# Deep Learning for Computer Vision

## Andrii Liubonko

Grammarly*

* The opinions expressed in this presentation and on the following slides are solely those of the presenter and not necessarily those of Grammarly

# Logistics

**4 units**

**3 types of homework:**
- paper review   [20 % of FINAL SCORE]
- notebooks      [30 % of FINAL SCORE]
- mini-project   [50 % of FINAL SCORE]

**important dates:**
*31  January, 23:59*

***course repo:***
https://github.com/lyubonko/ucu2022cv

# Overview of the course

**Unit I**
[T] Intro
[P] pytorch

**Unit II**
[T] CNNs in depth, Object Detection
[P] simple nets

**Unit III**
[T] Attention in CV, Transformers
[P] classification

**Unit IV**
[T] Generative models, Diffusion Models
[P] project structure, stable diffusion

"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

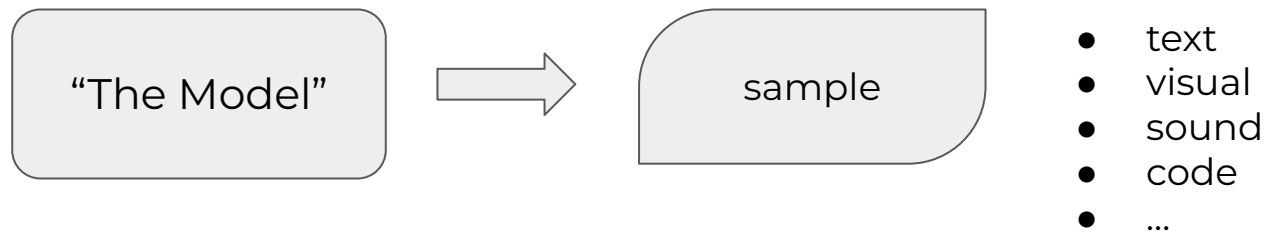"Creating noise from data is easy; creating data from noise is generative modeling"
© 2011.13456

# Content of today's lecture

- Generative modeling

  - appetiser
  - approaches

- Diffusion models

  - intro
  - tale of marvelous Gaussians
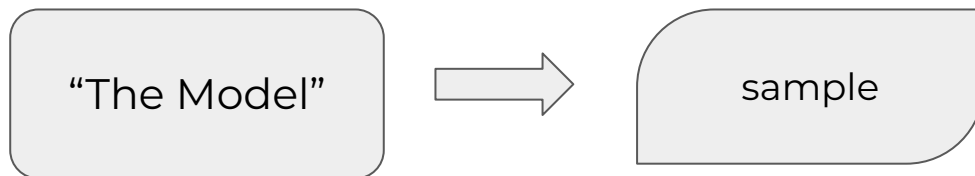  - more details

- Stable Diffusion

# Generative modeling
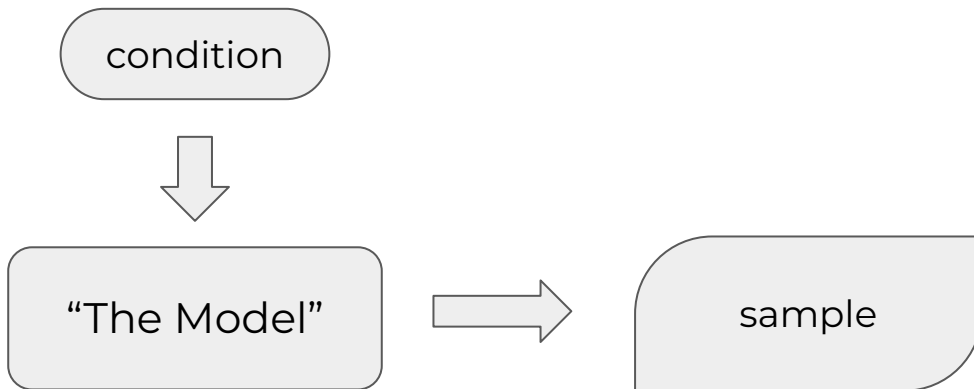
- uncontrollable (unconditional) generation:

"The Model" → sample

- text
- visual
- sound
- code
- …

# Generative modeling

- uncontrollable (unconditional) generation:

"The Model" → sample

- text
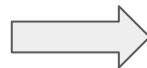- visual
- sound
- code
- …

- controllable (conditional) generation:

condition → "The Model" → sample

# Generative modeling

# Generative modeling



condition:

"The Model"

# Generative modeling

Control over multiple conditions:



condition:

+

condition:

**Texture Description**

| upper clothing texture ∨ | lower clothing texture ∨ | outer clothing texture ∨ |
| pure color ∨ | denim ∨ | stripe/spline ∨ |

"The Model" →

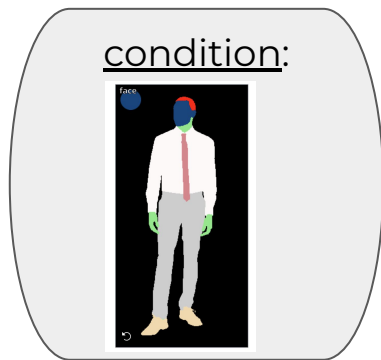https://huggingface.co/spaces/CVPR/drawings-to-human

# Generative modeling

Control over multiple conditions:



condition:

+

condition:

**Texture Description**
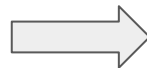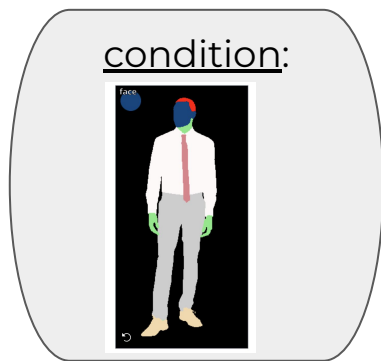
| upper clothing texture | lower clothing texture | outer clothing texture |
| pure color | pure color | plaid/lattice |

"The Model"

→

https://huggingface.co/spaces/CVPR/drawings-to-human

# Generative modeling



condition:

"a cover for the cosmopolitan magazine"

→

"The Model"

→

# Generative modeling



condition:

"Christmas at Lviv"

→

"The Model"

→

https://huggingface.co/spaces/stabilityai/stable-diffusion

# Generative modeling [for text]

condition:
(prompt)

"Once upon a time,"

$\Rightarrow$

"The Model"

$\Rightarrow$

sample

---

⚡ **Text Generation demo**

using **gpt2**

📝 Text Generation                                    Examples ⌄

Once upon a time, he went out of his way to not be caught. Though he had been arrested with his head set upon the railing when the storm struck, he did not get his last glimpse of the night sky, where a strange woman and

https://huggingface.co/tasks/text-generation

# Generative modeling [for text]

- controllable (conditional) generation:

# Content of today's lecture

- Generative modeling

  - appetiser
  - approaches

- Diffusion models

  - intro
  - tale of marvelous Gaussians
  - more details

- Stable Diffusion

**EBM:** Approximate Maximum likelihood

$x$ → Energy $E(x)$ → R

**GAN:** Adversarial training

$x'$ | $x$ → Discriminator $D(x)$ → 0/1 | $z$ → Generator $G(z)$ → $x'$

**VAE:** Maximize variational lower bound

$x$ → Encoder $q_\phi(z|x)$ → $z$ → Decoder $p_\theta(x|z)$ → $x'$

**Flow-based Model:** Invertible transform of distributions

$x$ → Flow $f(x)$ → $z$ → Inverse $f^{-1}(z)$ → $x'$

**Diffusion Model:** Gradually add Gaussian noise and then reverse

$x$ ⇄ $z_1$ ⇄ $z_2$ ⋯ ⋯ ⋯ ⇄ $z_T$

**Autoregressive model:** Learn conditional of each variable given past

$x^0$ → $x^1$ → $x^2$ → $x^3$ ⋯ ⋯ ⋯ $x^D$

Kevin P. Murphy
"Probabilistic Machine Learning:
Advanced Topics", 2023

originally based on
What are Diffusion Models? | Lil'Log

**EBM:** Approximate Maximum likelihood — Energy $E(x)$

**GAN:** Adversarial training — Discriminator $D(x)$, Generator $G(z)$

**VAE:** Maximize variational lower bound — Encoder $q_\phi(z|x)$, Decoder $p_\theta(x|z)$

**Flow-based Model:** Invertible transform of distributions — Flow $f(x)$, Inverse $f^{-1}(z)$

**Diffusion Model:** Gradually add Gaussian noise and then reverse

**Autoregressive model:** Learn conditional of each variable given past

Kevin P. Murphy
"Probabilistic Machine Learning: Advanced Topics", 2023

originally based on
What are Diffusion Models? | Lil'Log
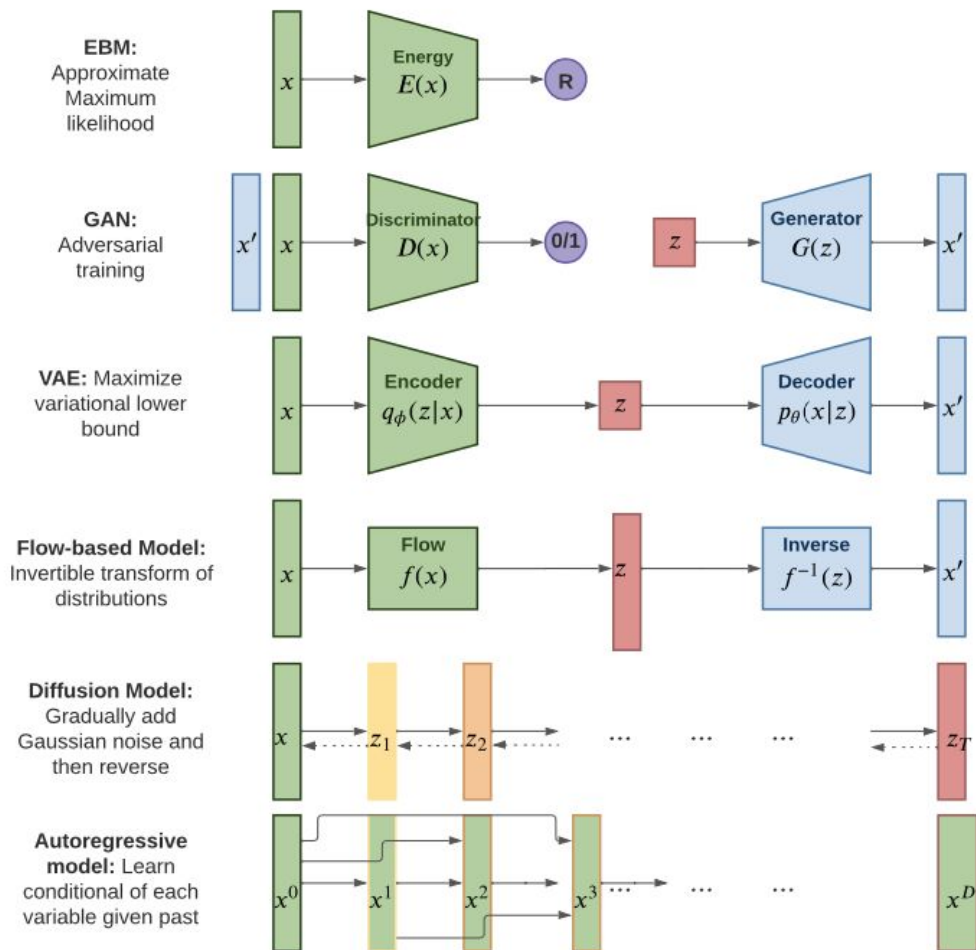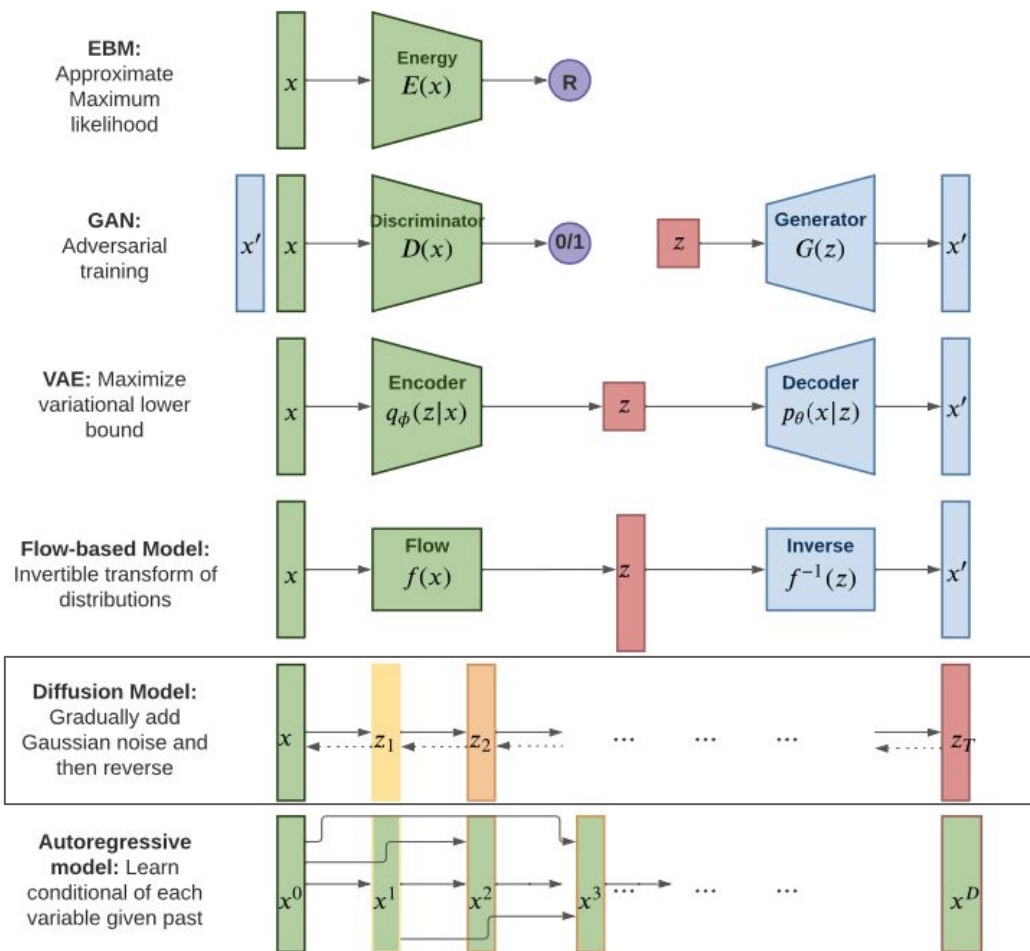
# Content of today's lecture

- Generative modeling

  - appetiser
  - approaches

- Diffusion models

  - intro
  - tale of marvelous Gaussians
  - more details

- Stable Diffusion

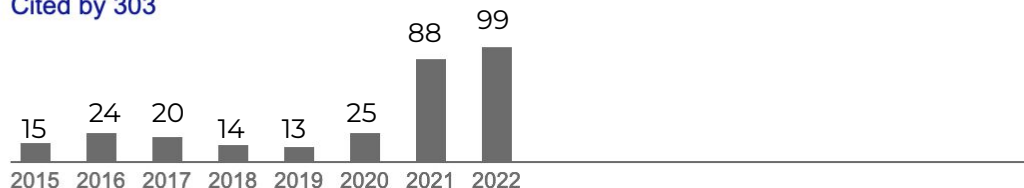# Deep unsupervised learning using nonequilibrium thermodynamics

| | |
|---:|:---|
| Authors | Jascha Sohl-Dickstein, Eric A Weiss, Niru Maheswaranathan, Surya Ganguli |
| Publication date | 2015/3/12 |
| Journal | International Conference on Machine Learning |
| Description | A central problem in machine learning involves modeling complex data-sets using highly flexible families of probability distributions in which learning, sampling, inference, and evaluation are still analytically or computationally tractable. Here, we develop an approach that simultaneously achieves both flexibility and tractability. The essential idea, inspired by non-equilibrium statistical physics, is to systematically and slowly destroy structure in a data distribution through an iterative forward diffusion process. We then learn a reverse diffusion process that restores structure in data, yielding a highly flexible and tractable generative model of the data. This approach allows us to rapidly learn, sample from, and evaluate probabilities in deep generative models with thousands of layers or time steps, as well as to compute conditional and posterior probabilities under the learned model. We additionally release an open source reference implementation of the algorithm. |

Total citations

| 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 |
|------|------|------|------|------|------|------|------|
| 15 | 24 | 20 | 14 | 13 | 25 | 88 | 99 |

| | |
|---|---|
| 15 | 2015 |

[1503.03585] Deep Unsupervised Learning using Nonequilibrium Thermodynamics
#foundational, #images,
Jascha Sohl-Dickstein [Stanford University]

| | |
|---|---|
| 24 | 2016 |

[2006.11239] Denoising Diffusion Probabilistic Models #foundational, #images,
#connecting score based method (1907.05600)
Jonathan Ho [UC Berkeley]

| | |
|---|---|
| 20 | 2017 |

[2011.13456] Score-Based Generative Modeling through Stochastic
Differential Equations code blog #foundational, #extending to Stochastic stuff

| | |
|---|---|
| 14 | 2018 |

Yang Song, Jascha Sohl-Dickstein [Stanford University + Google Brain]

| | |
|---|---|
| 13 | 2019 |

[2105.05233] Diffusion Models Beat GANs on Image Synthesis **code**
#images, [openAI]
Yang Song, Jascha Sohl-Dickstein [Stanford University + Google Brain]

| | |
|---|---|
| 25 | 2020 |

[2205.14217] Diffusion-LM Improves Controllable Text Generation
#language, #Stanford

| | |
|---|---|
| 88 | 2021 |

[2208.04202] Analog Bits: Generating Discrete Data using Diffusion Models
with Self-Conditioning

| | |
|---|---|
| 99 | 2022 |

#images, #language, #Google
Hinton

🤗 Diffusers provides pretrained diffusion models across multiple modalities, such as vision and audio, and serves as a modular toolbox for inference and training of diffusion models.
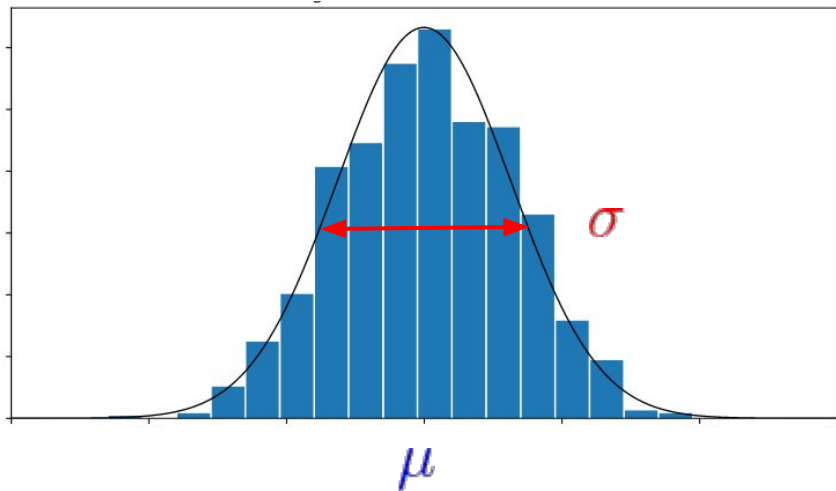
More precisely, 🤗 Diffusers offers:

- State-of-the-art diffusion pipelines that can be run in inference with just a couple of lines of code (see src/diffusers/pipelines).
- Various noise schedulers that can be used interchangeably for the prefered speed vs. quality trade-off in inference (see src/diffusers/schedulers).
- Multiple types of models, such as UNet, that can be used as building blocks in an end-to-end diffusion system (see src/diffusers/models).
- Training examples to show how to train the most popular diffusion models (see examples).

https://github.com/huggingface/diffusers

- Generative modeling

- Diffusion models

  - intro
  - **tale of marvelous Gaussians**
  - details

- Application to NLP domain

  - [2205.14217] Diffusion-LM Improves Controllable Text Generation
  - [2208.04202] Analog Bits: Generating Discrete Data using Diffusion Models with Self-Conditioning
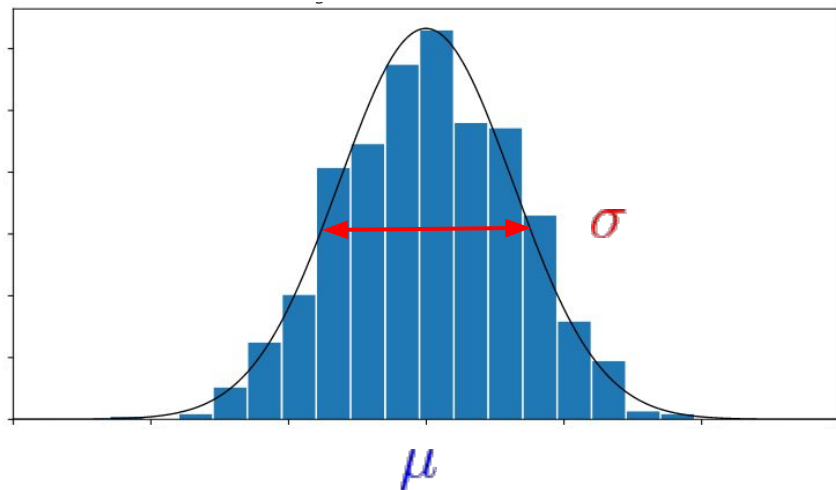  - [2208.00638] Composable Text Control Operations in Latent Space with Ordinary Differential Equations

# Gaussian distribution

$$\mathcal{N}(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]$$

# Gaussian distribution

$$\mathcal{N}(x;\mu,\sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]$$



## The Unreasonable Effectiveness of Mathematics in the Natural Sciences

Richard Courant Lecture in Mathematical Sciences delivered at New York University, May 11, 1959

EUGENE P. WIGNER

Princeton University

"*and it is probable that there is some secret here which remains to be discovered.*" (C. S. Peirce)

There is a story about two friends, who were classmates in high school, talking about their jobs. One of them became a statistician and was working on population trends. He showed a reprint to his former classmate. The reprint started, as usual, with the Gaussian distribution and the statistician explained to his former classmate the meaning of the symbols for the actual population, for the average population, and so on. His classmate was a bit incredulous and was not quite sure whether the statistician was pulling his leg. "How can you know that?" was his query. "And what is this symbol here?" "Oh," said the statistician, "this is $\pi$." "What is that?" "The ratio of the circumference of the circle to its diameter." "Well, now you are pushing your joke too far," said the classmate, "surely the population has nothing to do with the circumference of the circle."
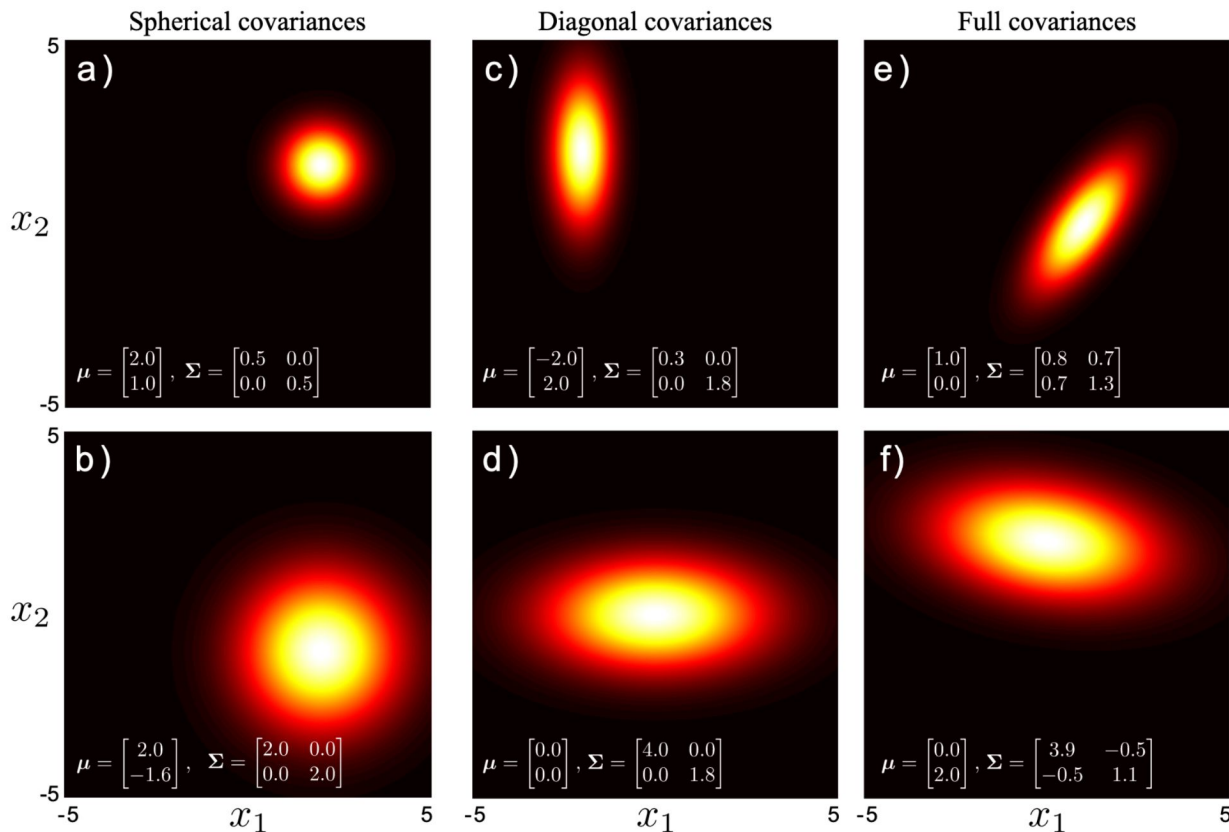
# Multivariate Gaussian distribution

$$\mathcal{N}(\mathbf{x}; \mu, \mathbf{\Sigma}) = \frac{1}{|\mathbf{\Sigma}|^{1/2} (2\pi)^{D/2}} \exp\left[ -\frac{1}{2} (\mathbf{x} - \mu)^T \cdot \mathbf{\Sigma}^{-1} \cdot (\mathbf{x} - \mu) \right]$$

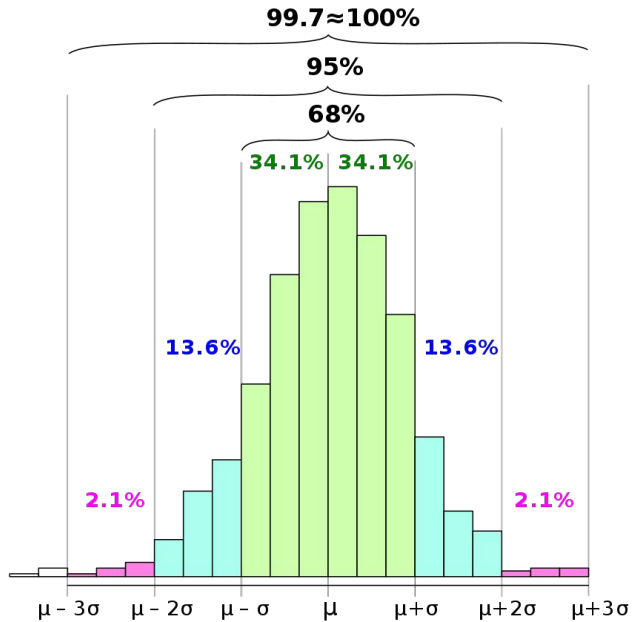**μ** is *[D x 1]* vector  that describes the position of the distribution

**Σ** is *[D x D]* positive definite matrix
(implying that $\mathbf{z}^T \cdot \mathbf{\Sigma} \cdot \mathbf{z}$ is positive for any real vector **z**)
and describes the shape of the distribution

$$\mathbf{\Sigma}_{spher} = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix} \qquad \mathbf{\Sigma}_{diag} = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \qquad \mathbf{\Sigma}_{full} = \begin{bmatrix} \sigma_{11}^2 & \sigma_{12}^2 \\ \sigma_{21}^2 & \sigma_{22}^2 \end{bmatrix}$$
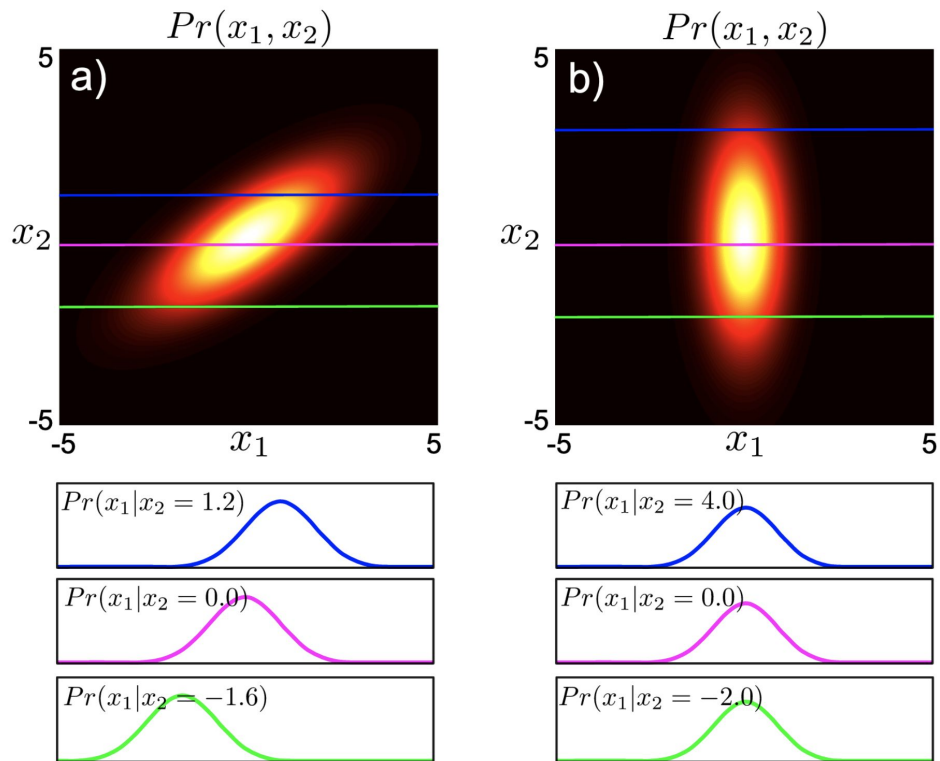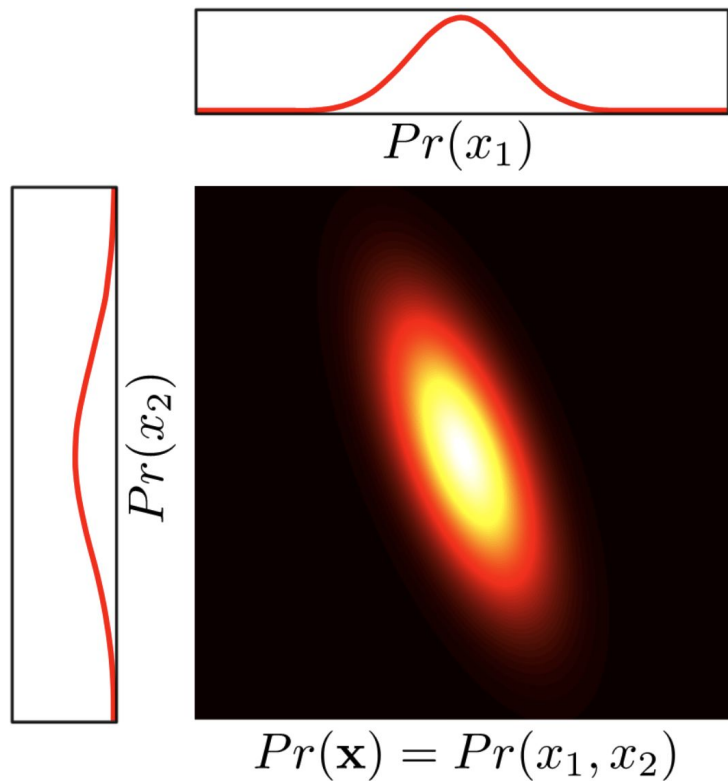


Simon J.D. Prince, "Computer Vision: Models, Learning, and Inference"

# Multivariate Gaussian distribution



| n = 1 | 99.73 % |
|---|---|
| n = 10 | 97.33 % |
| n = 100 | 76.31 % |
| n = 300 | 44.44 % |
| n = 1000 | 6.70 % |

# Conditional distributions (of multivariate Gaussian)



If we take any multivariate normal distribution, fix a subset of the variables, and look at the distribution of the remaining variables, this distribution will also take the form of a normal. The mean of this new normal depends on the values that we fixed the subset to, but the covariance is always the same.

Simon J.D. Prince, "Computer Vision: Models, Learning, and Inference"

# Marginal distributions (of multivariate Gaussian)



$$Pr(x_1)$$

$$Pr(x_2)$$

$$Pr(\mathbf{x}) = Pr(x_1, x_2)$$

The marginal distribution of any subset of variables in a normal distribution is also normally distributed. In other words, if we sum over the distribution in any direction, the remaining quantity is also normally distributed. To find the mean and the covariance of the new distribution, we can simply extract the relevant entries from the original mean and covariance matrix.

Simon J.D. Prince, "Computer Vision: Models, Learning, and Inference"

# Linear transformations of variables



$Pr(\mathbf{x})$

a)

$Pr(\mathbf{y})$

b) $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$

The form of the multivariate normal is preserved under linear transformations

**y = Ax + b**

If the original distribution was

$$P(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mu, \Sigma)$$

then the transformed variable **y** is distributed as:

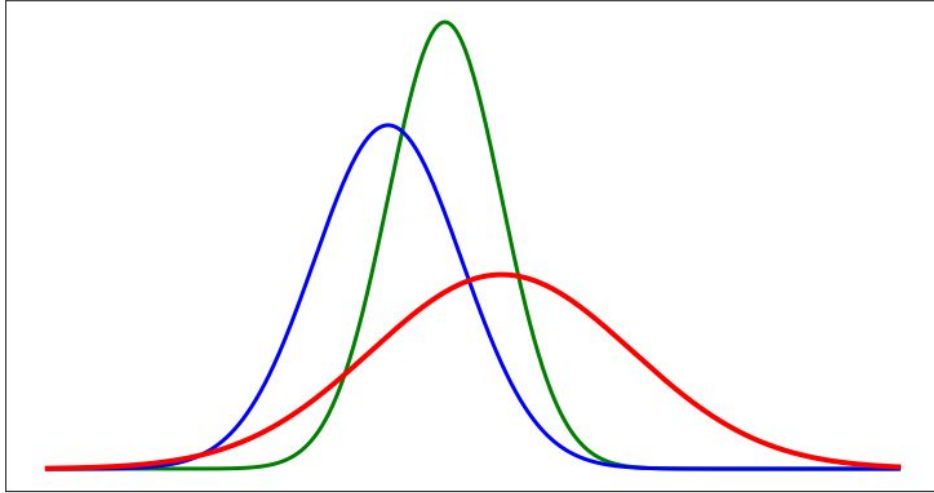$$Pr(\mathbf{y}) = \mathcal{N}(\mathbf{y}; \mathbf{A}\mu + \mathbf{b}, \mathbf{A}\Sigma\mathbf{A}^{T})$$

This relationship provides a simple method to draw samples from a normal distribution with mean μ and covariance Σ. We first draw a sample x from a standard normal distribution (μ = 0 and Σ = I) and then apply the transform y = $\Sigma^{1/2}$x + μ          ["*Reparameterization Trick*"]

Simon J.D. Prince, "Computer Vision:  Models, Learning, and Inference"

# Product of two Gaussians



The product of any two normals **N1** and **N2** is proportional to a third normal **N3** distribution, with a mean between the two original means and a variance that is smaller than either of the original distributions.

# Sum of two Gaussians



The sum of any two normals **N1** and **N2** is proportional to a third normal **N3** distribution, with a mean as a sum of two original means and a variance that is also sum of original variances

# Other nice properties



- KL divergence between two Gaussians has a nice form;
- the (inverse)Fourier transform of a Gaussian is another Gaussian
- the convolution of two Gaussians is another Gaussian

# Content of today's lecture

- Generative modeling

  - appetiser
  - approaches

- Diffusion models

  - intro
  - tale of marvelous Gaussians
  - more details

- Stable Diffusion
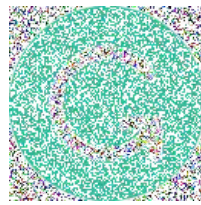
Structure ⟶ Uniform distribution

Structure ←———————— Uniform distribution

- Destroy all structure in data distribution using diffusion process

- Learn reversal of diffusion process

- *Reverse diffusion process* is the *model* of the data

$$q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

$$\mathbf{x}_0 \longrightarrow \cdots \longrightarrow \mathbf{x}_{t-1} \longrightarrow \mathbf{x}_t \longrightarrow \cdots \longrightarrow \mathbf{x}_T$$

$$q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

$$\mathbf{x} \in \mathbb{R}^D$$

[3, 128, 128] => D = 50k

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

Gaussian noise

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

$\mathbf{x} \in \mathbb{R}^D$

[3, 128, 128] => D = 50k

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

Gaussian noise

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$\mathbf{x}_0 \quad \cdots \quad \mathbf{x}_{t-1} \quad \mathbf{x}_t \quad \cdots \quad \mathbf{x}_T$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

$\mathbf{x} \in \mathbb{R}^2$

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

Gaussian noise

[-0.18, -0.49]          [-0.20, -0.62]          [0.12, 0.54]

$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$

$q(\mathbf{x}_t|\mathbf{x}_{t-1})$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

$$0 < \beta_1 < \beta_2 < ... < \beta_T < 1$$

$\mathcal{N}(\mathbf{x}; \mathbf{x}_0, \mathbf{0})$ - - - - - - - - - - - limiting cases - - - - - - - - - - - $\mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$

Gaussian noise

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

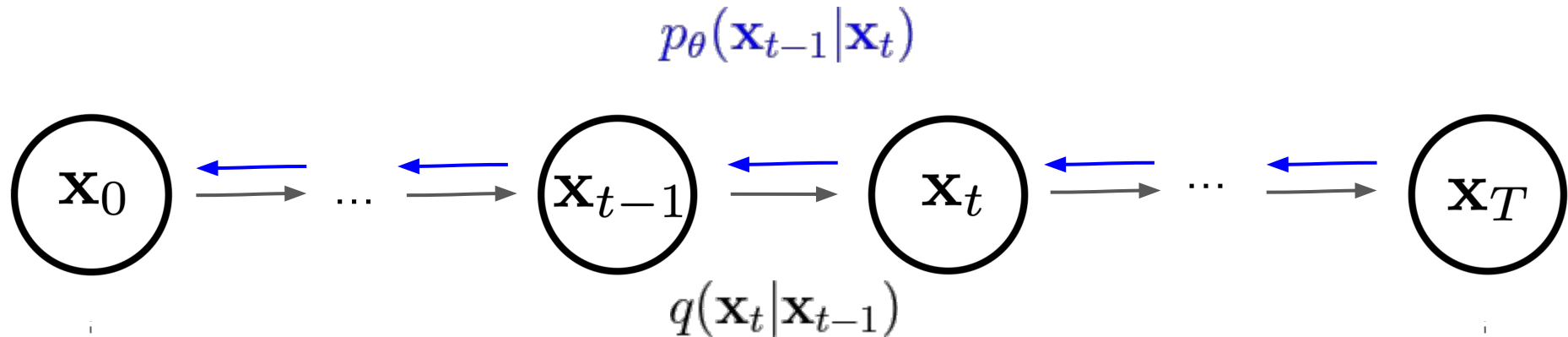$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

$$\bar{\alpha}_t := \prod_{s=1}^{T} \alpha_s \quad \alpha_t := 1 - \beta_t$$

$\mathcal{N}(\mathbf{x}; \mathbf{x}_0, \mathbf{0})$ ----- limiting cases ----- $\mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$

Gaussian noise

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$$

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$$
$$= \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$
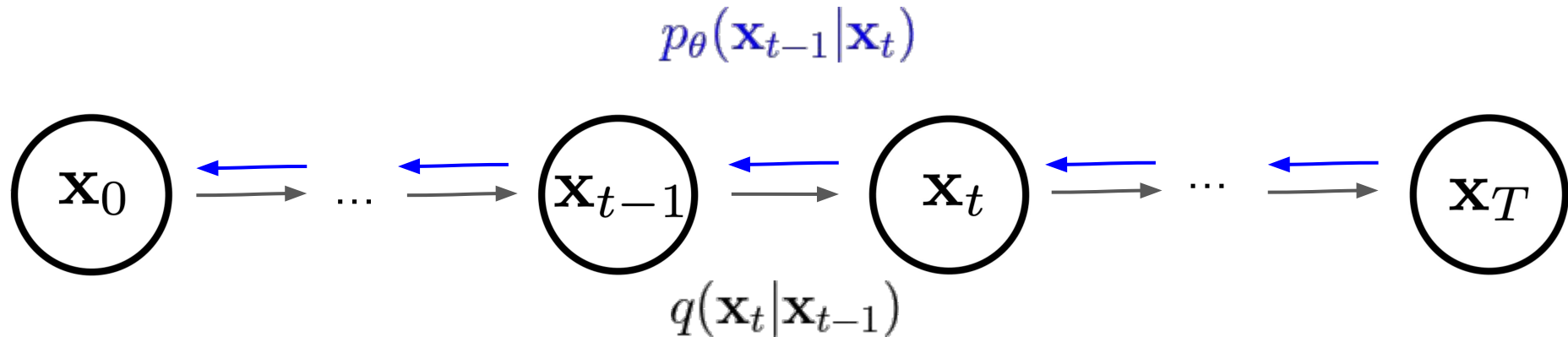
$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$$

$$= \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

Trainable network

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$
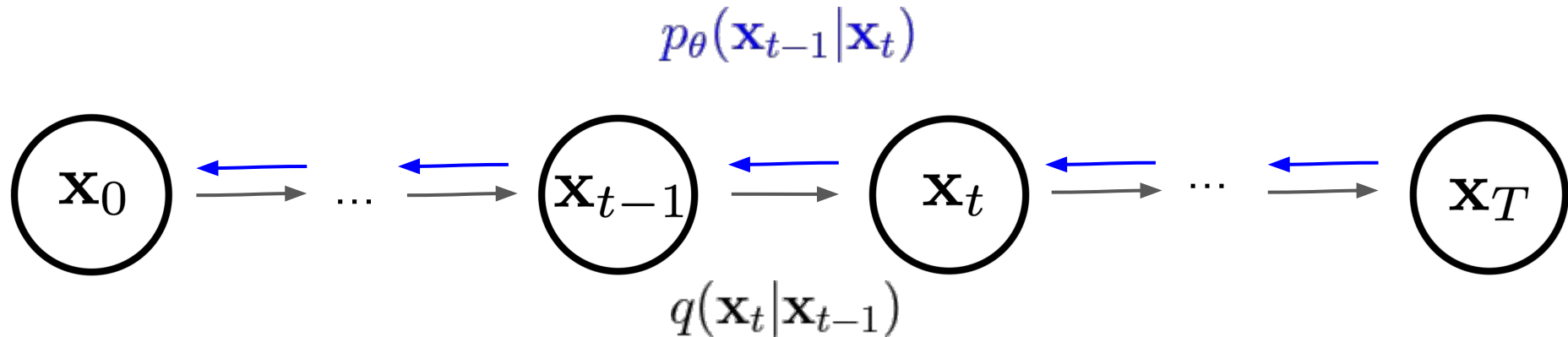
$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$$

$$= \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\alpha_t}\left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \cdot \underbrace{\epsilon(\mathbf{x}_t, t)}\right)$$

Trainable network

**Algorithm 1** Training

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

**Algorithm 2** Sampling

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

# Conditional generation

To explicit incorporate class information into the diffusion process, Dhariwal & Nichol (2021) trained a classifier $f_\phi(y|\mathbf{x}_t, t)$ on noisy image and use gradients $\nabla_\mathbf{x} \log f_\phi(y|\mathbf{x}_t, t)$ to guide the diffusion sampling process toward the target class label .

---

**Algorithm 2** Classifier guided DDIM sampling, given a diffusion model $\epsilon_\theta(x_t)$, classifier $f_\phi(y|x_t)$, and gradient scale $s$.

---

Input: class label $y$, gradient scale $s$
$x_T \leftarrow$ sample from $\mathcal{N}(0, \mathbf{I})$
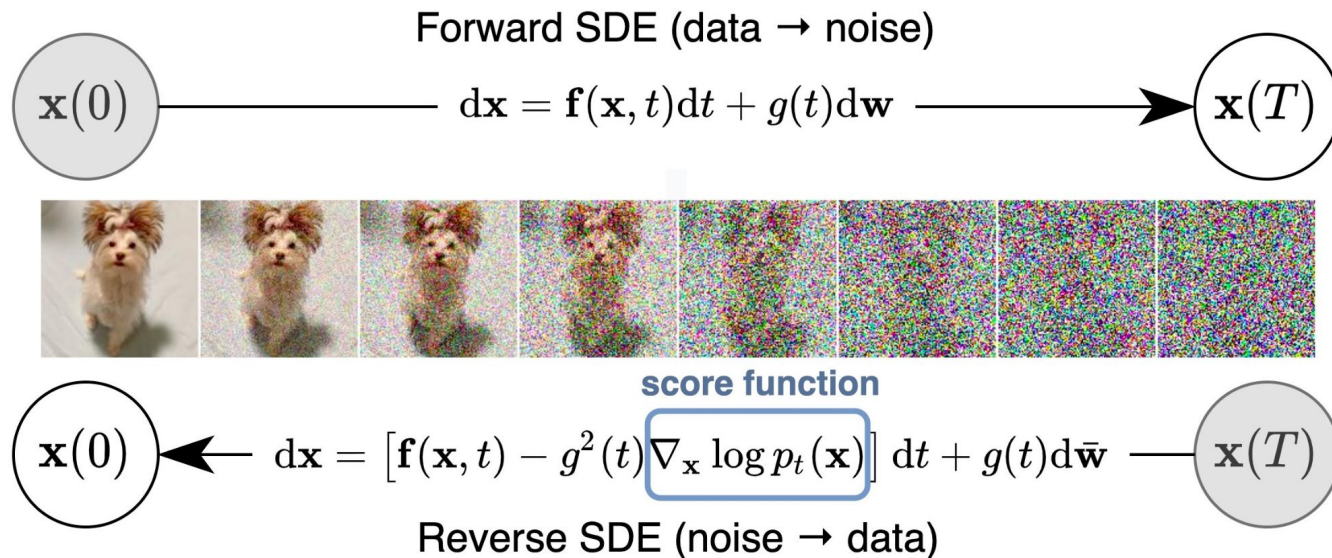**for all** $t$ from $T$ to 1 **do**
$\quad \hat{\epsilon} \leftarrow \epsilon_\theta(x_t) - \sqrt{1 - \bar{\alpha}_t} \, \nabla_{x_t} \log f_\phi(y|x_t)$
$\quad x_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \left( \frac{x_t - \sqrt{1-\bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}$
**end for**
**return** $x_0$

---

# Score-based formulation



Forward SDE (data → noise)

$$\mathrm{d}\mathbf{x} = \mathbf{f}(\mathbf{x}, t)\mathrm{d}t + g(t)\mathrm{d}\mathbf{w}$$

$\mathbf{x}(0)$      $\mathbf{x}(T)$

**score function**

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

$\mathbf{x}(0)$      $\mathbf{x}(T)$

Reverse SDE (noise → data)

[Generative Modeling by Estimating Gradients of the Data Distribution | Yang Song](#)

# Content of today's lecture

- Generative modeling

  - appetiser
  - approaches

- Diffusion models

  - intro
  - tale of marvelous Gaussians
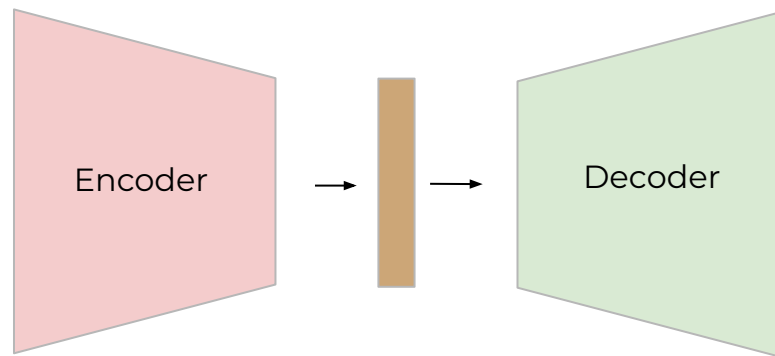  - more details

- Stable Diffusion

# Stable Diffusion



[2112.10752] High-Resolution Image Synthesis with Latent Diffusion Models

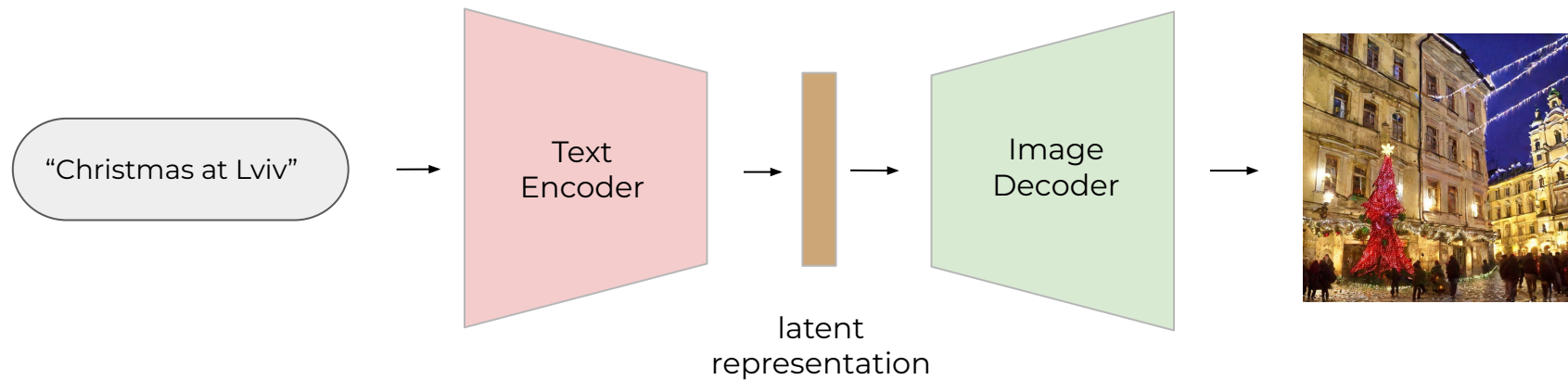GitHub - CompVis/stable-diffusion: A latent text-to-image diffusion model

Official Announcement:
Stable Diffusion launch announcement — Stability AI

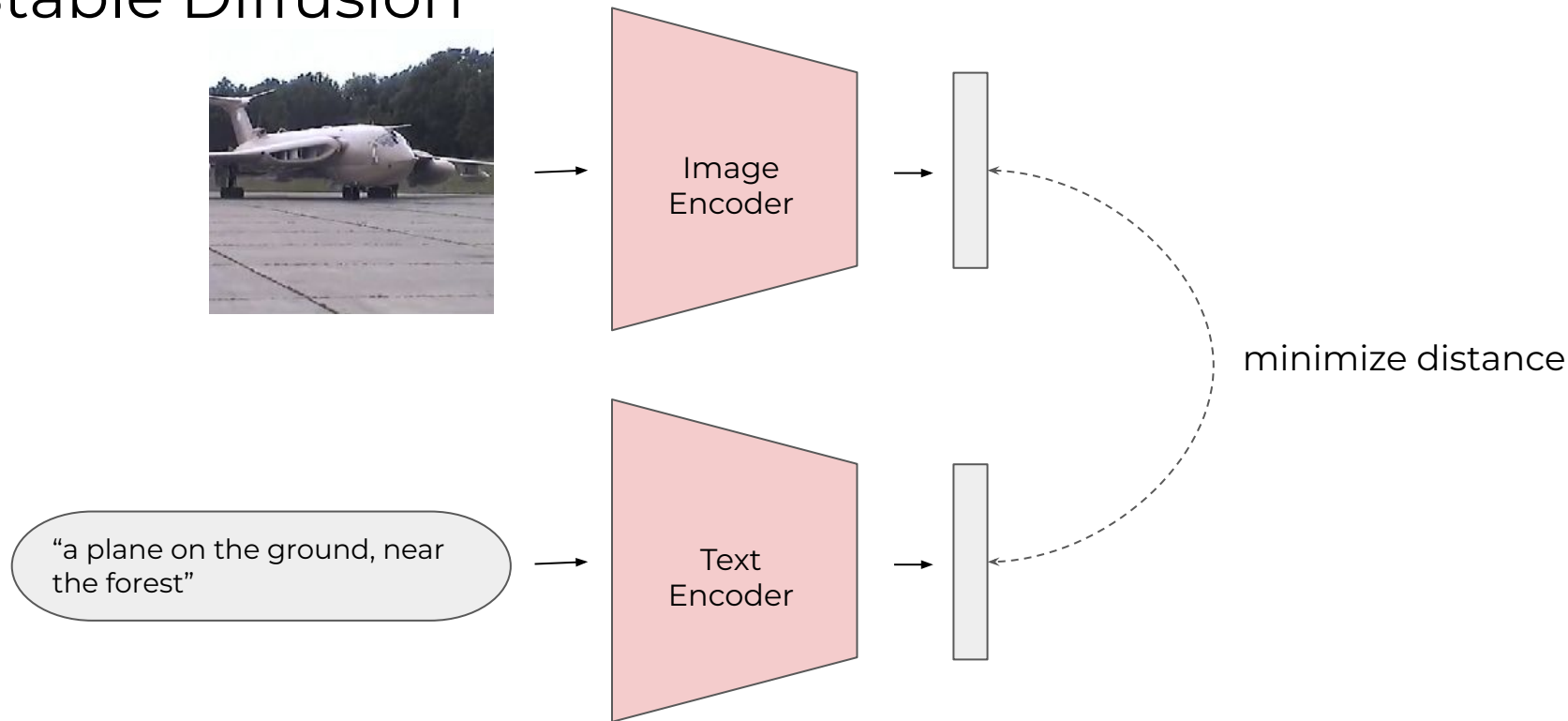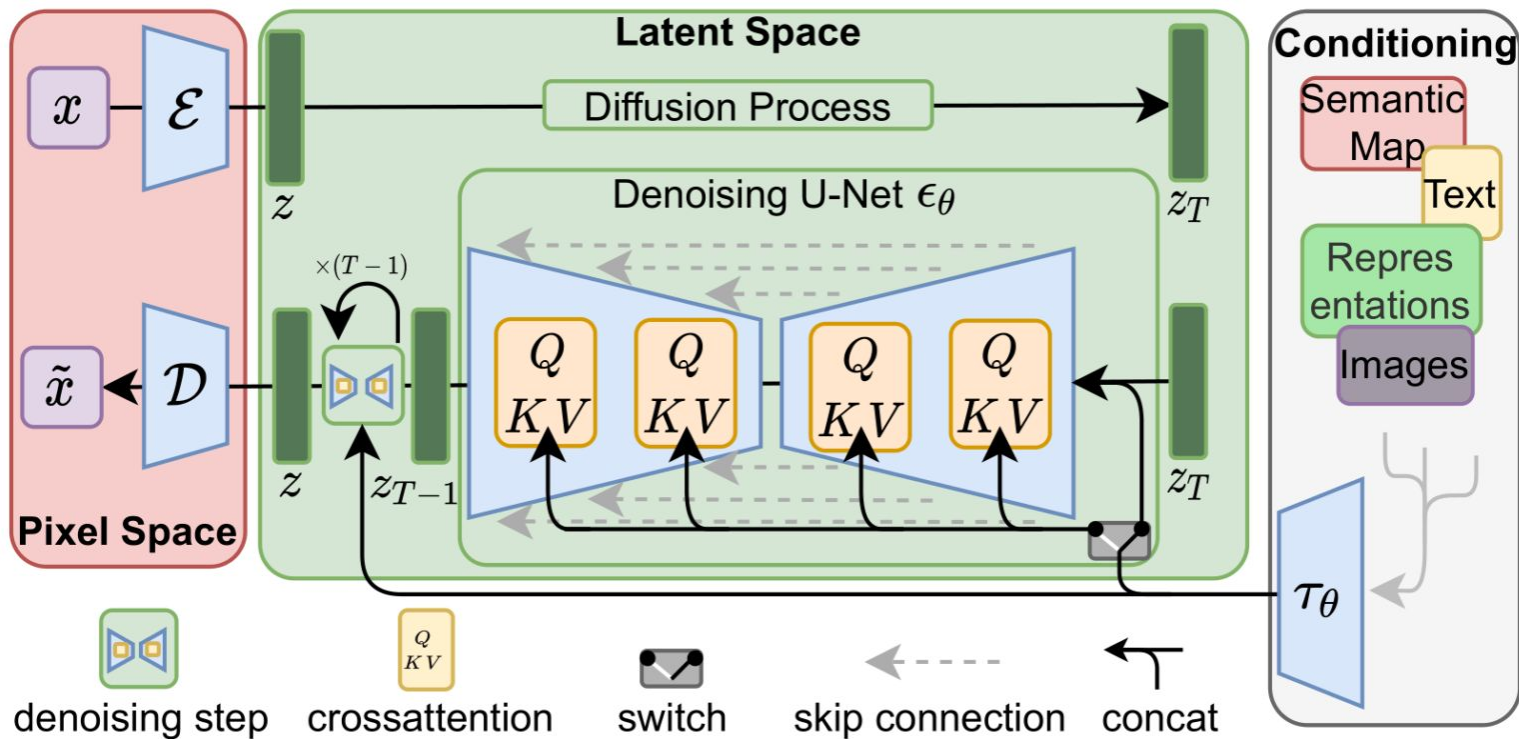# Stable Diffusion

# Stable Diffusion

# Stable Diffusion



minimize distance

"a plane on the ground, near the forest"

GitHub - openai/CLIP: Contrastive Language-Image Pretraining
GitHub - mlfoundations/open_clip: An open source implementation of CLIP.

# Stable Diffusion



[2112.10752] High-Resolution Image Synthesis with Latent Diffusion Models

# Summary

- We are witnessing a revolution in computer-assisted creativity.

- Generation won't only be limited to images, already now there are video, 3D generation, and more is coming.