

Separation between NN and kernel

- For approximation and optimization, neural network has no advantage over kernel. Why NN gives better performance: **generalization**.
- [Allen-Zhu and Li '20] Construct a class of functions \mathcal{F} such that $y = f(x)$ for some $f \in \mathcal{F}$:
 - no kernel is sample-efficient; *for all kernel*
 - Exists a neural network that is sample-efficient.

Separation between NN and kernel

Defn Kernel method is a linear method
with an embedding $\phi: \mathbb{R}^d \rightarrow \mathcal{H}$ Hilbert space
 \Rightarrow It turns an element $f \in \mathcal{H}$ into a
prediction function

$$y = \langle f, \phi(x) \rangle$$

The method uses n samples, $\{x_i\}_{i=1}^n$, $x \in \mathbb{R}^d$
observe $\{y_i\}_{i=1}^n$

$$f \in \text{span} \left(\phi(x_i) \right)_{i=1}^n, \quad i \in [n]$$

e.g. $\underset{f}{\text{argmin}} \quad \frac{1}{n} \sum_{i=1}^n (y_i - \langle \phi(x_i), f \rangle)^2 + \lambda \|f\|_{\mathcal{H}}^2$

Separation between NN and kernel

Thus : \exists a class of functions $\mathcal{C} \subseteq \{c: \mathbb{R}^d \rightarrow \mathbb{R}\}$
and a distribution μ over \mathbb{R}^d s.t.

i) \forall kernel method, $\forall c \in \mathcal{C}$
given $y_i = c(x_i)$

if $\mathbb{E}_{x \sim \mu} [(c(x) - \langle f, \phi(x) \rangle)^2] \leq \frac{1}{9}$
then $n \geq 2^{d-1}$

ii) \exists simple procedure s.t. it can output
the true c as long as $n \geq d$
this procedure can be simulated/approximated
by NN + GD

Separation between NN and kernel

Pf: M : unit of $\{-1, 1\}^d$, 2^d elements
 $\mathcal{C} = \{C_S(x) = \prod_{i \in S} x_i, S \subset \{1, \dots, d\}\}$

part ii) choose a basis: $\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ \vdots \\ 1 \end{pmatrix}, \dots, \begin{pmatrix} 1 \\ \vdots \\ -1 \end{pmatrix}$
 e_1, e_2, \dots, e_d

observe $y_i = C(e_i)$, if $i \in S \Rightarrow y_i = 1$
& $i \notin S \Rightarrow y_i = -1$

\Rightarrow know whether i is in S or not

\Rightarrow identify $S \Rightarrow C_S$

$$\mathcal{C}(d) = \{1, \dots, d\}$$

Separation between NN and kernel

Part i) ϕ is a basis for $\{f: \{-1, 1\}^d \rightarrow \mathbb{R}\}$
with distribution

$$\mathbb{E}_{X \sim \mu} [\phi_S(X) \cdot \phi_{S'}(X)] = \begin{cases} 0 & \text{if } S \neq S' \\ 1 & \text{if } S = S' \end{cases} \quad \text{or symmetry}$$

Goal: to compute

$$\mathbb{E}_{X \sim \mu} [(C_{S^*}(X) - \langle f, \phi(X) \rangle)^2]$$

since $f \in \text{span}(\phi(x_i))_{i=1}^n$

$$\Rightarrow f = \sum_{i=1}^n a_i \phi(x_i)$$

$$x \mapsto \sum_{S \in \mathcal{C}(d)} \lambda_{i,S} (\phi_S(x))$$

Separation between NN and kernel

$$\Rightarrow \mathbb{E}_{x \sim \mu} \left[(C_{S^*}(x) - \langle f, \Phi(x) \rangle)^2 \right]$$

$$= \mathbb{E}_{x \sim \mu} \left[(C_{S^*}(x) - \sum_{S \in \mathcal{A}} \sum_{i=1}^n \alpha_i \lambda_{i,S} C_S(x))^2 \right]$$

$$= \left(1 - \sum_{i=1}^n \alpha_i \lambda_{i,S^*} \right)^2 + \sum_{S \neq S^*} \left(\sum_i \alpha_i \lambda_{i,S} \right)^2$$

by assumption error $\leq \frac{1}{9}$

$$\Rightarrow \left(1 - \sum_{i=1}^n \alpha_i \lambda_{i,S^*} \right)^2 \leq \frac{1}{9}$$

$$\& \sum_{S \neq S^*} \left(\sum_i \alpha_i \lambda_{i,S} \right)^2 \leq \frac{1}{9}$$

$\Rightarrow n \geq 2^{d-1}$

Notations: $\Lambda : 2^d \times n$ ($n \leq 2^d$)

$$\Lambda_{s,i} = \lambda_{i,s}$$

$$A : n \times 2^d$$

$$A_{i,s^*} = a_{i,s^*}$$

$$\Omega = \Lambda A : 2^d \times 2^d \text{ of rank-} n$$

$$(1) \quad \underline{(1 - \Omega_{s^*, s^*})^2 \leq \frac{1}{9} \Rightarrow \Omega_{s^*, s^*} \geq \frac{2}{3}}$$

$$\sum_{s \neq s^*} \Omega_{s, s^*}^2 \leq \frac{1}{9}$$

$$2^d \quad \left(\text{Diagram showing a large circle with two smaller circles inside, one labeled } \leq \frac{1}{9} \text{ and the other } \geq \frac{2}{3} \right)$$

$$\Omega = \text{diag}(\Omega) + \Omega' \quad \text{off-diagonal}$$

$$\|\Omega'\|_F^2 \leq \frac{2^d}{9}$$

$$= \sum \text{eig}(\Omega')^2$$

$$\Rightarrow \Omega' \text{ has at most } \frac{2^d}{4} \text{ eigenvalues } > \frac{2}{3}$$

\Rightarrow Consider subspace with eigenvalues $< \frac{2}{3}$, which has dimension at least $\frac{3}{4} \cdot 2^d$

$\forall x \in \text{this subspace}$

$$\begin{aligned} \|\Omega x\|_2 &= \|\text{diag}(\Omega)x + \Omega'x\|_2 \\ &\geq \|\text{diag}(\Omega)x\|_2 - \|\Omega'x\|_2 \\ &> \frac{2}{3}x - \frac{2}{3}x = 0 \end{aligned}$$

$$\Rightarrow \text{rank}(\Omega) \geq \frac{3}{4} \cdot 2^d \Rightarrow n \geq \frac{3}{4} \cdot 2^d \quad \square$$

Convolutional Neural Networks



Multi-layer Neural Network

$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

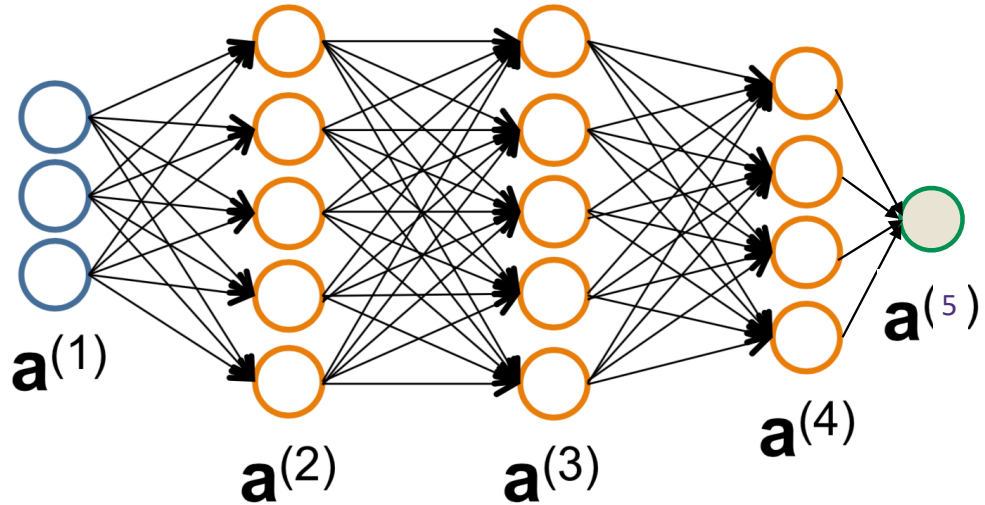
⋮

$$z^{(l+1)} = \Theta^{(l)} a^{(l)}$$

$$a^{(l+1)} = g(z^{(l+1)})$$

⋮

$$\hat{y} = a^{(L+1)}$$



$$L(y, \hat{y}) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

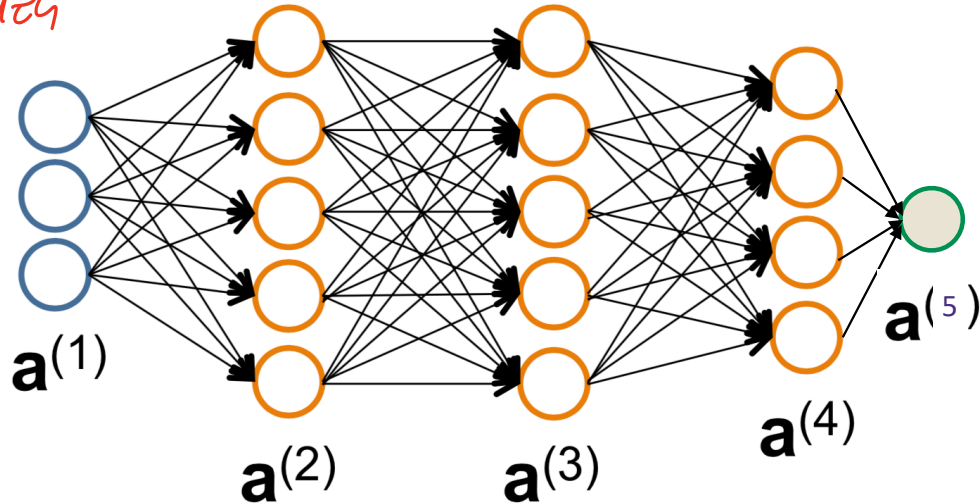
Binary
Logistic
Regression

Neural Network Architecture

The neural network architecture is defined by the number of layers, and the number of nodes in each layer, but also by allowable edges.

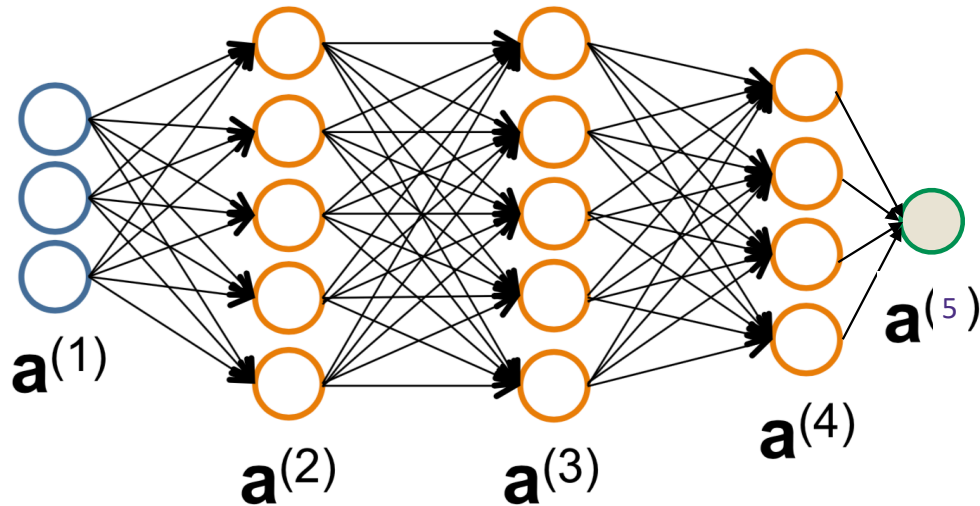
width

depth



Neural Network Architecture

The neural network architecture is defined by the number of layers, and the number of nodes in each layer, but also by **allowable edges**.



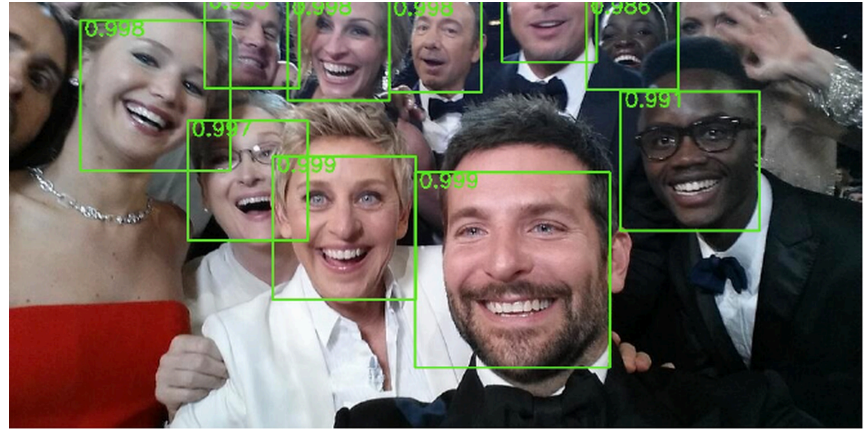
We say a layer is **Fully Connected (FC)** if all linear mappings from the current layer to the next layer are permissible.

$$\mathbf{a}^{(k+1)} = g(\Theta \mathbf{a}^{(k)}) \quad \text{for any } \Theta \in \mathbb{R}^{n_{k+1} \times n_k}$$

A lot of parameters!! $n_1 n_2 + n_2 n_3 + \cdots + n_L n_{L+1}$

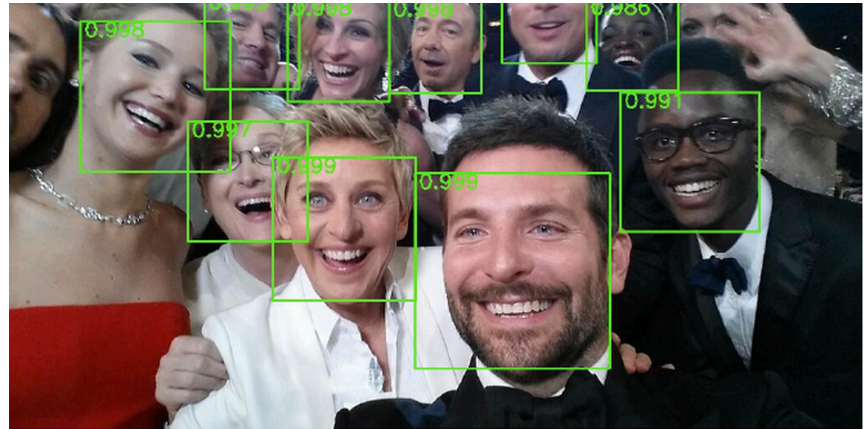
Neural Network Architecture

Objects are often **localized in space** so to find the faces in an image, not every pixel is important for classification—makes sense to drag a window across an image.

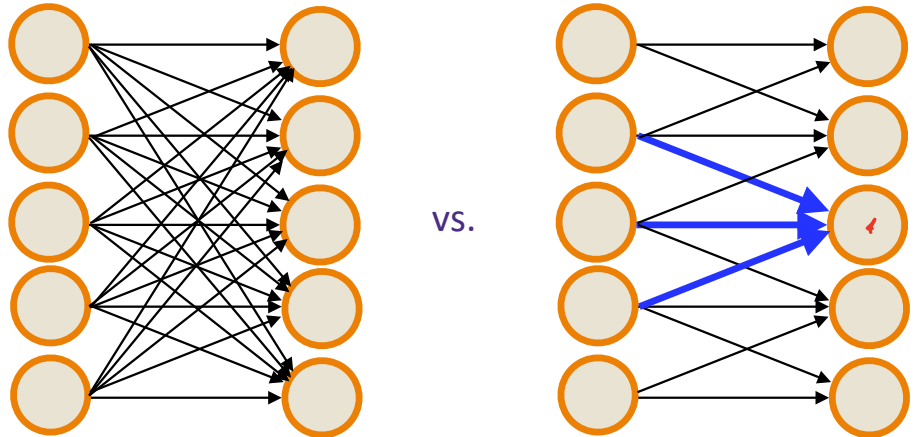


Neural Network Architecture

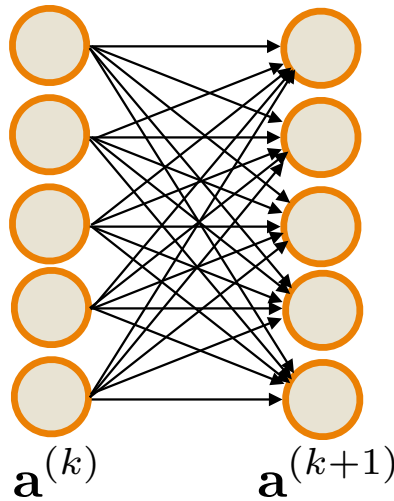
Objects are often **localized in space** so to find the faces in an image, not every pixel is important for classification—makes sense to drag a window across an image.



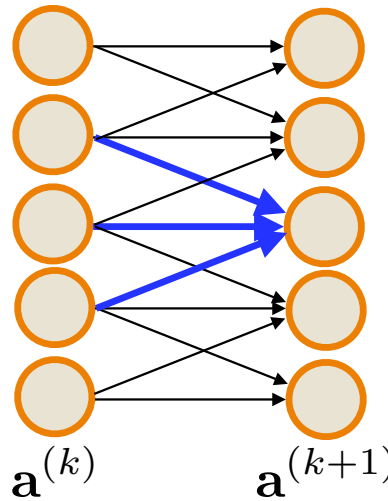
Similarly, to identify edges or other local structure, it makes sense to only look at **local information**



Neural Network Architecture



vs.



hardcode

$$\begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & \Theta_{0,2} & \Theta_{0,3} & \Theta_{0,4} \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & \Theta_{1,3} & \Theta_{1,4} \\ \Theta_{2,0} & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & \Theta_{2,4} \\ \Theta_{3,0} & \Theta_{3,1} & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ \Theta_{4,0} & \Theta_{4,1} & \Theta_{4,2} & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix}$$

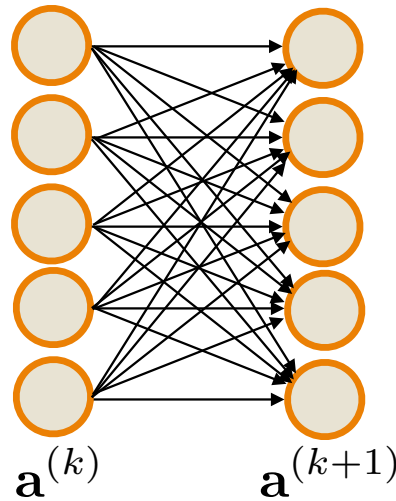
$$\begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & 0 & 0 & 0 \\ \Theta_{1,0} & \Theta_{1,1} & 0 & 0 & 0 \\ 0 & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & 0 \\ 0 & 0 & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ 0 & 0 & 0 & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix}$$

Parameters: n^2

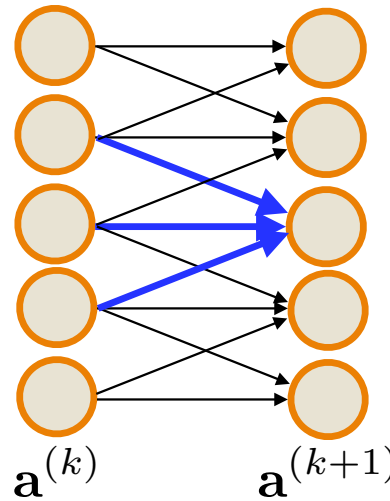
$3n - 2$

$$\mathbf{a}_i^{(k+1)} = g \left(\sum_{j=0}^{n-1} \Theta_{i,j} \mathbf{a}_j^{(k)} \right)$$

Neural Network Architecture



vs.



Mirror/share local weights everywhere
(e.g., structure equally likely to be anywhere in image)

$$\begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & \Theta_{0,2} & \Theta_{0,3} & \Theta_{0,4} \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & \Theta_{1,3} & \Theta_{1,4} \\ \Theta_{2,0} & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & \Theta_{2,4} \\ \Theta_{3,0} & \Theta_{3,1} & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ \Theta_{4,0} & \Theta_{4,1} & \Theta_{4,2} & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix}$$

Parameters: n^2

$$\begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & 0 & 0 & 0 \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & 0 & 0 \\ 0 & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & 0 \\ 0 & 0 & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ 0 & 0 & 0 & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix}$$

$3n - 2$

$$\begin{bmatrix} \theta_1 & \theta_2 & 0 & 0 & 0 \\ \theta_0 & \theta_1 & \theta_2 & 0 & 0 \\ 0 & \theta_0 & \theta_1 & \theta_2 & 0 \\ 0 & 0 & \theta_0 & \theta_1 & \theta_2 \\ 0 & 0 & 0 & \theta_0 & \theta_1 \end{bmatrix}$$

3

$$\mathbf{a}_i^{(k+1)} = g \left(\sum_{j=0}^{n-1} \Theta_{i,j} \mathbf{a}_j^{(k)} \right)$$

$$\mathbf{a}_i^{(k+1)} = g \left(\sum_{j=0}^{m-1} \theta_j \mathbf{a}_{i+j}^{(k)} \right)$$

Neural Network Architecture

Fully Connected (FC) Layer

$$\begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & \Theta_{0,2} & \Theta_{0,3} & \Theta_{0,4} \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & \Theta_{1,3} & \Theta_{1,4} \\ \Theta_{2,0} & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & \Theta_{2,4} \\ \Theta_{3,0} & \Theta_{3,1} & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ \Theta_{4,0} & \Theta_{4,1} & \Theta_{4,2} & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix}$$

$$\mathbf{a}_i^{(k+1)} = g \left(\sum_{j=0}^{n-1} \Theta_{i,j} \mathbf{a}_j^{(k)} \right)$$

Convolutional (CONV) Layer (1 filter)

$$\begin{bmatrix} \theta_1 & \theta_2 & 0 & 0 & 0 \\ \theta_0 & \theta_1 & \theta_2 & 0 & 0 \\ 0 & \theta_0 & \theta_1 & \theta_2 & 0 \\ 0 & 0 & \theta_0 & \theta_1 & \theta_2 \\ 0 & 0 & 0 & \theta_0 & \theta_1 \end{bmatrix} \quad m=3$$

$$\mathbf{a}_i^{(k+1)} = g \left(\sum_{j=0}^{m-1} \theta_j \mathbf{a}_{i+j}^{(k)} \right) = g([\theta * \mathbf{a}^{(k)}]_i)$$

Convolution*

$\theta = (\theta_0, \dots, \theta_{m-1}) \in \mathbb{R}^m$ is referred to as a “filter”

Example (1d convolution)

$$(\theta * x)_i = \sum_{j=0}^{m-1} \theta_j x_{i+j}$$

stride = 1

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
|---|---|---|---|---|

Input $x \in \mathbb{R}^n$

| | | |
|---|---|---|
| 1 | 0 | 1 |
|---|---|---|

Filter $\theta \in \mathbb{R}^m$

| | | |
|--|--|--|
| | | |
|--|--|--|

Output $\theta * x$

$\in \mathbb{R}^{n-2}$

Example (1d convolution)

$$(\theta * x)_i = \sum_{j=0}^{m-1} \theta_j x_{i+j}$$

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
|---|---|---|---|---|

Input $x \in \mathbb{R}^n$

| | | |
|---|---|---|
| 1 | 0 | 1 |
|---|---|---|

Filter $\theta \in \mathbb{R}^m$

| | | | | |
|------------------------|------------------------|------------------------|---|---|
| 1 <small>x1</small> | 1 <small>x0</small> | 1 <small>x1</small> | 0 | 0 |
|------------------------|------------------------|------------------------|---|---|

| | | |
|---|--|--|
| 2 | | |
|---|--|--|

Output $\theta * x$

Example (1d convolution)

$$(\theta * x)_i = \sum_{j=0}^{m-1} \theta_j x_{i+j}$$

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
|---|---|---|---|---|

Input $x \in \mathbb{R}^n$

| | | |
|---|---|---|
| 1 | 0 | 1 |
|---|---|---|

Filter $\theta \in \mathbb{R}^m$

| | | | | |
|---|-----------------|-----------------|-----------------|---|
| 1 | 1 _{x1} | 1 _{x0} | 0 _{x1} | 0 |
|---|-----------------|-----------------|-----------------|---|

| | | |
|---|---|--|
| 2 | 1 | |
|---|---|--|

Output $\theta * x$

Example (1d convolution)

$$(\theta * x)_i = \sum_{j=0}^{m-1} \theta_j x_{i+j}$$

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
|---|---|---|---|---|

Input $x \in \mathbb{R}^n$

| | | |
|---|---|---|
| 1 | 0 | 1 |
|---|---|---|

Filter $\theta \in \mathbb{R}^m$

| | | | | |
|---|---|-----------------|-----------------|-----------------|
| 1 | 1 | 1 _{x1} | 0 _{x0} | 0 _{x1} |
|---|---|-----------------|-----------------|-----------------|

| | | | | |
|---|---|---|---|---|
| 1 | 2 | 1 | 1 | 0 |
|---|---|---|---|---|

Output $\theta * x$

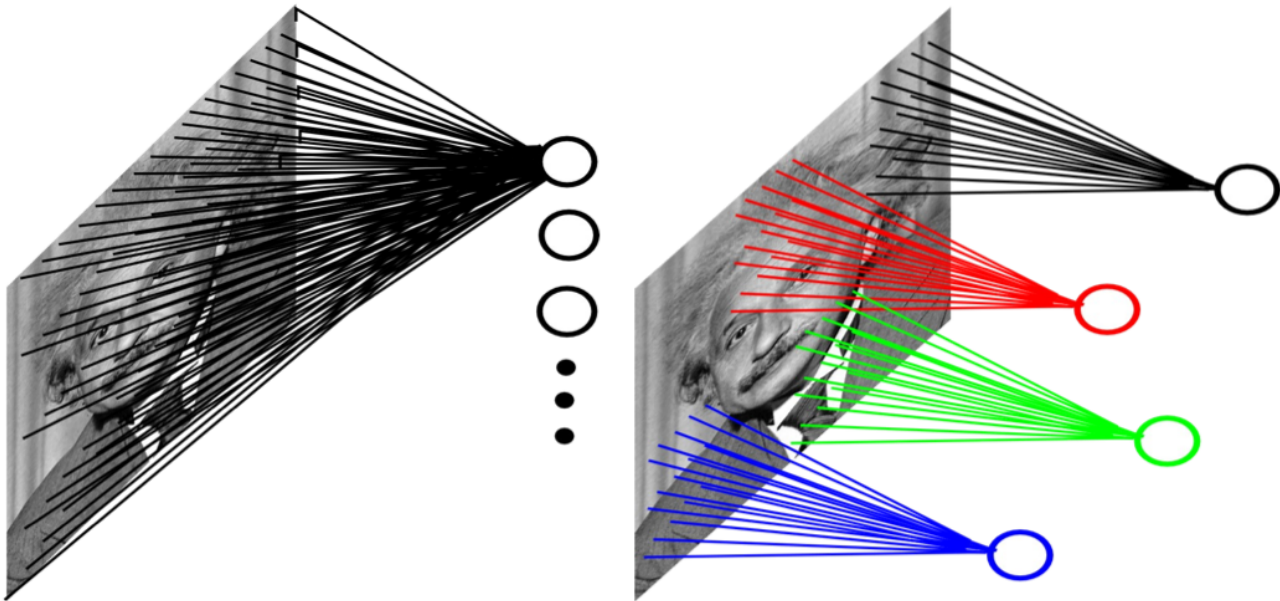
padding

| | | | | | | |
|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|

2d Convolution Layer

■ Example: 200x200 image

- ▶ Fully-connected, 400,000 hidden units = 16 billion parameters
- ▶ Locally-connected, 400,000 hidden units 10x10 fields = 40 million params
- ▶ Local connections capture local dependencies



Convolution of images (2d convolution)

$$(I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n)$$

4

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Image I

| | | |
|---|---|---|
| 1 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 1 |

Filter K

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Image

| | | |
|---|--|--|
| 4 | | |
| | | |
| | | |

Convolved
Feature

$$I * K$$

$$(n-2), (n-2)$$

Convolution of images

$$(I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n)$$

Image I

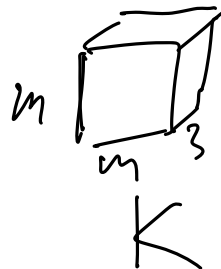
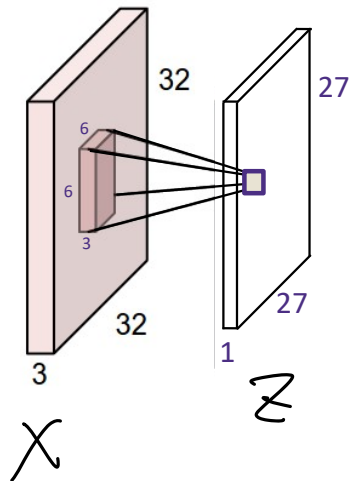


hand
crafted

NN: learned

| Operation | Filter K | Convolved Image $I * K$ |
|----------------------------------|----------------------------------------------------------------------------------|-------------------------|
| Edge detection | $\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$ | |
| | $\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$ | |
| | $\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$ | |
| Sharpen | $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ | |
| Box blur (normalized) | $\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ | |
| Gaussian blur (approximation) | $\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$ | |

Stacking convolved images



$$r=3$$

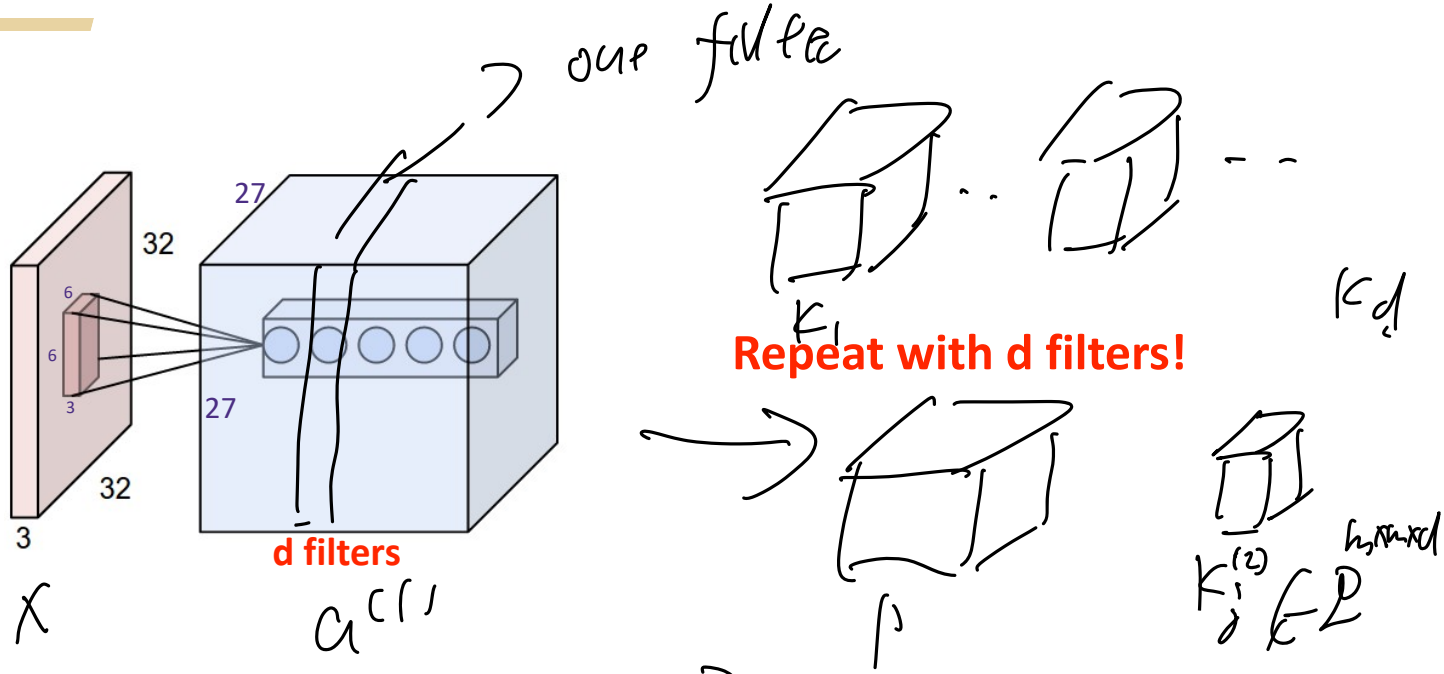
c_{channel}

$$x \in \mathbb{R}^{n \times n \times \underline{r}}$$

$\mathcal{D} \mathcal{G} \mathcal{B}$

$$Z = \sum_{\alpha=1}^r X[:, :, \alpha] * K[:, :, \alpha]$$

Stacking convolved images

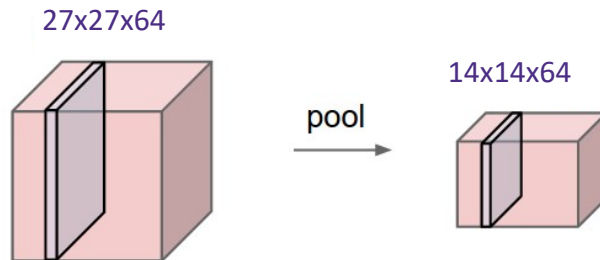
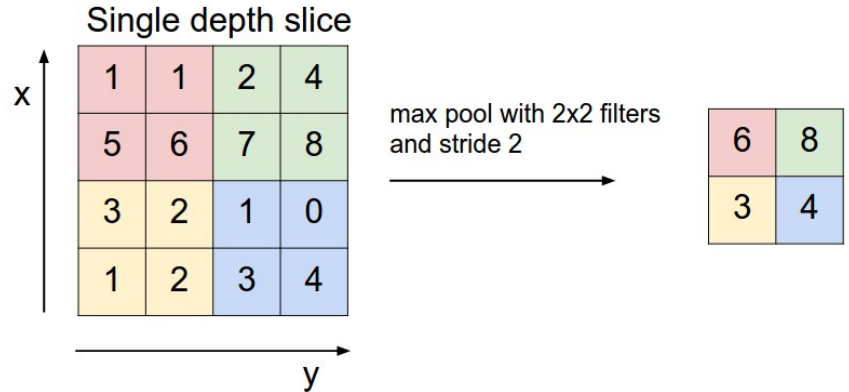


$$A_{\vec{i}}^{(1)} = \gamma \left(\sum_{\alpha=1}^3 X[\vec{i}, : \alpha] * K^{(1)}[\vec{i}, : \alpha] \right)$$

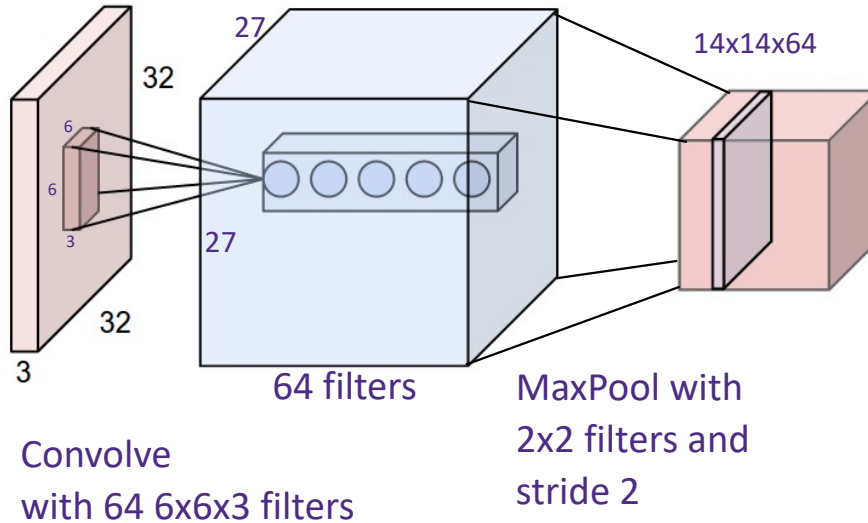
$$A_{\vec{j}}^{(2)} = \gamma \left(\sum_{\vec{i}=1}^d A_{\vec{i}}^{(1)}[\vec{i}, : \vec{j}] * K^{(2)}[\vec{i}, : \vec{j}] \right)$$

Pooling

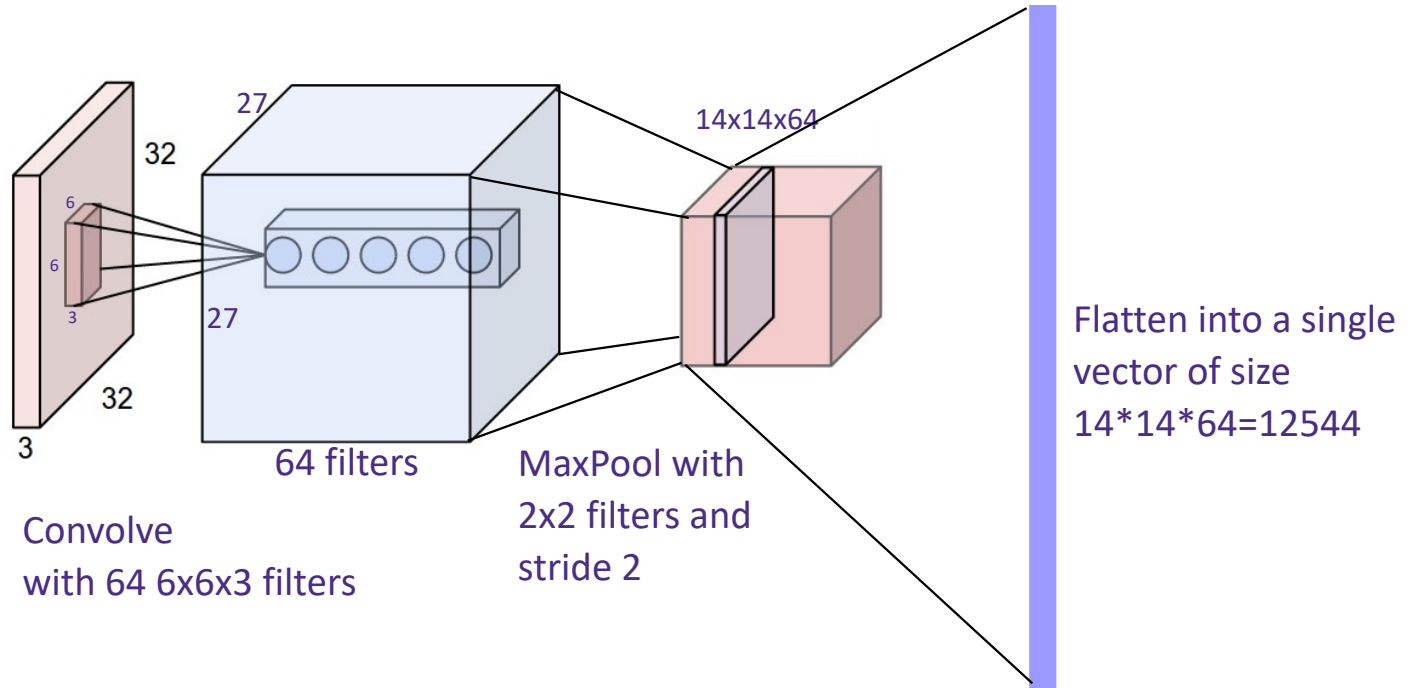
Pooling reduces the dimension and can be interpreted as “This filter had a high response in this general region”



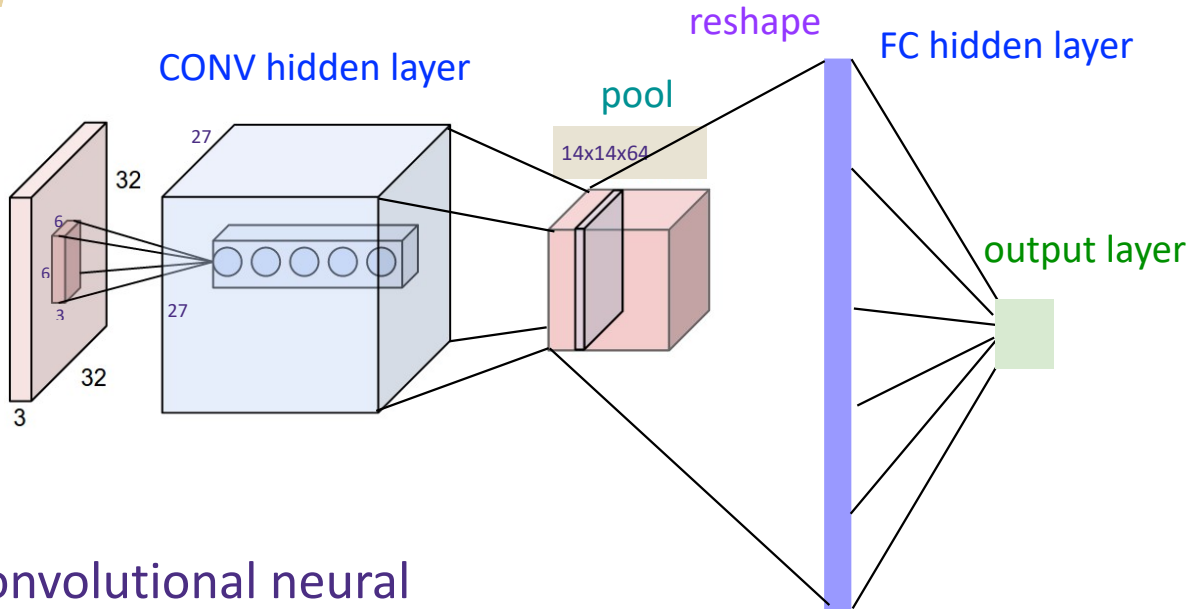
Pooling Convolution layer



Flattening

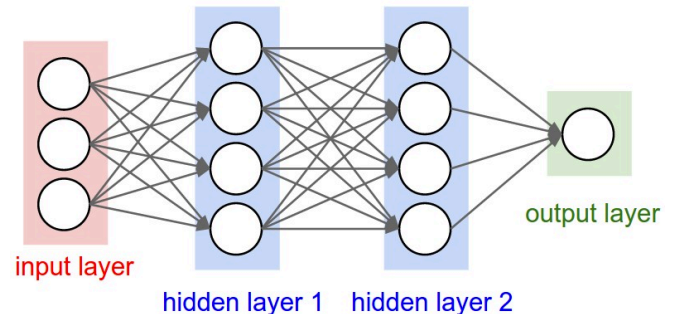


Training Convolutional Networks

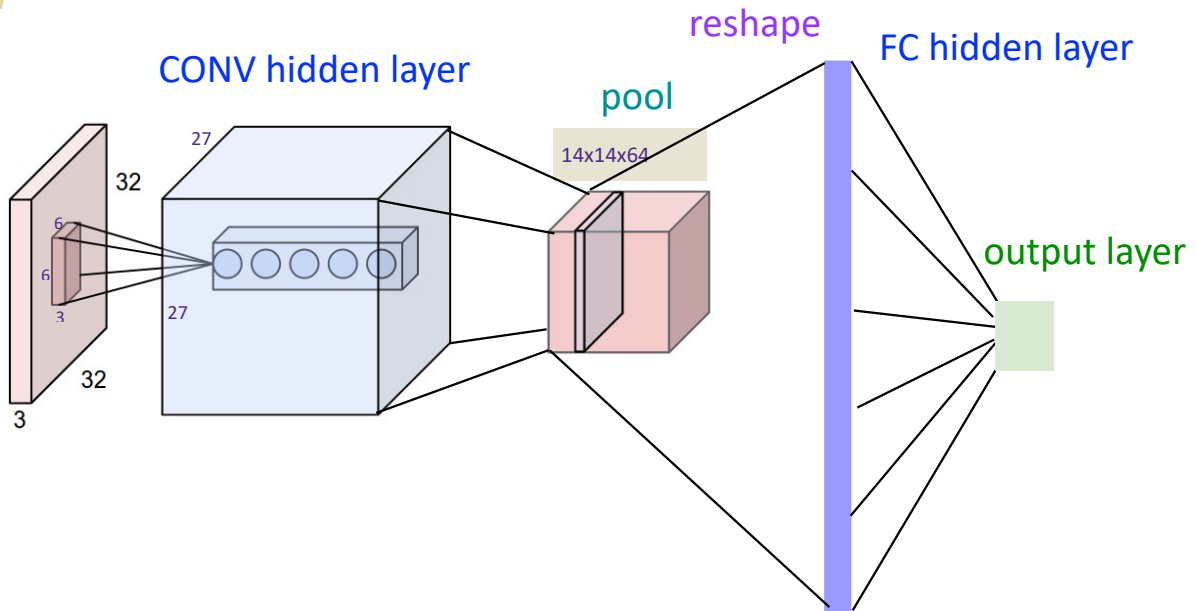


Recall: Convolutional neural networks (CNN) are just regular fully connected (FC) neural networks with some connections removed.

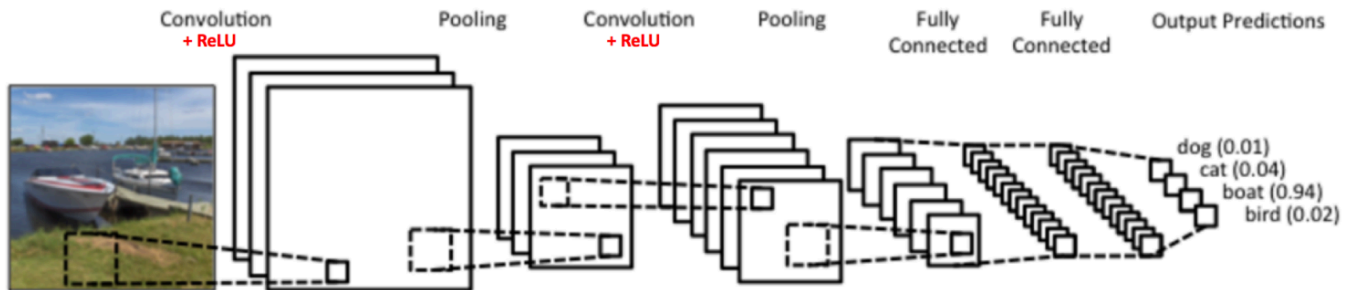
Train with SGD!

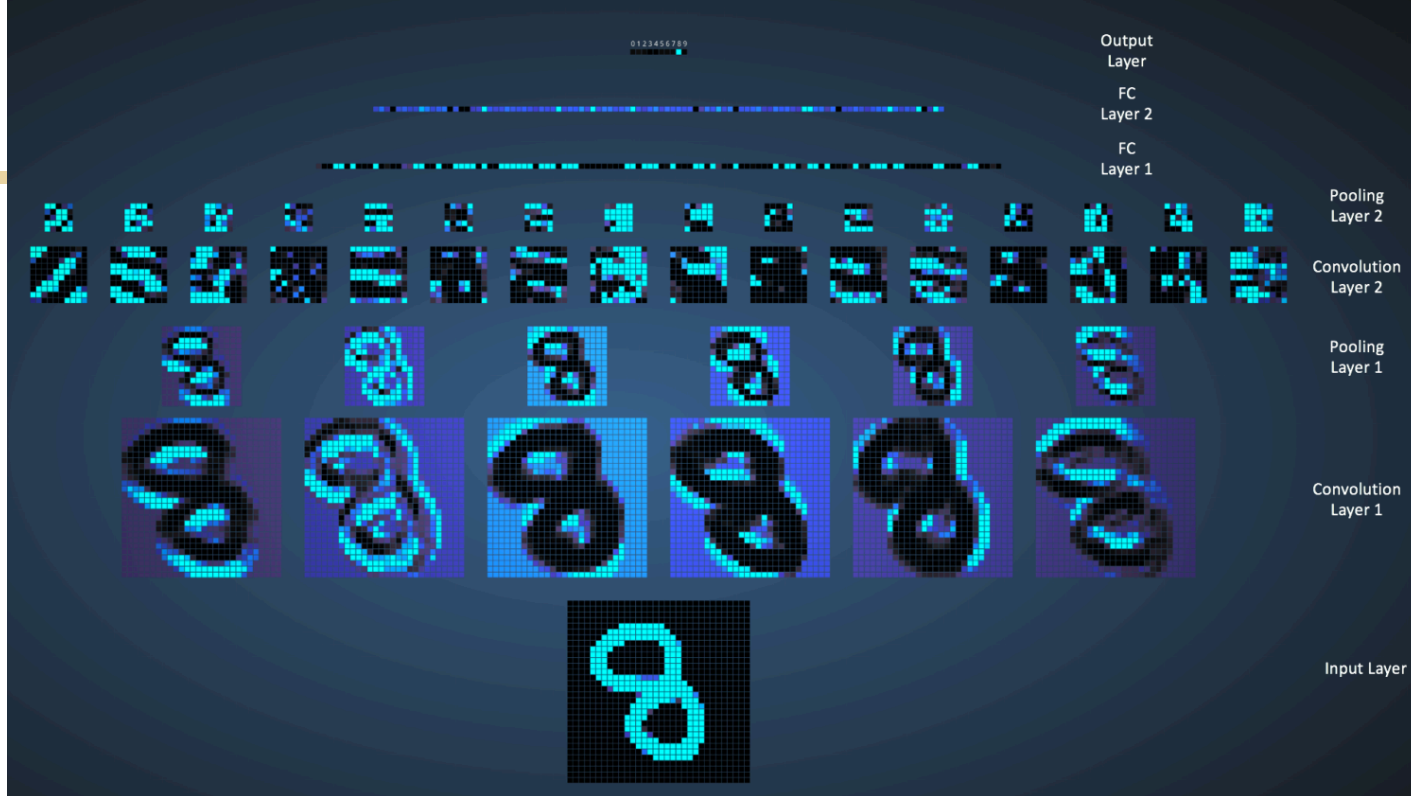


Training Convolutional Networks

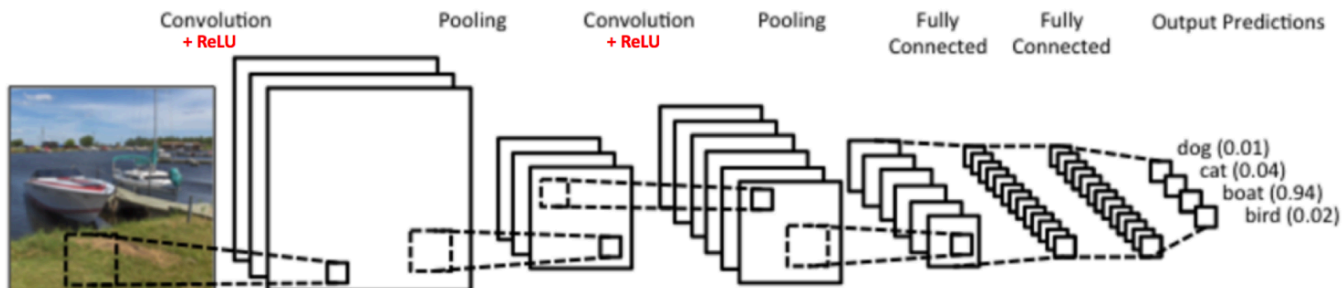


Real example network: LeNet





Real example network: LeNet



Famous CNNs



ImageNet Dataset

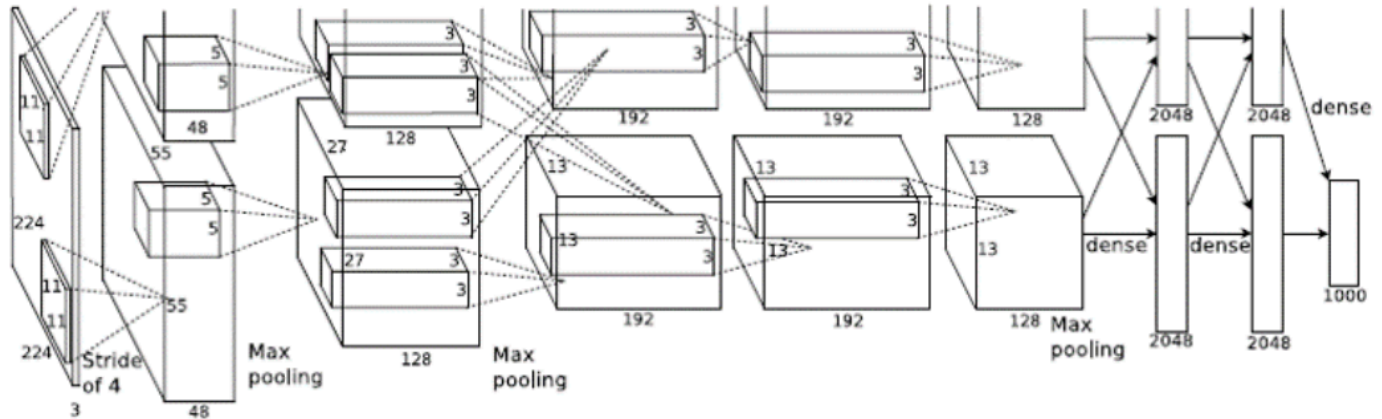
~14 million images, 20k classes



Deng et al. "Imagenet: a large scale hierarchical image database" '09

AlexNet

Breakthrough on ImageNet: ~the beginning of deep learning era



Krizhevsky, Sutskever, Hinton “ImageNet Classification with Deep Convolutional Neural Networks”, NIPS 2012.

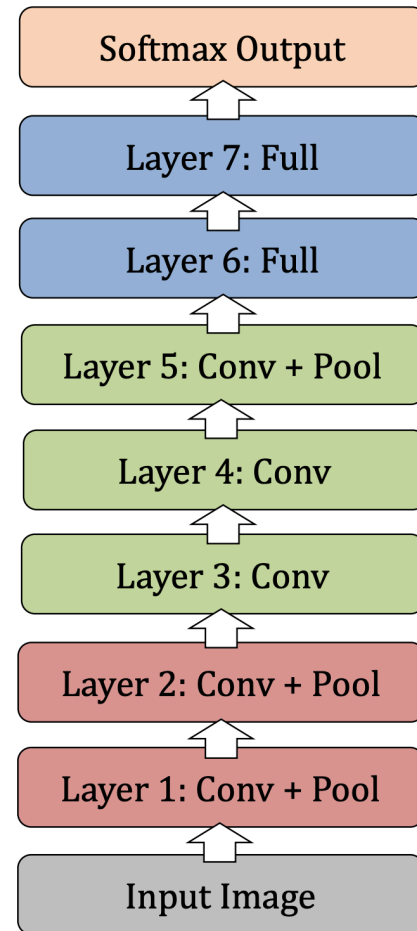
AlexNet

8 layers, ~60M parameters

Top5 error: 18.2%

Techniques used:

ReLU activation, overlapping pooling, dropout, ensemble (create 10 patches by cropping and average the predictions), data-augmentation (intensity of RGB channels)



[From Rob Fergus' CIFAR 2016 tutorial]

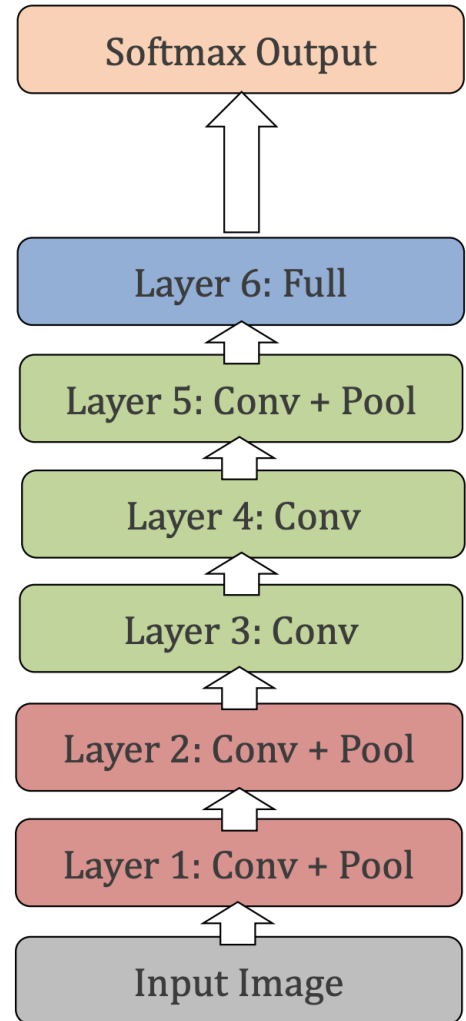
AlexNet

Remove top fully-connected layer 7

Drop ~16 million parameters

1.1% drop in performance

[From Rob Fergus' CIFAR 2016 tutorial]



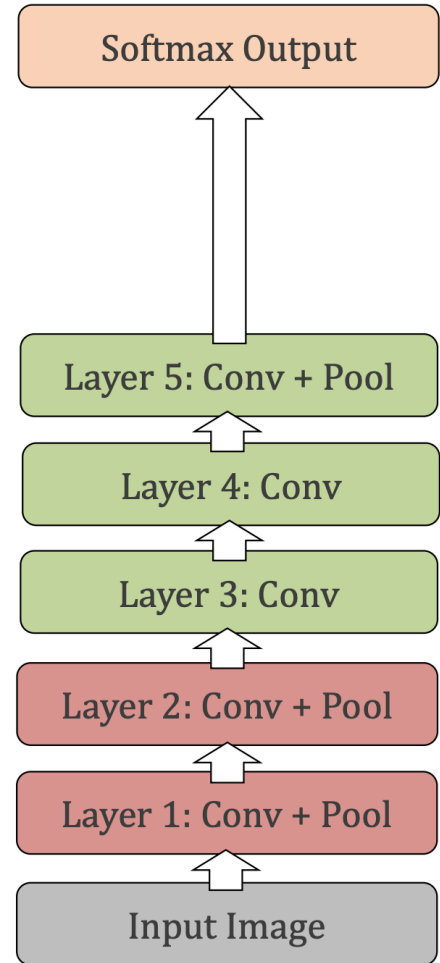
AlexNet

Remove both fully connected
layers 6 and 7

Drop ~50 million parameters

5.7% drop in performance

[From Rob Fergus' CIFAR 2016 tutorial]

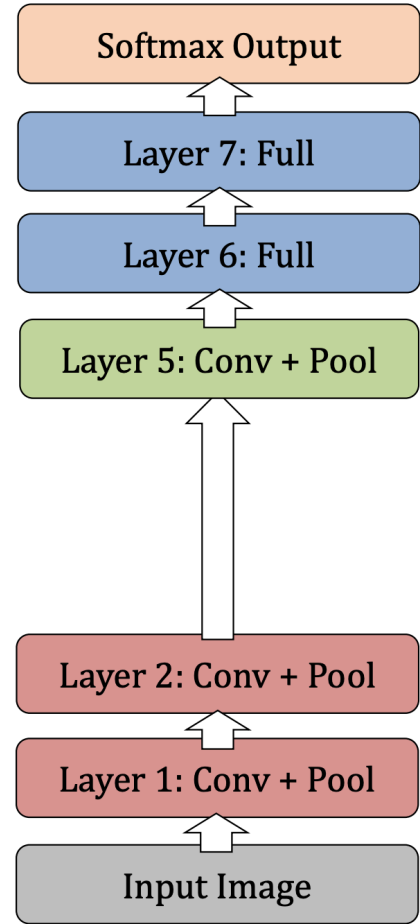


AlexNet

Remove upper convolutio / feature extractor layers (layer 3 and 4)

Drop ~1 million parameters

3% drop in performance



[From Rob Fergus' CIFAR 2016 tutorial]

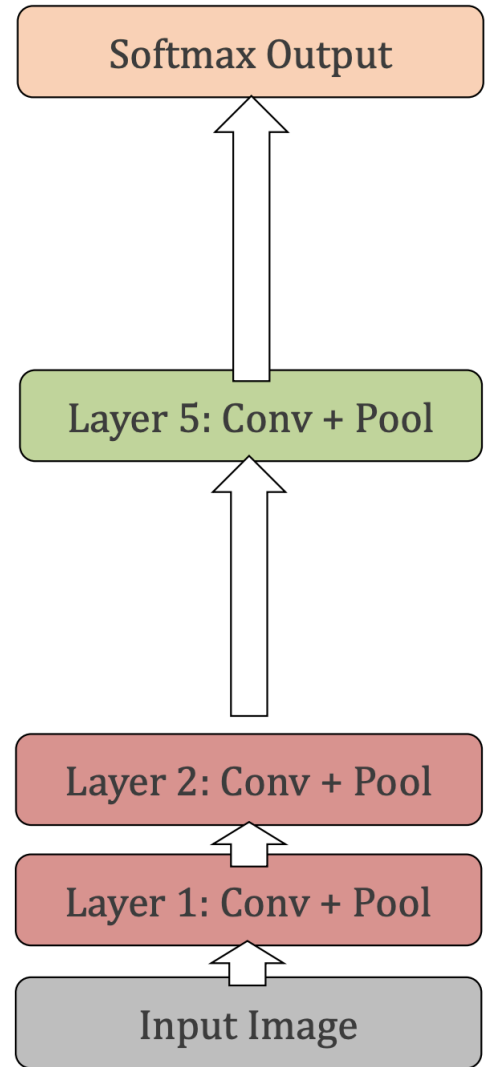
AlexNet

Remove top fully connected layer
6,7 and upper convolution layers
3,4.

33.5% drop in performance.

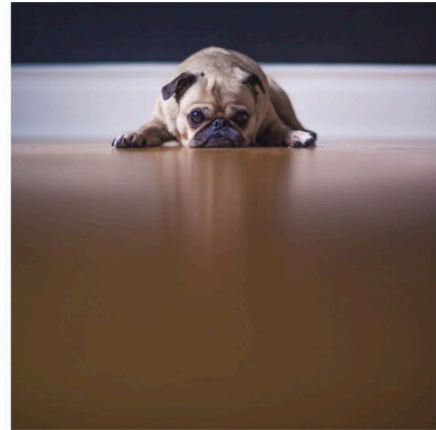
Depth of the network is the key.

[From Rob Fergus' CIFAR 2016 tutorial]



GoogLeNet

Motivation: multiscale nature of images

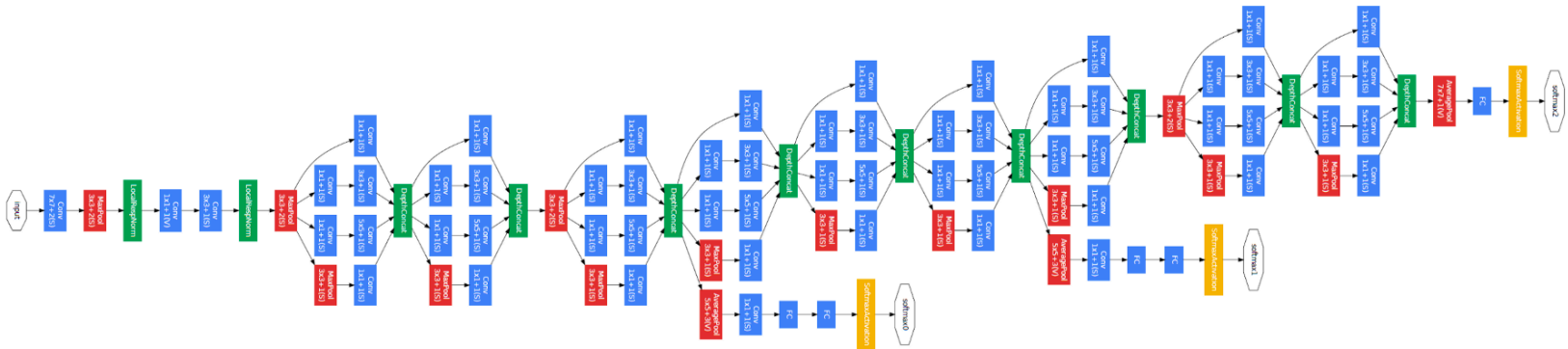


Large kernel for global features, and **smaller kernel** for local features.

Idea: have multiple different-size kernels at any layer.

[Going Deep with Convolutions, Szegedy et al. '14]

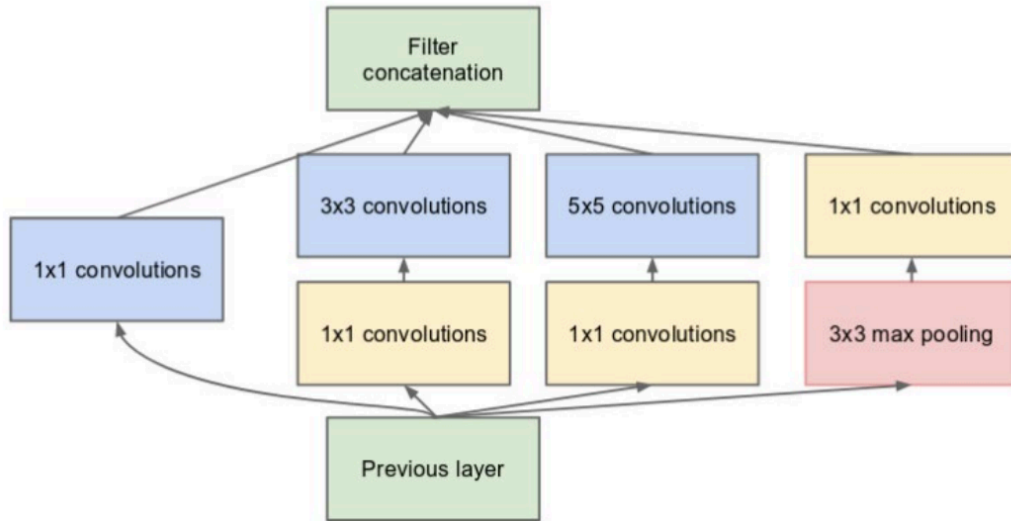
GoogLeNet



Large kernel for global features, and **smaller kernel** for local features.

Idea: have multiple different-size kernels at any layer.

Inception Module



Multiple filter scales at each layer

Dimensionality reduction to keep computational requirements down

[Going Deep with Convolutions, Szegedy et al. '14]

Residual Networks

Motivation: extremely deep nets are hard to train (gradient explosion/vanishing)

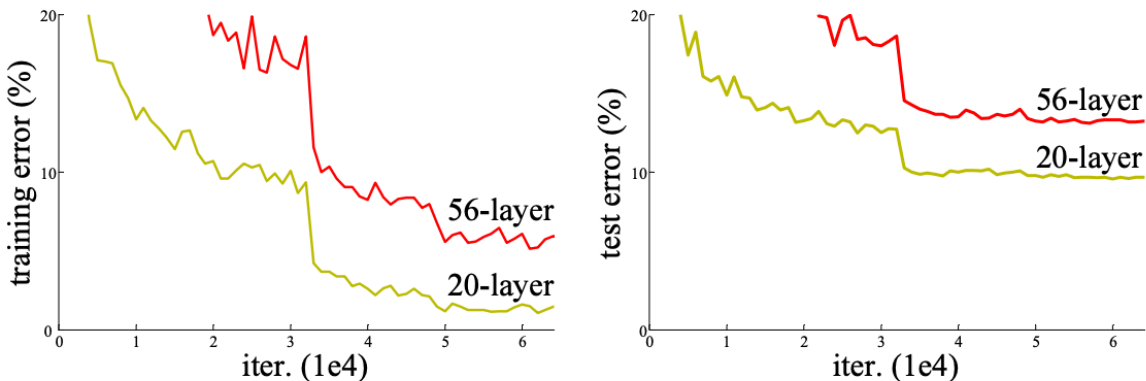
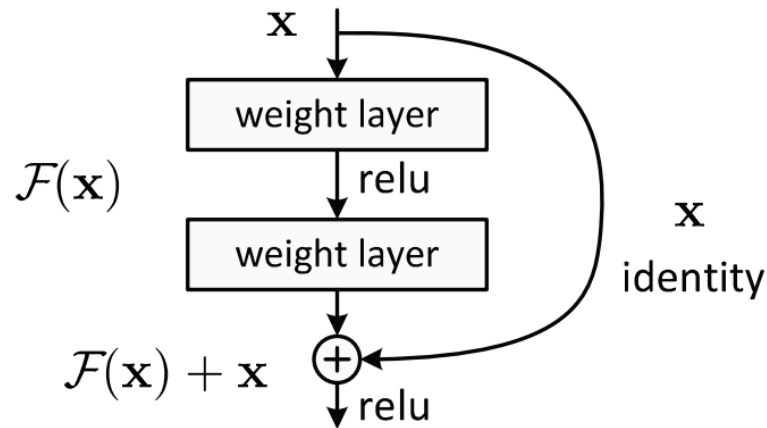


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

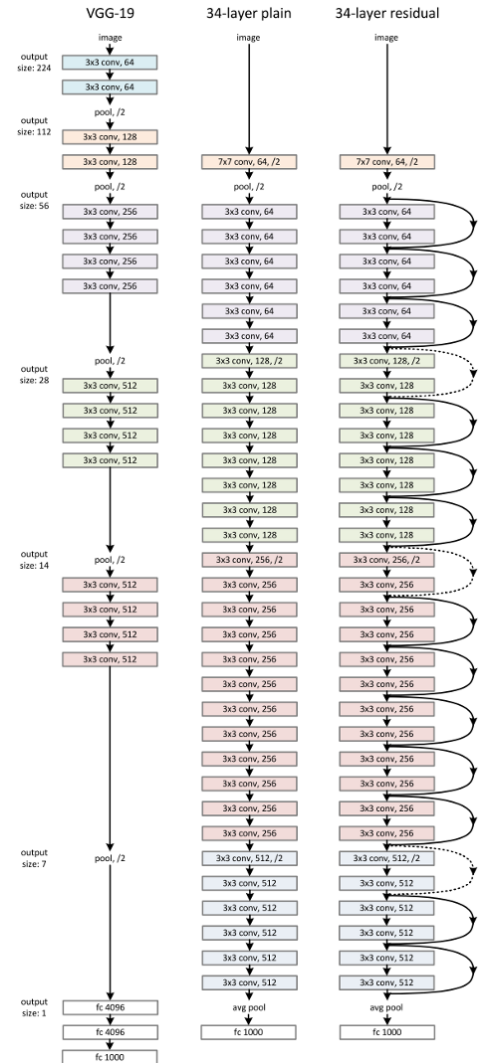
Residual Networks

Idea: identity shortcut, skip one or more layers.

Justification: network can easily simulate shallow network ($F \approx 0$), so performance should not degrade by going deeper.



- [He, Zhang, Ren, Sun, '16]

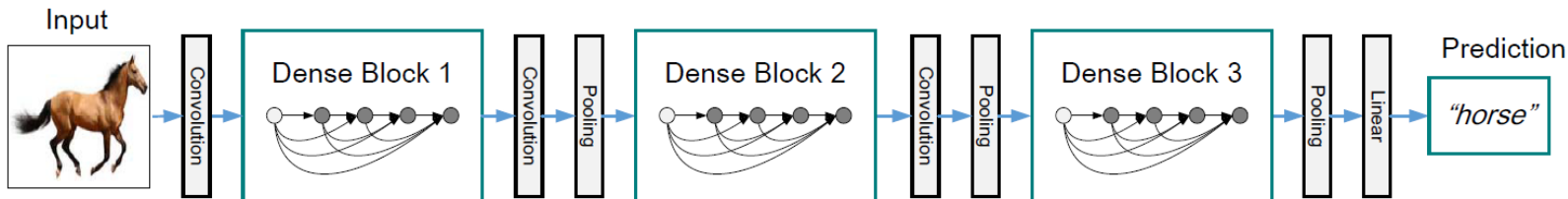
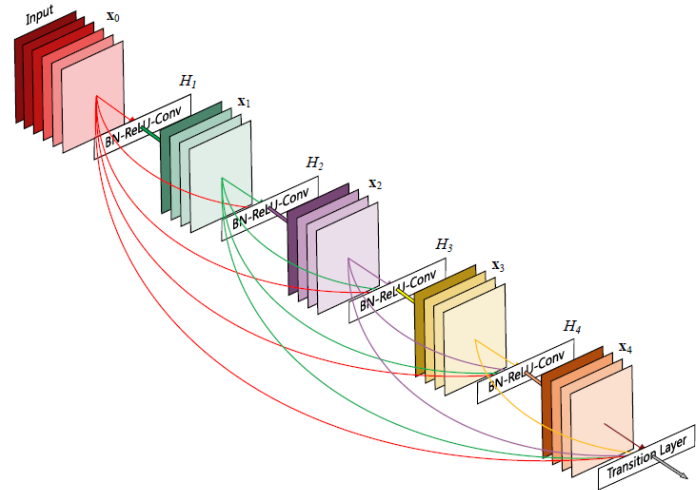


Densely Connected Network

Idea: explicit forward output of layer to all future layers (by concatenation)

Intuition: helps vanishing gradients, encourage reuse features (reduce parameter count)

Issues: network maybe too wide, need to be careful about memory consumption



[He, Zhang, Ren, Sun, '16]

Neural Architecture / Hyper-Parameter Search

Many design choices:

- Number of layers, width, kernel size, pooling, connections, etc.
- Normalization, learning rate, batch size, etc.

Strategies:

- Grid search
- Random search [Bergstra & Bengio '12]
- Bandit-based [Li et al. '16]
- Gradient-based (DARTS) [Liu et al. '19]
- Neural tangent kernel [Xu et al. '21]
- ...