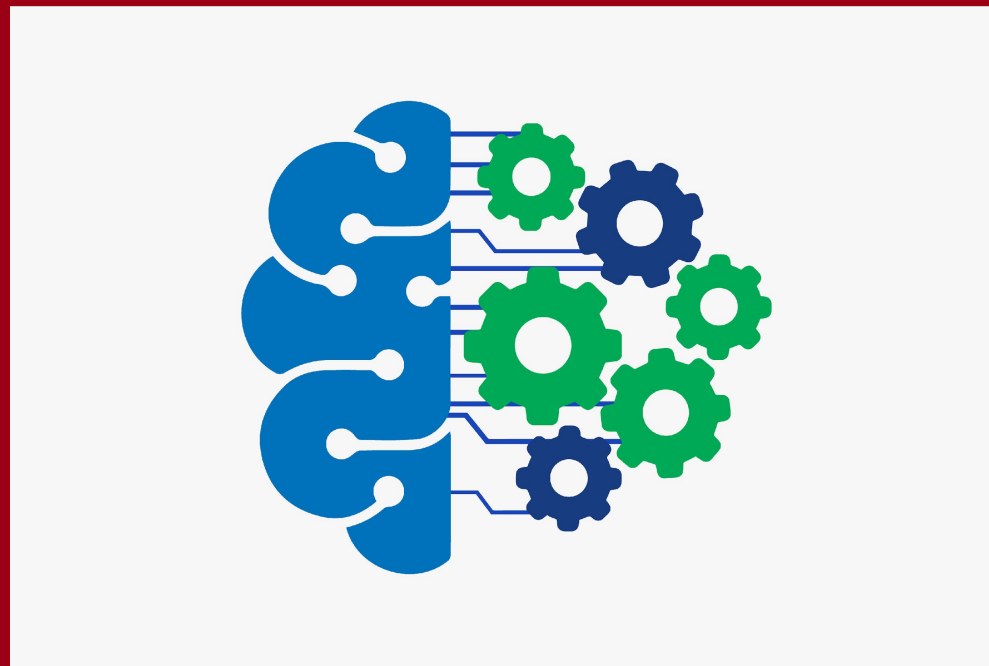




DIPARTIMENTO
DI INGEGNERIA
DELL'INFORMAZIONE



UNIVERSITÀ
DEGLI STUDI
DI PADOVA



Deep Learning: Advanced Approaches

Machine Learning 2022-23

Deep Learning: Advanced Approaches

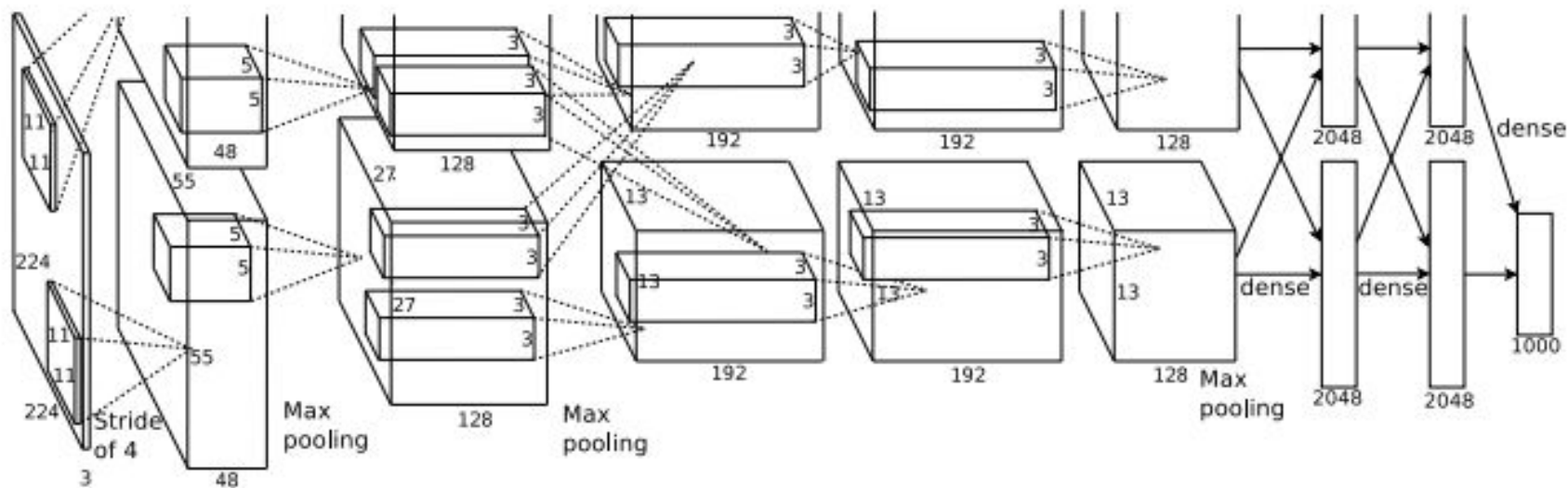
1. *Advanced CNN schemes*: Residual networks, skip connections, auto-encoders
2. *Generative models*: Generative Adversarial Networks (GAN)
3. *Modeling temporal information*: Recurrent Neural Networks (RNN) and Long-Short Term Memory (LSTM) (*not part of the course*)



Advanced CNN Models

- ❑ We'll see some relatively recent advanced architectures
- ❑ Some new concepts will be briefly introduced:
 - Residual Networks
 - Auto-Encoders
 - Skip Connections
 - Transposed Convolutions

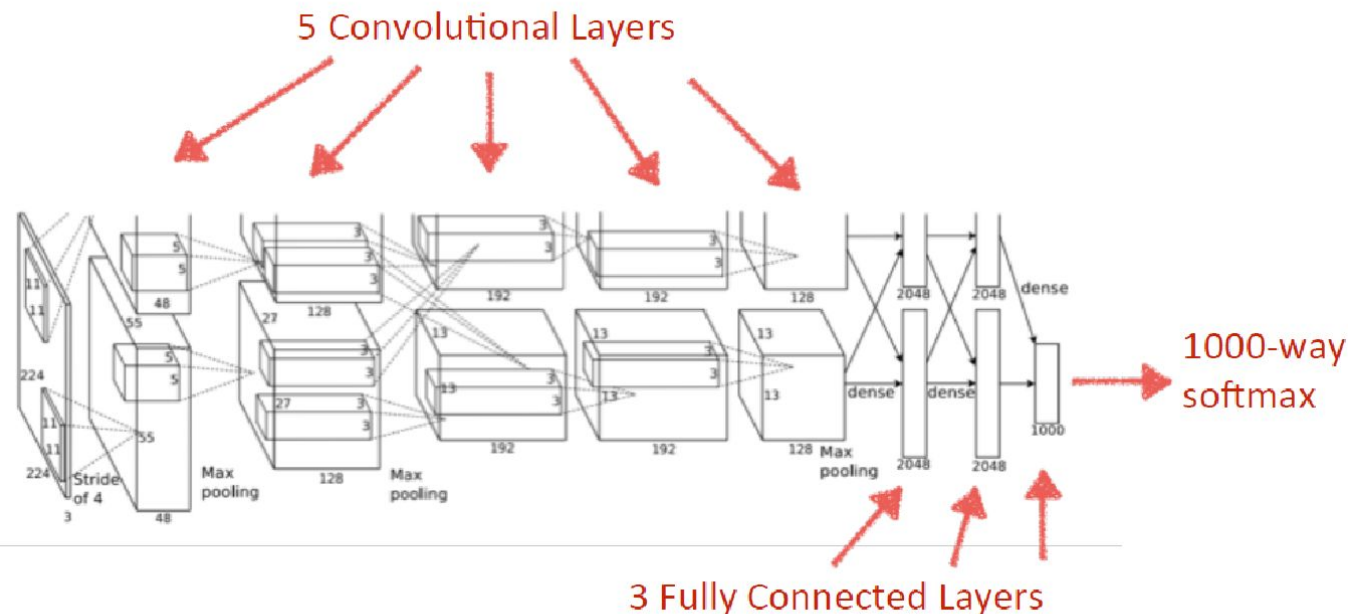
«Historical» Perspective: AlexNet (2012)



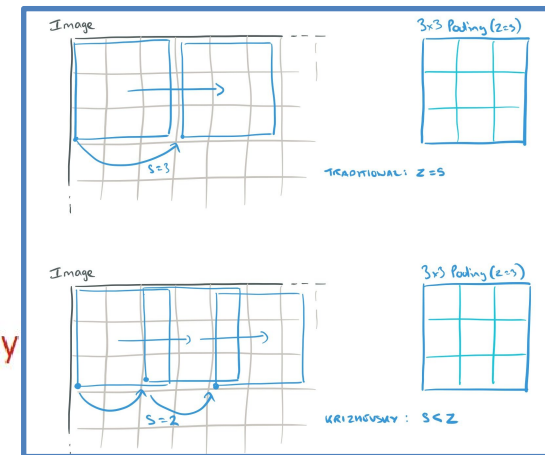
- ❑ **AlexNet** [3]: First Deep Learning approach outperforming “classic” methods (i.e., outperforming SVM or RF)
- ❑ Exploits 11x11, 5x5, 3x3, convolutions, max pooling, dropout, data augmentation, ReLU activations, SGD with momentum
- ❑ Split in 2 pipelines since it was trained with 2 GPUs (for 6 days)
 - According to Nvidia the DGX-2 server released in 2018 can train it in 18 mins!!!
- ❑ Complex but quite “standard” model



AlexNet: the Network

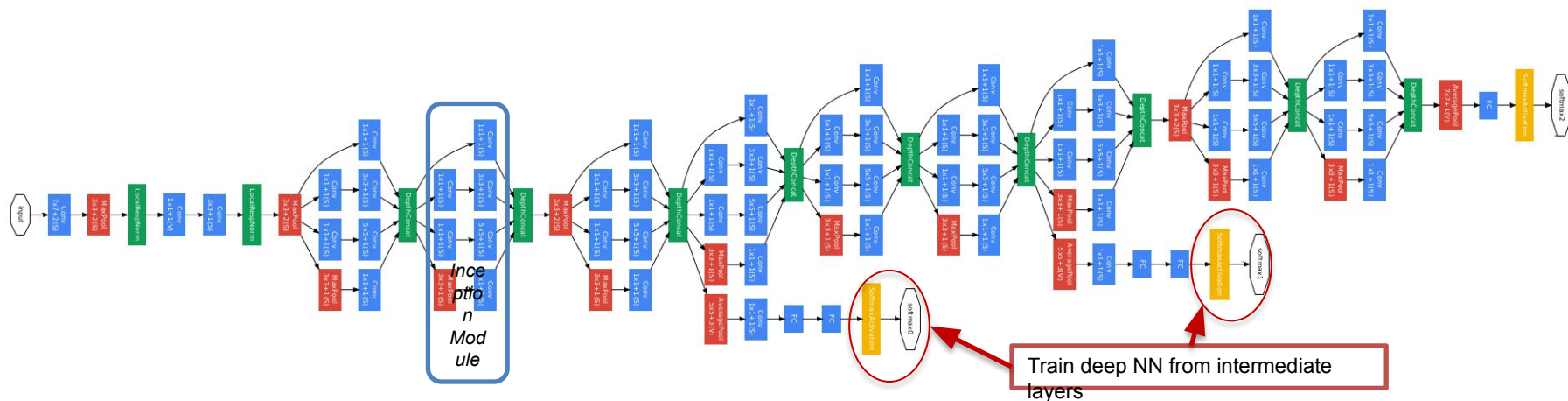


overlapping pooling



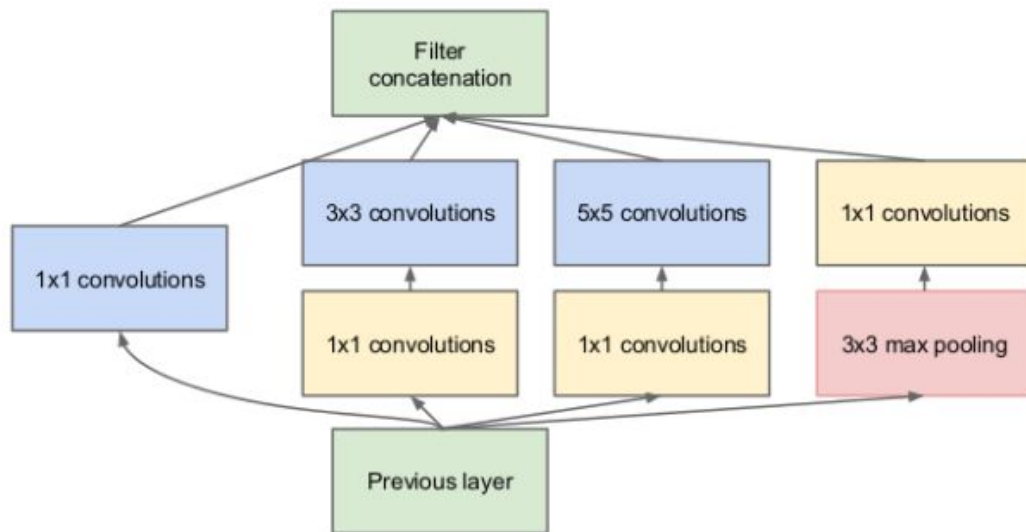
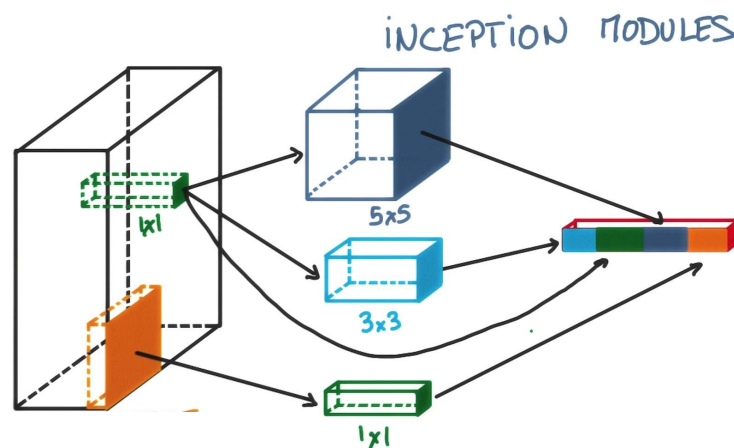
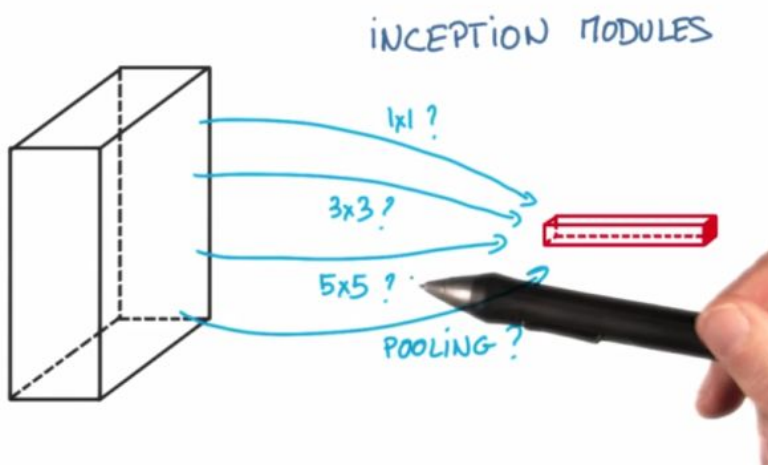
- ❑ 5 convolutional layers, 3 fully connected ones
- ❑ Many feature maps for each layer
- ❑ 650K neurons, 60M parameters
- ❑ Rectified Linear Units (ReLU) activations, **overlapping pooling**, dropout trick
- ❑ Training with randomly extracted 224x224 patches for more data

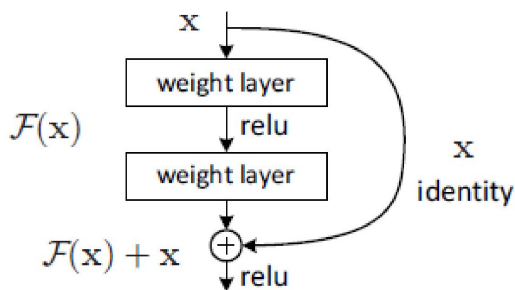
GoogleNet (Inception V1)



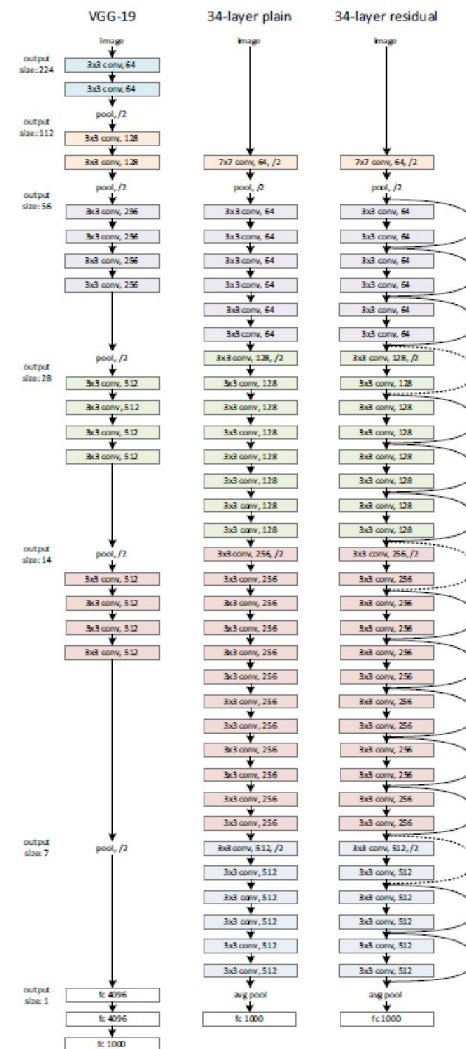
- ❑ Released in 2014, 1st method very close to human level performance
- ❑ Implemented a novel element: *the inception module*
 - This module performs multiple small convolutions with different sizes in parallel
- ❑ The network is a 22 layers deep CNN but reduced the number of parameters from 60M of AlexNet to 4M

The Inception Module

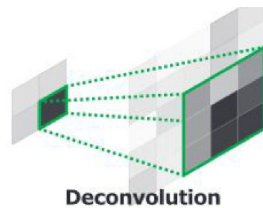
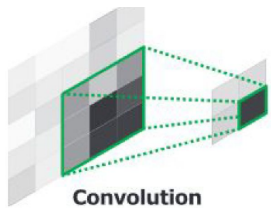
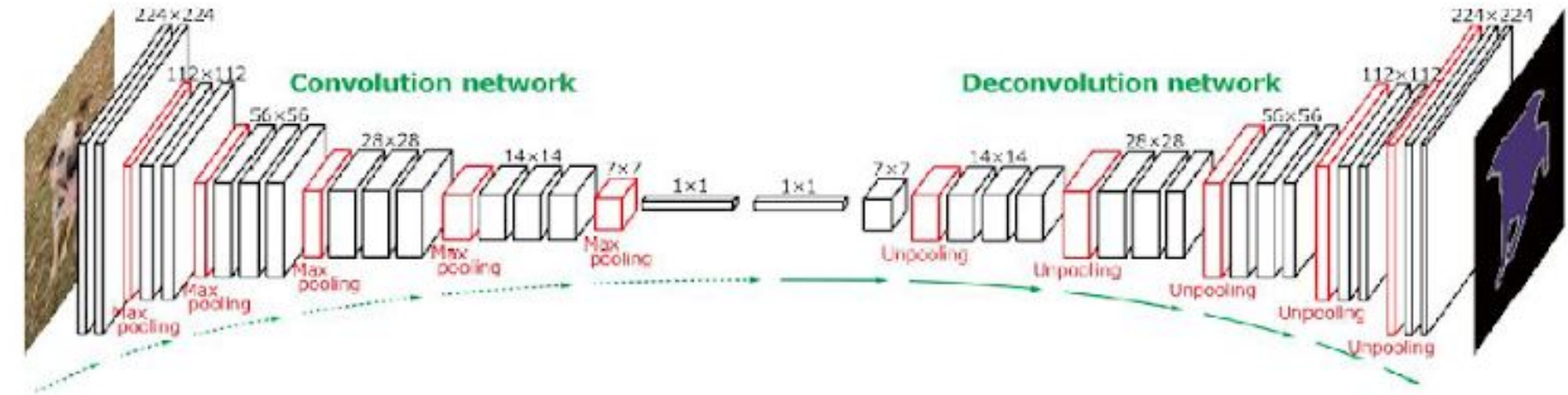




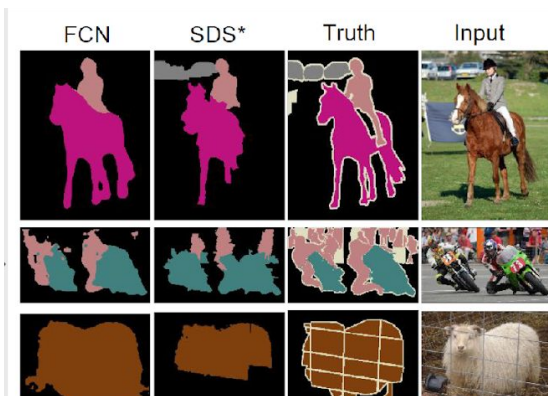
- ❑ Residual Neural Network [4] introduced in 2015 a novel architecture with “*skip connections*”
- ❑ Idea: try to estimate the residual w.r.t the previous estimation instead of the function itself
- ❑ Thanks to this technique they were able to train a NN with 152 layers with reasonable complexity
- ❑ Was able to beat human-level performance on image classification tasks



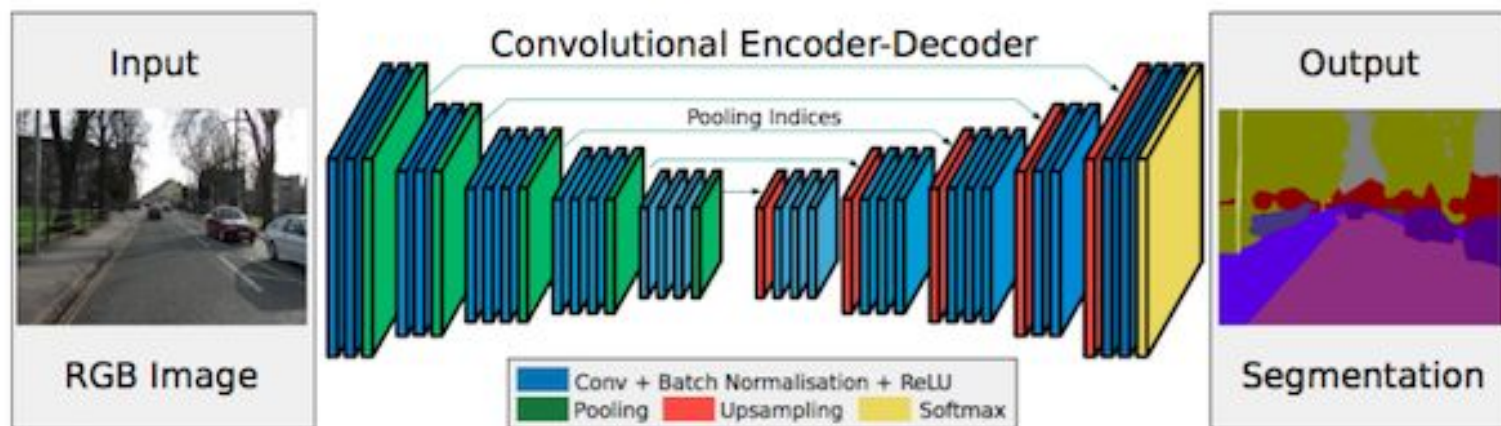
Upsampling : Transposed Convolutions



- ❑ In some applications the output has the same or even larger size than the input (e.g., semantic segmentation, denoising)
- ❑ Convolutional layers connect multiple input activations within a filter window to a single activation
- ❑ Transposed convolutions associate a single input activation with multiple outputs
- ❑ Use transposed convolutions for upsampling

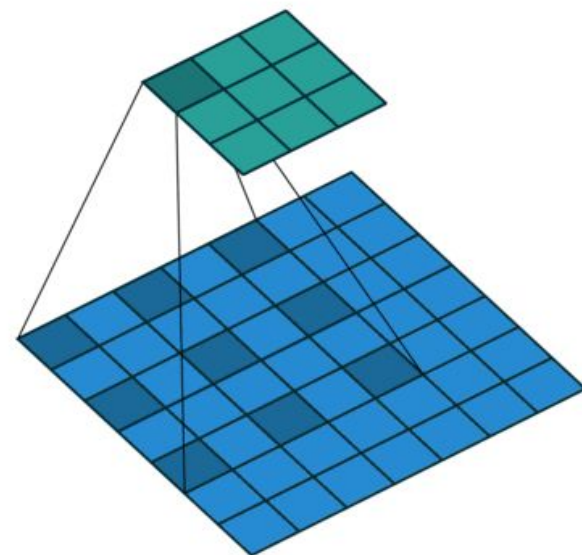


Encoder-Decoder Architectures

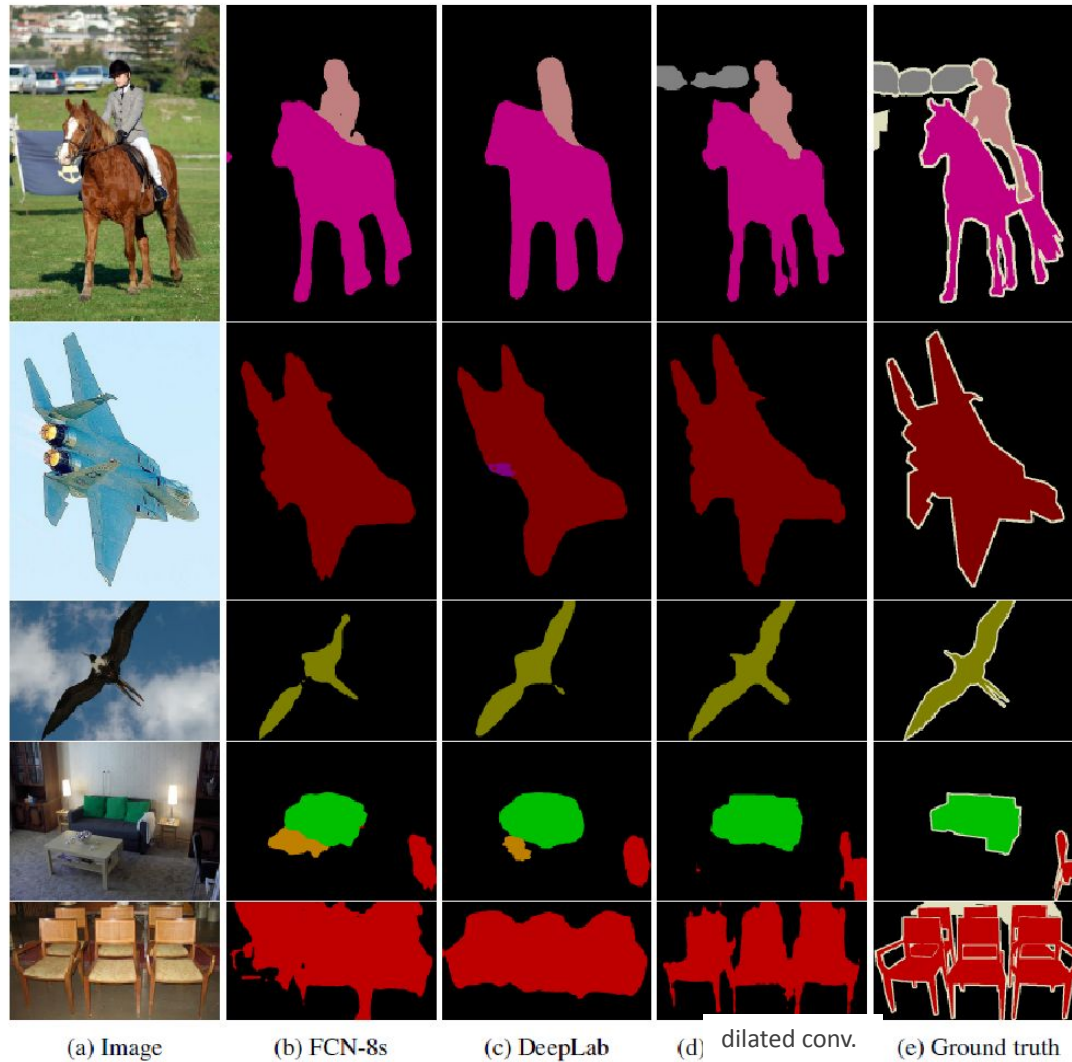


- ❑ The network is made of 2 parts, an **encoder** and a **decoder**
- ❑ A "**compressed**" description of the input data is created at the middle layers by the **encoder**
- ❑ The **decoder** expands it into the final result
- ❑ Maxpooling indices can be transferred to decoder to improve the reconstruction
- ❑ FCN and SegNet are among the first encoder-decoder architectures
 - Fully Convolutional Networks for Semantic Segmentation (2014)
 - A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation (*SegNet*) (2015)

Dilated Convolutions



- ❑ Large convolutions have a wide receptive field but requires a lot of parameters
- ❑ Use dilated (*atrous*) convolutions, to increase the field of view without increasing the spatial dimensions
- ❑ The convolution works on samples spaced apart with a regular step instead of over each single sample in the window.





GANs

- ❑ **Generative**
 - Learn a generative model
- ❑ **Adversarial**
 - Trained in an adversarial setting
- ❑ **Networks**
 - Use Deep Neural Networks

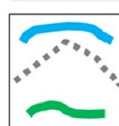
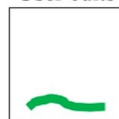
Introduced by Goodfellow et al in 2014 [2]

Generative Models

Which one is Computer generated?



User edits



Generated images



Why Generative Models?

- ❑ We have only seen discriminative models so far...

- Given a vector \mathbf{x} , predict a label \mathbf{y} → The model estimates $P(\mathbf{y}|\mathbf{x})$

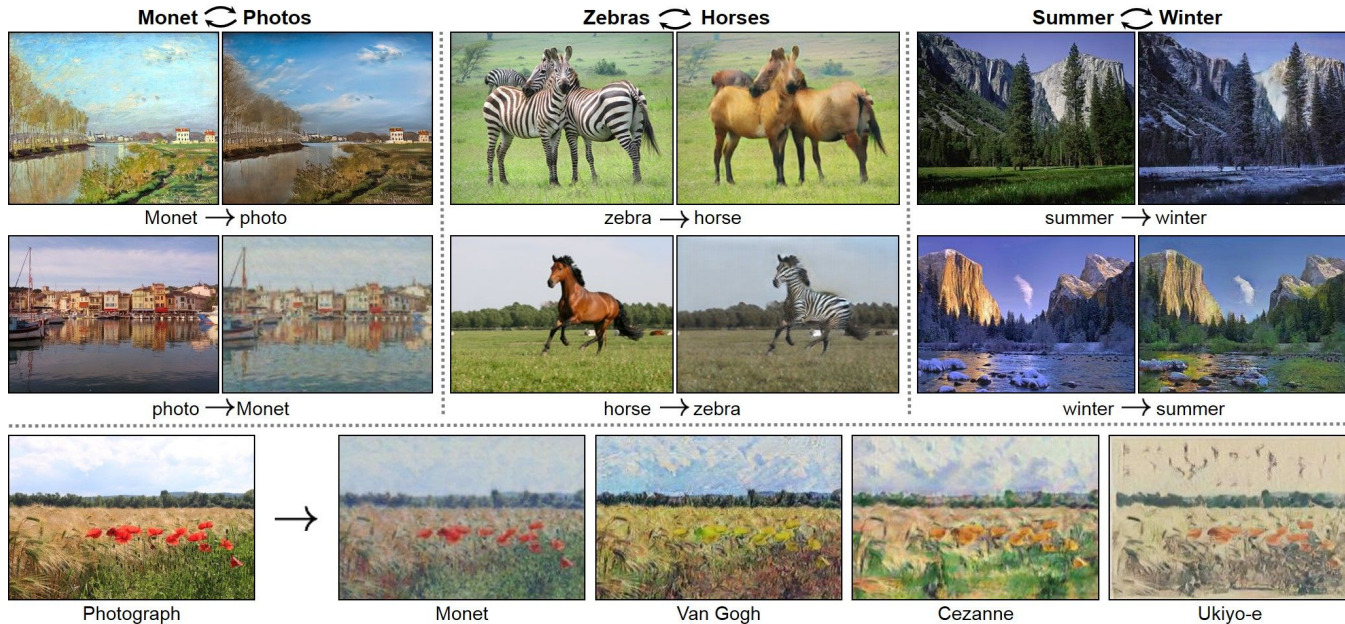
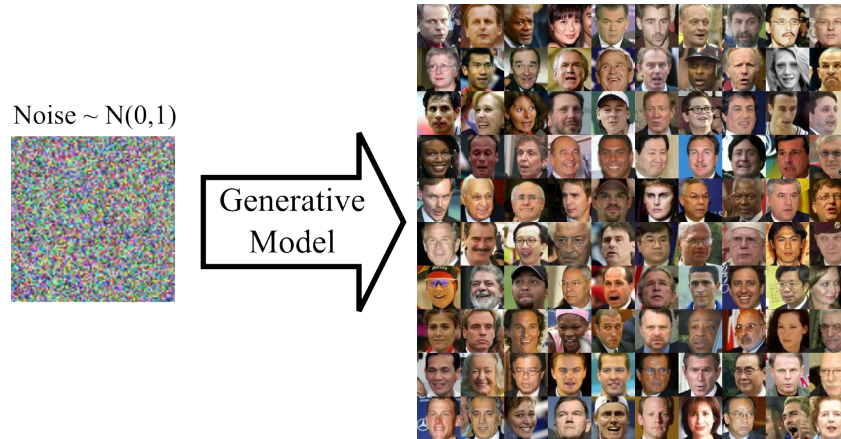
- ❑ Discriminative models can't model $P(\mathbf{x})$, i.e. the probability of seeing a certain data sample

- Thus, can't sample from $P(\mathbf{x})$, i.e. *can't generate new data samples*

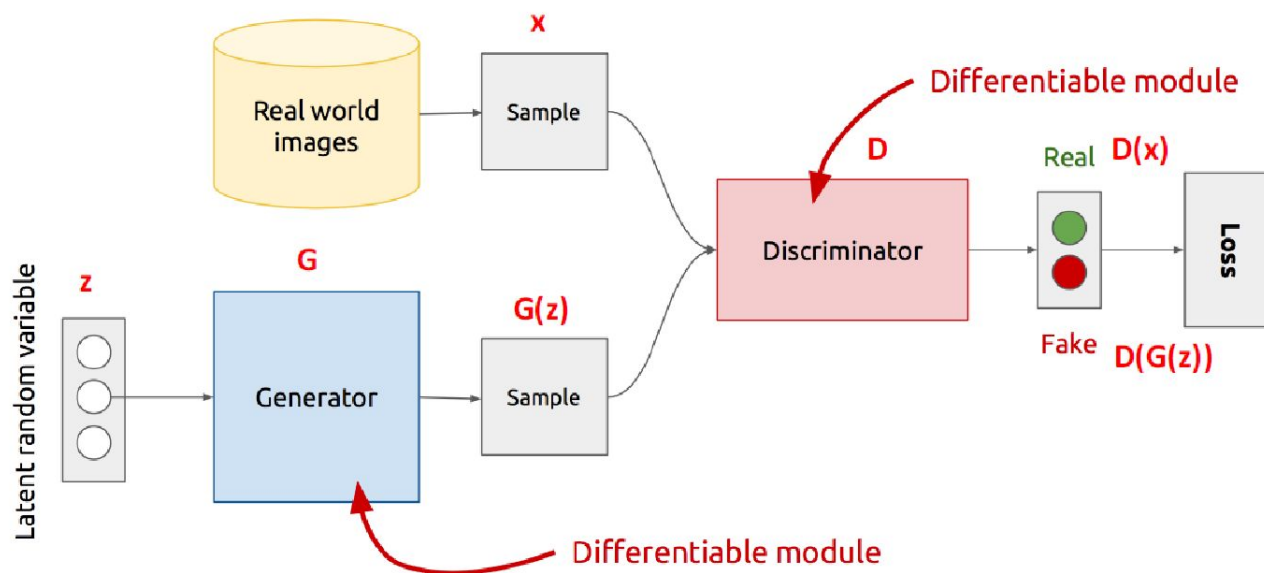
- ❑ Generative models can model $P(\mathbf{x})$

- Can *generate new data samples*

Generative Models: Examples



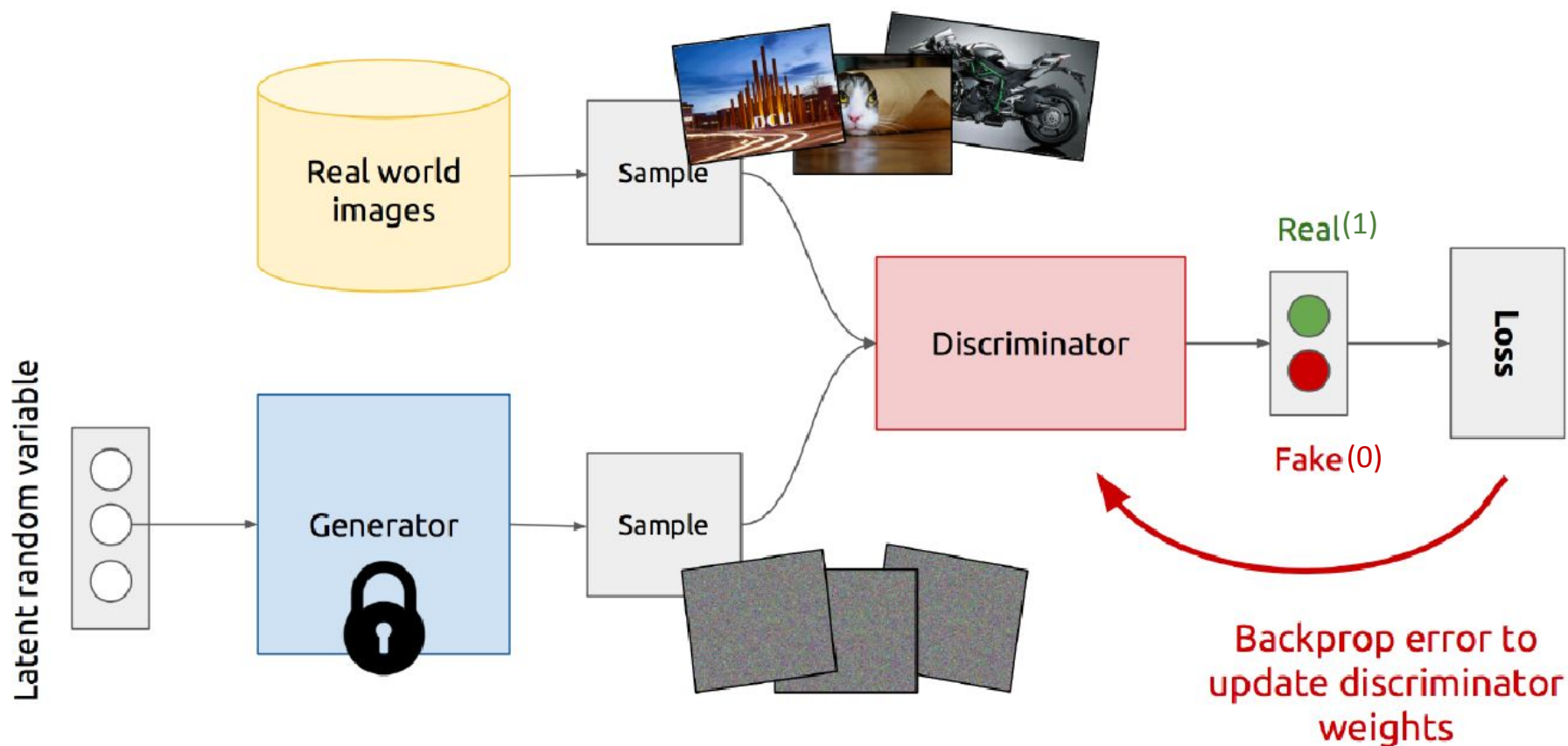
Generative Adversarial Networks (GAN)



A GAN is composed of two sub-networks:

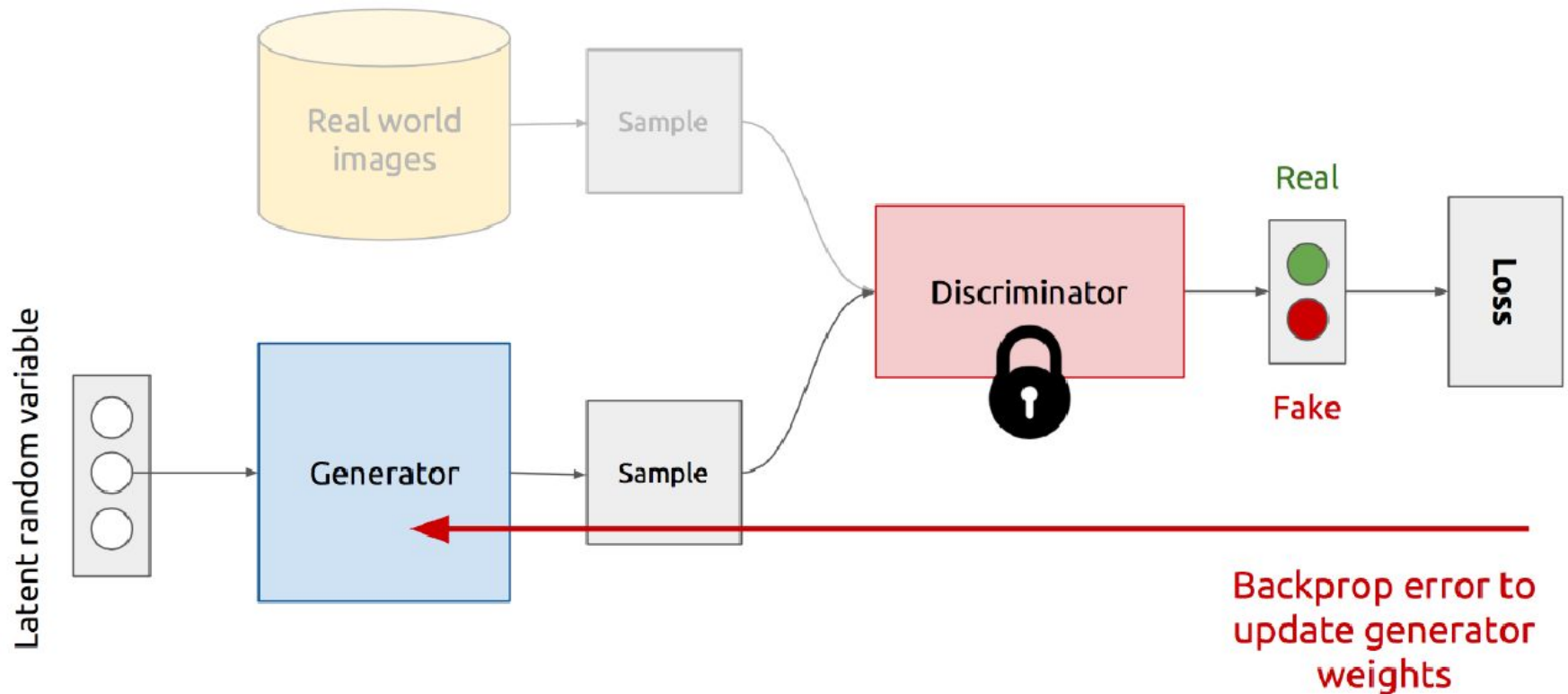
1. **Generator (G):** generate fake samples, tries to fool the Discriminator
 2. **Discriminator (D):** tries to distinguish between real and fake samples
- ❑ Train them against each other (in practice we alternate between training **Generator** and **Discriminator**)
 - ❑ Repeat this and we get better **Generator** and **Discriminator**

Discriminator Training



Target: *Minimize* discriminator loss

Generator Training



Target: *Maximize* discriminator loss



Loss Function

$$\min_G \max_D V(D, G)$$

- It is formulated as a **minimax game**, where:
 - The Discriminator is trying to maximize its reward $V(\mathbf{D}, \mathbf{G})$ (or minimize its loss)
 - The Generator is trying to minimize Discriminator's reward (or maximize its loss)

$$V(D, G) = \underbrace{\mathbb{E}_{x \sim p(x)} [\log D(x)]}_{\text{true samples}} + \underbrace{\mathbb{E}_{z \sim q(z)} [\log(1 - D(G(z)))]}_{\text{fake (generated) samples}}$$

- The Nash equilibrium of this particular game is achieved at:
 - $P_{data}(x) = P_{gen}(x) \quad \forall x$
 - $D(x) = \frac{1}{2} \quad \forall x$

Training Algorithm

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(x^{(i)}) + \log (1 - D(G(z^{(i)}))) \right].$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)}))).$$

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

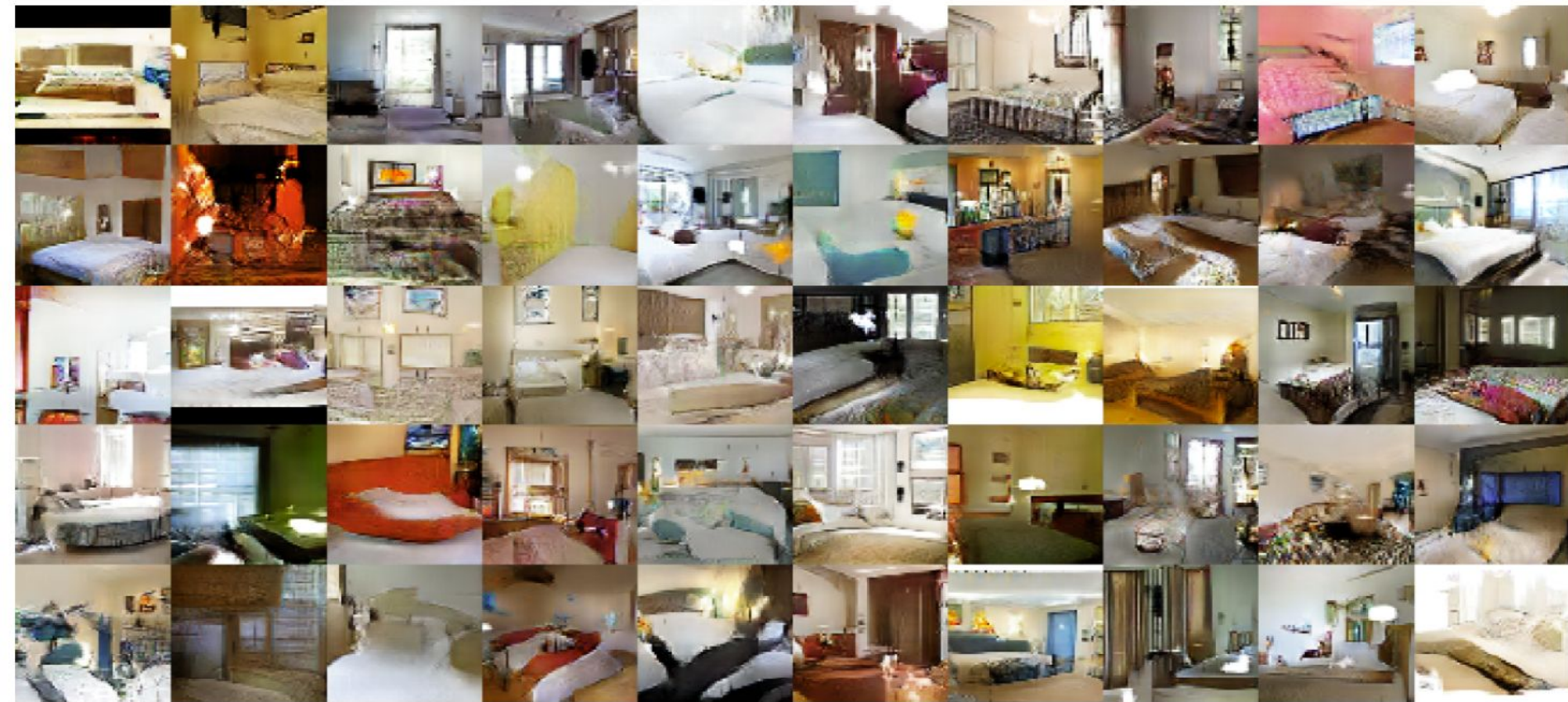
Discriminator
updates

Generator
updates



DIPARTIMENTO
DI INGEGNERIA
DELL'INFORMAZIONE

Examples (1): Generated Images



Examples (2)



Condition image

Target pose sequence

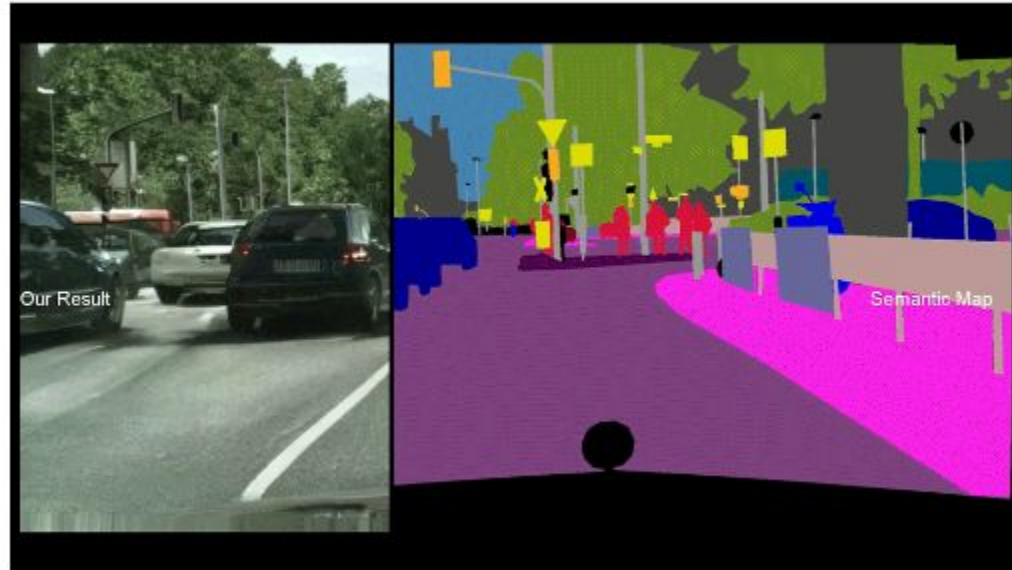
Refined results



(b) Handbag images (input) & **Generated** shoe images (output)



Examples (3)





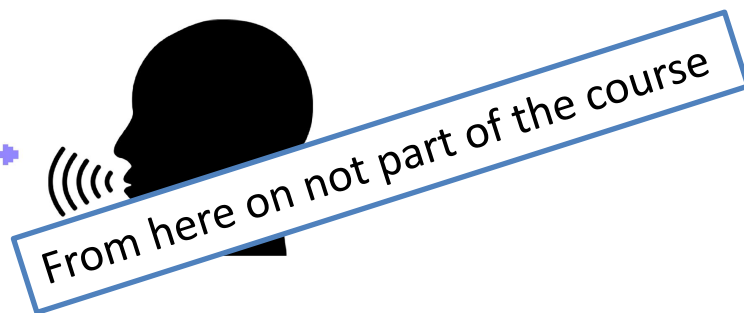
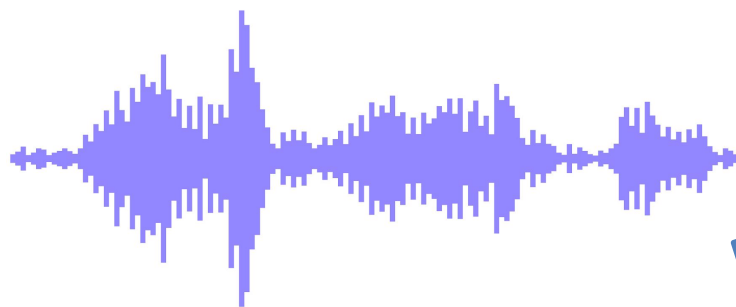
Many Other Approaches....

- ❑ This was just a quick overview of some relatively recent results
 - For ICT students more approaches will be presented in computer vision, neural networks and deep learning and many other courses....
- ❑ Huge amount of resources is currently spent on Deep Learning research
- ❑ Many other schemes exist, and every month there is a new one outperforming previous results



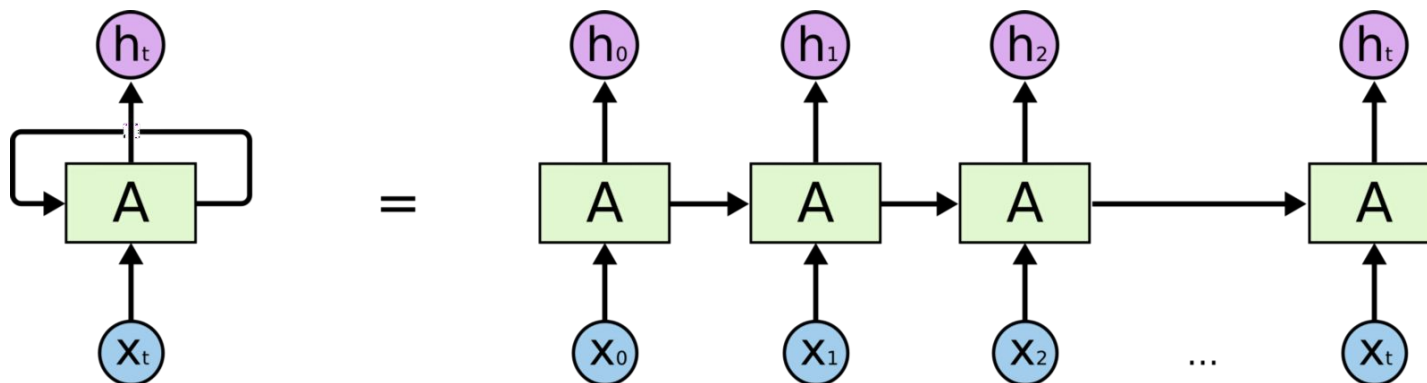
- ❑ End of the course material
- ❑ RNN/LSTM slides only for personal interest

Exploit Temporal Information



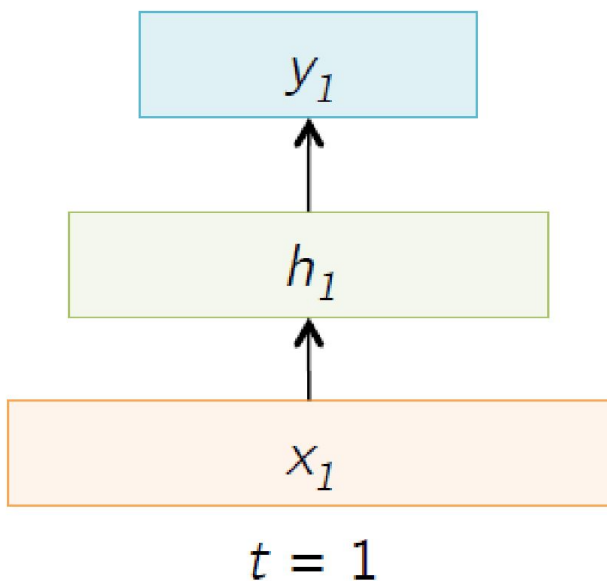
- ❑ Not all problem data can be fitted into a representation with fixed-length inputs and outputs !
- ❑ Problems such as **speech recognition** or **time-series prediction** require a system able to store and use context information
- ❑ Example:
 - Sequence of bits: output YES if the number of 1s is even, else NO
 - e.g., "1000010101" → YES (4 ones), "100011" → NO (3 ones), ...
 - Hard/Impossible to choose a fixed context window
 - There can always be a new sample longer than anything seen

Recurrent Neural Networks (RNN)

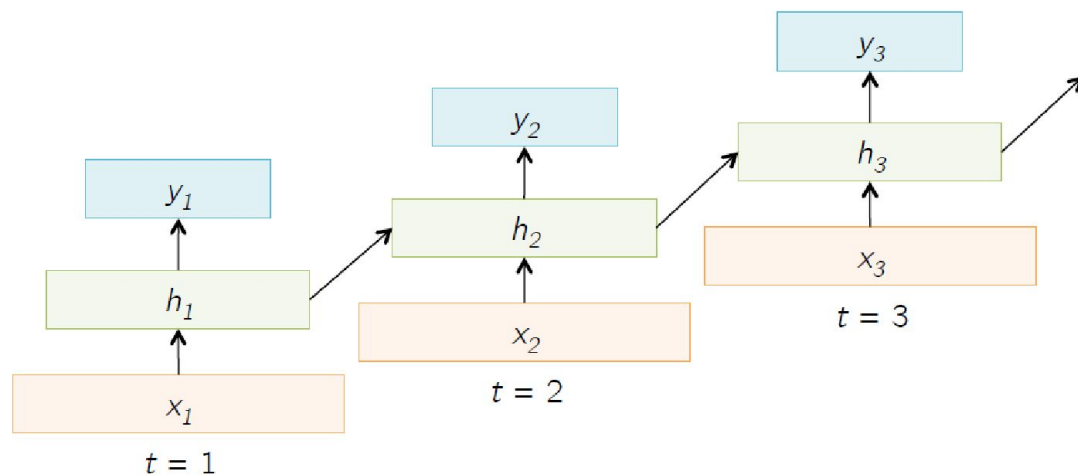


- ❑ Recurrent Neural Networks (RNN) take the previous outputs or hidden states as inputs
- ❑ The composite input at time t has some historical information about the happenings at times $t' < t$
- ❑ RNNs are useful as their intermediate values (state) can store information about past inputs for a time that is not fixed a priori

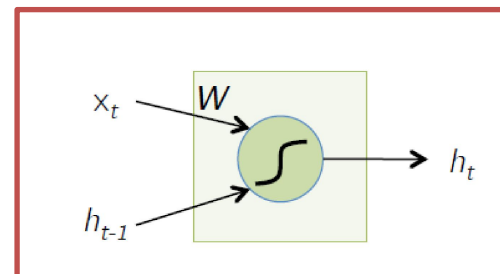
Feedforward vs Recurrent Networks



Sample Feedforward Network

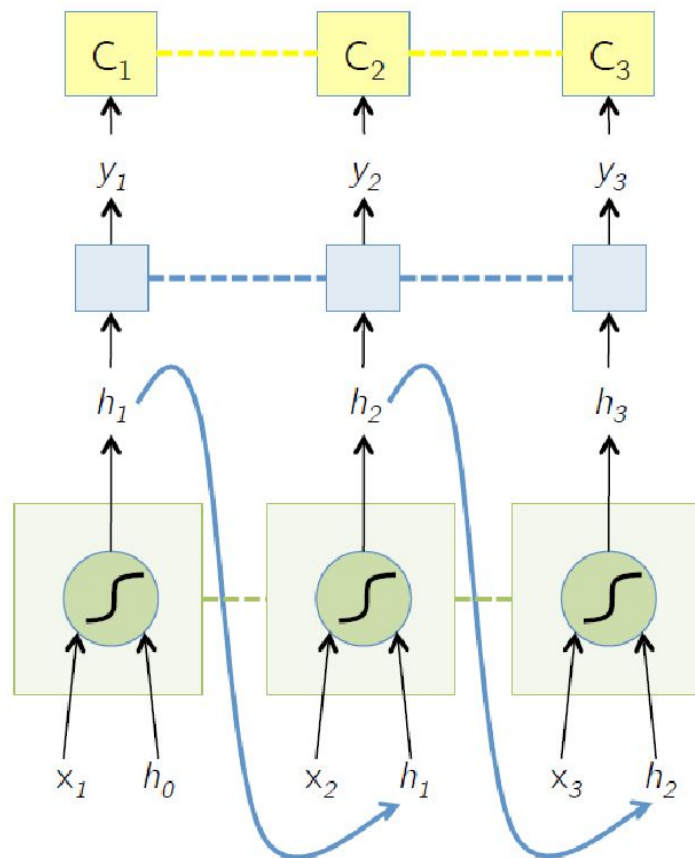


Sample Recurrent Network



Basic RNN cell

Basic RNN model



$$y_t = F(h_t)$$

$$C_t = \text{Loss}(y_t, \text{GT}_t)$$

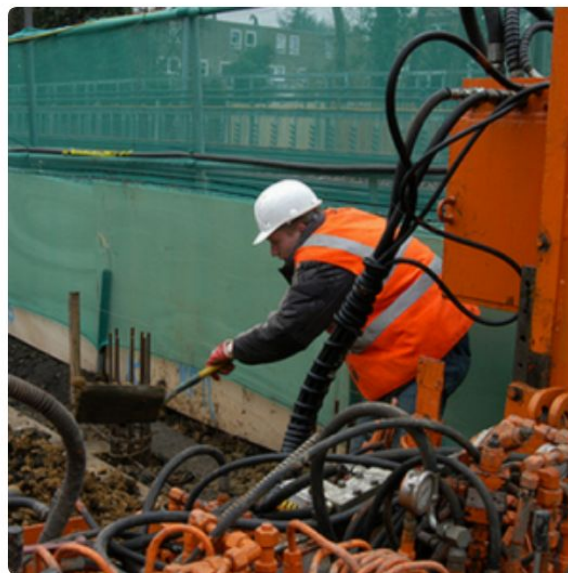
----- indicates shared weights

- Note that the weights are **shared** over time!
- Essentially, copies of the RNN cell are made over time (unrolling/unfolding), with different inputs at different time steps !

Example: Image Captioning (1)



"man in black shirt is playing guitar."



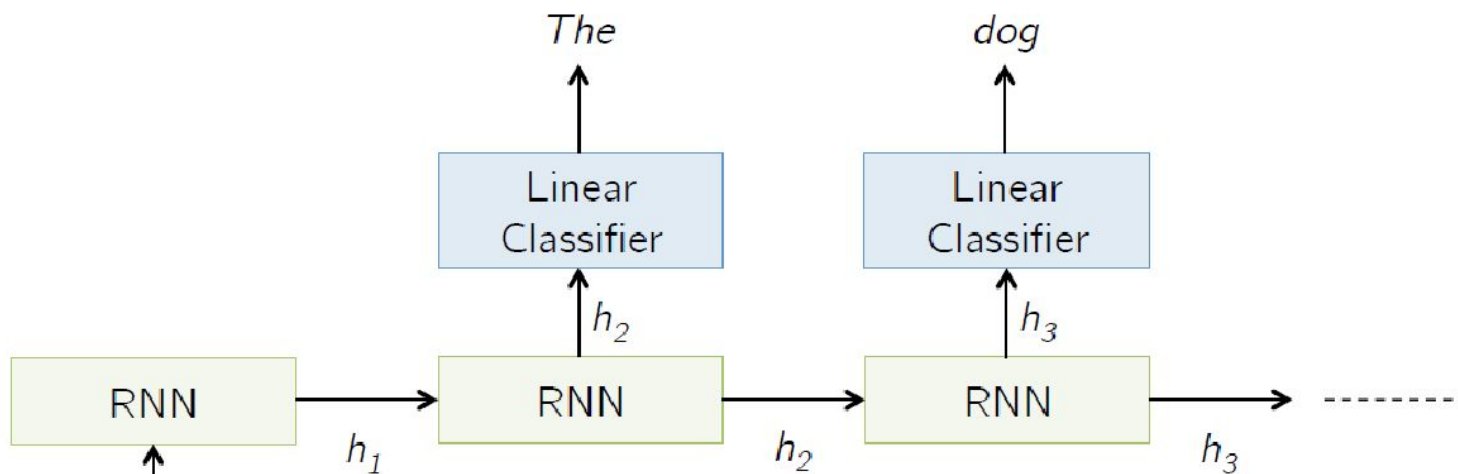
"construction worker in orange safety vest is working on road."



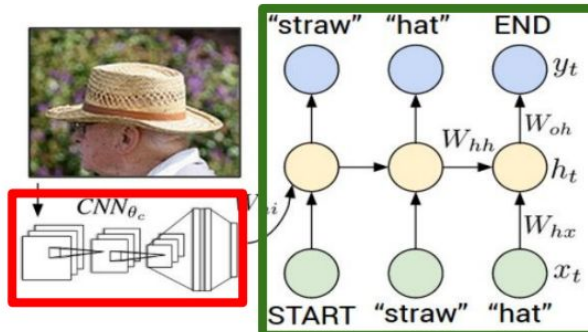
"two young girls are playing with lego toy."

- Given an image produce a sentence describing its content
- **Input:** Image features (e.g., output of a CNN)
- **Output:** Multiple words (e.g., one sentence)

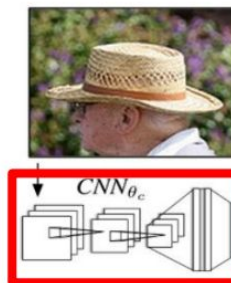
Example: Image Captioning (2)



Recurrent Neural Network



Convolutional Neural Network



Example: Image Captioning (3)

A person riding a motorcycle on a dirt road.



Two dogs play in the grass.



A herd of elephants walking across a dry grass field.



A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



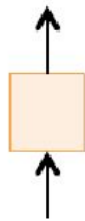
A close up of a cat laying on a couch.





Input-Output Scenarios

Single - Single



Feed-forward Network

Single - Multiple

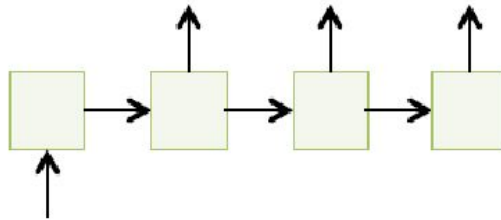
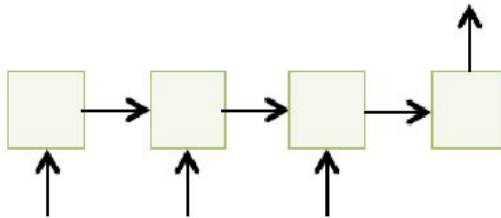


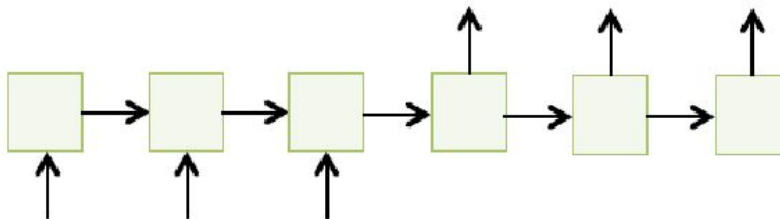
Image Captioning

Multiple - Single



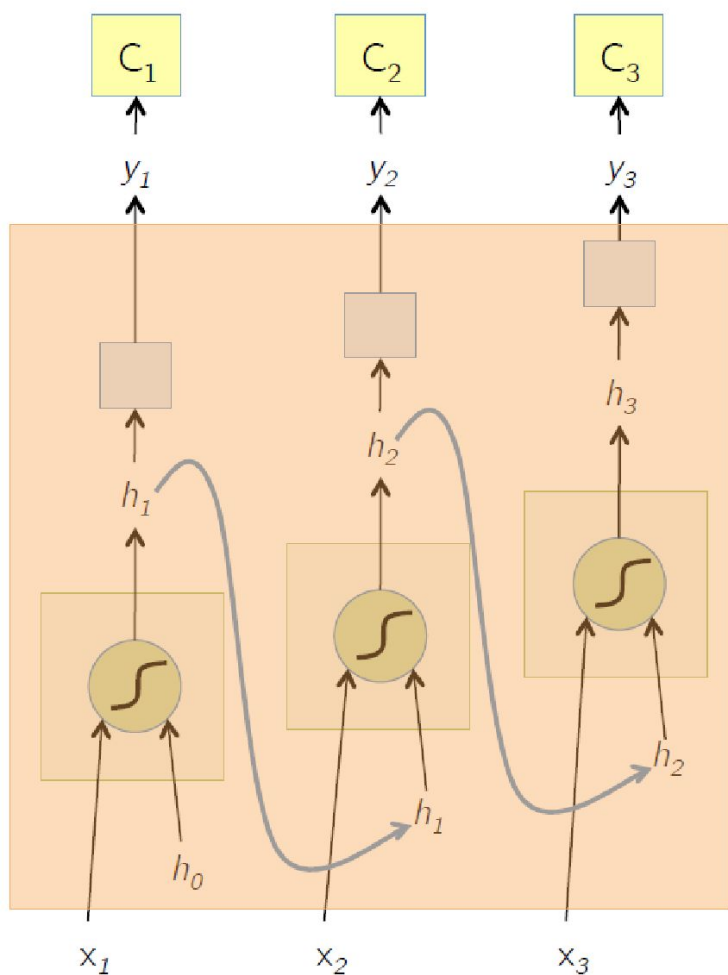
Sentiment Classification

Multiple - Multiple



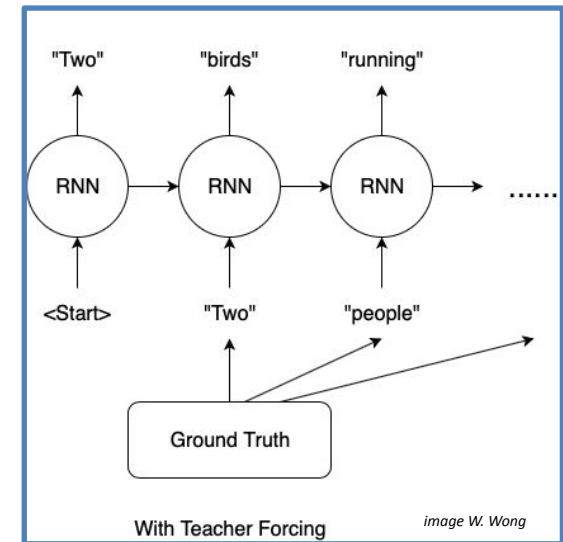
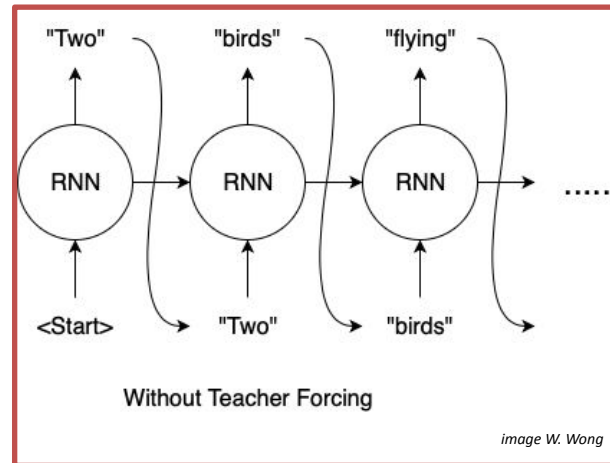
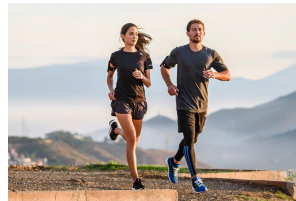
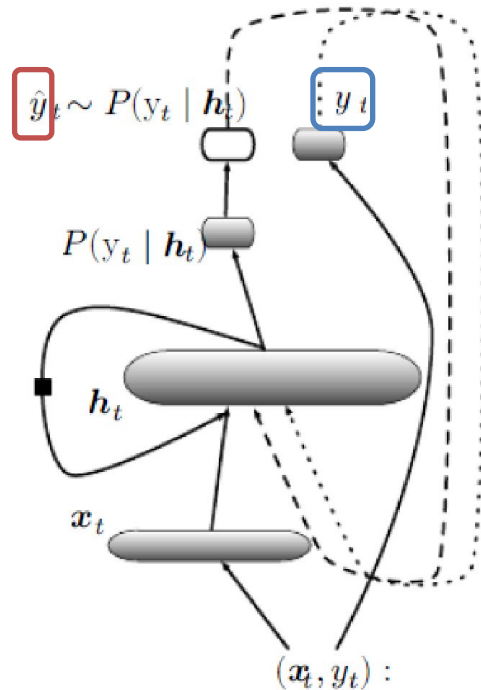
Translation

BackPropagation Through Time (BPTT)



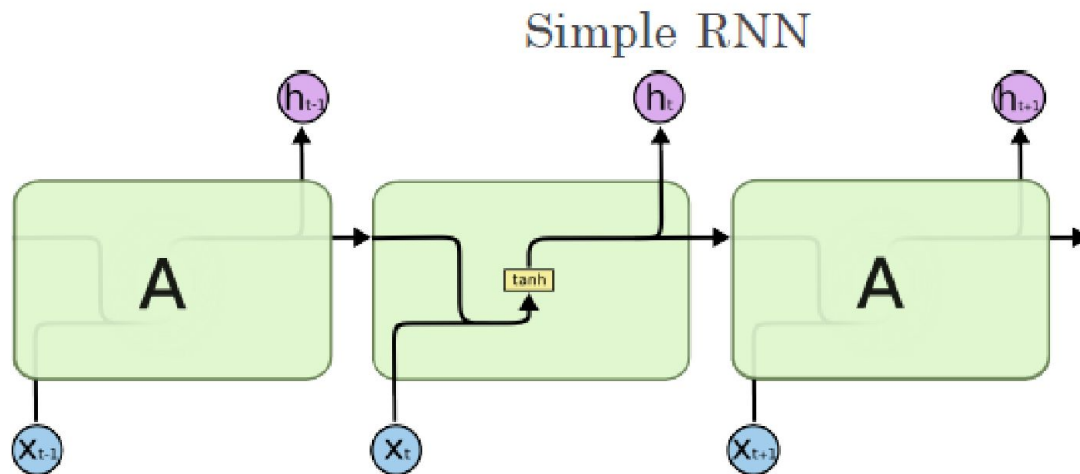
- ❑ One of the methods used to train RNNs
- ❑ The unfolded network (used during forward pass) is treated as one big feed-forward network!
- ❑ This unfolded network accepts the whole time series as input
- ❑ The weight updates are computed for each copy in the unfolded network (using standard BackPropagation), then summed (or averaged) and finally applied to the RNN weights
- ❑ Training RNN is challenging

Teacher Forcing



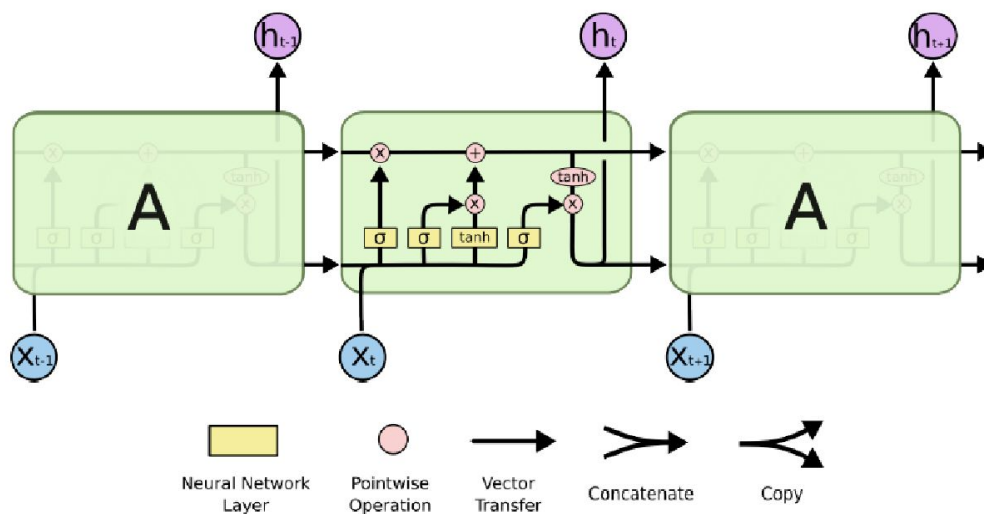
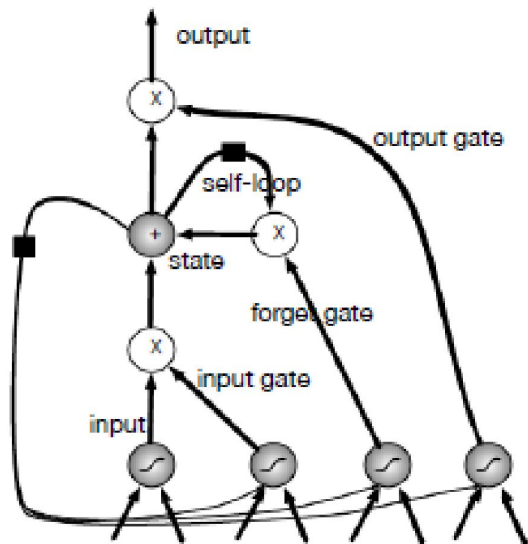
- ❑ If the output is used for the hidden state it is possible to choose if using **network's output** or **ground truth labels**
- ❑ With GT labels (**teacher forcing**) the model is easier to train
- ❑ But generalization properties can be poor

Model Long Time Temporal Relationship



- Baseline RNNs are good for short time temporal relationships
- But they are not able to capture long-time relationships since the gradients vanish or explode
- Also in some applications (e.g., word recognition) a way of "*forgetting*" the state is needed

Long-Short-Term-Memory (LSTM)



<p>Cell state: Taken in input from previous instant and forwarded to the next</p>	<p>Forget gate: If set to 0 can delete the content of the cell state</p>	<p>Input gate: Enable contribution of new input (the gate can enable/disable the input)</p>	<p>Output: Filtered version of state controlled by output gate</p>



References

- [1]: S Hochreiter, J Schmidhuber, "[Long short-term memory](#)", Neural computation 9 (8), 1735-1780, 1997
- [2]: Goodfellow, J Pouget-Abadie, M Mirza, B Xu, D Warde-Farley, S Ozair, "[Generative adversarial nets](#)", Advances in neural information processing systems, 2672-2680, 2014
- [3] Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "[Imagenet classification with deep convolutional neural network](#)" Advances in neural information processing systems, 2012
- [4] He, K., Zhang, X., Ren, S., & Sun, J. (2016). "[Deep residual learning for image recognition](#)", In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 770-778).
- [5] Yu, F., & Koltun, V. (2015). "[Multi-scale context aggregation by dilated convolutions](#)", arXiv preprint arXiv:1511.07122.

The papers can be downloaded from elearning