

Why AI needs Deep Generative Modeling?

Jakub M. Tomczak

What is **intelligence?**

What is **intelligence?**

...

INFORMATION, INTELLIGENCE AND ARTIFICIAL INTELLIGENCE

What is **intelligence**?

...

What is **artificial intelligence**?

INFORMATION, INTELLIGENCE AND ARTIFICIAL INTELLIGENCE

What is **intelligence**?

...

What is **artificial intelligence**?



INFORMATION, INTELLIGENCE AND ARTIFICIAL INTELLIGENCE

What is **intelligence**?

...

What is **artificial intelligence**?

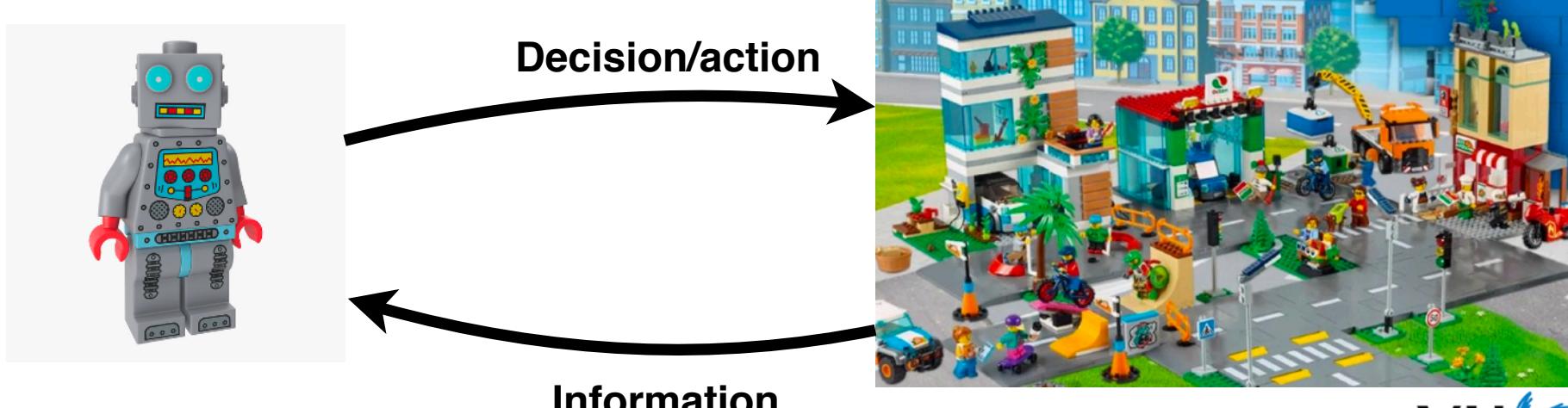


INFORMATION, INTELLIGENCE AND ARTIFICIAL INTELLIGENCE

What is **intelligence**?

...

What is **artificial intelligence**?

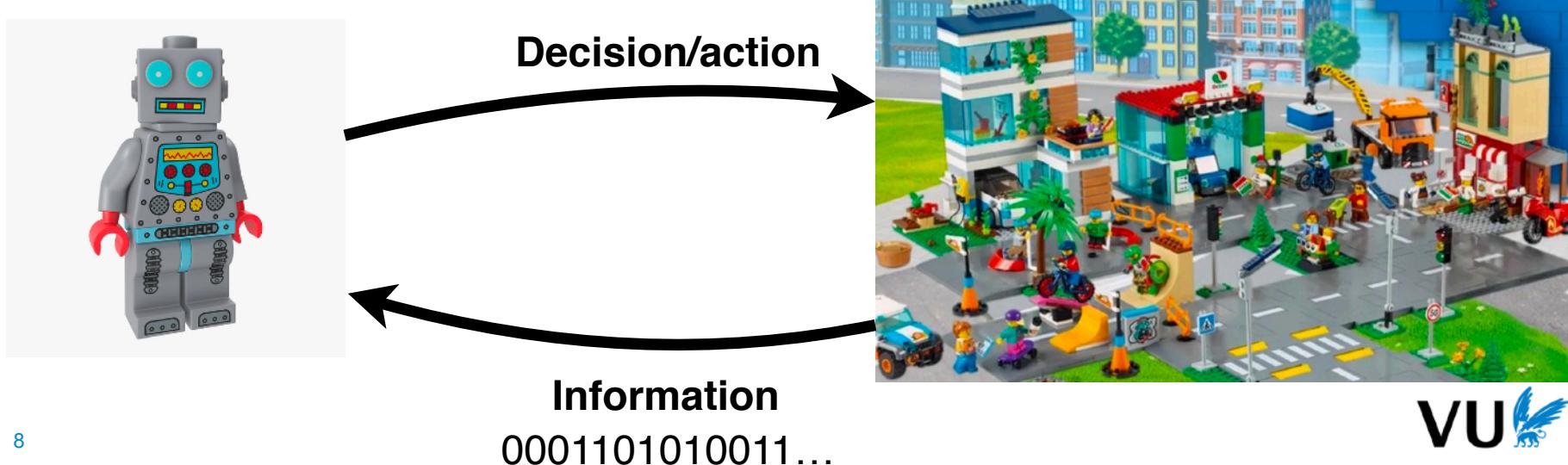


INFORMATION, INTELLIGENCE AND ARTIFICIAL INTELLIGENCE

What is **intelligence**?

...

What is **artificial intelligence**?



What is **artificial intelligence**?

- **Information** processing
- **Information** storing
- **Information** transmission



What is **artificial intelligence**?

- **Information** processing
- **Information** storing
- **Information** transmission
- **Decision** making



INFORMATION, INTELLIGENCE AND ARTIFICIAL INTELLIGENCE

What is **artificial intelligence**?

- **Information** processing
 - **Information** storing
 - **Information** transmission
 - **Decision** making
- } Learning
Knowledge representation
Models...



What is **artificial intelligence**?

- **Information** processing
 - **Information** storing
 - **Information** transmission
 - **Decision** making
- } Learning
Knowledge representation
Models...



The question is how to formalize the problem of AI?

INFORMATION, INTELLIGENCE AND ARTIFICIAL INTELLIGENCE

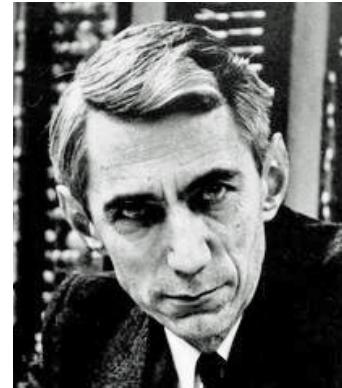
Information (a quick recap)



Claude Shannon

Information (a quick recap)

We have a random source of data x .



Claude Shannon

Information (a quick recap)

We have a random source of data x .

We can quantify the **uncertainty** of this source by calculating **the entropy**:



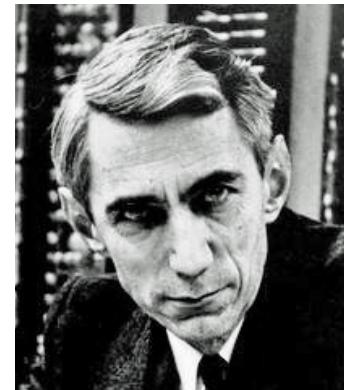
Claude Shannon

$$\mathbb{H}[x] = - \sum_x p(x) \log p(x)$$

Information (a quick recap)

We have a random source of data x .

We can quantify the **uncertainty** of this source by calculating **the entropy**:



Claude Shannon

$$\mathbb{H}[x] = - \sum_x p(x) \log p(x)$$

Entropy is max if all x 's are equiprobable.

Entropy is min if the probability of one value is 1.

Information (a quick recap)

We have a random source of data x .

We can quantify the **uncertainty** of this source by calculating **the entropy**:



Claude Shannon

$$\mathbb{H}[x] = - \sum_x p(x) \log p(x)$$

Optimal message length \approx the entropy.

Information (a quick recap)

We have two random sources: x and y .

We can quantify the **uncertainty** of them by calculating
the joint entropy:



Claude Shannon

$$\mathbb{H}[x, y] = - \sum_{x,y} p(x, y) \log p(x, y)$$

or **the conditional entropy**:

$$\mathbb{H}[y | x] = - \sum_{x,y} p(x, y) \log p(y | x)$$

Mutual Information (a quick recap)

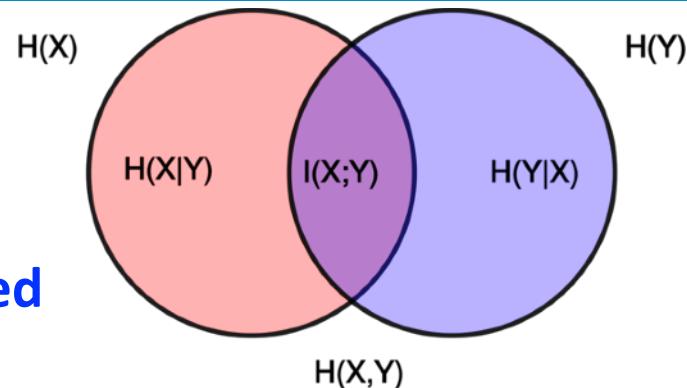
We have two random sources: x and y .

Mutual Information (a quick recap)

We have two random sources: x and y .

We can quantify how much **information is shared
by the two sources:**

$$I[x; y] = H[y] - H[y | x]$$

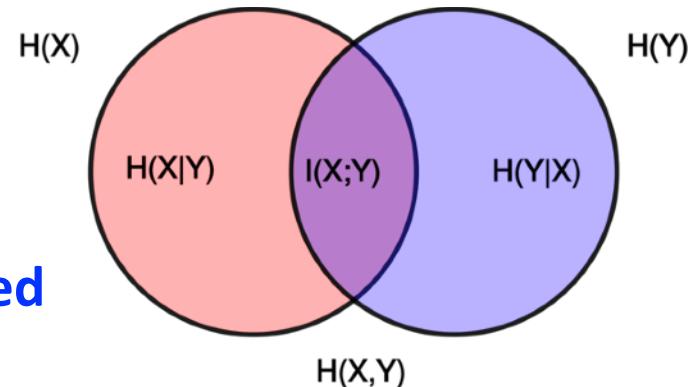


INFORMATION, INTELLIGENCE AND ARTIFICIAL INTELLIGENCE

Mutual Information (a quick recap)

We have two random sources: x and y .

We can quantify how much **information is shared by the two sources:**



$$I[x; y] = H[y] - H[y | x]$$

or **how much knowing one source reduces uncertainty about the other.**

We have two random sources: x (e.g., images) and y (e.g., decisions).

We have two random sources: x (e.g., images) and y (e.g., decisions).

We have also a **model** m (a representation of a world).

We have two random sources: x (e.g., images) and y (e.g., decisions).

We have also a **model** m (a representation of a world).

The **goal** of AI is to **maximize** the **mutual information** between (x, y) and m :

$$\mathbb{I}[(x, y); m] = \mathbb{H}[x, y] - \mathbb{H}[x, y | m]$$

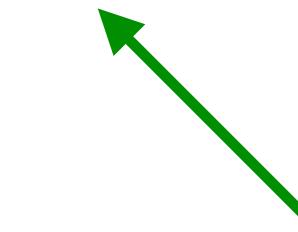
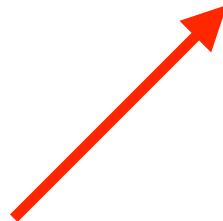
We have two random sources: x (e.g., images) and y (e.g., decisions).

We have also a **model** m (a representation of a world).

The **goal** of AI is to **maximize** the **mutual information** between (x, y) and m :

$$\mathbb{I}[(x, y); m] = \mathbb{H}[x, y] - \mathbb{H}[x, y | m]$$

Entropy of the world
(model has no influence on that)



That's the “real” goal!

The **goal** of AI is to **maximize** the **mutual information** between (x, y) and m

(or minimize $\mathbb{H}[x, y \mid m]$, i.e., minimize uncertainty of the world):

$$\mathbb{H}[x, y \mid m] = \sum_{x,y,m} p(x, y, m) [\log p(y \mid x, m) + \log p(x \mid m)]$$

The **goal** of AI is to **maximize** the **mutual information** between (x, y) and m

(or minimize $\mathbb{H}[x, y | m]$, i.e., minimize uncertainty of the world):

$$\mathbb{H}[x, y | m] = \sum_{x,y,m} p(x, y, m) [\log p(y | x, m) + \log p(x | m)]$$



A model for
decision making



A model for
understanding
the world.

The **goal** of AI is to **maximize** the **mutual information** between (x, y) and m (or minimize $\mathbb{H}[x, y \mid m]$, i.e., minimize uncertainty of the world).

In order to achieve that, AI should focus on learning **two models**:

- **A model for decision making:** $p(y \mid x, m)$
- **A model for understanding the world:** $p(x \mid m)$

WHAT HAPPENS IF WE LEARN ONLY DECISION MAKING

The bulk of AI is focused on the decision making part **only**!

WHAT HAPPENS IF WE LEARN ONLY DECISION MAKING

The bulk of AI is focused on the decision making part **only**!

Example: Let's say we have a model that is well trained.



$$p(y = \text{cat}|\mathbf{x}) = 0.90$$

$$p(y = \text{dog}|\mathbf{x}) = 0.05$$

$$p(y = \text{horse}|\mathbf{x}) = 0.05$$

WHAT HAPPENS IF WE LEARN ONLY DECISION MAKING

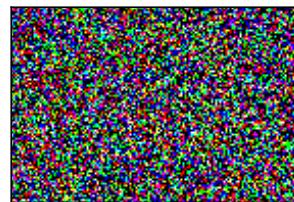
The bulk of AI is focused on the decision making part **only**!

Example: Let's say we have a model that is well trained.



$$\begin{aligned} p(y = \text{cat}|\mathbf{x}) &= 0.90 \\ p(y = \text{dog}|\mathbf{x}) &= 0.05 \\ p(y = \text{horse}|\mathbf{x}) &= 0.05 \end{aligned}$$

+



noise

=



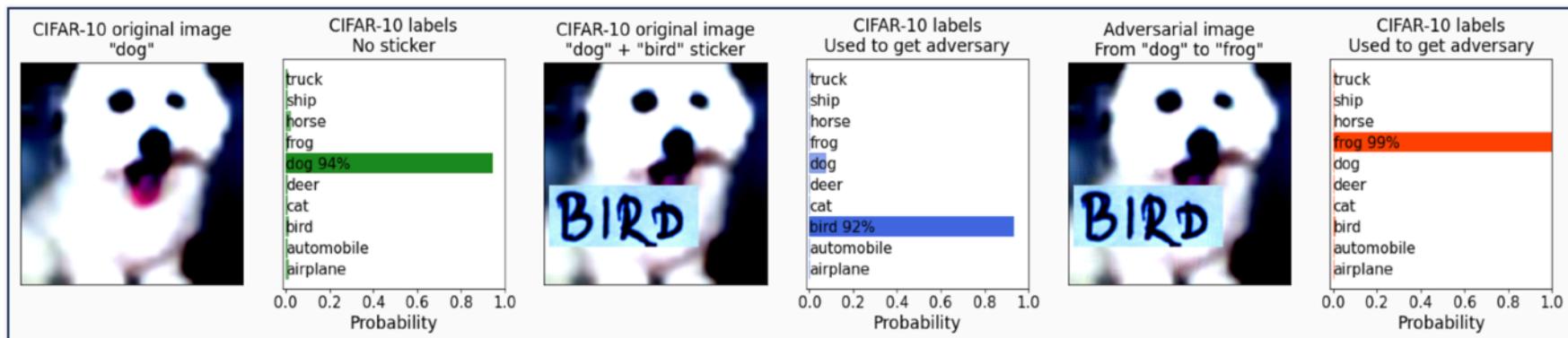
$$\begin{aligned} p(y = \text{cat}|\mathbf{x}) &= 0.05 \\ p(y = \text{dog}|\mathbf{x}) &= 0.05 \\ p(y = \text{horse}|\mathbf{x}) &= 0.90 \end{aligned}$$

But after adding a little noise it could fail completely...

IS LEARNING CLASSIFIERS ENOUGH?

Let's assume we have a perfectly trained neural net.

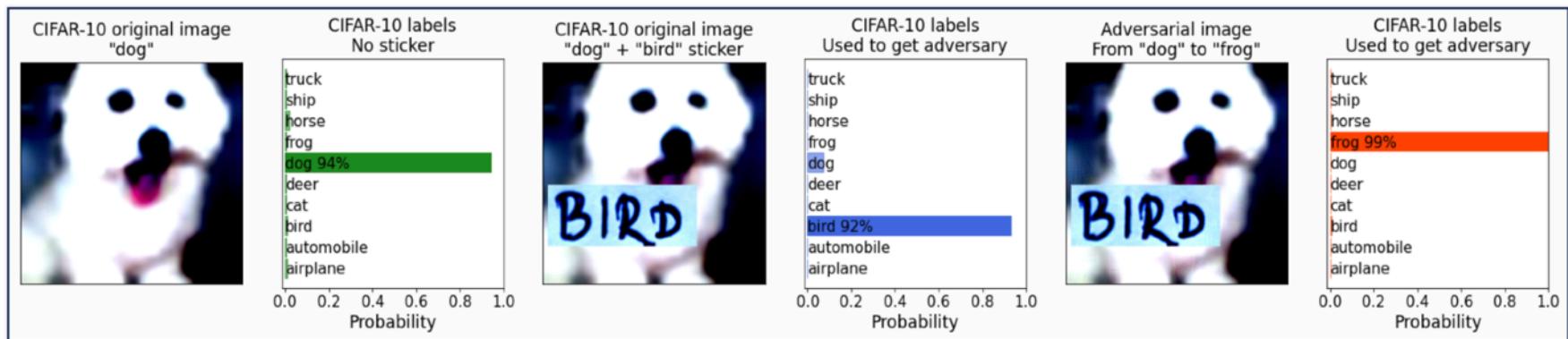
What happens if we add (adversarial) noise to an image?



IS LEARNING CLASSIFIERS ENOUGH?

Let's assume we have a perfectly trained neural net.

What happens if we add (adversarial) noise to an image?



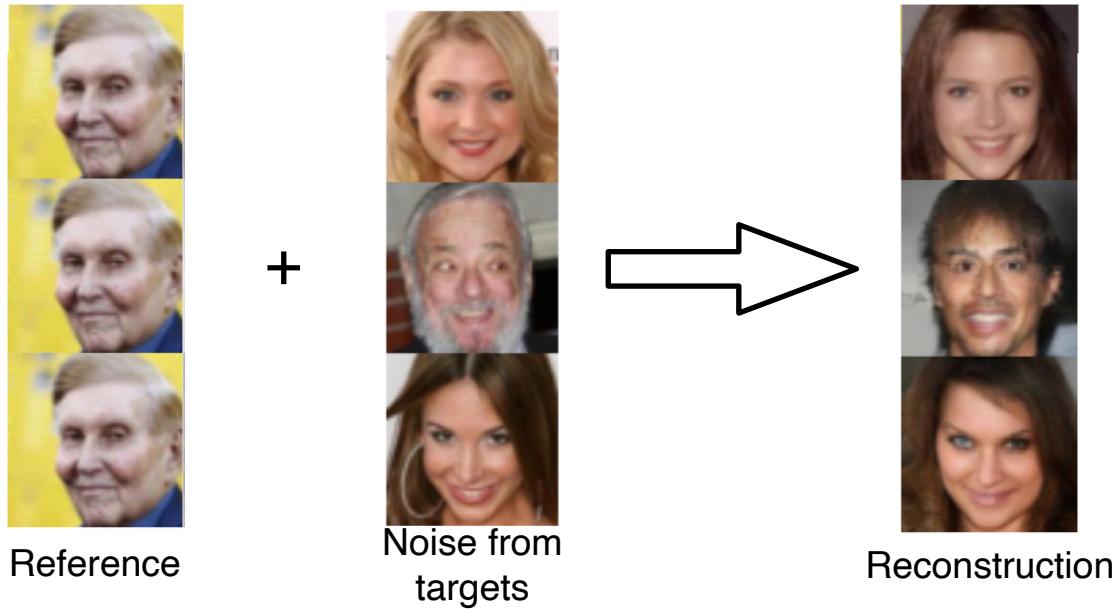
It fails completely...

³³S. Fort, "Pixels still beat text: Attacking the OpenAI CLIP model with text patches and adversarial pixel perturbations", [\[Link\]](#)

IS LEARNING CLASSIFIERS ENOUGH?

Let's assume we have a perfectly trained neural net.

What happens if we add (adversarial) noise to an image?

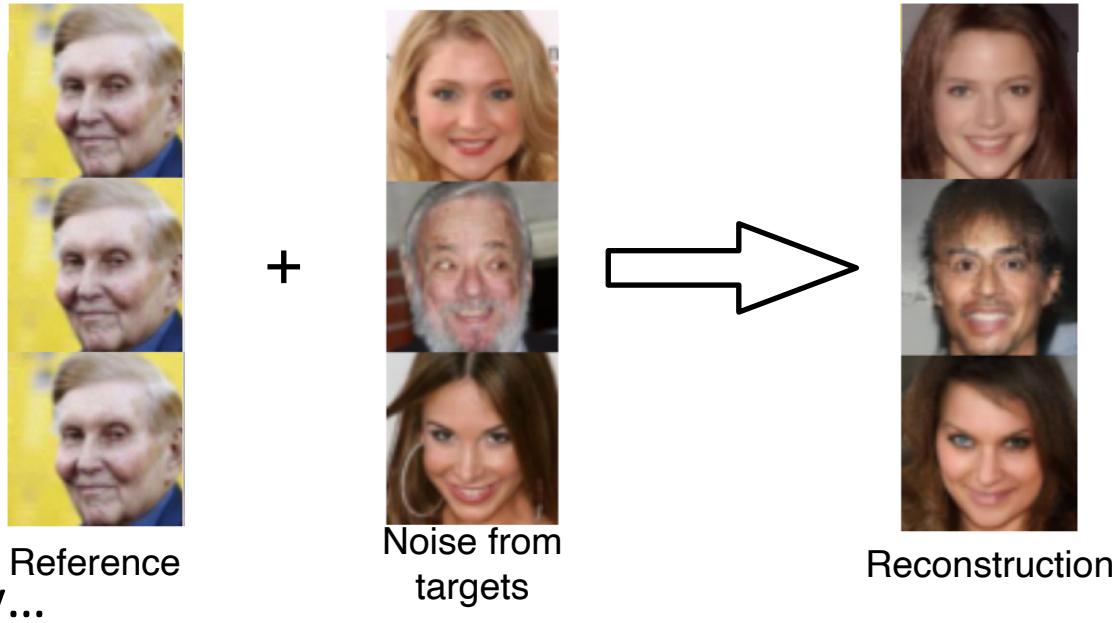


³⁴A. Kuzina, M. Welling, J.M. Tomczak, “Diagnosing Vulnerability of Variational Auto-Encoders to Adversarial Attacks”, [\[Link\]](#)

IS LEARNING CLASSIFIERS ENOUGH?

Let's assume we have a perfectly trained neural net.

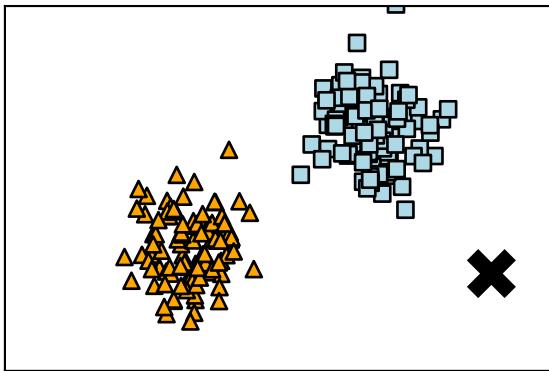
What happens if we add (adversarial) noise to an image?



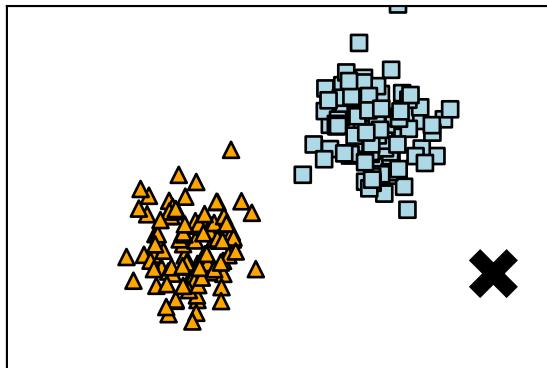
It fails completely...

³⁵A. Kuzina, M. Welling, J.M. Tomczak, “Diagnosing Vulnerability of Variational Auto-Encoders to Adversarial Attacks”, [\[Link\]](#)

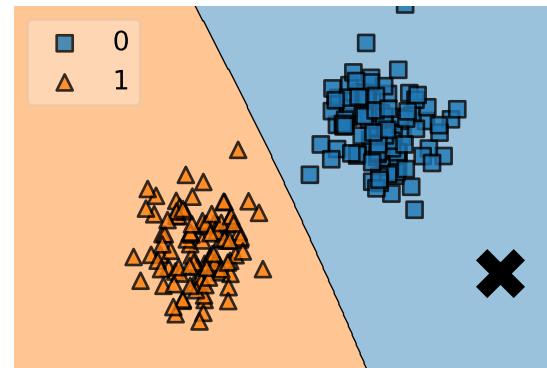
DEEP GENERATIVE MODELING: WHY DO WE NEED THEM?



DEEP GENERATIVE MODELING: WHY DO WE NEED THEM?



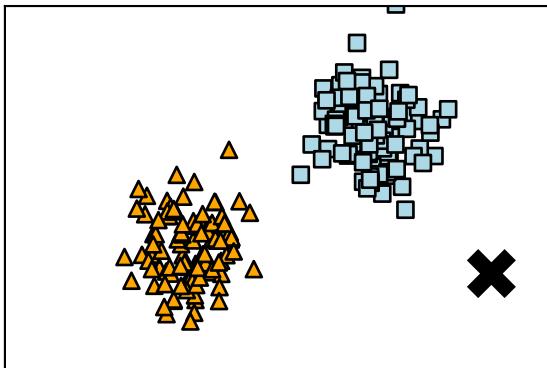
Data



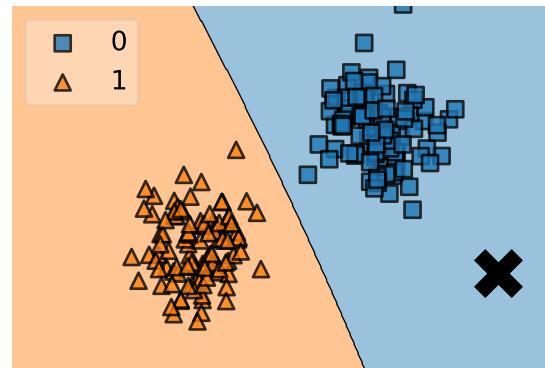
$p(y|\mathbf{x})$

$p(\text{blue}|\mathbf{x})$ is high
= certain decision!

DEEP GENERATIVE MODELING: WHY DO WE NEED THEM?

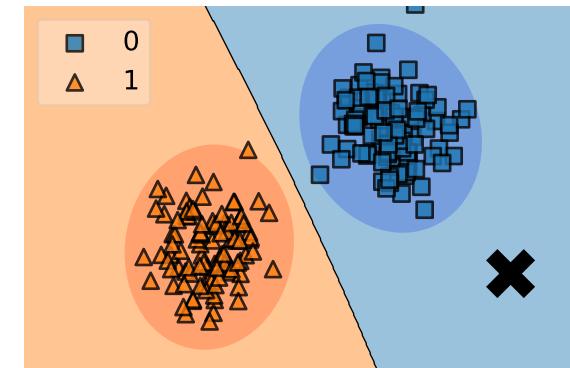


Data



$p(y|\mathbf{x})$

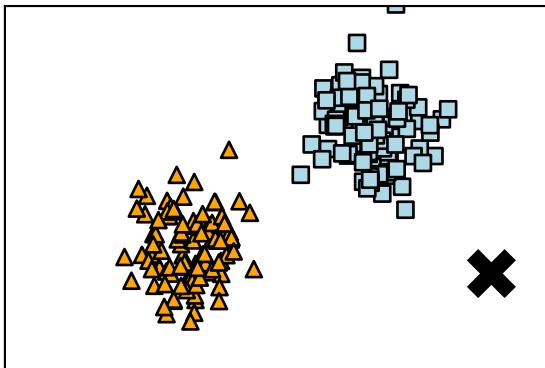
$p(\text{blue}|\mathbf{x})$ is high
= certain decision!



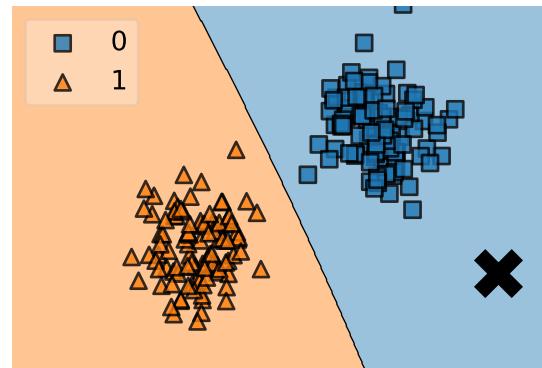
$p(\mathbf{x}, y) = p(y|\mathbf{x}) p(\mathbf{x})$

$p(\text{blue}|\mathbf{x})$ is high
and $p(\mathbf{x})$ is low
= uncertain decision!

DEEP GENERATIVE MODELING: WHY DO WE NEED THEM?

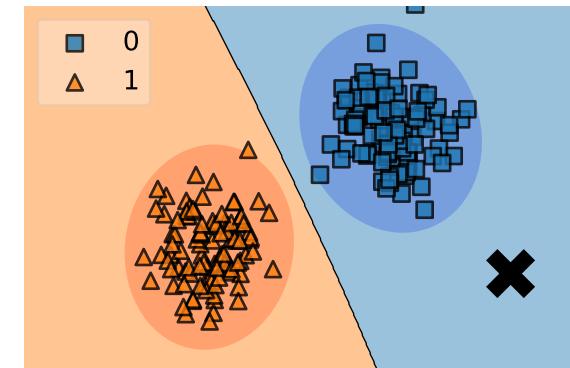


Data



$p(y|x)$

$p(\text{blue}|\mathbf{x})$ is high
= certain decision!



$p(\mathbf{x}, y) = p(y|\mathbf{x}) p(\mathbf{x})$

$p(\text{blue}|\mathbf{x})$ is high
and $p(\mathbf{x})$ is low
= uncertain decision!

Thus, learning the conditional is only a part of the story!
How can we learn $p(\mathbf{x})$?

DEEP GENERATIVE MODELING: WHY DO WE NEED THEM?

We clearly see that training a neural network (i.e., a conditional distribution):

$$p(y | \mathbf{x}) = \text{softmax} (\text{NN}(\mathbf{x}))$$

is **not enough!**



Granny Smith	0.1%
iPod	99.7%
library	0.0%
pizza	0.0%
toaster	0.0%
dough	0.0%

DEEP GENERATIVE MODELING: WHY DO WE NEED THEM?

We clearly see that training a neural network (i.e., a conditional distribution):

$$p(y | \mathbf{x}) = \text{softmax} (\text{NN}(\mathbf{x}))$$

is **not enough!**

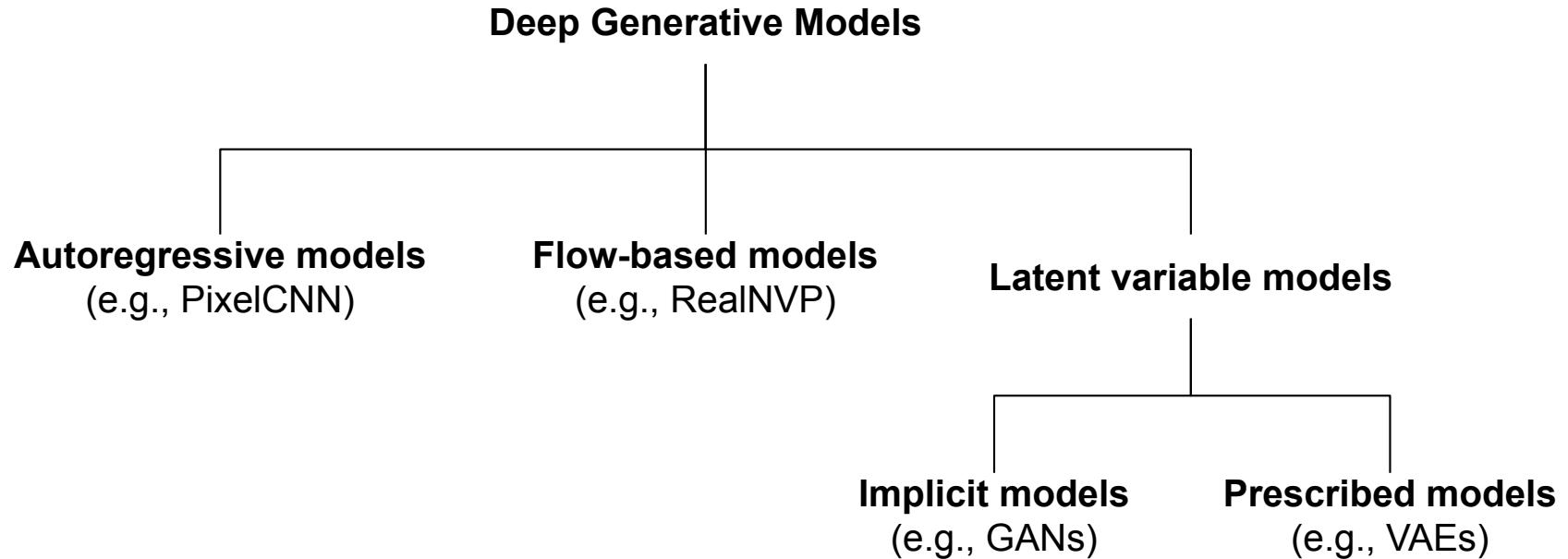
What can we do then?

Or, how to modify the **wrong certainty?**

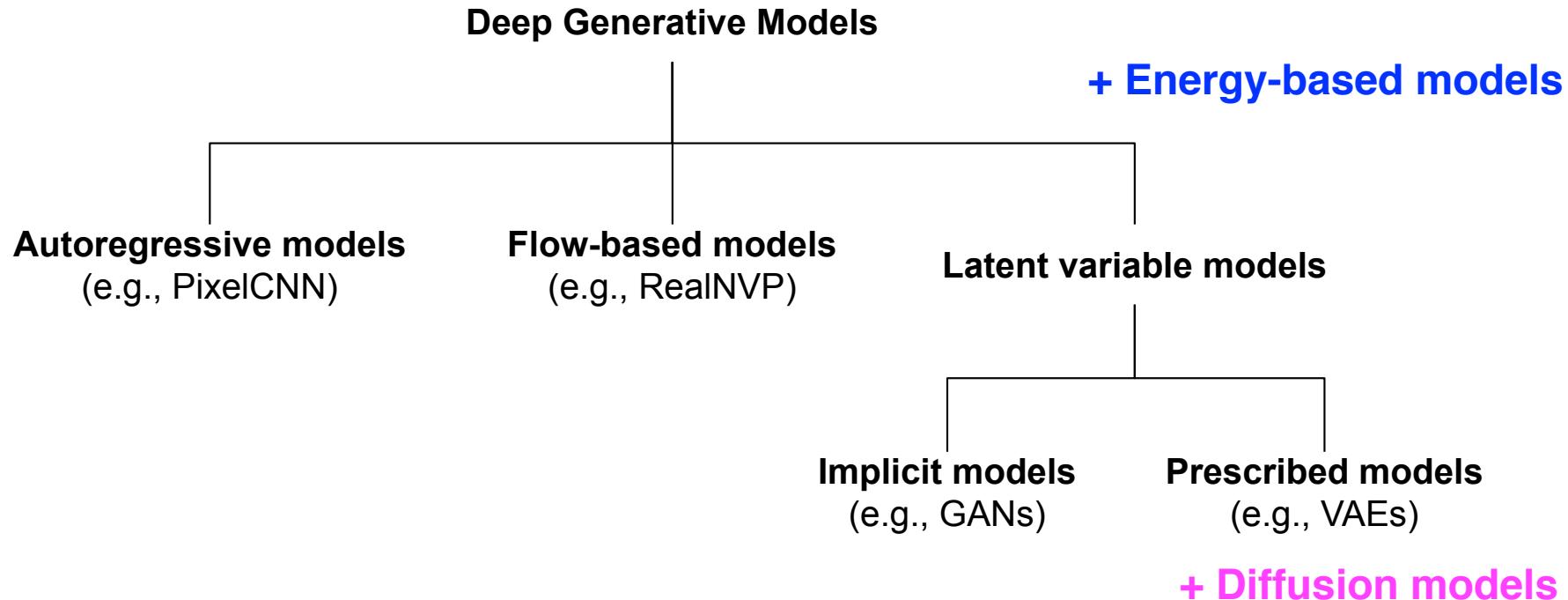


Granny Smith	0.1%
iPod	99.7%
library	0.0%
pizza	0.0%
toaster	0.0%
dough	0.0%

DEEP GENERATIVE MODELING: HOW WE CAN FORMULATE IT?



DEEP GENERATIVE MODELING: HOW WE CAN FORMULATE IT?



DEEP GENERATIVE MODELING: HOW WE CAN FORMULATE IT?

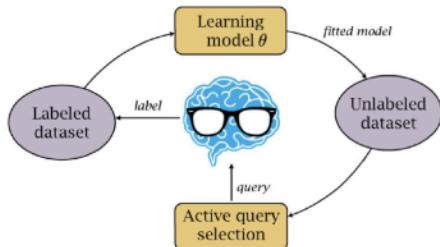
Generative models	Training	Likelihood	Sampling	Lossy compression	Lossless compression
Autoregressive models	stable	exact	slow	no	yes
Flow-based models	stable	exact	fast/slow	no	yes
Implicit models	unstable	no	fast	no	no
Prescribed model	stable	approximate	fast	yes	no

DEEP GENERATIVE MODELING: WHERE CAN WE USE IT?

“ i want to talk to you . ”
“ want to be with you . ”
“ i do n’t want to be with you . ”
i do n’t want to be with you .
she did n’t want to be with him .

he was silent for a long moment .
he was silent for a moment .
it was quiet for a moment .
it was dark and cold .
there was a pause .
it was my turn .

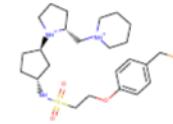
Text analysis



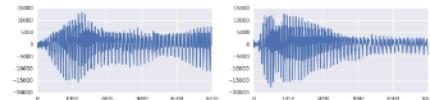
Active Learning



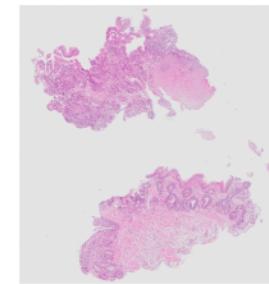
Image analysis



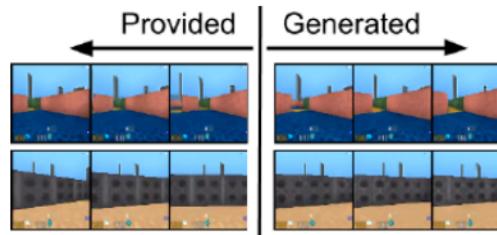
Graph analysis



Audio analysis



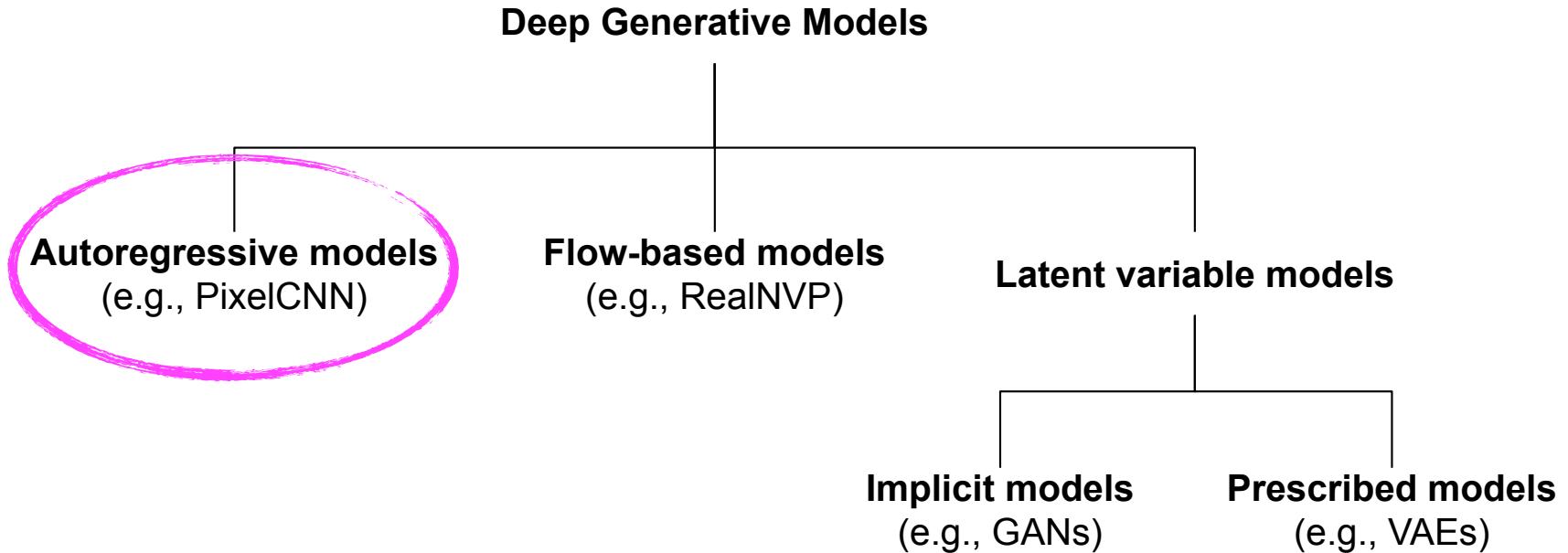
Medical data



Reinforcement Learning

and more...

DEEP GENERATIVE MODELING: HOW WE CAN FORMULATE IT?



AUTOREGRESSIVE MODELS (ARMS)

We are interested in modeling

$$p(\mathbf{x})$$

where $\mathbf{x} \in \{0, 1, \dots, 255\}^{D \times 3}$ is an RGB image (for instance).

AUTOREGRESSIVE MODELS (ARMS)

We are interested in modeling

$$p(\mathbf{x})$$

where $\mathbf{x} \in \{0, 1, \dots, 255\}^{D \times 3}$ is an RGB image (for instance).

We can use the **product rule**:

$$p(\mathbf{x}) = p(x_1) \prod_{d=2}^D p(x_d | \mathbf{x}_{<d})$$

where $\mathbf{x}_{<d} = [x_1, x_2, \dots, x_{d-1}]^\top$

AUTOREGRESSIVE MODELS (ARMS)

We are interested in modeling

$$p(\mathbf{x})$$

where $\mathbf{x} \in \{0, 1, \dots, 255\}^{D \times 3}$ is an RGB image (for instance).

We can use the **product rule**:

$$p(\mathbf{x}) = p(x_1) \prod_{d=2}^D p(x_d | \mathbf{x}_{<d})$$

where $\mathbf{x}_{<d} = [x_1, x_2, \dots, x_{d-1}]^\top$

Example:

$$p(\mathbf{x}) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2)$$

AUTOREGRESSIVE MODELS (ARMS)

We are interested in modeling

$$p(\mathbf{x})$$

where $\mathbf{x} \in \{0, 1, \dots, 255\}^{D \times 3}$ is an RGB image (for instance).

We can use the **product rule**:

$$p(\mathbf{x}) = p(x_1) \prod_{d=2}^D p(x_d | \mathbf{x}_{<d})$$

where $\mathbf{x}_{<d} = [x_1, x_2, \dots, x_{d-1}]^\top$

Training objective:

$$\ln p(\mathbf{x}) = \ln p(x_1) + \sum_{d=2}^D \ln p(x_d | \mathbf{x}_{<d})$$

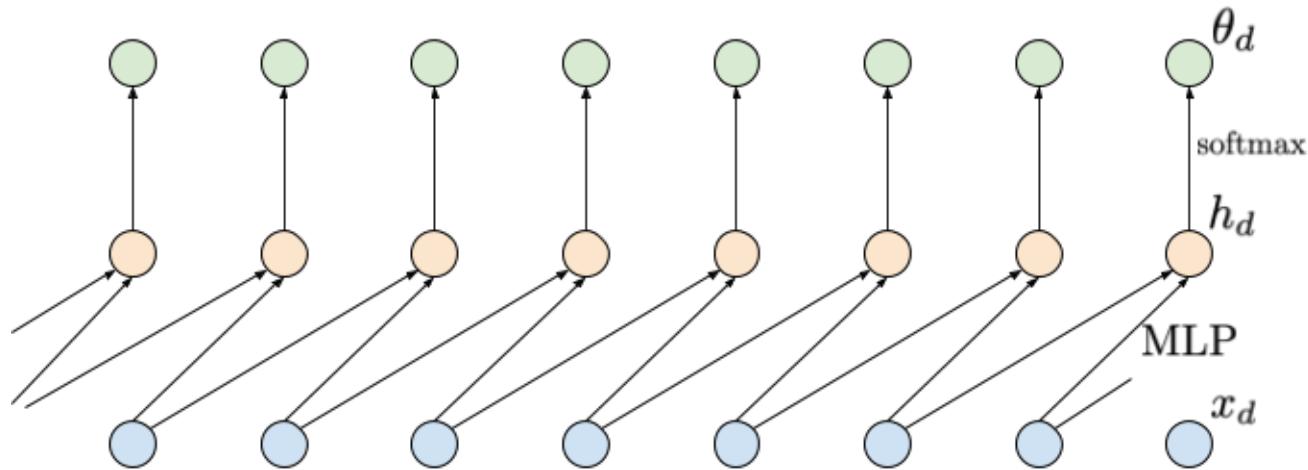
AUTOREGRESSIVE MODELS (ARMS)

Question: How we can model the conditionals $p(x_d | \mathbf{x}_{<d})$ **efficiently**?

AUTOREGRESSIVE MODELS (ARMS)

Question: How we can model the conditionals $p(x_d | \mathbf{x}_{<d})$ **efficiently**?

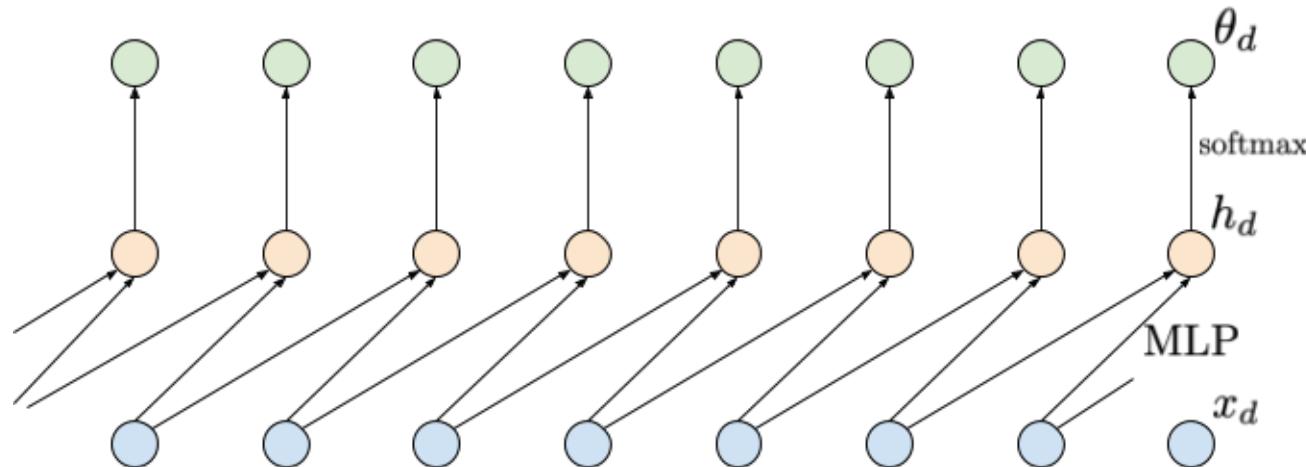
Approach 1: Finite memory



AUTOREGRESSIVE MODELS (ARMS)

Question: How we can model the conditionals $p(x_d | \mathbf{x}_{<d})$ **efficiently**?

Approach 1: Finite memory



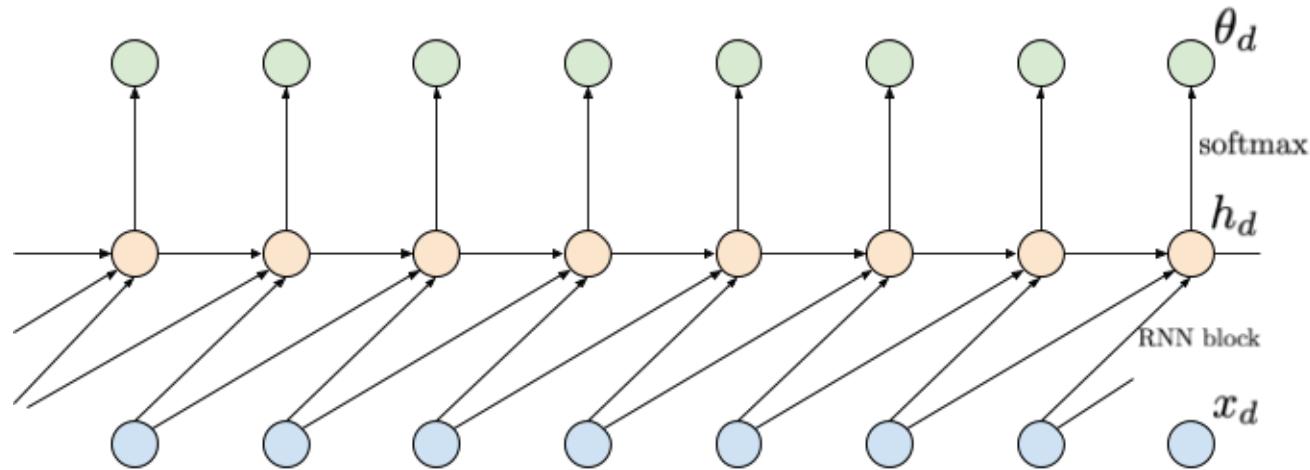
Easy!

Limited dependencies!
How many we should take?

AUTOREGRESSIVE MODELS (ARMS)

Question: How we can model the conditionals $p(x_d | \mathbf{x}_{<d})$ **efficiently**?

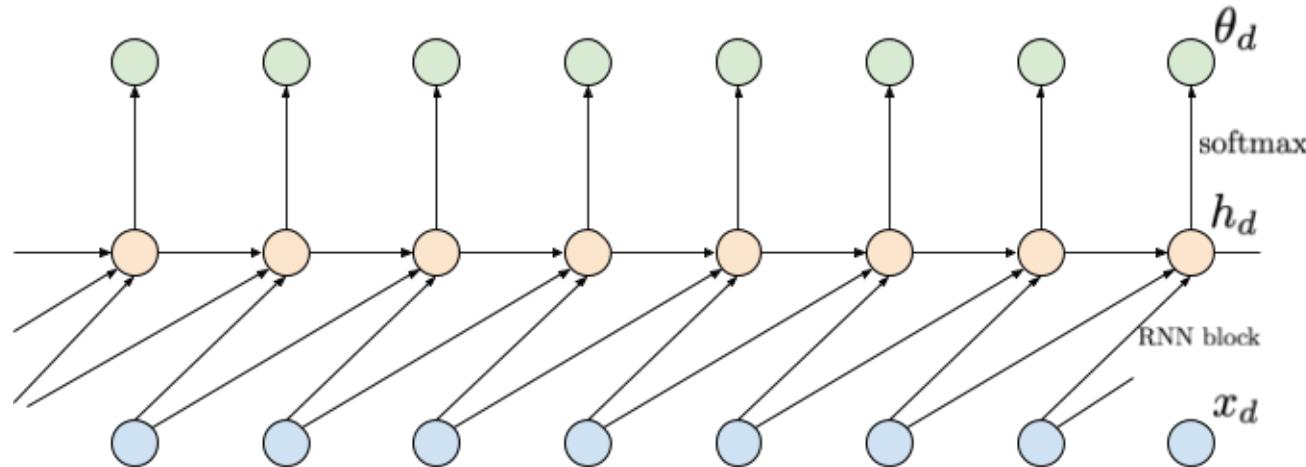
Approach 2: Long-range memory with RNNs



AUTOREGRESSIVE MODELS (ARMS)

Question: How we can model the conditionals $p(x_d | \mathbf{x}_{<d})$ **efficiently**?

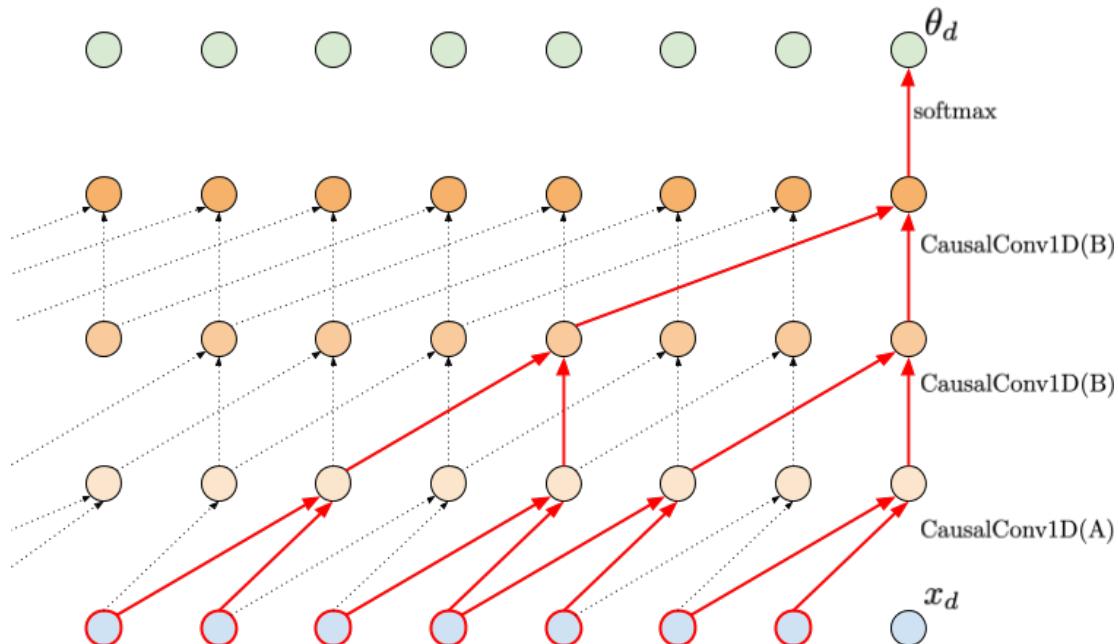
Approach 2: Long-range memory with RNNs



AUTOREGRESSIVE MODELS (ARMS)

Question: How we can model the conditionals $p(x_d | \mathbf{x}_{<d})$ **efficiently**?

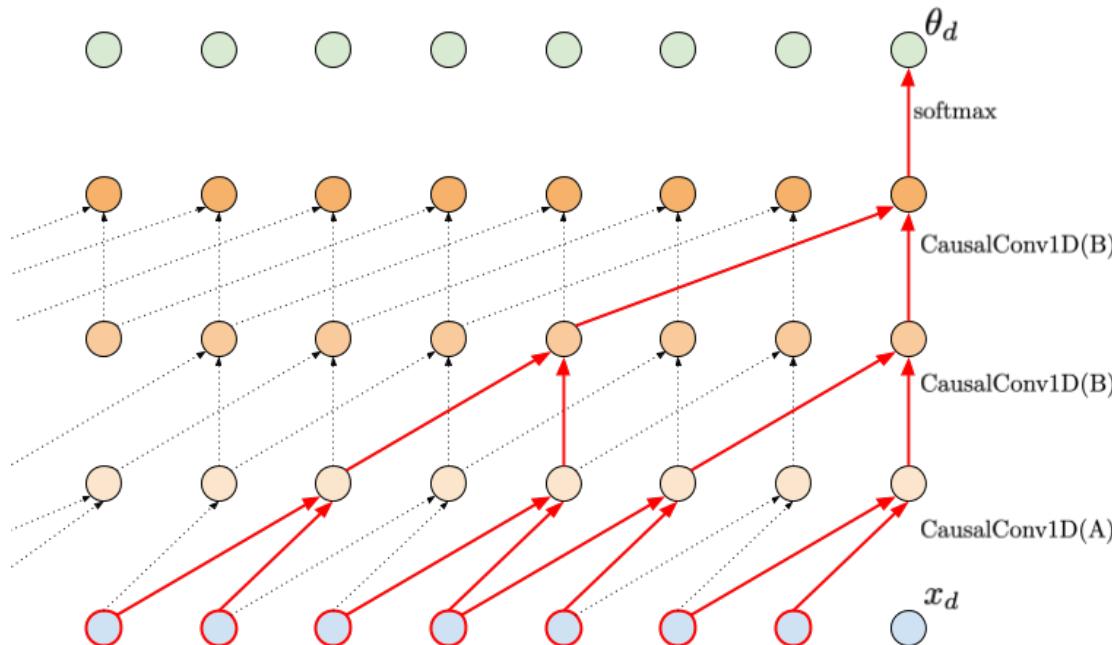
Approach 3: Long-range memory with CNNs



AUTOREGRESSIVE MODELS (ARMS)

Question: How we can model the conditionals $p(x_d | \mathbf{x}_{<d})$ **efficiently**?

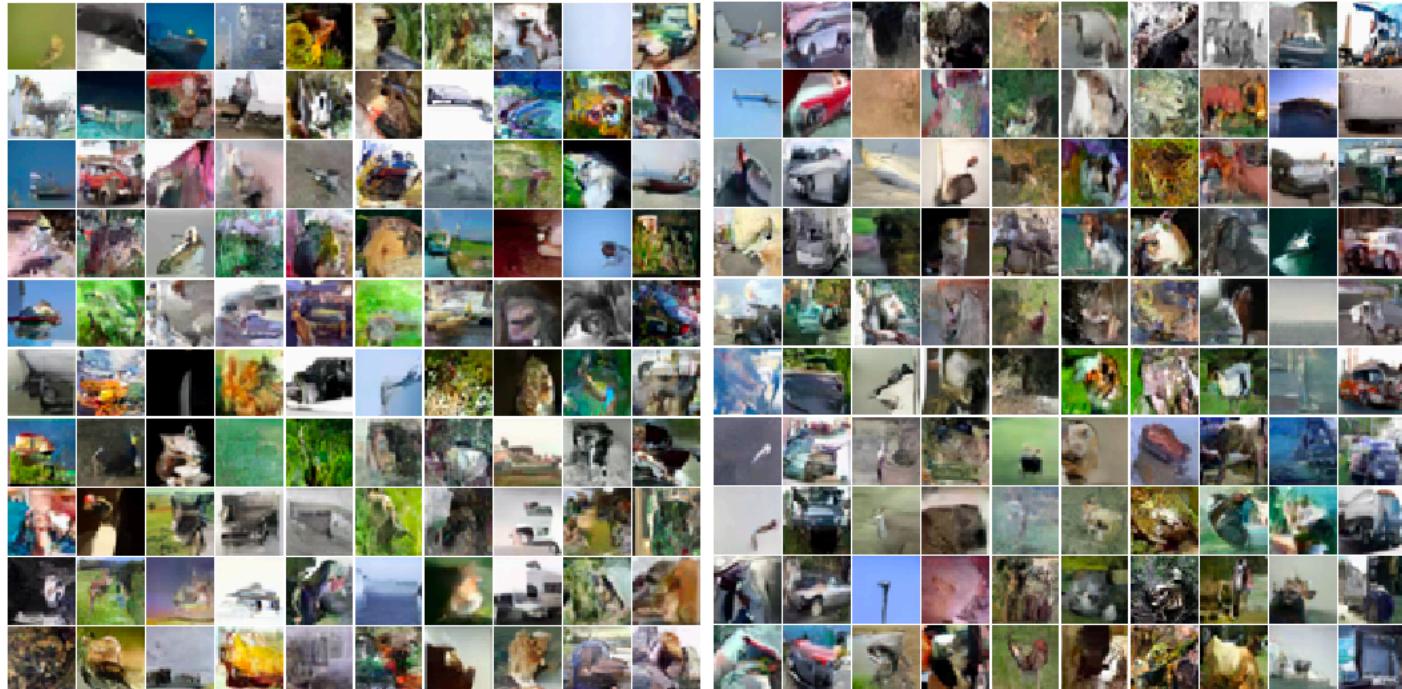
Approach 3: Long-range memory with CNNs



Easy!
Long-range dependencies!
No training issues!

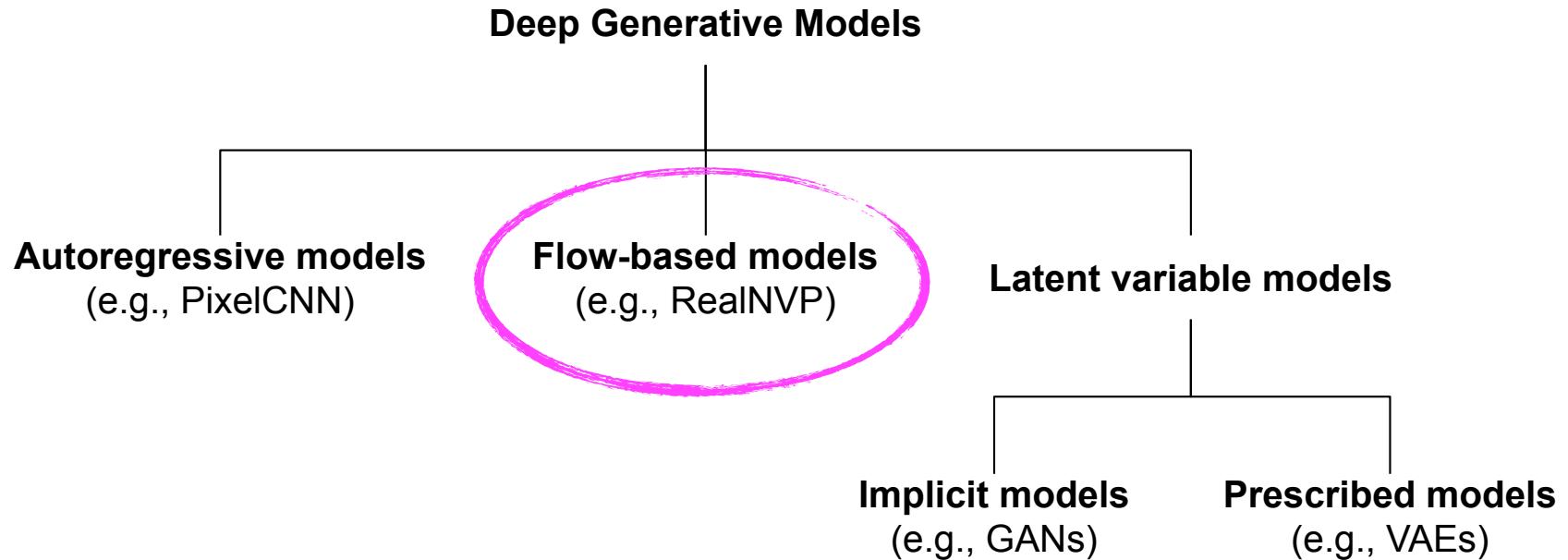
Slow generation

AUTOREGRESSIVE MODELS (ARMS)



Samples from a PixelCNN

DEEP GENERATIVE MODELING: HOW WE CAN FORMULATE IT?



FLows (Flow-based Models)

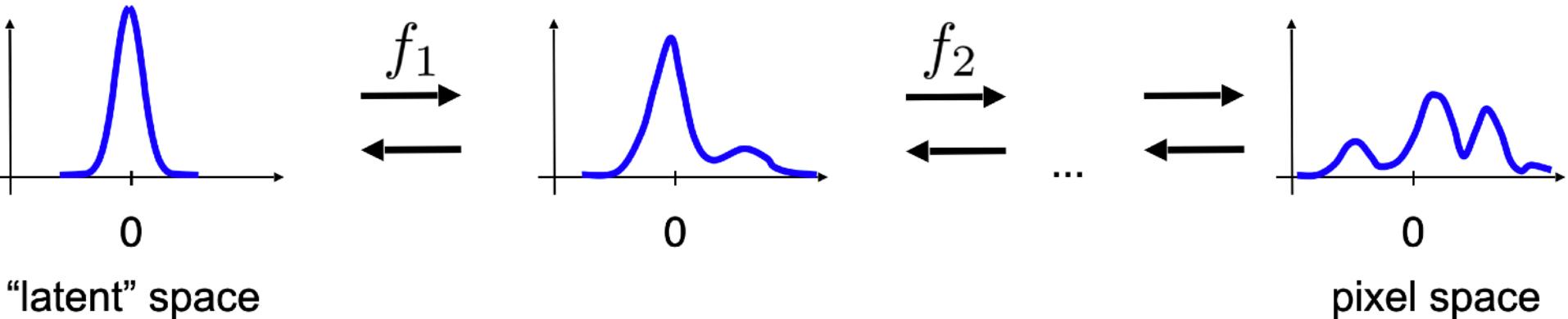
We change a random variable \mathbf{x} to another random variable \mathbf{z} using **invertible** transformations, $\mathbf{x}, \mathbf{z} \in \mathbb{R}^D$:

$$p(\mathbf{x}) = \pi(\mathbf{z}_0 = f^{-1}(\mathbf{x})) \prod_{i=1}^K \left| \mathbf{J}_{f_i}(z_{i-1}) \right|^{-1}$$

FLows (Flow-based Models)

We change a random variable \mathbf{x} to another random variable \mathbf{z} using **invertible** transformations, $\mathbf{x}, \mathbf{z} \in \mathbb{R}^D$:

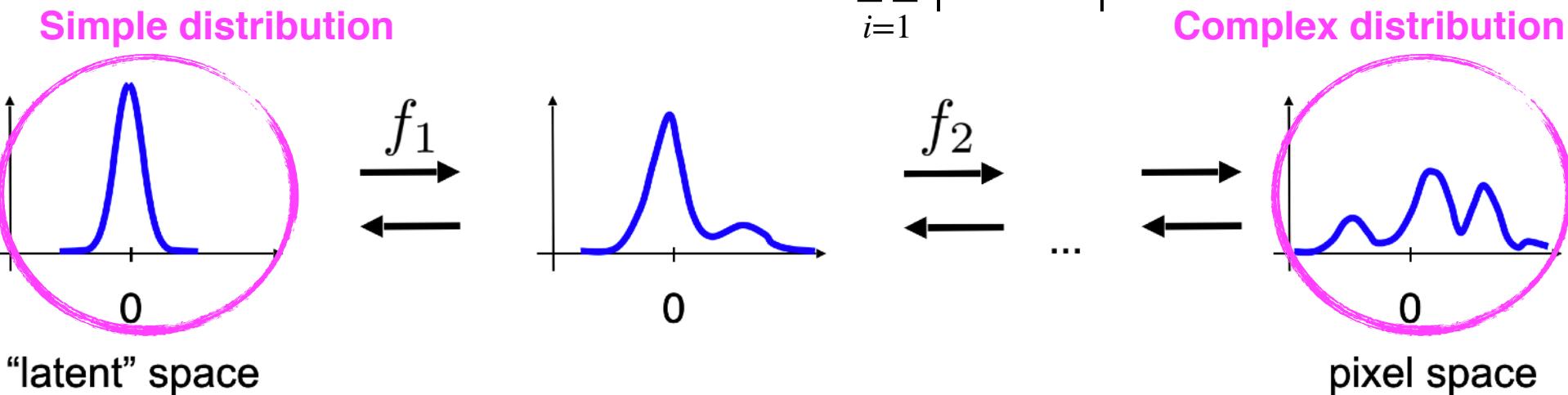
$$p(\mathbf{x}) = \pi(\mathbf{z}_0 = f^{-1}(\mathbf{x})) \prod_{i=1}^K \left| \mathbf{J}_{f_i}(z_{i-1}) \right|^{-1}$$



FLows (FLOW-BASED MODELS)

We change a random variable \mathbf{x} to another random variable \mathbf{z} using **invertible** transformations, $\mathbf{x}, \mathbf{z} \in \mathbb{R}^D$:

$$p(\mathbf{x}) = \pi(\mathbf{z}_0 = f^{-1}(\mathbf{x})) \prod_{i=1}^K \left| \mathbf{J}_{f_i}(z_{i-1}) \right|^{-1}$$

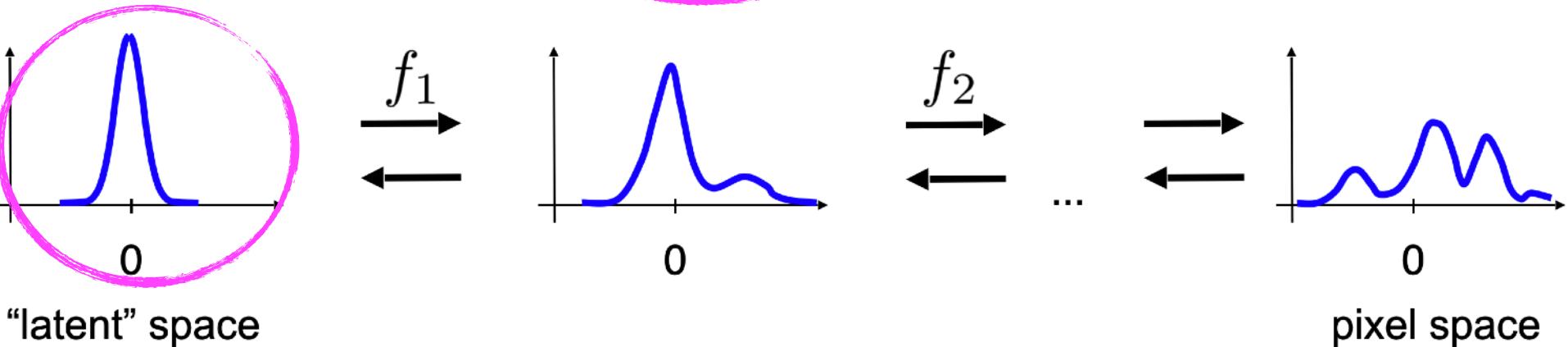


FLows (FLOW-BASED MODELS)

We change a random variable \mathbf{x} to another random variable \mathbf{z} using **invertible** transformations, $\mathbf{x}, \mathbf{z} \in \mathbb{R}^D$:

Known, e.g., Gaussian

$$p(\mathbf{x}) = \pi(\mathbf{z}_0 = f^{-1}(\mathbf{x})) \prod_{i=1}^K \left| \mathbf{J}_{f_i}(z_{i-1}) \right|^{-1}$$

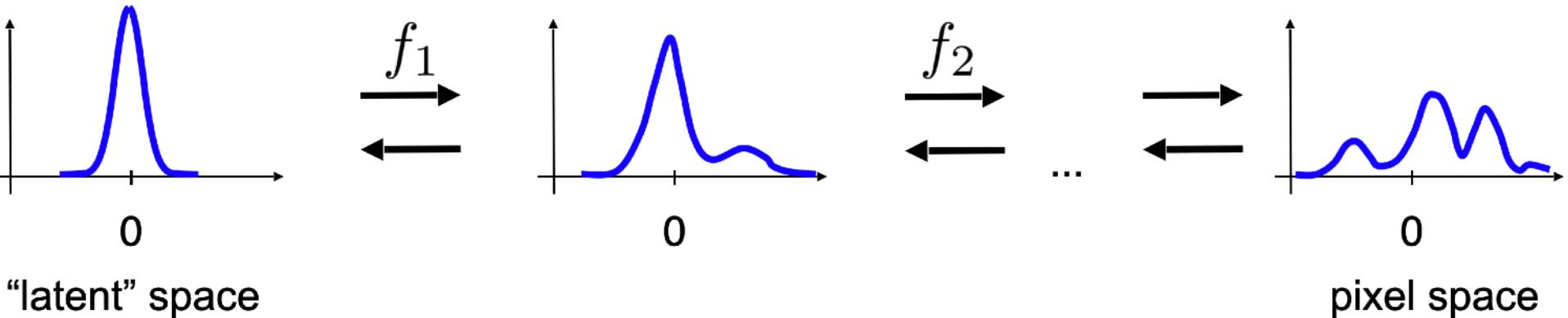


FLows (Flow-based Models)

We change a random variable \mathbf{x} to another random variable \mathbf{z} using **invertible** transformations, $\mathbf{x}, \mathbf{z} \in \mathbb{R}^D$:

Jacobian must be tractable

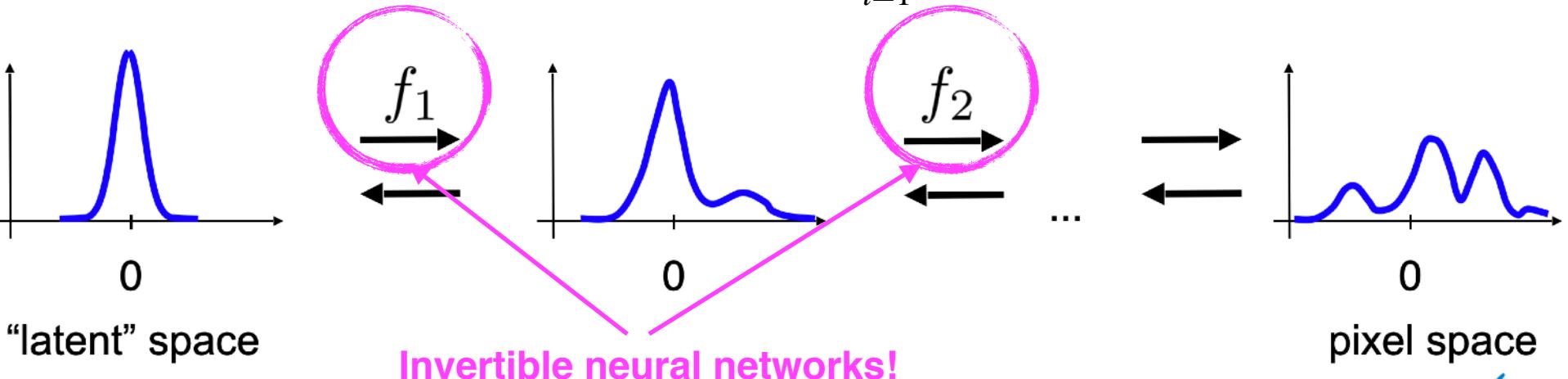
$$p(\mathbf{x}) = \pi(\mathbf{z}_0 = f^{-1}(\mathbf{x})) \prod_{i=1}^K \left| \mathbf{J}_{f_i}(z_{i-1}) \right|^{-1}$$



FLows (FLOW-BASED MODELS)

We change a random variable \mathbf{x} to another random variable \mathbf{z} using **invertible** transformations, $\mathbf{x}, \mathbf{z} \in \mathbb{R}^D$:

$$p(\mathbf{x}) = \pi(\mathbf{z}_0 = f^{-1}(\mathbf{x})) \prod_{i=1}^K \left| \mathbf{J}_{f_i}(z_{i-1}) \right|^{-1}$$



FLows (Flow-based Models)

We change a random variable \mathbf{x} to another random variable \mathbf{z} using **invertible** transformations, $\mathbf{x}, \mathbf{z} \in \mathbb{R}^D$:

$$p(\mathbf{x}) = \pi(\mathbf{z}_0 = f^{-1}(\mathbf{x})) \prod_{i=1}^K \left| \mathbf{J}_{f_i}(z_{i-1}) \right|^{-1}$$

Training objective:

$$\ln p(\mathbf{x}) = \ln \pi(\mathbf{z}_0 = f^{-1}(\mathbf{x})) - \sum_{i=1}^K \ln \left| \mathbf{J}_{f_i}(z_{i-1}) \right|$$

FLows (Flow-based Models): INVERTIBLE LAYERS

Two main components

1) Coupling layer:

$$\begin{aligned}\mathbf{y}_a &= \mathbf{x}_a \\ \mathbf{y}_b &= \exp\left(s\left(\mathbf{x}_a\right)\right) \odot \mathbf{x}_b + t\left(\mathbf{x}_a\right)\end{aligned}$$

is invertible by design:

$$\mathbf{x}_b = \left(\mathbf{y}_b - t(\mathbf{y}_a)\right) \odot \exp\left(-s(\mathbf{y}_a)\right)$$

$$\mathbf{x}_a = \mathbf{y}_a$$

2) Permutation layer

FLows (Flow-based Models): INVERTIBLE LAYERS

Two main components

1) Coupling layer:

$$\begin{aligned}\mathbf{y}_a &= \mathbf{x}_a \\ \mathbf{y}_b &= \exp\left(s(\mathbf{x}_a)\right) \odot \mathbf{x}_b + t(\mathbf{x}_a)\end{aligned}$$

is invertible by design:

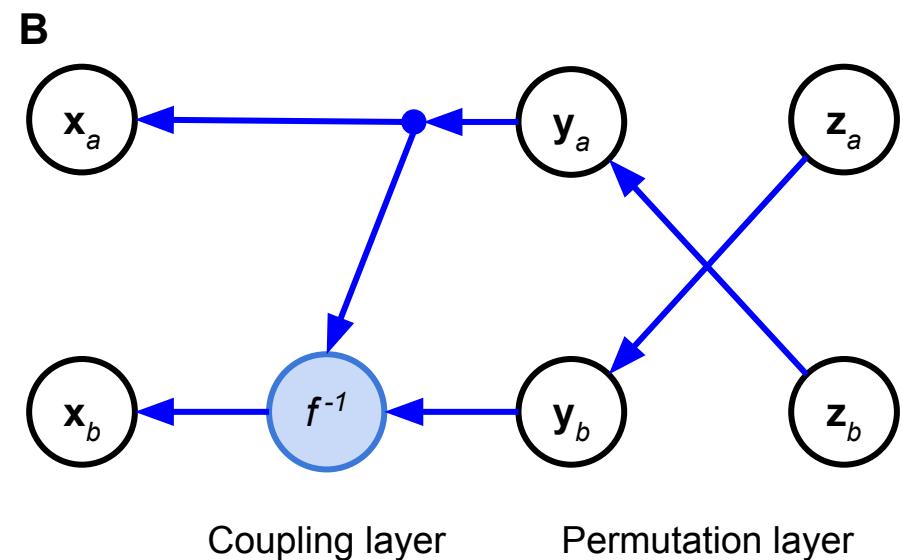
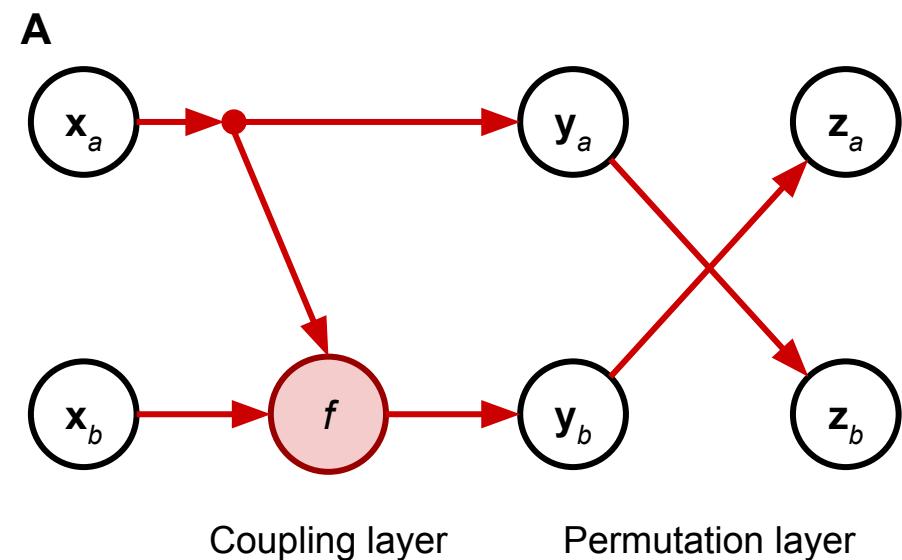
$$\begin{aligned}\mathbf{x}_b &= (\mathbf{y}_b - t(\mathbf{y}_a)) \odot \exp(-s(\mathbf{y}_a)) \\ \mathbf{x}_a &= \mathbf{y}_a\end{aligned}$$

Jacobian is tractable!

$$\det(\mathbf{J}) = \prod_{j=1}^{D-d} \exp\left(s(\mathbf{x}_a)\right)_j = \exp\left(\sum_{j=1}^{D-d} s(\mathbf{x}_a)_j\right)$$

2) Permutation layer $\det(\mathbf{J}) = 1$

FLows (Flow-based Models): INVERTIBLE LAYERS



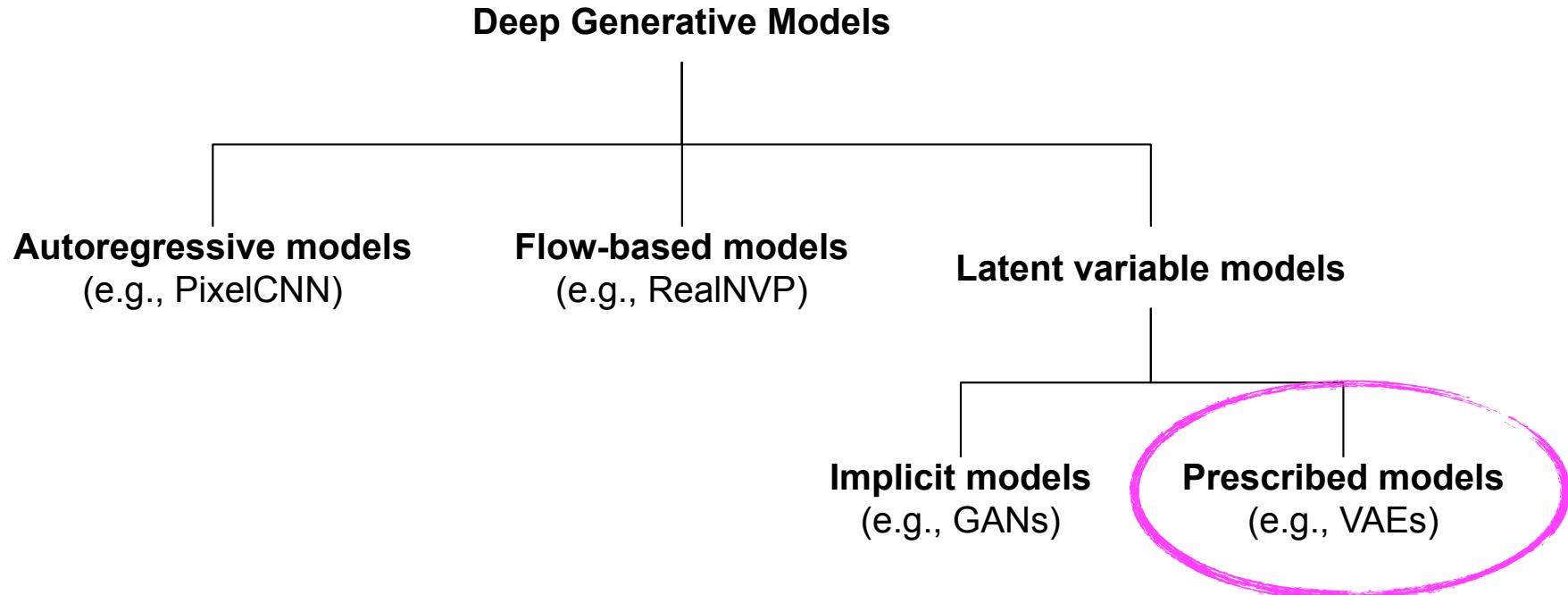
A: Forward pass. **B:** Inverse pass.

FLows (FLOW-BASED MODELS)



⁷⁰Kingma, D.P., and Prafulla D. "Glow: Generative flow with invertible 1x1 convolutions." *NeurIPSx 2018*

DEEP GENERATIVE MODELING: HOW WE CAN FORMULATE IT?

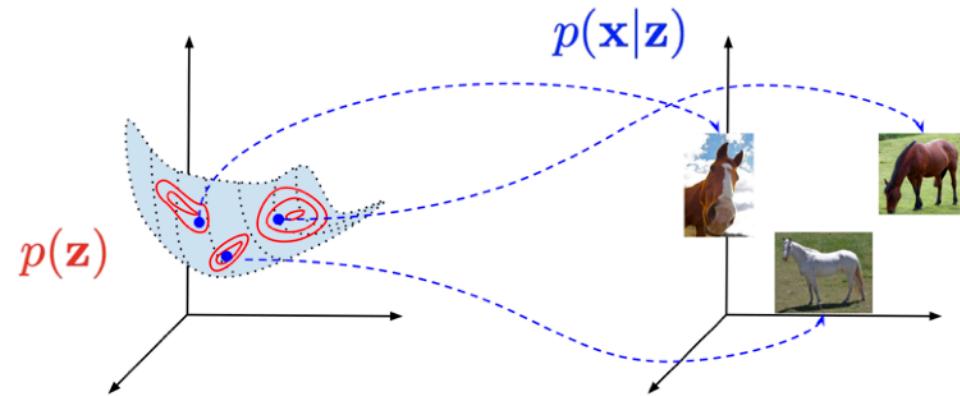


VARIATIONAL AUTO-ENCODERS

Let's consider a **latent variable model** where we distinguish:

- **latent variables** $\mathbf{z} \in \mathcal{Z}^M$
- **observable variables** $\mathbf{x} \in \mathcal{X}^D$

Latent variables lie on a **low-dimensional manifold**.



Generative process:

1. $\mathbf{z} \sim p(\mathbf{z})$
2. $\mathbf{x} \sim p(\mathbf{x} | \mathbf{z})$

VARIATIONAL AUTO-ENCODERS

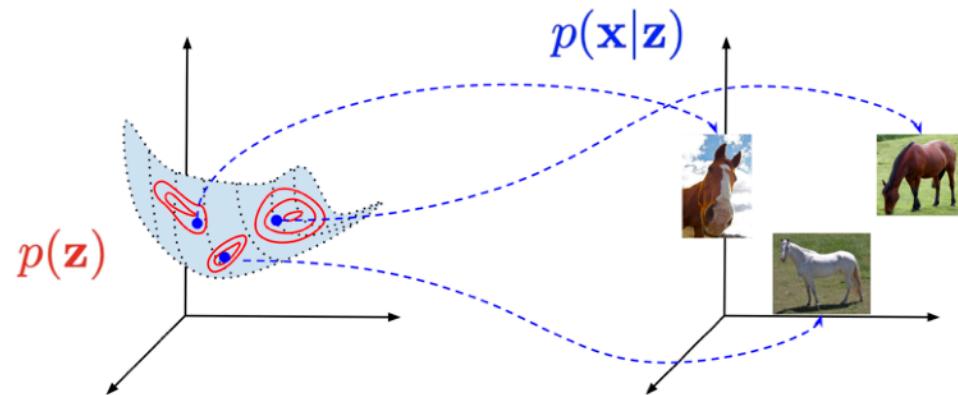
Let's consider a **latent variable model** where we distinguish:

- **latent variables** $\mathbf{z} \in \mathcal{Z}^M$
- **observable variables** $\mathbf{x} \in \mathcal{X}^D$

Latent variables lie on a **low-dimensional manifold**.

The objective function:

$$\ln p(\mathbf{x}) = \ln \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z}$$



Generative process:

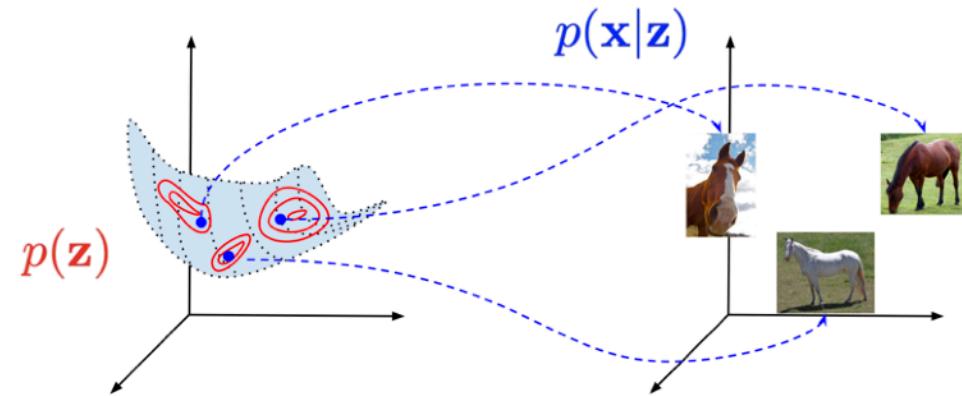
1. $\mathbf{z} \sim p(\mathbf{z})$
2. $\mathbf{x} \sim p(\mathbf{x} | \mathbf{z})$

VARIATIONAL AUTO-ENCODERS

Let's consider a **latent variable model** where we distinguish:

- **latent variables** $\mathbf{z} \in \mathcal{Z}^M$
- **observable variables** $\mathbf{x} \in \mathcal{X}^D$

Latent variables lie on a **low-dimensional manifold**.



Generative process:

1. $\mathbf{z} \sim p(\mathbf{z})$
2. $\mathbf{x} \sim p(\mathbf{x} | \mathbf{z})$

The objective function:

$$\ln p(\mathbf{x}) = \ln \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z}$$

The integral is intractable...

VARIATIONAL AUTO-ENCODERS

$$\begin{aligned}\ln p(\mathbf{x}) &= \ln \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z})]\end{aligned}$$

VARIATIONAL AUTO-ENCODERS

$$\begin{aligned}\ln p(\mathbf{x}) &= \ln \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z})]\end{aligned}$$

Variational posteriors

VARIATIONAL AUTO-ENCODERS

$$\begin{aligned}\ln p(\mathbf{x}) &= \ln \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z})]\end{aligned}$$

VARIATIONAL AUTO-ENCODERS

$$\begin{aligned}\ln p(\mathbf{x}) &= \ln \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \quad \text{Jensen's inequality} \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z})]\end{aligned}$$

VARIATIONAL AUTO-ENCODERS

$$\begin{aligned}\ln p(\mathbf{x}) &= \ln \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \quad \text{Jensen's inequality} \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z})]\end{aligned}$$

VARIATIONAL AUTO-ENCODERS

$$\begin{aligned}\ln p(\mathbf{x}) &= \ln \int p(\mathbf{x} \mid \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} \mid \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} \mid \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} \mid \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} \mid \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln p(\mathbf{x} \mid \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z})]\end{aligned}$$

VARIATIONAL AUTO-ENCODERS

$$\begin{aligned}\ln p(\mathbf{x}) &= \ln \int p(\mathbf{x} \mid \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} \mid \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} \mid \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} \mid \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} \mid \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln p(\mathbf{x} \mid \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z})]\end{aligned}$$

VARIATIONAL AUTO-ENCODERS

$$\begin{aligned}\ln p(\mathbf{x}) &= \ln \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z} \\ &= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z} \\ &= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\ &\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right]\end{aligned}$$

Reconstruction error “Regularization” term

$$= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z})]$$

VARIATIONAL AUTO-ENCODERS

$$\ln p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\ln p(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\ln q_\phi(\mathbf{z}|\mathbf{x}) - \ln p(\mathbf{z})]$$

ELBO: Evidence Lower Bound

VARIATIONAL AUTO-ENCODERS

$$\ln p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x})} [\ln q_\phi(\mathbf{z} | \mathbf{x}) - \ln p(\mathbf{z})]$$

We consider **amortized inference**: $q_\phi(\mathbf{z} | \mathbf{x})$

In other words, a single parameterization for each new input \mathbf{x} .

VARIATIONAL AUTO-ENCODERS

$$\ln p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x})} [\ln q_\phi(\mathbf{z} | \mathbf{x}) - \ln p(\mathbf{z})]$$

We consider **amortized inference**: $q_\phi(\mathbf{z} | \mathbf{x})$

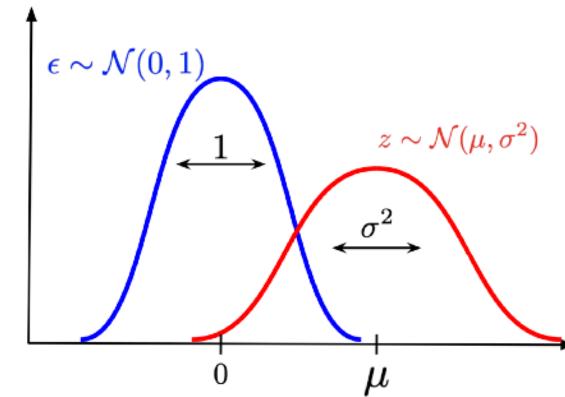
In other words, a single parameterization for each new input \mathbf{x} .

Moreover, we use **reparameterization trick**:

Every Gaussian variable could be defined as:

$$z = \mu + \sigma \cdot \varepsilon$$

where $\varepsilon \sim \mathcal{N}(0, 1)$



VARIATIONAL AUTO-ENCODERS

$$\ln p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x})} [\ln q_\phi(\mathbf{z} | \mathbf{x}) - \ln p(\mathbf{z})]$$

We consider **amortized inference**: $q_\phi(\mathbf{z} | \mathbf{x})$

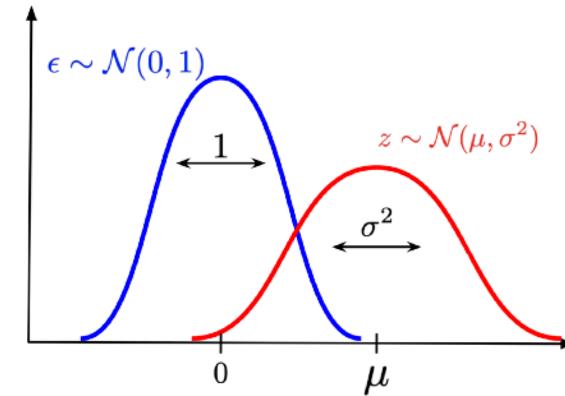
In other words, a single parameterization for each new input \mathbf{x} .

Moreover, we use **reparameterization trick**:

It reduces the variance of the gradients.

It allows to get randomness outside \mathbf{z} .

$$z = \mu + \sigma \cdot \epsilon$$

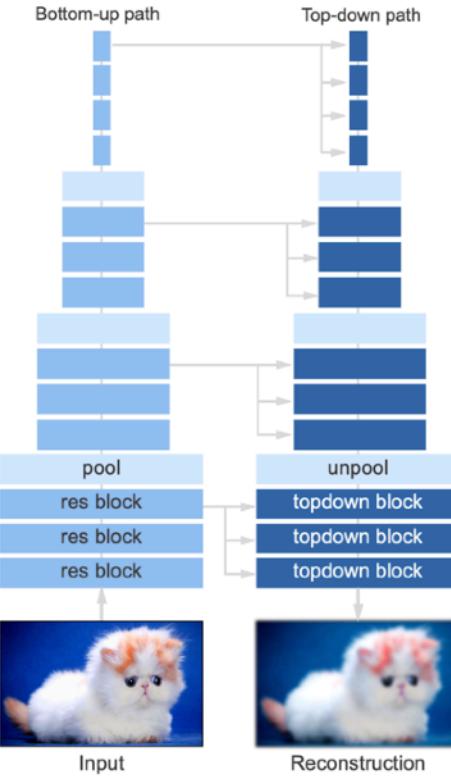


⁸⁶Kingma, D.P., and Welling, M.. "Auto-encoding variational bayes." *ICLR 2014*

VARIATIONAL AUTO-ENCODERS



Generations



Very Deep VAE

⁸⁷ Child, R. "Very Deep VAEs Generalize Autoregressive Models and Can Outperform Them on Images." *ICLR 2021*

VARIATIONAL AUTO-ENCODERS

i) selfVAE - downscale - 31v1



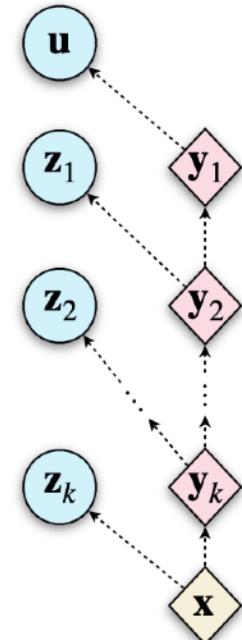
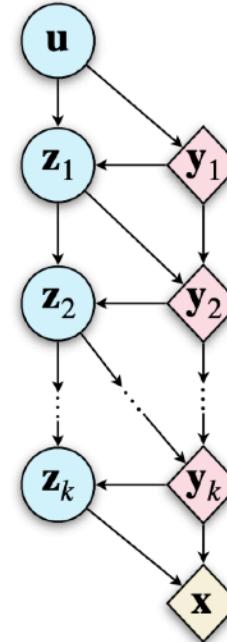
Generations

ii) selfVAE - sketch



i) Generative Model

Hierarchical VAE



CONCLUSION

- Here: **the likelihood-based generative models.**



CONCLUSION

- Here: **the likelihood-based generative models.**
- We **skipped** Generative Adversarial Nets & others.



CONCLUSION

- Here: **the likelihood-based generative models**.
- We **skipped** Generative Adversarial Nets & others.
- Why generative modeling?

$$p(\mathbf{x}, y) = p(y | \mathbf{x}) p(\mathbf{x})$$



CONCLUSION

- Here: **the likelihood-based generative models.**
- We **skipped** Generative Adversarial Nets & others.
- Why generative modeling?

$$p(\mathbf{x}, y) = p(y | \mathbf{x}) p(\mathbf{x})$$

- Important directions:
 - Better uncertainty quantification
 - New parameterization (new neural networks)
 - Out-of-Distribution
 - Continual learning



BLOG ABOUT DEEP GENERATIVE MODELING

If you are interested in going deeper into deep generative modeling, please take a look at my blog: [\[Blog\]](#)

- **Intro:** [\[Link\]](#)
- **ARMs:** [\[Link\]](#)
- **Flows:** [\[Link\]](#), [\[Link\]](#)
- **VAEs:** [\[Link\]](#), [\[Link\]](#)
- **Hybrid modeling:** [\[Link\]](#)

THANK YOU FOR YOUR ATTENTION

Jakub M. Tomczak
Computational Intelligence group
Vrije Universiteit Amsterdam

Webpage: <https://jmtomczak.github.io/>

Github: <https://github.com/jmtomczak>

Twitter: <https://twitter.com/jmtomczak>