# Deep Learning and Applications

DSA 5204 • Lecture 4
Dr Low Yi Rui (Aaron)
Department of Mathematics

**NUS**
National University
of Singapore

# Last Time

We discussed deep fully connected neural networks

$$h^{(i+1)} = \text{ReLU}\big(W^{(i+1)}h^{(i)} + b^{(i+1)}\big), \qquad h^{(0)} = x$$
$$\hat{y} = w^T h^{(\ell)} + c$$

They do not generally perform better than other ML methods, even with abundant data

Today, we will look at deep convolutional neural networks, which is arguably what put deep learning ahead in modern ML.

# Demo: Permutation Invariance of Fully-Connected NNs

# Permutations

A **permutation** of $n$ objects is a one-to-one transformation on these objects.

In other words, it is a **bijection** on $\{1, 2, \ldots, n\}$

$$\{1,2,3,4,5,6,7,8,9\}$$

$$p \qquad p(1) = 3, p(2) = 9 \ldots$$

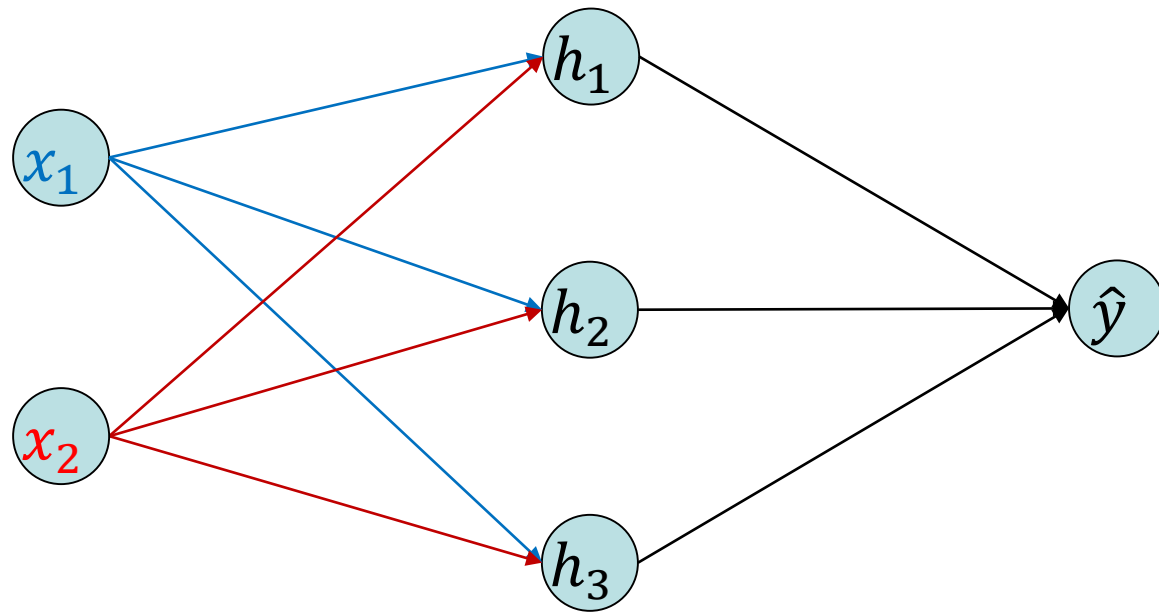$$\{3,9,1,5,4,6,8,7,2\}$$
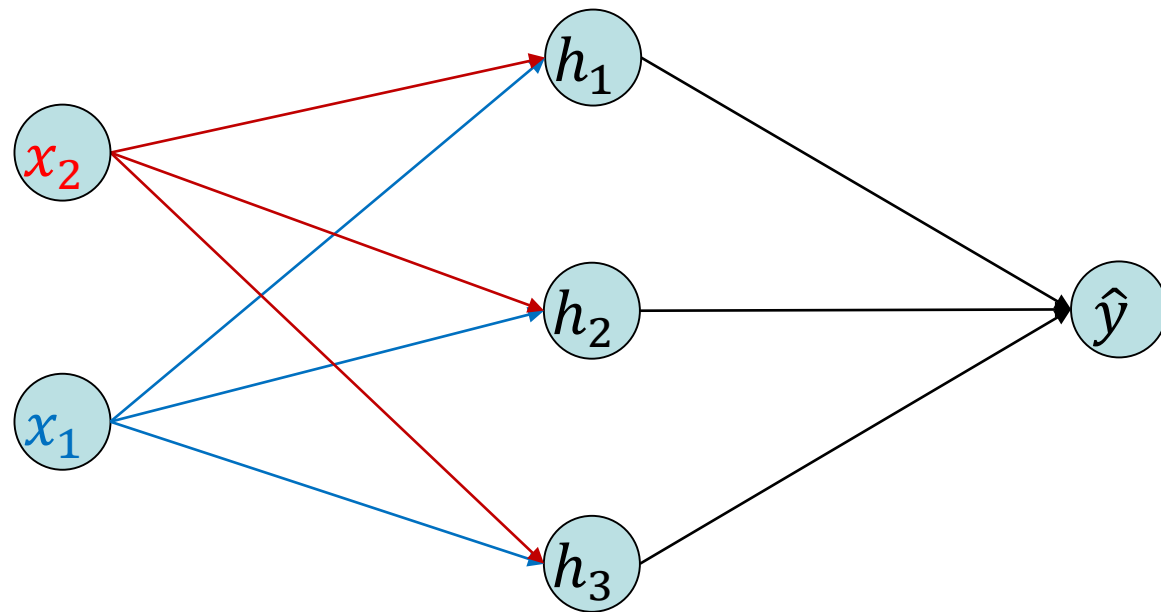
# Permutation Invariance of Hypothesis Spaces

Observe that the FCNN hypothesis space $\mathcal{H}$ has the following invariance property:

*Suppose $f \in \mathcal{H}$, then for any permutation $p$ on the indices, define $f_p(x_1, \dots, x_d) \equiv f\left(x_{p(1)}, \dots, x_{p(d)}\right)$, then the function $f_p \in \mathcal{H}$ as well*
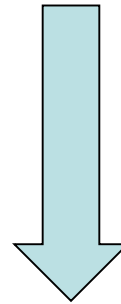
# FCNN and Permutation Invariance

FCNN does not care about permuting the signal's components since if we can fit one permutation, we can fit any permutation!

Is this sensible?

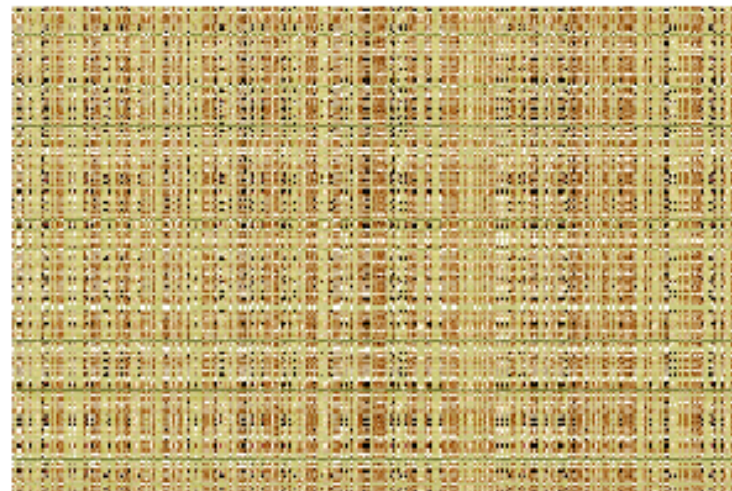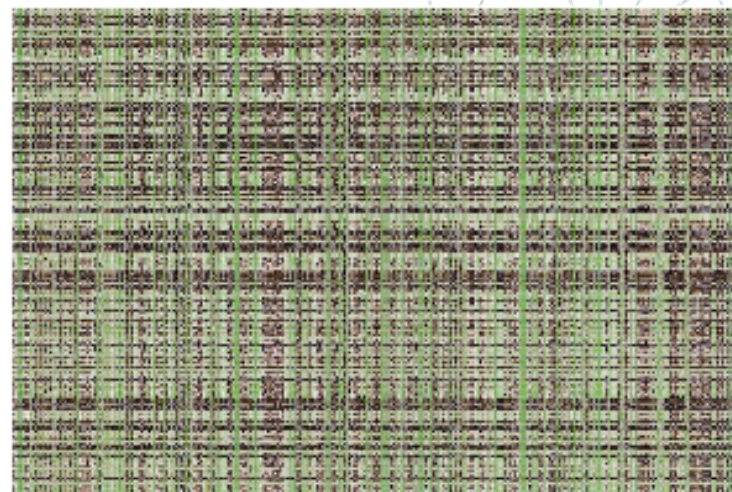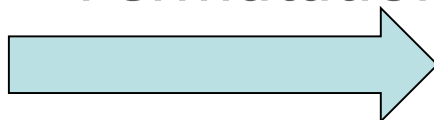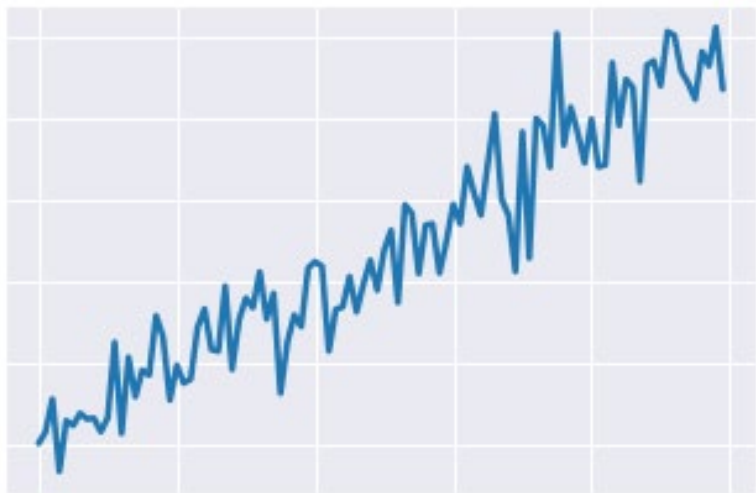| | Cement (component 1)(kg in a m^3 mixture) | Blast Furnace Slag (component 2)(kg in a m^3 mixture) | Fly Ash (component 3)(kg in a m^3 mixture) | Water (component 4)(kg in a m^3 mixture) | Superplasticizer (component 5)(kg in a m^3 mixture) | Coarse Aggregate (component 6)(kg in a m^3 mixture) | Fine Aggregate (component 7)(kg in a m^3 mixture) | Age (day) |
|---|---|---|---|---|---|---|---|---|
| 0 | 540.0 | 0.0 | 0.0 | 162.0 | 2.5 | 1040.0 | 676.0 | 28 |
| 1 | 540.0 | 0.0 | 0.0 | 162.0 | 2.5 | 1055.0 | 676.0 | 28 |
| 2 | 332.5 | 142.5 | 0.0 | 228.0 | 0.0 | 932.0 | 594.0 | 270 |
| 3 | 332.5 | 142.5 | 0.0 | 228.0 | 0.0 | 932.0 | 594.0 | 365 |
| 4 | 198.6 | 132.4 | 0.0 | 192.0 | 0.0 | 978.4 | 825.5 | 360 |

## Random Permutation

| | Fly Ash (component 3)(kg in a m^3 mixture) | Cement (component 1)(kg in a m^3 mixture) | Fine Aggregate (component 7)(kg in a m^3 mixture) | Superplasticizer (component 5)(kg in a m^3 mixture) | Age (day) | Coarse Aggregate (component 6)(kg in a m^3 mixture) | Water (component 4)(kg in a m^3 mixture) | Blast Furnace Slag (component 2)(kg in a m^3 mixture) |
|---|---|---|---|---|---|---|---|---|
| 0 | 0.0 | 540.0 | 676.0 | 2.5 | 28 | 1040.0 | 162.0 | 0.0 |
| 1 | 0.0 | 540.0 | 676.0 | 2.5 | 28 | 1055.0 | 162.0 | 0.0 |
| 2 | 0.0 | 332.5 | 594.0 | 0.0 | 270 | 932.0 | 228.0 | 142.5 |
| 3 | 0.0 | 332.5 | 594.0 | 0.0 | 365 | 932.0 | 228.0 | 142.5 |
| 4 | 0.0 | 198.6 | 825.5 | 0.0 | 360 | 978.4 | 192.0 | 132.4 |

Random Permutation

Random
Permutation

# Limitation of FCNN: The permutation invariance property loses spatial/temporal structure in the data!

# The Convolution Operation

# Tracking a Spaceship

**Suppose we are tracking the location of a spaceship with a laser sensor**

**The sensor provides a single output**

$$x(t) \approx \text{Position of spaceship at time } t$$

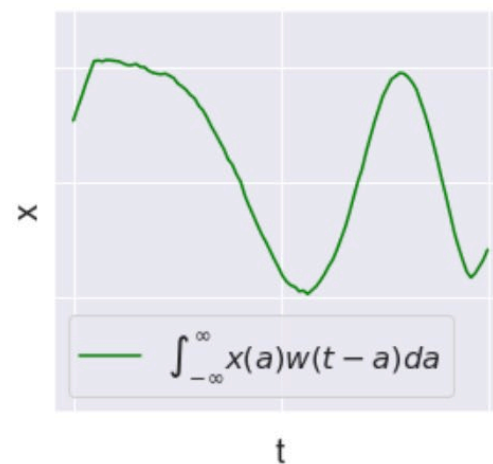**The sensor is noisy!**

# Averaged Measurement

To obtain a less noisy measurement, we average measurements received previously:

$$s(t) = \int_{-\infty}^{+\infty} x(a)w(t-a)\,da$$

For example, we can take

$$w(t-a) = \begin{cases} e^{-(t-a)} & t \geq a \\ 0 & t < a \end{cases}$$

# Convolution of Functions

The previous weighted average is in fact a special case of convolutions.

Given two real-valued functions $x(t)$ and $w(t)$, their convolution is denoted by

$$s(t) = x * w(t)$$

$$s(t) = \int_{-\infty}^{+\infty} x(a)w(t-a)da$$

**The convolution operation is symmetric (commutative)**

$$x * w = w * x$$

**But, we usually distinguishes them**
- $x(t)$ is called the **signal** or **input**
- $w(t)$ is called the **kernel** or **filter**
- The output $s = x * w$ is sometimes called the **feature map**

# Discrete Convolutions

In reality, measurements can only be carried out at discrete time steps

$$t = \ldots, 0,1,2, \ldots$$

In this case, it is useful to define the discrete convolution
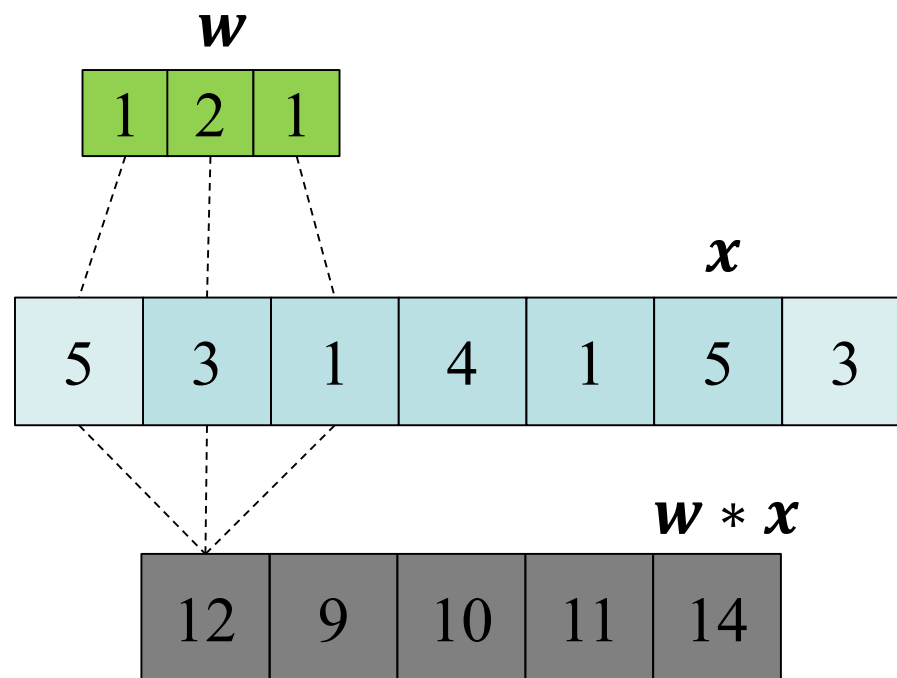
$$s(t) = (x * w)(t) = \sum_{-\infty}^{+\infty} x(a)w(t - a)$$

# Finite Convolutions and Boundary Conditions

Practical signals are finite, so need to truncate

- Circular convolutions

Practical signals are finite, so need to truncate

- Valid convolutions



$w$

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

$x$

| 3 | 1 | 4 | 1 | 5 |
|---|---|---|---|---|

$w * x$

| 9 | 10 | 11 |
|---|---|---|

# Finite Convolutions and Boundary Conditions

Practical signals are finite, so need to truncate
- Zero-Padded convolutions with "same" padding

# 2D Discrete (Finite) Convolutions

For image/vision applications, signals come in the form of a 2D matrix. In this case, we define the 2D convolution operation given an input image $I$ and a kernel $K$

$$S(i,j) = (I * K)(i,j)$$

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

As before, we have $I * K = K * I$

# Cross-correlation vs Convolution

**A similar concept to convolution is the cross-correlation:**

- Convolution

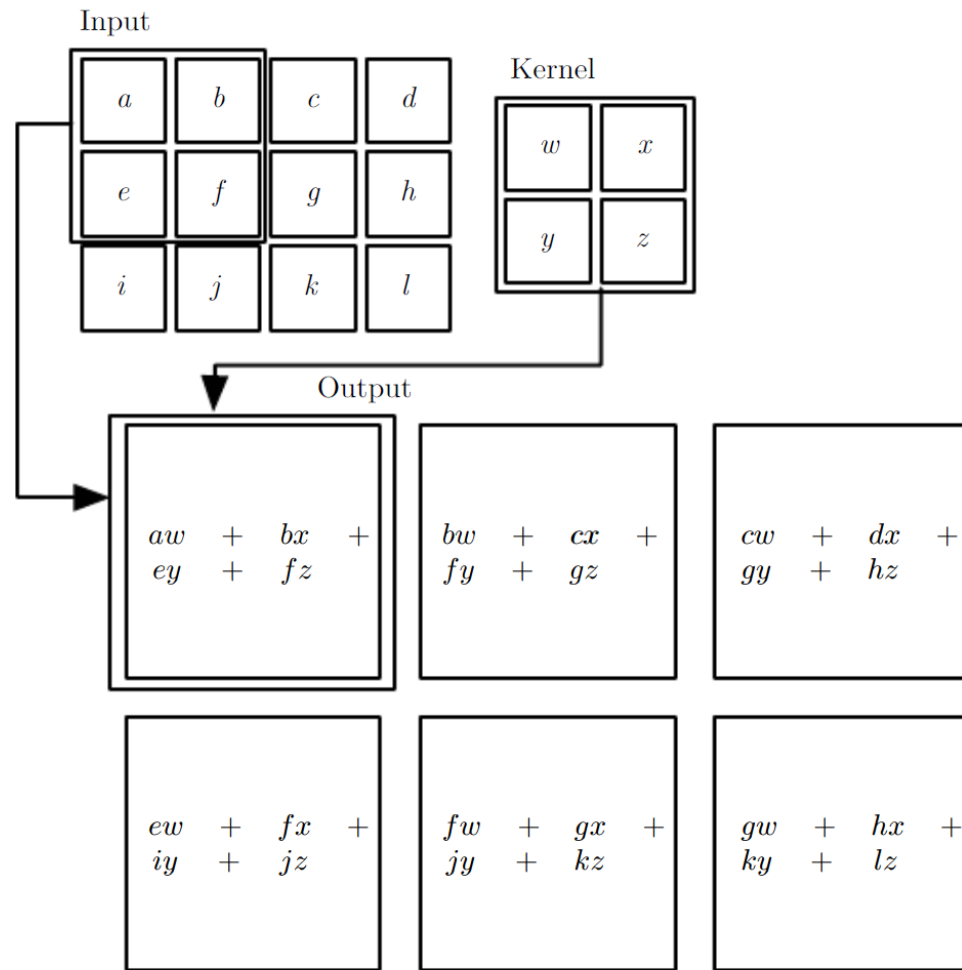$$S(i,j) = \sum_m \sum_n I(i - m, j - n)K(m,n)$$

- Cross-correlation
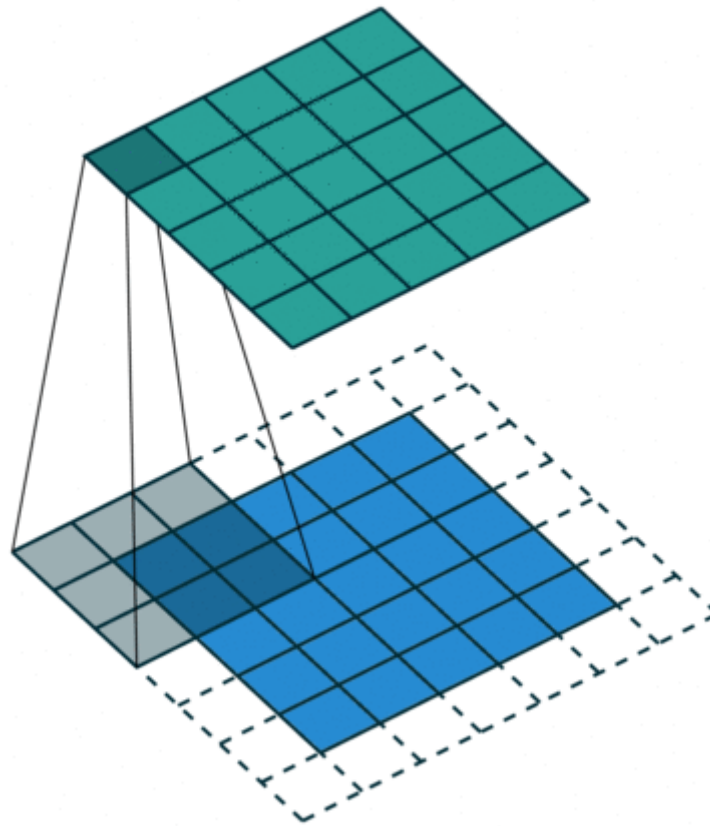
$$S(i,j) = \sum_m \sum_n I(i + m, j + n)K(m,n)$$

**Value is the same, but cross-correlation loses commutativity.**

**Deep learning libraries usually implement cross-correlation but calls it convolution**

# Example: 2D "Valid" Convolution

# Example: 2D "Same" Convolution

# Convolution as a Linear Operation

**Consider 1D convolution (zero-padded, same)**

$$w = (w_1, w_2, w_3) \text{ and } x = (x_1, \dots, x_5)$$

$$w * x = \underbrace{\begin{pmatrix} w_2 & w_3 & 0 & 0 & 0 \\ w_1 & w_2 & w_3 & 0 & 0 \\ 0 & w_1 & w_2 & w_3 & 0 \\ 0 & 0 & w_1 & w_2 & w_3 \\ 0 & 0 & 0 & w_1 & w_2 \end{pmatrix}}_{A(w)} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix}$$

**Convolutions are still linear operations:** $w * x = A(w)x$

# Convolution Neural Networks

**The basic idea of convolutional neural networks (CNN) is to replace the linear operation in FCNN by convolutions:**

$$h^{(i+1)} = g\big(W^{(i+1)}h^{(i)} + b^{(i+1)}\big)$$

$$h^{(i+1)} = g\big(w^{(i+1)} * h^{(i)} + b^{(i+1)}\big)$$
$$= g\big(A\big(w^{(i+1)}\big)h^{(i)} + b^{(i+1)}\big)$$

# Why Convolutions?

# Motivation 0: Convolutions are effective feature extractors



$$w_1 * x$$

$$w_1$$

$$x$$

$$w_1 * x$$

$$w_2$$

# Motivation 1: Sparse Interactions

**Due to the spatial locality of convolutions, not every node interact with every node!**

Convolution
$$A(\boldsymbol{w}) = \begin{pmatrix} w_2 & w_3 & 0 & 0 & 0 \\ w_1 & w_2 & w_3 & 0 & 0 \\ 0 & w_1 & w_2 & w_3 & 0 \\ 0 & 0 & w_1 & w_2 & w_3 \\ 0 & 0 & 0 & w_1 & w_2 \end{pmatrix}$$

Fully Connected
$$W = \begin{pmatrix} W_{11} & W_{12} & W_{13} & W_{14} & W_{15} \\ W_{21} & W_{22} & W_{23} & W_{24} & W_{25} \\ W_{31} & W_{32} & W_{33} & W_{34} & W_{35} \\ W_{41} & W_{42} & W_{43} & W_{44} & W_{45} \\ W_{51} & W_{52} & W_{53} & W_{54} & W_{55} \end{pmatrix}$$

Convolution

Fully Connected

Convolution

Receptive Field

Fully Connected

# Implicit Connectivity

**Although for each layer, the dependence is local.**
**For multi-layer CNNs, signals far apart are connected implicitly**

# Computational Efficiency

**Image size:** $m \times n$    **Kernel size:** $k \times l$

**Computational cost for one matrix multiplication:**
- CNN: $\mathcal{O}(k \times l \times m \times n)$
- FCNN: $\mathcal{O}(m^2 \times n^2)$

**This represents savings if** $k, l \ll m, n$**, which is often the case!**

# Motivation 2: Parameter Sharing

**Parameter sharing refers to using the same trainable parameters for more than one place in a ML model**

**An alternative name is tied weights**

$$A(\boldsymbol{w}) = \begin{pmatrix} w_2 & w_3 & 0 & 0 & 0 \\ w_1 & w_2 & w_3 & 0 & 0 \\ 0 & w_1 & w_2 & w_3 & 0 \\ 0 & 0 & w_1 & w_2 & w_3 \\ 0 & 0 & 0 & w_1 & w_2 \end{pmatrix} \quad \text{(Only 3 degrees of freedom)}$$
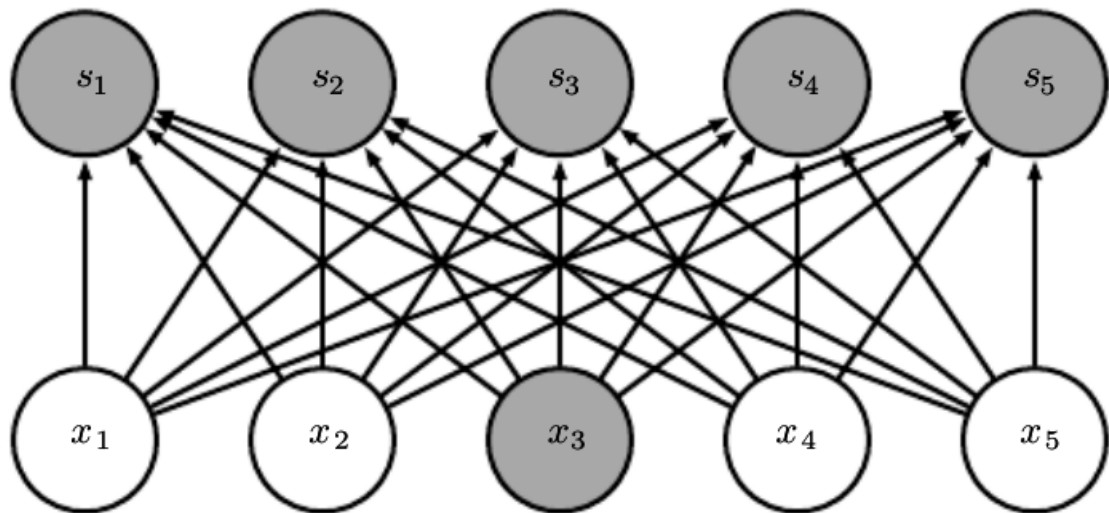
Convolution

Fully Connected

# Memory Efficiency

**Image size:** $m \times n$    **Kernel size:** $k \times l$

**Storage size for trained parameters:**
$$\mathcal{O}(k \times l)$$

**Without parameter sharing, need:**

$$\mathcal{O}(m \times n \times k \times l)$$

# Motivation 3: Equivariance and Invariance

Let $f, g$ be functions. We say that $f$ is equivariant with respect to $g$ if

$$f\big(g(x)\big) = g\big(f(x)\big)$$

We say that $f$ is invariant with respect to $g$ if

$$f\big(g(x)\big) = f(x)$$

# Example: Equivariance

**Let**

- $g(x) = \lambda x$ be a positive scaling function, where $\lambda \geq 0$

- $f(x) = \mathrm{ReLU}(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$

**Then, $f$ is equivariant with respect to $g$**

# Example: Invariance

**Let**
- $g(\boldsymbol{x}) = G\boldsymbol{x}$ where $G$ is an orthogonal matrix ($G^{-1} = G^T$)
- $f(\boldsymbol{x}) = \|\boldsymbol{x}\|^2$ be the norm function

**Then, $f$ is invariant with respect to $g$**

**Let**
- $g(\boldsymbol{x})_j = x_{j-\tau}$ be a translation by $\tau$
- $f(\boldsymbol{x}) = \mathbf{1}^T \boldsymbol{x} + c$ be an affine transformation

**Then $f$ is invariant with respect to $g$**

# General Group Equivariance and Invariance*

More generally, let $f: \mathbb{R}^d \to \mathbb{R}^{d'}$ be a function and $\mathcal{G}$ be a group (of symmetries).

For any $g \in \mathcal{G}$ we denote by $\phi_g$ and $\psi_g$ some group actions on $\mathbb{R}^d$ and $\mathbb{R}^{d'}$ respectively

- We say that $f$ is **equivariant** under $\mathcal{G}$ if for any $g \in \mathcal{G}$ we have

$$f \circ \phi_g = \psi_g \circ f$$

- Invariance is a special case of equivariance if $\psi_g = \mathrm{id}$ (identity map)

# Convolutions and Translations

Let $T_\tau$ be a translation of a signal $x$ in time by $\tau$

$$T_\tau(x)(t) = x(t - \tau)$$

Then, for any kernel $w$ we have

$$w * T_\tau(x) = T_\tau(w * x)$$

Because

$$[w * T_\tau(x)](t) = \int w(s)\big(T_\tau(x)\big)(t - s)ds$$
$$= \int w(s)x(t - s - \tau)ds$$
$$= \int w(s)x[(t - \tau) - s]ds = T_\tau(w * x)$$

In other words, convolutions are equivariant with respect to translations.

# Example: Discrete 1D Convolutions

Let

$$x = (1, 2, 3, 4, 5)$$
$$w = (1, -2, 1)$$

Define translation $T(x) = (x_3, x_4, x_5, x_1, x_2)$, then equivariance holds under circular convolution.

What about other paddings?

# Example: Element-wise Nonlinearities

**Instead of convolutions, we consider element-wise nonlinear activations.**

$$x = (1, 2, 3, 4, 5)$$
$$\sigma(x) = \big(\sigma(1), \dots, \sigma(5)\big)$$

**Then** $\sigma\big(T(x)\big) = T\big(\sigma(x)\big)$

# Compositions of Equivariant Transformations

**Let $f_1, f_2$ be equivariant with respect $g$, then so is $f_2 \circ f_1$:**

$$f_2 \circ f_1\big(g(x)\big) =$$
$$= f_2\left(f_1(g(x))\right) = f_2\left(g(f_1(x))\right)$$
$$= g\left(f_2(f_1(x))\right) = g\big(f_2 \circ f_1(x)\big)$$

**This means that each convolution layer**

$$h^{(i+1)} = \sigma\big(w * h^{(i)} + b\big)$$

**is equivariant with respect to translations**

# Building an Invariant Mapping

Let $f_1, \ldots, f_\ell$ be a collection of functions equivariant with respect to $g$

Let $F$ be another function that is invariant with respect to $g$

Then, observe that $F \circ f_\ell \circ \cdots \circ f_1$ is invariant with respect to $g$

In other words, a deep CNN

$$h^{(i+1)} = \sigma\big(W * h^{(i)} + b\big) \qquad h^{(0)} = x$$
$$\hat{y} = \mathbf{1}^T h^{(\ell)} + c$$

is invariant with respect to translations!

# Importance of Translation Invariance

**Let** $f^* : \text{Images} \rightarrow \{\text{dog}, \text{cat}\}$



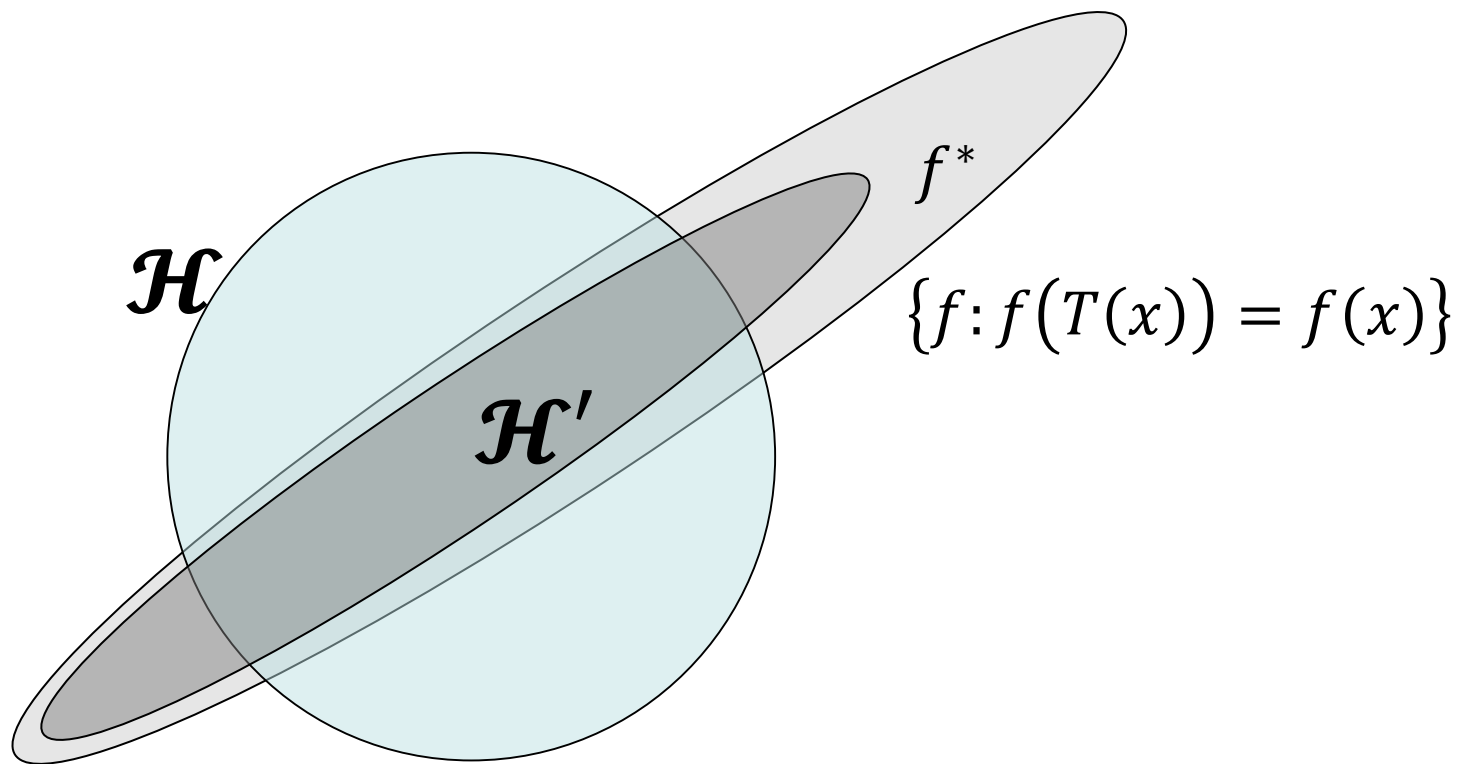These three images are translations of the same image, but they all give the label "dog".

In other words, we must have
$$f^*\big(T(I)\big) = f^*(I)$$

for this classification problem!

Using a CNN to solve this problem places precisely this **prior**

$\mathcal{H}$

$\mathcal{H}'$

$f^*$

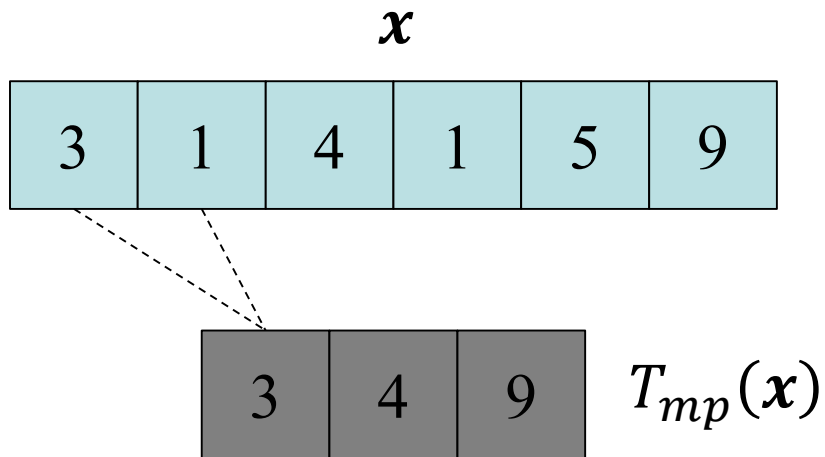$\{f : f(T(x)) = f(x)\}$

# Pooling Layers

**Pooling is another type of operation that builds approximate invariance to small translations/deformations**

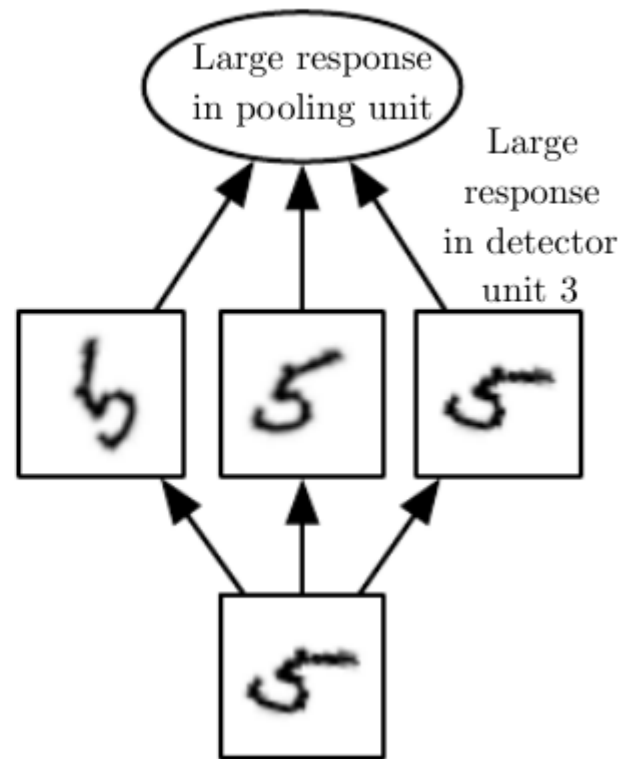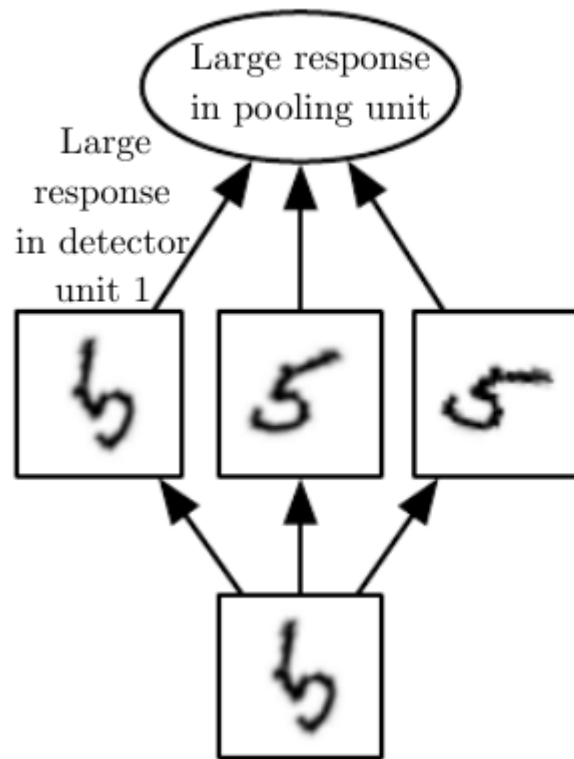**Simplest form is max pooling, but there are also other types, e.g. average pooling**
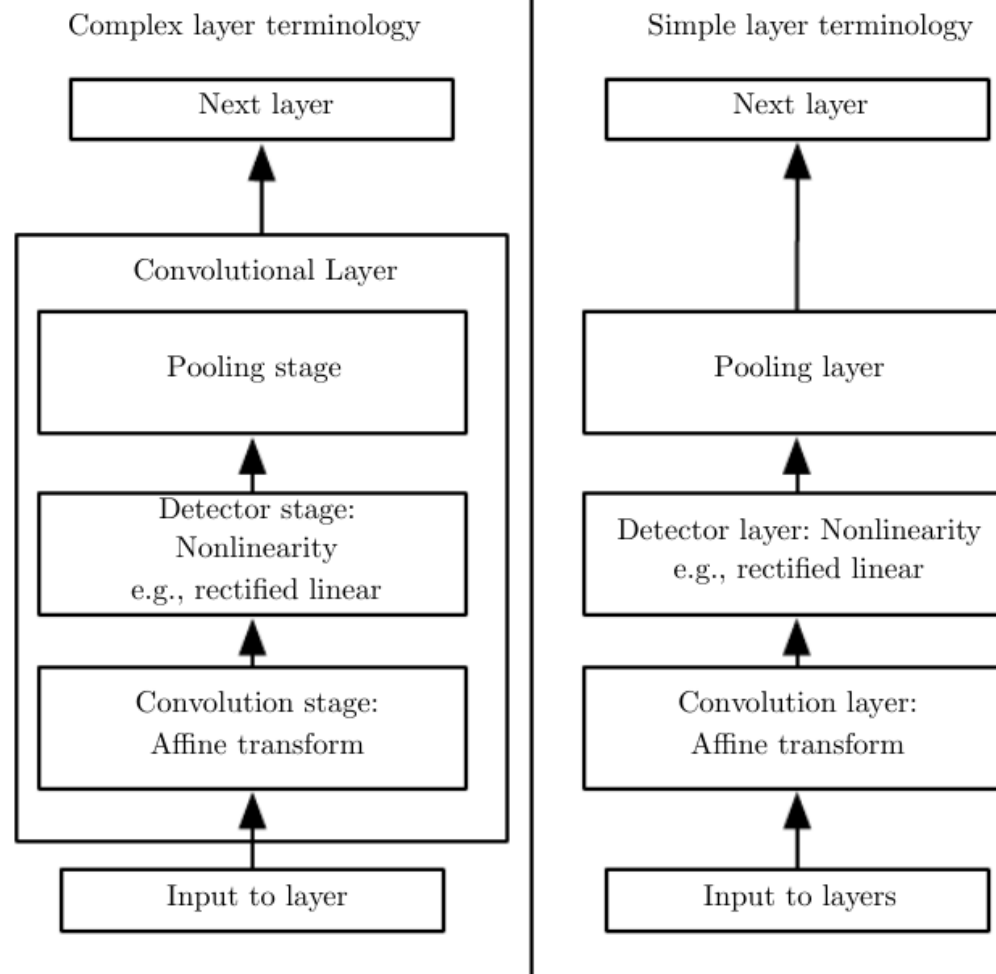
Max pooling In 1D with stride $p$:

$$T_{mp}(\boldsymbol{x})_k = \max_{i=kp,\dots,(k+1)p} x_i$$



$\boldsymbol{x}$

| 3 | 1 | 4 | 1 | 5 | 9 |
|---|---|---|---|---|---|

| 3 | 4 | 9 | $T_{mp}(\boldsymbol{x})$ |
|---|---|---|---|

$I$

| 3 | 1 | 4 | 1 |
|---|---|---|---|
| 5 | 9 | 2 | 6 |
| 5 | 3 | 5 | 8 |
| 9 | 7 | 9 | 3 |

$T_{mp}(I)$

| 9 | 6 |
|---|---|
| 9 | 9 |

# The Basic Structure of CNNs

# The Basic Structure of CNNs



Conv 1
$3 \times 3$ Filters

Conv 2
$3 \times 3$ Filters

Flatten

Max Pool
$2 \times 2$

Max Pool
$2 \times 2$

FCNN

$$\begin{pmatrix} 0.9 \\ 0.1 \end{pmatrix}$$

# Variants of Convolutions

**There are many variants to the basic conv and pooling layers introduced so far, including:**

- Strided convolutions
- Up and down-sampling operations
- Unshared convolutions
- Channel-wise/separable convolutions

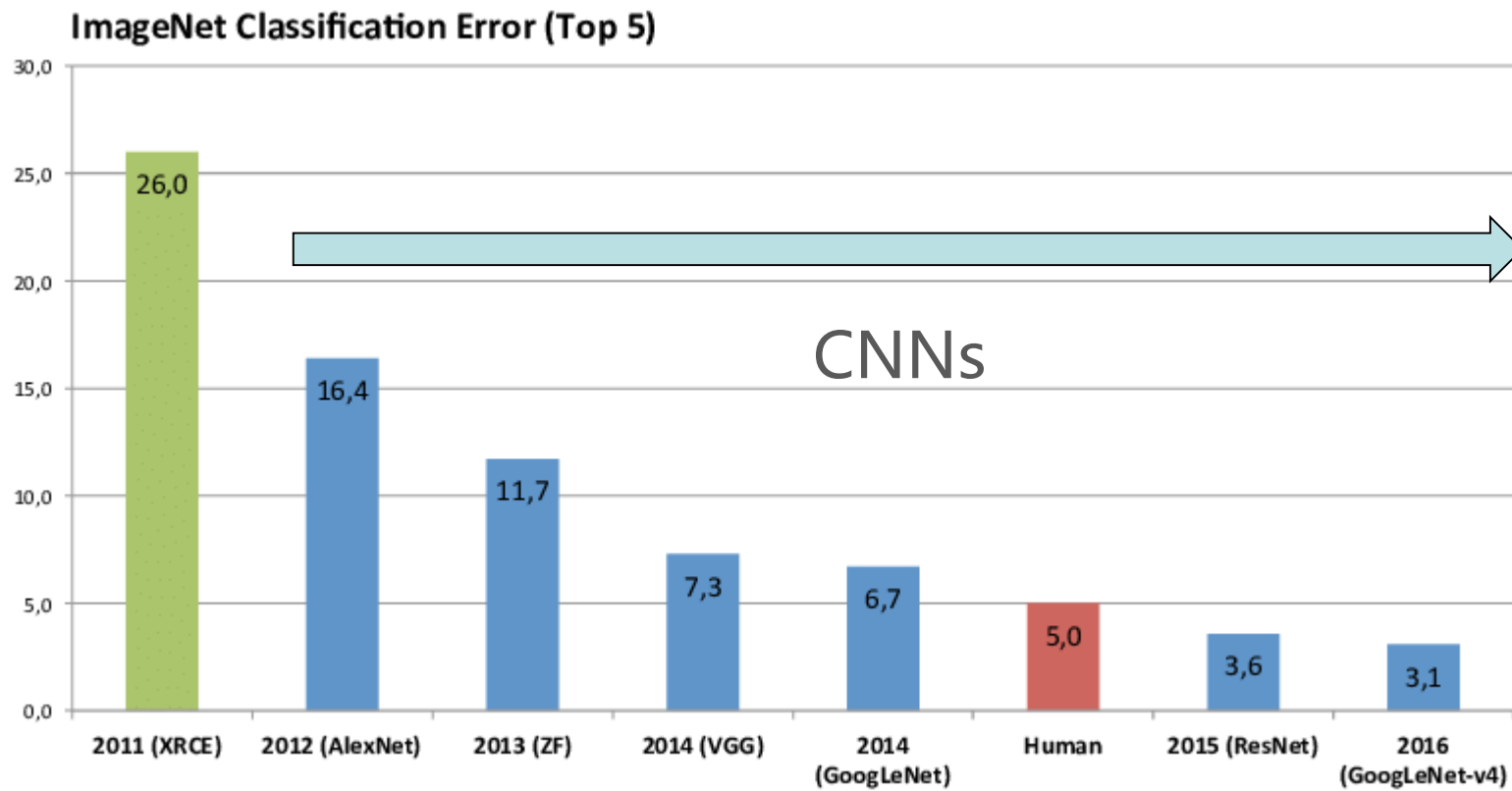**Refer to Chapter 9 in the deep learning book for more details**

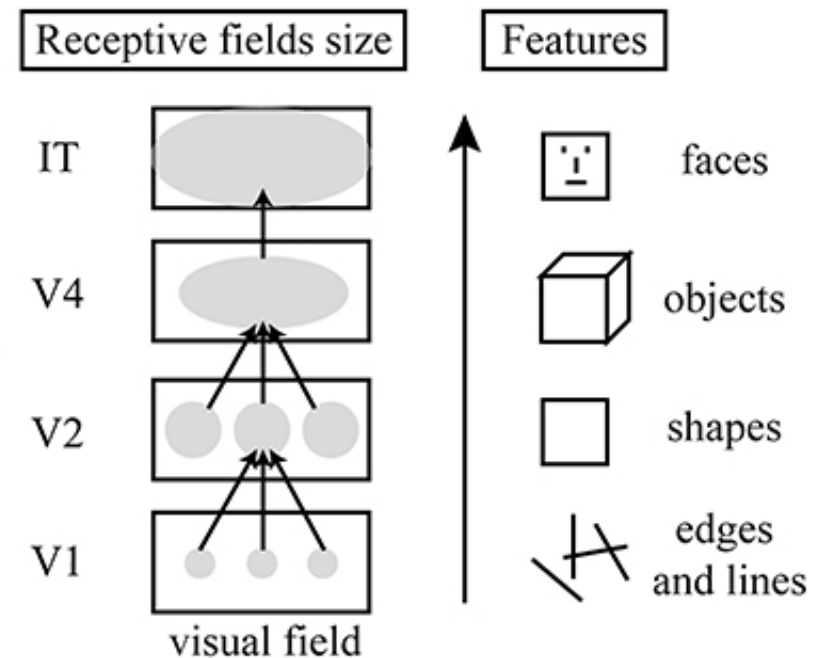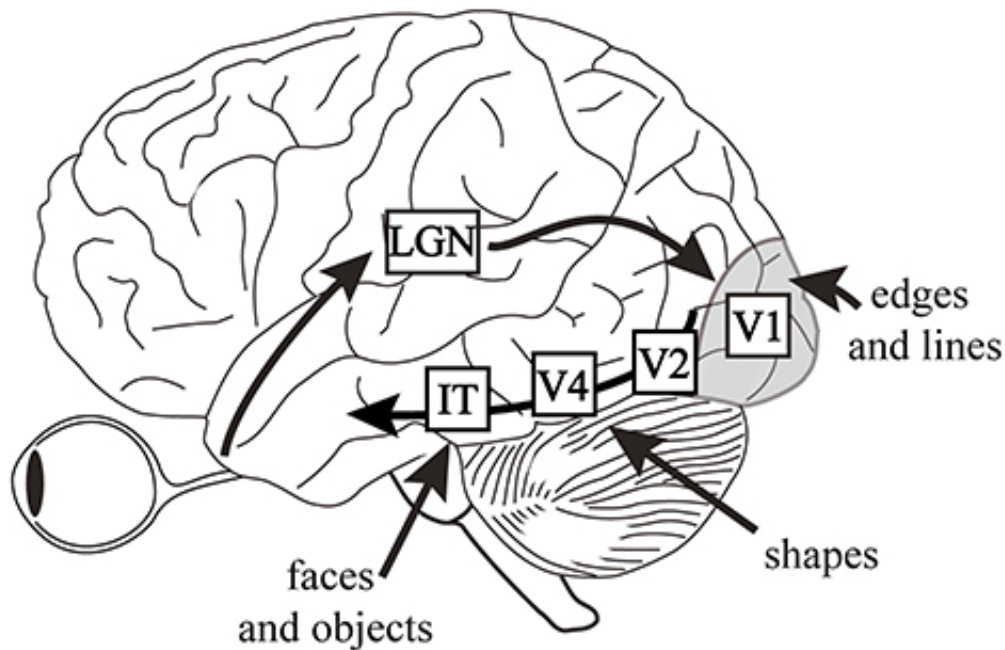# Historical notes on CNN

# The ImageNet Challenge

A large benchmark problem consisting of ~1M images, 20k categories



*http://www.image-net.org/*

# History of CNNs



ImageNet Classification Error (Top 5)

# Demo: CNNs for MNIST and Pneumonia Classification

# Summary

In this lecture, we introduced

- Limitations of FCNNs
- Convolution as a replacement for general linear transformations
- Motivations of using convolutions for image processing
- CNN based on convolution (and pooling)