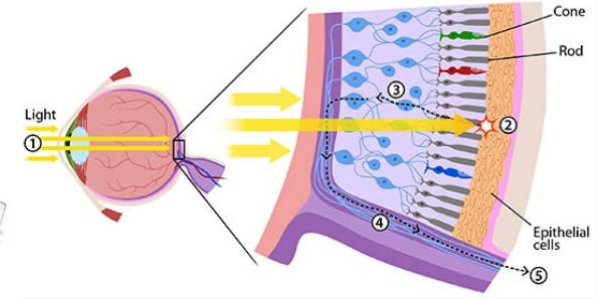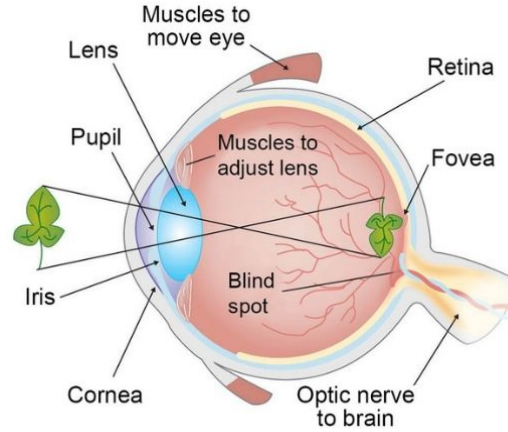# How do Computers see?

SURF

# How do Humans see?

SURF

# Mammalian Visual System

Cones and rods picks up photons and propagate the signal to the back of our brain

Along the way, the signal gets processed in stages



SURF

# Mammalian Visual System

Cones and rods picks up photons and propagate the signal to the back of our brain

Along the way, the signal gets processed in stages

The extracted information gets aggregated and formed into an image
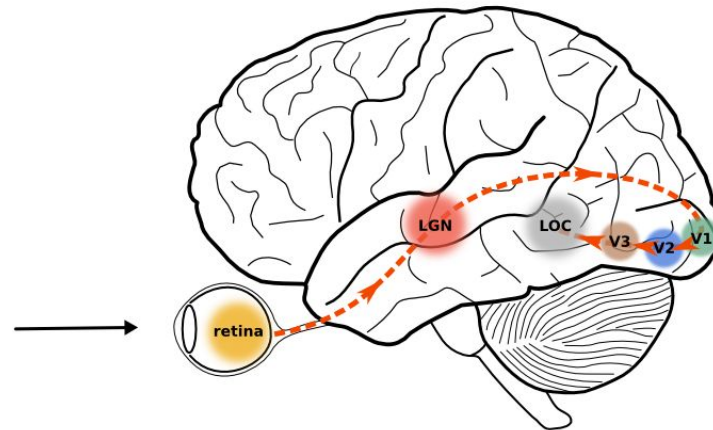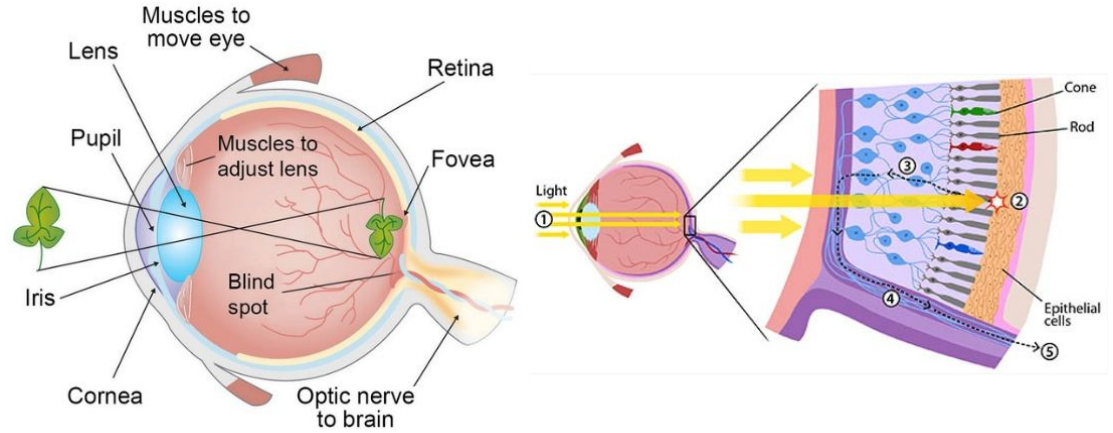


SURF

# Mammalian Visual System

Cones and rods picks up photons and propagate the signal to the back of our brain

Along the way, the signal gets processed in stages

The extracted information gets aggregated and formed into an image
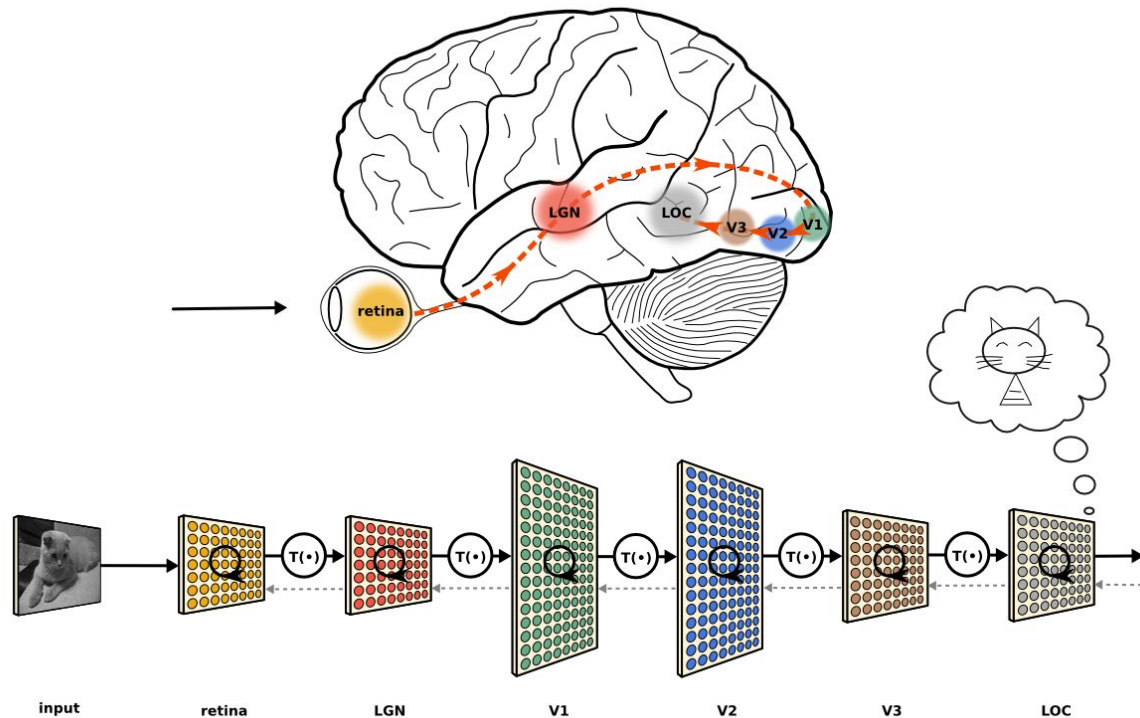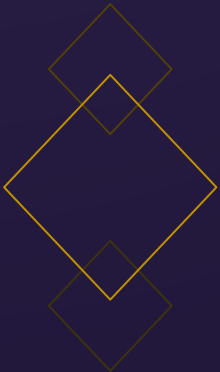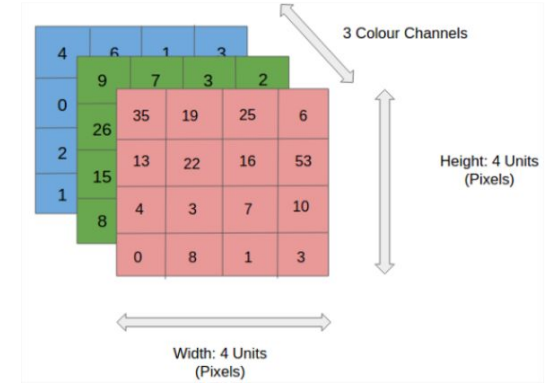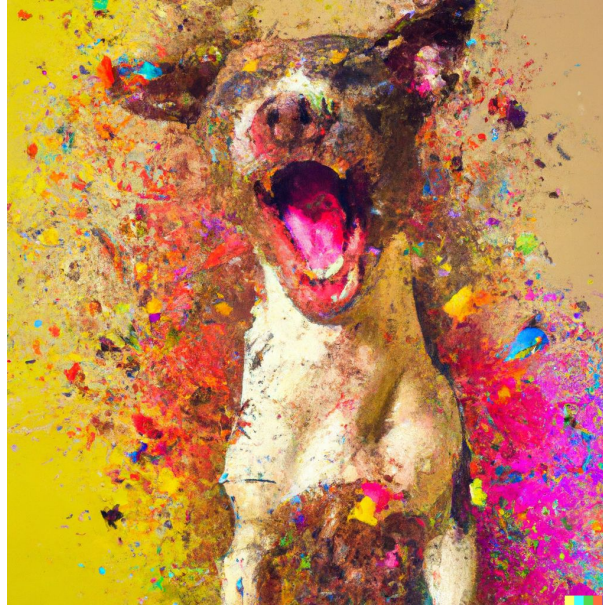
# 02. Convolutional Neural Networks

# Image Representation
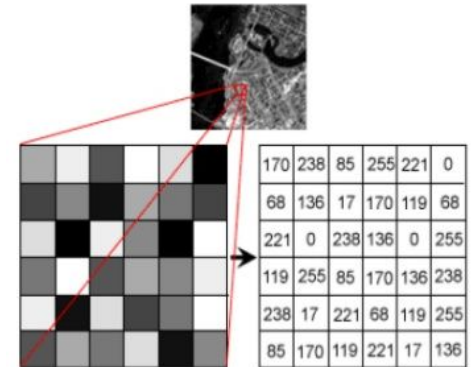
**01.** Is a matrix or grid of intensity values

**02.** integers [0, 255] or float points [0,1]

**03.** Each element in the matrix is a **pixel**

**04.** Can have 1 greyscale channel or multiple colour **channels: RGB**

# Image Representation

**01.** Is a matrix or grid of intensity values

**02.** integers [0, 255] or float points [0,1]

**03.** Each element in the matrix is a **pixel**

**04.** Can have 1 greyscale channel or multiple colour **channels: RGB**





3 Colour Channels

Height: 4 Units (Pixels)

Width: 4 Units (Pixels)

| 4 | 6 | 1 | 3 |
| 0 | 9 | 7 | 3 | 2 |
| 2 | 26 | 35 | 19 | 25 | 6 |
| 1 | 15 | 13 | 22 | 16 | 53 |
| | 8 | 4 | 3 | 7 | 10 |
| | | 0 | 8 | 1 | 3 |



| 170 | 238 | 85 | 255 | 221 | 0 |
| 68 | 136 | 17 | 170 | 119 | 68 |
| 221 | 0 | 238 | 136 | 0 | 255 |
| 119 | 255 | 85 | 170 | 136 | 238 |
| 238 | 17 | 221 | 68 | 119 | 255 |
| 85 | 170 | 119 | 221 | 17 | 136 |

SURF

# Correlated structures



$x_1$
$x_2$
$\vdots$
$x_{784}$

0
1
$\vdots$

Images are **correlated spatially**

Strong correlation between neighbouring pixels

Shifting an object will not alter its essence

# Correlated structures



$x_1$    0
$x_2$    1
$\vdots$
$x_{784}$

Images are **correlated spatially**

Strong correlation between neighbouring pixels

Shifting an object will not alter its essence

# Correlated structures



$x_1$
$x_2$
$\vdots$
$x_{784}$

0
1
$\vdots$

Images are **correlated spatially**

Strong correlation between neighbouring pixels

Shifting an object will not alter its essence

# Correlated structures



$x_1$   0

$x_2$   1

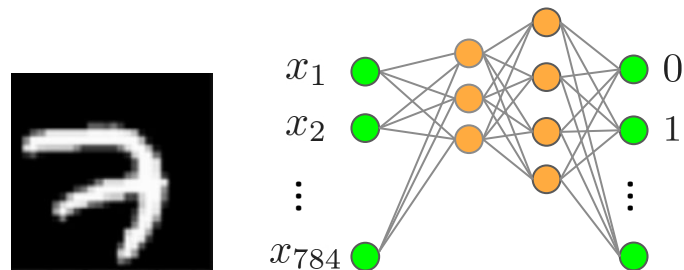$\vdots$

$x_{784}$

Images are **correlated spatially**

Strong correlation between neighbouring pixels

Shifting an object will not alter its essence

# Correlated structures



$x_1$  0
$x_2$  1
$\vdots$
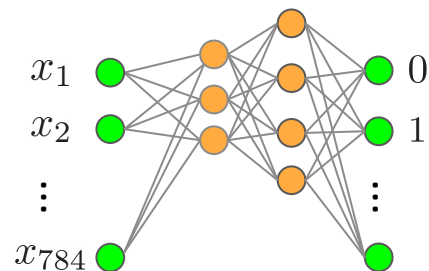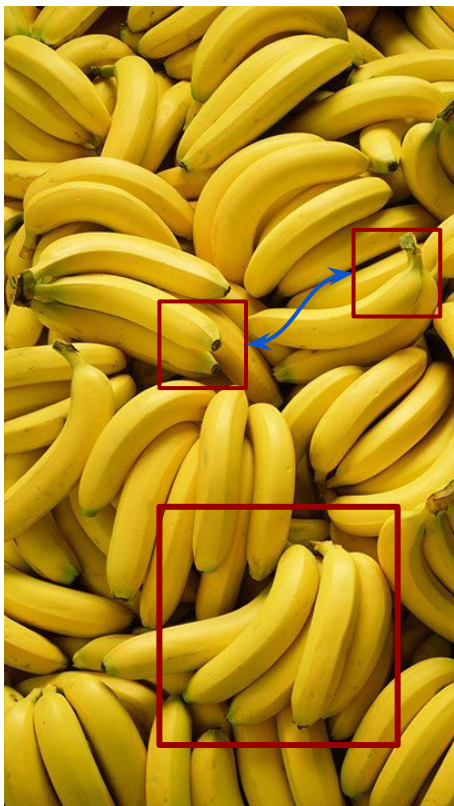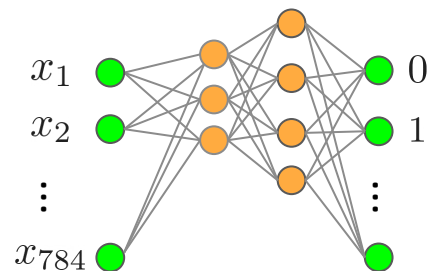$x_{784}$

Images are **correlated spatially**

Strong correlation between neighbouring pixels

Shifting an object will not alter its essence

# Correlated structures



$x_1$   0
$x_2$   1
$\vdots$
$x_{784}$

Detect the features across an image and aggregate them

Form complete features higher up in the **hierarchy**
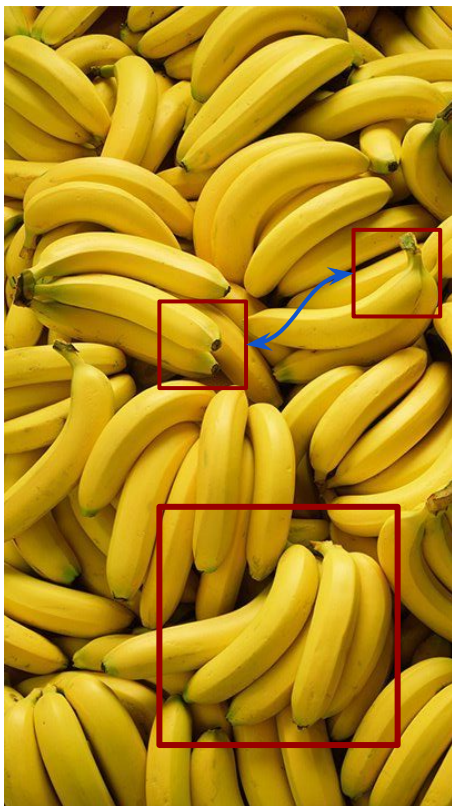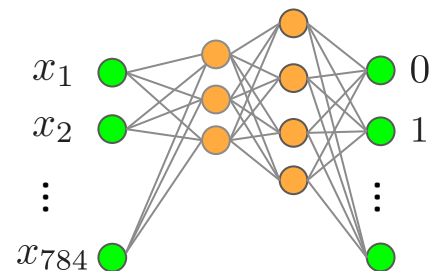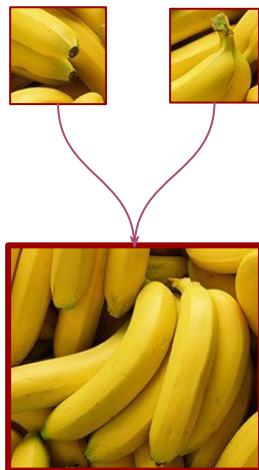
Images are **correlated spatially**

Strong correlation between neighbouring pixels

Shifting an object will not alter its essence

# ImageNet Challenge

**Computer Vision Benchmark**    **1.4M Images, 1000 classes**

**Image Classification**    **Difficult until 2012**

# ImageNet Challenge

**Computer Vision Benchmark**  **1.4M Images, 1000 classes**

**Image Classification**  **Difficult until 2012**



ImageNet Classification Error (Top 5)

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 2011 (XRCE) | 2012 (AlexNet) | 2013 (ZF) | 2014 (VGG) | 2014 (GoogLeNet) | Human | 2015 (ResNet) | 2016 (GoogLeNet-v4) |
| 26,0 | 16,4 | 11,7 | 7,3 | 6,7 | 5,0 | 3,6 | 3,1 |

# Convolutions

**Kernel** (filter): small matrix that we use to convolve an image

**Convolution:**

An operation that "blends" one function with another.

| Operation | Kernel | Image result |
|---|---|---|
| **Identity** | $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ | |
| **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) | $\frac{1}{9}\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ | |

SURF

# Convolutions

**Kernel** (filter): small matrix that we use to convolve an image

**Convolution:**

An operation that "blends" one function with another.

A filter applies a **convolution** operation on an image



| Operation | Kernel | Image result |
|---|---|---|
| **Identity** | $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ | |
| **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) | $\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ | |

**SURF**

# Convolutions

**Kernel** (filter): small matrix that we use to convolve an image

**Convolution:**

An operation that "blends" one function with another.

A filter applies a **convolution** operation on an image

# Convolutions
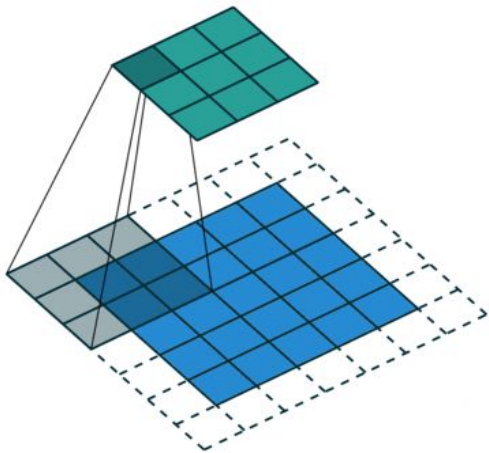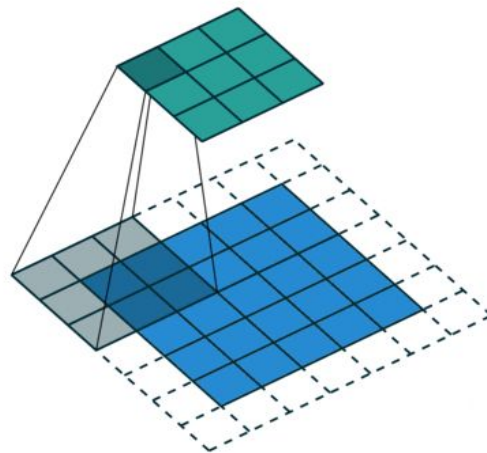
**Kernel** (filter): small matrix that we use to convolve an image

**Convolution:**

An operation that "blends" one function with another.

A filter applies a **convolution** operation on an image

# Convolutions Compute



**Source layer (image)**

**Convolutional Kernel (filter)**

**Destination layer**

$(-1×2) + (0×2) + (1×6) +$
$(-2×4) + (0×3) + (2×4) +$
$(-1×3) + (0×9) + (1×4) = 5$

# Convolutions Compute

| 2 | 2 | 6 | 8 | 2 | 0 | 1 | 2 |
|---|---|---|---|---|---|---|---|
| 4 | 3 | 4 | 5 | 1 | 9 | 6 | 3 |
| 3 | 9 | 4 | 4 | 7 | 7 | 6 | 9 |
| 1 | 3 | 4 | 6 | 8 | 2 | 2 | 1 |
| 8 | 4 | 6 | 2 | 3 | 1 | 8 | 8 |
| 5 | 8 | 9 | 0 | 1 | 0 | 2 | 3 |
| 9 | 2 | 6 | 6 | 3 | 6 | 2 | 1 |
| 9 | 8 | 8 | 2 | 6 | 3 | 4 | 5 |

**Source layer (image)**

**Convolutional Kernel (filter)**

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

**Destination layer**

| | | | | | |
|---|---|---|---|---|---|
| 5 | | | | | |

$$(-1\times2) + (0\times2) + (1\times6) +$$
$$(-2\times4) + (0\times3) + (2\times4) +$$
$$(-1\times3) + (0\times9) + (1\times4) = 5$$

The kernel is **shifted** across the image and produces a point value
The step size in which it shifts is the **stride**
The output is always smaller, we use **padding** to preserve dimensions

SURF

# Convolutions Compute



2 2 6 8 2 0 1 2
4 3 4 5 1 9 6 3
3 9 4 4 7 7 6 9
1 3 4 6 8 2 2 1
8 4 6 2 3 1 8 8
5 8 9 0 1 0 2 3
9 2 6 6 3 6 2 1
9 8 8 2 6 3 4 5

**Source layer (image)**

**Convolutional Kernel (filter)**

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

**Destination layer**

5

(-1×2) + (0×2) + (1×6) +
(-2×4) + (0×3) + (2×4) +
(-1×3) + (0×9) + (1×4) = 5

The kernel is **shifted** across the image and produces a point value
The step size in which it shifts is the **stride**
The output is always smaller, we use **padding** to preserve dimensions

The image shrinks according to

Output_size = inputSize -(KernelSize - 1)

SURF

# Convolutions

**Source layer (image)**

| 2 | 2 | 6 | 8 | 2 | 0 | 1 | 2 |
|---|---|---|---|---|---|---|---|
| 4 | 3 | 4 | 5 | 1 | 9 | 6 | 3 |
| 3 | 9 | 4 | 4 | 7 | 7 | 6 | 9 |
| 1 | 3 | 4 | 6 | 8 | 2 | 2 | 1 |
| 8 | 4 | 6 | 2 | 3 | 1 | 8 | 8 |
| 5 | 8 | 9 | 0 | 1 | 0 | 2 | 3 |
| 9 | 2 | 6 | 6 | 3 | 6 | 2 | 1 |
| 9 | 8 | 8 | 2 | 6 | 3 | 4 | 5 |

**Kernel (filter)**

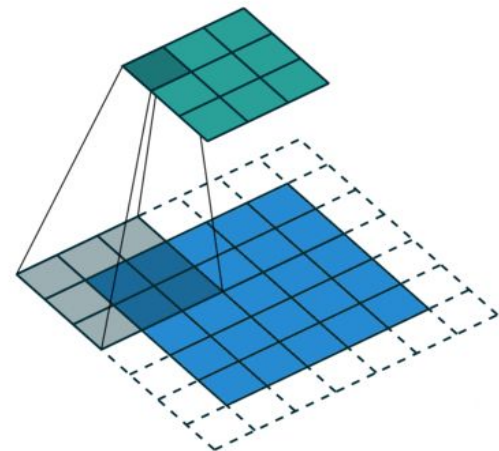| $w_1$ | $w_2$ | $w_3$ |
|-------|-------|-------|
| $w_4$ | $w_5$ | $w_6$ |
| $w_7$ | $w_8$ | $w_9$ |

**Feature map (activation map)**

|   | 5 |   |   |   |   |   |
|---|---|---|---|---|---|---|
|   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |

SURF

# Convolutions

**Source layer (image)**

| 2 | 2 | 6 | 8 | 2 | 0 | 1 | 2 |
|---|---|---|---|---|---|---|---|
| 4 | 3 | 4 | 5 | 1 | 9 | 6 | 3 |
| 3 | 9 | 4 | 4 | 7 | 7 | 6 | 9 |
| 1 | 3 | 4 | 6 | 8 | 2 | 2 | 1 |
| 8 | 4 | 6 | 2 | 3 | 1 | 8 | 8 |
| 5 | 8 | 9 | 0 | 1 | 0 | 2 | 3 |
| 9 | 2 | 6 | 6 | 3 | 6 | 2 | 1 |
| 9 | 8 | 8 | 2 | 6 | 3 | 4 | 5 |

**Kernel (filter)**

| $w_1$ | $w_2$ | $w_3$ |
|---|---|---|
| $w_4$ | $w_5$ | $w_6$ |
| $w_7$ | $w_8$ | $w_9$ |

**Feature map (activation map)**

The feature map grid shows the value $5$ in the second row, second column.

SURF

# Convolutions

**Source layer (image)**

| 2 | 2 | 6 | 8 | 2 | 0 | 1 | 2 |
|---|---|---|---|---|---|---|---|
| 4 | 3 | 4 | 5 | 1 | 9 | 6 | 3 |
| 3 | 9 | 4 | 4 | 7 | 7 | 6 | 9 |
| 1 | 3 | 4 | 6 | 8 | 2 | 2 | 1 |
| 8 | 4 | 6 | 2 | 3 | 1 | 8 | 8 |
| 5 | 8 | 9 | 0 | 1 | 0 | 2 | 3 |
| 9 | 2 | 6 | 6 | 3 | 6 | 2 | 1 |
| 9 | 8 | 8 | 2 | 6 | 3 | 4 | 5 |

**Kernel (filter)**

| $w_1$ | $w_2$ | $w_3$ |
|-------|-------|-------|
| $w_4$ | $w_5$ | $w_6$ |
| $w_7$ | $w_8$ | $w_9$ |

**Feature map (activation map)**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | 5 | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

**How do we know which kernels to use?**

Kernels are learnt and initialized randomly
during training the CNN learns spatial features

SURF

# Convolutions

**Source layer (image)**

| 2 | 2 | 6 | 8 | 2 | 0 | 1 | 2 |
|---|---|---|---|---|---|---|---|
| 4 | 3 | 4 | 5 | 1 | 9 | 6 | 3 |
| 3 | 9 | 4 | 4 | 7 | 7 | 6 | 9 |
| 1 | 3 | 4 | 6 | 8 | 2 | 2 | 1 |
| 8 | 4 | 6 | 2 | 3 | 1 | 8 | 8 |
| 5 | 8 | 9 | 0 | 1 | 0 | 2 | 3 |
| 9 | 2 | 6 | 6 | 3 | 6 | 2 | 1 |
| 9 | 8 | 8 | 2 | 6 | 3 | 4 | 5 |

**Kernel (filter)**

| $w_1$ | $w_2$ | $w_3$ |
|---|---|---|
| $w_4$ | $w_5$ | $w_6$ |
| $w_7$ | $w_8$ | $w_9$ |

**Feature map (activation map)**

| | 5 | | | | | |
|---|---|---|---|---|---|---|

**How do we know which kernels to use?**

**Fully connected**

Kernels are learnt and initialized randomly during training the CNN learns spatial features

$$\sum_{i \in \text{image}}^{W \times H \times C} \mathbf{x}_i \mathbf{w}_i$$

SURF

# Convolutions

## Source layer (image)

| 2 | 2 | 6 | 8 | 2 | 0 | 1 | 2 |
|---|---|---|---|---|---|---|---|
| 4 | 3 | 4 | 5 | 1 | 9 | 6 | 3 |
| 3 | 9 | 4 | 4 | 7 | 7 | 6 | 9 |
| 1 | 3 | 4 | 6 | 8 | 2 | 2 | 1 |
| 8 | 4 | 6 | 2 | 3 | 1 | 8 | 8 |
| 5 | 8 | 9 | 0 | 1 | 0 | 2 | 3 |
| 9 | 2 | 6 | 6 | 3 | 6 | 2 | 1 |
| 9 | 8 | 8 | 2 | 6 | 3 | 4 | 5 |

## Kernel (filter)

| $w_1$ | $w_2$ | $w_3$ |
|-------|-------|-------|
| $w_4$ | $w_5$ | $w_6$ |
| $w_7$ | $w_8$ | $w_9$ |

## Feature map (activation map)



**How do we know which kernels to use?**

Kernels are learnt and initialized randomly during training the CNN learns spatial features

**Fully connected**
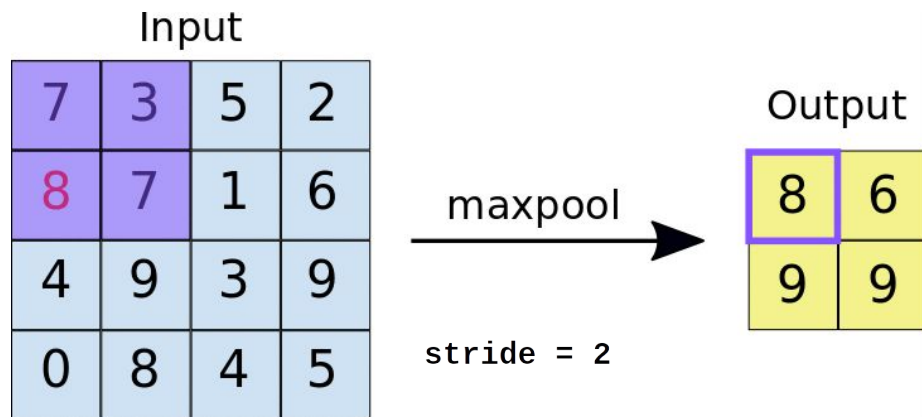
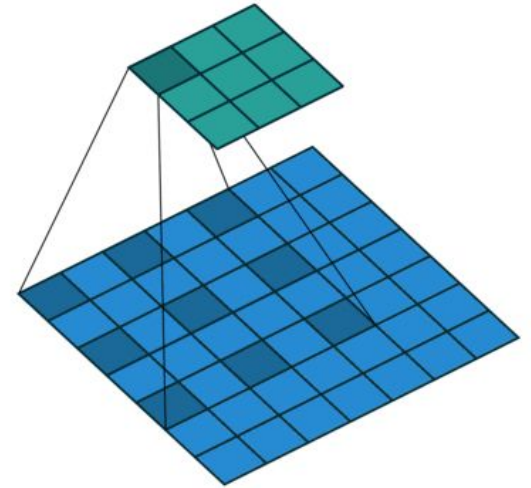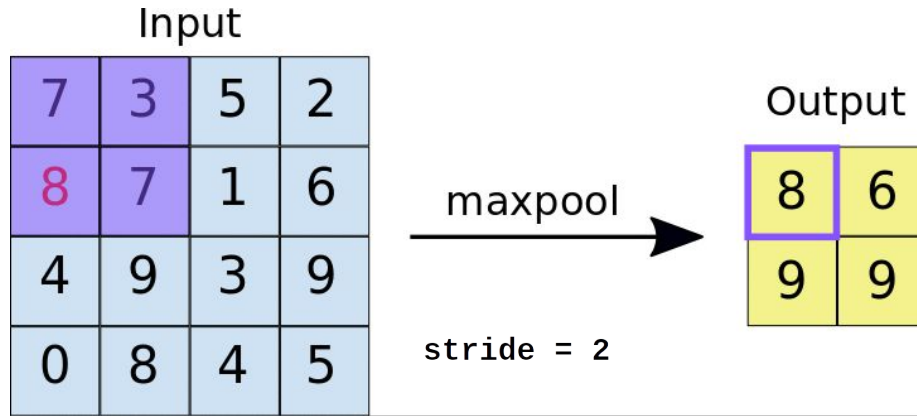$$\sum_{i \in \text{image}}^{W \times H \times C} \mathbf{x}_i \mathbf{w}_i$$

SURF

# Convolutions

**Source layer (image)**

| 2 | 2 | 6 | 8 | 2 | 0 | 1 | 2 |
|---|---|---|---|---|---|---|---|
| 4 | 3 | 4 | 5 | 1 | 9 | 6 | 3 |
| 3 | 9 | 4 | 4 | 7 | 7 | 6 | 9 |
| 1 | 3 | 4 | 6 | 8 | 2 | 2 | 1 |
| 8 | 4 | 6 | 2 | 3 | 1 | 8 | 8 |
| 5 | 8 | 9 | 0 | 1 | 0 | 2 | 3 |
| 9 | 2 | 6 | 6 | 3 | 6 | 2 | 1 |
| 9 | 8 | 8 | 2 | 6 | 3 | 4 | 5 |

**Feature map (activation map)**

| $w_1$ | $w_2$ | $w_3$ |
|---|---|---|
| $w_4$ | $w_5$ | $w_6$ |
| $w_7$ | $w_8$ | $w_9$ |

**Kernel (filter)**

**How do we know which kernels to use?**

Kernels are learnt and initialized randomly during training the CNN learns spatial features

**Fully connected**

$$\sum_{i \in \text{image}}^{W \times H \times C} \mathbf{x}_i \mathbf{w}_i$$

**Locally connected shared weights**

$$\sum_{i \in 3 \times 3}^{W \times H \times C} \mathbf{x}_i \mathbf{w}_i$$

SURF

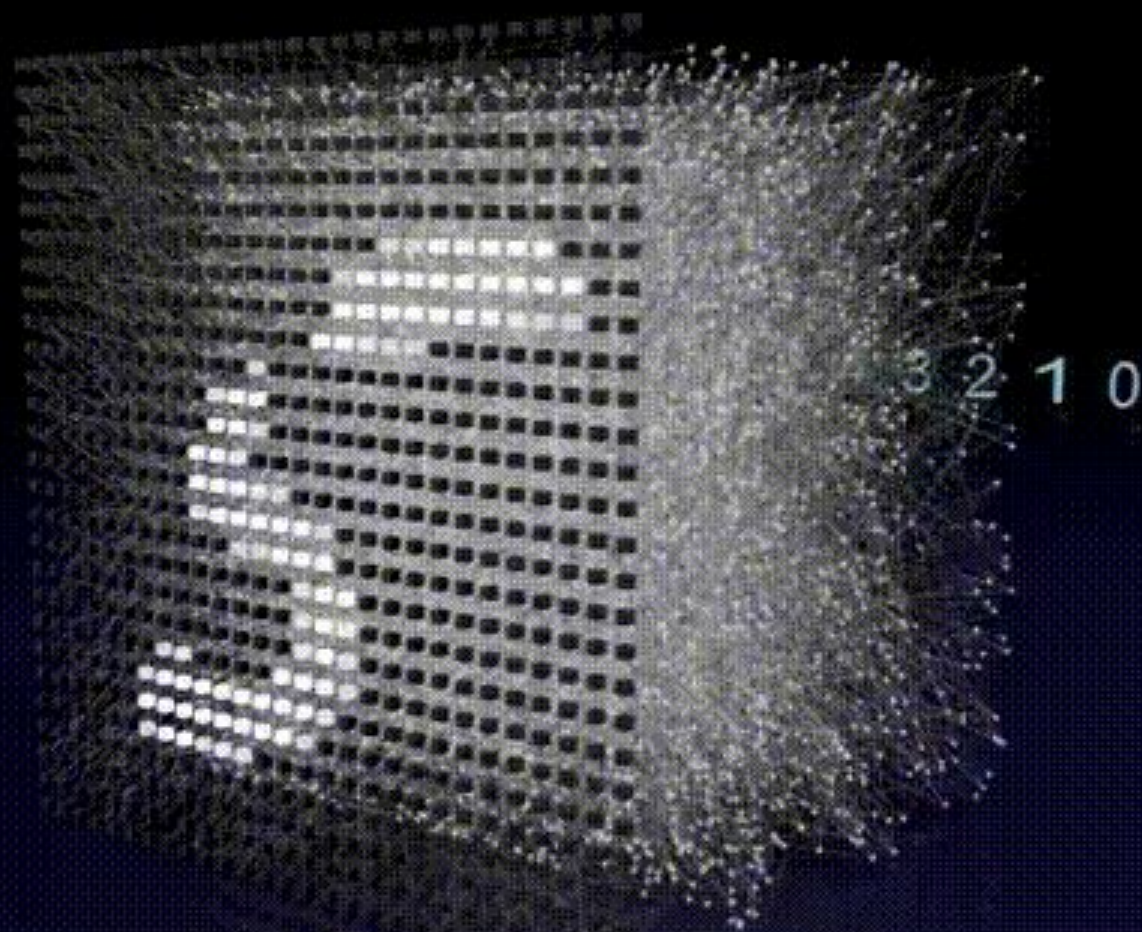# Pooling and dilated convolutions



We could discard information gradually by max pooling
Keep the strongest signal, works better than average pooling
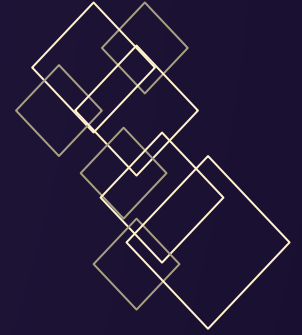Nowadays, we use a bigger stride for dimensionality reduction

# Pooling and dilated convolutions



We could discard information gradually by max pooling
Keep the strongest signal, works better than average pooling
Nowadays, we use a bigger stride for dimensionality reduction

Dilated convolutions can be used if you expect your images to have information which are spatially far from each other

**03.** CNN Cases
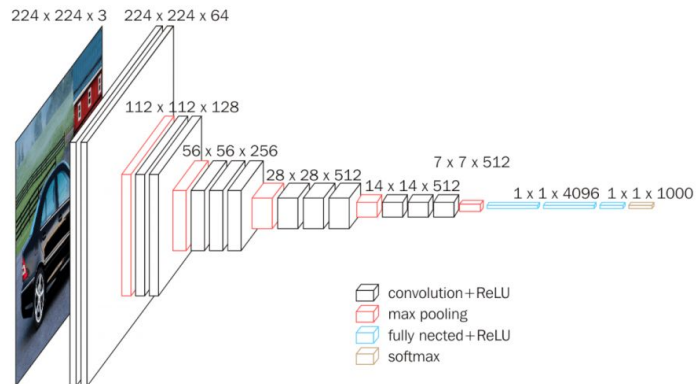
# AlexNet (2012)

# AlexNet (2012)



Images become smaller and number of filters increases as we go deeper

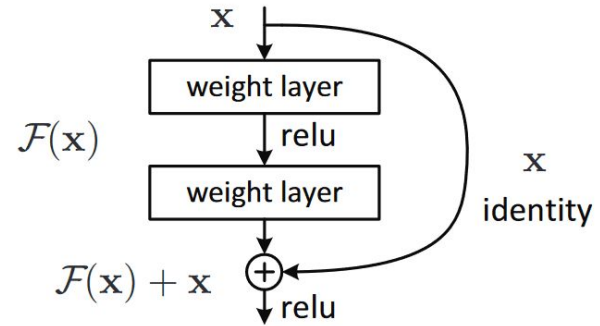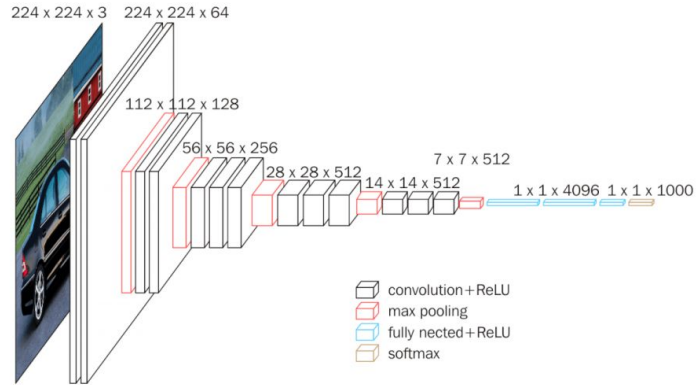8 layers with dropout and ReLU activations trained for 6 days on 2 GPUs

# VGG16 and ResNet (2014, 2015)



ImageNet Classification Error (Top 5)

- 2011 (XRCE): 26,0
- 2012 (AlexNet): 16,4
- 2013 (ZF): 11,7
- 2014 (VGG): 7,3
- 2014 (GoogLeNet): 6,7
- Human: 5,0
- 2015 (ResNet): 3,6
- 2016 (GoogLeNet-v4): 3,1

Images become smaller and number of filters increases as we go deeper

8 layers with dropout and ReLU activations trained for 6 days on 2 GPUs

SURF

# VGG16 and ResNet (2014, 2015)



224 x 224 x 3    224 x 224 x 64

112 x 112 x 128

56 x 56 x 256

28 x 28 x 512

14 x 14 x 512

7 x 7 x 512

1 x 1 x 4096    1 x 1 x 1000

convolution+ReLU
max pooling
fully nected+ReLU
softmax

VGG16: construct large and deep models (120M params)

Vanishing gradient problem

SURF

# VGG16 and ResNet (2014, 2015)



224 x 224 x 3  224 x 224 x 64

112 x 112 x 128

56 x 56 x 256

28 x 28 x 512

14 x 14 x 512

7 x 7 x 512

1 x 1 x 4096  1 x 1 x 1000

convolution+ReLU
max pooling
fully nected+ReLU
softmax

$\mathbf{x}$

weight layer

$\mathcal{F}(\mathbf{x})$          relu

weight layer

$\mathbf{x}$
identity

$\mathcal{F}(\mathbf{x}) + \mathbf{x}$   relu

VGG16: construct large and deep models (120M params)
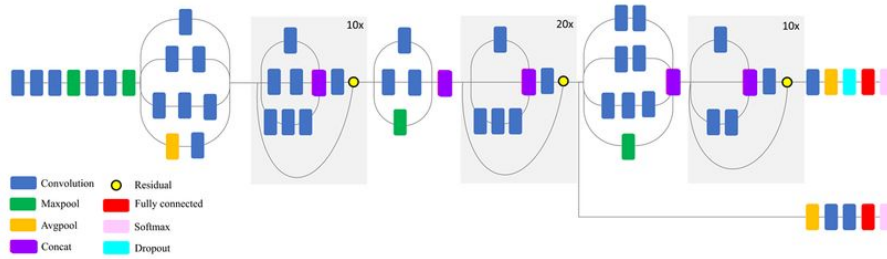
Vanishing gradient problem

SURF

# ResNet (2015)



Convolution  Residual
Maxpool     Fully connected
Avgpool     Softmax
Concat      Dropout

$\mathcal{F}(\mathbf{x})$

weight layer

relu

weight layer

$\mathbf{x}$
identity

$\mathcal{F}(\mathbf{x}) + \mathbf{x}$

relu

# ResNet (2015)



$$\mathcal{F}(x)$$

$$\mathcal{F}(x) + x$$

Convolution   Residual
Maxpool   Fully connected
Avgpool   Softmax
Concat   Dropout
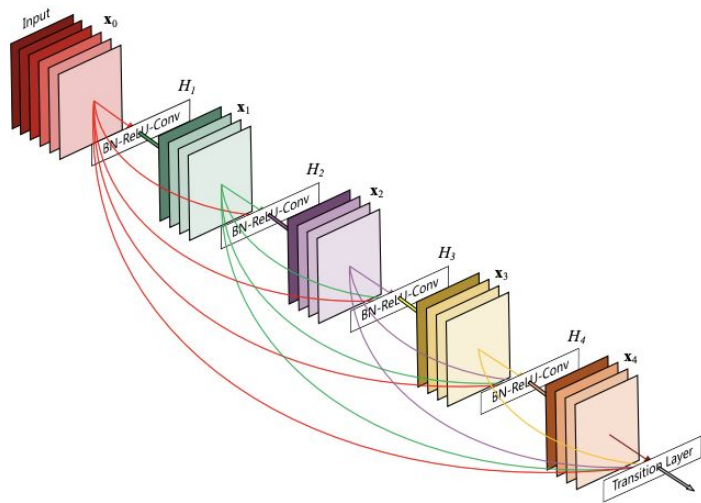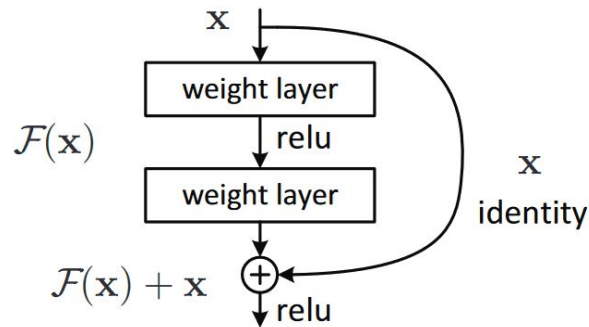
x
weight layer
relu
weight layer
x
identity
relu

ResNets: construct large and deep models with skip connections, able to train up to 152 layers (!)

Higher abstraction and less nuisance from vanishing or exploding gradients
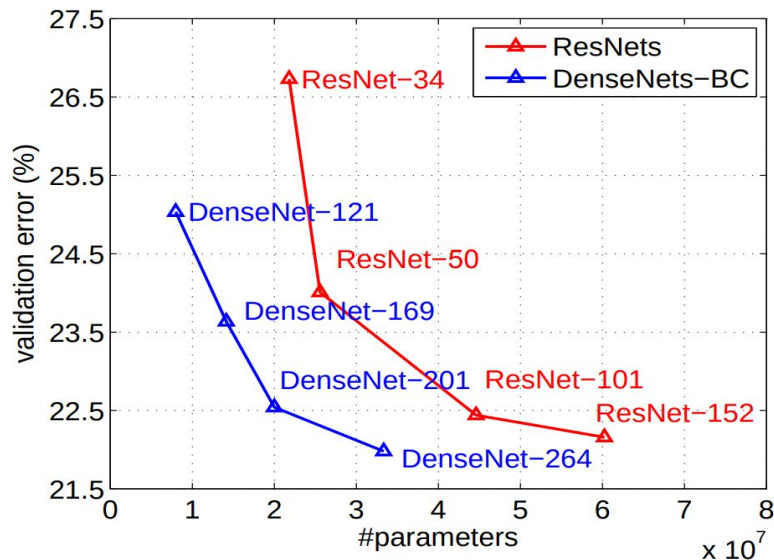
SURF

# ResNet (2015)



ResNets: construct large and deep models with skip connections, able to train up to 152 layers (!)

Higher abstraction and less nuisance from vanishing or exploding gradients

Early layers learn features that get progressively more abstract, we can **preserve** and **control** the flow of information with **skip connections**
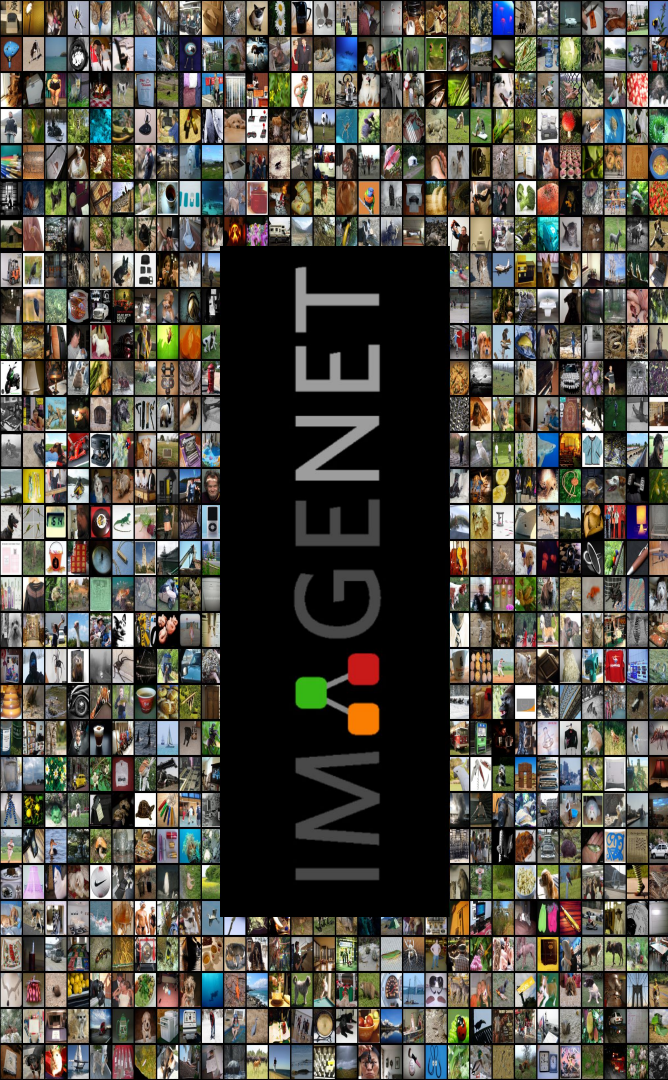
# DenseNet (2016)



ResNets: construct large and deep models with skip connections, able to train up to 152 layers (!)

Higher abstraction and less nuisance from vanishing or exploding gradients

Early layers learn features that get progressively more abstract, we can **preserve** and **control** the flow of information with **skip connections**
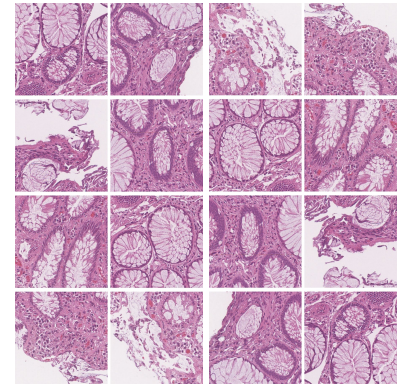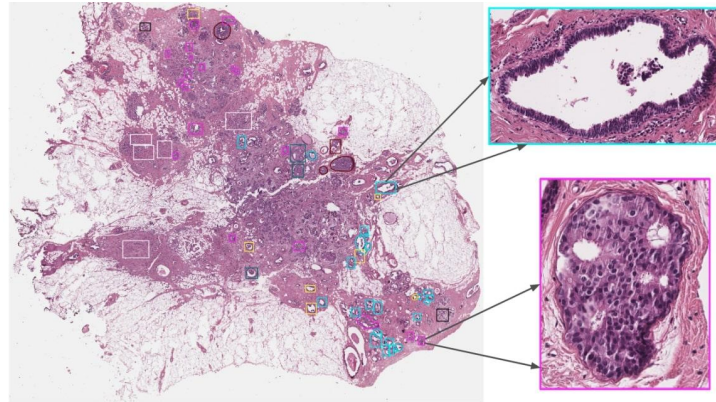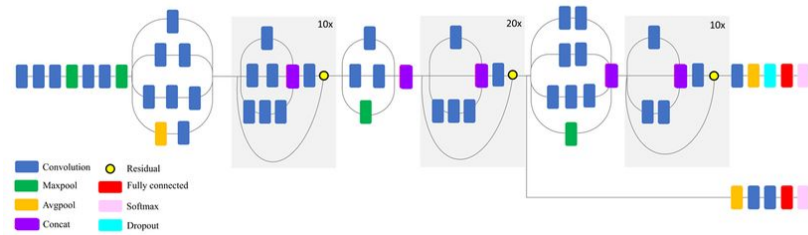
# DenseNet (2016)



ResNets: construct large and deep models with skip connections, able to train up to 152 layers (!)
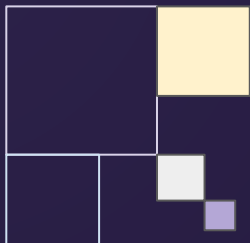
Higher abstraction and less nuisance from vanishing or exploding gradients

Early layers learn features that get progressively more abstract, we can **preserve** and **control** the flow of information with **skip connections**

# Transfer Learning

# Thank You

**High Performance Machine Learning Group**

SURF