

# Uczenie Maszynowe 2

## Modele generatywne

dr hab. Piotr Duda, prof. AGH i PCz



# Plan wykładu

## **Wykład 1**

- **Organiczone Maszyny Boltzmann**

## **Wykład 2**

- **Autoenkodery, VAE, WAE**

## **Wykład 3**

- **GAN, Normalizing Flows**

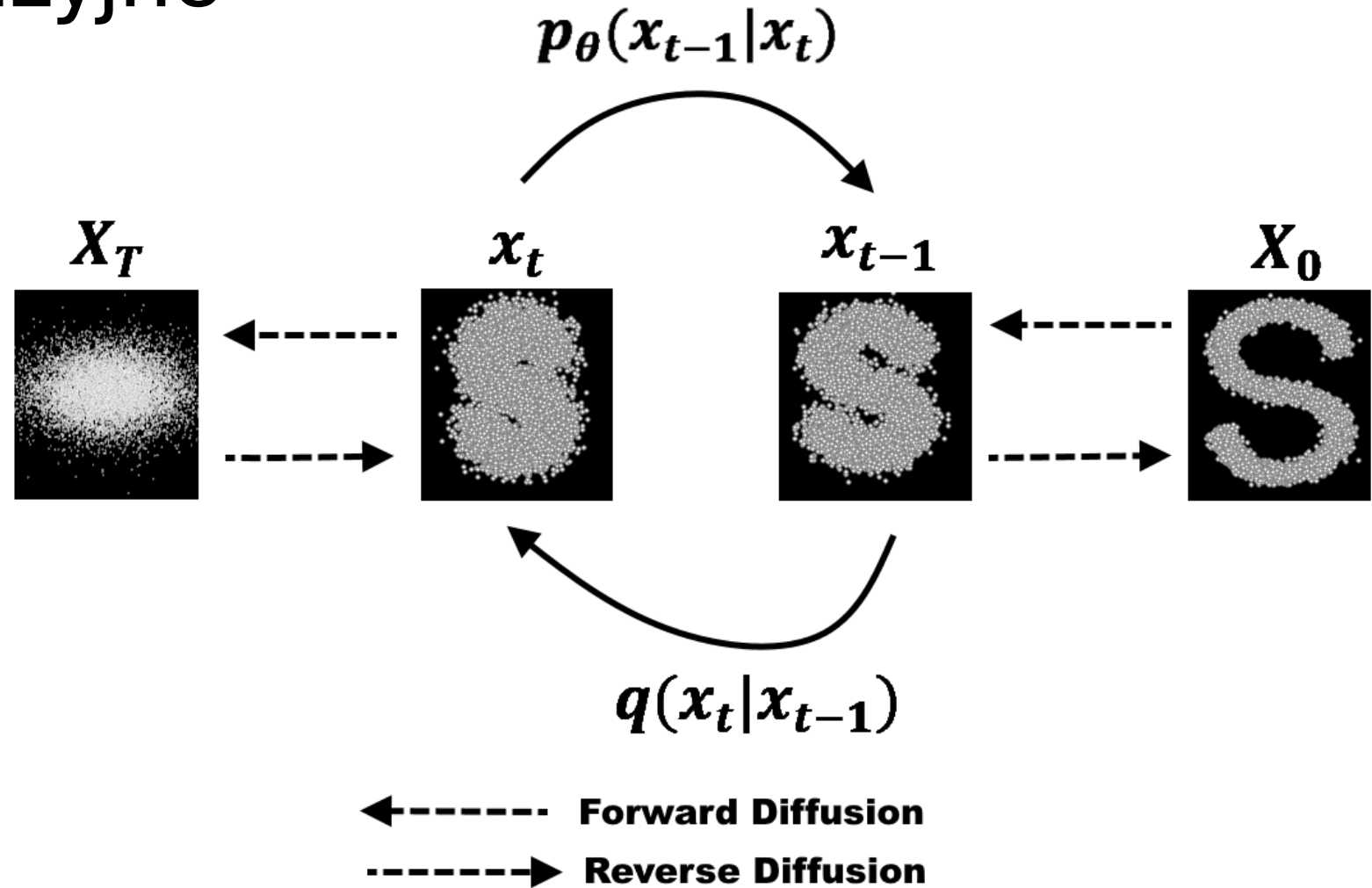
## **Wykład 4**

- **Podejście dyfuzyjne, Przegląd „aktualnych” modeli**

# Wykład 4 – Modele generatywne

- Podejście dyfuzyjne
- Transformery
- Podsumowanie metod
- Przegląd „aktualnych” rozwiązań
- Kolejny krok...

# Modele dyfuzyjne



# Koder

Pierwsze przekształcenie:

$$z_1 = \sqrt{1 - \beta_1}x + \sqrt{\beta_1}\epsilon_1$$

Kolejne przekształcenia:

$$z_t = \sqrt{1 - \beta_t}z_{t-1} + \sqrt{\beta_t}\epsilon_t,$$

dla  $t=2,3,\dots,T$ ,  $\beta_t \in [0,1]$ ,  $\epsilon_t \sim N(0,1)$ .

# Koder

$$P(z_1|x), P(z_2|z_1), \dots, P(z_T|z_{T-1})$$

$$P(z_1, \dots, z_T|x) = P(z_1|x) \prod_{i=2}^T P(z_i|z_{i-1})$$

# Diffusion kernel

$$z_t = \sqrt{\alpha_t}x + \sqrt{1 - \alpha_t}\epsilon,$$

gdzie

$$\alpha_t = \prod_{i=1}^t 1 - \beta_i,$$

a  $\epsilon \sim N(0,1)$

# Prawdopodobieństwo...

...wygenerowania bardziej zaszumionego obrazu

$$P(z_1|x) \Rightarrow N(\sqrt{1-\beta_1}x, \beta_1 I)$$

$$P(z_t|z_{t-1}) \Rightarrow N(\sqrt{1-\beta_t}z_{t-1}, \beta_t I)$$

... wygenerowania mniej zaszumionego obrazu

$$P(z_{t-1}|z_t) = \frac{P(z_t|z_{t-1})P(z_{t-1})}{P(z_t)}$$



# Dekoder

$$P(z_T) \Rightarrow N(0, I)$$

$$P(z_{t-1}|z_t, \varphi_t) \Rightarrow N(f_t(z_t, \varphi_t), \sigma_t^2 I)$$

$$P(x|z_1, \varphi_1) \Rightarrow N(f_1(z_1, \varphi_1), \sigma_1^2 I)$$

# Uczenie modelu

$$\operatorname{argmax}_{\varphi_t} \sum_{i=1}^n \log P(x_i | \varphi_t)$$

Założmy, że  $z_{1..T}$  to wielowymiarowa zmienna losowa pochodząca z rozkładu łącznego zmiennych  $z_1, \dots, z_T$ . Wówczas

$$\int P(z_{1..T} | x) \log \left( \frac{P(x, z_{1..T} | \varphi_{1,..,T})}{P(z_{1..T} | x)} \right) dz_{1..T} \leq \log P(x | \varphi_{1,..,T})$$

# Uczenie modelu

$$\int P(z_{1..T}|x) \log \left( \frac{P(x, z_{1..T} | \varphi_{1,..,T})}{P(z_{1..T}|x)} \right) dz_{1..T} \approx$$
$$\mathbb{E}_{P(z_1|x)} \log(P(x|z_1, \varphi_1)) - \sum_{i=2}^T \mathbb{E}_{P(z_i|x)} D_{KL}(P(z_{i-1}|z_i, x) || P(z_{i-1}|z_i, \varphi_i))$$

# Uczenie modelu

## Diffusion loss function

$$L[\phi_{1...T}] = \sum_{i=1}^I \left( \overbrace{-\log \left[ \text{Norm}_{\mathbf{x}_i} \left[ \mathbf{f}_1[\mathbf{z}_{i1}, \phi_1], \sigma_1^2 \mathbf{I} \right] \right]}^{\text{reconstruction term}} + \sum_{t=2}^T \frac{1}{2\sigma_t^2} \left\| \underbrace{\frac{1 - \alpha_{t-1}}{1 - \alpha_t} \sqrt{1 - \beta_t} \mathbf{z}_{it} + \frac{\sqrt{\alpha_{t-1}} \beta_t}{1 - \alpha_t} \mathbf{x}_i}_{\text{target, mean of } q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x})} - \underbrace{\mathbf{f}_t[\mathbf{z}_{it}, \phi_t]}_{\text{predicted } \mathbf{z}_{t-1}} \right\|^2 \right),$$

# Modyfikacja diffusion loss

Pamiętając, że

$$z_t = \sqrt{\alpha_t}x + \sqrt{1 - \alpha_t}\epsilon,$$

mamy

$$x = \frac{1}{\sqrt{\alpha_t}}z_t + \frac{\sqrt{1 - \alpha_t}}{\sqrt{\alpha_t}}\epsilon$$

# Uczenie modelu

## Diffusion loss function

$$L[\phi_{1...T}] = \sum_{i=1}^I \left( -\log \left[ \text{Norm}_{\mathbf{x}_i} \left[ \mathbf{f}_1[\mathbf{z}_{i1}, \phi_1], \sigma_1^2 \mathbf{I} \right] \right] \right. \\ \left. + \sum_{t=2}^T \frac{1}{2\sigma_t^2} \left\| \left( \frac{1}{\sqrt{1-\beta_t}} \mathbf{z}_{it} - \frac{\beta_t}{\sqrt{1-\alpha_t}\sqrt{1-\beta_t}} \epsilon_{it} \right) - \mathbf{f}_t[\mathbf{z}_{it}, \phi_t] \right\|^2 \right).$$

Pamiętając, że

$$z_t = \sqrt{1 - \beta_t} z_{t-1} + \sqrt{\beta_t} \epsilon_t$$

mamy

$$z_{t-1} = \frac{1}{\sqrt{1 - \beta_t}} z_t + \frac{\sqrt{\beta_t}}{\sqrt{1 - \beta_t}} \epsilon_t$$

Pamiętając, że

$$z_t = \sqrt{1 - \beta_t} z_{t-1} + \sqrt{\beta_t} \epsilon_t$$

mamy

$$z_{t-1} = \frac{1}{\sqrt{1 - \beta_t}} z_t + \frac{\sqrt{\beta_t}}{\sqrt{1 - \beta_t}} \epsilon_t$$

A skoro  $z_{t-1} = f_t(z_t, \varphi_t)$  to

$$z_{t-1} = \frac{1}{\sqrt{1 - \beta_t}} z_t + \frac{\sqrt{\beta_t}}{\sqrt{1 - \beta_t}} g_t(z_t, \varphi_t)$$



# Uczenie modelu

Diffusion loss function

$$\begin{aligned} L[\phi_{1...T}] &= \sum_{i=1}^I \sum_{t=1}^T \left\| \mathbf{g}_t[\mathbf{z}_{it}, \phi_t] - \epsilon_{it} \right\|^2 \\ &= \sum_{i=1}^I \sum_{t=1}^T \left\| \mathbf{g}_t \left[ \sqrt{\alpha_t} \cdot \mathbf{x}_i + \sqrt{1 - \alpha_t} \cdot \epsilon_{it}, \phi_t \right] - \epsilon_{it} \right\|^2, \end{aligned}$$

# Uczenie modelu

---

**Algorithm 18.1:** Diffusion model training

---

**Input:** Training data  $\mathbf{x}$

**Output:** Model parameters  $\phi_t$

**repeat**

**for**  $i \in \mathcal{B}$  **do** // For every training example index in batch  
         $t \sim \text{Uniform}[1, \dots, T]$  // Sample random timestep  
         $\epsilon \sim \text{Norm}[\mathbf{0}, \mathbf{I}]$  // Sample noise  
         $\ell_i = \left\| \mathbf{g}_t \left[ \sqrt{\alpha_t} \mathbf{x}_i + \sqrt{1 - \alpha_t} \epsilon, \phi_t \right] - \epsilon \right\|^2$  // Compute individual loss

    Accumulate losses for batch and take gradient step

**until** converged

---

# Generowanie danych

---

**Algorithm 18.2:** Sampling

---

**Input:** Model,  $\mathbf{g}_t[\bullet, \phi_t]$

**Output:** Sample,  $\mathbf{x}$

$\mathbf{z}_T \sim \text{Norm}_{\mathbf{z}}[\mathbf{0}, \mathbf{I}]$  // Sample last latent variable

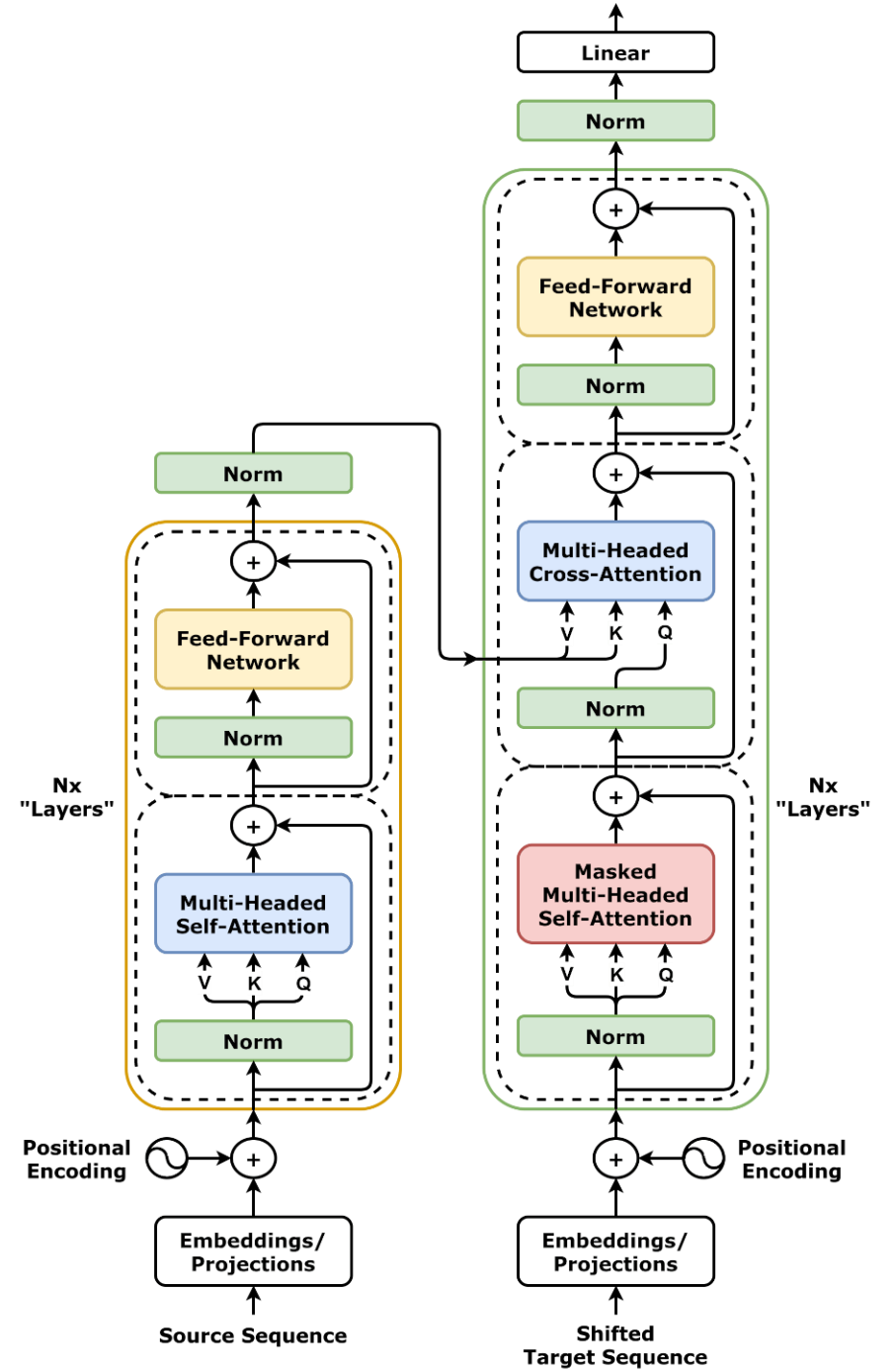
**for**  $t = T \dots 2$  **do**

$\hat{\mathbf{z}}_{t-1} = \frac{1}{\sqrt{1-\beta_t}}\mathbf{z}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}\sqrt{1-\beta_t}}\mathbf{g}_t[\mathbf{z}_t, \phi_t]$  // Predict previous latent variable  
     $\epsilon \sim \text{Norm}_{\epsilon}[\mathbf{0}, \mathbf{I}]$  // Draw new noise vector  
     $\mathbf{z}_{t-1} = \hat{\mathbf{z}}_{t-1} + \sigma_t \epsilon$  // Add noise to previous latent variable

$\mathbf{x} = \frac{1}{\sqrt{1-\beta_1}}\mathbf{z}_1 - \frac{\beta_1}{\sqrt{1-\alpha_1}\sqrt{1-\beta_1}}\mathbf{g}_1[\mathbf{z}_1, \phi_1]$  // Generate sample from  $\mathbf{z}_1$  without noise

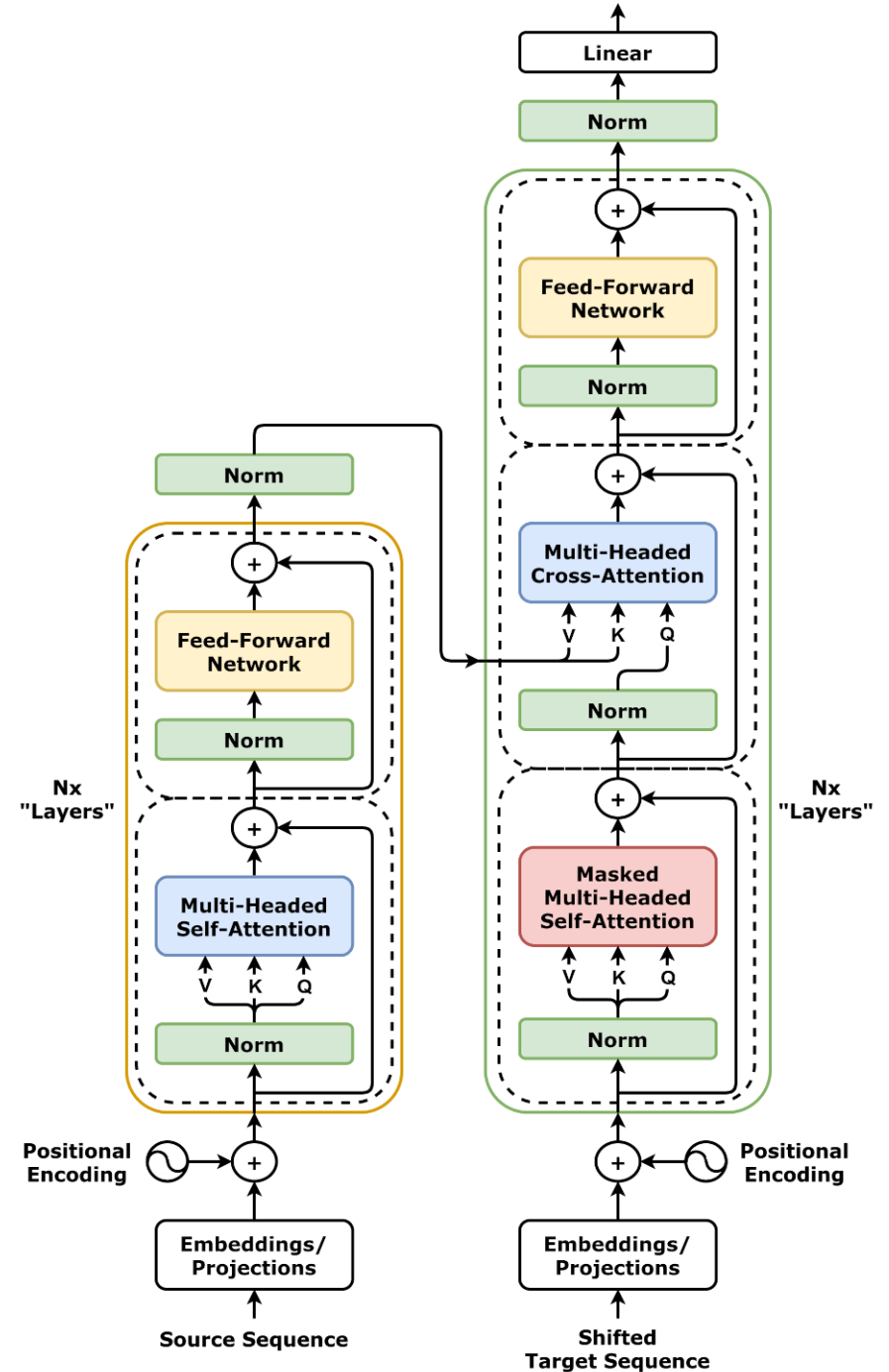
---

# Transformery



# Transformery

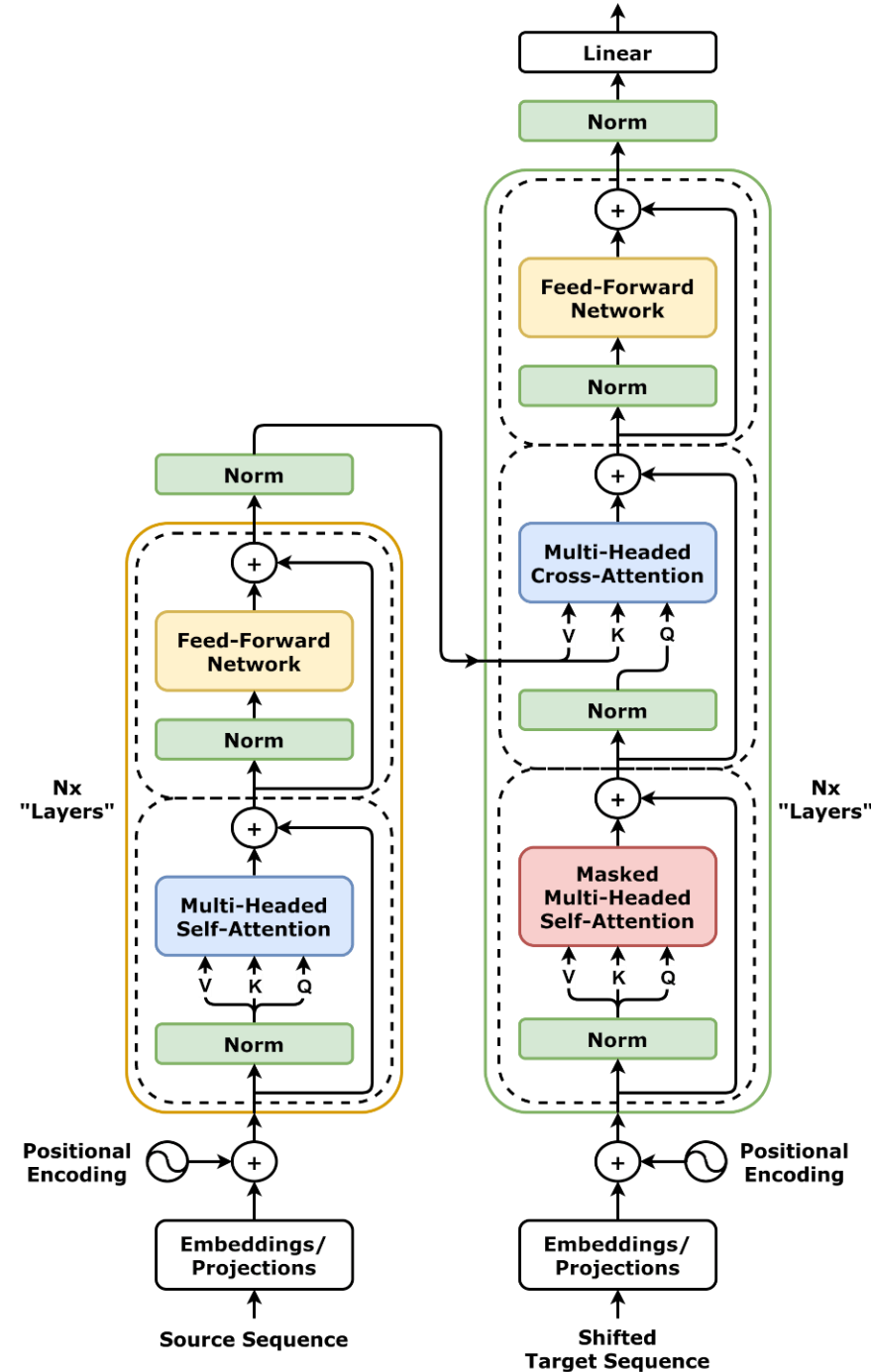
## Struktura transformera Encoder-Decoder



# Transformery

## Struktura transformera Encoder-Decoder

## Elementy struktur

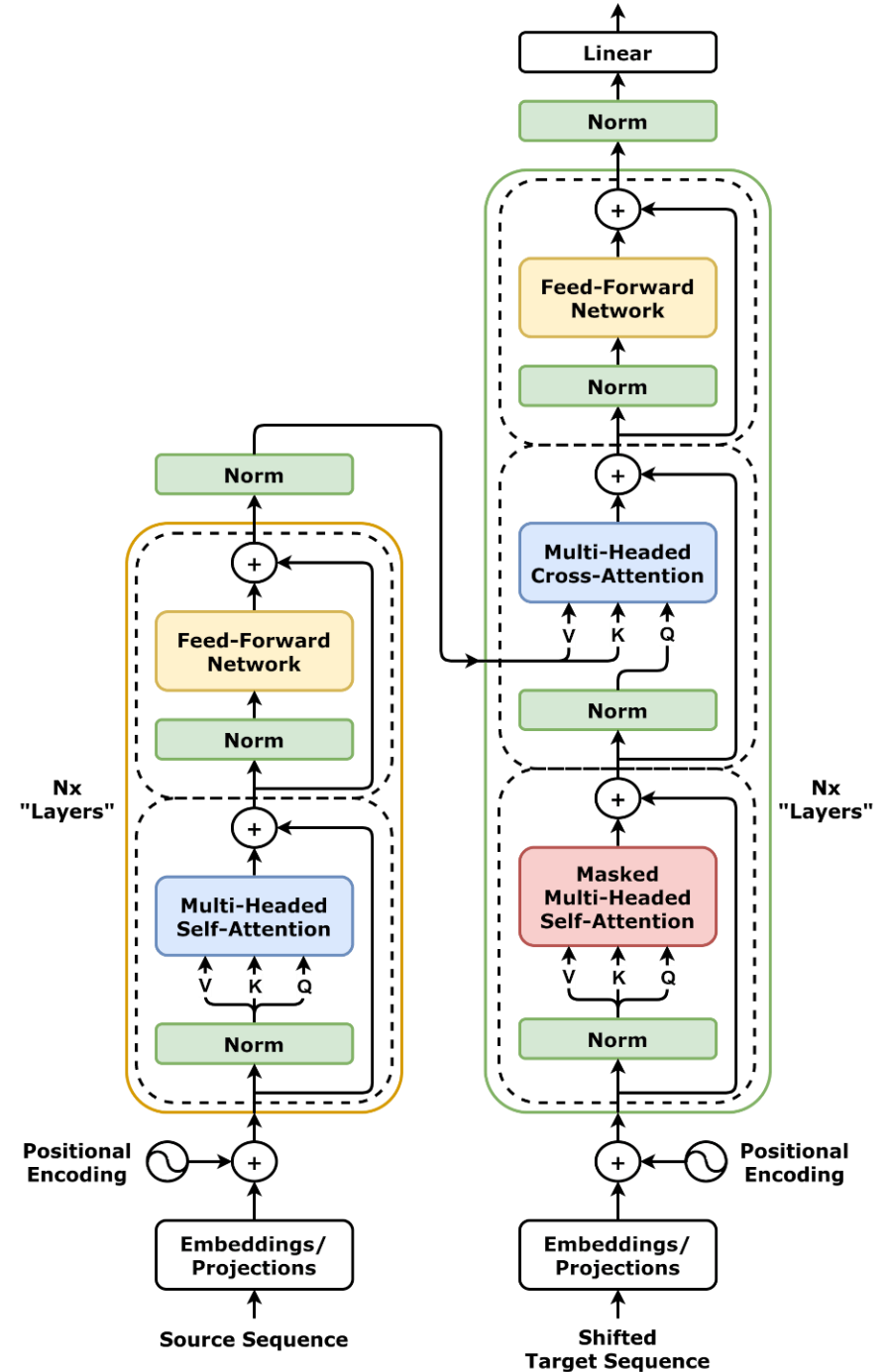


# Transformery

## Struktura transformera Encoder-Decoder

### Elementy struktur:

- Multi-head attention

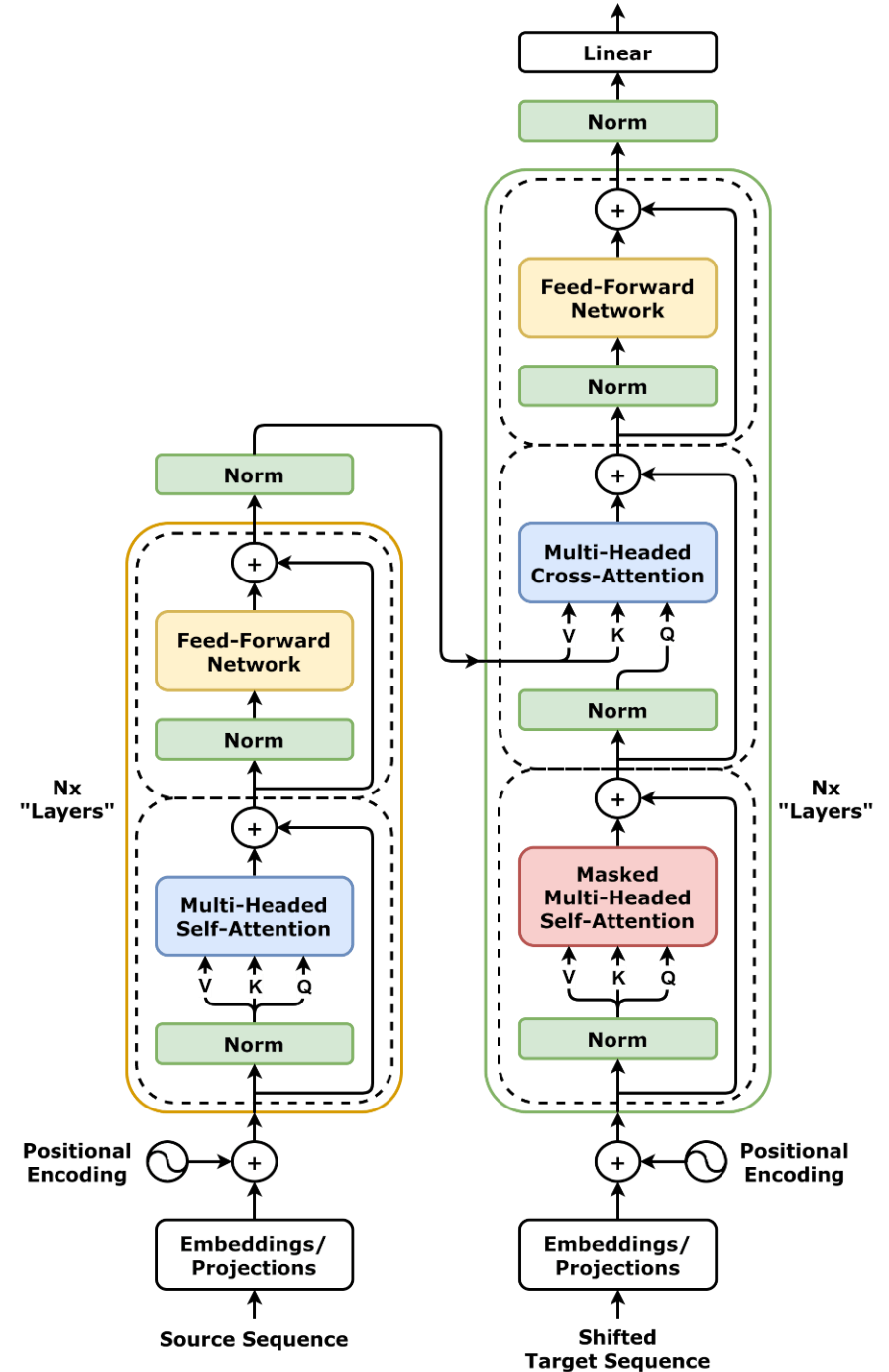


# Transformery

## Struktura transformera Encoder-Decoder

### Elementy struktur:

- Multi-head attention
- Embeddingi





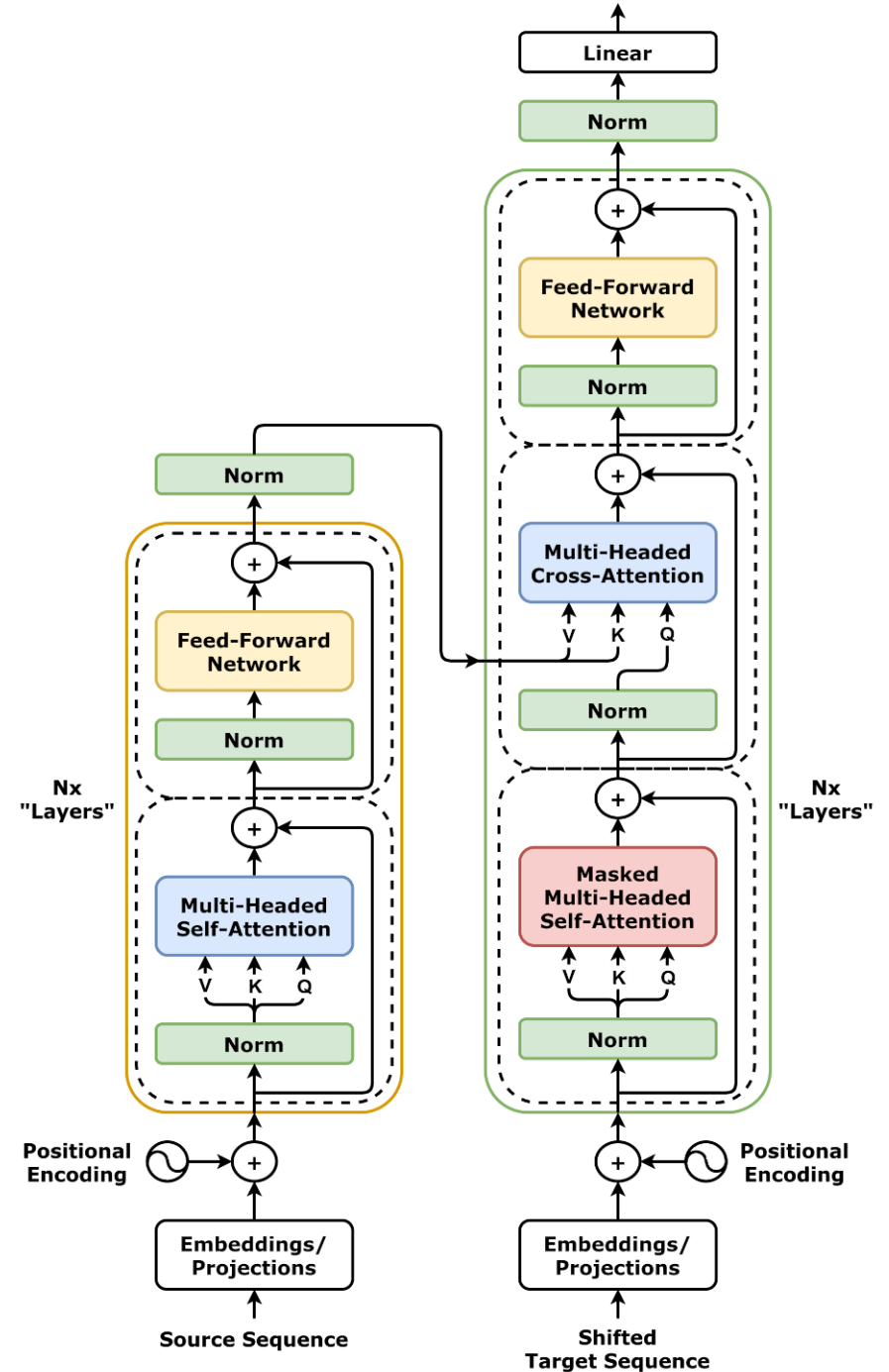
# Transformery

## Struktura transformera Encoder-Decoder

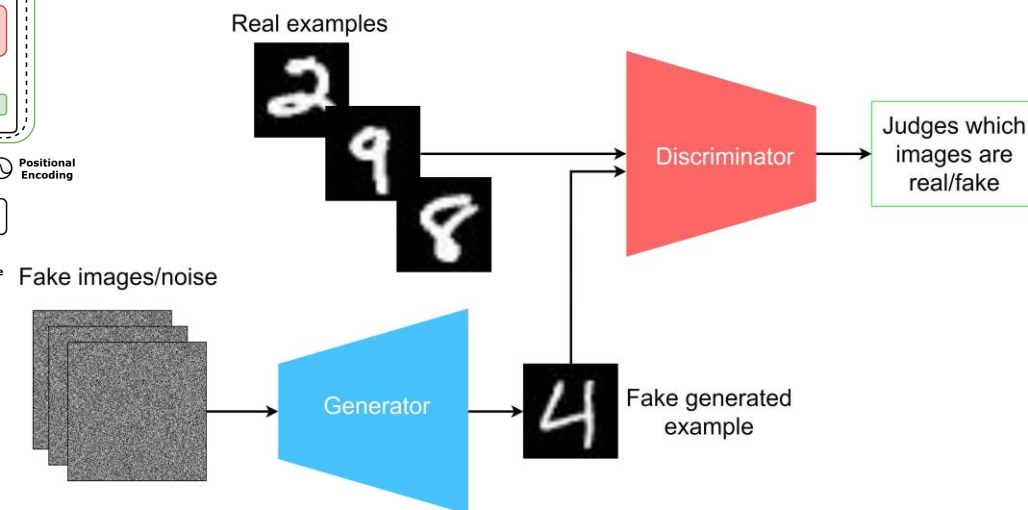
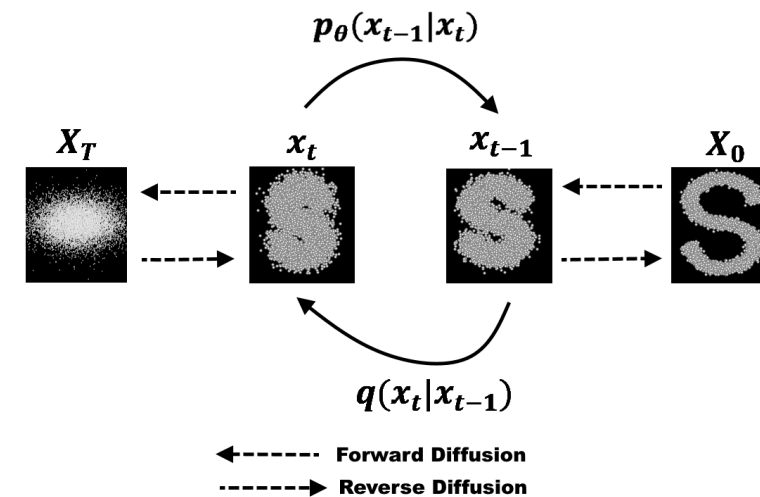
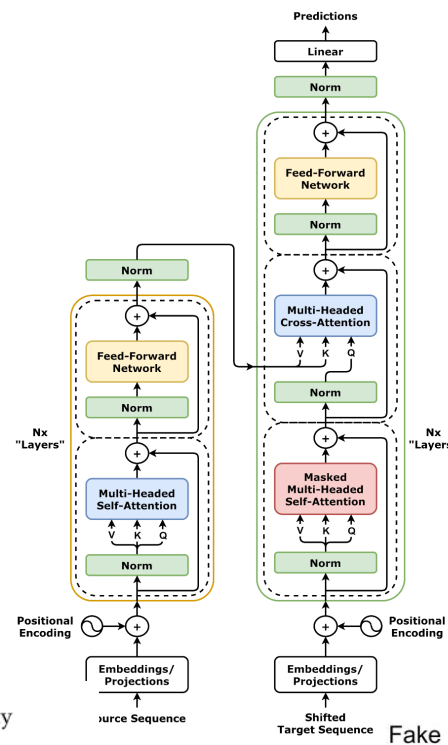
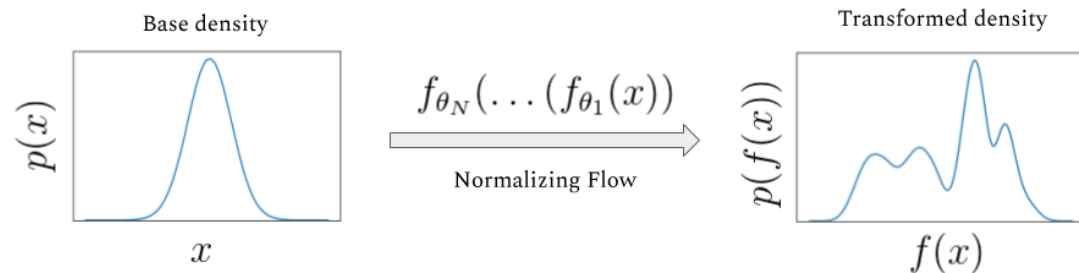
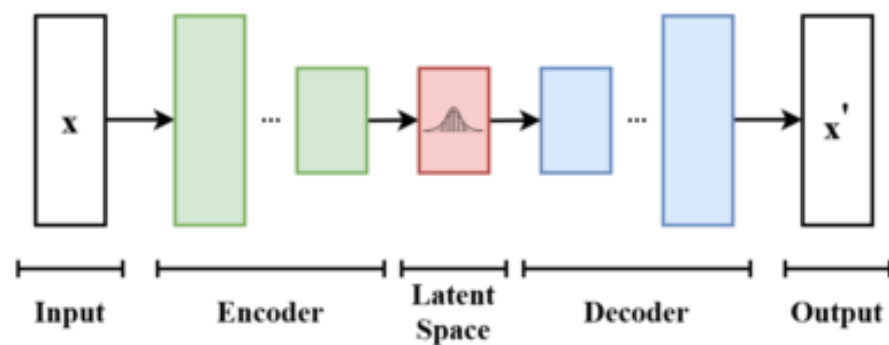
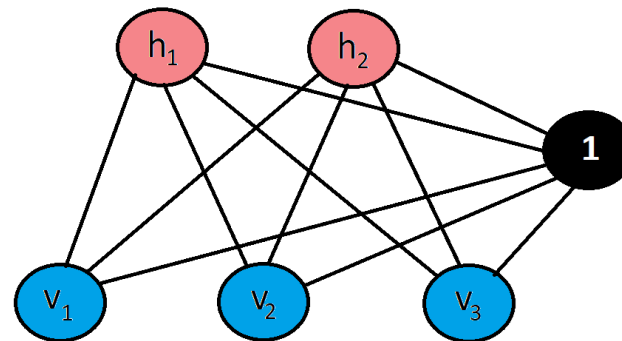
### Elementy struktur:

- Multi-head attention
- Embeddingi

### Uczenie



# Podsumowanie



# Porównanie podejść generacyjnych

## zalety - wady

Metoda	Zalety	Wady
<b>Modele dyfuzyjne</b>	<ul style="list-style-type: none"><li>- Wysoka jakość generowanych próbek</li><li>- Stabilność treningu</li><li>- Umożliwiają iteracyjną kontrolę nad generacją</li></ul>	<ul style="list-style-type: none"><li>- Wysokie wymagania obliczeniowe</li><li>- Powolny proces generacji (wiele kroków samplingu)</li></ul>
<b>Normalizing Flows</b>	<ul style="list-style-type: none"><li>- Łatwość obliczenia prawdopodobieństwa danych</li><li>- Odwracalność modelu</li></ul>	<ul style="list-style-type: none"><li>- Trudność w skalowaniu do bardziej złożonych danych</li><li>- Wymagają bardzo dokładnych parametrów</li></ul>
<b>GAN</b>	<ul style="list-style-type: none"><li>- Szybka generacja danych</li><li>- Wysoka jakość próbek, zwłaszcza obrazów</li></ul>	<ul style="list-style-type: none"><li>- Trudność w stabilnym trenowaniu (problem zbalansowania generatora i dyskryminatora)</li><li>- Mniej intuicyjna kontrola</li></ul>
<b>VAE</b>	<ul style="list-style-type: none"><li>- Łatwość implementacji</li><li>- Naturalne modelowanie zmiennych latentnych</li></ul>	<ul style="list-style-type: none"><li>- Generowane próbki mogą być mniej realistyczne niż GAN</li><li>- Rozmyte wyniki dla obrazów</li></ul>
<b>WAE</b>	<ul style="list-style-type: none"><li>- Lepsza jakość niż klasyczne VAE</li><li>- Stabilniejszy trening</li></ul>	<ul style="list-style-type: none"><li>- Wciąż słabsza jakość generacji niż GAN</li><li>- Problemy z kosztami obliczeń dla dużych zbiorów danych</li></ul>
<b>RBM</b>	<ul style="list-style-type: none"><li>- Prosta implementacja dla małych zbiorów danych</li><li>- Dobre do modelowania nieliniowych relacji</li></ul>	<ul style="list-style-type: none"><li>- Problemy ze skalowalnością</li><li>- Słaba wydajność przy bardziej złożonych danych</li></ul>

# Porównanie podejść generacyjnych

## Zastosowania

**Modele dyfuzyjne:** Nadają się do aplikacji wymagających wysokiej jakości i elastyczności generacji, np. edycji obrazów.

**Normalizing Flows:** Używane, gdy obliczenie dokładnego prawdopodobieństwa danych jest istotne (np. w analizie anomalii).

**GAN:** Idealne do generowania realistycznych obrazów. Popularne w grafice komputerowej i sztuce cyfrowej.

**VAE:** Sprawdzają się w analizie danych i uczeniu reprezentacji, gdzie zmienne latentne są ważne.

**WAE:** Dobry wybór dla kompromisu między prostotą a jakością generowanych danych.

**RBM:** Raczej przestarzałe, dobre dla małych zbiorów danych lub nauki podstawowych zasad generatywnych.

# Porównanie podejść generacyjnych

## Implementacje

Metoda	implementacja	Krótki opis
Modele dyfuzyjne	Stable Diffusion / DALL-E 2	Modele do generacji obrazów wysokiej jakości z opisów tekstowych.
Normalizing Flows	RealNVP / Glow	RealNVP wprowadził odwracalne transformacje, a Glow poprawił ich efektywność.
GAN	StyleGAN / BigGAN	StyleGAN jest używany do realistycznej generacji obrazów twarzy, BigGAN do dużych zbiorów danych.
VAE	Beta-VAE	Modyfikacja klasycznego VAE z regulacją wpływu zmiennej latentnej na generację.
WAE	WAE-GAN / WAE-MMD	Połączenie idei autoenkoderów z dystansem Wassersteina dla lepszej jakości wyników.
RBM	Deep Belief Networks (DBN)	Rozszerzenie RBM do głębszych struktur, używane w analizie danych.

# Obszary zastosowań a modele

Zastosowanie	Model / Metoda
Generowanie obrazów	Stable Diffusion
	StyleGAN
Generowanie tekstu	GPT
	VAE-LSTM
Generowanie dźwięku / mowy	WaveNet
	DiffWave
Generowanie muzyki	MuseNet
	Jukebox
Obróbka i edycja obrazów	DALL-E
	CycleGAN

# Obszary zastosowań a modele

Zastosowanie	Model / Metoda
Uczenie reprezentacji danych	Beta-VAE
	Deep Belief Networks
Wykrywanie anomalii	Glow
	VAE
Tworzenie 3D	NeRF
	GAN3D
Generowanie wideo	Video Diffusion Models
	MoCoGAN
Symulacje w nauce	Schrodinger Bridge Models
	PhysicsGAN

# Przegląd implementacji

Stable Diffusion 3.5

<https://stability.ai/news/introducing-stable-diffusion-3-5>

StyleGAN

<https://github.com/NVlabs/stylegan>

GPT

<https://github.com/openai/gpt-2>

WaveNet

<https://deepmind.google/discover/blog/wavenet-a-generative-model-for-raw-audio/>



# Przegląd implementacji

MuseNet

<https://openai.com/index/musenet/>

Jukebox

<https://openai.com/index/jukebox/>

DALL-E

<https://openai.com/index/dall-e-3/>

CycleGAN

<https://junyanz.github.io/CycleGAN/>

# Od czego zacząć?

## Documentations

 Search across all docs

### • Hub

Host Git-based models, datasets and Spaces on the Hugging Face Hub.

### • Hub Python Library

Client library for the HF Hub: manage repositories from your Python runtime.

### • Inference API (serverless)

Experiment with over 200k models easily using the serverless tier of Inference Endpoints.

### • Accelerate

### • Transformers

State-of-the-art ML for Pytorch, TensorFlow, and JAX.

### • Datasets

Access and share datasets for computer vision, audio, and NLP tasks.

### • Huggingface.js

A collection of JS libraries to interact with Hugging Face, with TS types included.

### • Inference Endpoints (dedicated)

Easily deploy models to production on dedicated, fully managed infrastructure.

### • Optimum

### • Diffusers

State-of-the-art diffusion models for image and audio generation in PyTorch.

### • Gradio

Build machine learning demos and other web apps, in just a few lines of Python.

### • Transformers.js

State-of-the-art Machine Learning for the web. Run Transformers directly in your browser, with no need for a server.

### • PEFT

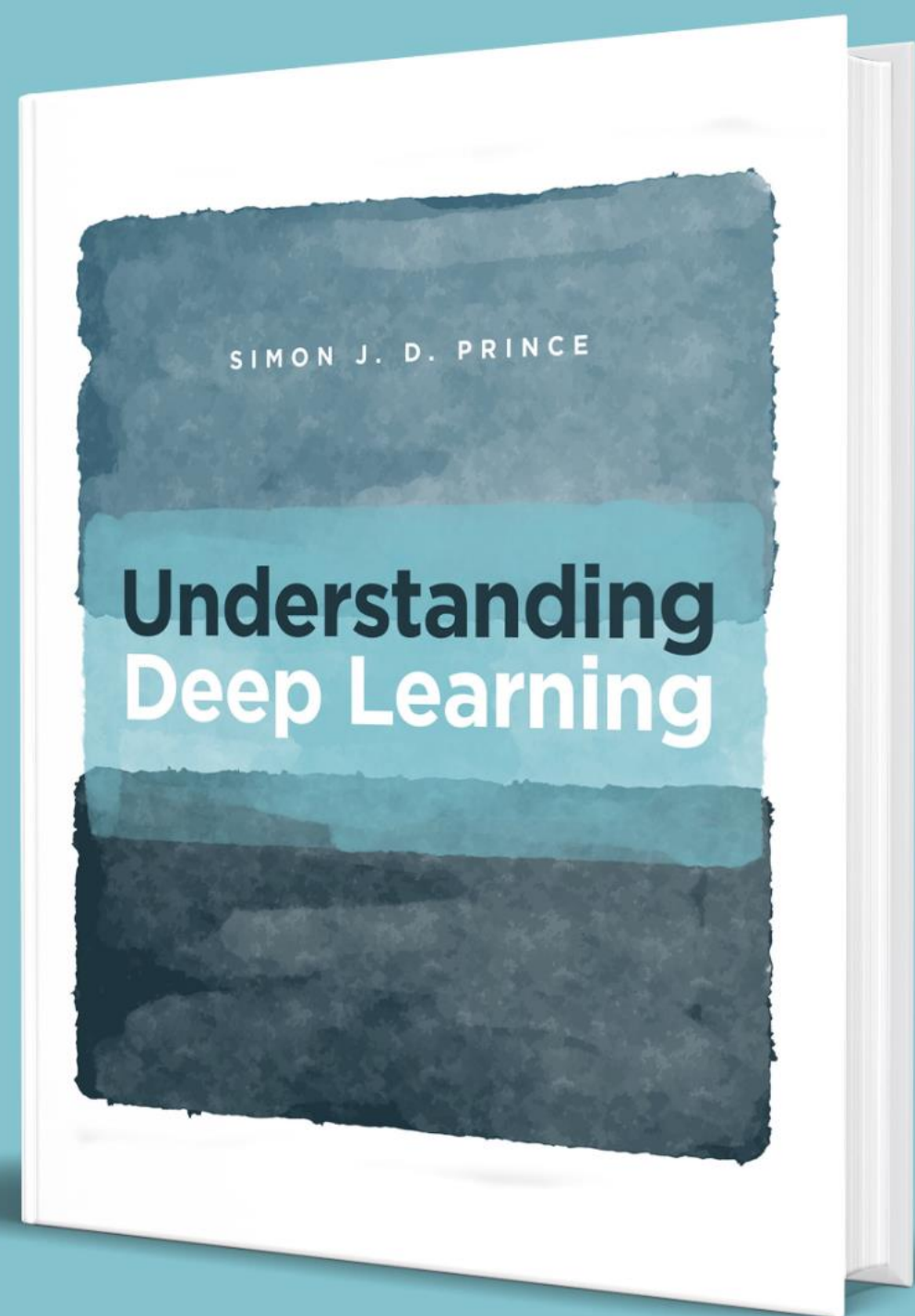
Parameter efficient finetuning methods for large models.

### • AWS Trainium & Inferentia

# Gdzie szukać teorii?

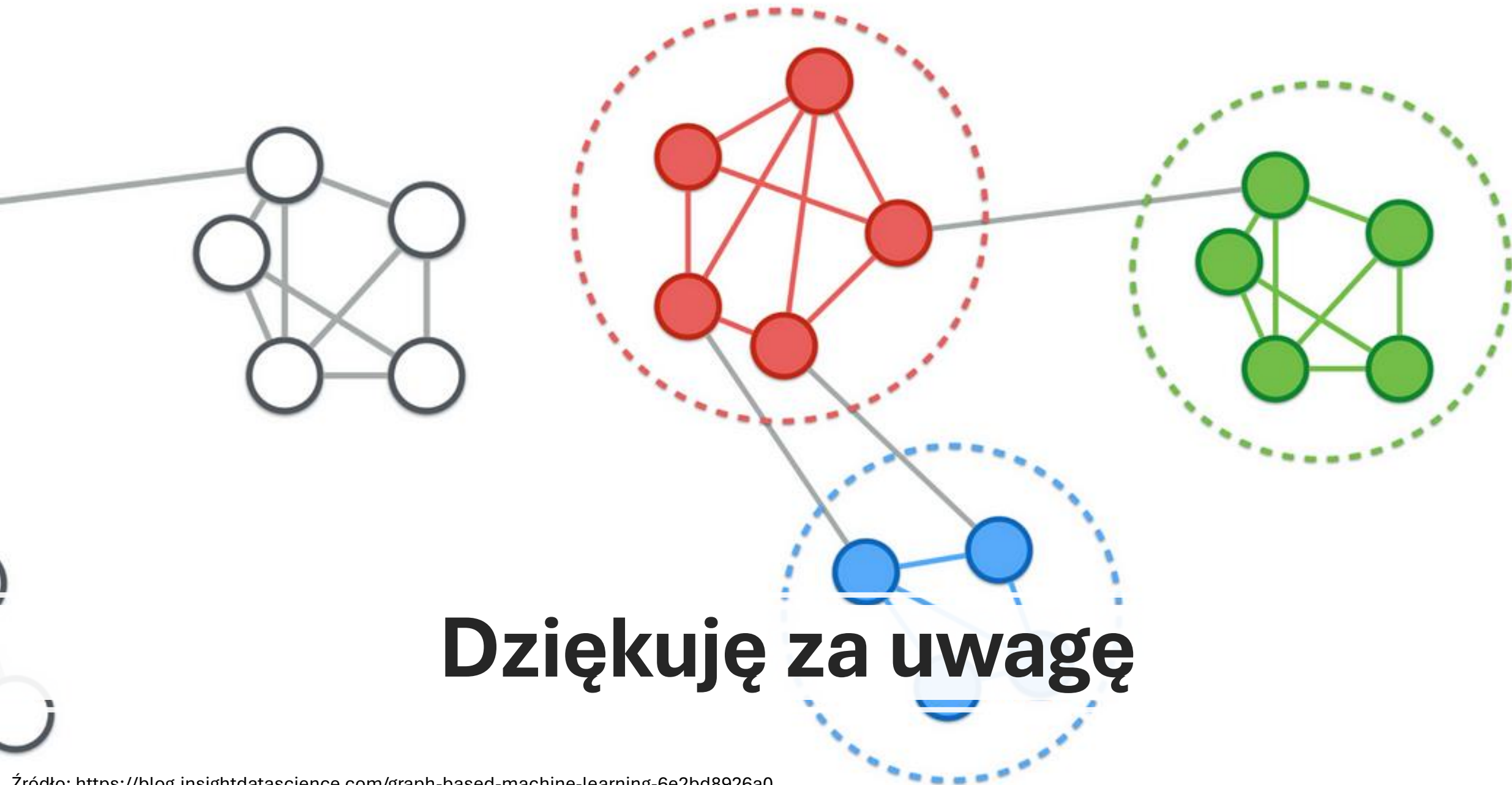
Simon JD. Prince,  
*Understanding deep learning.*  
MIT press, 2023.

<https://udlbook.github.io/udlbook/>



Pytania / Komentarze / Uwagi





**Dziękuję za uwagę**