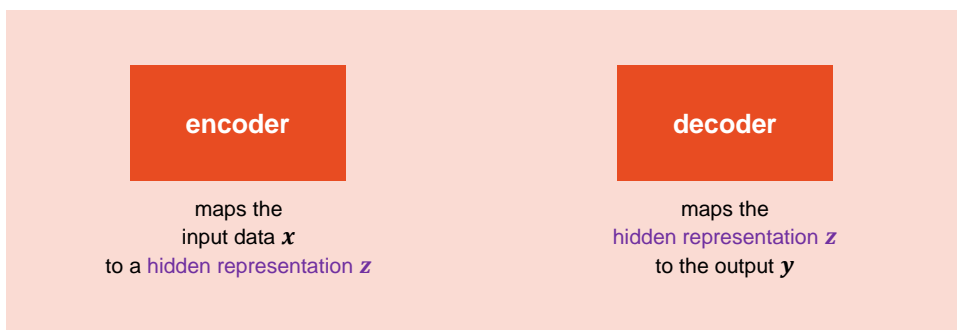# EE-559
# Deep Learning

## What's on today?

- **Deterministic autoencoders**: on discovering structure within the data
- **Variational autoencoders**: on using a continuous latent space
- **Vector quantised-variational autoencoder**: on using discrete latents
- **Diffusion probabilistic models**: on learning noise removal
- **Instruction-following diffusion models**: on editing with language
- **Exercises**: autoencoders and variational autoencoders

# Deterministic autoencoders
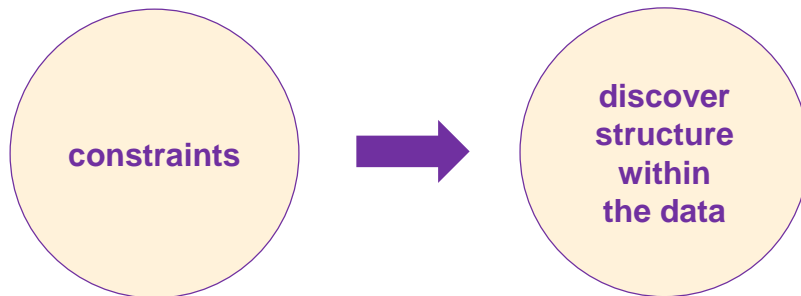
## Auto-associative neural network

| encoder | decoder |
|---|---|
| maps the input data $x$ to a hidden representation $z$ | maps the hidden representation $z$ to the output $y$ |

*trained* to generate an output $y$ that is as close as possible to $x$

**Concepts**:
Same number of input and output units; internal representation $z(x)$ of each *new* input.
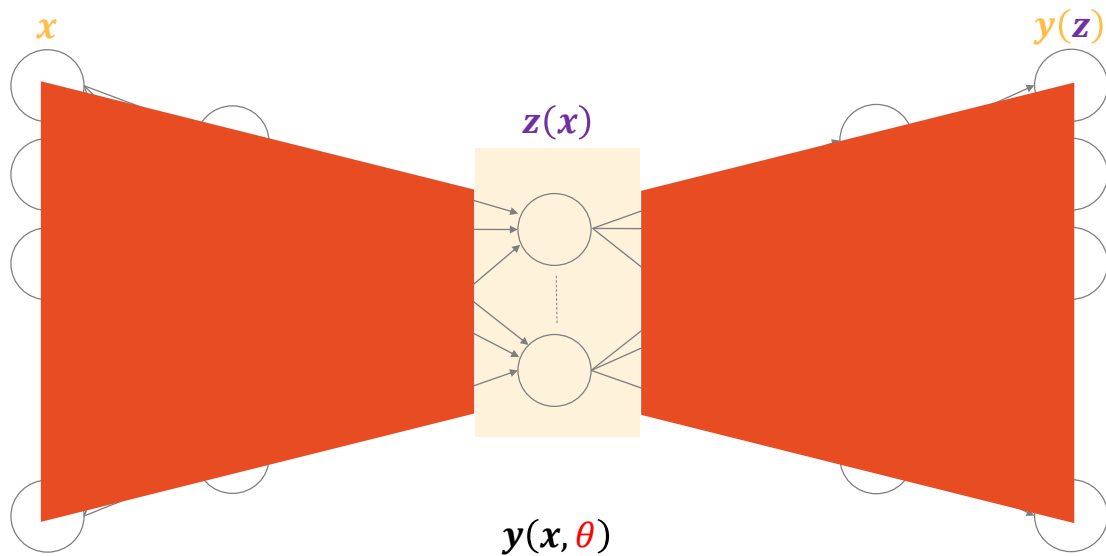
# Autoencoder



**constraints** → **discover structure within the data**

**Concept**:
Non-linear form of Principal Component Analysis (PCA)

# Autoencoder: network diagram



restrict the dimensionality of $z$
sparse representation for $z$

# Autoencoder

$x$             $z(x)$                $y(z)$

$y(x, \theta)$

# Autoencoder: training

predict missing values

remove additive noise

**encoder**        **decoder**

$x$      $z$      $y$

(internal) **representation** to be used in a downstream task (after training)

**Concepts**: Self-supervision; modify the training process to *undo* corruptions of the input vector $x$; dimensionality of subspace to be defined *before* training the network.

# Denoising autoencoders

$$x_n \xrightarrow{\text{noise}} \widetilde{x}_n$$

examples

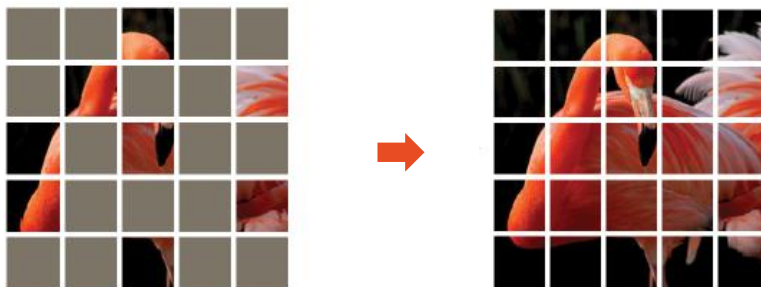set a random amount $\rho$ of input variables to $0$     $0 \le \rho \le 1$

add *independent* zero-mean Gaussian noise to each input variable

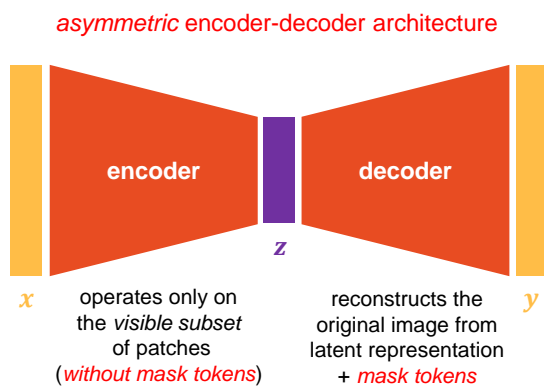$$E(\theta) = \sum_{n=1}^{N} \left\| y(\widetilde{x}_n, \theta) - x_n \right\|^2$$

**Concept**:
Network encouraged to learn (aspects of) the structure of the data by learning to denoise the input data.
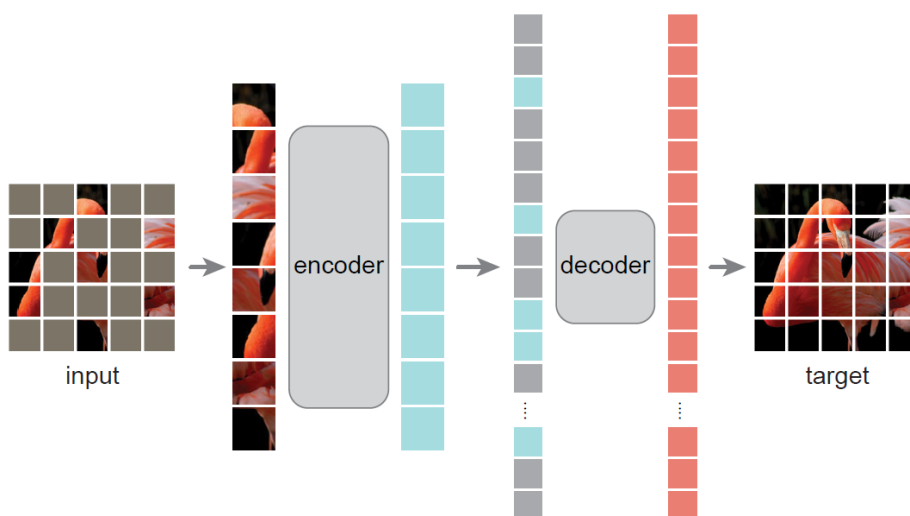
# Masked autoencoders



arXiv:2111.06377

# Masked autoencoders

*asymmetric* encoder-decoder architecture



$x$ — operates only on the *visible subset* of patches (*without mask tokens*)

$z$

reconstructs the original image from latent representation + *mask tokens* — $y$

**Concepts**: Pass only randomly selected input patches (as low as 25% of an image); each mask token is augmented with positional encoding.
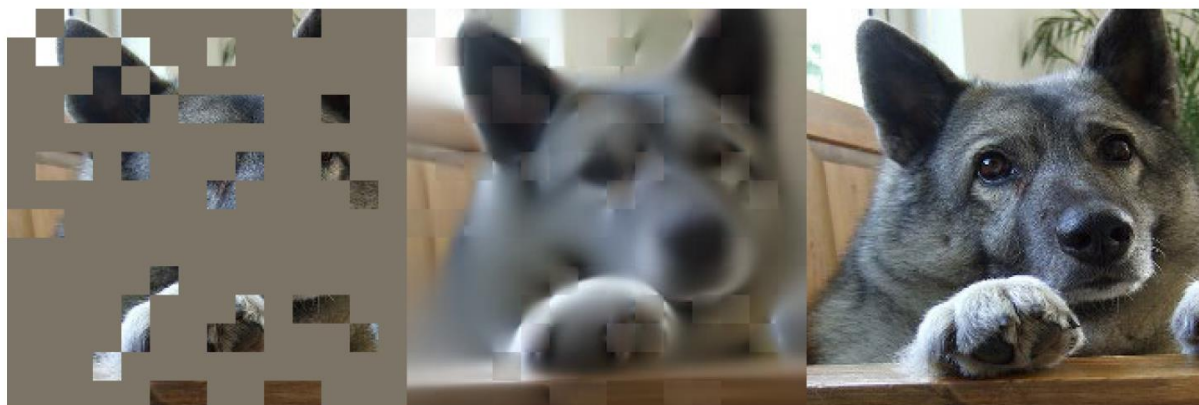
# Masked autoencoders



input → encoder → decoder → target

arXiv:2111.06377

# Masked autoencoder: example

# Masked autoencoder: example

# Sparse autoencoders

regularizer to encourage sparse representation

$$\tilde{E}(\theta) = E(\theta) + \lambda \sum_{k=1}^{K} |z_k|$$

un-regularized error

regularizer applied to the unit activations of a hidden layer

**Concept**:
To constrain the internal representation with a regularizer.

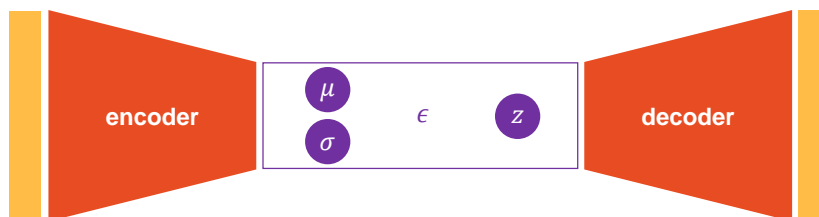# Variational autoencoders

# VAE

continuous, structured latent space



encoder  $\mu$  $\sigma$  →  $z$  decoder

mapping input data into the *parameters* of a
probability distribution (e.g. *mean* & *variance* of a Gaussian)

arXiv:1312.6114

# VAE: flow of gradients during training



encoder  $\mu$  $\sigma$   $\epsilon$   $z$  decoder

$\epsilon$ Gaussian **random variable** with *zero mean* and *unit variance*

$z = \sigma\epsilon + \mu$    reparametrization trick (**replaces** direct sample of $z$)

$z$ Gaussian distribution with *mean* $\mu$ and *variance* $\sigma^2$

arXiv:1312.6114

# **Vector Quantised-Variational Autoencoder**

## Neural discrete representation learning

discrete latent representation

prior is learnt
(rather than static)

**Concepts**: Parameterization of the posterior distribution of (discrete) latent variables given an observation; vector quantisation is used to learn a discrete latent.
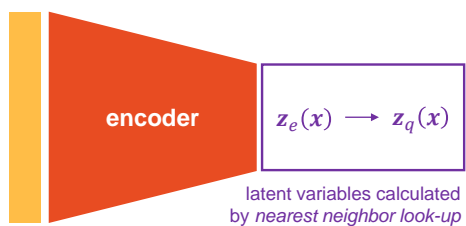
# VQ-VAE



**Concepts**: Categorical posterior and prior distributions;
samples (drawn from these distributions) index an embedding.

# Discrete latent variables

$e \in \mathbb{R}^{K \times D}$ latent embedding space

$K$ size of the discrete latent space

$D$ dimensionality of each latent embedding vector



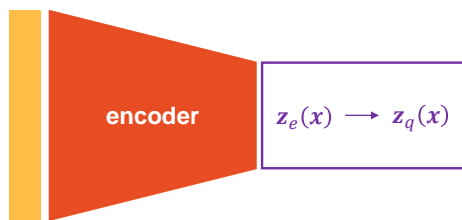posterior categorical distribution (one-hot)

output of the encoder network

$$q(z = k \mid x) = \begin{cases} 1 & \text{for } k = \text{argmin}_j ||z_e(x) - e_j||_2 \\ 0 & \text{otherwise} \end{cases}$$

embedding vector
$e_i \in \mathbb{R}^D \; i = 1, 2, \dots, K$

$$z_q(x) = e_k \qquad \text{where } k = \text{argmin}_j ||z_e(x) - e_j||_2$$

arXiv:1711.00937

# Learning

$e \in \mathbb{R}^{K \times D}$ latent embedding space
$K$ size of the discrete latent space
$D$ dimensionality of each latent embedding vector



$$z_q(x) = e_k \qquad \text{where } k = \text{argmin}_j ||z_e(x) - e_j||_2$$
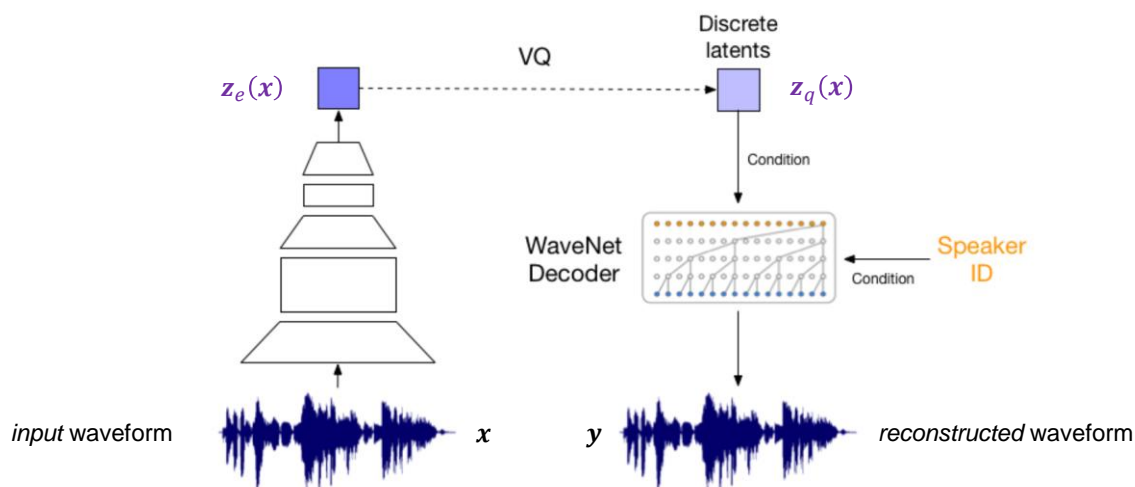
no gradient defined for this mapping!

copy gradients from decoder input $z_q(x)$
to encoder output $z_e(x)$

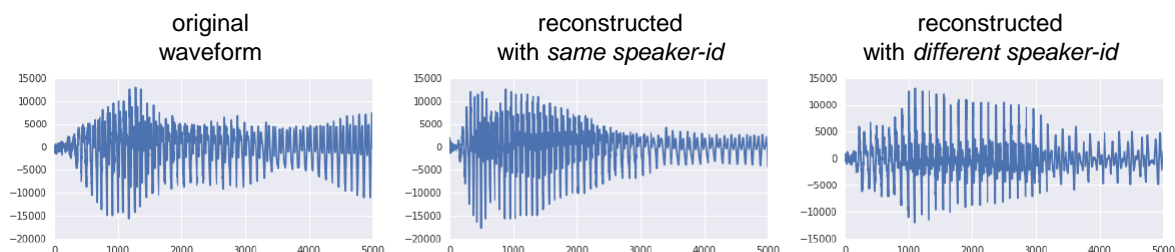arXiv:1711.00937

# VQ-VAE: example



arXiv:1711.00937

## VQ-VAE: example



the **contents** of the three waveforms are the **same**

# Diffusion probabilistic models

# Denoising diffusion probabilistic models



**encoder**

$z$

**decoder**

$x$

$y$

multi-step *noise process*
that transforms input
into a sample from a
Gaussian distribution

generate new data
from samples of a
Gaussian distribution

**Concepts**: Hierarchical variational autoencoder; encoder distribution is fixed (*defined by the noise process*);
only the generative distribution of the decoder is learned.

# Diffusion process



**encoder**

$z$

$x$

*Markov chain* that
gradually *adds noise* to the data
until the signal/information is destroyed

multi-step *noise process*
that transforms input
into a sample from a
Gaussian distribution

# Diffusion (forward) process

$$q(\boldsymbol{x}_{1:T}|\boldsymbol{x}_0) = \prod_{t=1}^{T} q(\boldsymbol{x}_t|\boldsymbol{x}_{t-1})$$

posterior

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

fixed Markov chain that
gradually adds Gaussian noise

$$\beta_1, \dots, \beta_T$$     variance schedule (the $\beta_i$ are fixed)
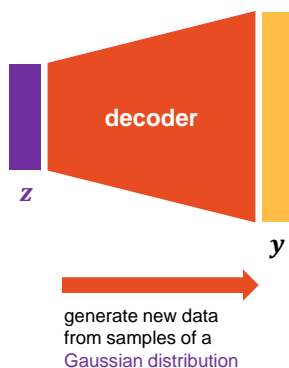
arXiv:2006.11239

# Reverse process

if diffusion consists of small amounts of Gaussian noise,
then it suffices to set the *sampling chain transitions*
to conditional Gaussians too



**decoder**

*z*

*y*

simple neural network parameterization!

generate new data
from samples of a
Gaussian distribution

arXiv:2006.11239

# Reverse process

$$p_\theta(\boldsymbol{x}_{0:T}) = p(\boldsymbol{x}_T) \prod_{t=1}^{T} p_\theta(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t)$$   **reverse process** (joint distribution)

$$p(\boldsymbol{x}_T) = \mathcal{N}(\boldsymbol{x}_T; \boldsymbol{0}, \mathbf{I})$$   start of the Markov chain

$$p_\theta(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t) = \mathcal{N}(\boldsymbol{x}_{t-1}; \mu_\theta(\boldsymbol{x}_t, t), \textstyle\sum_\theta(\boldsymbol{x}_t, t))$$   learned Gaussian transitions

$$\boldsymbol{x}_1, \boldsymbol{x}_2, \dots, \boldsymbol{x}_T$$   latents of the same dimension as the data $\boldsymbol{x}_0$

arXiv:2006.11239

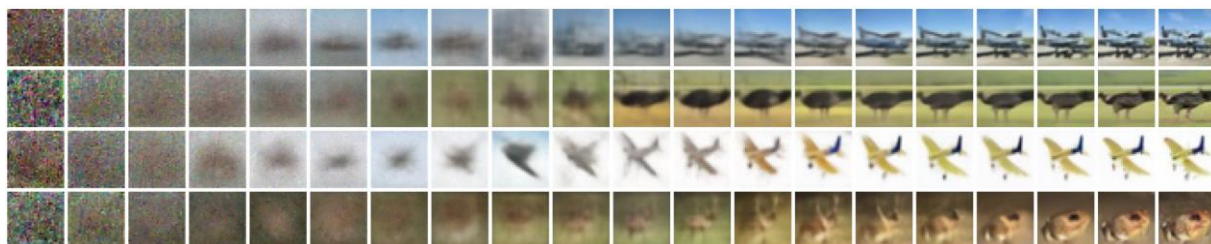# Diffusion probabilistic models

**Reverse process**:
parameterized Markov chain
trained using variational inference



arXiv:2006.11239

Reverse process: examples



arXiv:2006.11239

# Instruction-following diffusion models

# Text-to-image diffusion



*"Teddy bears swimming at the Olympics 400m Butterfly event"*

encoding text for image synthesis

image-text alignment

capturing complexity & compositionality of arbitrary natural language text inputs

arXiv:2205.11487

# Text-to-image diffusion



arXiv:2205.11487

# Anime quick response codes



https://arstechnica.com/information-technology/2023/06/redditor-creates-working-anime-qr-codes-using-stable-diffusion

# Anime quick response codes



https://arstechnica.com/information-technology/2023/06/redditor-creates-working-anime-qr-codes-using-stable-diffusion

# Editing images from human instructions
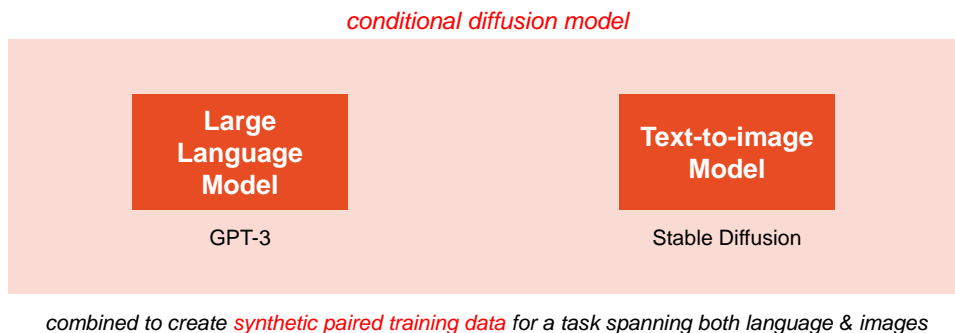
*conditional diffusion model*

| | |
|---|---|
| **Large Language Model** | **Text-to-image Model** |
| GPT-3 | Stable Diffusion |

*combined to create* *synthetic paired training data* *for a task spanning both language & images*

arXiv:2211.09800

# Synthetic multi-modal training data



(a) Generate text edits:

Input Caption: *"photograph of a girl riding a horse"* → GPT-3 → Instruction: *"have her ride a dragon"*
Edited Caption: *"photograph of a girl riding a dragon"*

(b) Generate paired images:

Input Caption: *"photograph of a girl riding a horse"*
Edited Caption: *"photograph of a girl riding a dragon"* → Stable Diffusion + Prompt2Prompt →

(c) Generated training examples:

*"convert to brick"*   *"Color the cars pink"*   *"Make it lit by fireworks"*   *"have her ride a dragon"*

arXiv:2211.09800

# Instruction-following diffusion model



*"turn her into a snake lady"*



arXiv:2211.09800

# Cross-attention control for editing

cross-attention maps **bind pixels** & **tokens** extracted from the prompt text

⬇

inject the cross-attention maps during the *diffusion process*
**controlling** which **pixels** attend to which **tokens** of the prompt text
during which *diffusion steps*

⬇         ⬇

change a single token's value in the prompt
(e.g. *cat* to *dog*)
fix the cross-attention maps
to preserve the scene composition

globally edit an image
(e.g. change the style)
add **new words** to the prompt
freeze the attention on previous tokens
allow **new attention** to flow to the **new tokens**

arXiv:2208.01626

# Cross-attention control for editing: example

*"Photo of a cat riding on a bicycle."*



| source image | *cat* > dog | *cat* > chicken | *cat* > squirrel | *cat* > elephant |

no need for model training, fine-tuning, extra data, or optimization

arXiv:2208.01626

# Practice exercises

# Today's practice exercises

• Autoencoders for anomaly detection

• Variational autoencoders for image generation

# What did we learn today?

• Deterministic autoencoders
• Variational autoencoders
• Vector quantised-variational autoencoder
• Diffusion probabilistic models
• Instruction-following diffusion models
• Exercises

# **EE-559**
# **Deep Learning**

andrea.cavallaro@epfl.ch