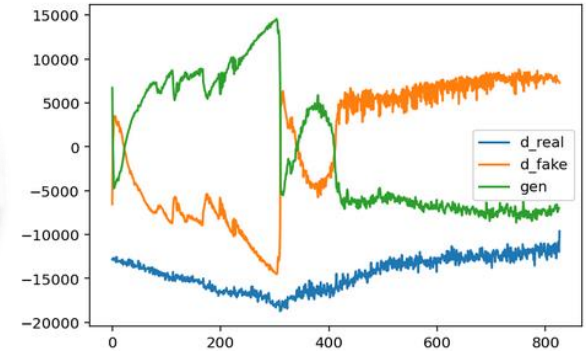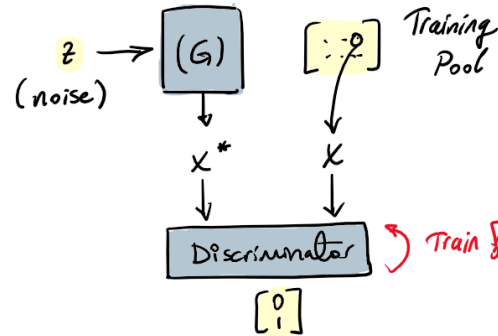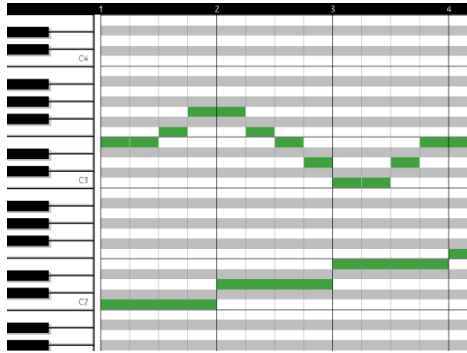# Data Driven Engineering I: Machine Learning for Dynamical Systems

**Introduction to Generative Learning: VAEs and GANs**

Institute of Thermal Turbomachinery
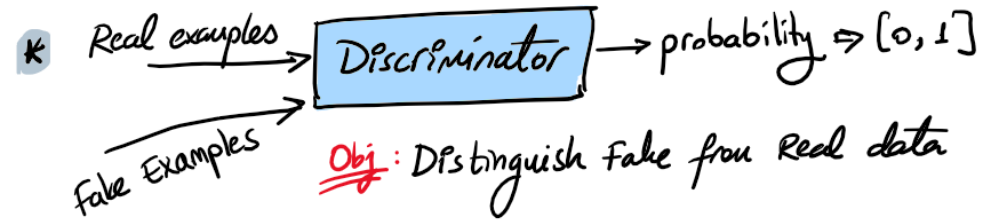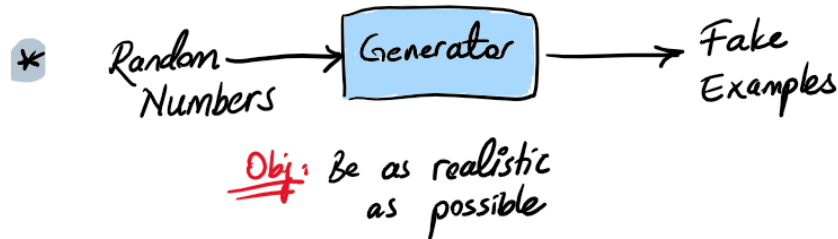Prof. Dr.-Ing. Hans-Jörg Bauer

# Generative Learning

* Generative Learning & Representation

* Latent Space $\Rightarrow$ Autoencoders $\rightarrow$ M' is easier to learn

* AE + Gaussian Sampling $\Rightarrow$ VAEs $\circ$ $\rightarrow$ Functional API
  $\rightarrow$ $\lambda$ Layer
  $\rightarrow$ Embeddings

* Music $\rightarrow$ MIDI $\rightarrow$ VAE, LSTM
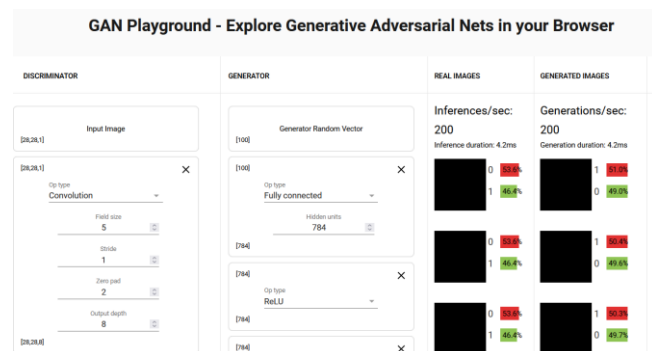
This Week : GANs for Music

# Generative Adversarial Networks: GANs

* Generative $\Rightarrow$ creating non-existing data

* Adversarial $\Rightarrow$ Competitive dynamics (game-like)

* Network $\Rightarrow$ Neural networks

* GANS (2014) → Generator
         → Discriminator

* Random Numbers → [Generator] → Fake Examples

  **Obj:** Be as realistic as possible

* Real examples → [Discriminator] → probability $\Rightarrow$ $[0, 1]$
  Fake Examples →

  **Obj:** Distinguish Fake from Real data

Everybody can dance now

Source Subject · Target Subject 1 · Target Subject 2 · Source Subject · Target Subject 1 · Target Subject 2

# Training GANs

* In MLP, we have a clear goal & measure

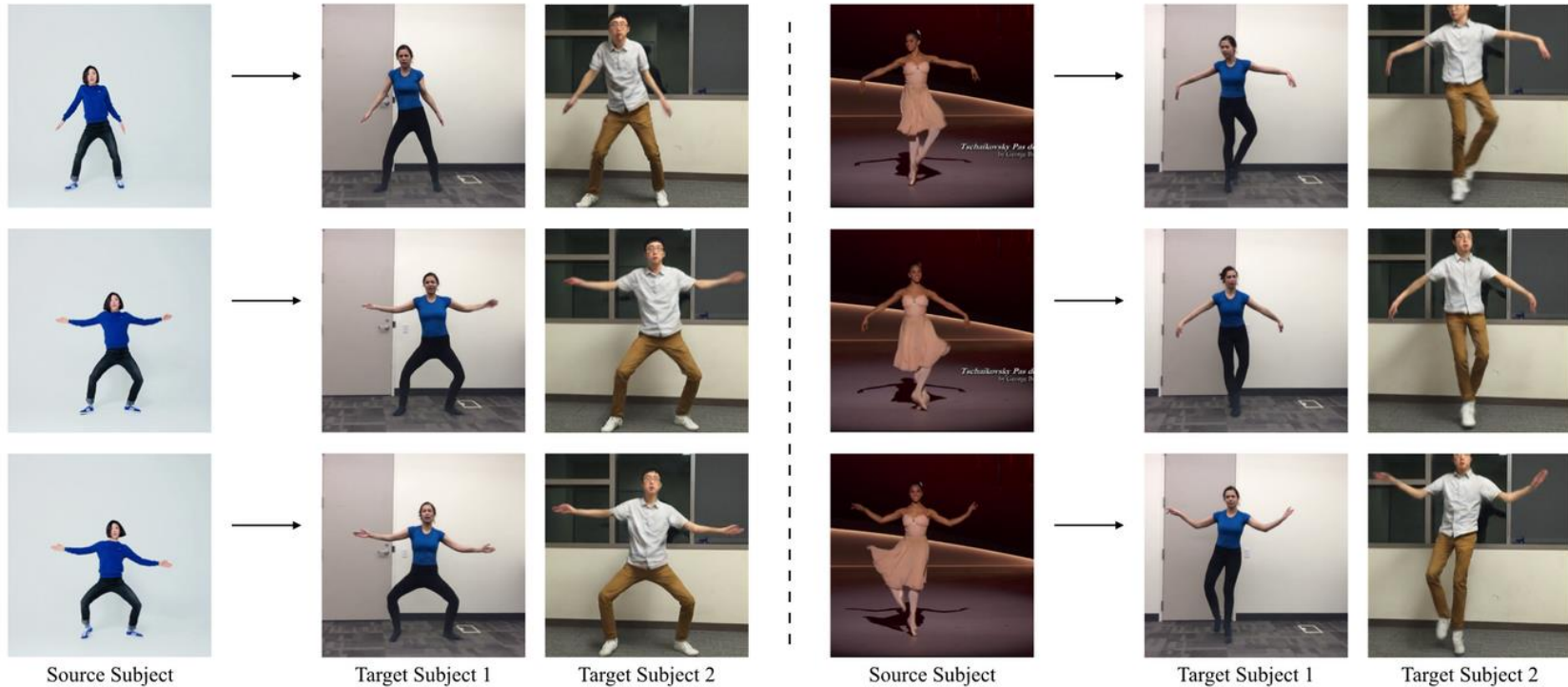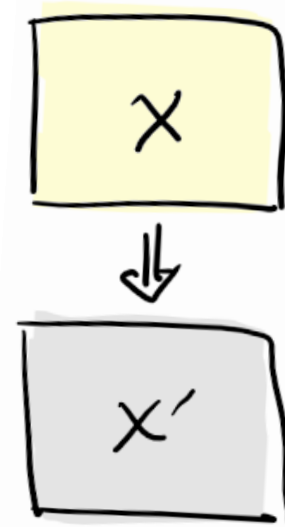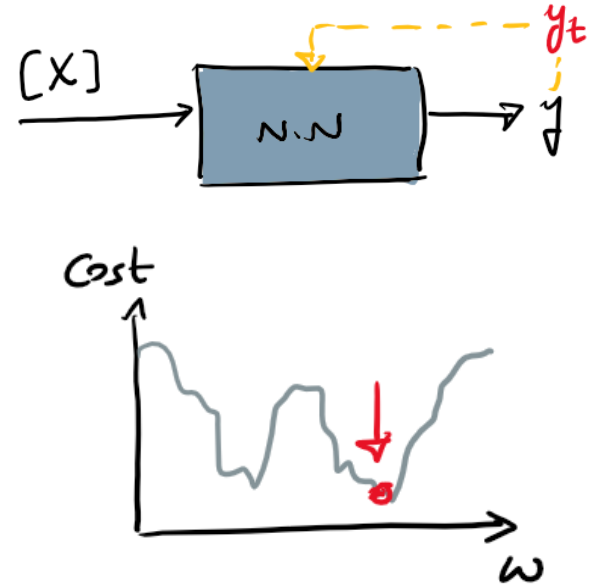  ~~Ex~~ Minimize Cross-entropy loss.

* In GANs, two networks have competing obj. ⚠

  $(G) \uparrow ; (D) \downarrow \, // \, (D) \uparrow ; (G) \downarrow$

# Training GANs

## Minmax Game

\* Adversarial ML

Game Theory
- \* Player 1 — Min.
- \* Player 2 — Max.

$w_1$

$w_2$

Nash Eq.

# Training GANs

**!** Nash Equilibrium := Point where neither "player" can improve their situation

- (G) := Fakes are indistinguishable from Real data

- (D) := at best random guess ( F/R ⟹ 1)

**!** In practice; ~ impossible to achieve Nash Eq.

it still works ...

# Training GANs

Loss Function        Likelihood $\Rightarrow$ Classification problem

$$\text{Value}(G, D, X, y) = \mathbb{E}\left(\log D(x)\right) + \mathbb{E}\left(1 - D(G(z))\right)$$

noise

log. probability D
correctly predicts reals

log. prob. D correctly predicts
fakes are fakes

Discriminator $\Rightarrow$ max. accuracy. of D

Generator $\Rightarrow$ min. accuracy of D

# Training Algorithm:

For each training do:

## # Train (D):

(1) Take a random real example from Training data, $x$

(2) Get a fake example from Generator, $x^*$

(3) Use Discriminator to classify $x$ & $x^*$.

(4) Compute the class. error.

(5) Backprop. error & update Discriminator trainable parameters.

$z$ (noise) $\rightarrow$ (G)

Training Pool

$x^*$          $x$

Discriminator  } Train

$\begin{bmatrix} 0 \\ 1 \end{bmatrix}$

Likelihood $\Rightarrow$ Classification problem

# Training Algorithm:

## #Train (G):

(6) Generate a new fake $X^*$.

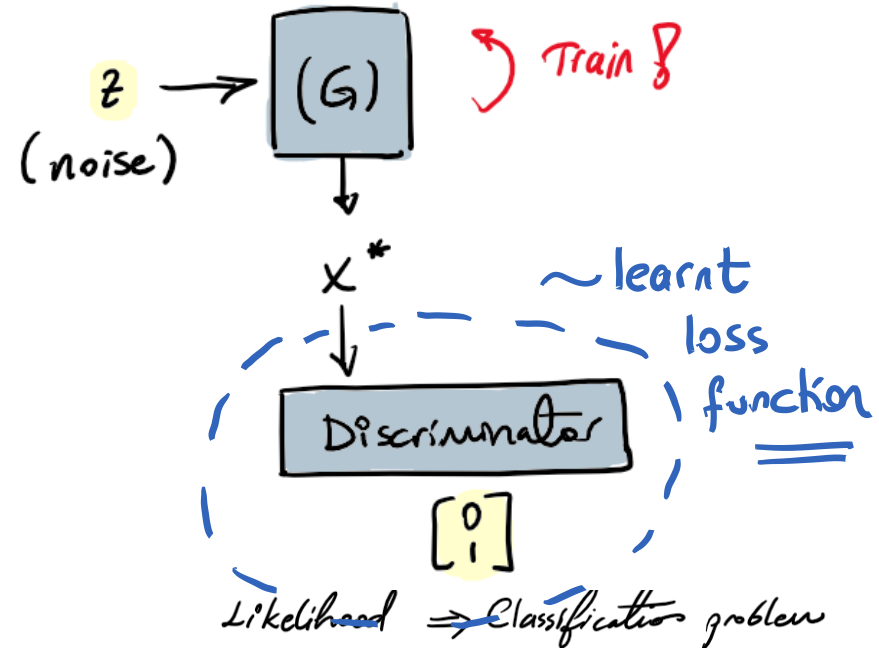(7) Use ~~Discriminator~~ to classify $X^*$.

(8). Compute the error.

(9). Update ~~Generator~~'s trainable parameters via backprop.

end for

$z$ (noise) $\rightarrow$ (G) $\quad$ Train

$\downarrow$

$X^*$

$\downarrow$

Discriminator $\quad$ ~learnt loss function

$\begin{bmatrix} 0 \\ 1 \end{bmatrix}$

Likelihood $\Rightarrow$ Classification problem

# Putting it together...

for i steps of training, do:

for j steps of D training, do:

→ Sample mini-batch of real $X; y=1$

→ Create mini-batch of fake $X^*; y=0$

→ Update D Model

$$\frac{1}{N} \sum^{N} \left[ \log D(x) + \log\left(1 - D(G(z))\right) \right]$$

→ Create a set of fake examples $X^\varepsilon, y=1$

→ Update G model

$$\frac{1}{N} \sum^{N} \log\left(1 - D(G(z))\right)$$

Likelihood $\Rightarrow$ Classification problem

# Training GANs
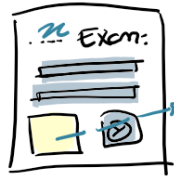
* Training is very difficult

- Disc. usually wins ⟹ make the game unfair ⸮

- Training with more Epoch. can make it worse.

- More data ⟹ more confusion ⟹ may worsen | Baby sitting ⸮

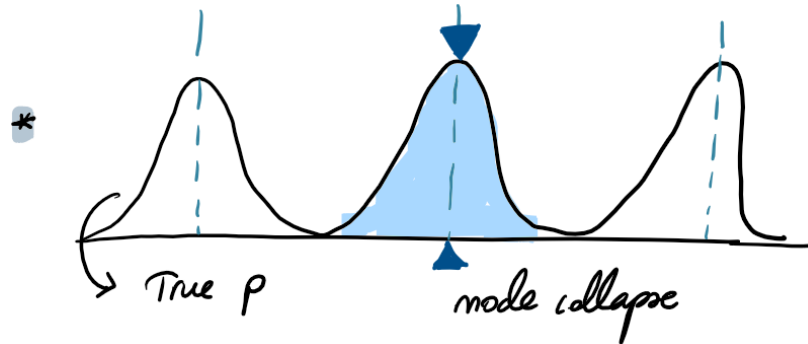- Cross-domain architectures may fail

# Training GANs

**Training problems #1: Mode Collapse**

" Abracadabra 🔔 „ ~ "I create as I speak 🔔

$$\begin{pmatrix} -4 \\ -2 \\ 0 \\ +2 \\ +4 \end{pmatrix}$$

$G \Rightarrow$ "0„ (✓)

$[-4, -2, 2, 4]$

True P

mode collapse

True values

Dr. Cihan Ates- Generative Learning II

Institute of Thermal Turbomachinery (ITS)

# Training GANs

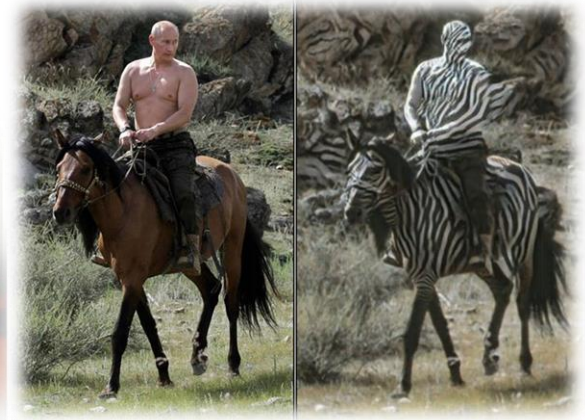Training Problem #2 : Over generalization

* Mods that should not exist, do exist.

* $\begin{bmatrix} -4 \\ -2 \\ 0 \\ +2 \\ +4 \end{bmatrix}$  $\boxed{G}$ $\nearrow$ $-2/3$
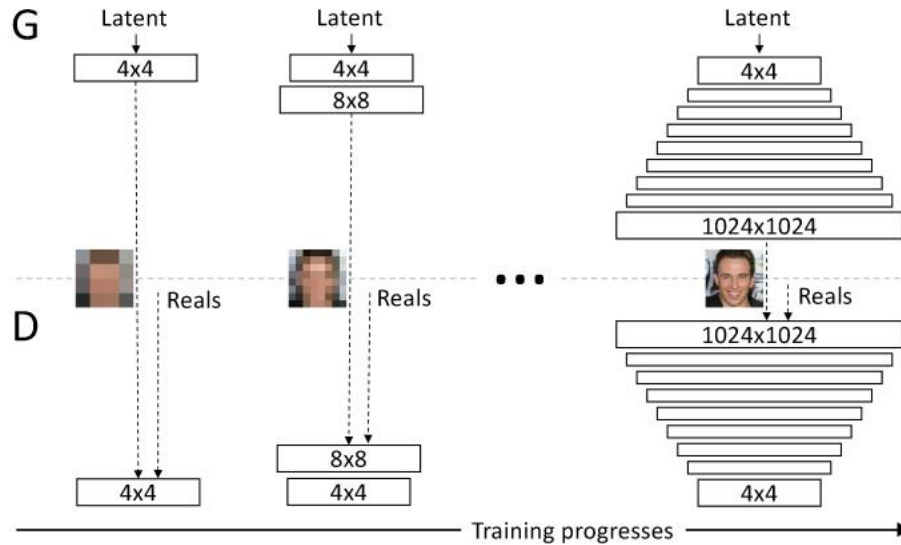$\searrow$ $1/2$
$\dots$
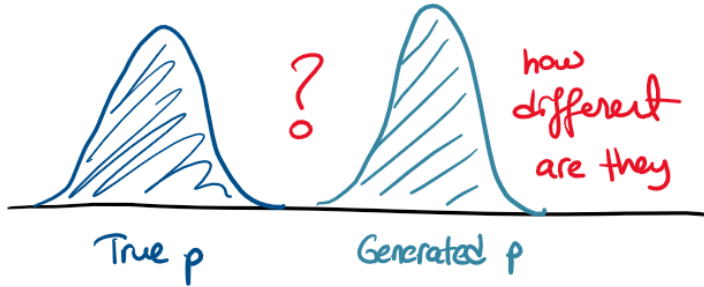
True values
[integers]

[real numbers]

* Image generation
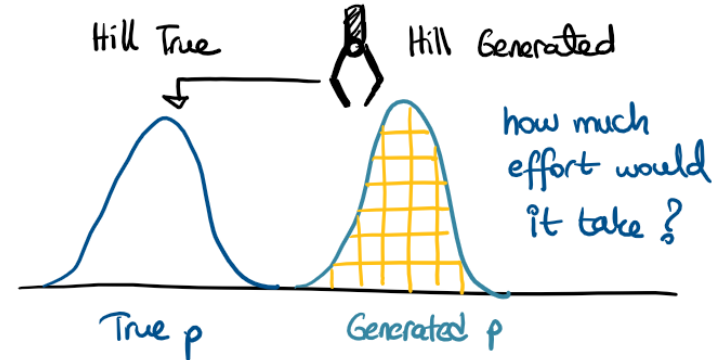
# Possible Remedies

① Growing the network gradually

# Possible Remedies

② Alternative loss definitions ⇒ Wasserstein GAN



True p          Generated p

? how different are they

\* Distance := (similarity) $\begin{pmatrix} TV\ distance \\ KL\ divergence \\ JS\ divergence \\ ... \end{pmatrix}$ ⇒ Earth-mover Distance

Hill True          Hill Generated

True p          Generated p

how much effort would it take?

# Possible Remedies

② Alternative loss definitions ⟹ Wasserstein GAN

**Obj:** Learn correct representation ⟹ "ρ"



$G_1(z)$     $G_2(z)$     Data X

How similar
they are?

- KL Divergence ⟹ no overlap ⟹ ∞

- JS Divergence ⟹ no overlap ⟹ finite, same value

- W Distance ⟹ "$G_2$ is better than $G_1$".

$$E\left(\log D(x)\right) + E\left(\log\left(1 - D\left(G(\overset{\text{noise}}{z})\right)\right)\right)$$

⇓

$$E\left(D(x)\right) - E\left(D(G(\overset{\text{noise}}{z}))\right)$$

Score ⟹ Regressor Problem

## Possible Remedies

② Alternative loss definitions ⇒ Wasserstein GAN

✱ Class. ⇒ $[0, 1]$

↳ Train... ⇒
$$0.9971 \quad \| \quad 0.0004$$
$$0.9965 \quad \| \quad 0.00013$$
$$0.9984 \quad \| \quad 0.00021$$
... ...

⇩ vanishing gradient problem

✱ Stabilized learning

# I/O in Music

① Symbolic Representation

② Audio

* Piano Rolls
  ↓
  "image like"

  Eg:

  * MuseGAN
  * MidiNET



⇒ freq.



① ② ③

↓↓↓ instruments
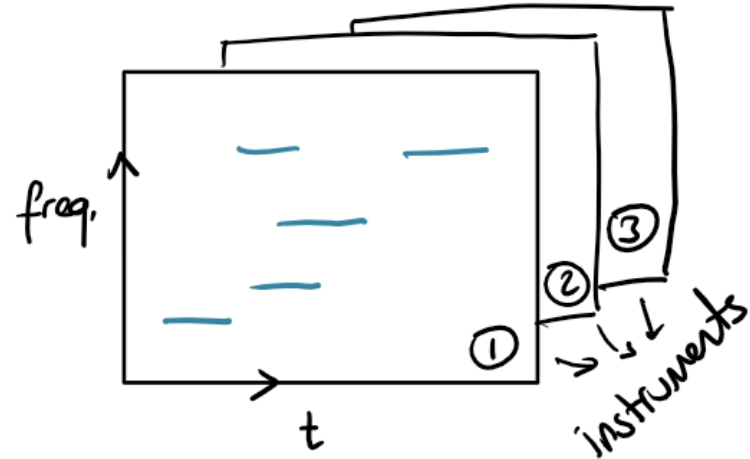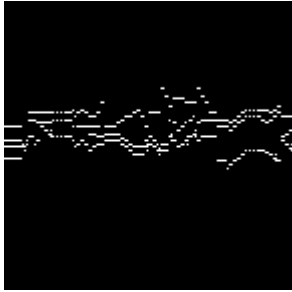
t

# Possible Remedies

## To Do List

(i) Create D model
- W $\Rightarrow$ linear act. function

(ii) Create G model
- Final layer $\Rightarrow$ tanh $[-1, 1]$

(iii) Add weight clipping in D for W

Ateliers & Saveurs in Montreal

colab

**Algorithm 1** WGAN, our proposed algorithm. All experiments in the paper used the default values $\alpha = 0.00005$, $c = 0.01$, $m = 64$, $n_{\text{critic}} = 5$.

---

**Require:** : $\alpha$, the learning rate. $c$, the clipping parameter. $m$, the batch size.
  $n_{\text{critic}}$, the number of iterations of the critic per generator iteration.
**Require:** : $w_0$, initial critic parameters. $\theta_0$, initial generator's parameters.
1: **while** $\theta$ has not converged **do**
2:    **for** $t = 0, ..., n_{\text{critic}}$ **do**
3:        Sample $\{x^{(i)}\}_{i=1}^{m} \sim \mathbb{P}_r$ a batch from the real data.
4:        Sample $\{z^{(i)}\}_{i=1}^{m} \sim p(z)$ a batch of prior samples.
5:        $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^{m} f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^{m} f_w(g_\theta(z^{(i)})) \right]$
6:        $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$
7:        $w \leftarrow \text{clip}(w, -c, c)$
8:    **end for**
9:    Sample $\{z^{(i)}\}_{i=1}^{m} \sim p(z)$ a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^{m} f_w(g_\theta(z^{(i)}))$
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_\theta)$
12: **end while**

---