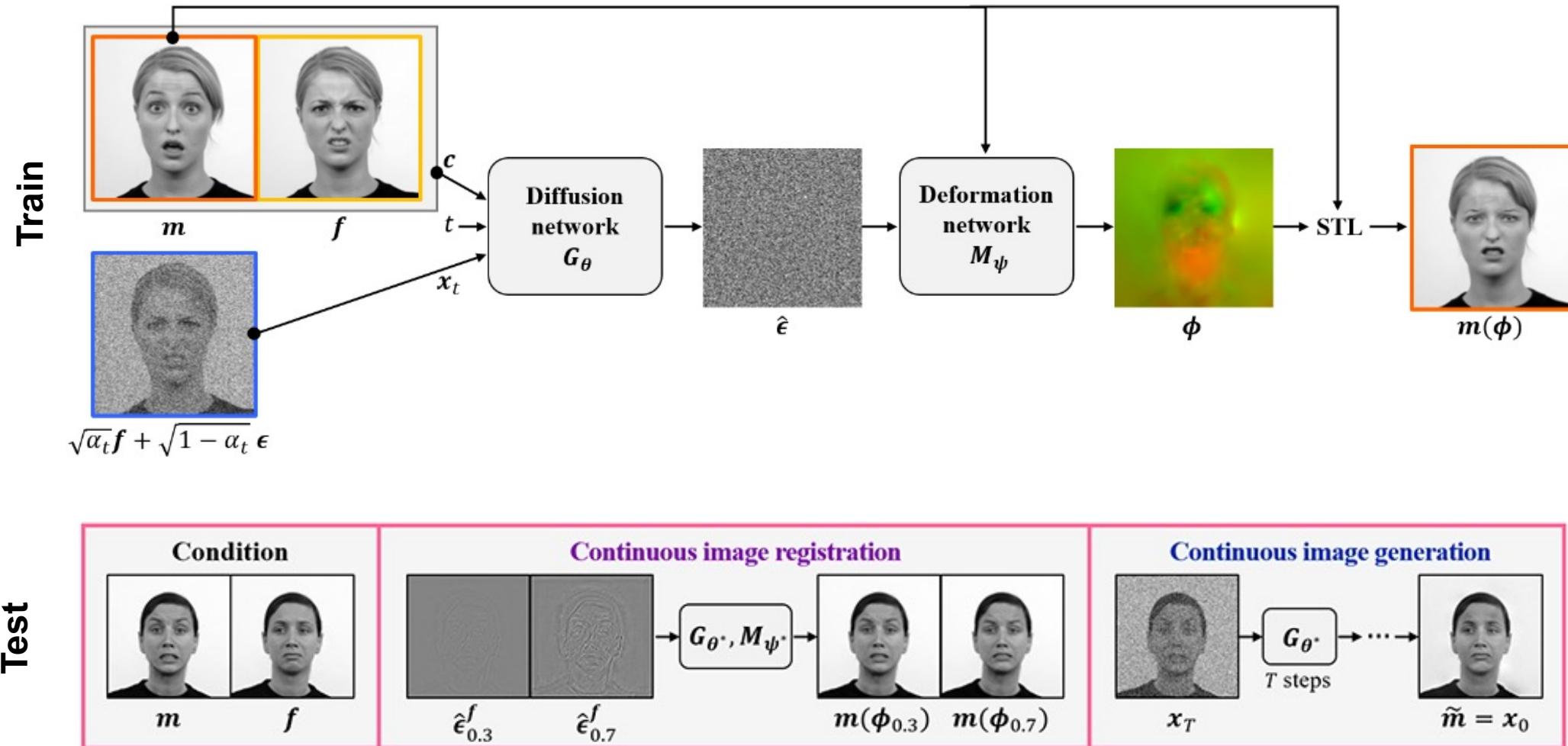


## Volume-wise Model

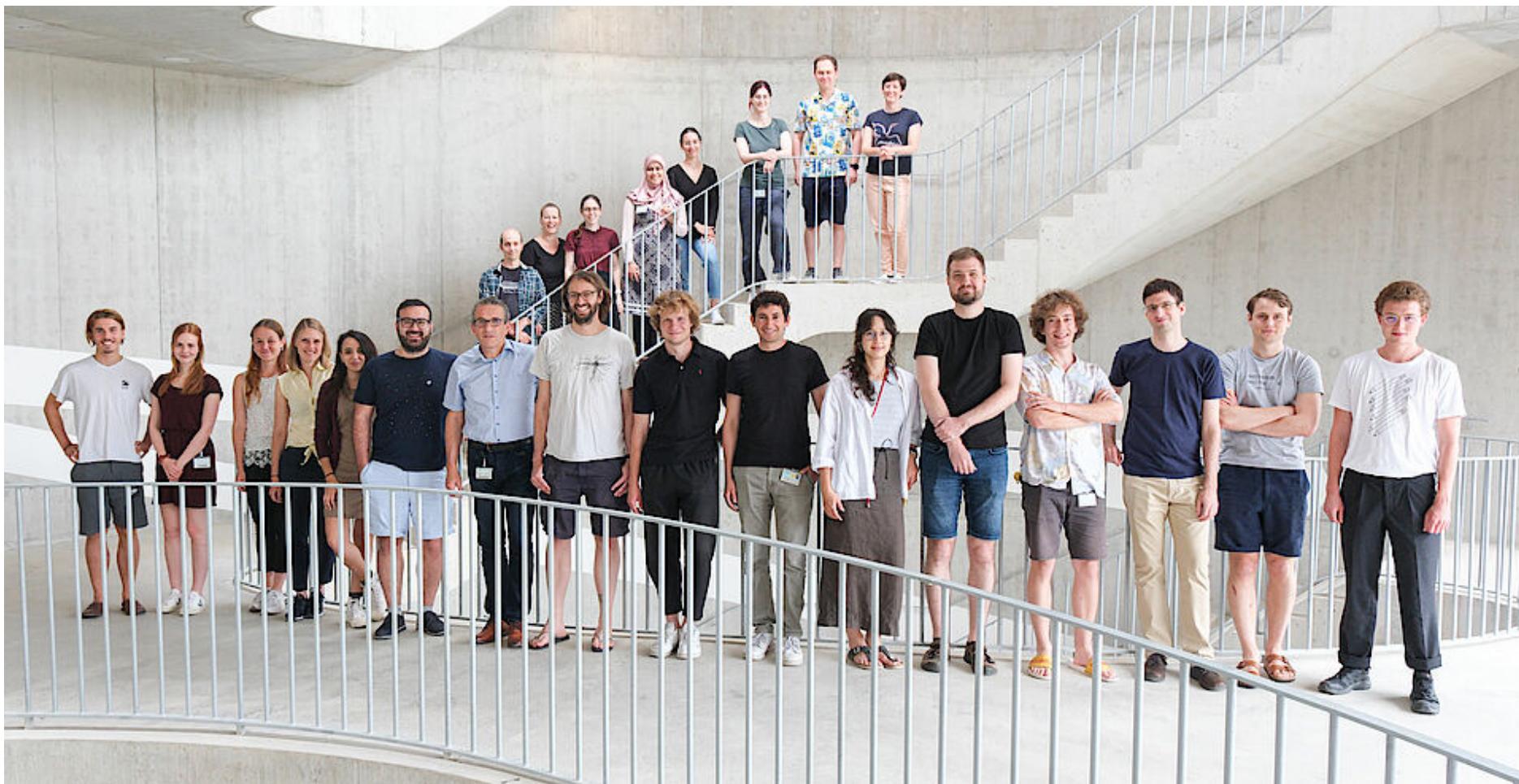
# **Image Registration**

Examples from the Community

# DiffuseMorph



# Thank you!



**Center for medical Image Analysis & Navigation (CIAN), Prof. Philippe C. Cattin**  
Universität Basel

# **Useful Key References, Gits to Watch, etc.**

## **Surveys**

- <https://arxiv.org/abs/2209.02646>
- <https://arxiv.org/abs/2209.00796>

## **Github**

- <https://github.com/heejkoo/Awesome-Diffusion-Models>

## **Tutorials**

- <https://cvpr2022-tutorial-diffusion-models.github.io>
- <https://huggingface.co/blog/annotated-diffusion>
- <https://huggingface.co/docs/diffusers>