



**girafe**  
ai

# Computer Vision with Deep Learning: an overview

**Radoslav Neychev**

- Different problems in Computer Vision
- Convolutional layer recap
- Main ideas of famous architectures in CV
- Datasets overview

# Different problems in Computer Vision – an overview

# Computer Vision tasks



- Classification
- Localization
- Detection
- Semantic Segmentation
- Instance Segmentation
- Image Generation
- Adversarial Attacks
- etc..

# Classification

**Input** image ( $C \times H \times W$ )



Cat

**Output**

class label (int)

# Classification

**Input image (C x H x W)**



Cat

## Application

- most image recognition tasks
- as a sub-task:
  - detection / localization
  - segmentation
  - prediction
  - ...

# Localization

**Input image (C x H x W)**



Cat

**Output**

class label (int)

BB\* coordinates (4 ints)

**Assumption**

one object per image

\* BB - bounding box

## Input image (C x H x W)



Cat

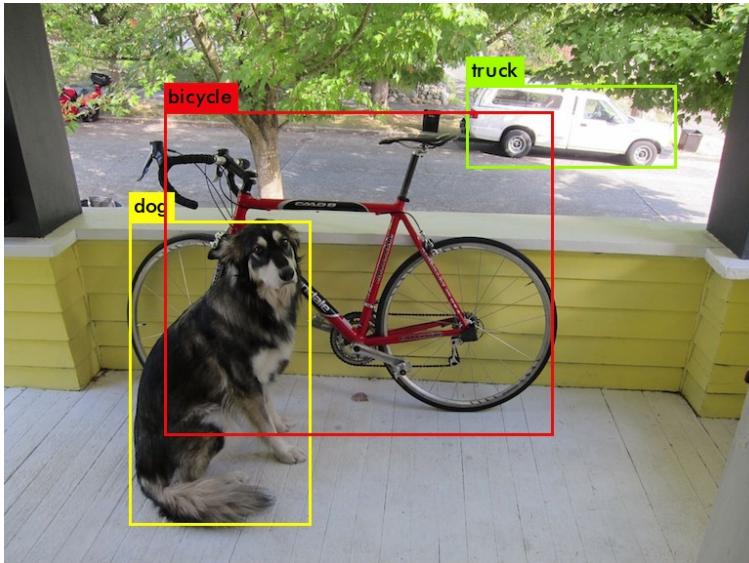
## Application

- smartphone interaction
- biometrics
- gesture recognition
- robotics
- object tracking
- ...

\* BB - bounding box

# Detection

**Input image (C x H x W)**



**Output**

class label (int)

confidence (0-1 float)

BB (4 ints)

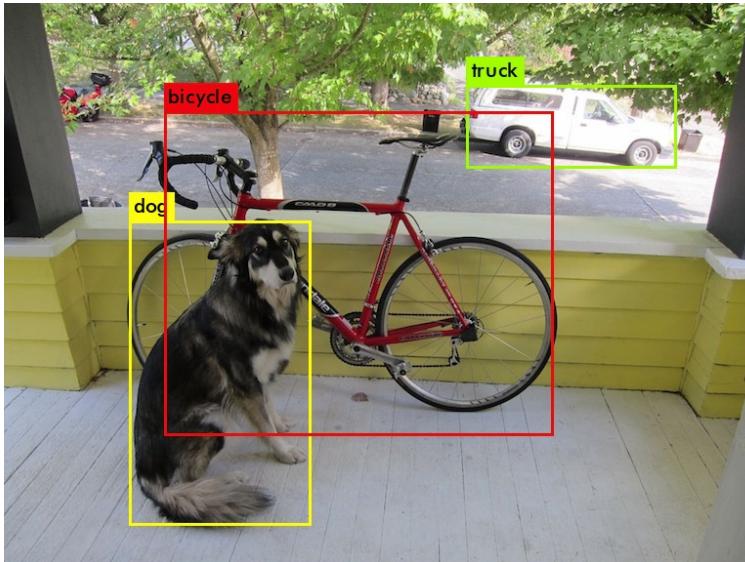
**Assumption**

multiple objects per image

} N times

# Detection

## Input image ( $C \times H \times W$ )



## Application

- biometrics
- surveillance
- robotics
- self-driving cars
- image enhancement
- ...

# Semantic Segmentation

**Input image (C x H x W)**

This image is CC0 public domain

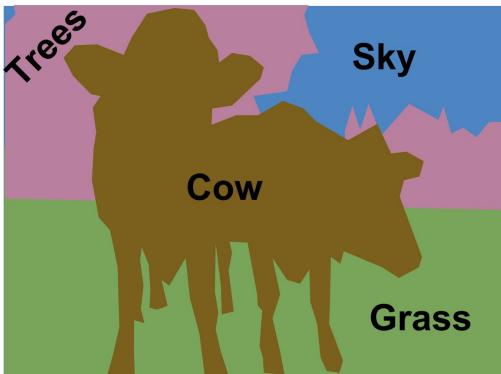


**Output**

**mask of classes (1 x H x W)**

**Assumption**

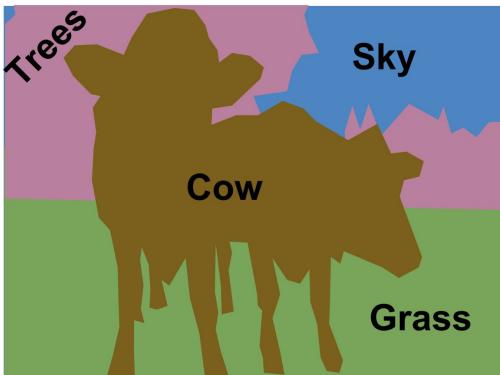
just class label for each pixel



# Semantic Segmentation

## Input image ( $C \times H \times W$ )

[This image is CC0 public domain](#)

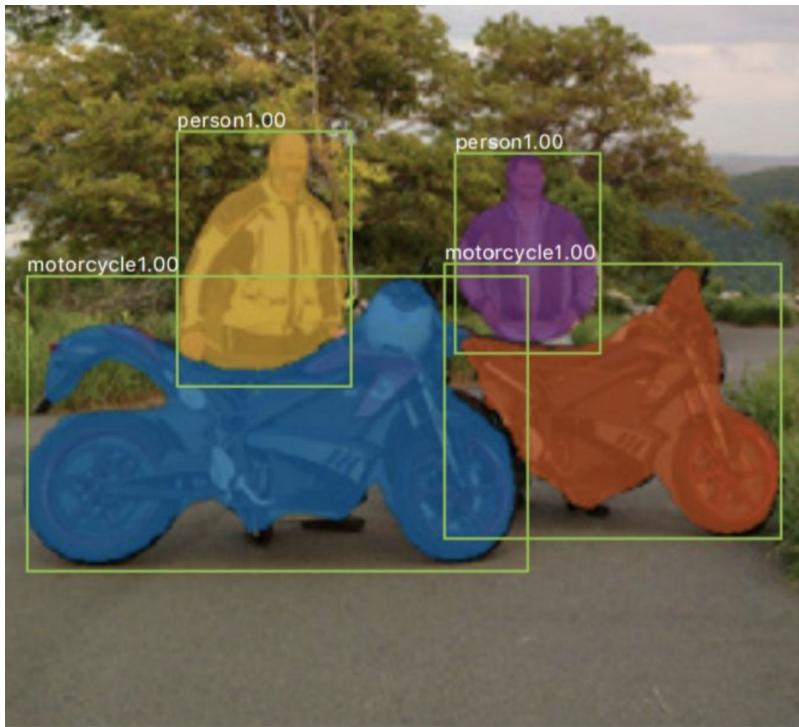


## Application

- biometrics
- robotics
- self-driving cars
- image enhancement
- image super-resolution
- video recognition
- ...

# Instance Segmentation

**Input image (C x H x W)**



**Output**

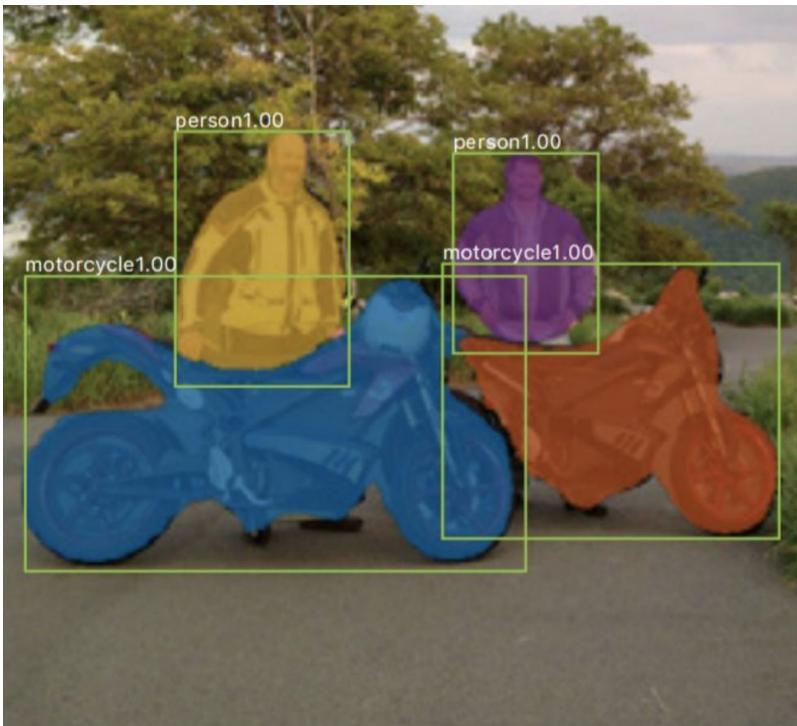
mask of class instances  
(detection + 1 x H x W)

**Assumption**

separated masks for each distinct object

# Instance Segmentation

## Input image (C x H x W)

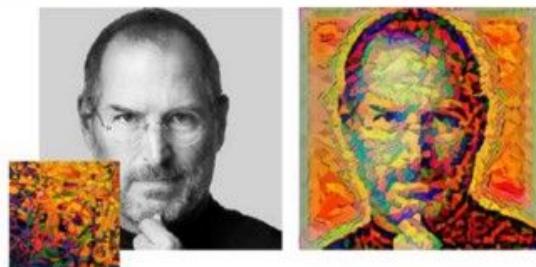


## Application

- biometrics
- robotics
- self-driving cars
- image enhancement
- image super-resolution
- video recognition
- ...

# Image Generation

**Input** image ( $C \times H \times W$ )



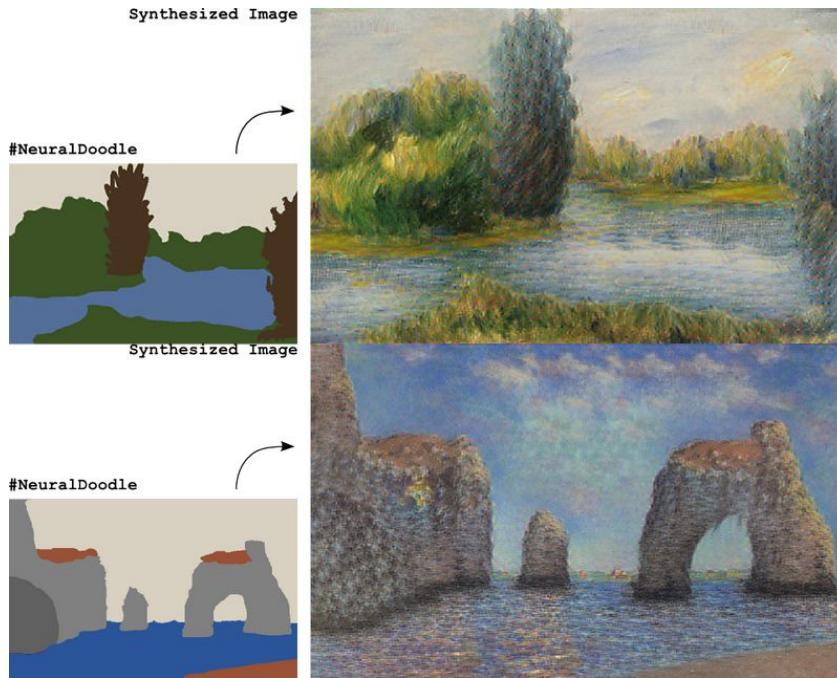
**Output**

image(s) of size  $C \times H \times W$

[image credit](#)

# Image Generation

## Input image ( $C \times H \times W$ )



## Application

- data augmentation
- dataset synthesis
- art generation
- ...

[image credit](#)

# Adversarial Attacks

**Input** image ( $C \times H \times W$ )



$+ .007 \times$



$x$   
“panda”  
57.7% confidence

$\text{sign}(\nabla_x J(\theta, x, y))$   
“nematode”  
8.2% confidence

=



$x +$   
 $\epsilon \text{sign}(\nabla_x J(\theta, x, y))$   
“gibbon”  
99.3 % confidence

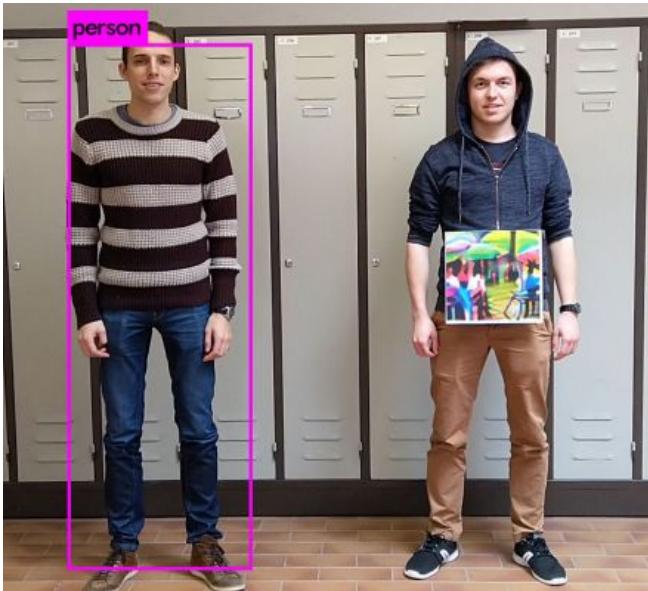
**Output**

image(s) of size  $C \times H \times W$

[image credit](#)

# Adversarial Attacks

**Input image (C x H x W)**



## Application

- data augmentation
- model testing
- surveillance protection
- ...

# Computer Vision tasks: Summary

What we have discussed today:

- Classification
- Localisation
- Detection
- Semantic Segmentation
- Instance Segmentation
- Image Generation
- Adversarial Attacks

# Computer Vision tasks: Summary

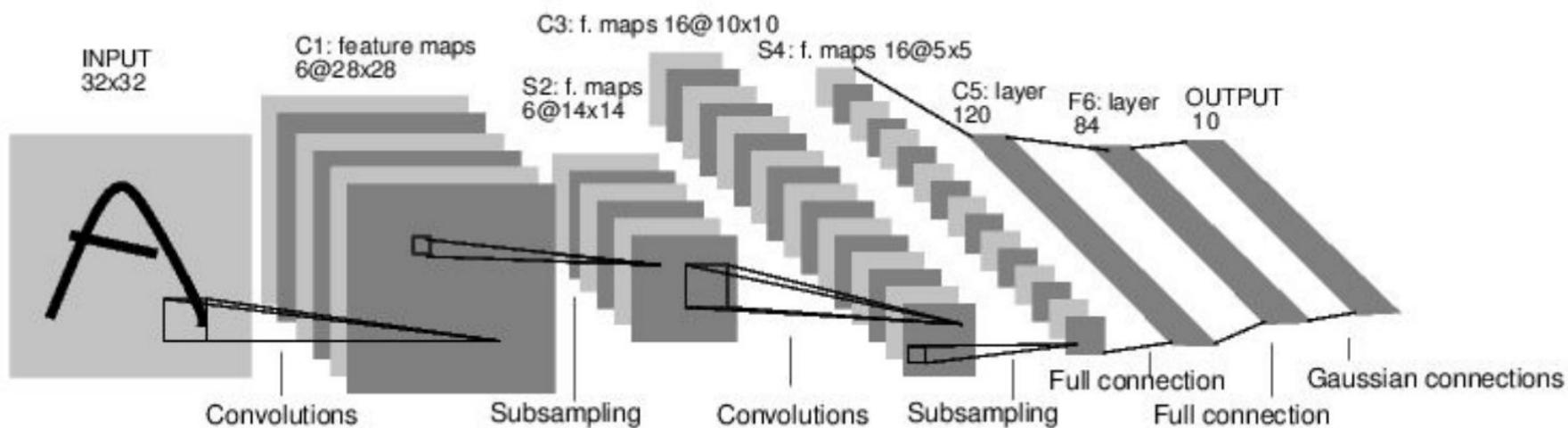
What we have discussed today:

- Classification
- Localisation
- Detection
- Semantic Segmentation
- Instance Segmentation
- Image Generation
- Adversarial Attacks

What else you can do:

- Image captioning
- Image encoding
- Image morphing
- Image upscaling
- Video object tracking
- Video interpolation
- Online learning
- Scene reconstruction
- ...

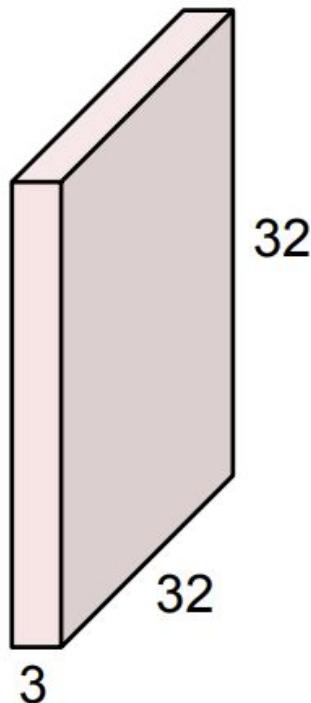
# Convolutional Layer – recap



[LeNet-5, LeCun 1998]

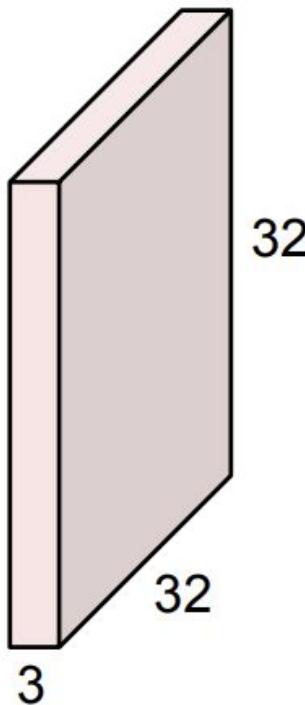
# Convolutional layer

32x32x3 image



# Convolutional layer

32x32x3 image



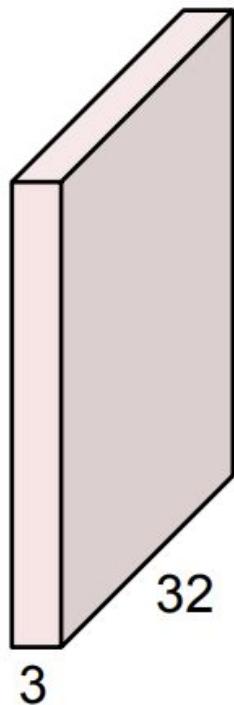
5x5x3 filter



**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolutional layer

32x32x3 image



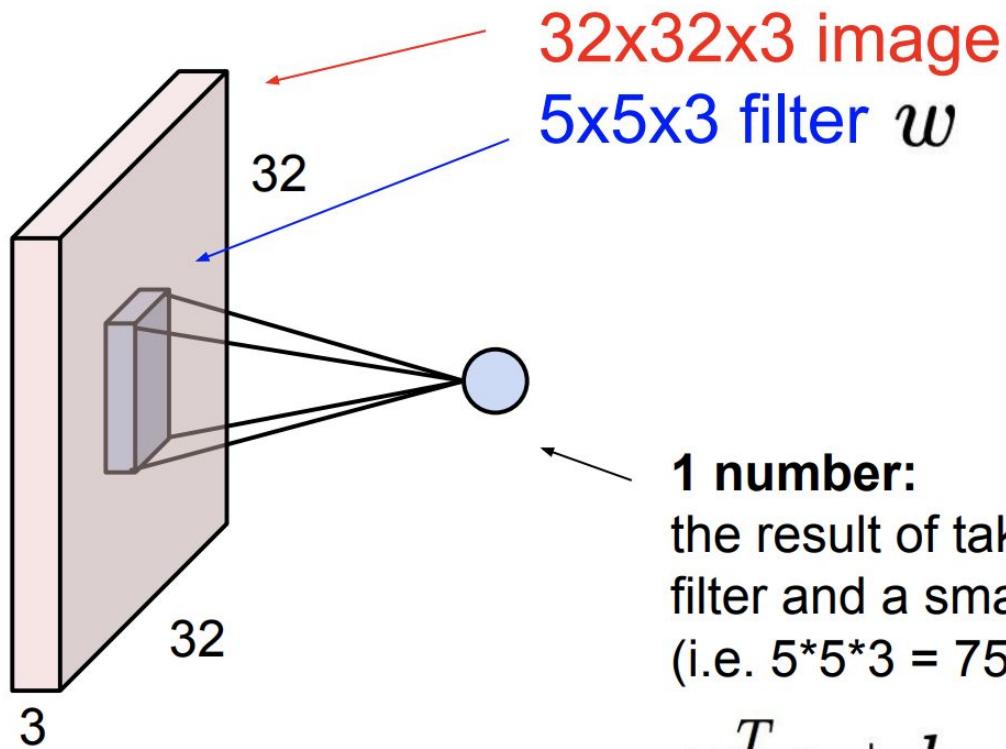
←  
5x5x3 filter



Filters extend the *depth* of the original image

**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

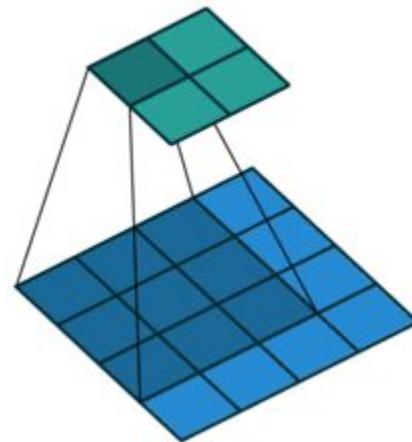
# Convolutional layer



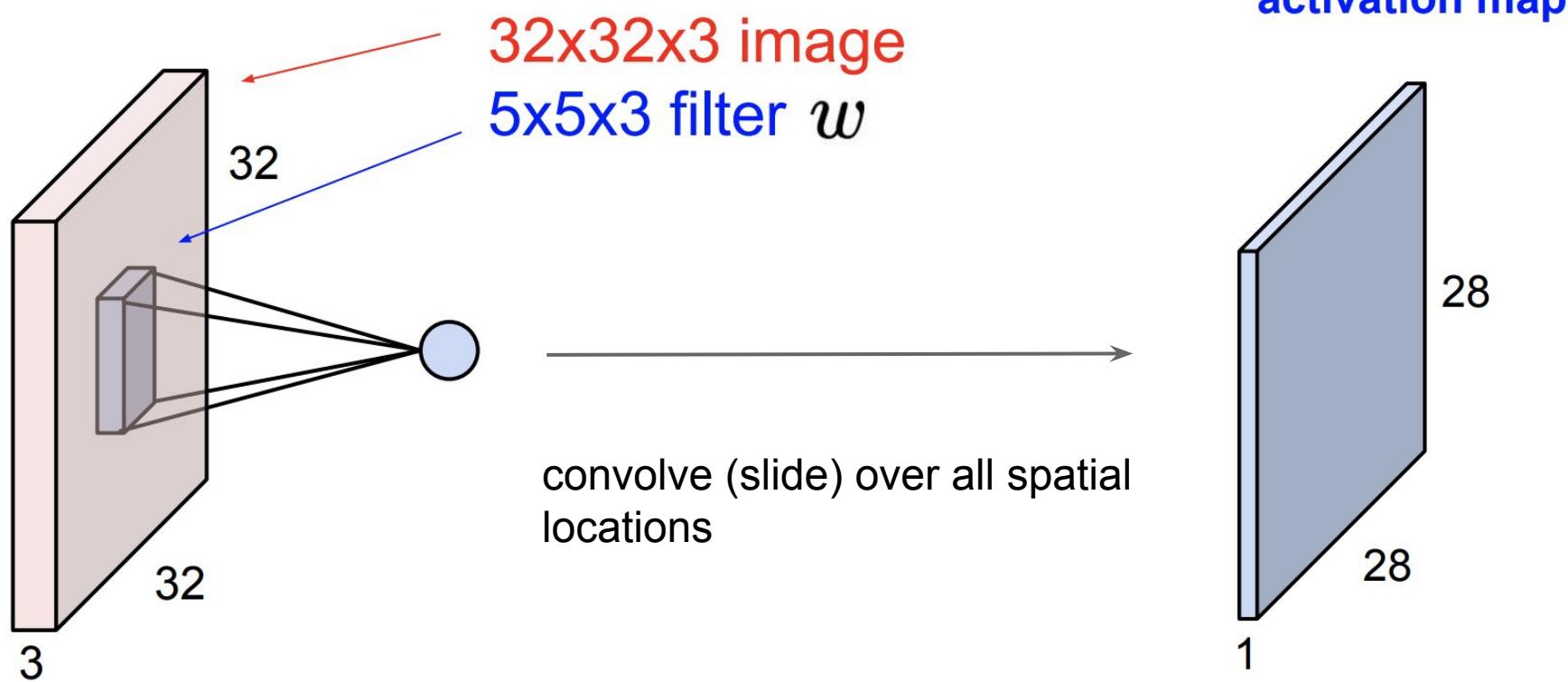
# Convolutional layer



# Convolutional layer

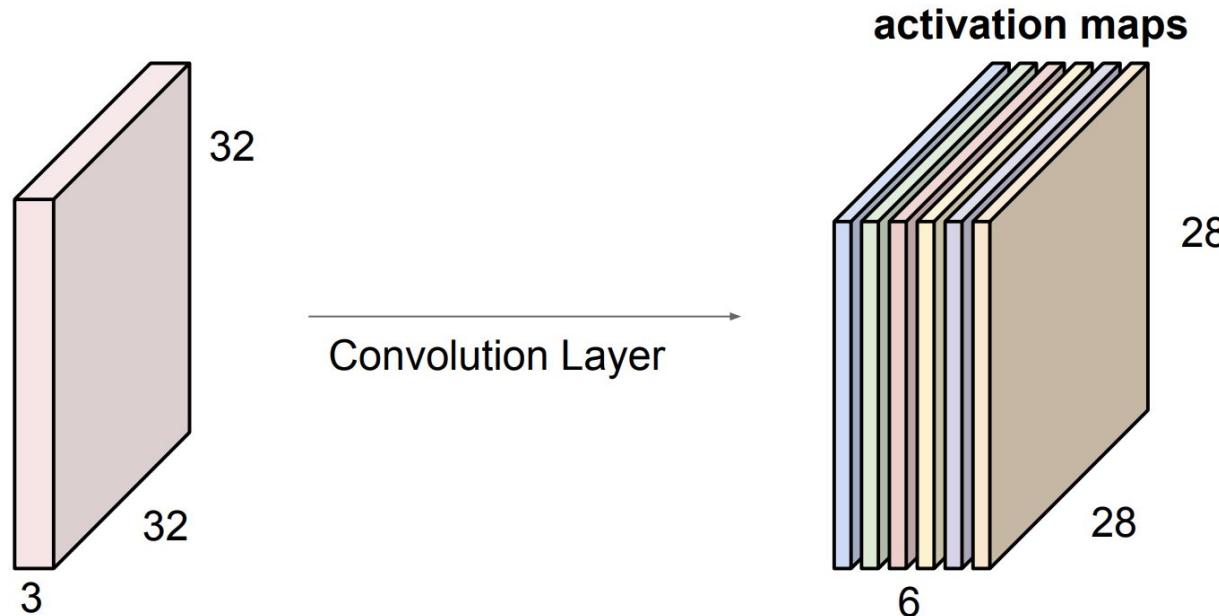


# Convolutional layer



# Convolutional layer

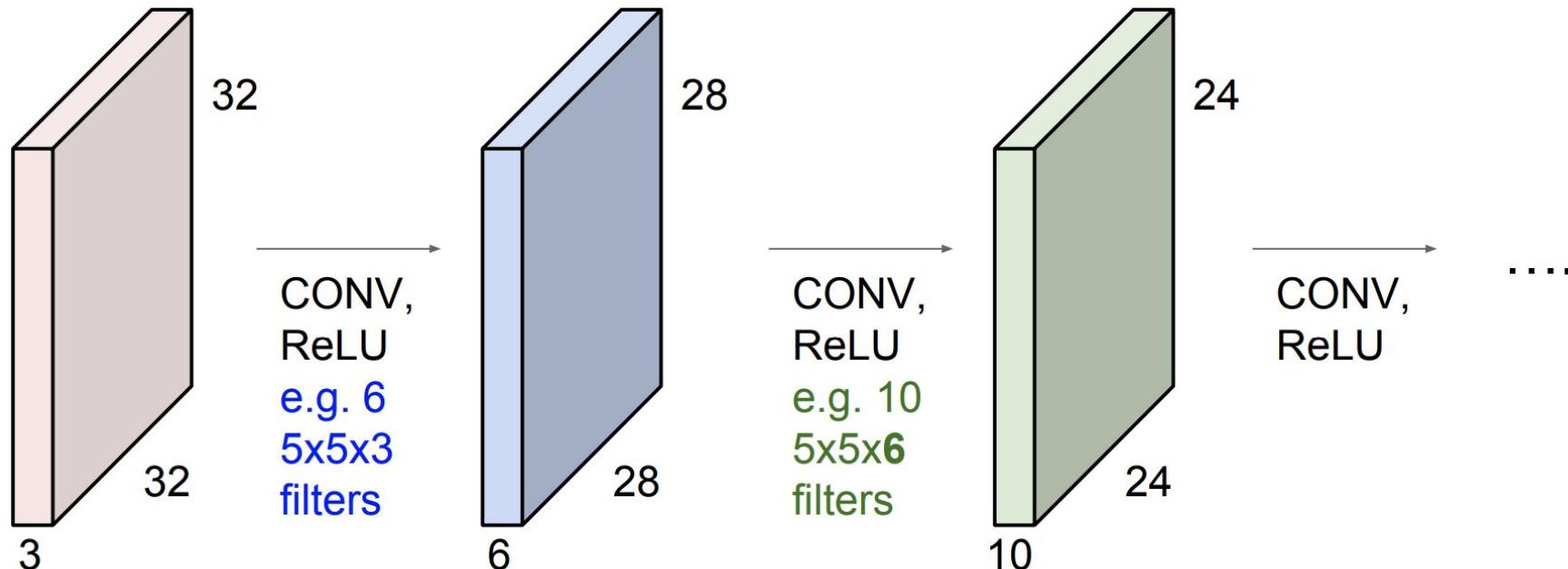
For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size 28x28x6!

# Convolutional layer

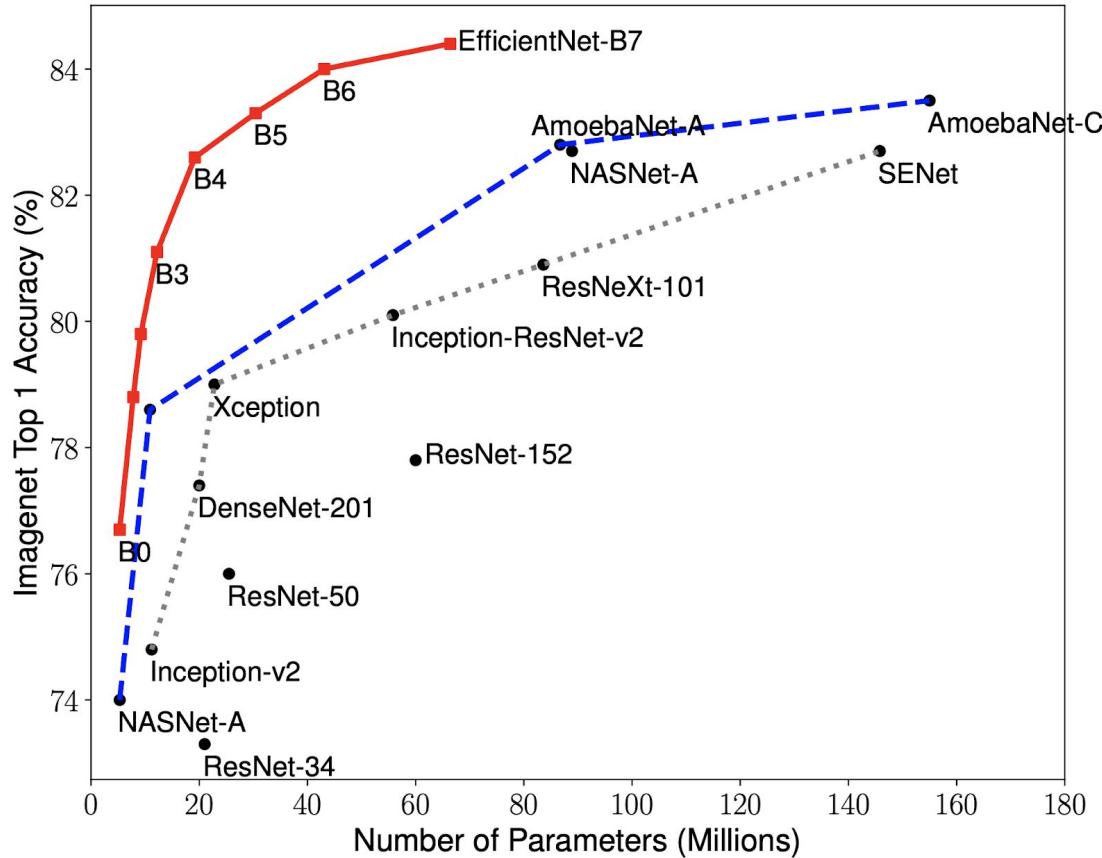
**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



# Famous NN Architectures in Computer Vision

# Overview

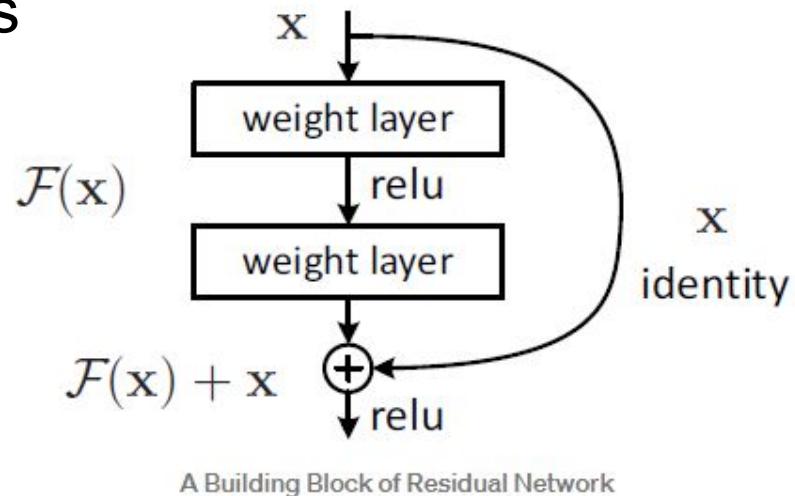
- Usually Deeper networks show better results
- But there are problems:
  - vanishing gradients
  - exploding gradients
  - biased activations
  - memory
  - speed
  - overfitting



[image credit](#)

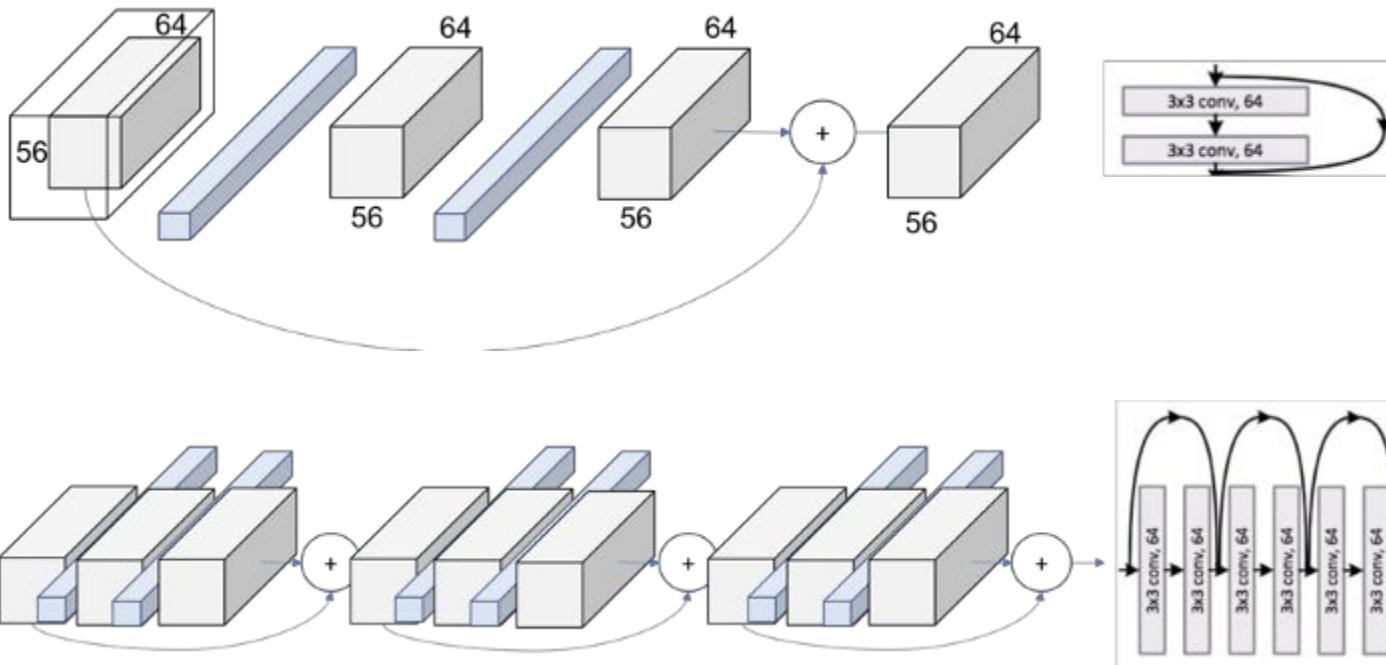
## Main feature – residual blocks:

- Allow training deeper networks
- Improve vanishing gradient problem
- Look like LSTM

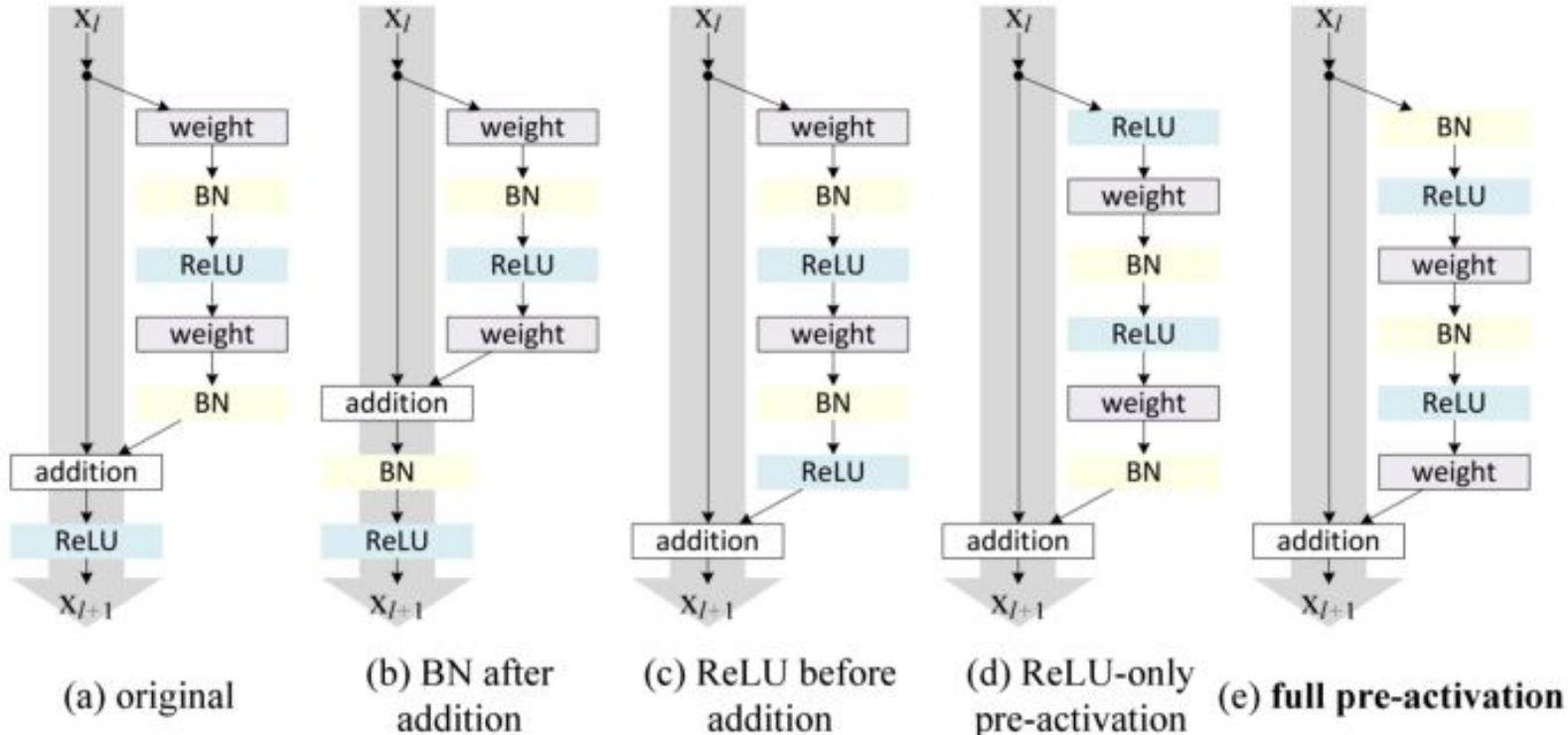


A Building Block of Residual Network

# ResNet [2015]

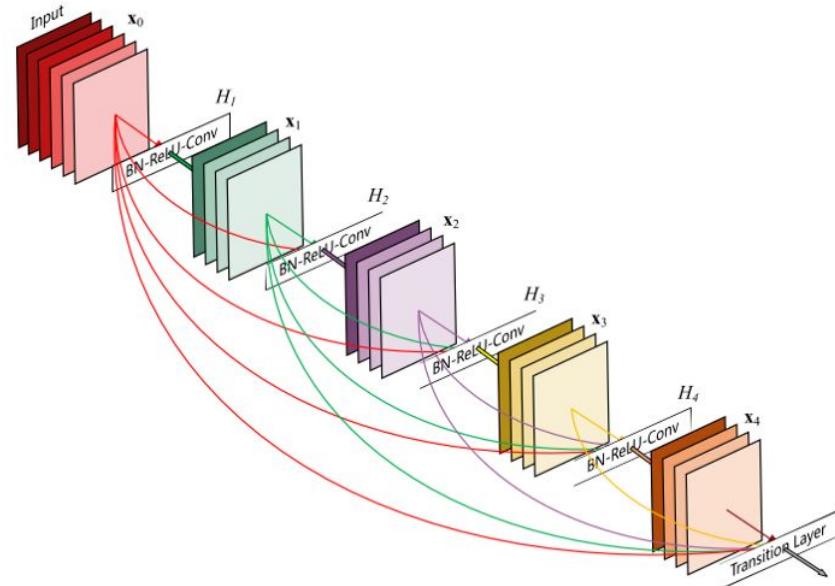


# ResNet [2015]

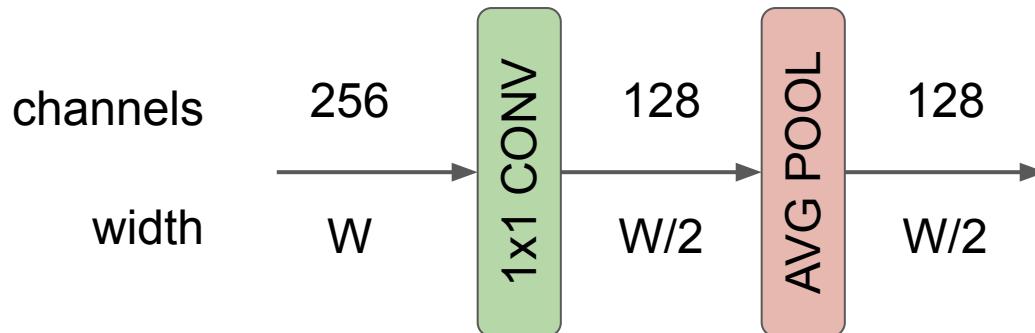


## Features:

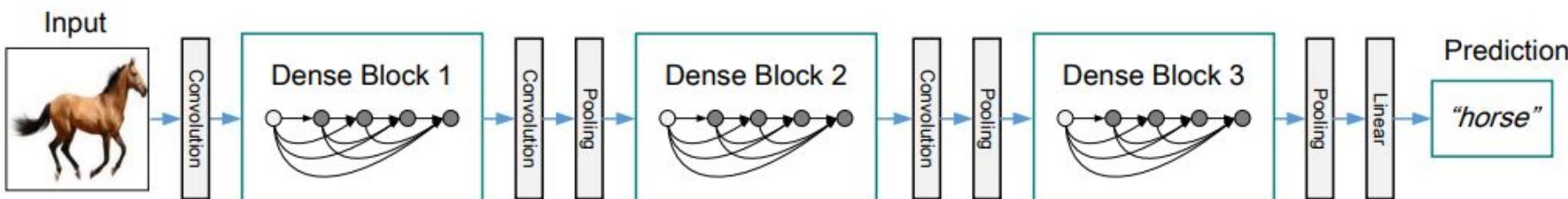
- Connects each layer to every other layer in a feed-forward fashion
- Concatenate outputs
- Improves vanishing gradient problem



# DenseNet [2016]



Transition block visualization

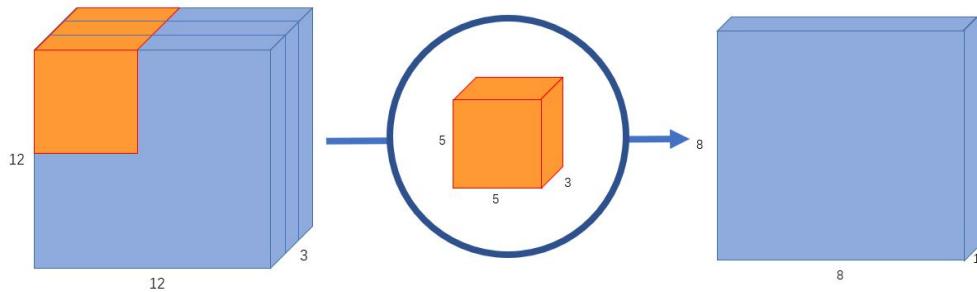


**Figure 2:** A deep DenseNet with three dense blocks. The layers between two adjacent blocks are referred to as transition layers and change feature-map sizes via convolution and pooling.

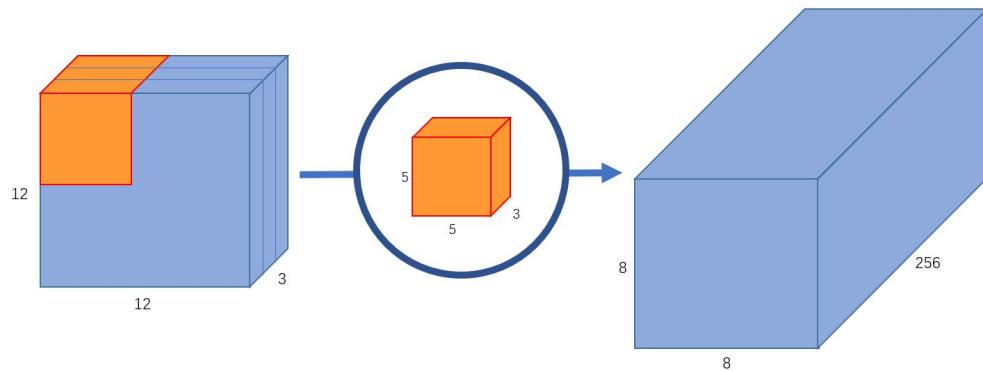
# Convolutions

## Normal convolution

- Uses all input channels:
- Is it redundant?



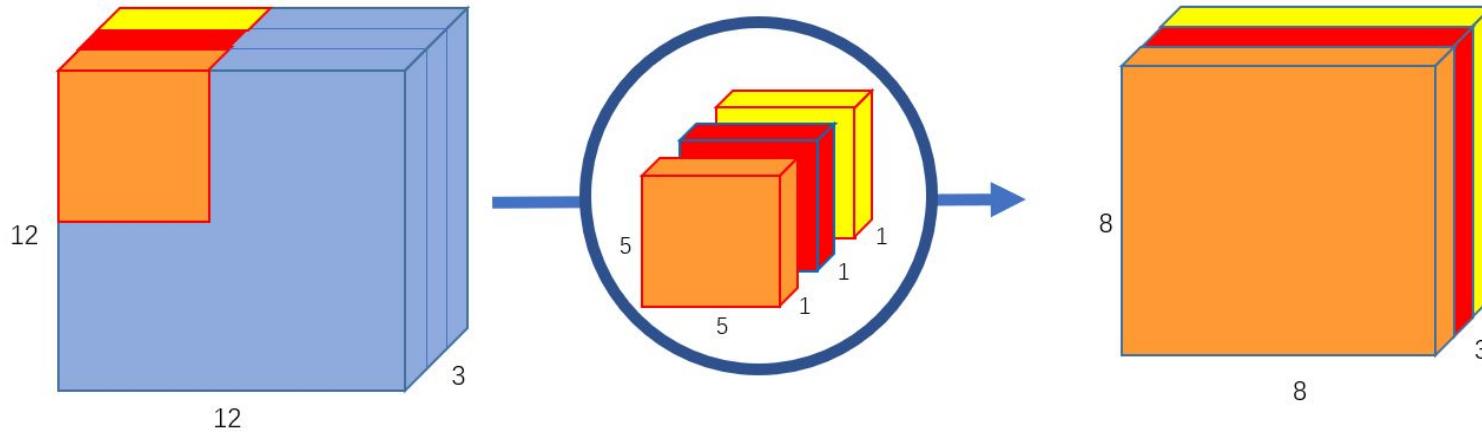
- Kernel may be factorized
  - channels: **depthwise**
  - spatially: **separable**



→ Less memory / FLOPs

Kernel size: C x K x K x H

# Convolutions: Depthwise

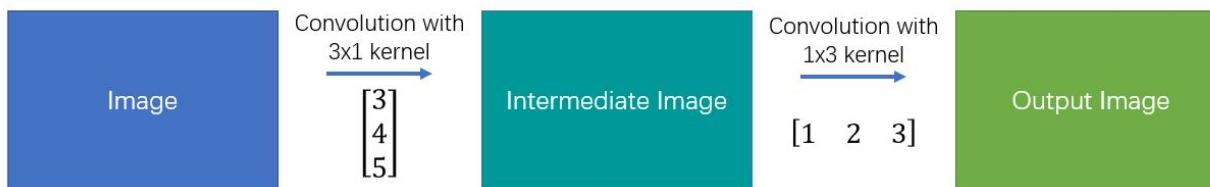


Keep channel number

Kernel size:  $C \times K \times K$

# Convolutions: Separable

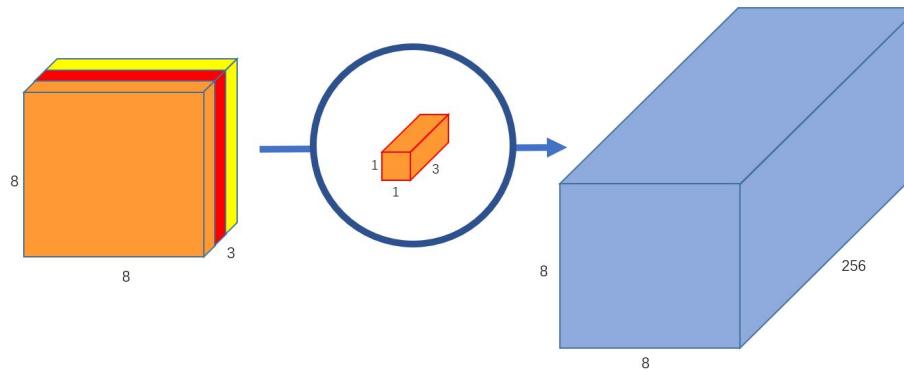
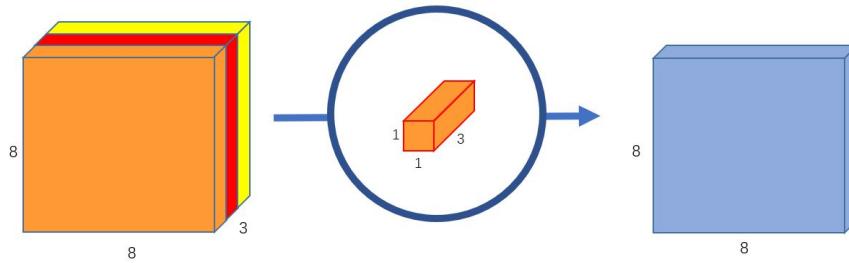
$$\begin{bmatrix} 3 & 6 & 9 \\ 4 & 8 & 12 \\ 5 & 10 & 15 \end{bmatrix} = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} \times [1 \quad 2 \quad 3]$$



Keep channel number

Filter size:  $C \times (K + K)$

# Convolutions: Pointwise



Change only channel number

Filter size: C x 1 x 1 x H

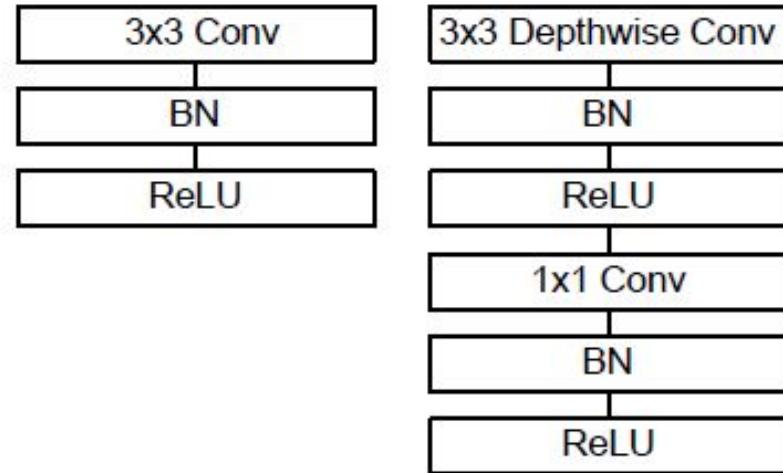
## Compare by kernel size:

- Normal:  $C \times K \times K \times H$
- Depthwise:  $C \times K \times K \rightarrow H$  times smaller
- Separable:  $C \times (K + K) \rightarrow 0.5 \times K \times H$  times smaller
- Pointwise:  $C \times 1 \times 1 \times H \rightarrow K \times K$  times smaller

Can combine them!

## Features:

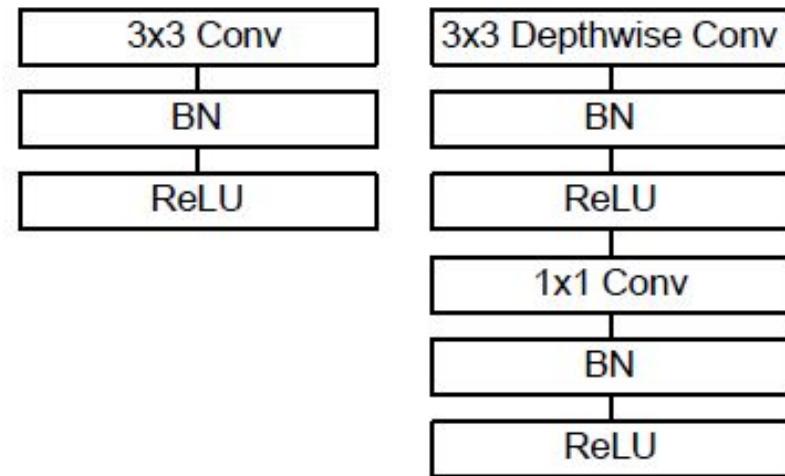
- Split convolution
  - “depthwise” part
  - “pointwise” part
- Less parameters
- Same performance
- Use in low-power devices



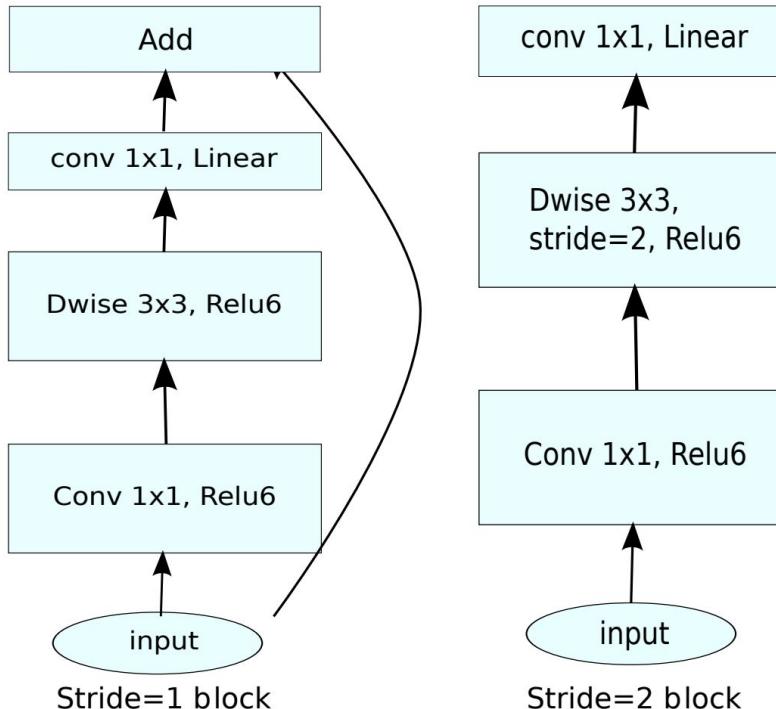
# MobileNet [2017]

Table 1. MobileNet Body Architecture

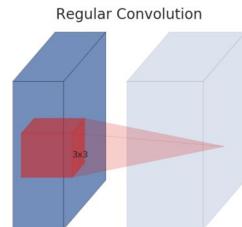
Type / Stride	Filter Shape	Input Size
Conv / s2	$3 \times 3 \times 3 \times 32$	$224 \times 224 \times 3$
Conv dw / s1	$3 \times 3 \times 32$ dw	$112 \times 112 \times 32$
Conv / s1	$1 \times 1 \times 32 \times 64$	$112 \times 112 \times 32$
Conv dw / s2	$3 \times 3 \times 64$ dw	$112 \times 112 \times 64$
Conv / s1	$1 \times 1 \times 64 \times 128$	$56 \times 56 \times 64$
Conv dw / s1	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 128$	$56 \times 56 \times 128$
Conv dw / s2	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 256$	$28 \times 28 \times 128$
Conv dw / s1	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 256$	$28 \times 28 \times 256$
Conv dw / s2	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 512$	$14 \times 14 \times 256$
5× Conv dw / s1	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
	$1 \times 1 \times 512 \times 512$	$14 \times 14 \times 512$
Conv dw / s2	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
Conv / s1	$1 \times 1 \times 512 \times 1024$	$7 \times 7 \times 512$
Conv dw / s2	$3 \times 3 \times 1024$ dw	$7 \times 7 \times 1024$
Conv / s1	$1 \times 1 \times 1024 \times 1024$	$7 \times 7 \times 1024$
Avg Pool / s1	Pool $7 \times 7$	$7 \times 7 \times 1024$
FC / s1	$1024 \times 1000$	$1 \times 1 \times 1024$
Softmax / s1	Classifier	$1 \times 1 \times 1000$



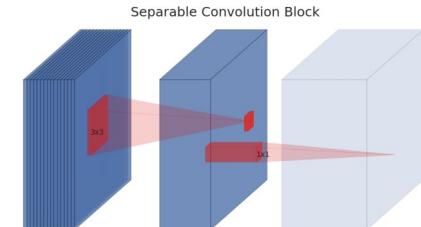
# MobileNet v2 [2018]



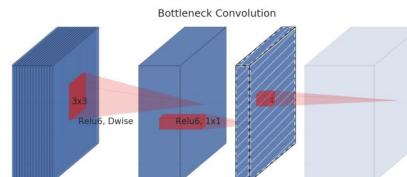
(a) Regular



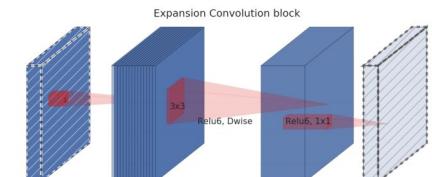
(b) Separable



(c) Separable with linear bottleneck



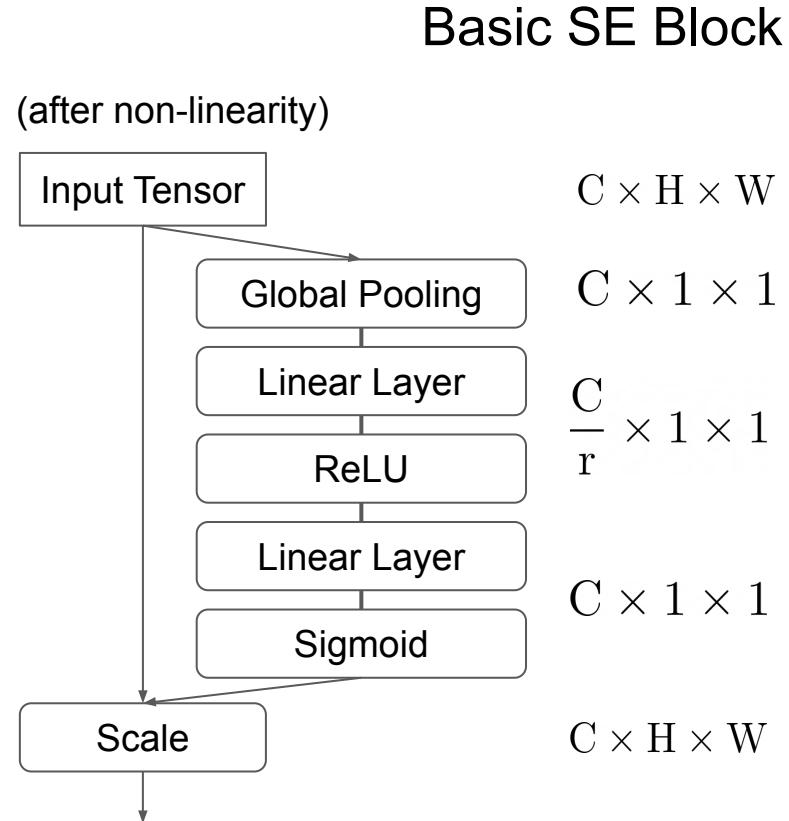
(d) Bottleneck with expansion layer



# Squeeze & Excitation Networks [2018]

## Features:

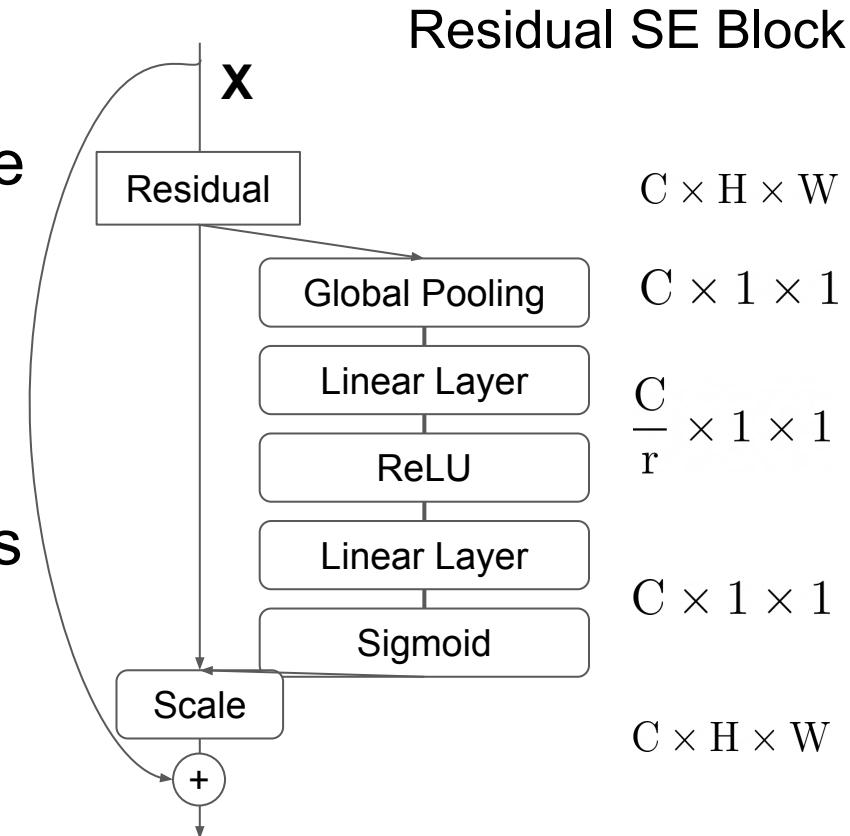
- Use global averages of feature maps for better classification
- Scale feature map channels
- “Excite” useful info
- Suppress redundant info
- Apply: deep CNNs’ later layers



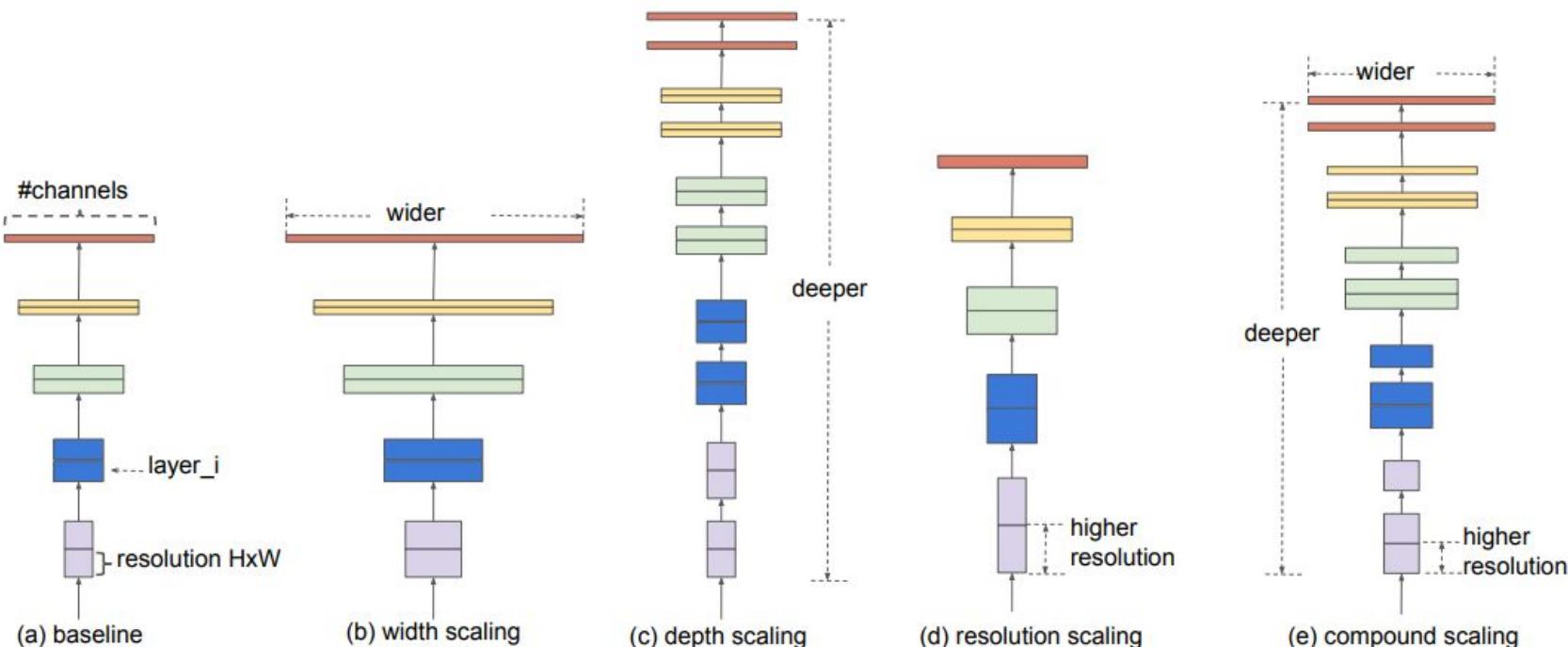
# Squeeze & Excitation Networks [2018]

## Features:

- Use global averages of feature maps for better classification
- Scale feature map channels
- “Excite” useful info
- Suppress redundant info
- Apply: deep CNNs’ later layers



# EfficientNet [2018]



**Figure 2. Model Scaling.** (a) is a baseline network example; (b)-(d) are conventional scaling that only increases one dimension of network width, depth, or resolution. (e) is our proposed compound scaling method that uniformly scales all three dimensions with a fixed ratio.

## Features:

- Scaling the model by:
  - depth (layers)
  - width (channels)
  - input resolution
- Compound scaling for better performance
- Neural architecture search

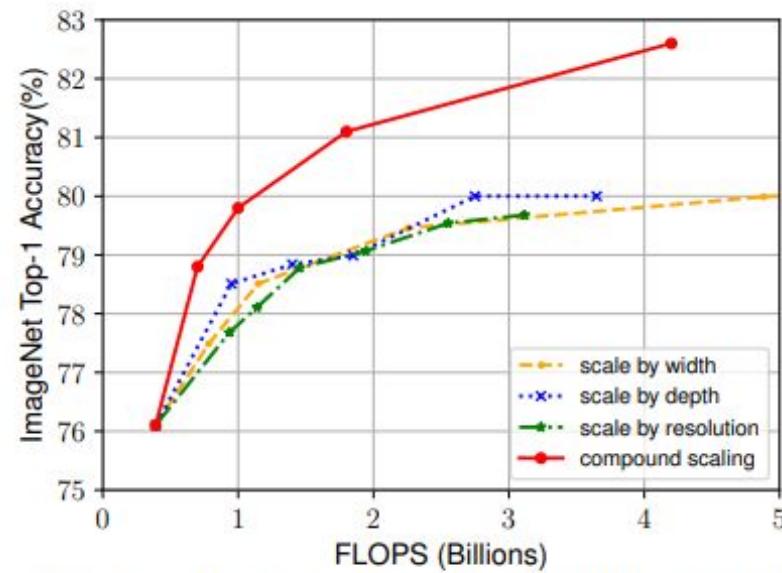


Figure 8. Scaling Up EfficientNet-B0 with Different Methods.

# Common Datasets

# Datasets

## CIFAR10 & CIFAR100

- 60K 32x32 colour images
- 10 classes
- 6K images per class  
or
- 100 classes
- 600 images per class

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



truck



## PASCAL VOC 2012

- (Visual Object Classes)
- developed since 2005
- 11,5k images
- 20 classes
- 31,5k detections
- 7k segmentations
- Images came from flickr.com

**Table 1** The VOC classes

Vehicles	Household	Animals	Other
Aeroplane	Bottle	Bird	Person
Bicycle	Chair	Cat	
Boat	Dining table	Cow	
Bus	Potted plant	Dog	
Car	Sofa	Horse	
Motorbike	TV/Monitor	Sheep	
Train			

The classes can be considered in a notional taxonomy



# Datasets

## ImageNet

- WordNet hierarchy
- 14.2m images
- 1000 classes
- 1m images with BBs
- diverse objects properties



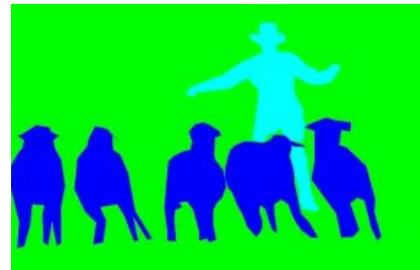
IMAGENET

<http://image-net.org/explore>

# Datasets

## MS COCO 2017

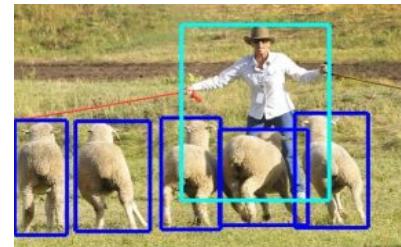
- 328k images
- 91 object categories
- 82 with more than 5,000 labeled instances
- 2,5m labeled instances
- 250k people with keypoints



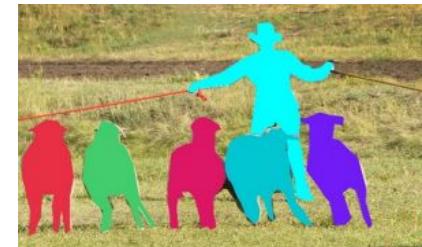
semantic segmentation



classification



object localization



instance segmentation

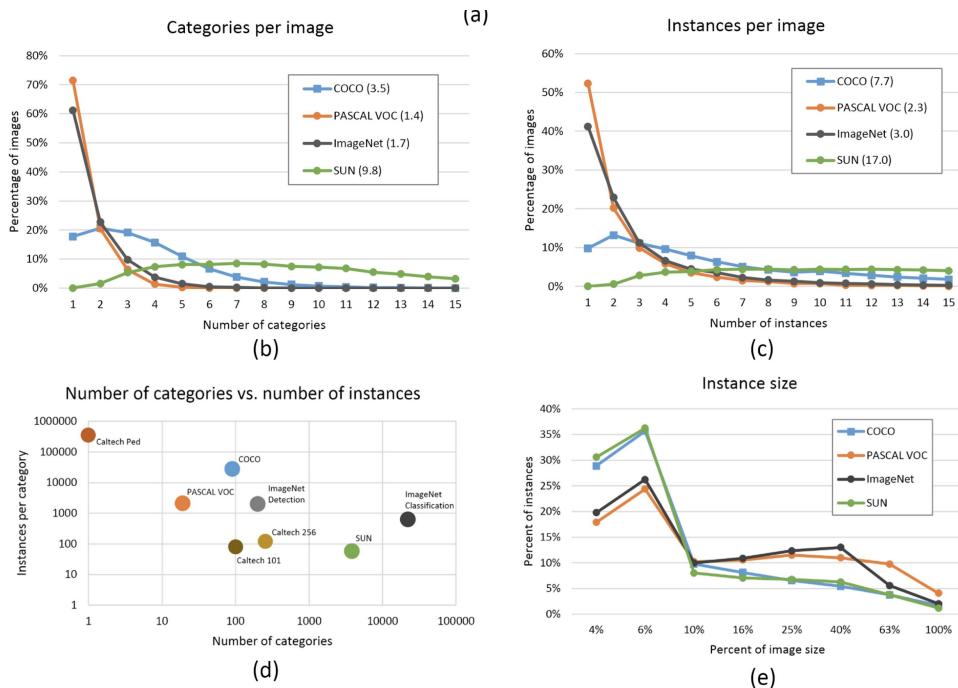
<https://arxiv.org/pdf/1405.0312.pdf>

<http://cocodataset.org/#explore>

# Datasets

## MS COCO 2017

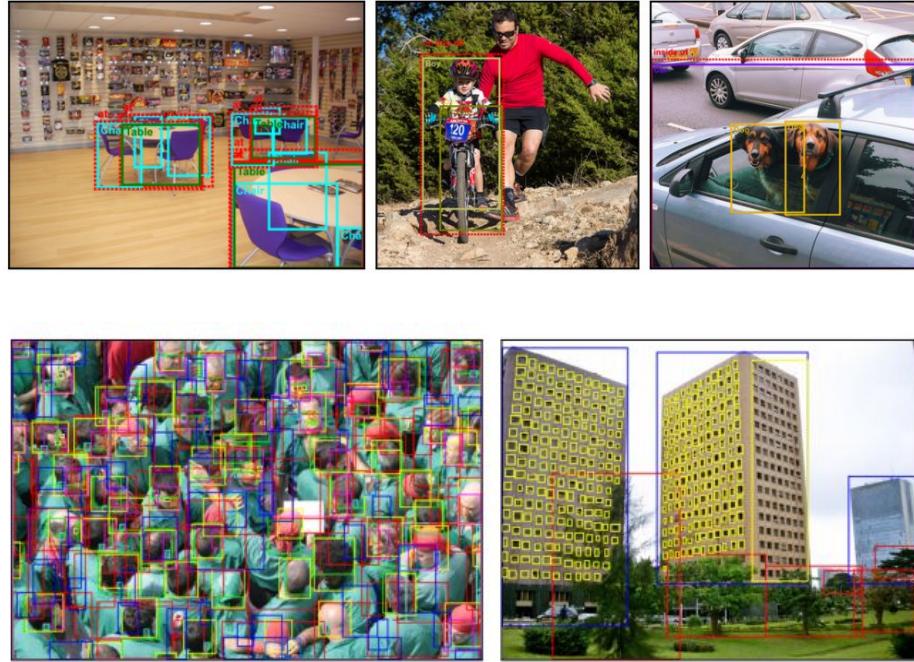
- 328k images
- 91 object categories
- 82 with more than 5,000 labeled instances
- 2.5m labeled instances
- 250k people with keypoints



<https://arxiv.org/pdf/1405.0312.pdf>  
<http://cocodataset.org/#explore>

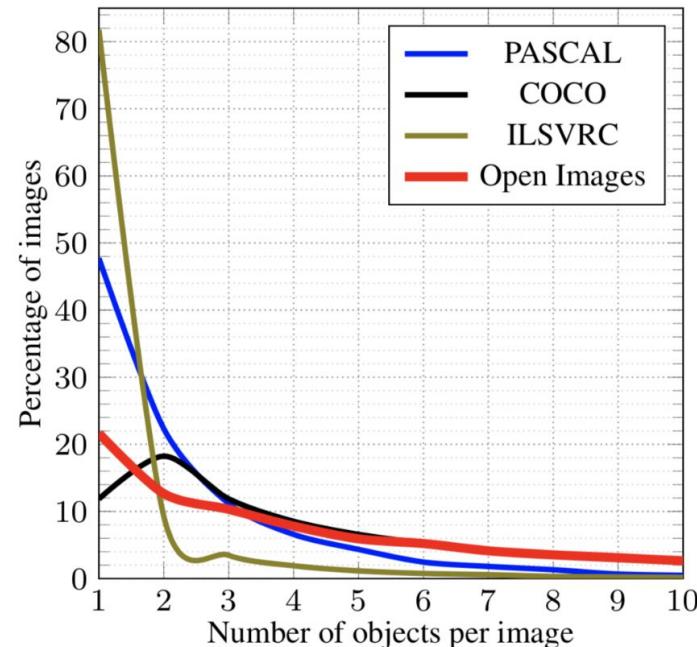
## Open Images

- Currently relevant challenge
- 15m boxes on 600 categories
- 2,8m instance segmentations on 350 categories
- 36,5m image-level labels on 20k categories
- 391k relationship annotations of 329 relationships



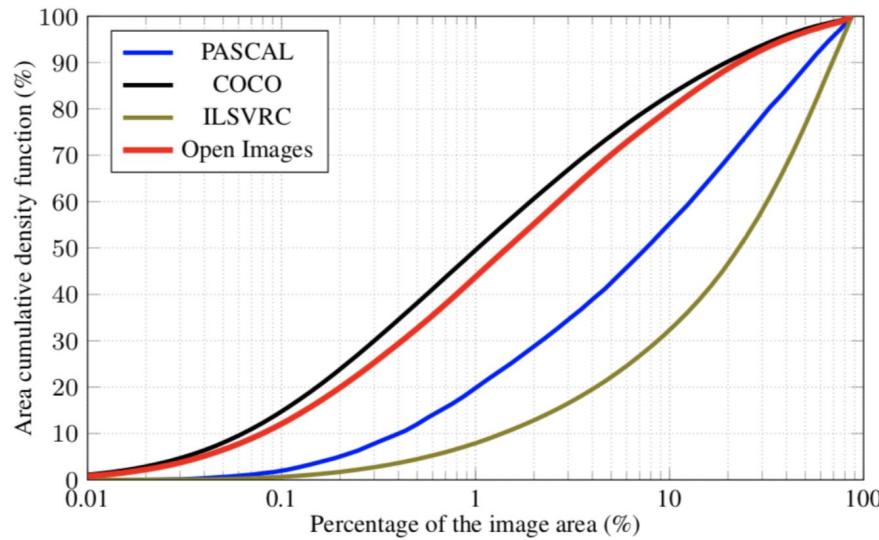
## Open Images

- Currently relevant challenge
- 15m boxes on 600 categories
- 2,8m instance segmentations on 350 categories
- 36,5m image-level labels on 20k categories
- 391k relationship annotations of 329 relationships



## Open Images

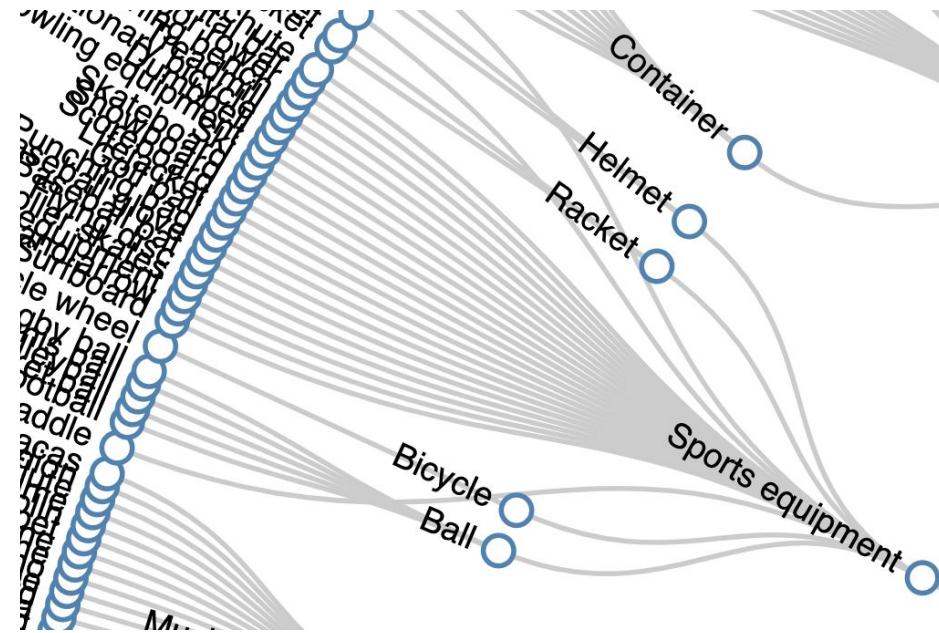
- Currently relevant challenge
- 15m boxes on 600 categories
- 2,8m instance segmentations on 350 categories
- 36,5m image-level labels on 20k categories
- 391k relationship annotations of 329 relationships



# Datasets

## Open Images

- Currently relevant challenge
- 15m boxes on 600 categories
- 2,8m instance segmentations on 350 categories
- 36,5m image-level labels on 20k categories
- 391k relationship annotations of 329 relationships



- Vanishing and exploding gradients are still present
- Careful model scaling is important
- Convolution layers can be factorized
- Datasets are getting better
- A lot of Computer Vision tasks yet to be solved