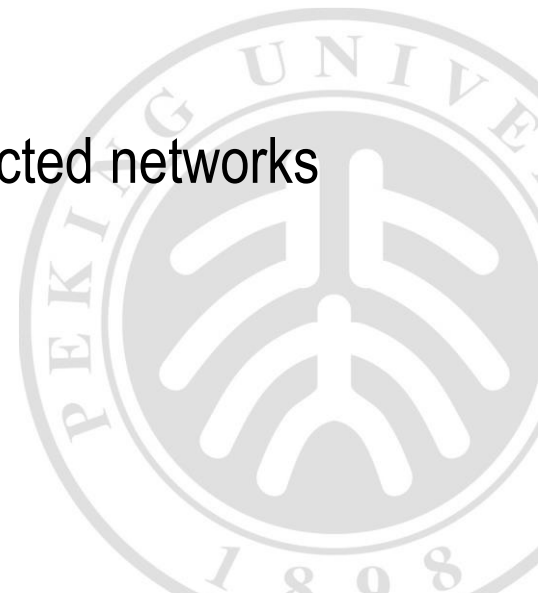- *Slides modified from Feifei Li.*

- Artificial Neurons

  - **Perceptrons** can do simple classification tasks

  - **Multi-Layer Non-Linear** Neural Network can approximate a wide range of functions

  - **Back-Propagation** with the chain rule to train the network

- Convolutional Neural Networks

  - Convolution operates on **local spatial** area

  - **Fewer parameters** compared with fully connected networks

  - **Automatically** extracts **features** from images

- CNN Architectures
  - AlexNet
  - VGG
  - GoogLeNet
  - ResNet
- Recurrent Neural Network
  - Vanilla RNN
  - Backpropagation through time
  - Long Short-Term Memory
- Beyond CNN and RNN
  - Unsupervised Learning
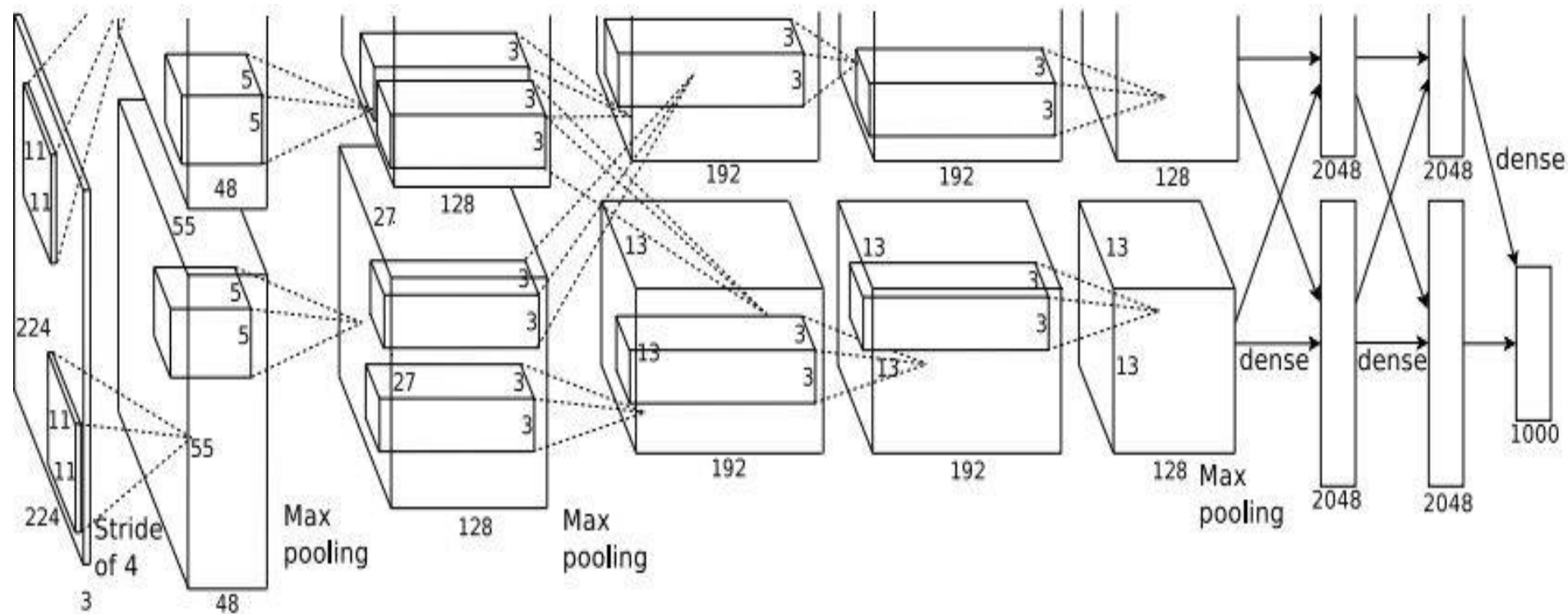  - Generative Adversarial Network

- **CNN Architectures**
  - **AlexNet**
  - **VGG**
  - **GoogLeNet**
  - **ResNet**
- Recurrent Neural Network
  - Vanilla RNN
  - Backpropagation through time
  - Long Short-Term Memory
- Beyond CNN and RNN
  - Unsupervised Learning
  - Generative Adversarial Network

## Architecture:

CONV1  MAX POOL1  NORM1  CONV2  MAX POOL2  NORM2  CONV3  CONV4

CONV5 MAX POOL3  FC6 FC7  FC8

**Full (simplified) AlexNet architecture:**

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

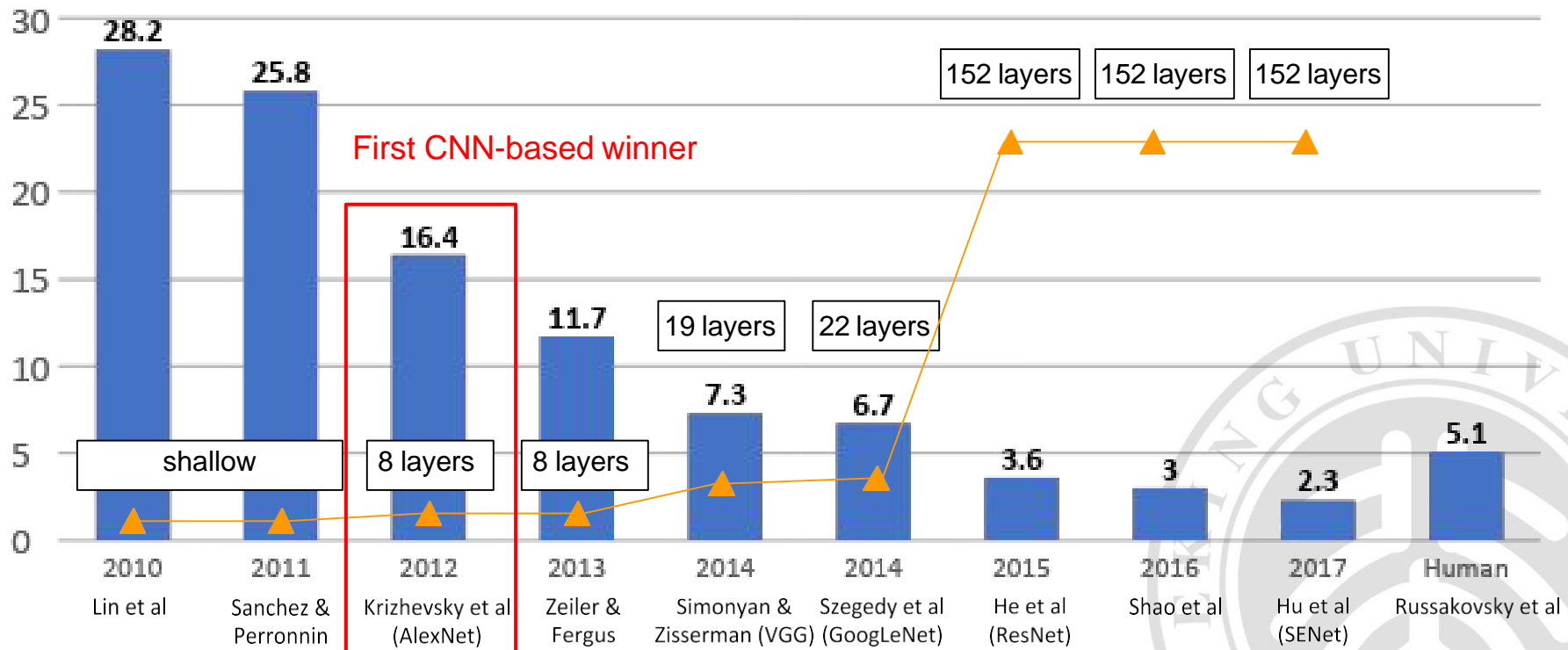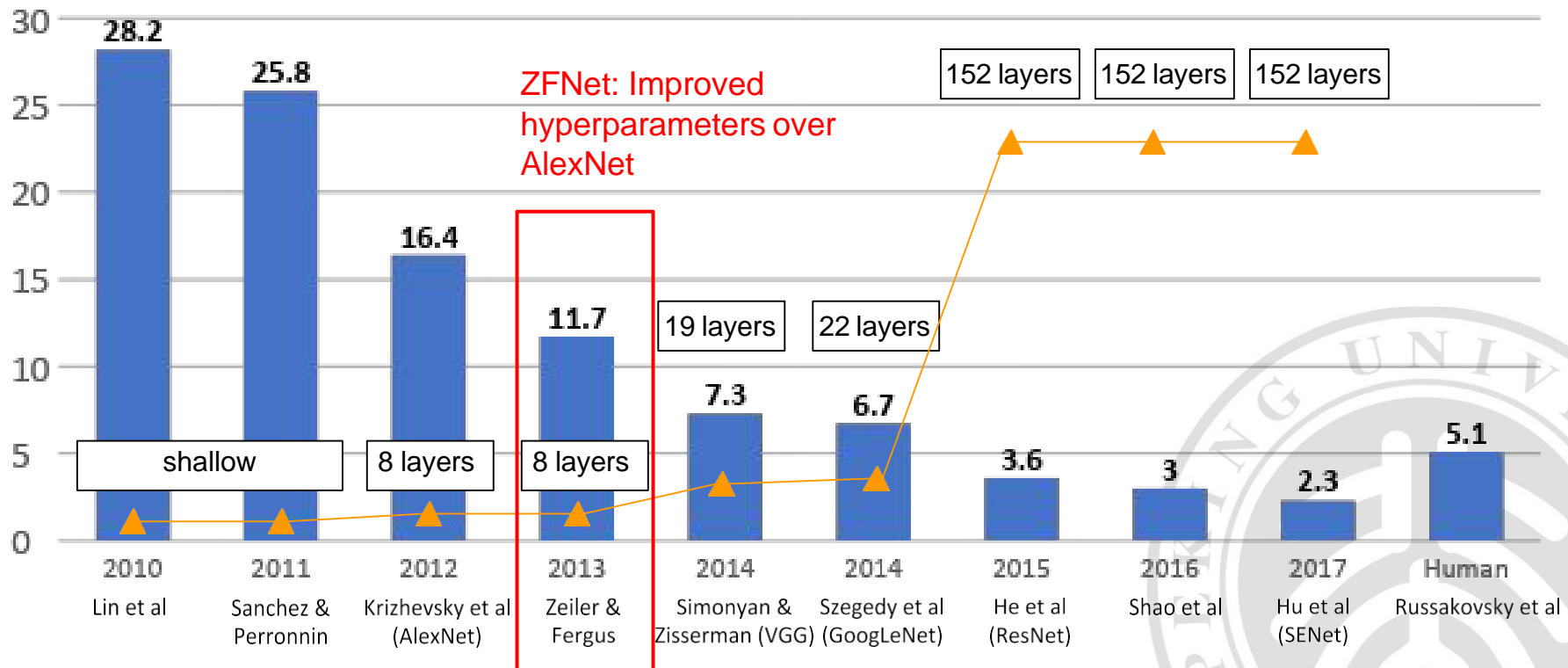[4096] FC7: 4096 neurons
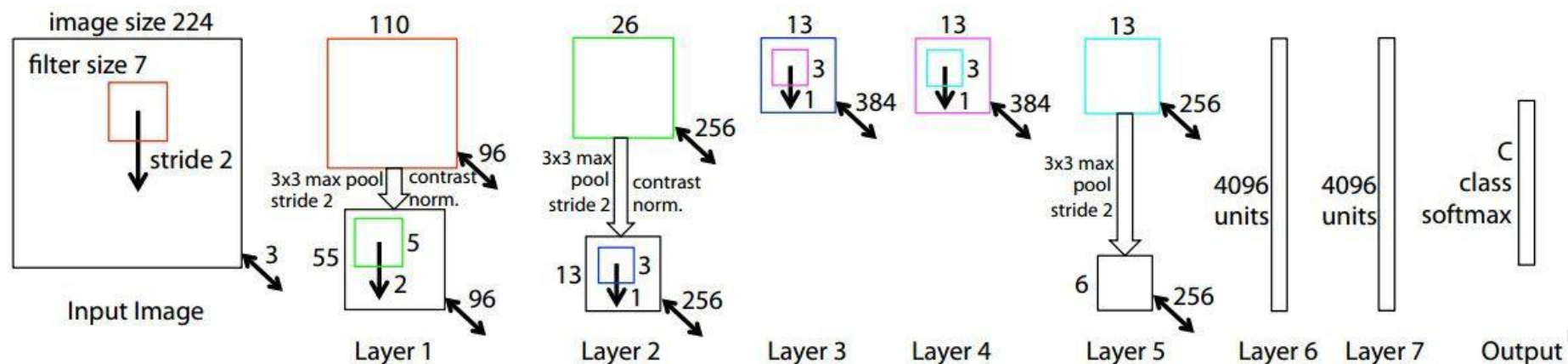
[1000] FC8: 1000 neurons (class scores)

**Details / Retrospectives:**
- First use of ReLU
- Used Norm layers (not common anymore)
- Heavy data augmentation
- Dropout 0.5
- Batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% →15.4%

**Full (simplified) AlexNet architecture:**

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)

**Details/Retrospectives:**

- First use of ReLU
- Used Norm layers (not common anymore)
- Heavy data augmentation
- Dropout 0.5
- Batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
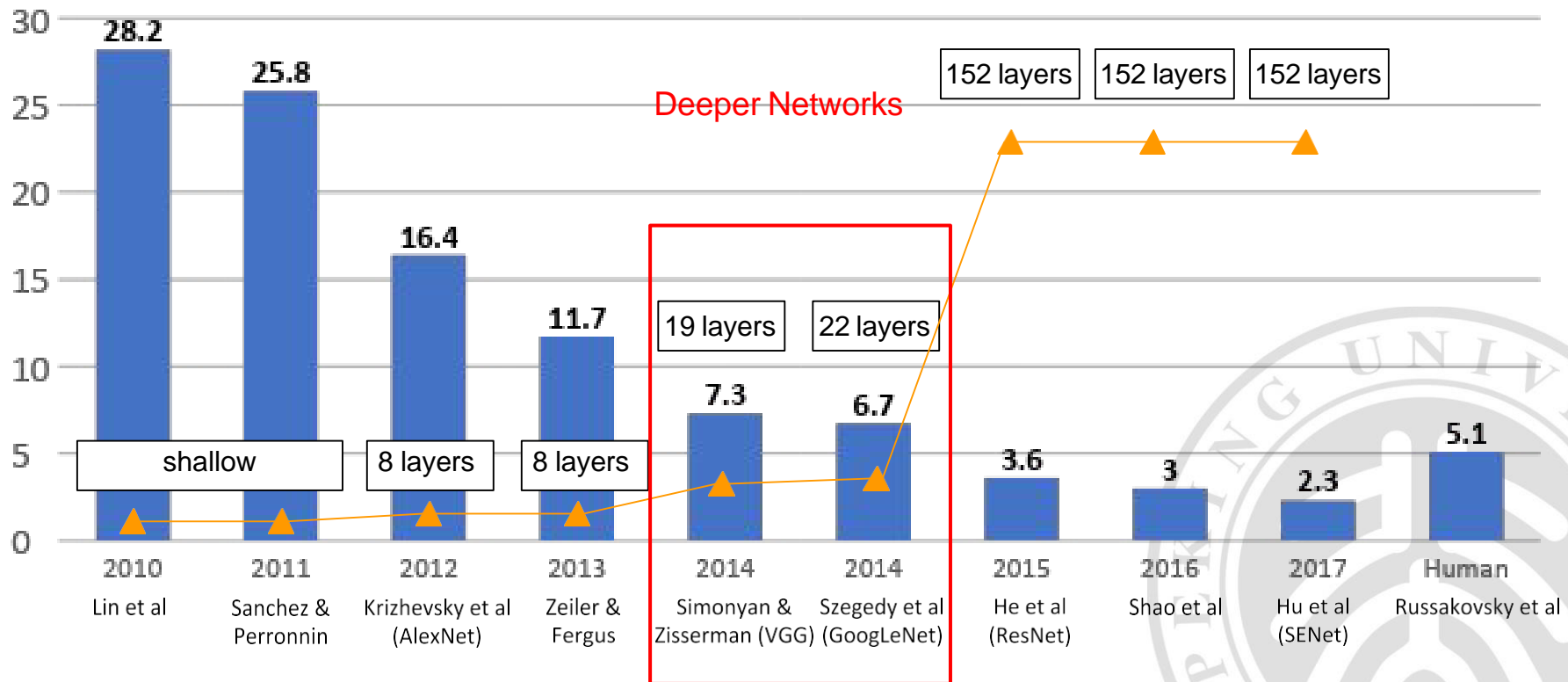- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% → 15.4%

AlexNet but:
CONV1: change from (11x11 stride 4) x48 to (7x7 stride 2) x96
CONV3,4,5: instead of 384, 384, 256 filters use 512, 1024, 512

ImageNet top 5 error: 16.4% → 11.7%

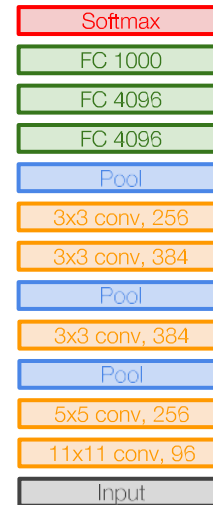## Small filters, deeper networks

8 layers (AlexNet)
→ 16 - 19 layers (VGG16Net)

Only 3x3 CONV: stride 1, pad 1 and 2x2
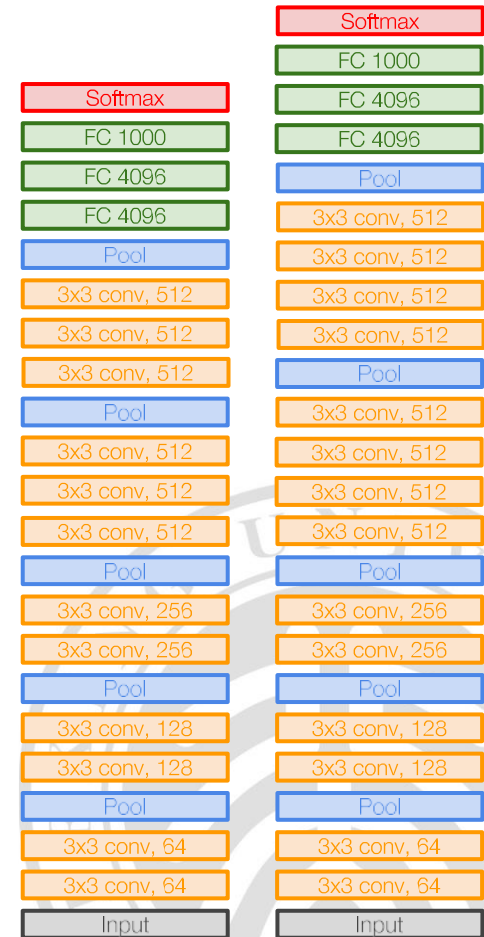MAX POOL: stride 2

11.7% top 5 error in ILSVRC'13  (ZFNet)
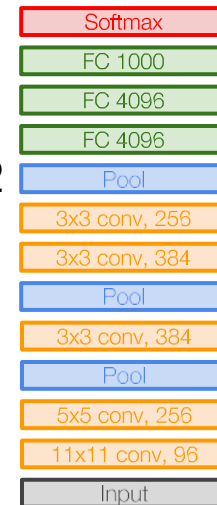→ 7.3% top 5 error in ILSVRC'14

| AlexNet |
|---|
| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 384 |
| Pool |
| 3x3 conv, 384 |
| Pool |
| 5x5 conv, 256 |
| 11x11 conv, 96 |
| Input |

| VGG16 |
|---|
| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv, 128 |
| Pool |
| 3x3 conv, 64 |
| 3x3 conv, 64 |
| Input |

| VGG19 |
|---|
| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv, 128 |
| Pool |
| 3x3 conv, 64 |
| 3x3 conv, 64 |
| Input |

## Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer

But deeper, more non-linearities

And fewer parameters: $3 * (3^2 C^2)$ vs. $7^2 C^2$ for C channels per layer
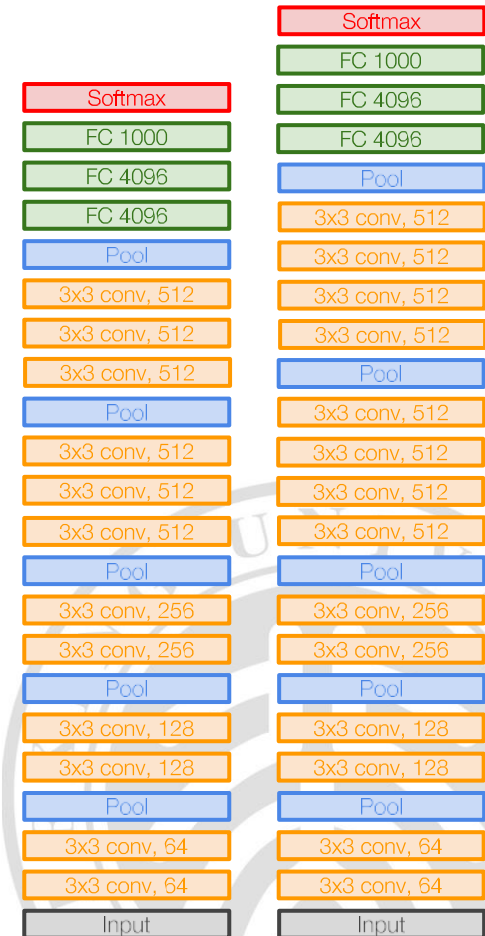


AlexNet          VGG16          VGG19

- **Details:**
  - ILSVRC'14 2nd in classification, 1st in localization
  - Similar training procedure as Krizhevsky 2012
  - No Local Response Normalisation (LRN)
  - Use VGG16 or VGG19 (VGG19 only slightly better, more memory)
  - Use ensembles for best results
  - FC7 features generalize well to other tasks

AlexNet

VGG16

VGG19

Deeper networks, with computational efficiency

- 22 layers
- Efficient "Inception" module
- No FC layers
- Only 5 million parameters!
  12x less than AlexNet
- ILSVRC'14 classification winner  (6.7% top 5 error)

Inception module

"Inception module":

Design a good local network topology (network within a network),
and then stack these modules on top of each other.



Inception module

Apply parallel filter operations on the input from previous layer:
- Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)
- Pooling operation (3x3)



Naive Inception module

Concatenate all filter outputs together depth-wise

Q: What is the problem with this?
[Hint: Computational complexity]

Q: What is the problem with this?
[Hint: Computational complexity]

Example:

Q3:What is output size after filter concatenation?

28x28x(128+192+96+256) = 28x28x672

| Filter concatenation |

28x28x128     28x28x192     28x28x96     28x28x256

| 1x1 conv, 128 | | 3x3 conv, 192 | | 5x5 conv, 96 | | 3x3 pool |

Module input: 28x28x256

| Input |

Naive Inception module

**Conv Ops:**
[1x1 conv, 128] 28x28x128x1x1x256
[3x3 conv, 192] 28x28x192x3x3x256
[5x5 conv, 96] 28x28x96x5x5x256
**Total: 854M ops**

Very expensive compute

Pooling layer also preserves feature depth, which means total depth after concatenation can only grow at every layer!

Q: What is the problem with this? [Hint: Computational complexity]

Example:

Q3:What is output size after filter concatenation?

28x28x(128+192+96+256) = **529k**

Filter concatenation

28x28x128    28x28x192    28x28x96    28x28x256

| 1x1 conv, 128 | 3x3 conv, 192 | 5x5 conv, 96 | 3x3 pool |

Module input: 28x28x256

Input

Naive Inception module

Solution: "bottleneck" layers that use 1x1 convolutions to reduce feature depth

1x1 CONV
with 32 filters

(Each filter has size 1x1x64, and performs a 64-dimensional dot product)

56

56

64

56

56

32

56

56

64

1x1 CONV
with 32 filters

Preserves spatial
dimensions, reduces depth!

Projects depth to lower
dimension (combination of
feature maps)

56

56

32

1x1 conv "bottleneck" layers



Naive Inception module

Inception module with dimension reduction

28x28x480

Filter
concatenation

28x28x128  28x28x192  28x28x96  28x28x64

| 1x1 conv, 128 | 3x3 conv, 192 | 5x5 conv, 96 | 1x1 conv, 64 |

28x28x64  28x28x64  28x28x256

| 1x1 conv, 64 | 1x1 conv, 64 | 3x3 pool |

Module input:
28x28x256

Previous Layer

Inception module with dimension reduction

Using same parallel layers as naive example, and adding "1x1 conv, 64 filter" bottlenecks:

**Conv Ops:**
[1x1 conv, 64] 28x28x64x1x1x256
[1x1 conv, 64] 28x28x64x1x1x256
[1x1 conv, 128] 28x28x128x1x1x256
[3x3 conv, 192] 28x28x192x3x3x64
[5x5 conv, 96] 28x28x96x5x5x64
[1x1 conv, 64] 28x28x64x1x1x256
**Total: 358M ops**

Compared to 854M ops for naive version, bottleneck can also reduce depth after pooling layer

Full GoogLeNet
architecture



Stem Network:
Conv-Pool-
2x Conv-Pool

Full GoogLeNet
architecture



Stacked Inception
Modules

Full GoogLeNet
architecture



Classifier output
(removed expensive FC layers!)
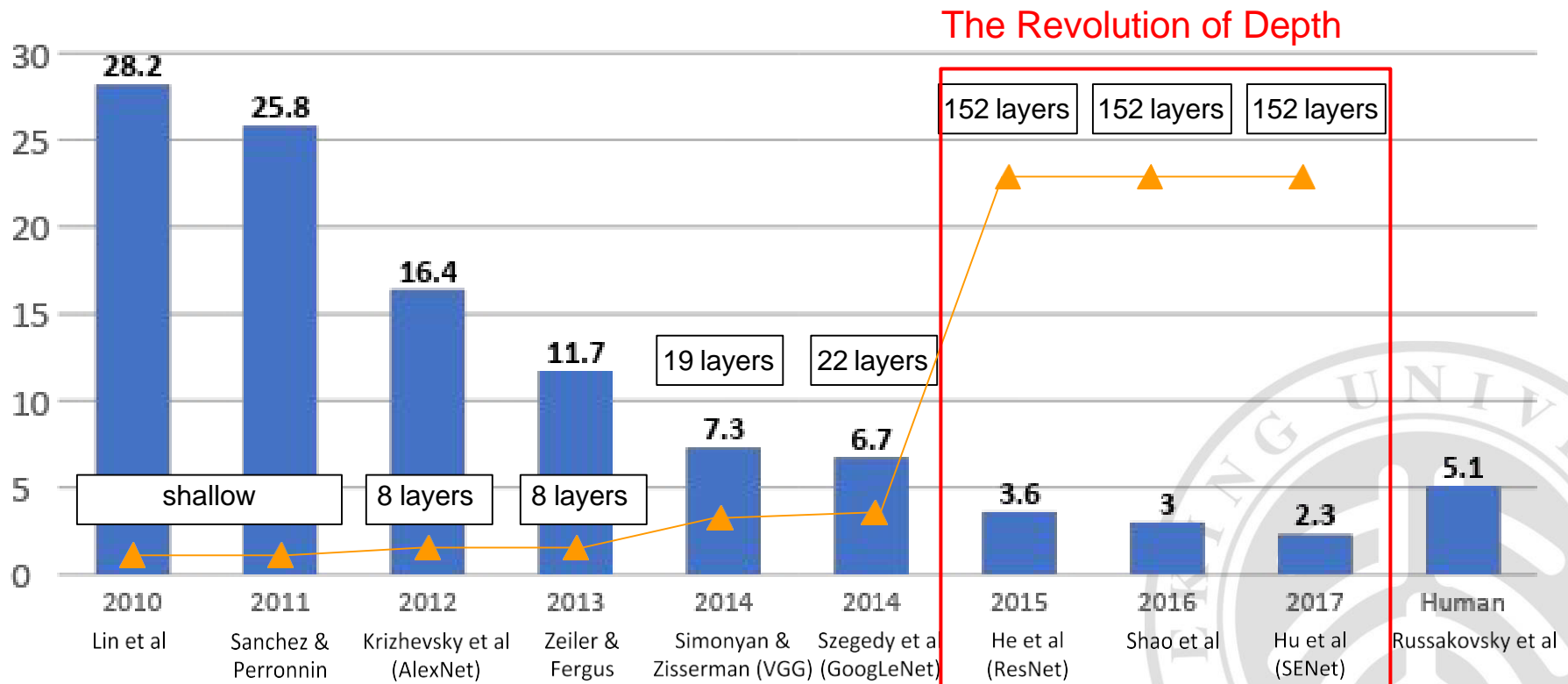
Full GoogLeNet architecture



Auxiliary classification outputs to inject additional gradient at lower layers
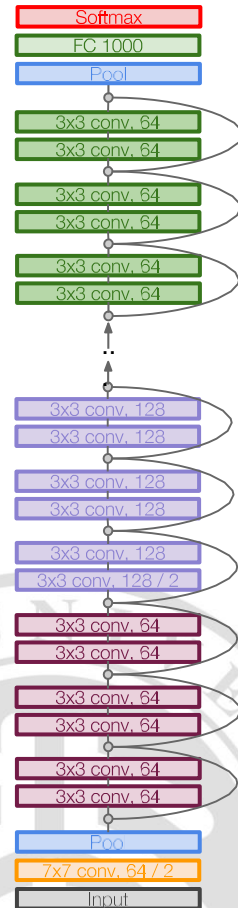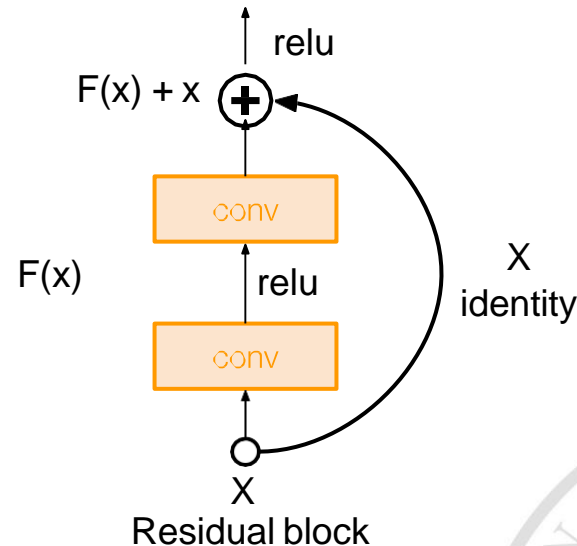(AvgPool-1x1Conv-FC-FC-Softmax)

Full GoogLeNet architecture



22 total layers with weights
(parallel layers count as 1 layer => 2 layers per Inception module. Don't count auxiliary output layers)

The Revolution of Depth

**Very deep networks using residual connections**

- 152-layer model for ImageNet
- ILSVRC'15 classification winner (3.57% top 5 error)
- Swept all classification and detection competitions in ILSVRC'15 and COCO'15!

relu

$F(x) + x$ ⊕

conv

$F(x)$  relu

conv

X identity

X

Residual block

Softmax
FC 1000
Pool
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
...
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128 / 2
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
Pool
7x7 conv, 64 / 2
Input

What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



56-layer model performs worse on both training and test error
→ The deeper model performs worse, but it's not caused by overfitting!

Hypothesis: the problem is an *optimization* problem, deeper models are harder to optimize

The deeper model should be able to perform at least as well as the shallower model.

A solution by construction is copying the learned layers from the shallower model and setting additional layers to identity mapping.

Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



"Plain" layers

Residual block

Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



$H(x) = F(x) + x$

"Plain" layers

Residual block

Use layers to fit residual $F(x) = H(x) - x$ instead of $H(x)$ directly

## Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
- Additional conv layer at the beginning



Residual block

For deeper networks
(ResNet-50+),
use "bottleneck" layer to
improve efficiency
(similar to GoogLeNet)

28x28x256
output

1x1 conv, 256 filters projects
back to 256 feature maps
(28x28x256)

3x3 conv operates over
only 64 feature maps

1x1 conv, 64 filters
to project to
28x28x64

1x1 conv, 256

3x3 conv, 64

1x1 conv, 64

28x28x256
input

Training ResNet in practice:

- Batch Normalization after every CONV layer
- Xavier 2/ initialization from He et al.
- SGD + Momentum (0.9)
- Learning rate: 0.1, divided by 10 when validation error plateaus
- Mini-batch size 256
- Weight decay of 1e-5
- No dropout used

## Experimental Results

- Able to train very deep networks without degrading (152 layers on ImageNet, 1202 on Cifar)
- Deeper networks now achieve lowing training error as expected
- Swept 1st place in all ILSVRC and COCO 2015 competitions

MSRA @ ILSVRC & COCO 2015 Competitions

- **1st places** in all five main tracks
  - ImageNet Classification: *"Ultra-deep"* (quote Yann) **152-layer** nets
  - ImageNet Detection: **16%** better than 2nd
  - ImageNet Localization: **27%** better than 2nd
  - COCO Detection: **11%** better than 2nd
  - COCO Segmentation: **12%** better than 2nd

ILSVRC 2015 classification winner (3.6%  top 5 error)
-- better than "human  performance"!
(Russakovsky 2014)

- **VGG, GoogLeNet, ResNet** all in wide use, available in model zoos

- ResNet current best default

- Trend towards extremely deep networks

- Significant research centers around design of layer / skip connections and improving gradient flow

- Efforts to investigate necessity of depth Vs. width and residual connections

- CNN Architectures
  - AlexNet
  - VGG
  - GoogLeNet
  - ResNet
- **Recurrent Neural Network**
  - **Vanilla RNN**
  - **Backpropagation through time**
  - **Long Short-Term Memory**
- Beyond CNN and RNN
  - Unsupervised Learning
  - Generative Adversarial Network

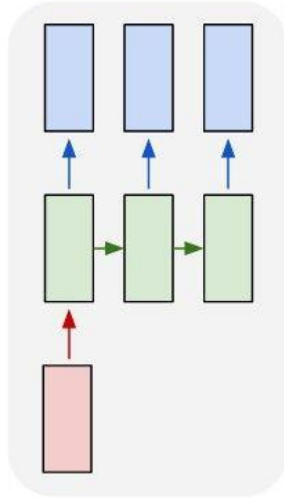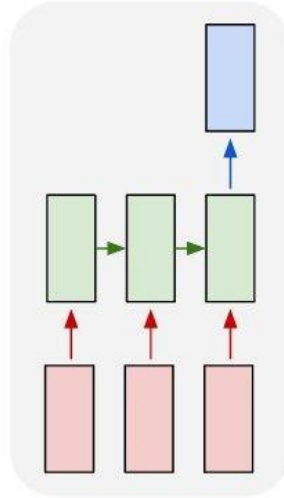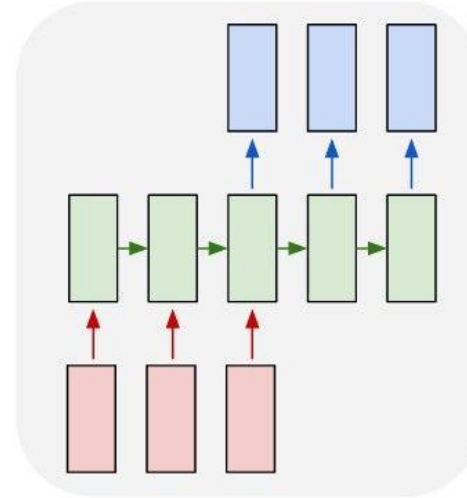one to one     one to many     many to one     many to many     many to many

e.g. **Image Captioning**
image → sequence of words

one to one one to many many to one many to many many to many

e.g. **Sentiment Classification**
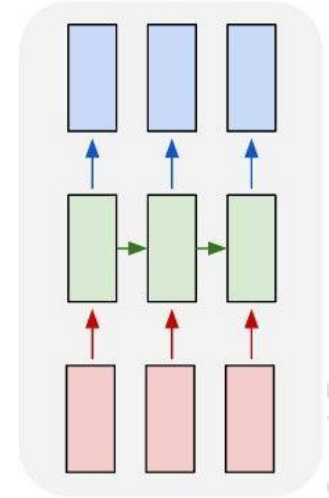sequence of words → sentiment

one to one    one to many    many to one    many to many    many to many
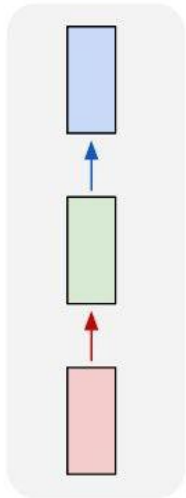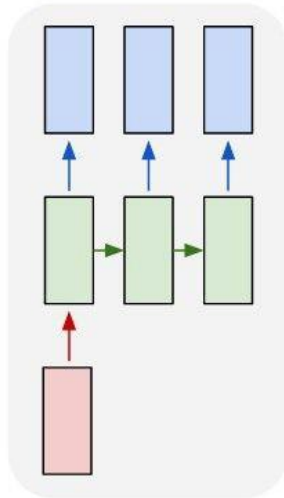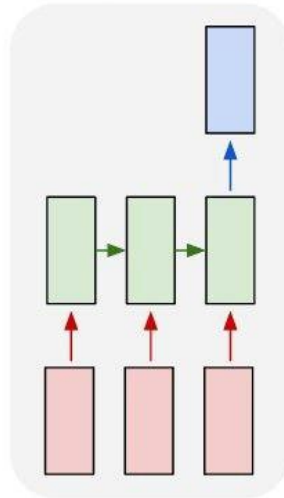
e.g. **Machine Translation**
seq of words → seq of words

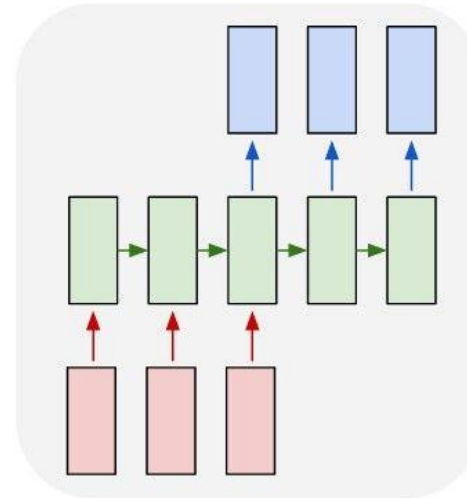e.g. **Video classification on frame level**

y

usually want to predict a
vector at some time steps

RNN

x

We can process a sequence of vectors **x** by applying a **recurrence formula** at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

new state

some function with parameters W

old state

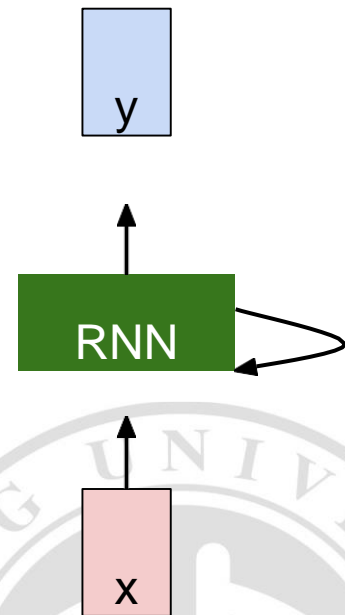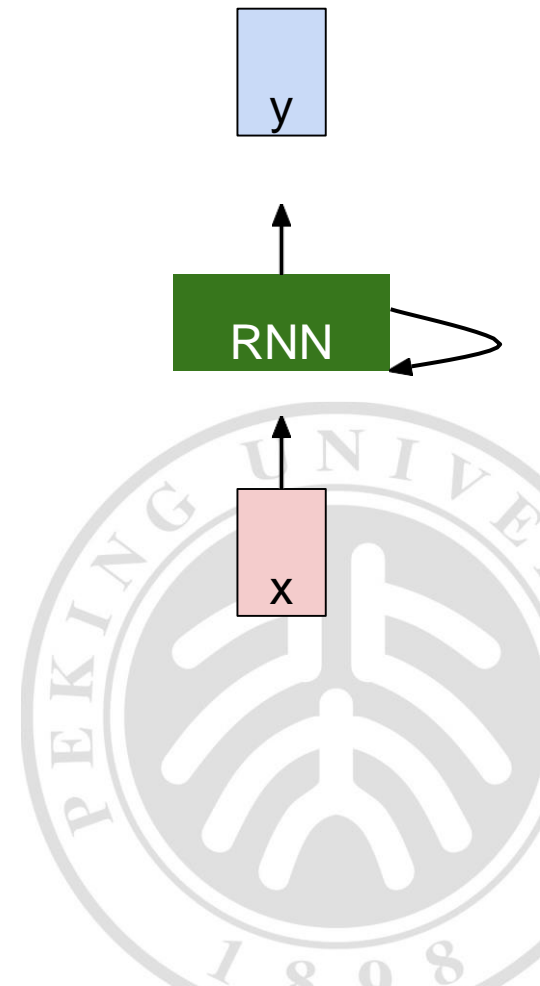input vector at some time step

y

RNN

x

We can process a sequence of vectors **x** by applying a **recurrence formula** at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

Notice: the same function and the same set of parameters are used at every time step.

The state consists of a single *"hidden"* vector **h**:
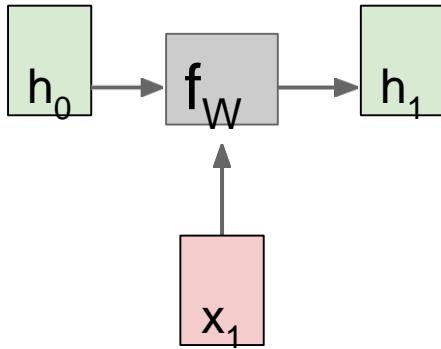
$$h_t = f_W(h_{t-1}, x_t)$$

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$

$$y_t = W_{hy} h_t$$

y

RNN

x

Sometimes called a "Vanilla RNN" or an "Elman RNN" after Prof. Jeffrey Elman

Re-use the same weight matrix at every time-step

**Many to one**: Encode input
sequence in a single vector



Sutskever et al, "Sequence to Sequence Learning with Neural Networks", NIPS 2014

**One to many:** Produce output sequence from single input vector

**Many to one**: Encode input sequence in a single vector



Sutskever et al, "Sequence to Sequence Learning with Neural Networks", NIPS 2014

Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient

Explain Images with Multimodal Recurrent Neural Networks, Mao et al.
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei
Show and Tell: A Neural Image Caption Generator, Vinyals et al.
Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.
Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

v

**Wih**

test image

y0

h0

x0
<START>

<START>

**before:**

h = tanh(Wxh * x + Whh * h)

**now:**

h = tanh(Wxh * x + Whh * h **+ Wih * v**)

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

y0    y1

h0 → h1

x0
<START>

straw

hat

<START>

sample!

sample
<END> token
=> finish.

*A cat sitting on a suitcase on the floor*



*A cat is sitting on a tree branch*



*A dog is running in the grass with a frisbee*



*A white teddy bear sitting in the grass*



*Two people walking on the beach with surfboards*



*A tennis player in action on the court*



*Two giraffes standing in a grassy field*



*A man riding a dirt bike on a dirt track*

*A woman is holding a cat in her hand*



*A person holding a computer mouse on a desk*



*A woman standing on a beach holding a surfboard*



*A bird is perched on a tree branch*



*A man in a baseball uniform throwing a ball*

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$h \in \mathbb{R}^n$      $W^l \; [n \times 2n]$

LSTM:

$W^l \; [4n \times 2n]$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

depth

time

Backpropagation from $h_t$ to $h_{t-1}$ multiplies by W

(actually $W_{hh}^T$)



$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$= \tanh\left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

$$= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

Computing gradient of $h_0$ involves many factors of W (and repeated tanh)

Largest singular value > 1: **Exploding gradients**

Largest singular value < 1: **Vanishing gradients**

**Gradient clipping**: Scale gradient if its norm is too big

```
grad_norm = np.sum(grad * grad)
if grad_norm > threshold:
    grad *= (threshold / grad_norm)
```

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

**Vanilla RNN**

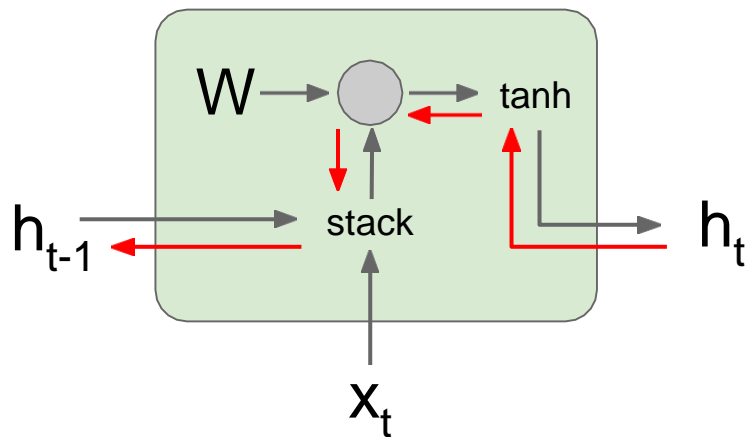$$h_t = \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

**LSTM**

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

## Uninterrupted gradient flow!



Similar to ResNet !

[Hochreiter et al., 1997]

- CNN Architectures
  - AlexNet
  - VGG
  - GoogLeNet
  - ResNet
- Recurrent Neural Network
  - Vanilla RNN
  - Backpropagation through time
  - Long Short-Term Memory
- **Beyond CNN and RNN**
  - **Unsupervised Learning**
  - **Generative Adversarial Network**

**Supervised Learning**

**Data**: (x, y)
x is data, y is label

**Goal**: Learn a *function* to map x → y

**Examples**: Classification, regression, object detection, semantic segmentation, image captioning, etc.



→ Cat

Classification

**Supervised Learning**

**Data**: (x, y)
x is data, y is label

**Goal**: Learn a *function* to map x → y

**Examples**: Classification, regression, object detection, semantic segmentation, image captioning, etc.



**DOG**, **DOG**, **CAT**

Object Detection

**Supervised Learning**

**Data**: (x, y)
x is data, y is label

**Goal**: Learn a *function* to map x $\rightarrow$ y

**Examples**: Classification, regression, object detection, semantic segmentation, image captioning, etc.



GRASS, CAT,
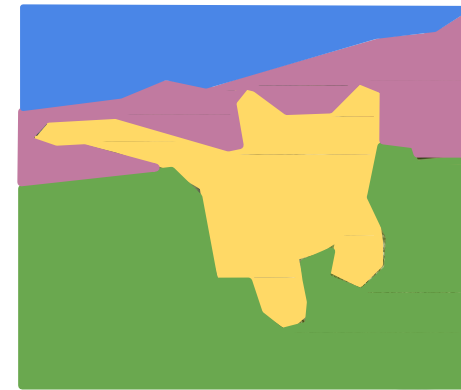TREE, SKY

Semantic Segmentation

**AIPKU**

## Supervised Learning

**Data**: (x, y)
x is data, y is label

**Goal**: Learn a *function* to map x $\rightarrow$ y

**Examples**: Classification, regression, object detection, semantic segmentation, image captioning, etc.



*A cat sitting on a suitcase on the floor*

Image Captioning

**Unsupervised Learning**

**Data**: x
Just data, no labels!

**Goal**: Learn some underlying hidden *structure* of the data

**Examples**: Clustering, dimensionality reduction, feature learning, density estimation, etc.

K-means clustering

**Unsupervised Learning**

**Data**: x
Just data, no labels!

**Goal**: Learn some underlying hidden *structure* of the data

**Examples**: Clustering, dimensionality reduction, feature learning, density estimation, etc.



3-d &rarr; 2-d

Principal Component Analysis
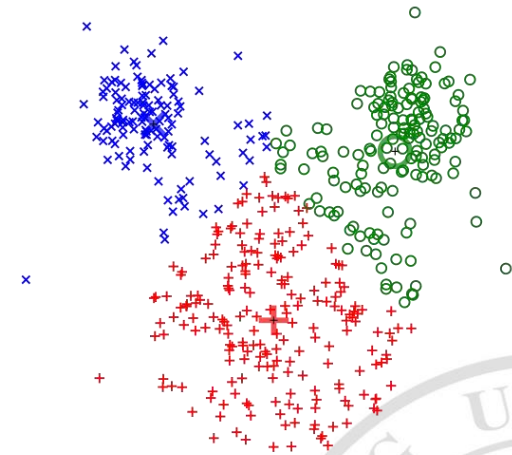(Dimensionality reduction)

**Unsupervised Learning**

**Data**: x
Just data, no labels!

**Goal**: Learn some underlying hidden *structure* of the data

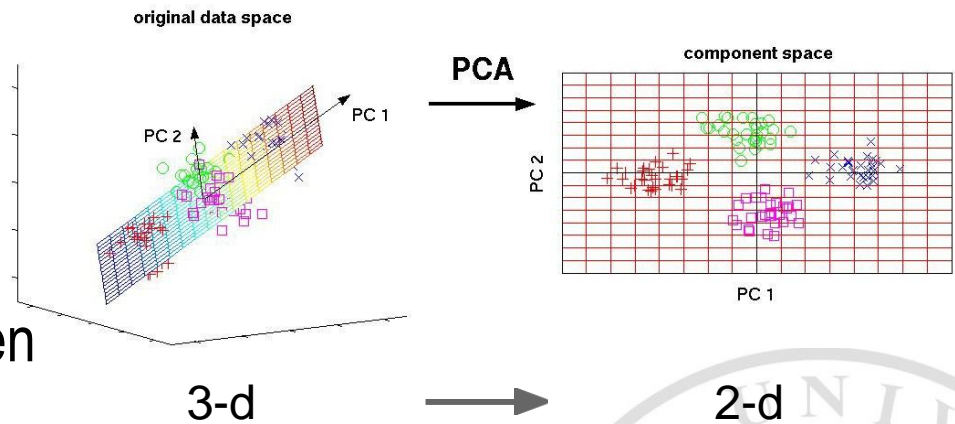**Examples**: Clustering, dimensionality reduction, feature learning, density estimation, etc.

L2 Loss function:
$$\|x - \hat{x}\|^2$$

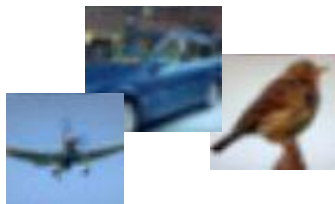Reconstructed input data — $\hat{x}$

Decoder

**Features** — $z$

Encoder

Input data — $x$

Reconstructed data

Encoder: 4-layer conv
Decoder: 4-layer upconv

Input data

Autoencoders
(Feature learning)

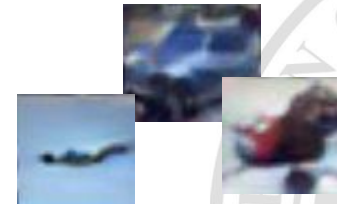- Given training data, generate new samples from same distribution
- **Several flavors:**
  - Explicit density estimation: explicitly define and solve for $p_{model}(x)$
  - Implicit density estimation: learn model that can sample from $p_{model}(x)$ w/o explicitly defining it

Training data ~ $p_{data}(x)$

Generated samples ~ $p_{model}(x)$

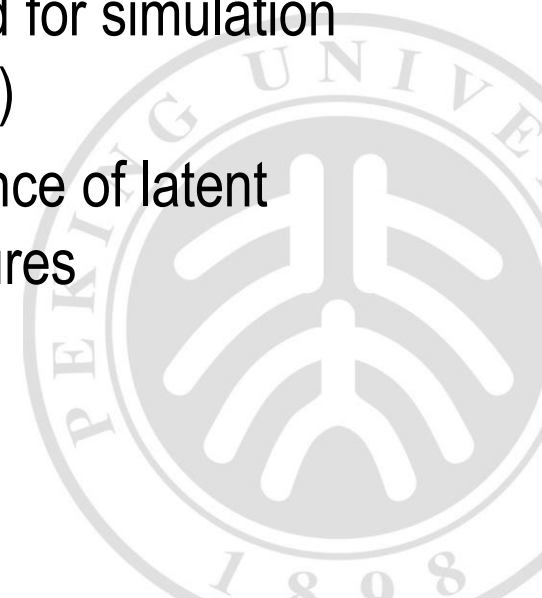Want to learn $p_{model}(x)$ similar to $p_{data}(x)$

- Realistic samples for artwork, super-resolution, colorization, etc.
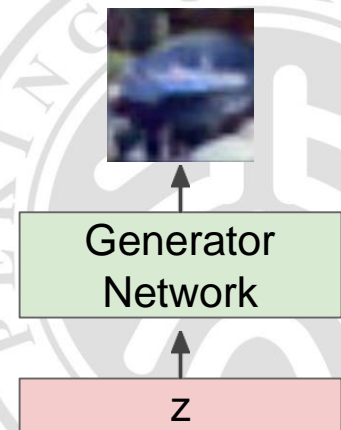


- Generative models of time-series data can be used for simulation and planning (reinforcement learning applications!)

- Training generative models can also enable inference of latent representations that can be useful as general features

- Problem: Want to sample from complex, high-dimensional training distribution. No direct way to do this!

- Solution: Sample from a simple distribution, e.g. random noise. Learn transformation to training distribution.

- Q: What can we use to represent this complex transformation?

- A: A neural network!

Output: Sample from
training distribution

Generator
Network

z

Input: Random noise

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

**Generator network**: try to fool the discriminator by generating real-looking images

**Discriminator network**: try to distinguish between real and fake images

**Generator network**: try to fool the discriminator by generating real-looking images

**Discriminator network**: try to distinguish between real and fake images
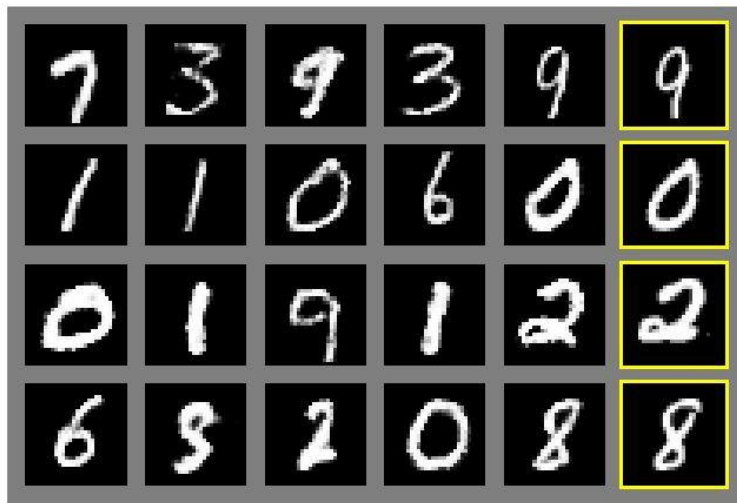
Train jointly in **minimax game**

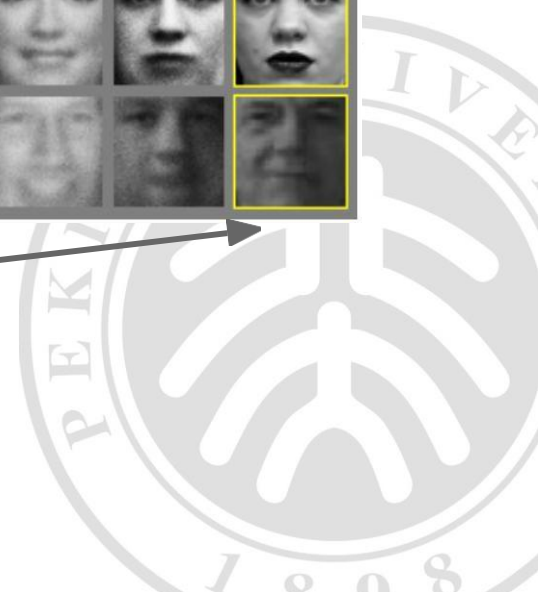<span style="color:blue">Discriminator outputs likelihood in (0,1) of real image</span>

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log \underbrace{D_{\theta_d}(x)} + \mathbb{E}_{z \sim p(z)} \log(1 - \underbrace{D_{\theta_d}(G_{\theta_g}(z))}) \right]$$

<span style="color:blue">Discriminator output for real data x</span>   <span style="color:blue">Discriminator output for generated fake data G(z)</span>

-   Discriminator ($\theta_d$) wants to **maximize objective** such that D(x) is close to 1 (real) and  D(G(z)) is close to 0 (fake)

-   Generator ($\theta_g$) wants to **minimize objective** such that D(G(z)) is close to 1 (discriminator is fooled into thinking generated G(z) is real)
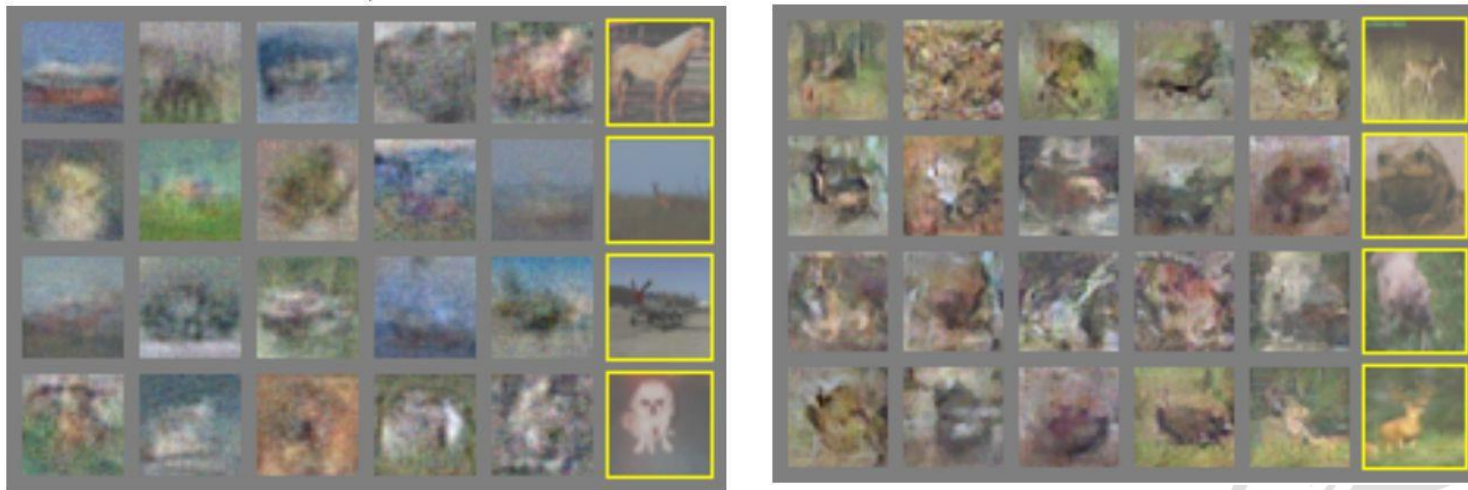
Generated samples



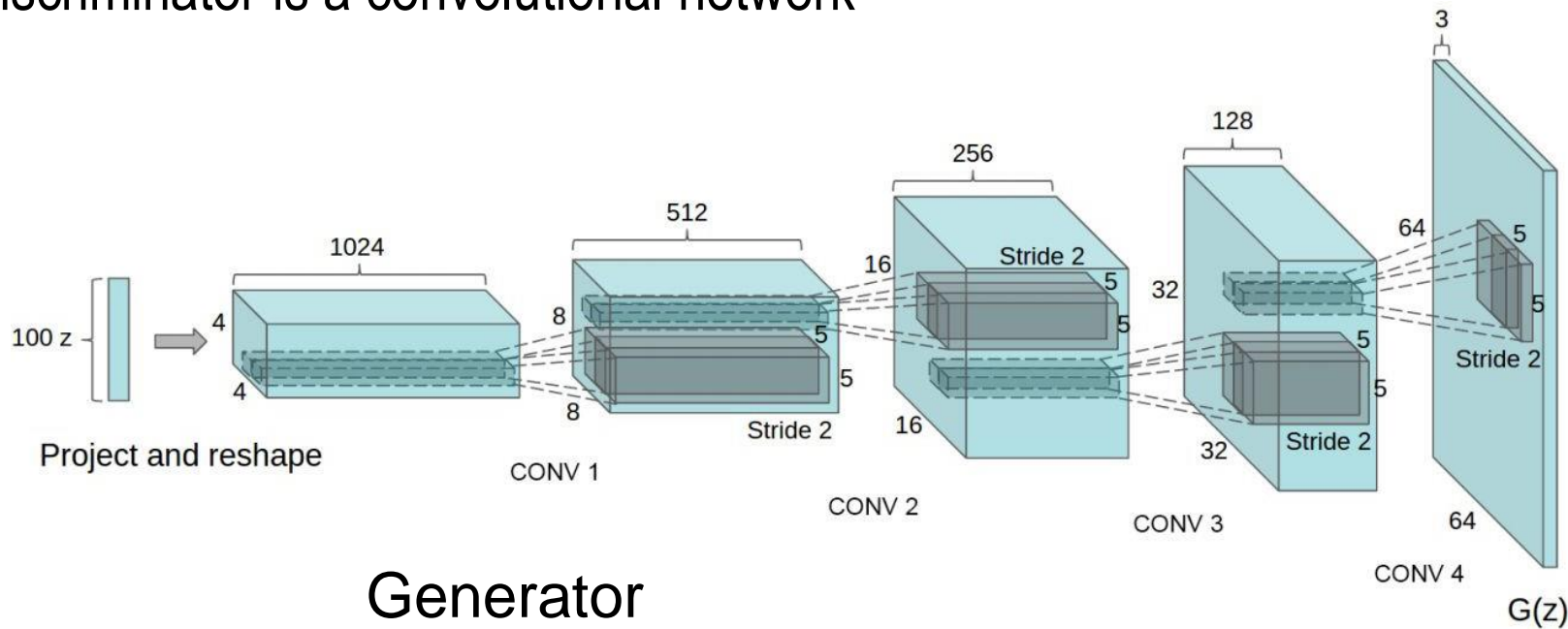Nearest neighbor from training set

Generated samples (CIFAR-10)



Nearest neighbor from training set

Generator is an upsampling network with fractionally-strided convolutions
Discriminator is a convolutional network



Generator

Samples from the model look much better!



Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

Smiling woman    Neutral woman    Neutral man

Samples from the model

Average Z vectors, do arithmetic

Smiling Man

Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

- Don't work with an explicit density function
- Take game-theoretic approach: learn to generate from training distribution through 2-player game
- Pros:
  - Beautiful, state-of-the-art samples!
- Cons:
  - Trickier / more unstable to train
- Active areas of research:
  - Better loss functions, more stable training (Wasserstein GAN, LSGAN, many others)
  - Conditional GANs, GANs for all kinds of applications