

# Deep Learning

---

Ulf Brefeld

SSAM!

build: June 25, 2024

Machine Learning Group

Leuphana University of Lüneburg

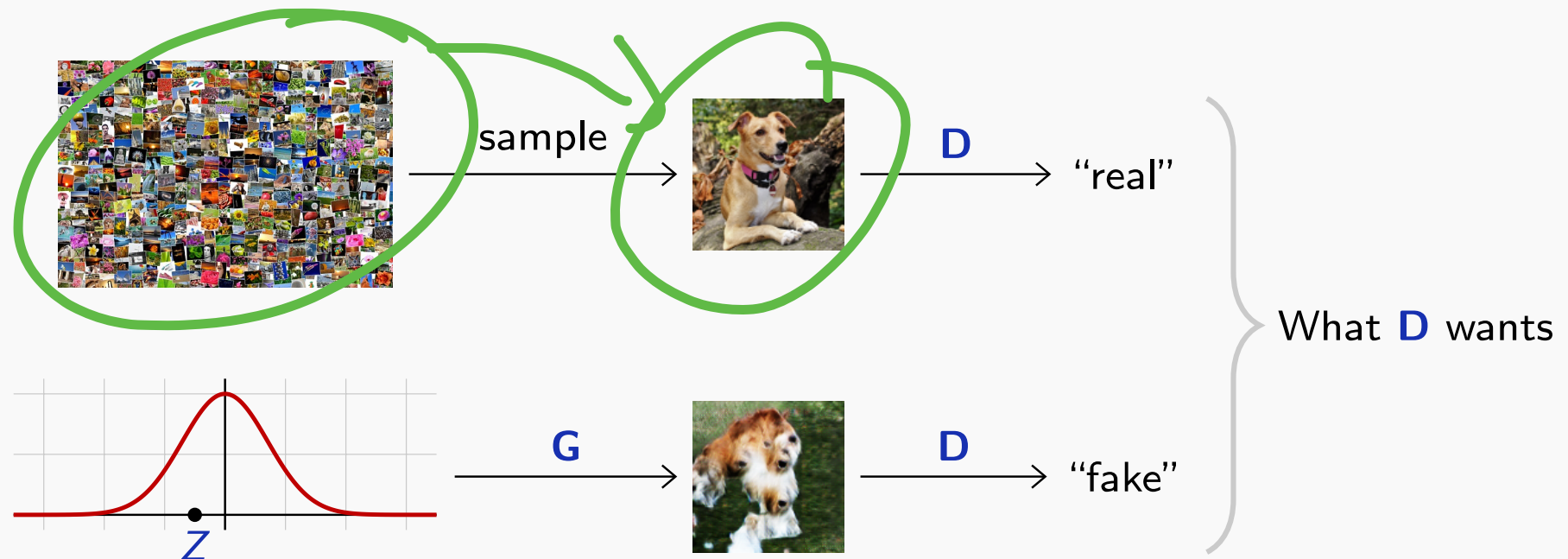
# Generative Adversarial Networks

---

A different approach to learn high-dimensional generative models are Generative Adversarial Networks (Goodfellow et al., 2014)

The idea behind GANs is to train two networks simultaneously:

- A discriminator  $\mathbf{D}$  to classify samples as *real* or *fake*
- A generator  $\mathbf{G}$  to create samples that fool  $\mathbf{D}$



It is adversarial since the two networks have antagonistic objectives.

A bit more formally, let  $\mathcal{X}$  be the signal space and  $D$  the dimension of the latent space

- The generator is trained so that [ideally] if it gets a random normal-distributed  $Z$  as input, it produces a sample following the data distribution as output

$$\mathbf{G} : \mathbb{R}^D \rightarrow \mathcal{X}$$

- The discriminator is trained so that if it gets a sample as input, it predicts if it is genuine.

$$\mathbf{D} : \mathcal{X} \rightarrow [0, 1]$$

If  $\mathbf{G}$  is fixed, to train  $\mathbf{D}$  given a set of real examples

$$x_n \sim \mu, n = 1, \dots, N$$

we create

$$z_n \sim \mathcal{N}(0, 1), 1 \leq n \leq N$$

build a binary data-set

$$\mathcal{D} = \left\{ \underbrace{(x_1, 1), \dots, (x_N, 1)}_{\text{real samples} \sim p_X}, \underbrace{(\mathbf{G}(z_1), 0), \dots, (\mathbf{G}(z_N), 0)}_{\text{fake samples} \sim p_G} \right\}$$

and minimize the binary cross-entropy

$$E(\mathbf{D}) = -\frac{1}{2N} \left( \sum_{n=1}^N \log \mathbf{D}(x_n) + \sum_{n=1}^N \log(1 - \mathbf{D}(\mathbf{G}(z_n))) \right)$$
$$\approx -\frac{1}{2} (\hat{\mathbb{E}}_{X \sim p_X} [\log \mathbf{D}(X)] + \hat{\mathbb{E}}_{X \sim p_G} [\log(1 - \mathbf{D}(X))])$$

where  $p_X$  is the true distribution of the data, and  $p_G$  is the distribution of  $\mathbf{G}(Z)$  with  $Z \sim \mathcal{N}(0, 1)$

CRF  $\rightarrow$  F&S  
gradient  
cliché  
↙

The situation is slightly more complicated since we also want to optimize  $\mathbf{G}$  to maximize  $\mathbf{D}$ 's loss.

Goodfellow et al. (2014) provide an analysis of the resulting equilibrium of that strategy.

$$E(\mathbf{D}) = \left( \mathbb{E}_{X \sim \mu} [\log \mathbf{D}(X)] + \mathbb{E}_{X \sim \mu_G} [\log(1 - \mathbf{D}(X))] \right)$$

is high if  $\mathbf{D}$  is doing a good job (low cross entropy), and low if  $\mathbf{G}$  fools  $\mathbf{D}$

Thus, the theoretically best  $\mathbf{G}^*$  that fools any  $\mathbf{D}$  is

$$\mathbf{G}^* = \operatorname{argmin}_{\mathbf{G}} \max_{\mathbf{D}} V(\mathbf{D}, \mathbf{G})$$

Similarly, the optimal discriminator for a given generator  $\mathbf{G}$  is

$$\mathbf{D}_{\mathbf{G}}^* = \operatorname{argmax}_{\mathbf{D}} V(\mathbf{D}, \mathbf{G})$$

The joint objective becomes

$$\mathbf{G}^* = \operatorname{argmin}_{\mathbf{G}} V(\mathbf{D}_{\mathbf{G}}^*, \mathbf{G})$$

We have

$$\begin{aligned} V(\mathbf{D}, \mathbf{G}) &= \mathbb{E}_{X \sim \mu} [\log \mathbf{D}(X)] + \mathbb{E}_{X \sim \mu_{\mathbf{G}}} [\log(1 - \mathbf{D}(X))] \\ &= \int_X \underbrace{\mu(x)}_{\text{pink}} \underbrace{\log \mathbf{D}(x)}_{\text{green}} + \underbrace{\mu_{\mathbf{G}}(x)}_{\text{pink}} \underbrace{\log(1 - \mathbf{D}(x))}_{\text{green}} dx \end{aligned}$$

Since

$$\operatorname{argmax}_d \mu(x) \log d + \underbrace{\mu_{\mathbf{G}}(x)}_{\text{green}} \log(1 - d) = \frac{\mu(x)}{\mu(x) + \mu_{\mathbf{G}}(x)}$$

and

$$\mathbf{D}_{\mathbf{G}}^* = \operatorname{argmax}_{\mathbf{D}} V(\mathbf{D}, \mathbf{G})$$

we get (if there is no additional regularization!)

$$\forall x, \mathbf{D}_{\mathbf{G}}^*(x) = \frac{\mu(x)}{\mu(x) + \mu_{\mathbf{G}}(x)}$$



$$d = \underset{?}{d}^* \quad \text{argmax}_{D(x)} \quad \underbrace{\mu(x)}_A \underbrace{\log D(x)}_d + \underbrace{\mu_G(x)}_B \underbrace{\log(1-D(x))}_d$$

$$\Rightarrow \frac{\partial}{\partial d} A \cdot \log(d) + B \cdot \log(1-d)$$

$$= A \frac{1}{d} + B \frac{1}{1-d} (-1)$$

$$= A \frac{1}{d} - B \frac{1}{1-d} \quad | \cdot d$$

$$= A - B \frac{d}{1-d}$$

$$= A - B \frac{1}{\frac{1}{d} - 1} \stackrel{!}{=} 0$$

$$\Rightarrow A = + B \frac{1}{\frac{1}{d} - 1} \quad | \cdot (\frac{1}{d} - 1)$$

$$\Rightarrow \frac{1}{\lambda} A - A = +B$$

$$\frac{1}{\lambda} A = A + B$$

$$\frac{1}{\lambda} = \frac{A+B}{A}$$

$$\lambda^* = \frac{A}{A+B}$$

$$\Rightarrow Y(x) = \frac{\mu_x(x)}{\mu_x(x) + \mu_G(x)}$$

So, since

$$\forall x, \mathbf{D}_{\mathbf{G}}^*(x) = \frac{\mu(x)}{\mu(x) + \mu_{\mathbf{G}}(x)}$$

we have

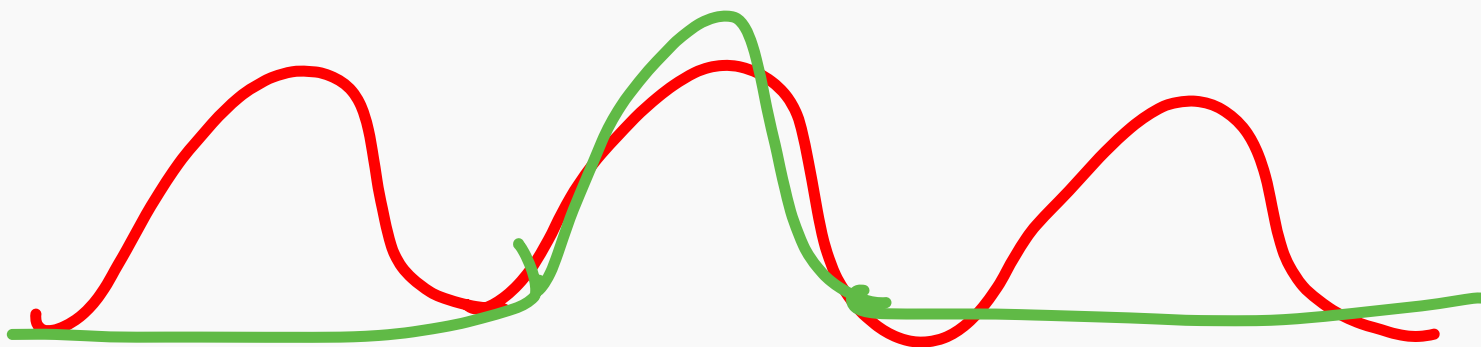
$$\begin{aligned} V(\mathbf{D}_{\mathbf{G}}^*, \mathbf{G}) &= \mathbb{E}_{X \sim \mu} [\log \mathbf{D}_{\mathbf{G}}^*(X)] + \mathbb{E}_{X \sim \mu_{\mathbf{G}}} [\log(1 - \mathbf{D}_{\mathbf{G}}^*(X))] \\ &= \mathbb{E}_{X \sim \mu} \left[ \log \frac{\mu(X)}{\mu(X) + \mu_{\mathbf{G}}(X)} \right] + \mathbb{E}_{X \sim \mu_{\mathbf{G}}} \left[ \log \frac{\mu_{\mathbf{G}}(X)}{\mu(X) + \mu_{\mathbf{G}}(X)} \right] \\ &= KL \left( \mu \parallel \frac{\mu + \mu_{\mathbf{G}}}{2} \right) + KL \left( \mu_{\mathbf{G}} \parallel \frac{\mu + \mu_{\mathbf{G}}}{2} \right) - \log 4 \\ &= JSD(\mu, \mu_{\mathbf{G}}) - \log 4 \end{aligned}$$

where  $JSD$  is the Jensen-Shannon Divergence, a dissimilarity measure between distributions.

Training a standard GAN often results in two pathological behaviors:

- Oscillations without convergence: we have no guarantee that the loss will actually decrease
- The infamous *mode collapse*, when  $G$  models very well a small part of the data, concentrating on a few modes

Additionally, performance is hard to assess, often being a “beauty contest”.

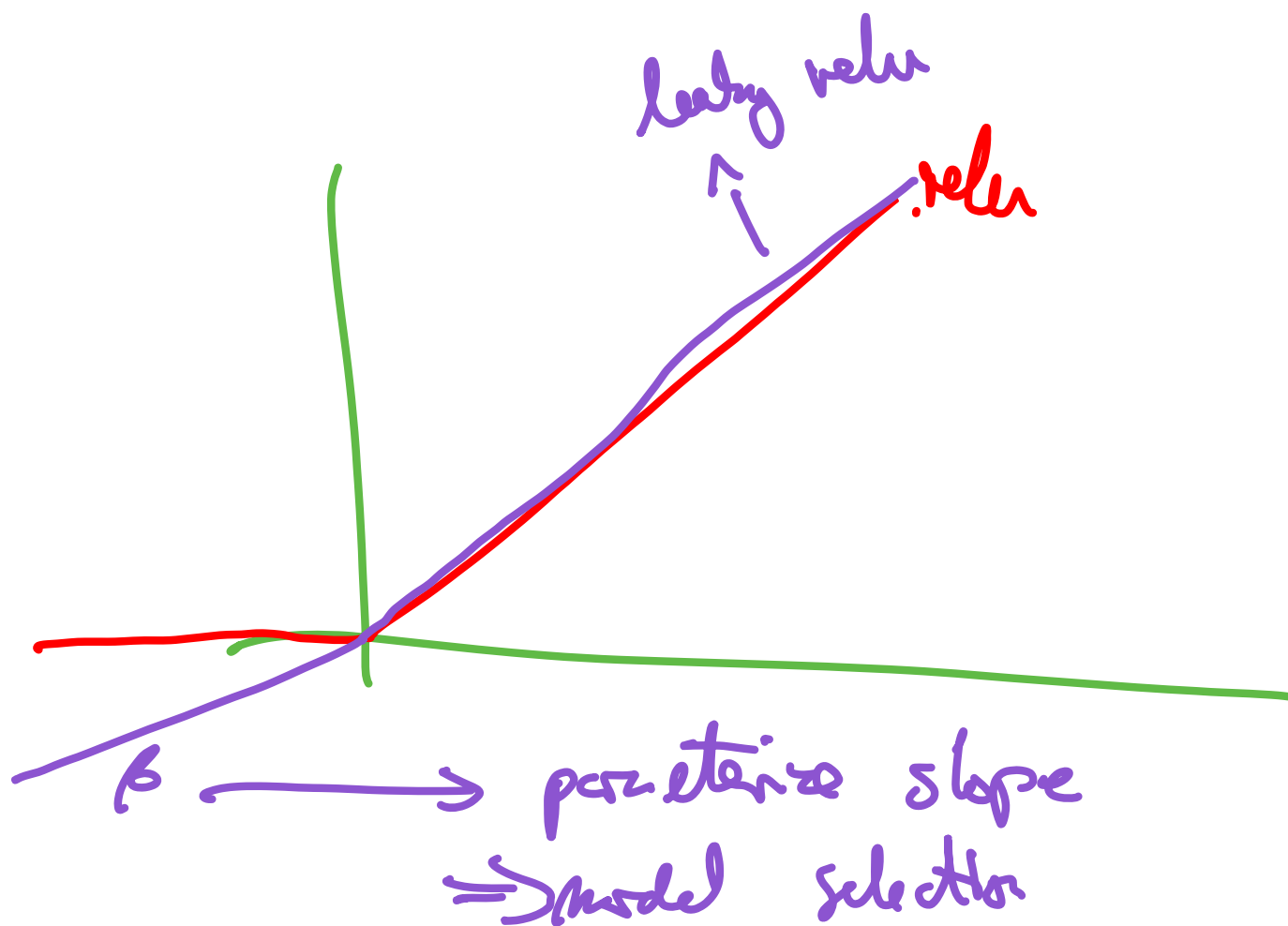


*We also encountered difficulties attempting to scale GANs using CNN architectures commonly used in the supervised literature. However, after extensive model exploration we identified a family of architectures that resulted in stable training across a range of datasets and allowed for training higher resolution and deeper generative models.*

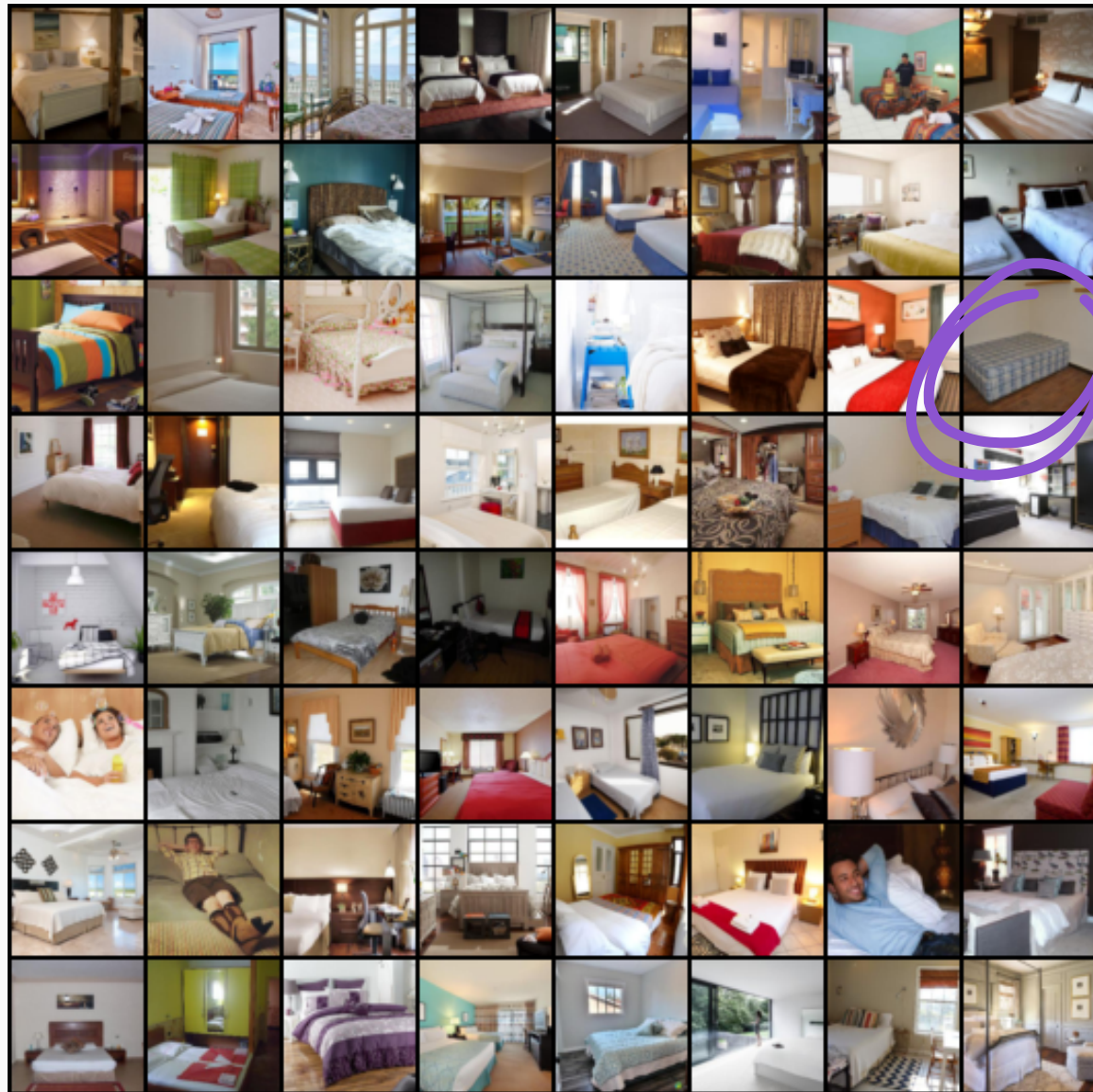
Radford et al. (2015)

Radford et al. (2015) adopted the following “rules”:*of - this,*

- Replace pooling with
  - Strided convolutions in **D**
  - Transposed convolutions in **G**
- Use BatchNorm in both **D** and **G**
- No fully connected layers
- Output of **G** uses tanh, then ReLU elsewhere
- LeakyReLUs in **D**



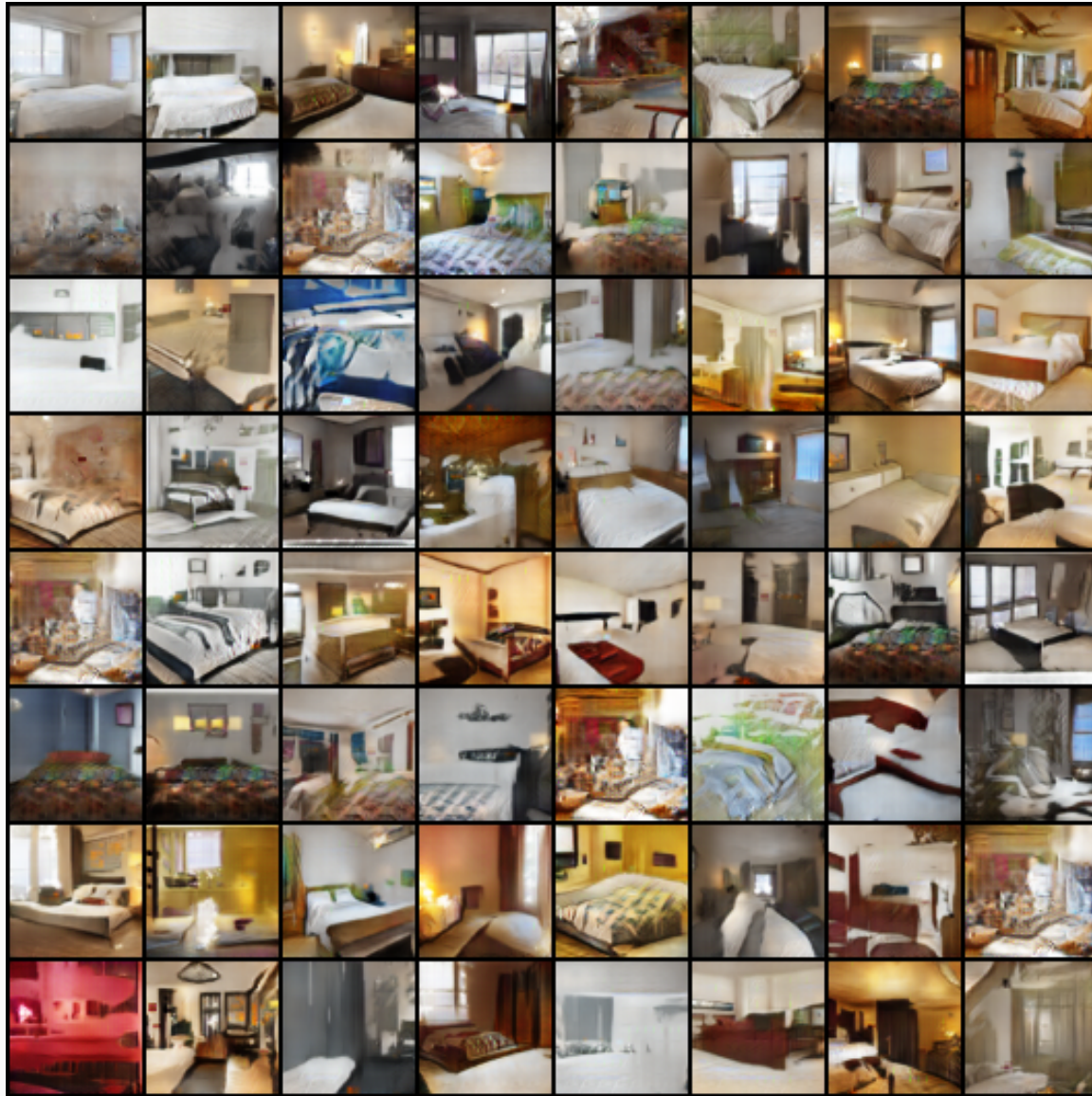
# “bedroom” class of the LSUN data set



Real samples



# “bedroom” class of the LSUN data set



Fake samples (20 epochs)

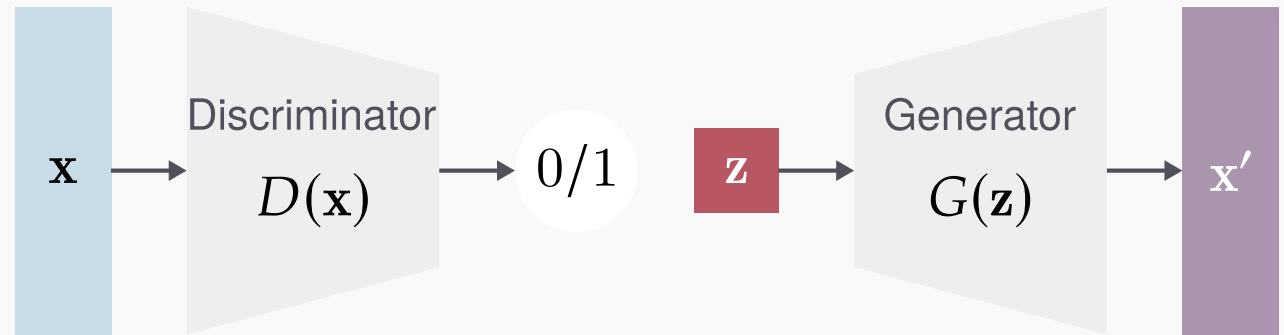
Replace discriminators by statistical test

$G \rightarrow$  data set  $\rightarrow D(\text{data set})$   
(real data set)  
 $\downarrow$   
sig. diff  
or not

# Deep Generative Modeling

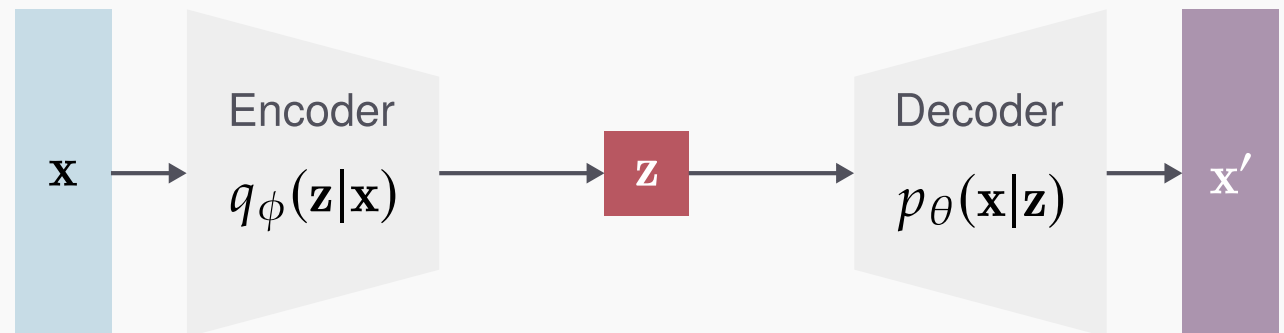
## GAN

Minimax  
classification  
loss



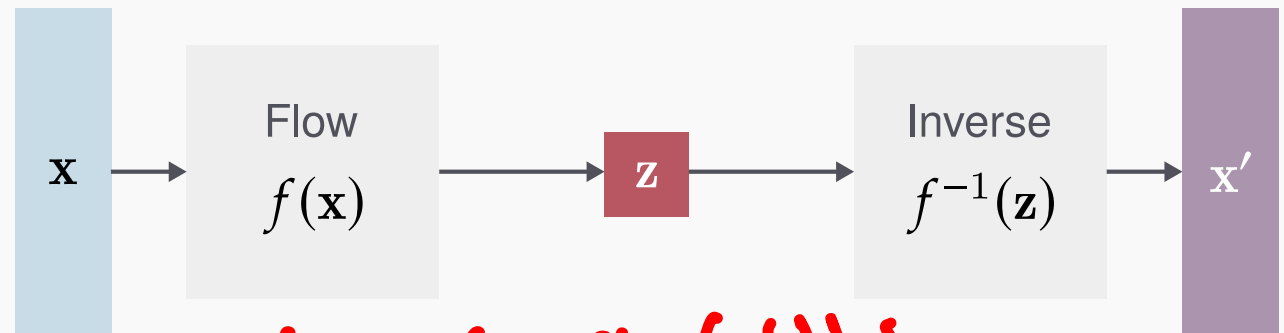
## VAE

Maximize ELBO



## Flow

Minimize NLL



$$f_n \circ f_{n-1} \circ \dots \circ f_1(\mathbf{x})$$

# Recurrent Neural Networks

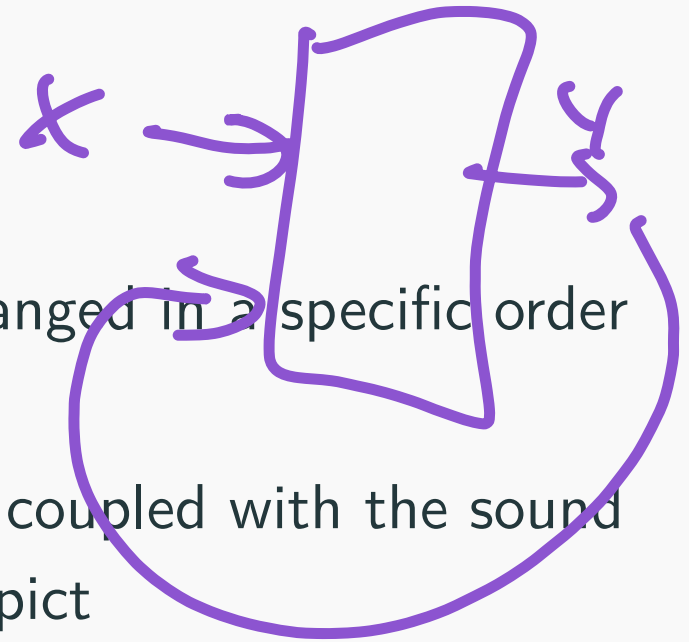
---

# TD- Gammon

Temporal or sequential data are different from simple vector data because the order in which they are arranged is part of the underlying phenomena they represent.

Examples include

- text, where words and letters are arranged in a specific order to convey a certain meaning
- video, where the sequence of images coupled with the sound signal determine the content they depict

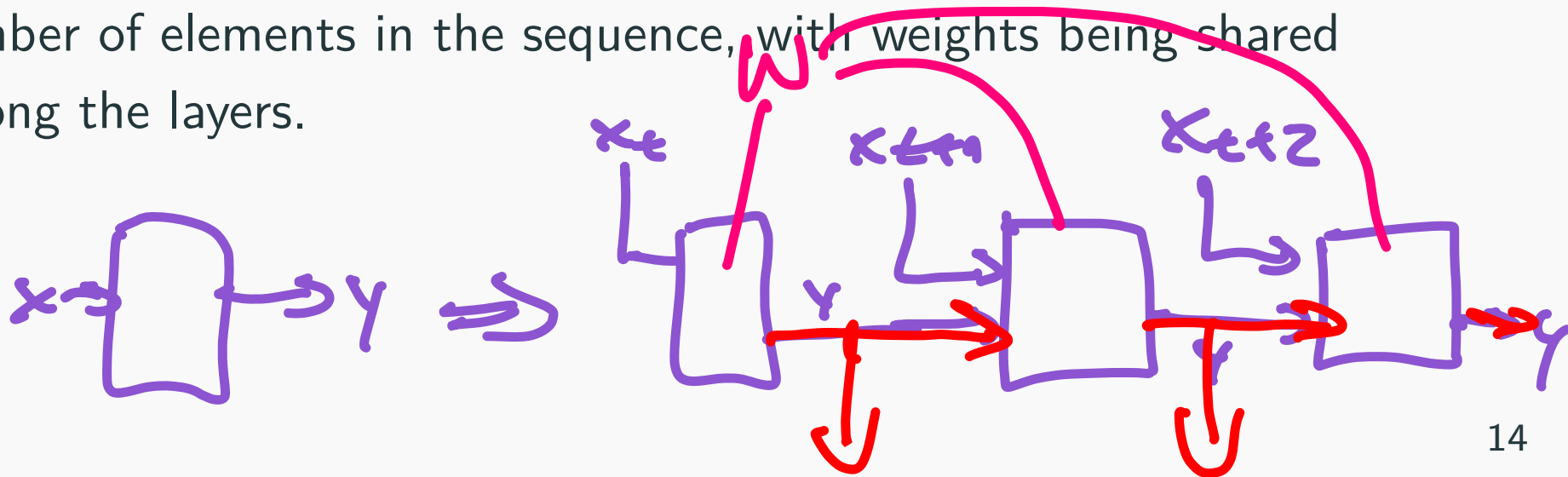


*Recurrent neural networks* are a class of neural networks designed specifically for handling sequential data.

Neural network architectures usually follow a specific “forward” direction: activations are always further from the inputs as it is transformed.

In recurrent architectures, however, this is not the case. The outputs of some neurons can be connected to the inputs of neurons that came before.

Such recurrent models can be seen as a MLP as deep as the number of elements in the sequence, with weights being shared among the layers.

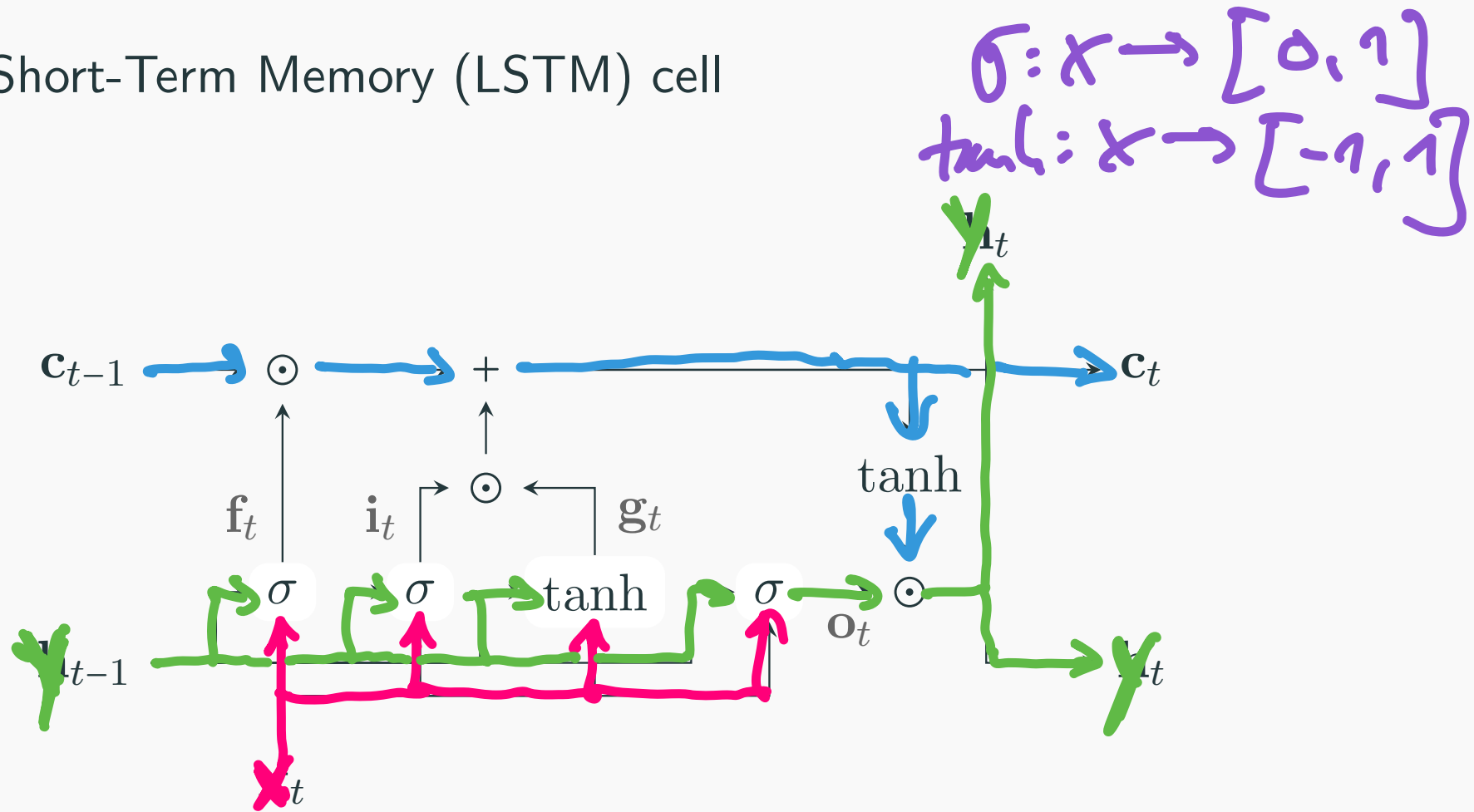


Similarly to regular deep architectures, connecting neurons in a recurrent fashion can easily cause vanishing gradients.

In fact, some of the first analyses on the vanishing gradients problem were conducted while studying recurrent architectures.

There are two main recurrent architectures designed to work with very long sequences: the *long short-term memory* (LSTM) and its simplification the *gated recurrent unit* (GRU).

# Long Short-Term Memory (LSTM) cell



The value of the cell state  $\mathbf{c}_{t-1}$  is increased by  $\mathbf{i}_t \odot \mathbf{g}_t$ . The forget gate  $\mathbf{f}_t$  allows the LSTM to easily reset the value of the cell. The output is typically  $\mathbf{h}_t$ , and  $\mathbf{c}_t$  is used internally.



The full updates are calculated as:

$$\mathbf{i}_t = \sigma \left( \mathbf{z}^{(ii)}(\mathbf{x}_t) + \mathbf{z}^{(hi)}(\mathbf{h}_{t-1}) \right) \quad (\text{input gate})$$

$$\mathbf{f}_t = \sigma \left( \mathbf{z}^{(if)}(\mathbf{x}_t) + \mathbf{z}^{(hf)}(\mathbf{h}_{t-1}) \right) \quad (\text{forget gate})$$

$$\mathbf{g}_t = \tanh \left( \mathbf{z}^{(ig)}(\mathbf{x}_t) + \mathbf{z}^{(hg)}(\mathbf{h}_{t-1}) \right) \quad (\text{full input gate})$$

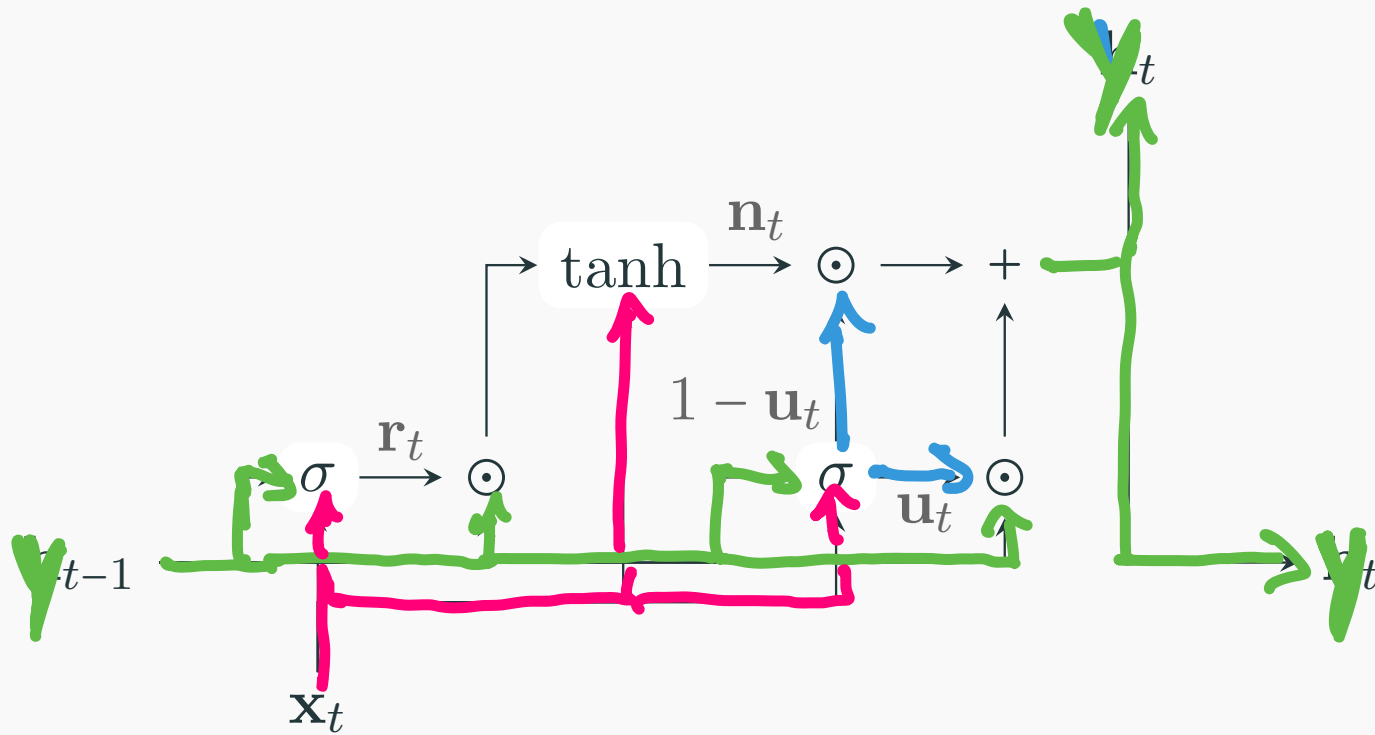
$$\mathbf{o}_t = \sigma \left( \mathbf{z}^{(io)}(\mathbf{x}_t) + \mathbf{z}^{(ho)}(\mathbf{h}_{t-1}) \right) \quad (\text{output gate})$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t \quad (\text{cell state})$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \quad (\text{output})$$

where  $\sigma$  is a sigmoid function and  $\mathbf{z}^{(\cdot)}(\cdot)$  are linear transformations.

# Gated Recurrent Unit (GRU) cell



$$\mathbf{r}_t = \sigma (\mathbf{z}_{ir}(\mathbf{x}_t) + \mathbf{z}_r(\mathbf{h}_{t-1})) \quad \text{(reset gate)}$$

$$\mathbf{u}_t = \sigma (\mathbf{z}_{iz}(\mathbf{x}_t) + \mathbf{z}_{hz}(\mathbf{h}_{t-1})) \quad \text{(update gate)}$$

$$\mathbf{n}_t = \tanh (\mathbf{z}_{in}(\mathbf{x}_t) + \mathbf{r}_t \odot \mathbf{z}_{hn}(\mathbf{h}_{t-1}))$$

$$\mathbf{h}_t = (1 - \mathbf{u}_t) \odot \mathbf{n}_t + \mathbf{u}_t \odot \mathbf{h}_{t-1} \quad \text{(output)}$$

# References

I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial networks. CoRR, abs/1406.2661, 2014.

A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. CoRR, abs/1511.06434, 2015.