

CME 296: Diffusion & Large Vision Models



Afshine Amidi & Shervine Amidi



Teaching staff



Afshine and Shervine

Teaching staff



Afshine

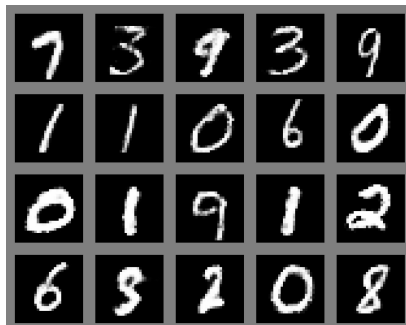
Centrale Paris ('16), MIT ('17)
Uber, Google, Netflix



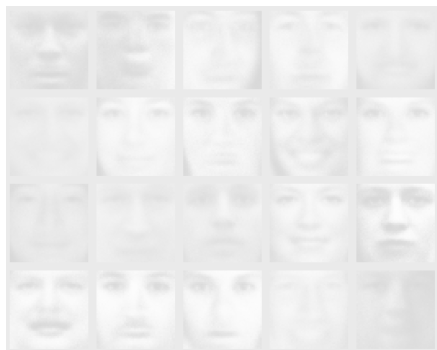
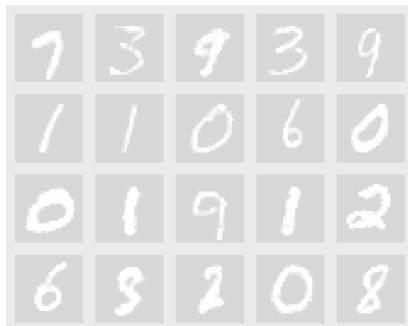
Shervine

Centrale Paris ('16), Stanford ('19)
Uber, Google, Netflix

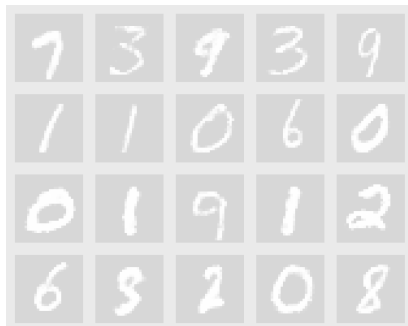
Motivation behind CME 296



Motivation behind CME 296



Motivation behind CME 296



?



CME 296 overview

Goals.

1. Understand **paradigm** behind how we generate images
2. Learn **how** image generation models are **trained & used**

CME 296 overview

Goals.

1. Understand **paradigm** behind how we generate images
2. Learn **how** image generation models are **trained & used**

Audience. Interested in state-of-the-art image generation models

1. Career goal
2. Personal project
3. Curiosity

Prerequisites

- **Linear algebra**

Vectors, matrices, gradient, divergence

Prerequisites

- **Linear algebra**

Vectors, matrices, gradient, divergence

- **Probability theory**

Bayes' rule, conditional/marginal probability, expectation, covariance, Gaussian distribution

Prerequisites

- **Linear algebra**

Vectors, matrices, gradient, divergence

- **Probability theory**

Bayes' rule, conditional/marginal probability, expectation, covariance, Gaussian distribution

- **Differential equations**

Ordinary/stochastic differential equations, numerical solvers

Prerequisites

- **Linear algebra**

Vectors, matrices, gradient, divergence

- **Probability theory**

Bayes' rule, conditional/marginal probability, expectation, covariance, Gaussian distribution

- **Differential equations**

Ordinary/stochastic differential equations, numerical solvers

- **Basics of machine learning**

Loss function, training, inference, neural network

Logistics

Date & time.

- Fridays from 3:30pm to 5:20pm
- Thornton 110

Logistics

Date & time.

- Fridays from 3:30pm to 5:20pm
- Thornton 110


Details about the class.

- 2 units, Letter or Credit/No credit
- Lectures are recorded
- 2 exams ([example](#) for the LLM class)
 - Midterm (50% grade), scheduled on Friday May 1st
 - Final exam (50% grade), scheduled on Monday June 8th

Material

Class website. cme296.stanford.edu

- Contains syllabus & logistics
- Slides and recordings will be posted there


 **CME 296** SYLLABUS **CHEATSHEET** FAQ


CME 296 - Diffusion & Large Vision Models

This course explores diffusion-based generative models for vision. You will study the foundations of diffusion, score matching and flow matching, modern architectures such as U-Nets and Diffusion Transformers, and methods for controllable image generation and evaluation. The course combines theory with practical insights into state-of-the-art generative models. Ideal for students with a background in linear algebra, probability, calculus and machine learning.

[Syllabus](#) [Cheatsheet](#) [Canvas](#)

Instructors

 **Afshine Amidi**
Instructor

 **Shervine Amidi**
Instructor

Logistics

Time: Fridays 3:30pm - 5:20pm
Location: Thornton 110
Grade: Letter or Credit/No Credit

VIP Cheatsheet

CME 296 – DIFFUSION & LARGE VISION MODELS https://cme296.stanford.edu

VIP Cheatsheet:
Diffusion & Large Vision Models

Afshine AMIDI and Shervine AMIDI
April 2, 2026

1 Generation paradigms

1.1 Discrete-time diffusion

Convention: Clean data is $t = 0$ and noisy data is $t = T$.

DDPM – A Denoising Diffusion Probabilistic Model (DDPM) learns to produce data via p_θ by reversing a known Markov diffusion process q that gradually adds noise to clean images.

Starting from a clean image $x_0 \sim p_0$, we obtain noisier images x_{t+1} using the current x_t with:

$$x_{t+1} = \sqrt{1 - \beta_t} x_t + \sqrt{\beta_t} \epsilon$$

with noise schedule β_t which can be rewritten as:

$$x_t = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon$$

with $\epsilon \sim \mathcal{N}(0, 1)$, $\alpha_t = 1 - \beta_t$, $\alpha_T = \prod_{s=1}^T \alpha_s$

KL divergence – Kullback-Leibler divergence (KL divergence) is a measure that compares two probability distributions p and q :

$$\text{KL}[p||q] = \int \log \left(\frac{p(x)}{q(x)} \right) p(x) dx = \mathbb{E}_p \left[\log \left(\frac{p}{q} \right) \right]$$

Variational formulation – Using the Evidence Lower Bound, we derive a tractable loss:

$$\mathbb{E}_{x_0}[\log(p_\theta(x_0))] \geq \text{ELBO} = - \sum_{s=1}^T \text{KL}(q(x_s | x_{s-1}, x_0) || p_\theta(x_s | x_{s-1})) + \text{extra terms}$$

Derivation: Use Jensen's inequality and simplify loss to the KL divergence term.

Training – The DDPM loss estimates the amount of noise ϵ to remove from the noisy image $x_t = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon$ at time t by minimizing the L_2 loss with respect to the actual noise removed ϵ :

$$\mathcal{L}_{\text{DDPM}} = \mathbb{E}_{x_0, \epsilon, t} [||\epsilon_\theta(x_t, t) - \epsilon||^2]$$

Sampling – We sample a less noisy image \hat{x}_{t-1} using the noised image x_t and the predicted noise $\epsilon_\theta(x_t, t)$ to remove between x_t and x_0 :

$$\hat{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(x_t, t) \right) + \sigma_t$$

DDIM – A Denoising Diffusion Implicit Model (DDIM) is a deterministic sampling method to generate new samples faster by dropping the Markov assumption.

Derivation: Use convenient marginal distribution that has the same DDPM loss.

1.2 Score-based diffusion

Convention: Clean data is $t = 0$ and noisy data is $t = T$.

Score – The score function of a distribution is equal to $\nabla_x \log(p(x))$ and indicates the direction where the log-likelihood is the highest.

Remark: Compared to $\nabla_x p(x)$, the score is easier to compute, appears naturally in mathematical derivations and is more numerically stable, which is why it is used in practice.

DDSM – Denoising Score Matching (DSM) is a training objective that learns the score function by predicting the gradient of the log data density from noisy inputs $x_t = x_0 + \sigma_t \epsilon$.

Training – Noise Conditional Score Networks (NCSN) trained with Score Matching with Langevin Dynamics (SMLD) estimates the score $\nabla_x \log(p_{\sigma_t}(x))$ at multiple noise levels σ_t :

$$\mathcal{L}_{\text{NCSN}} = \sum_{t=1}^T \lambda_t \mathbb{E}_x [||s_\theta(x, \sigma_t) - \nabla_x \log p_{\sigma_t}(x)||^2]$$

Sampling – Anisotropic Langevin Dynamics (ALD) is a sampling method that generates data while gradually reducing to progressively refine samples toward the data distribution:

$$x_t = x_{t-1} + \frac{\alpha_t}{2} s_\theta(x_{t-1}, \sigma_t) + \sqrt{\alpha_t} \epsilon_t$$

SDE formulation – The continuous formulation describes the forward diffusion process over time as a Stochastic Differential Equation (SDE) driven by a Wiener process, with time-dependent drift and diffusion terms that progressively perturb data with noise.

STANFORD UNIVERSITY 1 SPRING 2026

Link to PDF: <https://github.com/afshinea/stanford-cme-296-diffusion-large-vision-models>

Class communications

Canvas.

- Announcements
- Class discussions via Ed

Any other general inquiries / questions, email us:

- cme296-spr2526-staff@lists.stanford.edu
- afshine@stanford.edu, shervine@stanford.edu

Example of a slide in this class

Explanation of a CME 296 concept

Example of a slide in this class

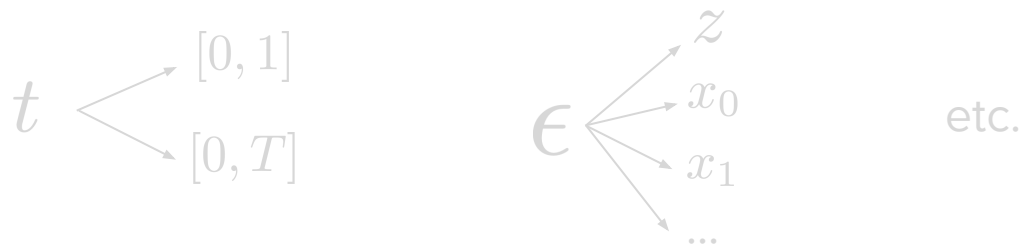
Explanation of a CME 296 concept

Source & suggested reading, if interested



Important note on conventions

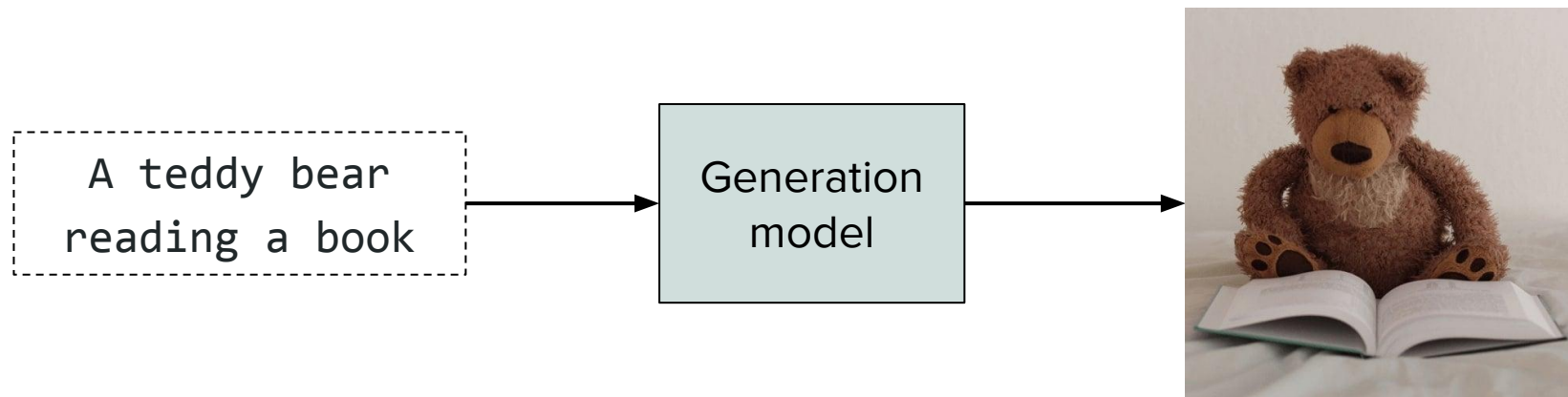
⚠ Field is full of different notations / conventions ⚠



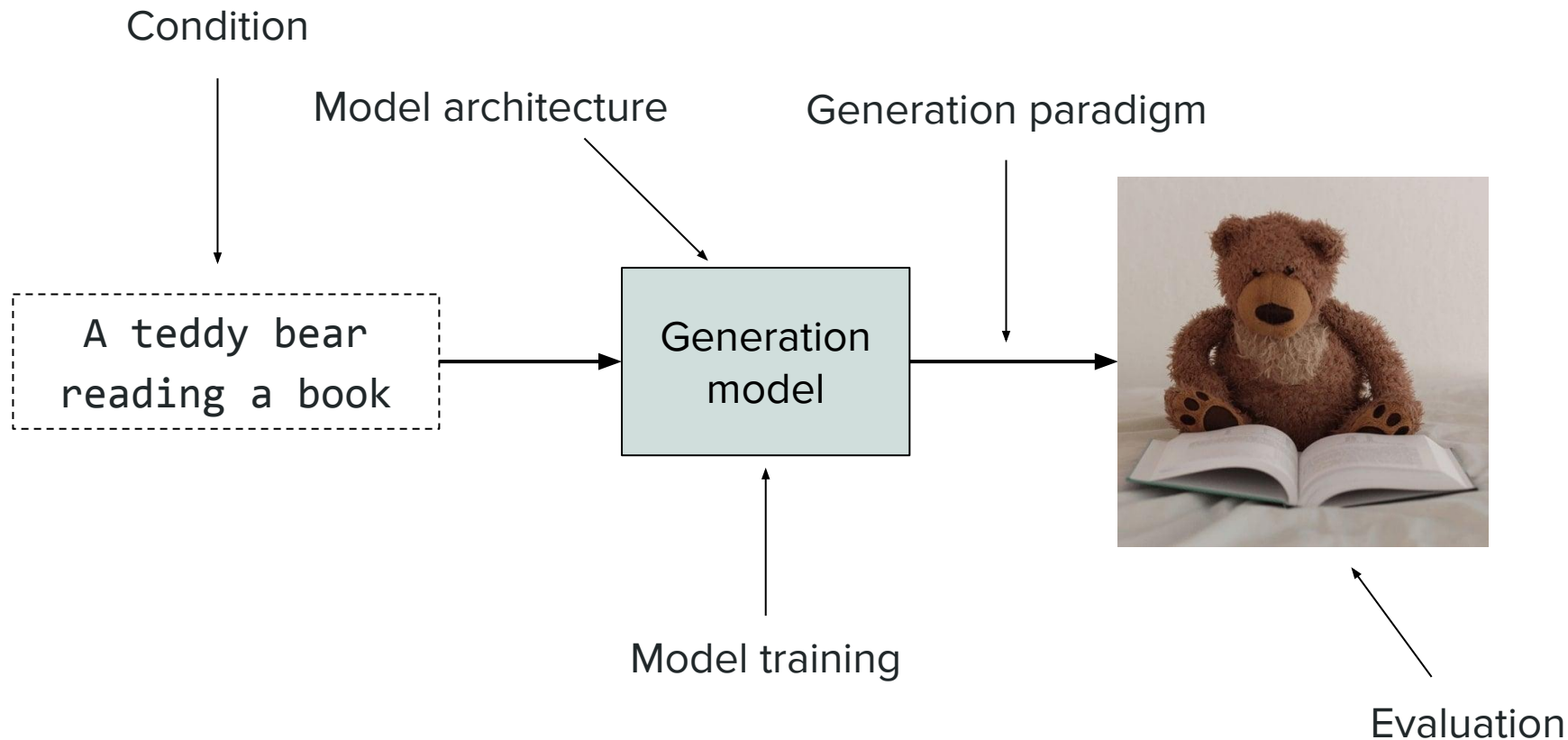
CME 296 will (try to):

1. Rely on the **most commonly-used** notations
2. **Stay consistent** throughout the class when that makes sense
3. **Call out** unintuitive change in notations

Outline of this class



Outline of this class

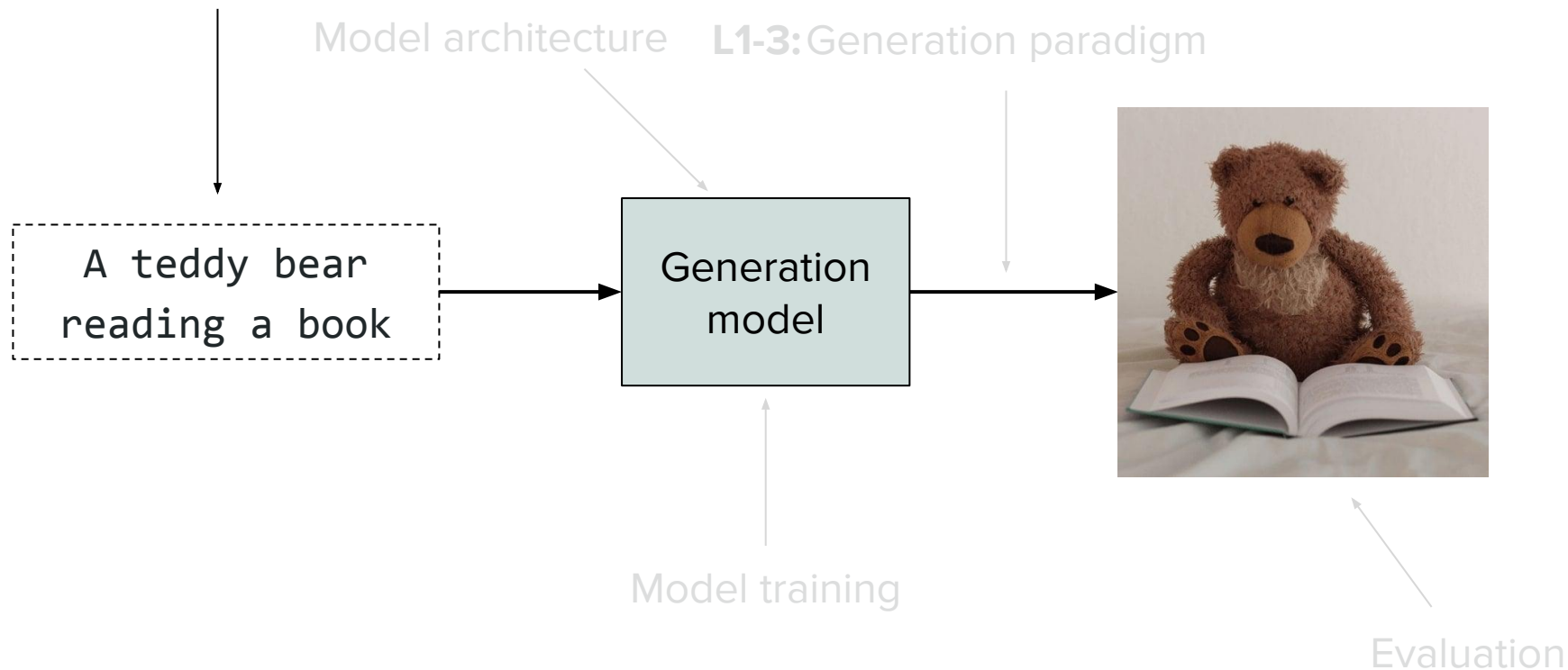


Outline of this class

L4: Condition

Model architecture

L1-3: Generation paradigm

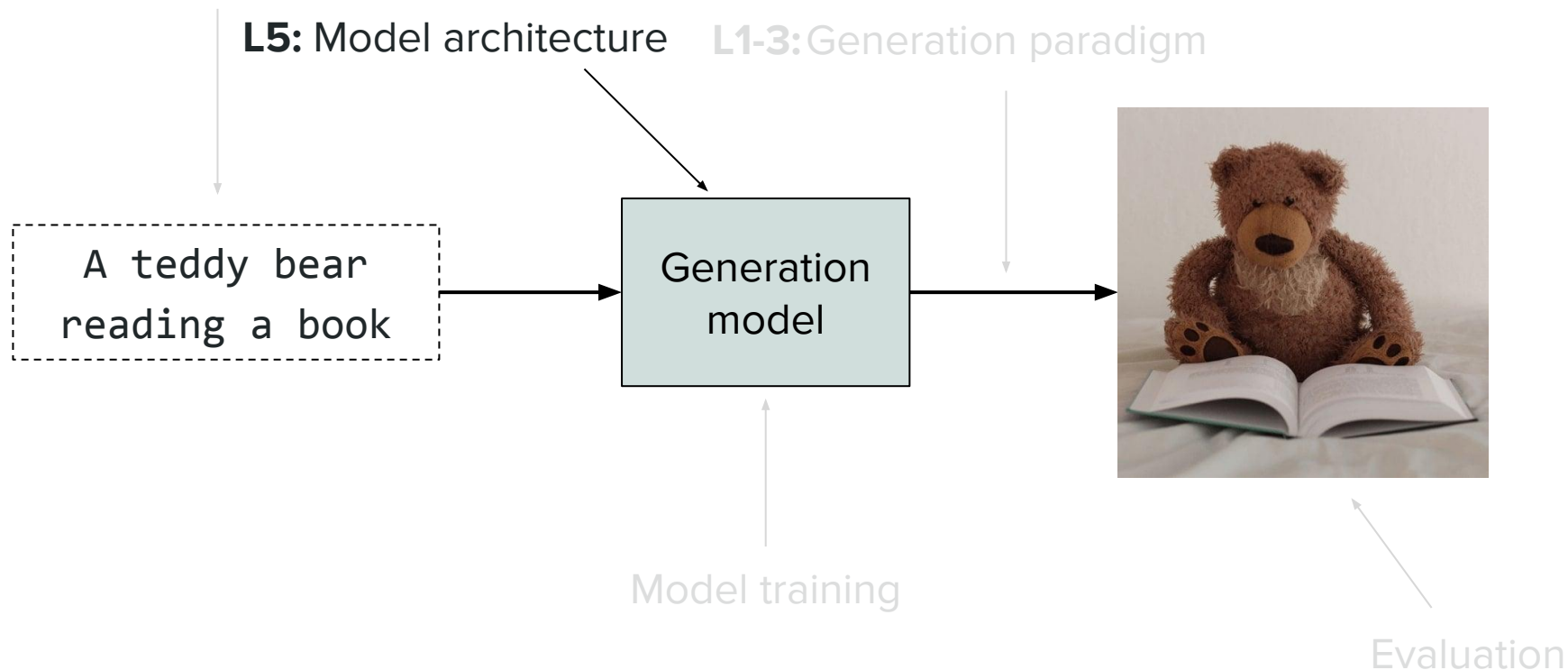


Outline of this class

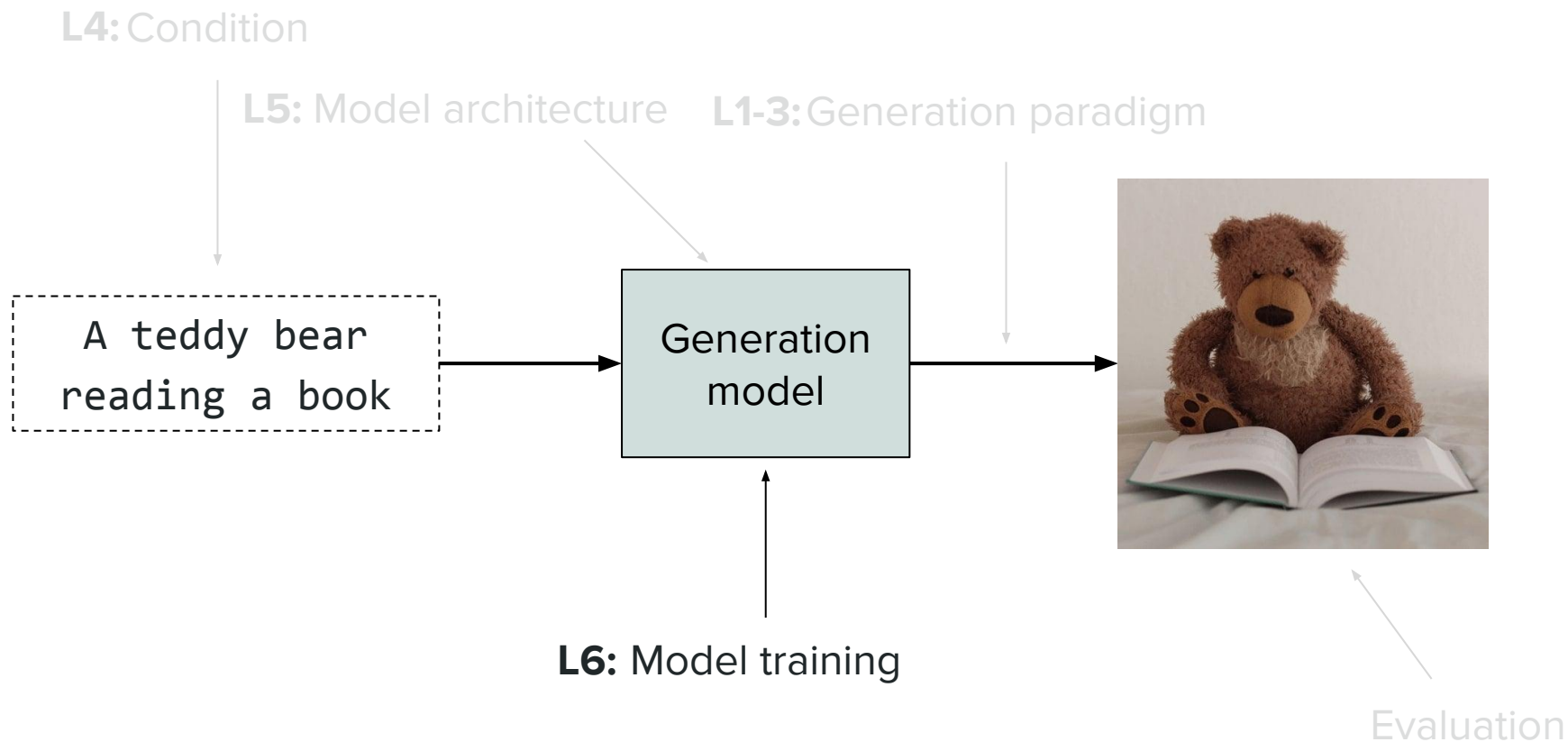
L4: Condition

L5: Model architecture

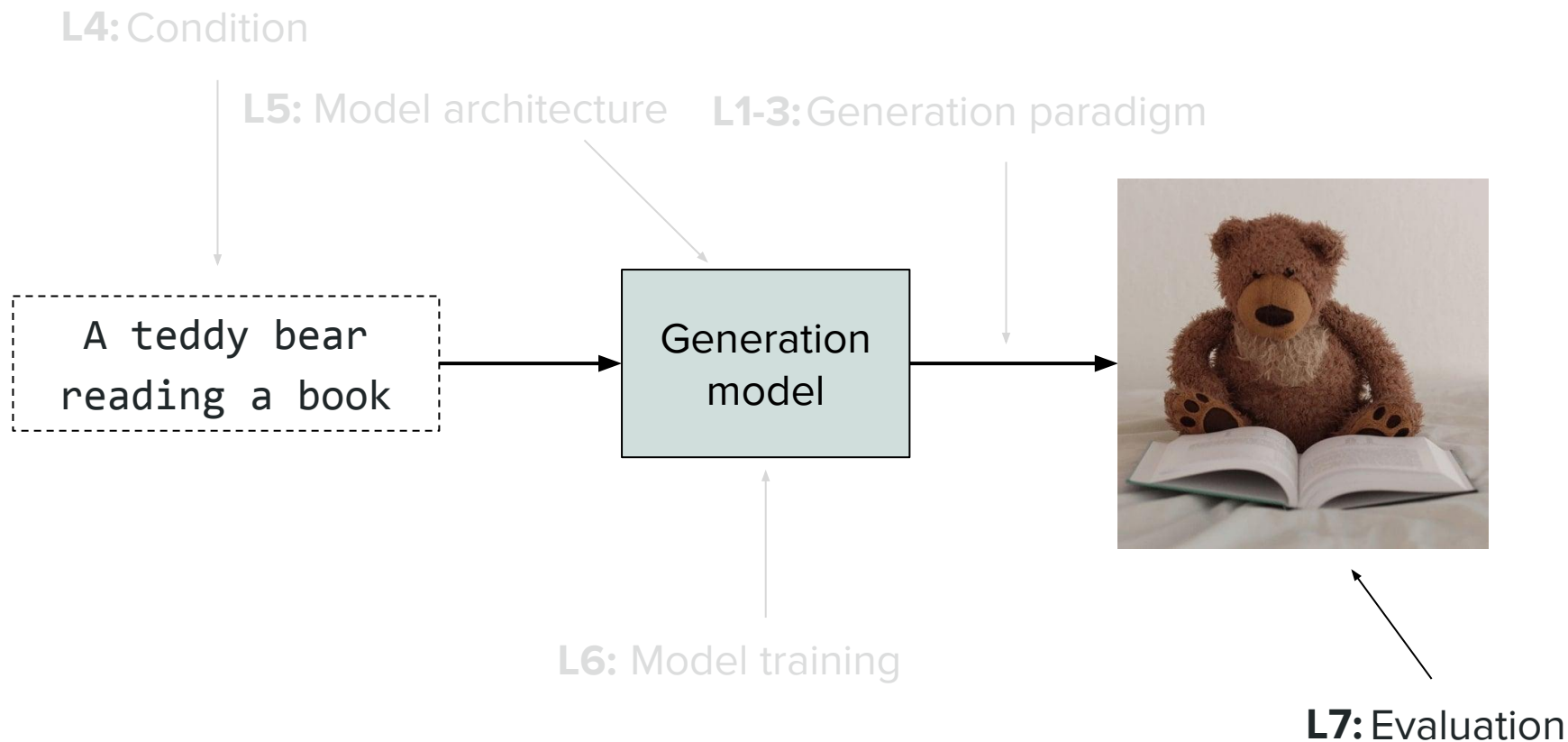
L1-3: Generation paradigm



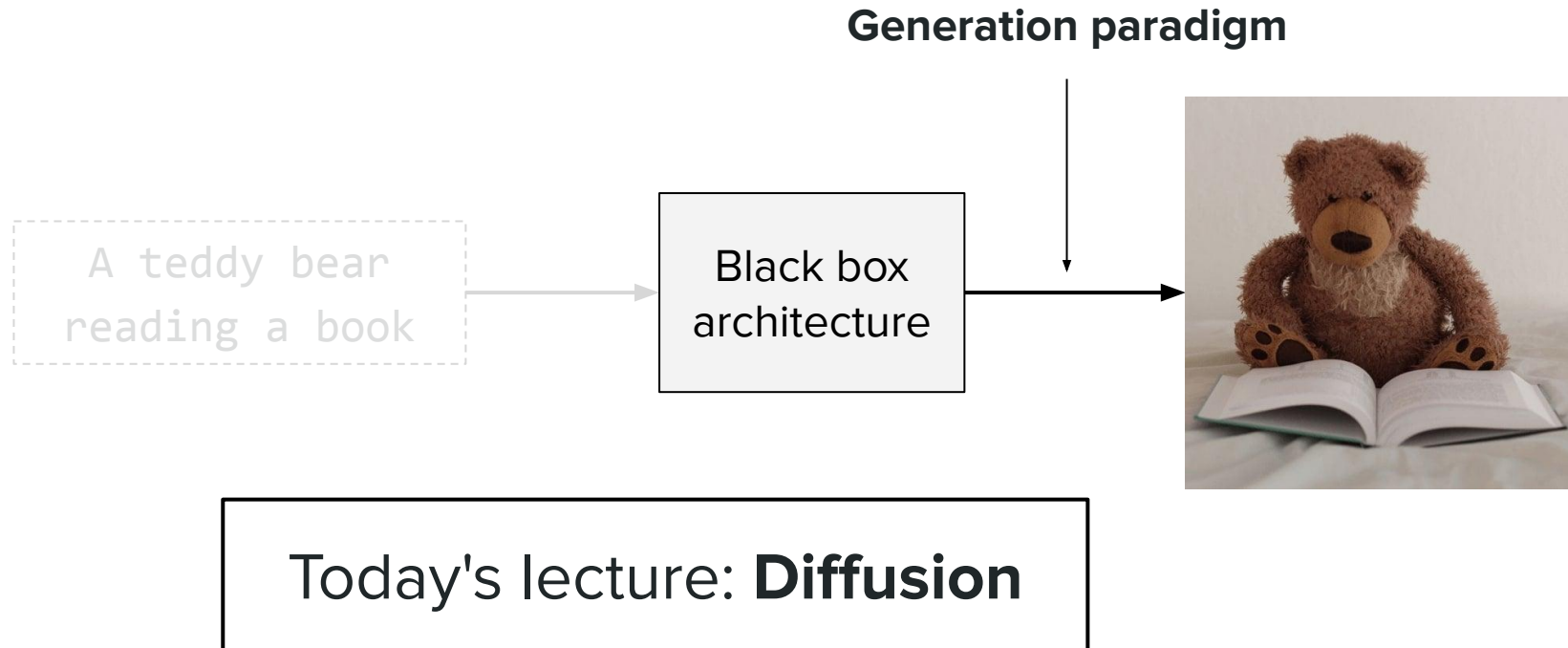
Outline of this class



Outline of this class



First three lectures: generation paradigms





Diffusion & Large Vision Models

Motivation

Forward / backwards process

Variational formulation

Training

Inference

Faster sampling

Problem formulation

Suppose we are given a **sample** of observations from p_{data}



Objective. Generate new images from a distribution of images p_{data}

Problem formulation

Suppose we are given a **sample** of observations from p_{data}

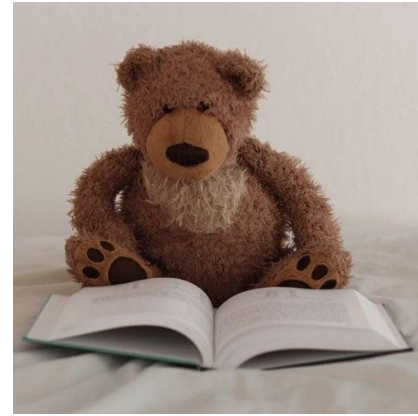
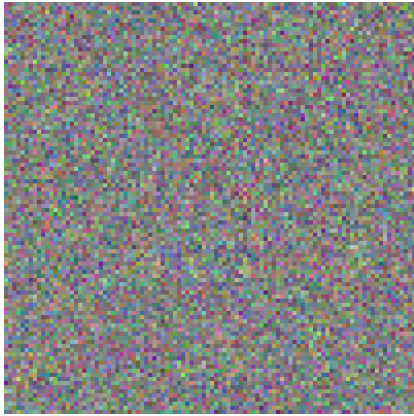


Objective. Generate new images from a distribution of images p_{data}



Unconditioned ← Focus of the 3 first lectures

Strategy: generate from noise



Why start from noise?

- Noise is **easy** to sample
- Introduces **randomness**
- Gaussian distribution has **nice properties**

Analogy: building a sculpture



The sculpture is already complete within the marble block, before I start my work. It is already there, I just have to chisel away the superfluous material.

–Michelangelo



Diffusion & Large Vision Models

Motivation

Forward / backwards process

Variational formulation

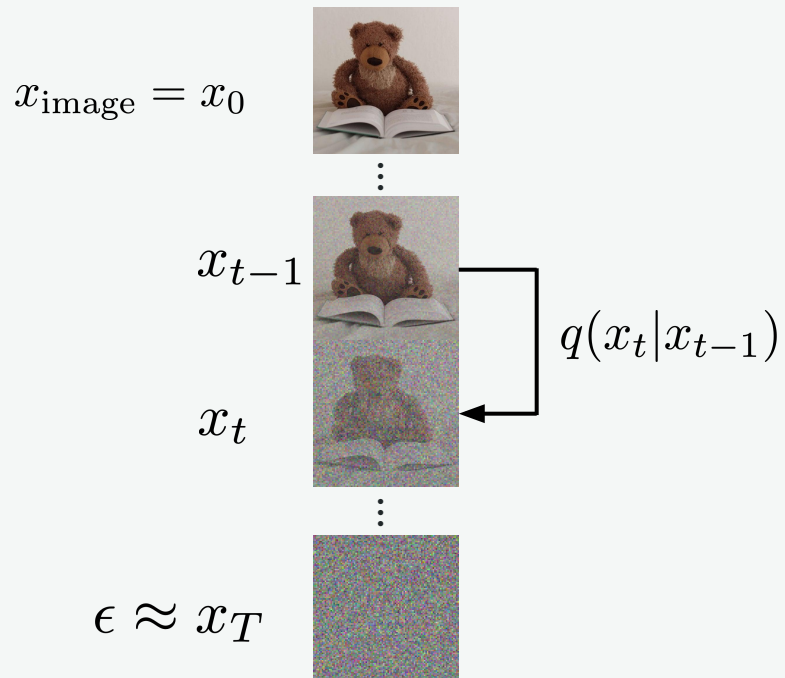
Training

Inference

Faster sampling

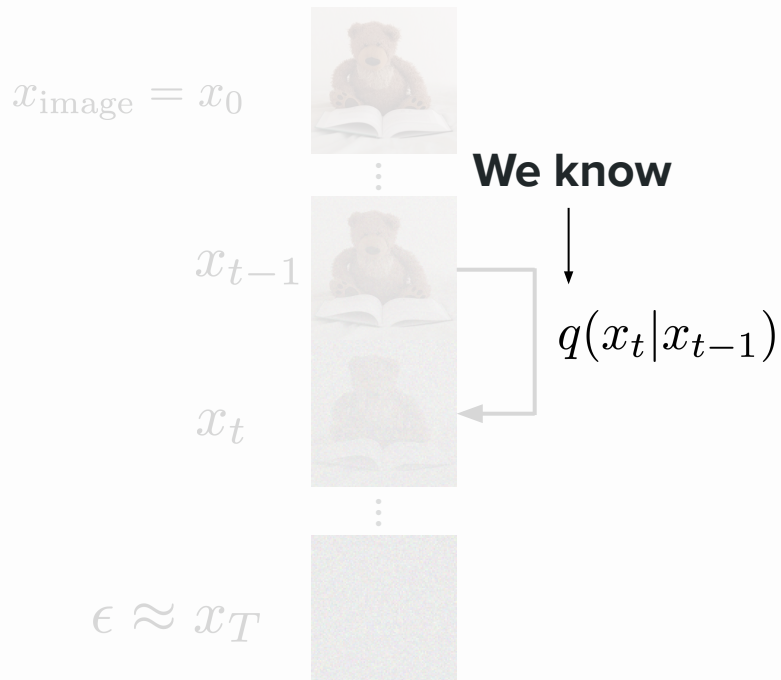
Intuition behind diffusion

1. Add noise (forward process)



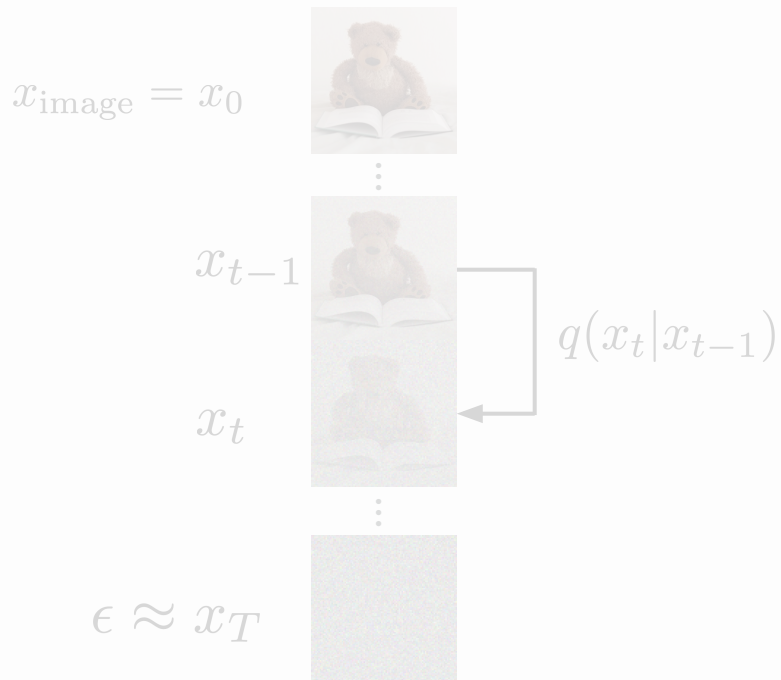
Intuition behind diffusion

1. Add noise (forward process)

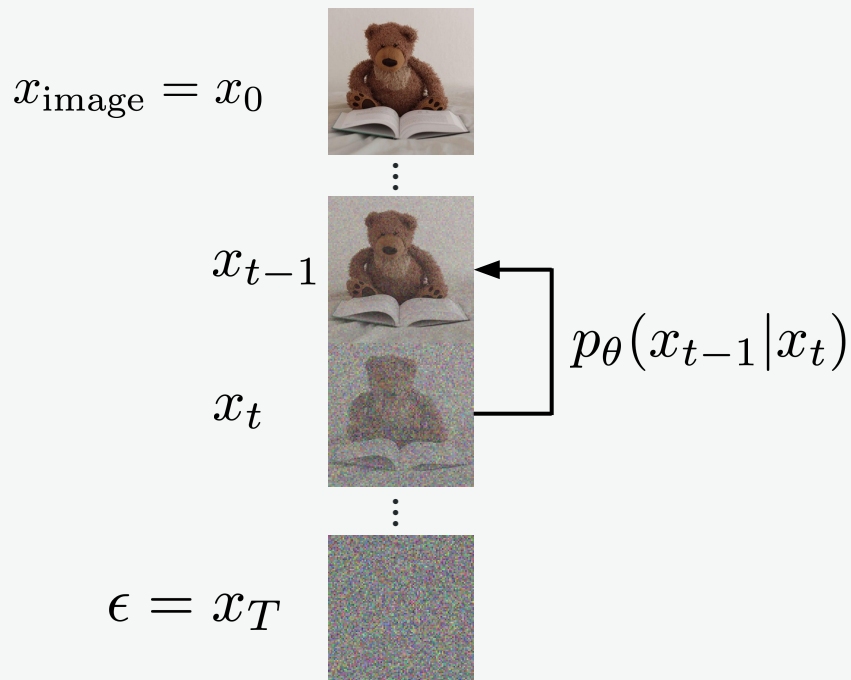


Intuition behind diffusion

1. Add noise (forward process)

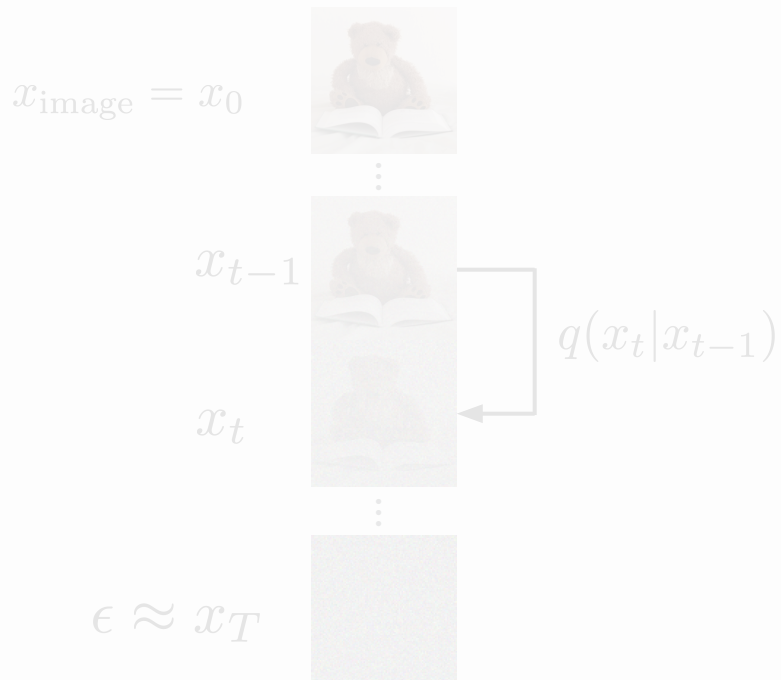


2. Learn to denoise it (reverse process)



Intuition behind diffusion

1. Add noise (forward process)



2. Learn to denoise it (reverse process)

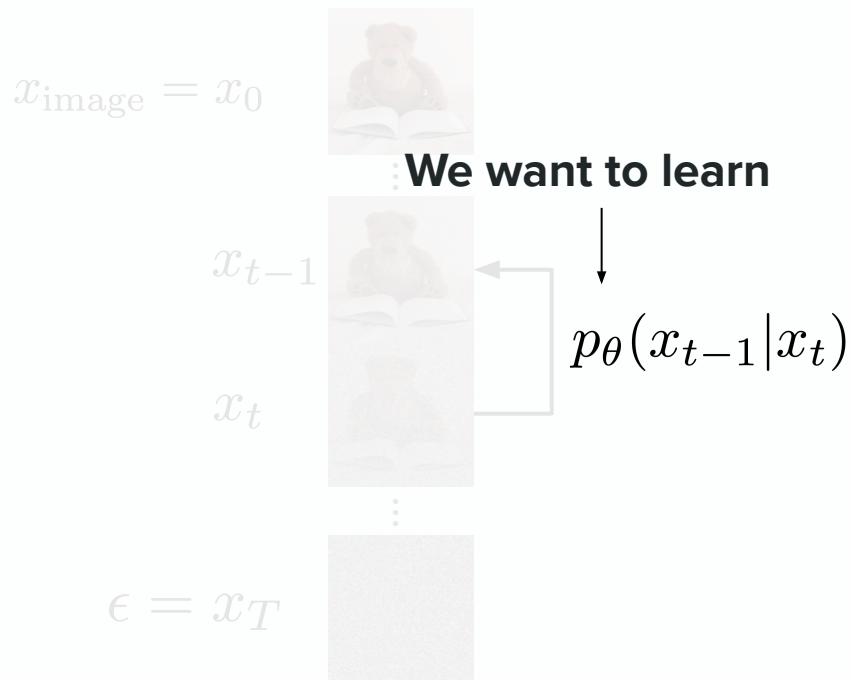
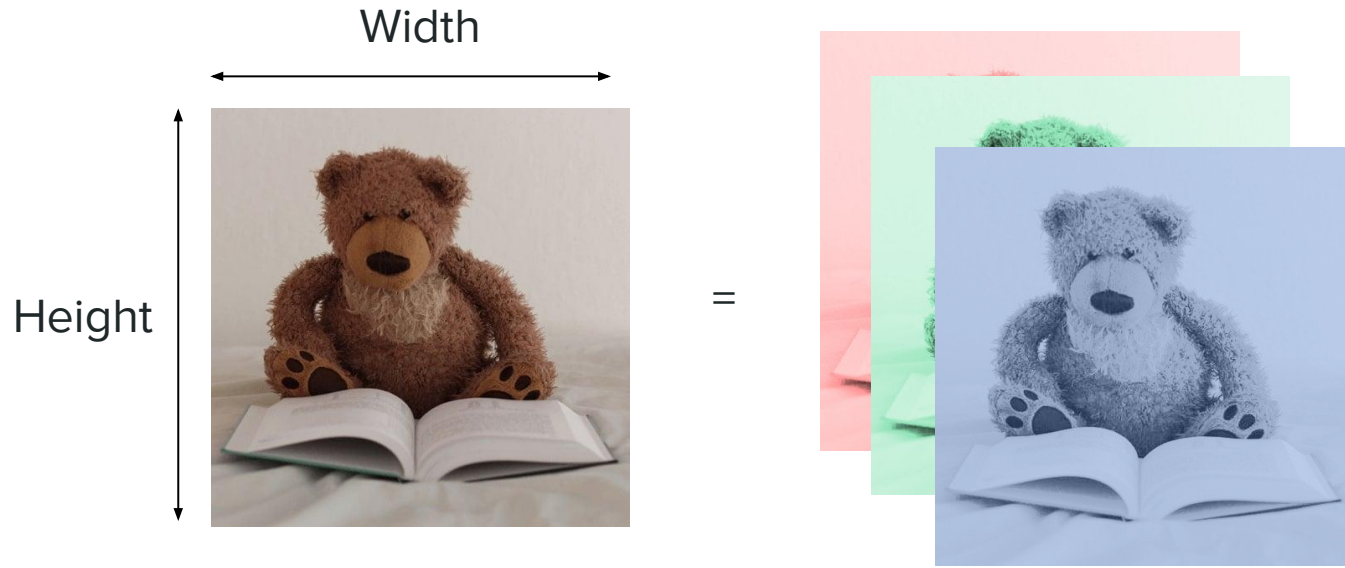
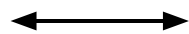
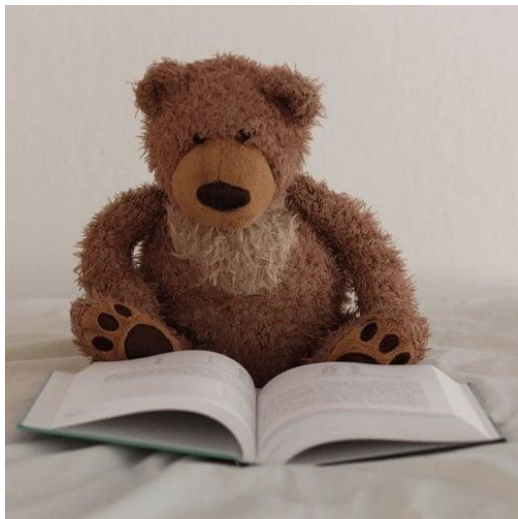


Image representation



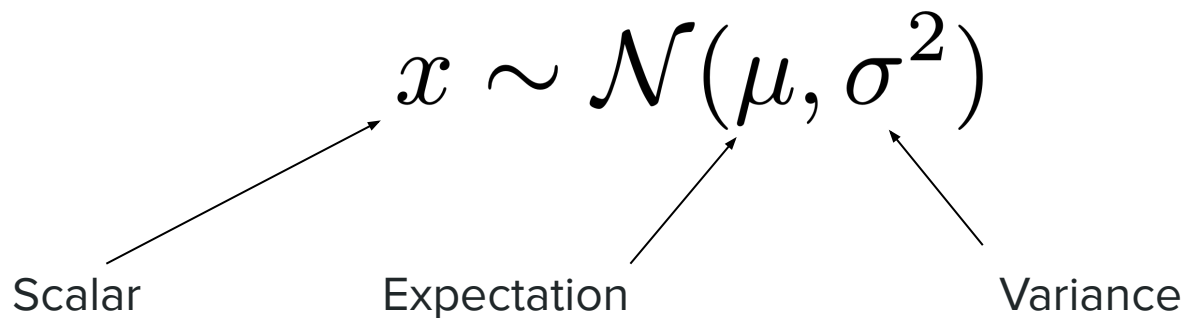
R **G** **B**
pixel = (0 to 255, 0 to 255, 0 to 255)

Image representation

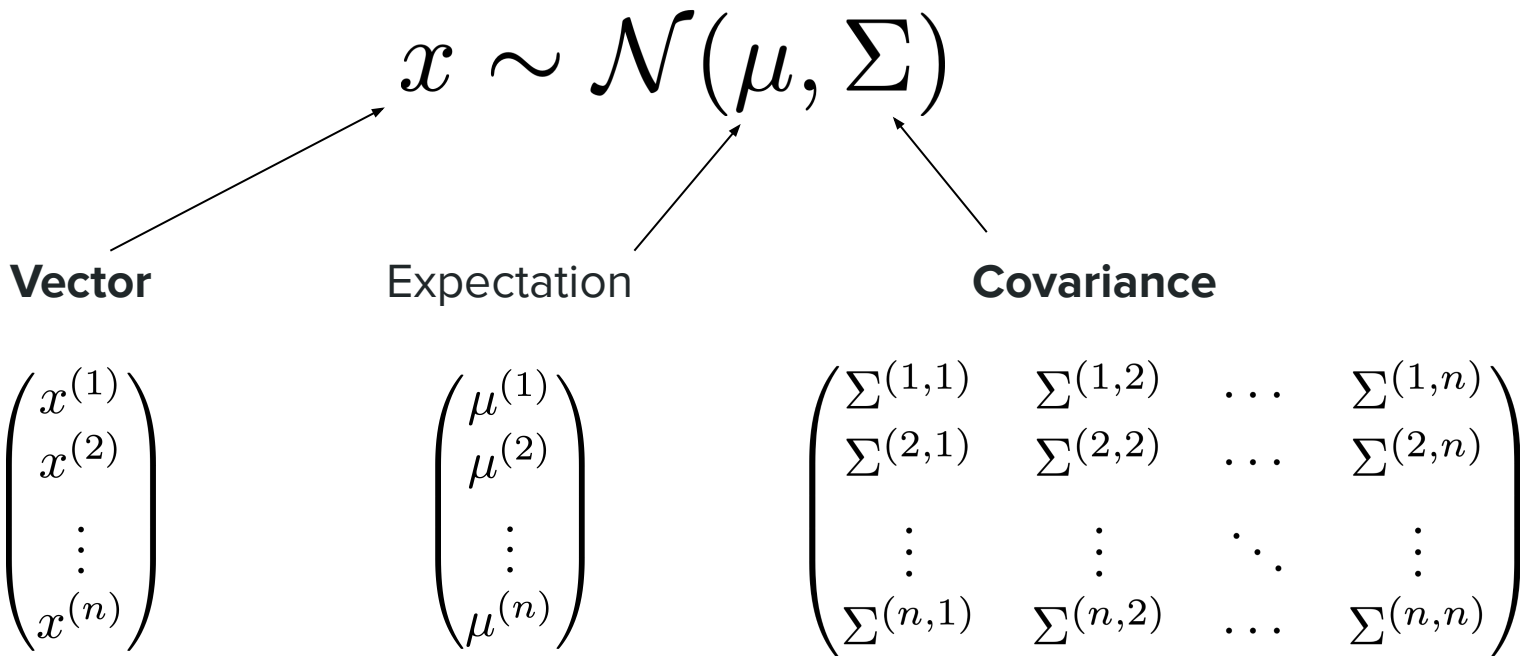


$$x_{\text{image}} = \begin{pmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{pmatrix}$$

Probability distribution in 1-D space



Probability distribution in high-dim space



Meaning of probability distribution in high-dim space

$$x \sim \mathcal{N}(\mu, \Sigma)$$

$$x^{(i)} \sim \mathcal{N}(\mu^{(i)}, \Sigma^{(i,i)}) \quad \text{and} \quad \text{Cov}(x^{(i)}, x^{(j)}) = \Sigma^{(i,j)}$$

Good news. We will mainly (if not only) be dealing with **isotropic** Gaussians:

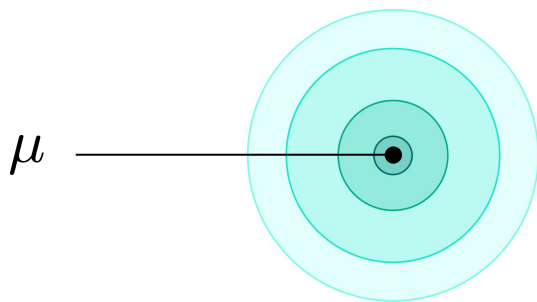
$$\Sigma = \sigma^2 I = \begin{pmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{pmatrix}$$

Meaning of probability distribution in high-dim space

$$x \sim \mathcal{N}(\mu, \Sigma)$$

$$x^{(i)} \sim \mathcal{N}(\mu^{(i)}, \Sigma^{(i,i)}) \quad \text{and} \quad \text{Cov}(x^{(i)}, x^{(j)}) = \Sigma^{(i,j)}$$

Good news. We will mainly (if not only) be dealing with **isotropic** Gaussians:



same variance
in all directions

Meaning of probability distribution in high-dim space

$$x \sim \mathcal{N}(\mu, \Sigma)$$

$$x^{(i)} \sim \mathcal{N}(\mu^{(i)}, \Sigma^{(i,i)}) \quad \text{and} \quad \text{Cov}(x^{(i)}, x^{(j)}) = \Sigma^{(i,j)}$$

Good news. We will mainly (if not only) be dealing with **isotropic** Gaussians:

$$p(x) = \frac{1}{(2\pi\sigma^2)^{d/2}} \exp\left(-\frac{\|x - \mu\|^2}{2\sigma^2}\right)$$

probability density
is known and "easy"

Noise representation

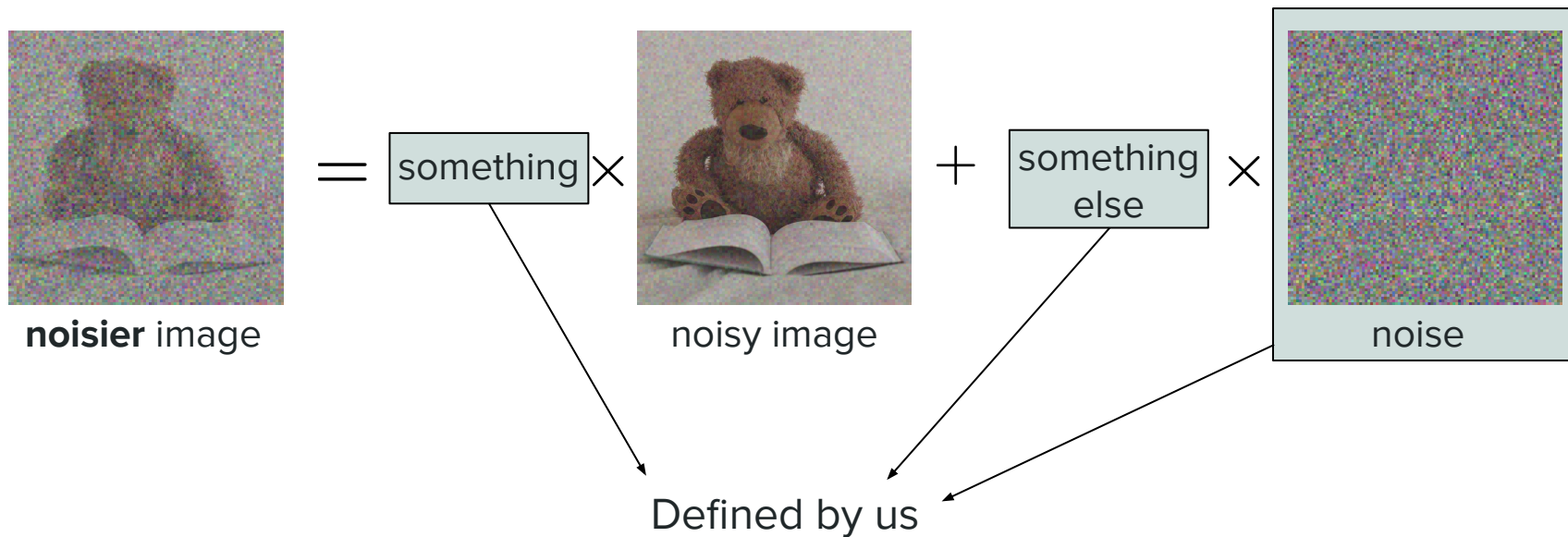


$$\longleftrightarrow \epsilon = \begin{pmatrix} \epsilon^{(1)} \\ \epsilon^{(2)} \\ \vdots \\ \epsilon^{(n)} \end{pmatrix} \sim \mathcal{N}(0, I)$$

In other words, $\epsilon^{(i)} \sim \mathcal{N}(0, 1)$ independent and identically distributed

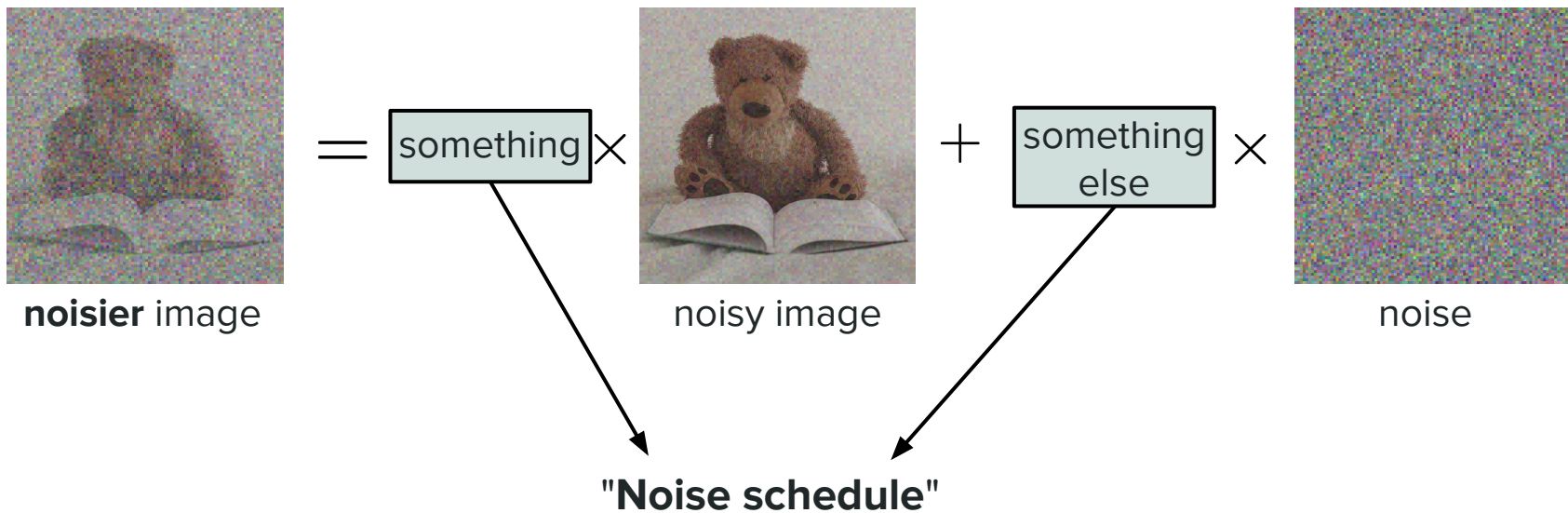
Process behind this

q is **chosen** by us and its **distribution** is **known** (reminder: it is stochastic!)



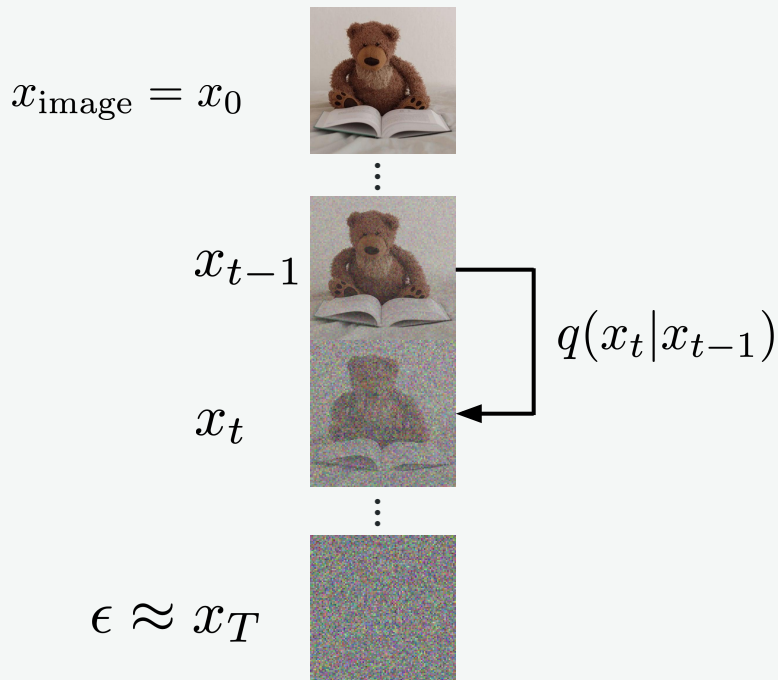
Process behind this

q is **chosen** by us and its **distribution** is **known** (reminder: it is stochastic!)



Mathematical formalism of forward process

1. Add noise (forward process)



Noisier image

Noise

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

Previous image

with $0 \leq \beta_1 < \beta_2 < \dots < \beta_T \leq 1$
noise schedule

Mathematical formalism of forward process

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



$$q(x_t | x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) I)$$

$$\text{with } \bar{\alpha}_t = \prod_{i=1}^t \alpha_i \quad \text{and} \quad \alpha_t = 1 - \beta_t$$

Derivation: Variance of sum of independent Gaussians is sum of respective variances.



Diffusion & Large Vision Models

Motivation

Forward / backwards process

Variational formulation

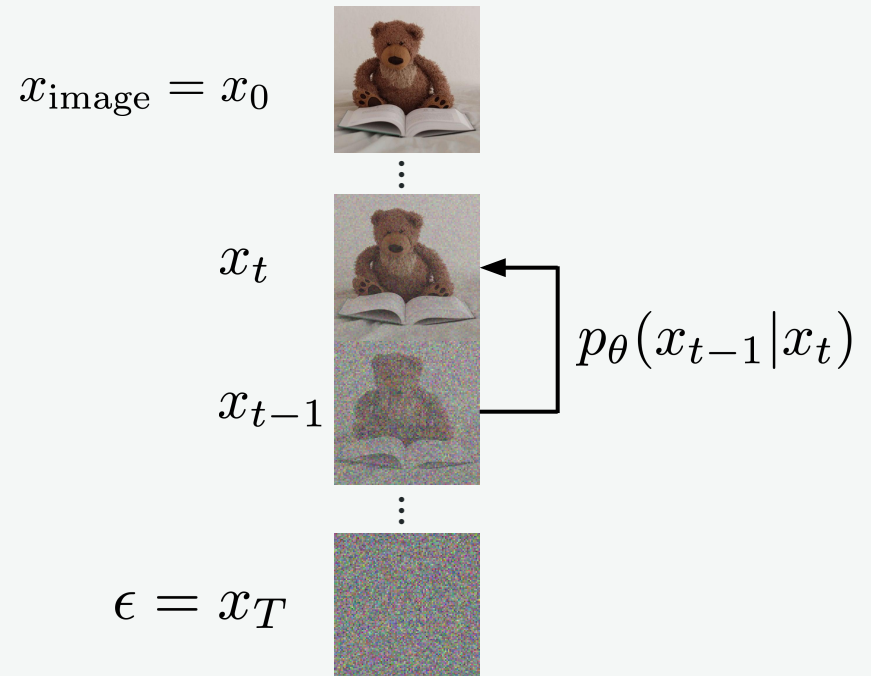
Training

Inference

Faster sampling

Learn reverse process

2. Learn to **denoise** it (reverse process)

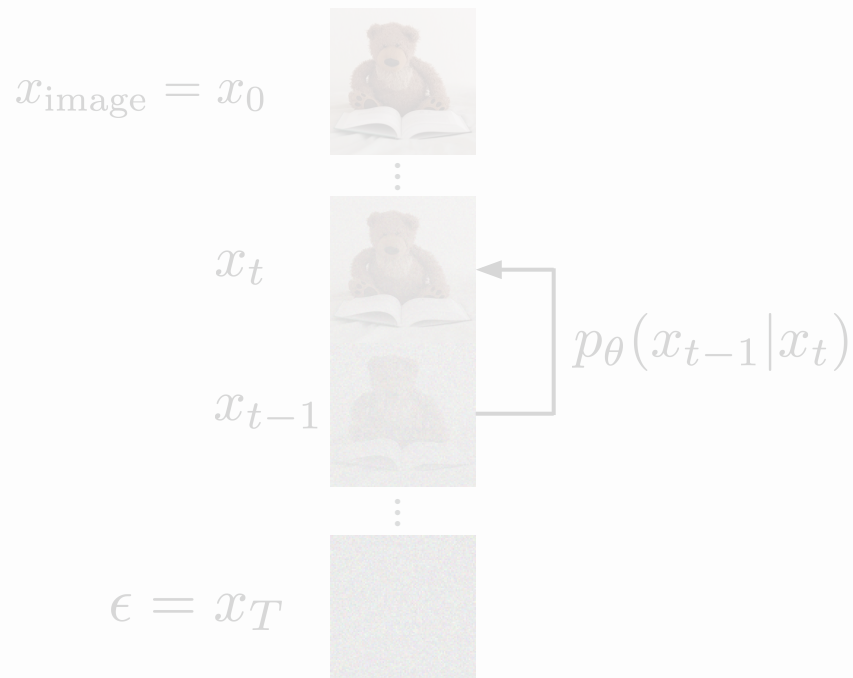


Learn reverse process

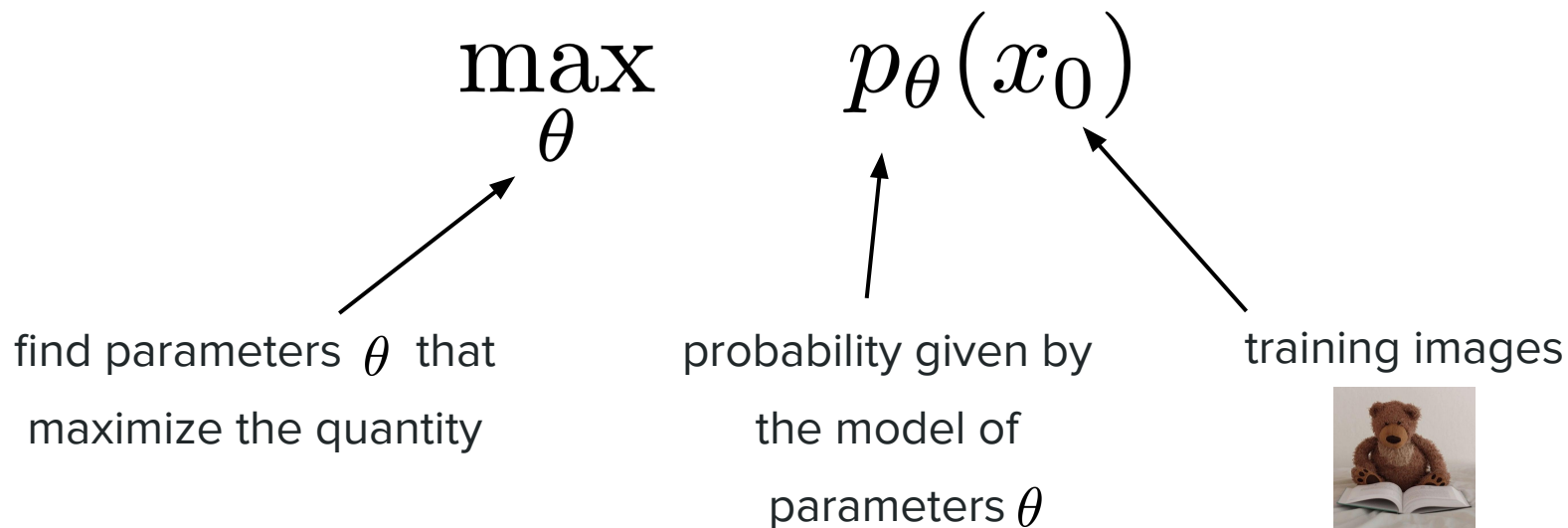
Objective. We want to learn:

$$p_{\theta}$$

2. Learn to denoise it (reverse process)



Objective



Objective

computationally stable +
nice properties from log

$$\max_{\theta} \log p_{\theta}(x_0)$$

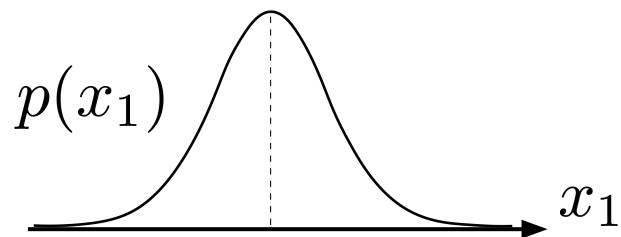
find parameters θ that
maximize the quantity

probability given by
the model of
parameters θ

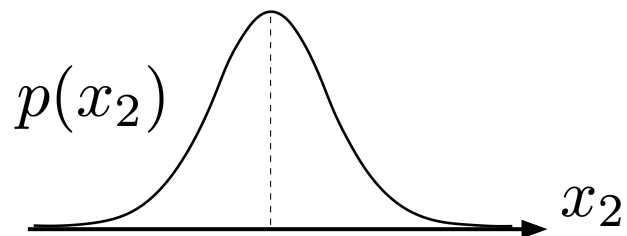
training images



Refresher: joint probability distribution



Probability density of x_1



Probability density of x_2

Refresher: joint probability distribution

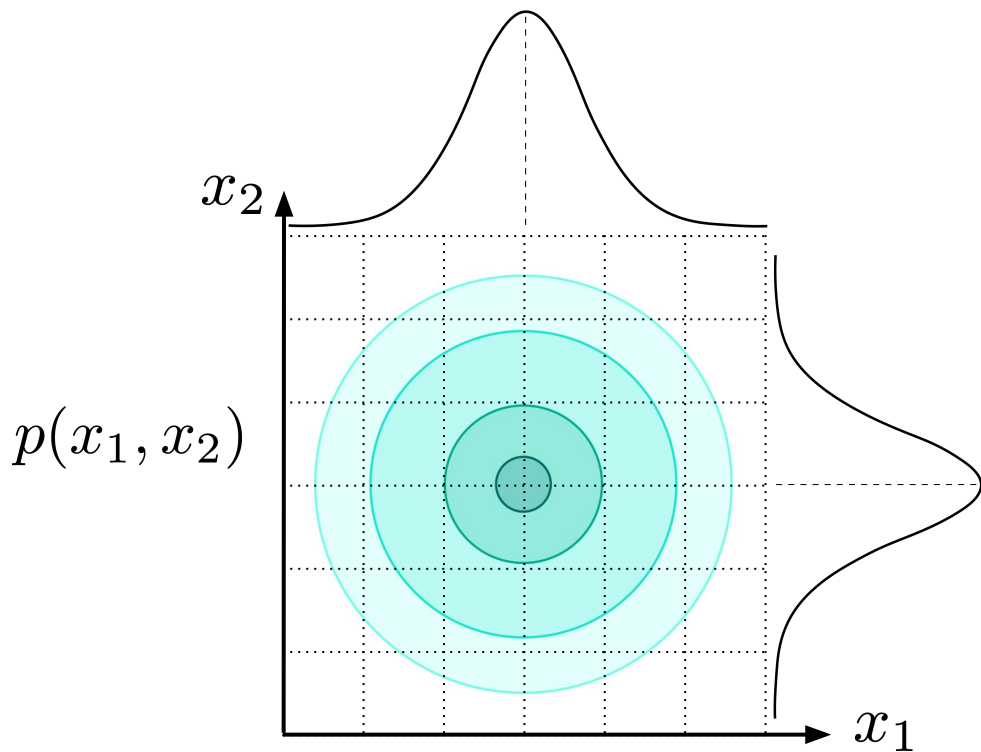
Joint probability distribution.

$$p(x_1, x_2) = p(x_1) \times p(x_2|x_1)$$



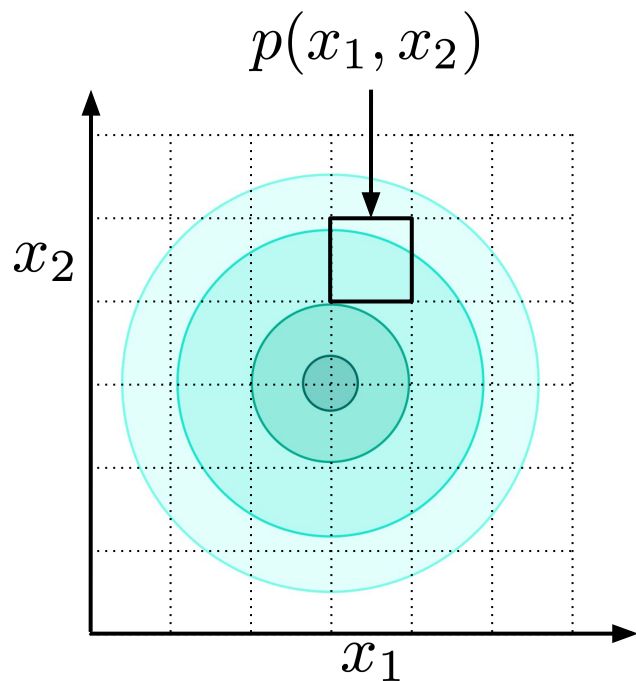
Conditional probability: if you know x_1 then what is the probability of x_2

Refresher: joint probability distribution



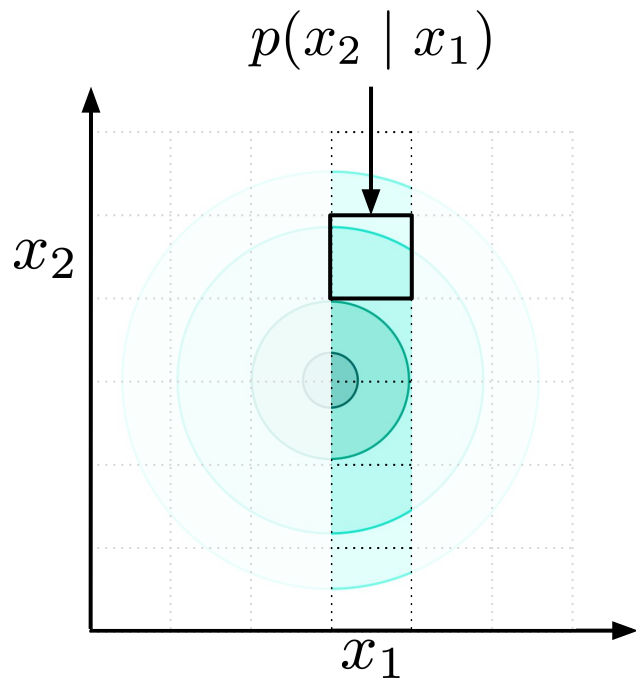
Joint probability $p(x_1, x_2)$ is the probability that we have x_1 **and** x_2

Refresher: joint probability distribution



Joint probability $p(x_1, x_2)$ is the probability that we have x_1 **and** x_2

Refresher: joint probability distribution



Conditional probability $p(x_2 | x_1)$
is the probability that we have x_2
knowing that we have x_1

Refresher: joint probability distribution

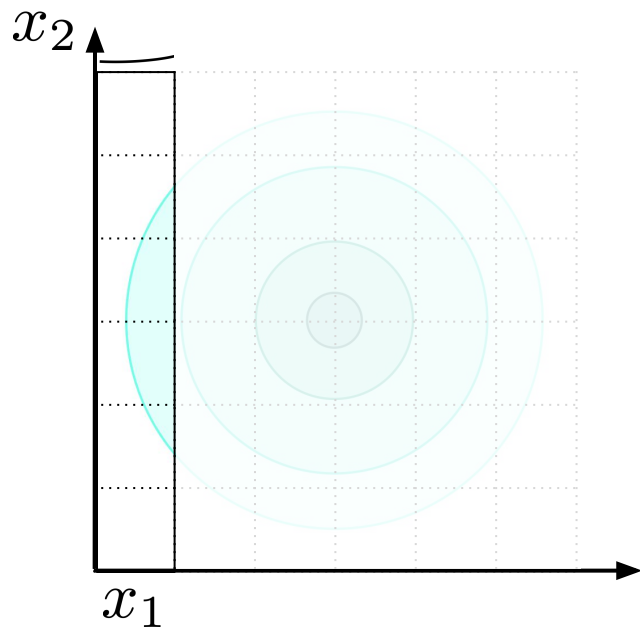
Marginalization.

$$p(x_1) = \int p(x_1, x_2) dx_2$$



Sum across all x_2

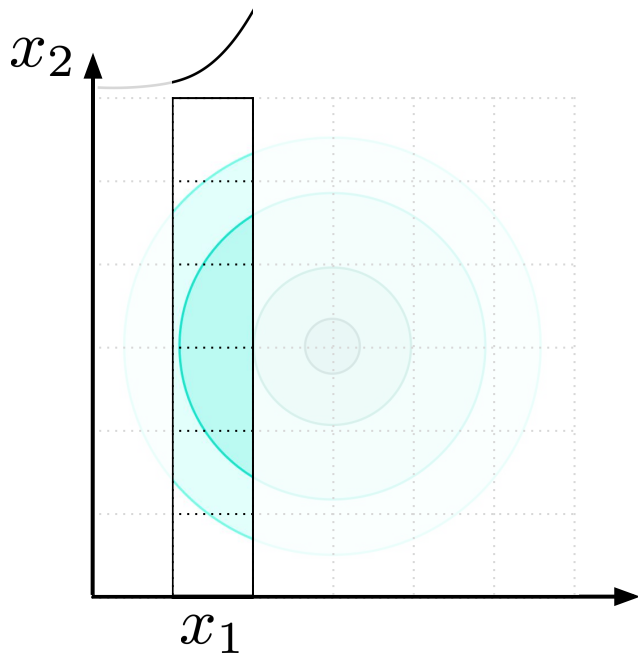
Refresher: joint probability distribution



$$p(x_1) = \int p(x_1, x_2) dx_2$$

↑
Sum across all x_2

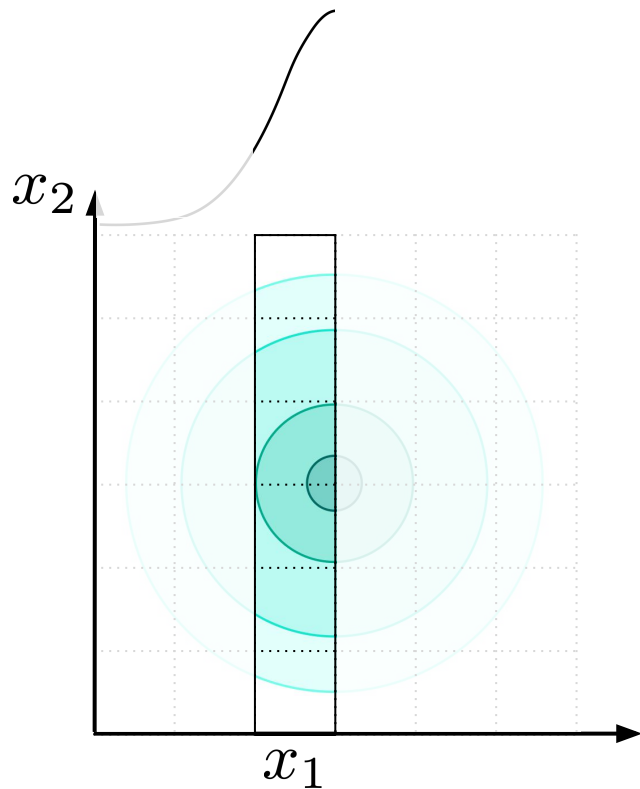
Refresher: joint probability distribution



$$p(x_1) = \int p(x_1, x_2) dx_2$$

↑
Sum across all x_2

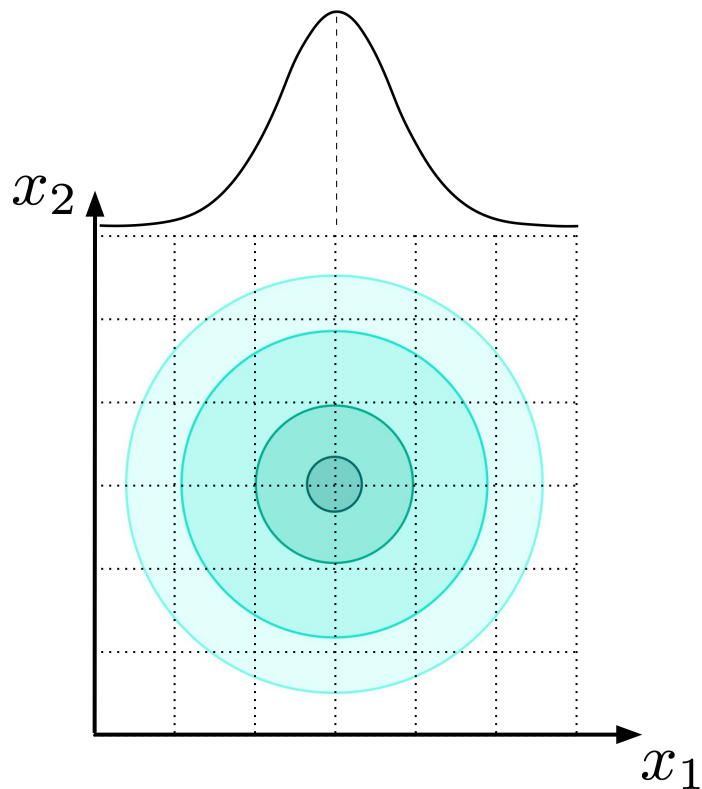
Refresher: joint probability distribution



$$p(x_1) = \int p(x_1, x_2) dx_2$$

↑
Sum across all x_2

Refresher: joint probability distribution



$$p(x_1) = \int p(x_1, x_2) dx_2$$



Sum across all x_2

Refresher: joint probability distribution

Joint probability distribution.

$$p(x_1, x_2, \dots, x_T) = p(x_1) \times p(x_2|x_1) \times \dots \times p(x_T|x_1, \dots, x_{T-1})$$

Marginalization.

$$p(x_1) = \int p(x_1, x_2, \dots, x_T) dx_2 \dots dx_T$$

Refresher: joint probability distribution

Joint probability distribution.

$$p(x_1, x_2, \dots, x_T) = p(x_1) \times p(x_2|x_1) \times \dots \times p(x_T|x_1, \dots, x_{T-1})$$

$$p(x_1, x_2, \dots, x_T) = p(x_{1:T})$$

Marginalization.

$$p(x_1) = \int p(x_1, x_2, \dots, x_T) dx_2 \dots dx_T$$

Refresher: joint probability distribution

Joint probability distribution (Simplified notation).

$$p(x_{1:T}) = p(x_1) \prod_{t=2}^T p(x_t | x_{1:t-1})$$

Marginalization (Simplified notation).

$$p(x_1) = \int p(x_{1:T}) dx_{2:T}$$

Problem

$$\log p_{\theta}(x_0) = ?$$

Derivation of a tractable loss

Strategy.

- 1 Derive **lower bound** for maximum (log-)likelihood estimation
- 2 **Expand terms** of lower bound
- 3 Show lower bound is **tractable**
- 4 Deduce **final loss function**

Step 1: Derivation of a lower bound

- 1 Derive **lower bound** for maximum (log-)likelihood estimation

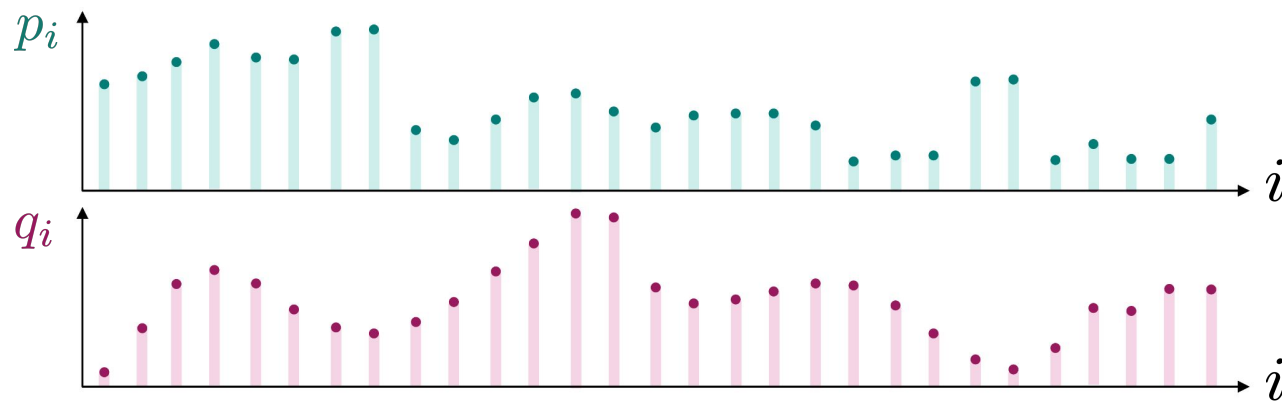
$$\mathbb{E}_{x_0 \sim q(x_0)} \left[\log p_\theta(x_0) \right] \geq \mathbb{E}_{x_{0:T} \sim q(x_{0:T})} \left[\log \frac{p_\theta(x_{0:T})}{q(x_{1:T} | x_0)} \right]$$

ELBO = **E**vidence **L**ower **B**ound

Derivation: Use Jensen's inequality on a convenient variational distribution

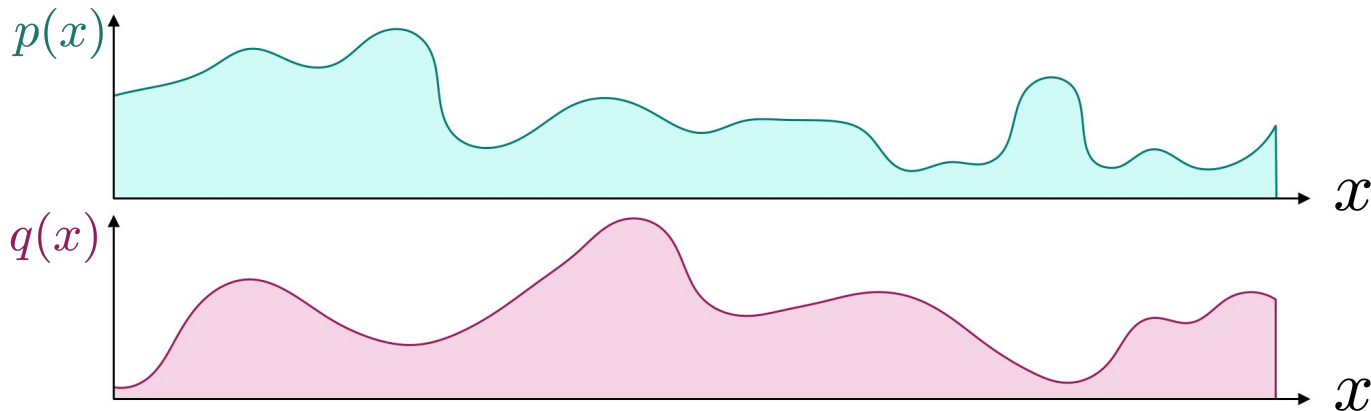
Refresher: KL divergence (discrete)

$$\text{KL}(P\|Q) = \sum_{i=1}^n p_i \log \left(\frac{p_i}{q_i} \right)$$



Refresher: KL divergence (continuous)

$$\text{KL}(P\|Q) = \int \log \left(\frac{p(x)}{q(x)} \right) p(x) dx = \mathbb{E}_{x \sim p} \left[\log \left(\frac{p(x)}{q(x)} \right) \right]$$



Step 3: Show lower bound is tractable

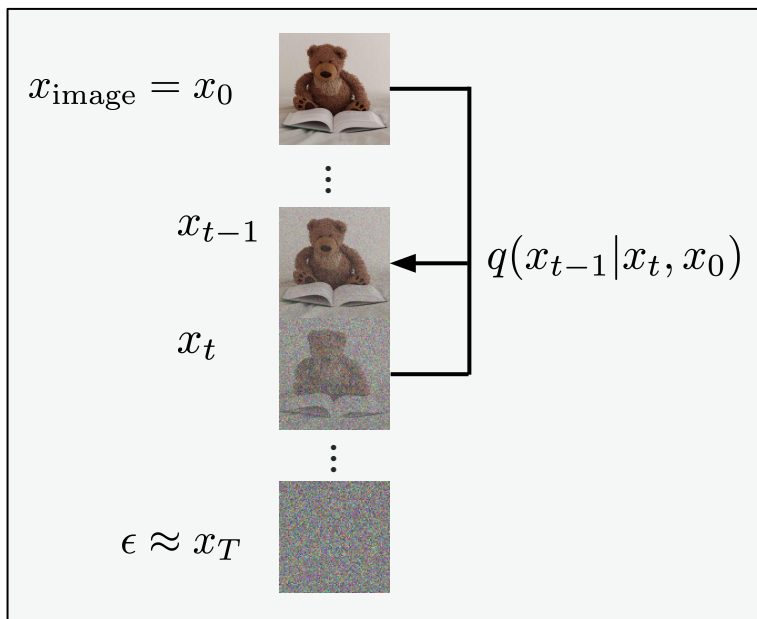
3 Show lower bound is **tractable**

a Show $q(x_{t-1} \mid x_t, x_0)$ is tractable

$$\text{ELBO} = - \sum_{t=2}^T \text{KL}(q(x_{t-1} \mid x_t, x_0) \parallel p_{\theta}(x_{t-1} \mid x_t)) + \text{some other terms}$$

Step 3: Show lower bound is tractable

3 a Show $q(x_{t-1} | x_t, x_0)$ is tractable



noisy image clean image

$q(x_{t-1} | x_t, x_0)$

less noisy image

Refresher: Bayes' rule

$$\text{Posterior } p(A|B) = \frac{\overset{\text{Likelihood}}{p(B|A)} \overset{\text{Prior}}{p(A)}}{\underset{\text{Marginal likelihood, aka evidence}}{p(B)}}$$

Step 3: Show lower bound is tractable

- 3 a Show $q(x_{t-1} | x_t, x_0)$ is tractable

$$q(x_{t-1} | x_t, x_0) = \mathcal{N}(\tilde{\mu}_t, \tilde{\beta}_t)$$

sampled noise used to go
from clean to noisy image

$$\text{with } \tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_t \right)$$

noisy image

Derivation: Compute density with Bayes' rule via $q(\cdot | x_0)$ + Markov property

Step 3: Show lower bound is tractable

3 Show lower bound is **tractable**

a Show $q(x_{t-1} | x_t, x_0)$ is tractable

b Show $p_\theta(x_{t-1} | x_t)$ is tractable

$$\text{ELBO} = - \sum_{t=2}^T \text{KL}(q(x_{t-1} | x_t, x_0) \| p_\theta(x_{t-1} | x_t)) + \text{some other terms}$$

Step 3: Show lower bound is tractable

- 3 b Show $p_\theta(x_{t-1} | x_t)$ is tractable

We assume that p_θ is Gaussian:

$$p_\theta(x_{t-1} | x_t) = \mathcal{N}(\mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

Justification: If $q(x_t | x_{t-1})$ Gaussian, then $q(x_{t-1} | x_t)$ is approx. Gaussian

Step 4: Deduce final loss function

4 Deduce **loss function**

$$\mathcal{L}_{\text{DDPM}} = \mathbb{E}_{t, x_0, \epsilon} \left[\left\| \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) - \epsilon \right\|^2 \right]$$

$$t \sim \mathcal{U}\{1, T\} \quad x_0 \sim q_0(x_0) \quad \epsilon \sim \mathcal{N}(0, I)$$

Derivation: Compute KL divergence between two Gaussians

Summary of what we have done

Strategy.

- 1 Derive **lower bound** for maximum (log-)likelihood estimation ← **ELBO!**
- 2 **Expand terms** of lower bound ← **Surfaced KL divergence**
- 3 Show lower bound is **tractable** ← **Bayes' rule + Gaussian properties**
- 4 Deduce **final loss function** ← **Incredibly simple noise prediction!**



Diffusion & Large Vision Models

Motivation

Forward / backwards process

Variational formulation

Training

Inference

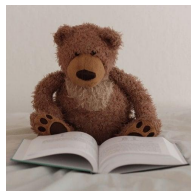
Faster sampling

Training recipe

1. Sample:

clean image

$$x_0 \sim q_0(x_0)$$



noise

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



time step

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

noised image

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$$



Training recipe

1. Sample:

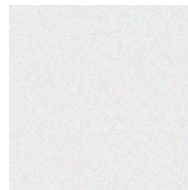
clean image

$$x_0 \sim q_0(x_0)$$



noise

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



time step

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

noised image

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$$



2. Use x_t and t to **predict** ϵ via $\epsilon_\theta(x_t, t)$

Training recipe

1. Sample:

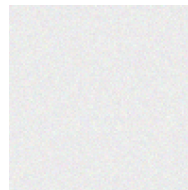
clean image

$$x_0 \sim q_0(x_0)$$



noise

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



time step

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

noised image

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$$



2. Use x_t and t to **predict** ϵ via $\epsilon_\theta(x_t, t)$

3. Compute **loss** $\mathcal{L} = \|\epsilon_\theta(x_t, t) - \epsilon\|^2$ and **backpropagate** through ϵ_θ



Diffusion & Large Vision Models

Motivation

Forward / backwards process

Variational formulation

Training

Inference

Faster sampling

Inference recipe

1. Sample:

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



2. Perform **iterative update** from x_t to x_{t-1}

Inference recipe

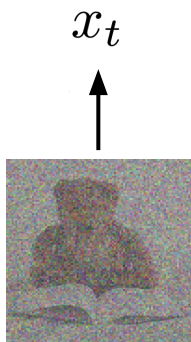
1. Sample:

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

noise



2. Perform **iterative update** from x_t to x_{t-1}



Inference recipe

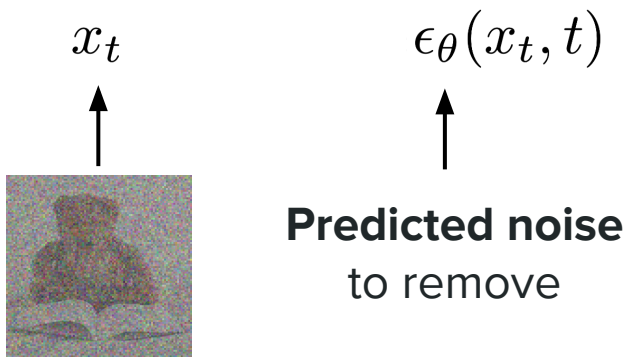
1. Sample:

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

noise



2. Perform **iterative update** from x_t to x_{t-1}



Inference recipe

1. Sample:

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

noise



2. Perform **iterative update** from x_t to x_{t-1}

$$x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t)$$



Predicted noise
to remove

Inference recipe

1. Sample:

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

noise



2. Perform **iterative update** from x_t to x_{t-1}

$$\frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(x_t, t) \right) + \sigma_t z$$



Predicted noise
to remove

Sampled from a
Normal distribution

Inference recipe

1. Sample:

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

noise



2. Perform **iterative update** from x_t to x_{t-1}

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



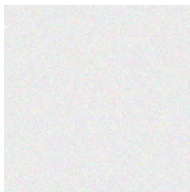
Predicted noise
to remove

Sampled from a
Normal distribution

Inference recipe

1. Sample:

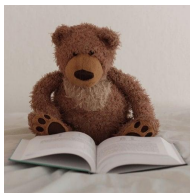
$$x_T \sim \mathcal{N}(0, I)$$



2. Perform **iterative update** from x_t to x_{t-1}

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

3. Obtain **final** image x_0



Landmark paper that is a must-read

DDPM = Denoising Diffusion Probabilistic Models

Denoising Diffusion Probabilistic Models

Jonathan Ho
UC Berkeley
jonathanho@berkeley.edu

Ajay Jain
UC Berkeley
ajayj@berkeley.edu

Pieter Abbeel
UC Berkeley
pabbeel@cs.berkeley.edu

Abstract

We present high quality image synthesis results using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive lossy decompression scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR10 dataset, we obtain an Inception score of 9.46 and a state-of-the-art FID score of 3.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/hojonathanho/diffusion>.



Diffusion & Large Vision Models

Motivation

Forward / backwards process

Variational formulation

Training

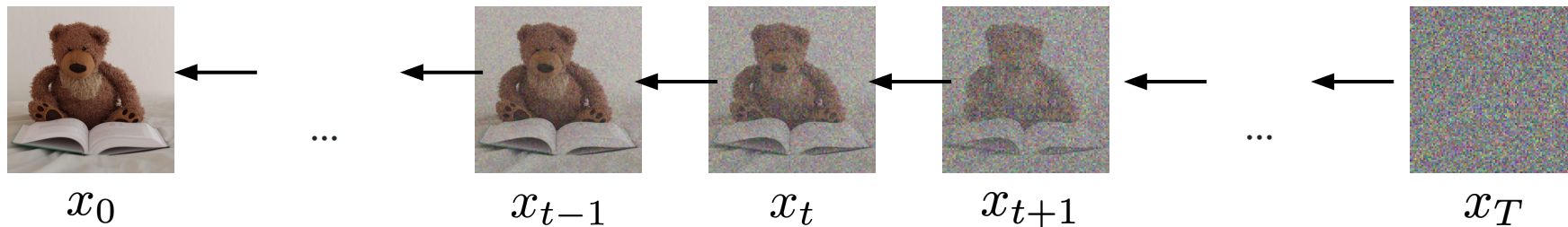
Inference

Faster sampling

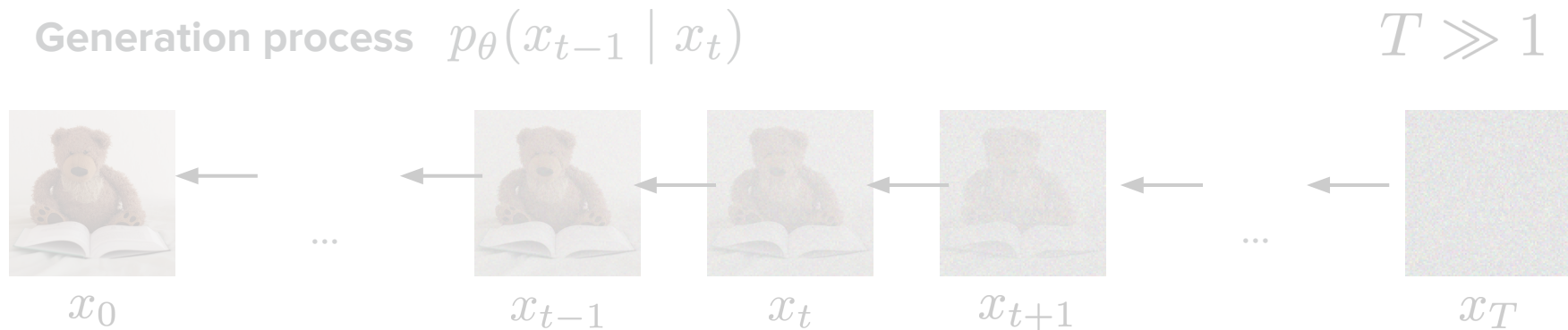
Limitations

Generation process $p_{\theta}(x_{t-1} | x_t)$

$T \gg 1$



Limitations



Orders of magnitude

- $O(T)$ more expensive than VAE/GAN
- $T = 1000 \rightarrow$ possibly spend \sim minutes per sample

Limitations

Generation process $p_{\theta}(x_{t-1} | x_t)$

$T \gg 1$



Orders of magnitude



- $O(T)$ more expensive than VAE/GAN
- $T = 1000 \rightarrow$ possibly spend ~minutes per sample

In practice, DDPM is too slow at inference time!

Mitigation attempts

Starting point. Recall iterative generative formula between two steps:

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

  **Predicted noise**
to remove

Sampled from a Normal distribution

Mitigation attempts

Attempt. What if we derive a formula by induction?

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

$$x_{t-2} = \frac{1}{\sqrt{\alpha_{t-1}}} \left(x_{t-1} - \frac{1 - \alpha_{t-1}}{\sqrt{1 - \bar{\alpha}_{t-1}}} \epsilon_{\theta}(x_{t-1}, t - 1) \right) + \sigma_{t-1} z$$

⋮

$$x_{\tau_{i-1}} = Ax_{\tau_i} + \sum_{s=\tau_{i-1}+1}^{\tau_i} B_s \epsilon_{\theta}(x_s, s) + C_s z_s$$

Mitigation attempts

Attempt. What if we derive a formula by induction?

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

$$x_{t-2} = \frac{1}{\sqrt{\alpha_{t-1}}} \left(x_{t-1} - \frac{1 - \alpha_{t-1}}{\sqrt{1 - \bar{\alpha}_{t-1}}} \epsilon_{\theta}(x_{t-1}, t-1) \right) + \sigma_{t-1} z$$

⋮

$$x_{\tau_{i-1}} = Ax_{\tau_i} + \sum_{s=\tau_{i-1}+1}^{\tau_i} B_s \epsilon_{\theta}(x_s, s) + C_s z_s$$

Problem: Still need to evaluate function several times. No progress.

Mitigation attempts

Attempt 2. What if we directly skip steps?

$$\tau_0 = 0 < \tau_1 < \dots < \tau_{S-1} < \tau_S = T$$

$$x_{\tau_{i-1}} = \sqrt{\frac{\bar{\alpha}_{\tau_{i-1}}}{\bar{\alpha}_{\tau_i}}} \left(x_{\tau_i} - \frac{1 - \bar{\alpha}_{\tau_i}/\bar{\alpha}_{\tau_{i-1}}}{\sqrt{1 - \bar{\alpha}_{\tau_i}}} \epsilon_{\theta}(x_{\tau_i}, \tau_i) \right) + \sigma_{\tau_i} z$$

Mitigation attempts

Attempt 2. What if we directly skip steps?

$$\tau_0 = 0 < \tau_1 < \dots < \tau_{S-1} < \tau_S = T$$

$$x_{\tau_{i-1}} = \sqrt{\frac{\bar{\alpha}_{\tau_{i-1}}}{\bar{\alpha}_{\tau_i}}} \left(x_{\tau_i} - \frac{1 - \bar{\alpha}_{\tau_i}/\bar{\alpha}_{\tau_{i-1}}}{\sqrt{1 - \bar{\alpha}_{\tau_i}}} \epsilon_{\theta}(x_{\tau_i}, \tau_i) \right) + \boxed{\sigma_{\tau_i} z}$$

Problem: mix large jumps + stochasticity = bad quality

Problem reformulation

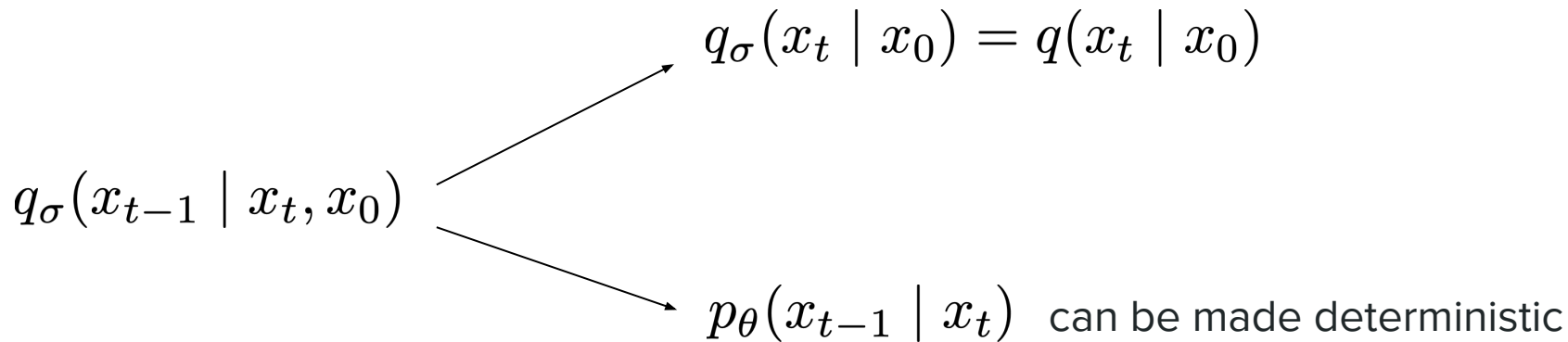
Follow-up idea. How to remove stochasticity while making these larger jumps?

Key property. Loss function only relies on marginals $q(x_t \mid x_0)$

Problem statement. Find q such that:

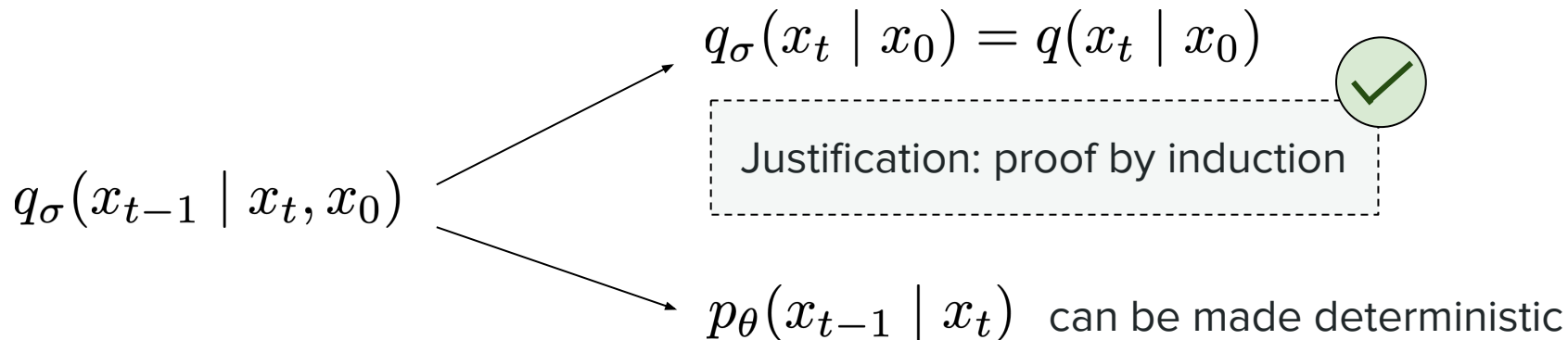
- marginal is the same as in DDPM
- generation process can be made deterministic

Problem reformulation



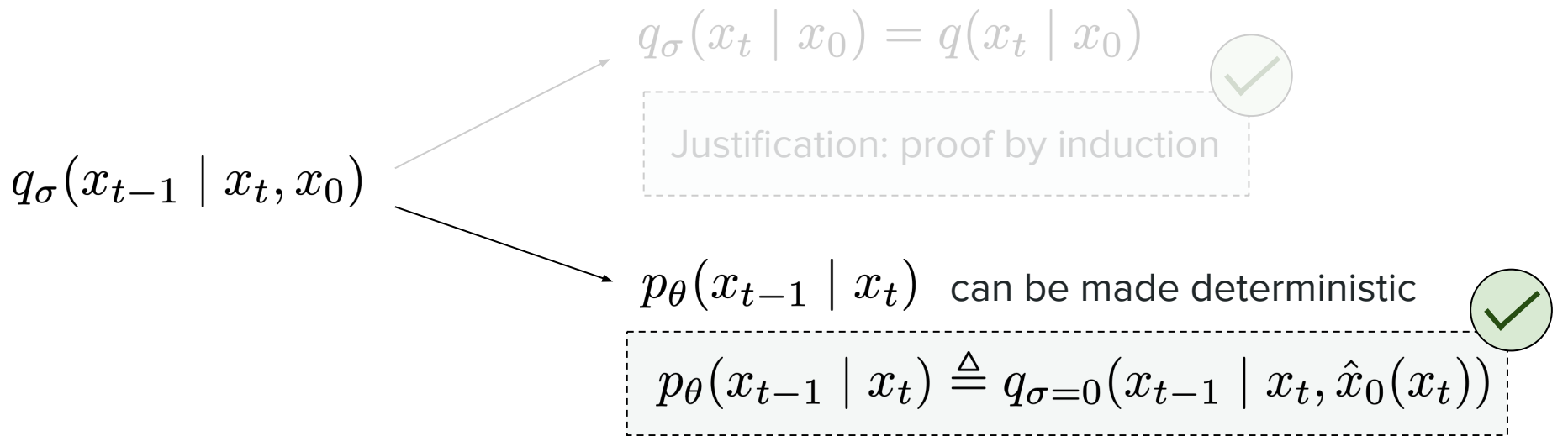
$$q_\sigma(x_{t-1} \mid x_t, x_0) = \mathcal{N}\left(\sqrt{\bar{\alpha}_{t-1}} x_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \frac{x_t - \sqrt{\bar{\alpha}_t} x_0}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2 I\right)$$

Problem reformulation



$$q_\sigma(x_{t-1} \mid x_t, x_0) = \mathcal{N}\left(\sqrt{\bar{\alpha}_{t-1}} x_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \frac{x_t - \sqrt{\bar{\alpha}_t} x_0}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2 I\right)$$

Problem reformulation



$$q_\sigma(x_{t-1} \mid x_t, x_0) = \mathcal{N}\left(\sqrt{\bar{\alpha}_{t-1}} x_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \cancel{\sigma_t^2}} \frac{x_t - \sqrt{\bar{\alpha}_t} x_0}{\sqrt{1 - \bar{\alpha}_t}}, \cancel{\sigma_t^2} I\right)$$

Why is DDIM called that way?

DDIM = **D**enoising **D**iffusion **I**mplicit **M**odels

$$\sigma_t = 0$$

$$\hat{x}_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(x_t, t)}{\sqrt{\bar{\alpha}_t}}$$

"our prediction
of the clean
image from time
t"



Derivation: use the marginal distribution formula

Why is DDIM called that way?

DDIM = **D**enoising **D**iffusion **I**mplicit **M**odels

$$\sigma_t = 0$$

$$\hat{x}_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(x_t, t)}{\sqrt{\bar{\alpha}_t}}$$

$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1}} \epsilon_\theta(x_t, t) + \emptyset$$

↑ "Our prediction of the clean image from time t"

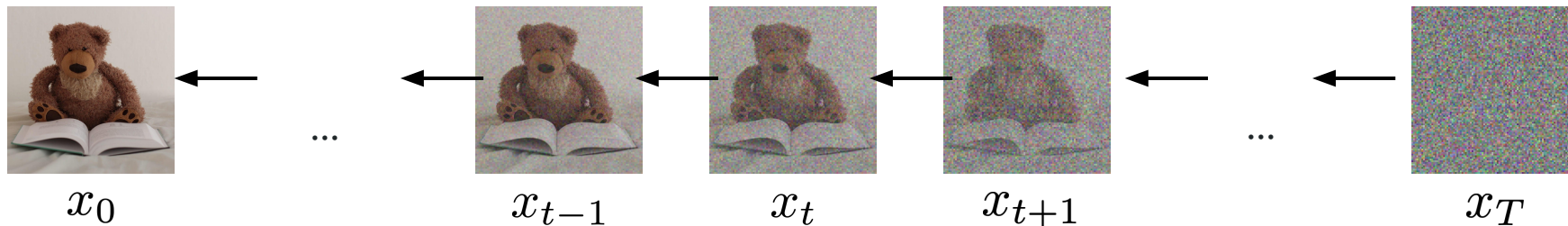
↑ "...that we noise to time t-1..."

↑ ... with no stochasticity!

Why is the DDIM formulation useful?

Update rule

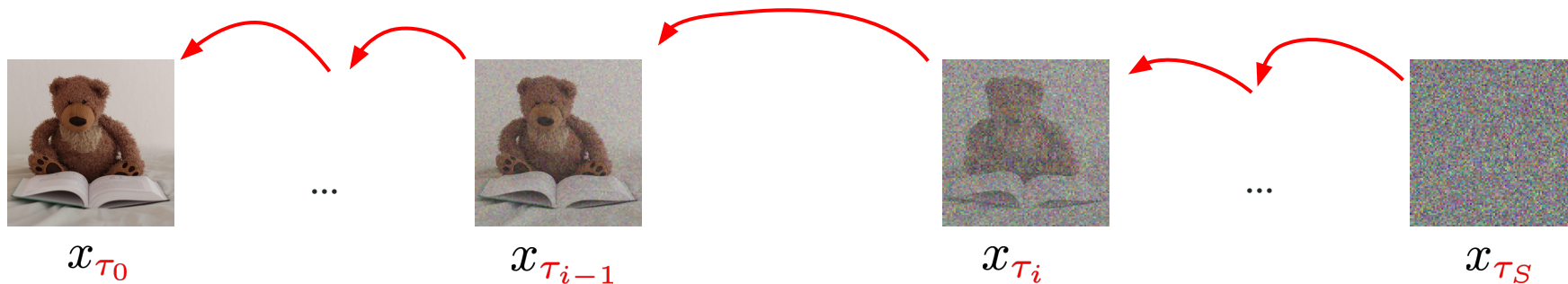
$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{x}_0(t) + \sqrt{1 - \bar{\alpha}_{t-1}} \epsilon_{\theta}(x_t, t)$$



Why is the DDIM formulation useful?

Update rule: **accelerated!**

$$x_{\tau_{i-1}} = \sqrt{\bar{\alpha}_{\tau_{i-1}}} \hat{x}_0(\tau_i) + \sqrt{1 - \bar{\alpha}_{\tau_{i-1}}} \epsilon_{\theta}(x_{\tau_i}, \tau_i)$$



$$\tau_0 = 0 < \tau_1 < \dots < \tau_{S-1} < \tau_S = T$$

$$S \ll T \quad \text{with} \quad \frac{T}{S} \triangleq \text{"speedup"}$$

Applications

Speedup

- Typically 10x-100x
- Need to trade off speed with quality

Applications

Speedup

- Typically 10x-100x
- Need to trade off speed with quality

Quality (based on CIFAR10 experiments)

Speed-up	1x	10x	20x	50x	100x
FID impact	baseline	+3%	+16%	+70%	+330%

Applications

Speedup

- Typically 10x-100x
- Need to trade off speed with quality

Quality (based on CIFAR10 experiments)

Speed-up	1x	10x	20x	50x	100x
FID impact	baseline	+3%	+16%	+70%	+330%

Discussion

- Skipping steps with inter-steps stochasticity → much worse quality
- Schedule for $\{\tau_i\}$ can be chosen as a hyperparameter

DDIM conclusion

Recipe

1. Choose "convenient" modeling assumptions compatible with DDPM optimization
2. Remove inter-step stochasticity = DDIM
3. Skip steps!

$$p_{\theta}(x_{t-1} \mid x_t) \longrightarrow p_{\theta}(x_{\tau_i-1} \mid x_{\tau_i})$$

Thank you for your attention!