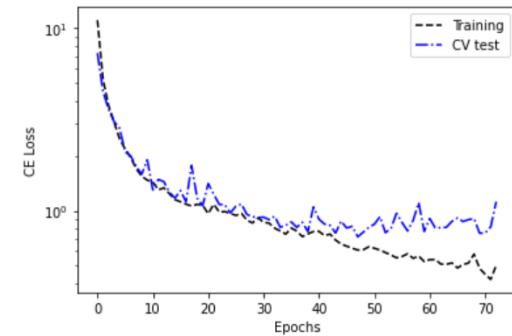
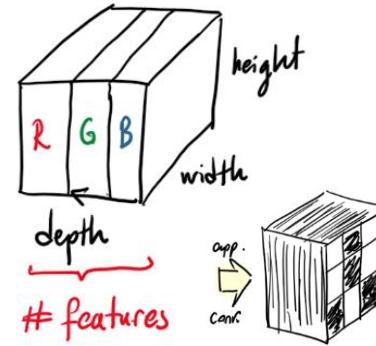
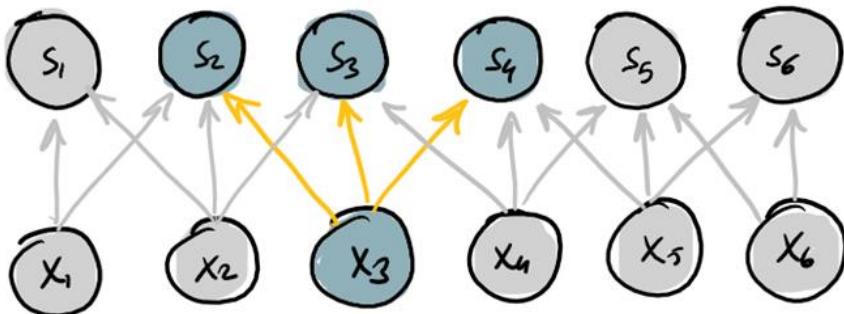


Data Driven Engineering II: Advaced Topics

Image processing and analysis

Institute of Thermal Turbomachinery
Prof. Dr.-Ing. Hans-Jörg Bauer



Summary of last week

- What is convolution
- Convolution operation in ANN
 - 1. What is a filter?
 - What do we learn in CNN?
 - Advantages compared to MLPs
- Padding, Stride, Pooling
- Rule of thumbs

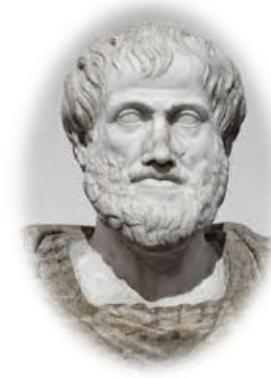
Case Study~

- ① Image classification
 - > augmentation

Outline of the week :

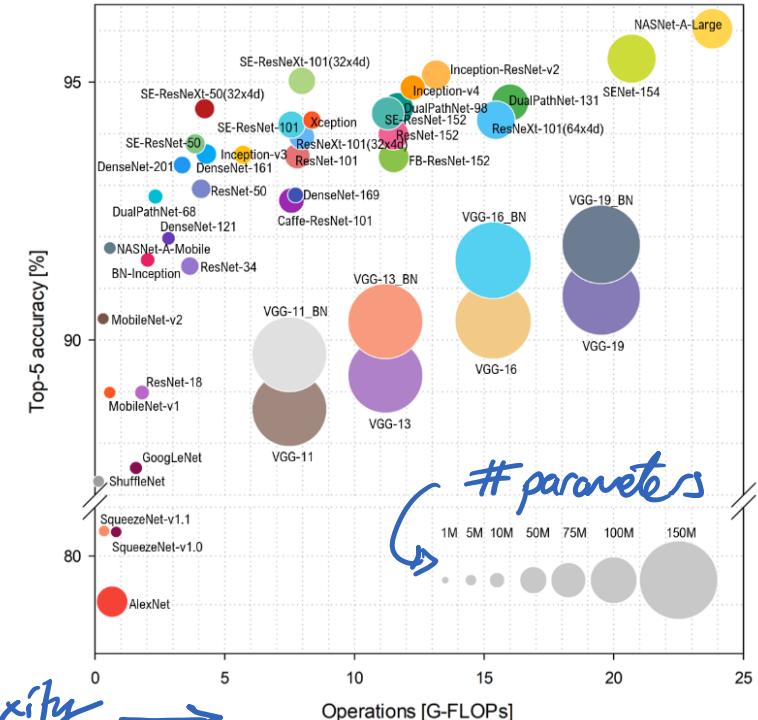
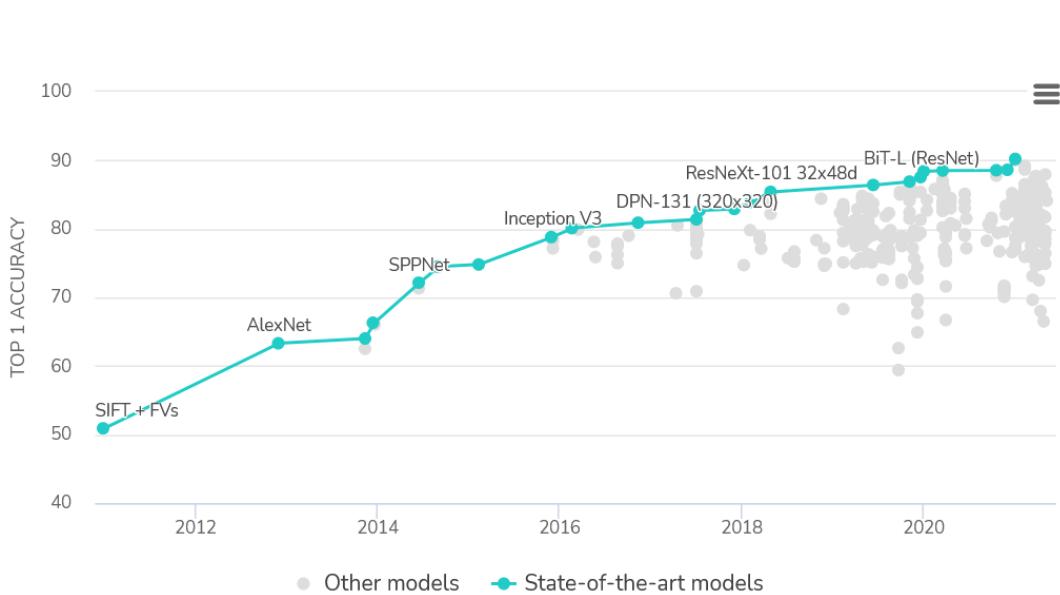
Conv. Neural Networks

- * What is CNN ?
- * Why convolution is useful ?
- * Where is it useful ?
- * CNN – How does it work ?
- * "Hall of fame" : Popular Arch- }
- * Transfer Learning with CNN



"The soul never thinks without a picture ."

where to begin ...



What is different?

- * Increase in depth
- * increase in width \Rightarrow more filters \Rightarrow filter hacks
- * Regularization & fine tuning
- * Training hacks
- * Data augmentation & tr. learning
- * Increased Image resolution

Who does not like classics?

LeNet-5 (1998)

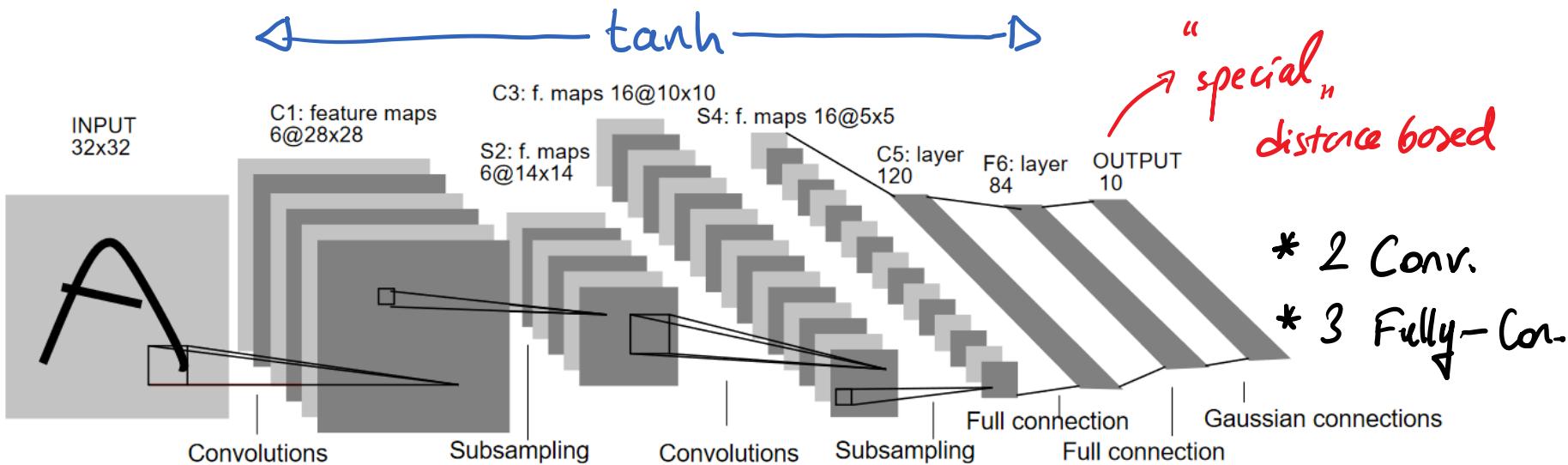
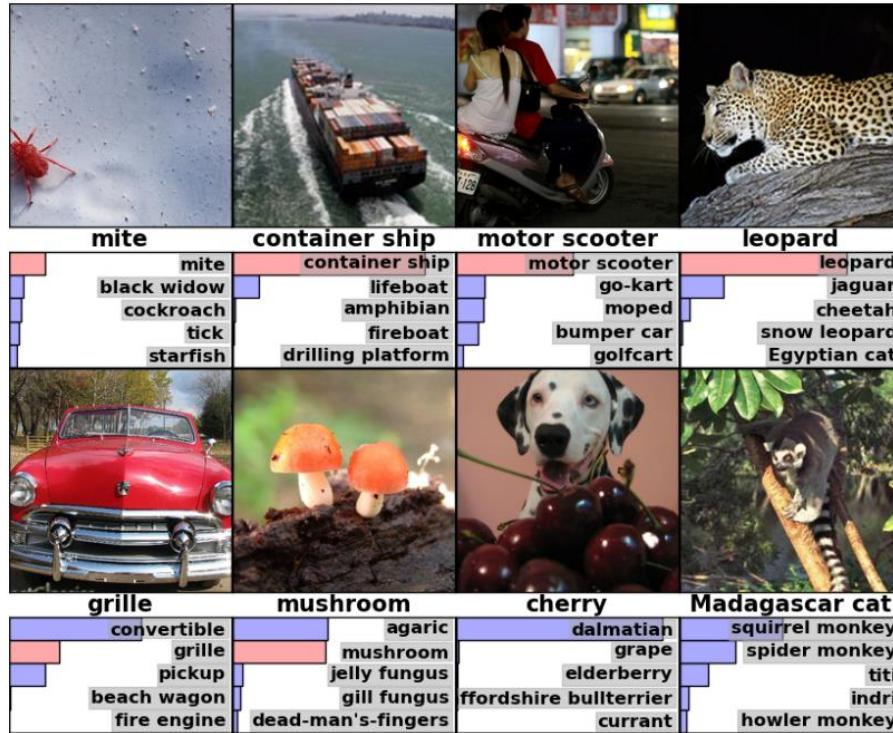


Image recognition – AlexNet (2012)



- Cited 128,948 times
- 1.2 Million images
- **5 layer network:**
 - 60 million parameters
 - **650k point neurons**

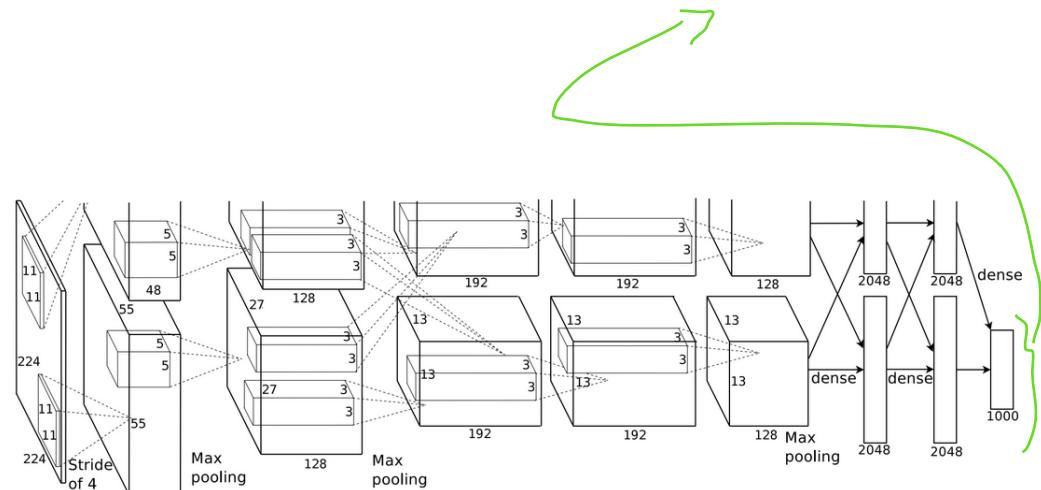
One small step for men, one giant leap for CNN

AlexNet (2012)

11x11 kernel makes so many parameters.
This brings the problem of the vanishing gradients.

Thousands classes, that's why it has thousands neurons.

- * Deeper := 5 Conv. + 3 Fully Con.
→ 60 M parameters
- * Stack conv. layers
- * Use ReLu + Dropouts
- * Data augmentation

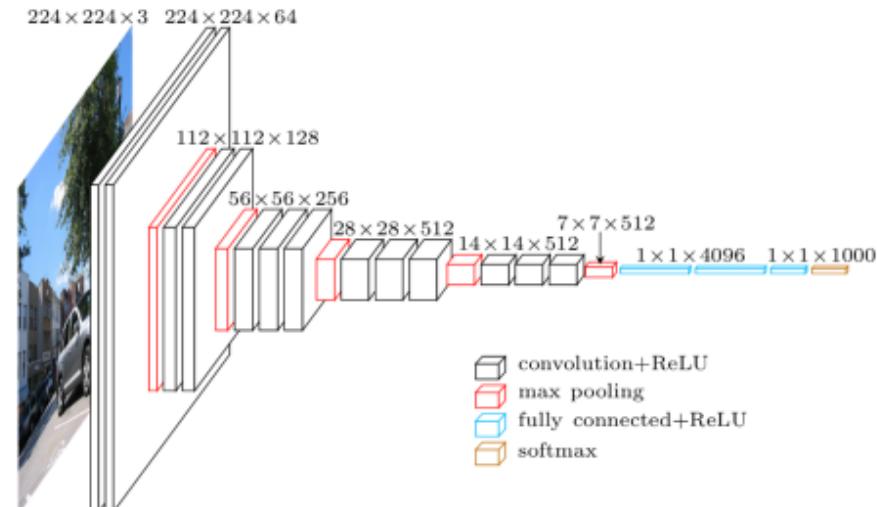


Deep Learning ...

VGG-16 (2014)

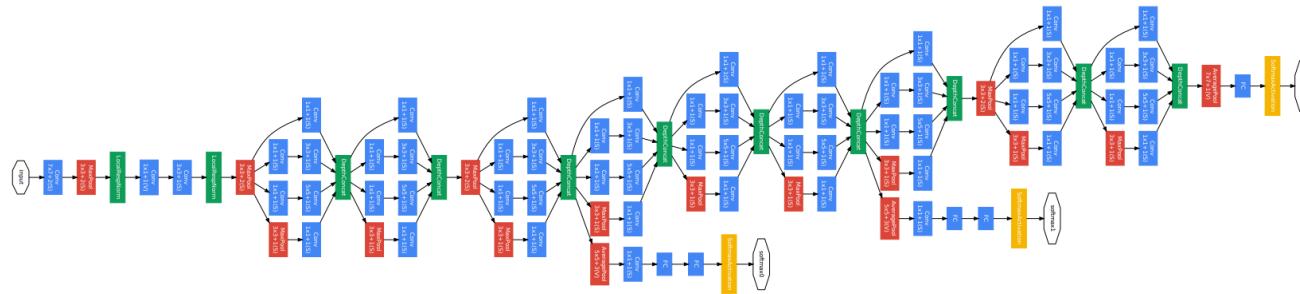
- * 13 Conv. + 3 Fully-connected.
- * 138 M parameters
- * ReLU + Smaller kernels ($2 \times 2, 3 \times 3$)

↳ Deeper variant "VGG-19"



Things get multi-layered :

inception-v1 // GoogleNet (2014)

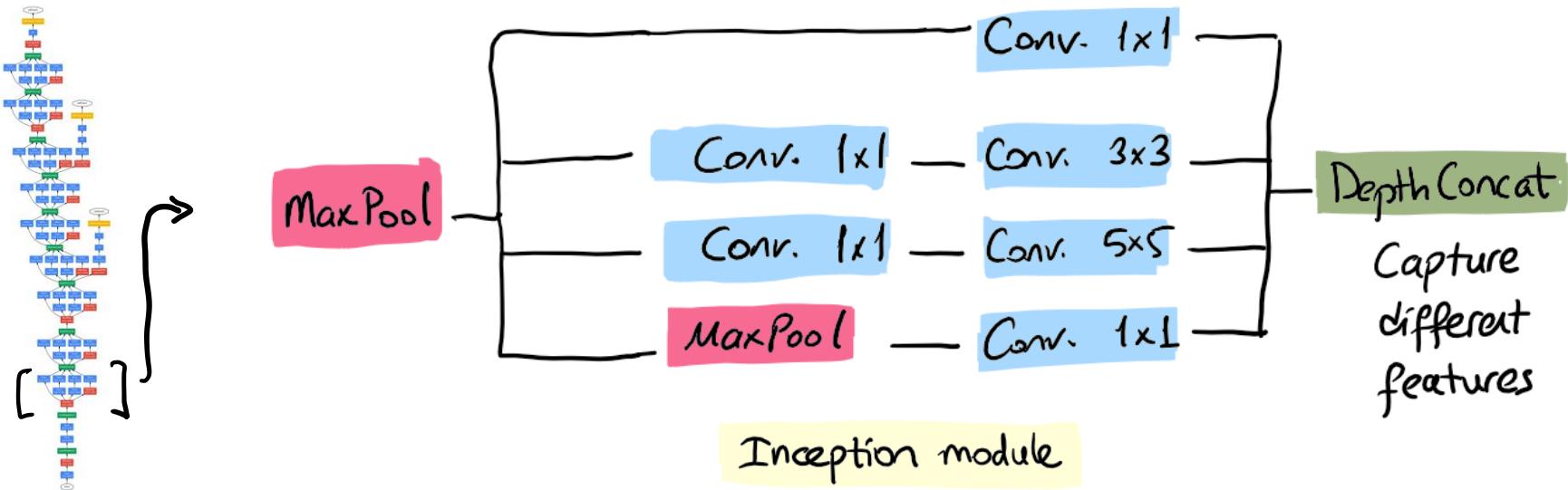


- * 22 Layers \Rightarrow 5M parameters
- * Inception module \Rightarrow "network in network",
enables much deeper network

Things get multi-layered :

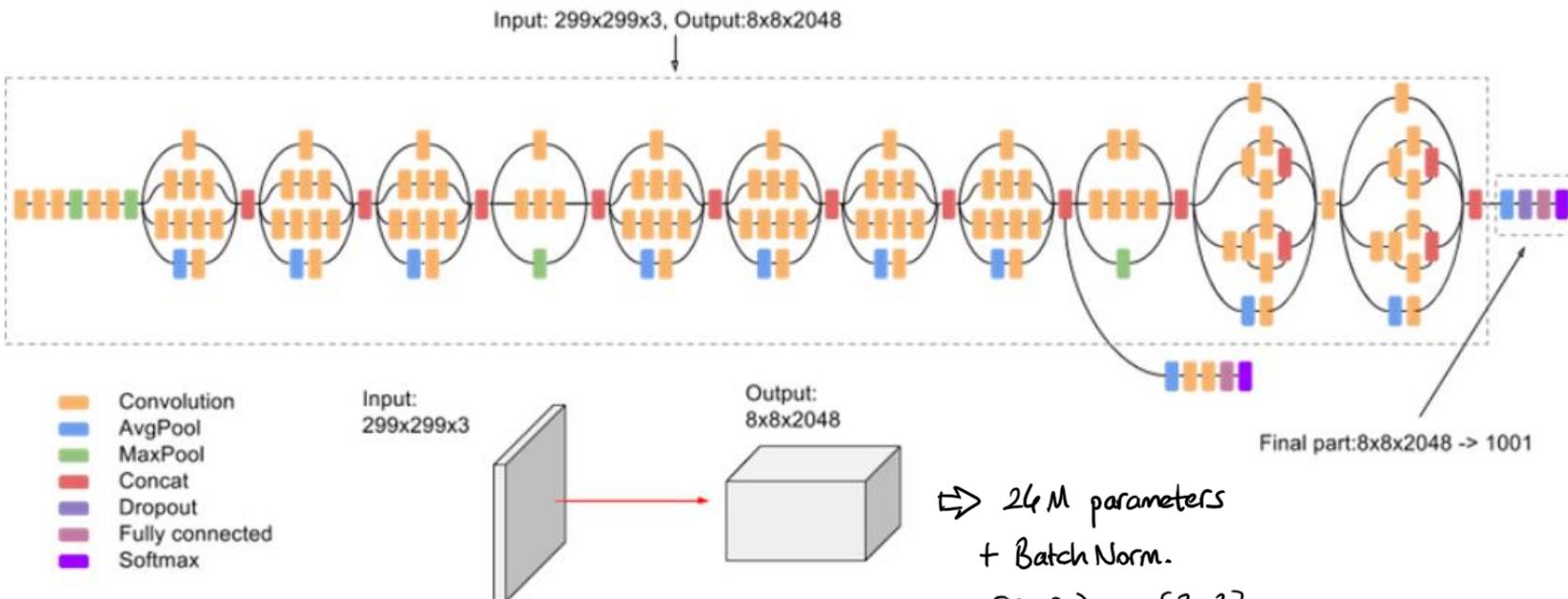
inception-v1 // GoogleNet (2014)

Now we have a multiple pipeline compared with the previous class.



Things get multi-layered : v3 (2015)

Instead of using the 3x3 filter size, they modified here.
Te models are getting bigger and bigger.



⇒ 26M parameters
+ Batch Norm.

- $[7 \times 7] \Rightarrow [3 \times 3]$
- $[5 \times 5] \Rightarrow [3 \times 3]$
- $[n \times n] \Rightarrow [1 \times n, n \times 1]$

(v4 → 2016)

A smart touch:

ResNet (2015)

* much deeper network \rightarrow 152 layers

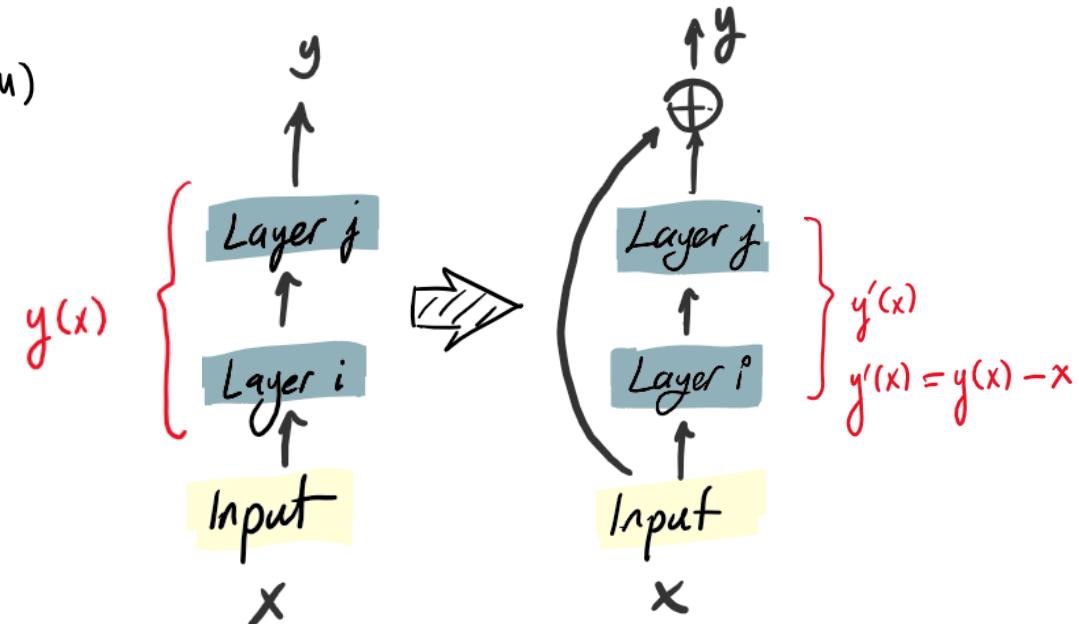
\hookrightarrow fewer parameters (26M)

! Exp./Var. Gradient problem

\hookleftarrow Skip Connections
+ Batch Norm

* Residual learning

* LSTM ?



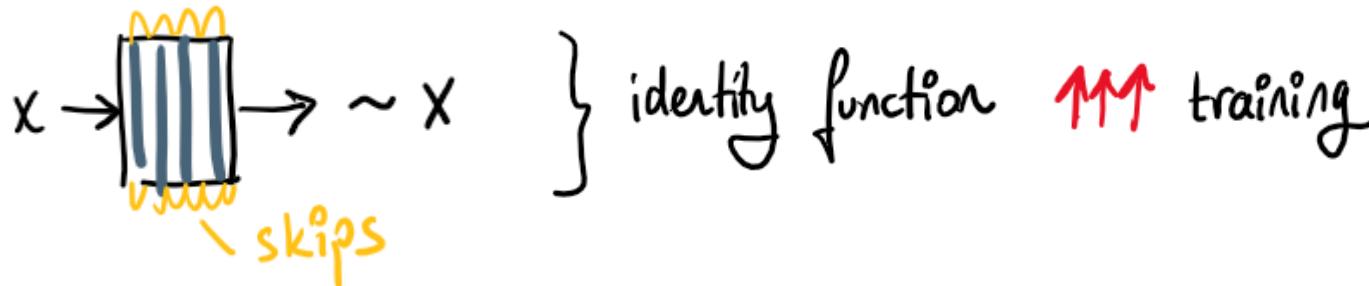
A smart touch:

ResNet (2015)

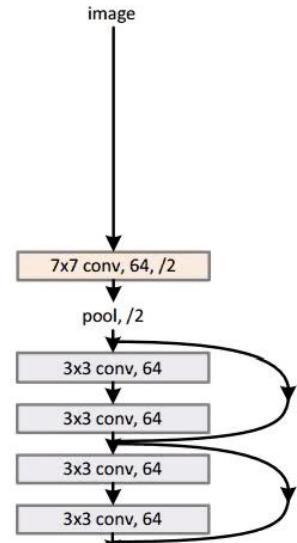
How does it help training?

① When network is initialized;

$$\bar{w} \rightarrow \phi$$



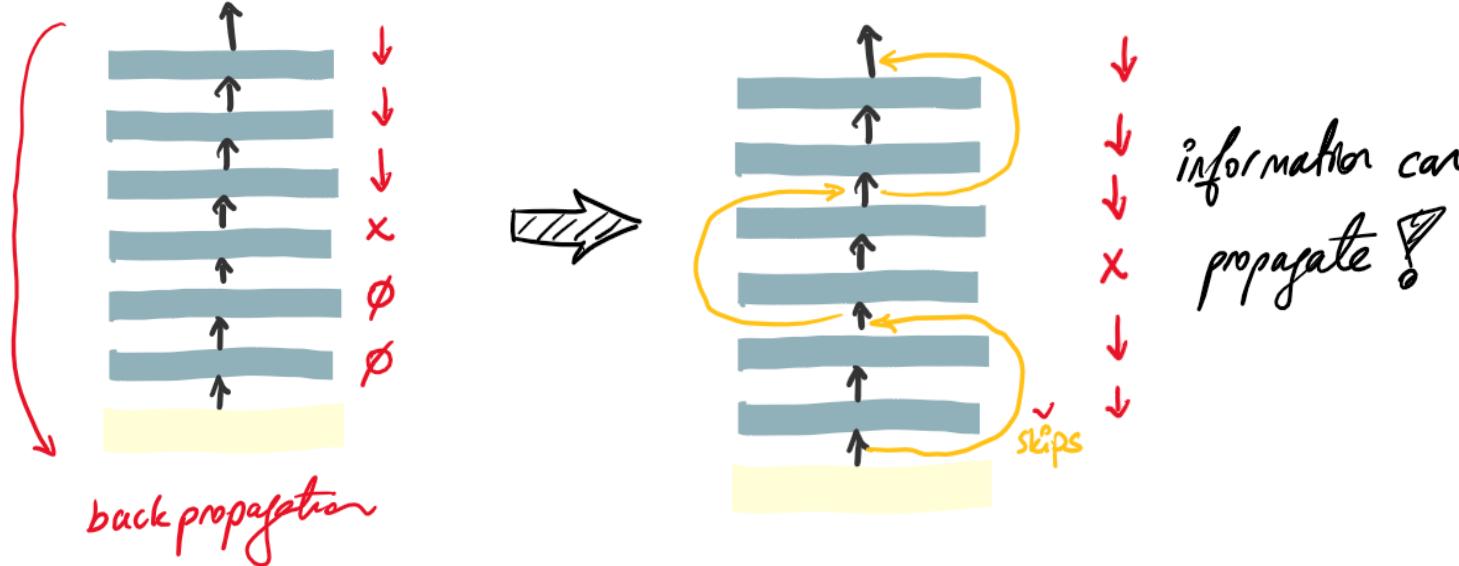
34-layer residual



A smart touch: ResNet (2015)

How does it help training?

② Learning bottlenecks



It looks nice . Lets use it everywhere ☺

DenseNet 2017

- * Use skip connections in dense blocks ↗ Re-use features with less parameters
- * use "concatenate", rather than adding ! Needs proper padding

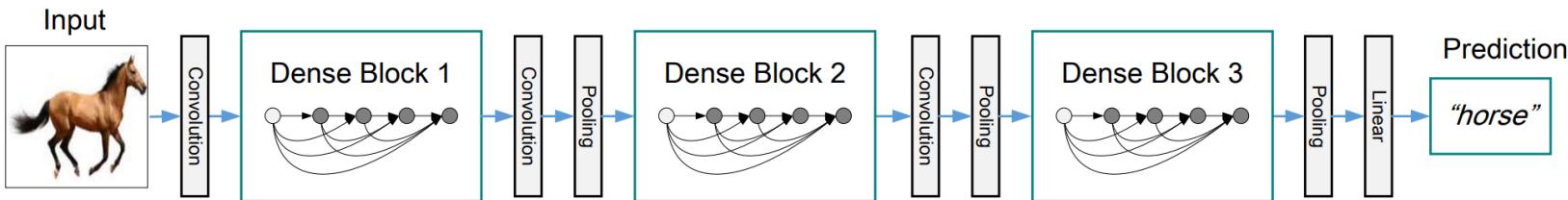
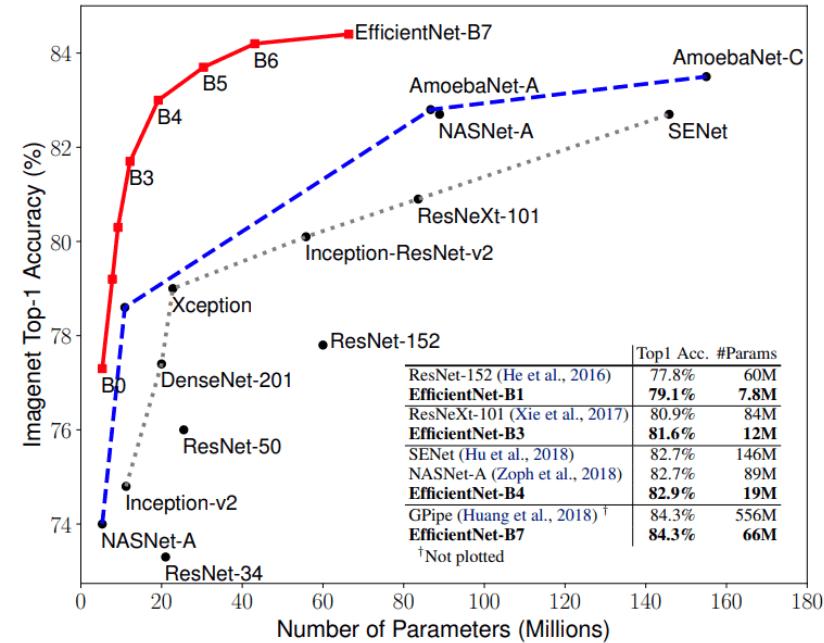
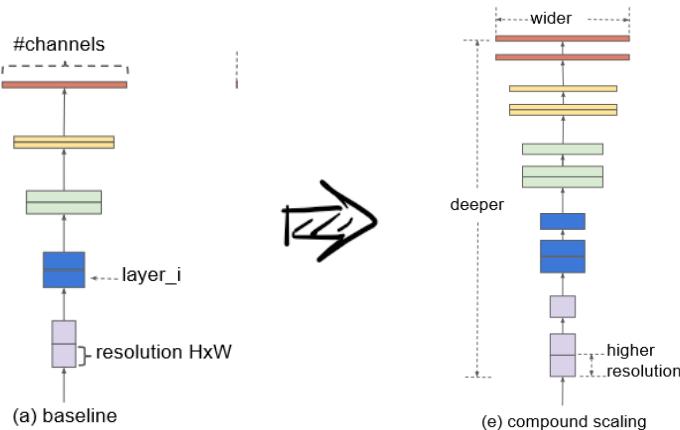


Figure 2: A deep DenseNet with three dense blocks. The layers between two adjacent blocks are referred to as transition layers and change feature-map sizes via convolution and pooling.

EfficientNet: "Rethinking model scaling"

* balancing network
↓
scaling method { depth
width
resolution



EfficientNet: "Rethinking model scaling"

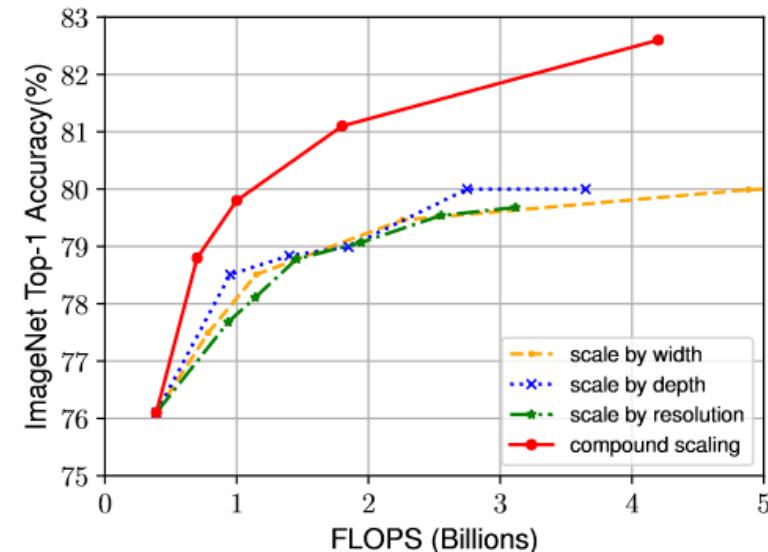
* balancing network {
 ↴
 scaling method depth
 width
 resolution

- depth = $d = \alpha^{\phi}$; $\alpha \beta^2 \gamma^2 \approx 2$ (FLOPs)
- width = $w = \beta^{\phi}$; $\alpha, \beta, \gamma > 1$
- res. = $r = \gamma^{\phi}$

(i) $\phi = 1$; do grid search on α, β, γ
 e.g. (1.2, 1.2, 1.15)

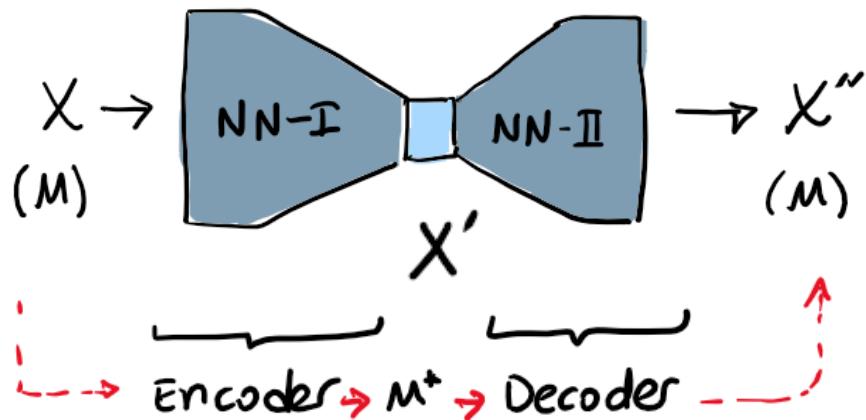
(ii) Fix α, β, γ ; scale ϕ with respect to hardware

⇒ Start small; gradually scale it



A useful Tool: Conv. Auto-encoders

* DDE-1 \Rightarrow "Autoencoders", \Rightarrow "PIV & GAN-Representation Projects,"

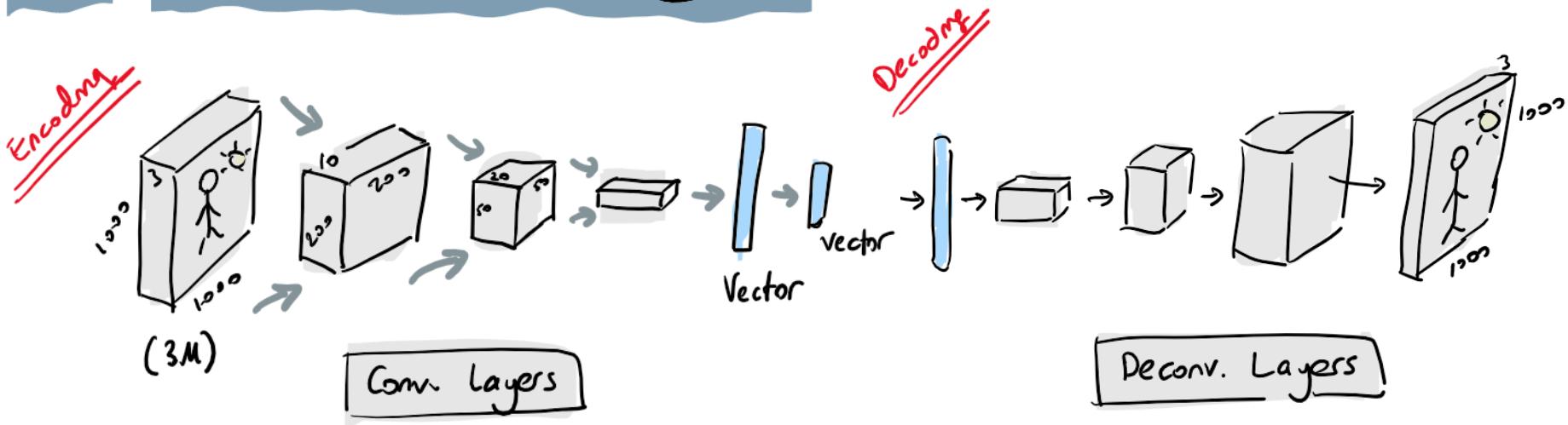


$(X - X'') := \frac{\text{Reconst.}}{\text{Error}}$

$M^* \ll M$

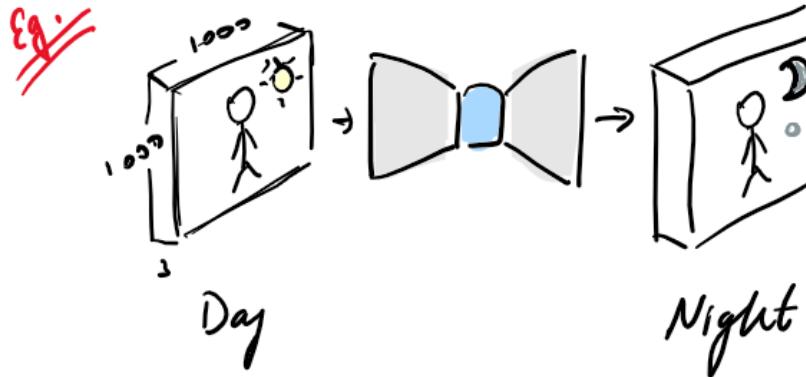
 Representation Learning

A useful Tool: Conv. Auto-encoders



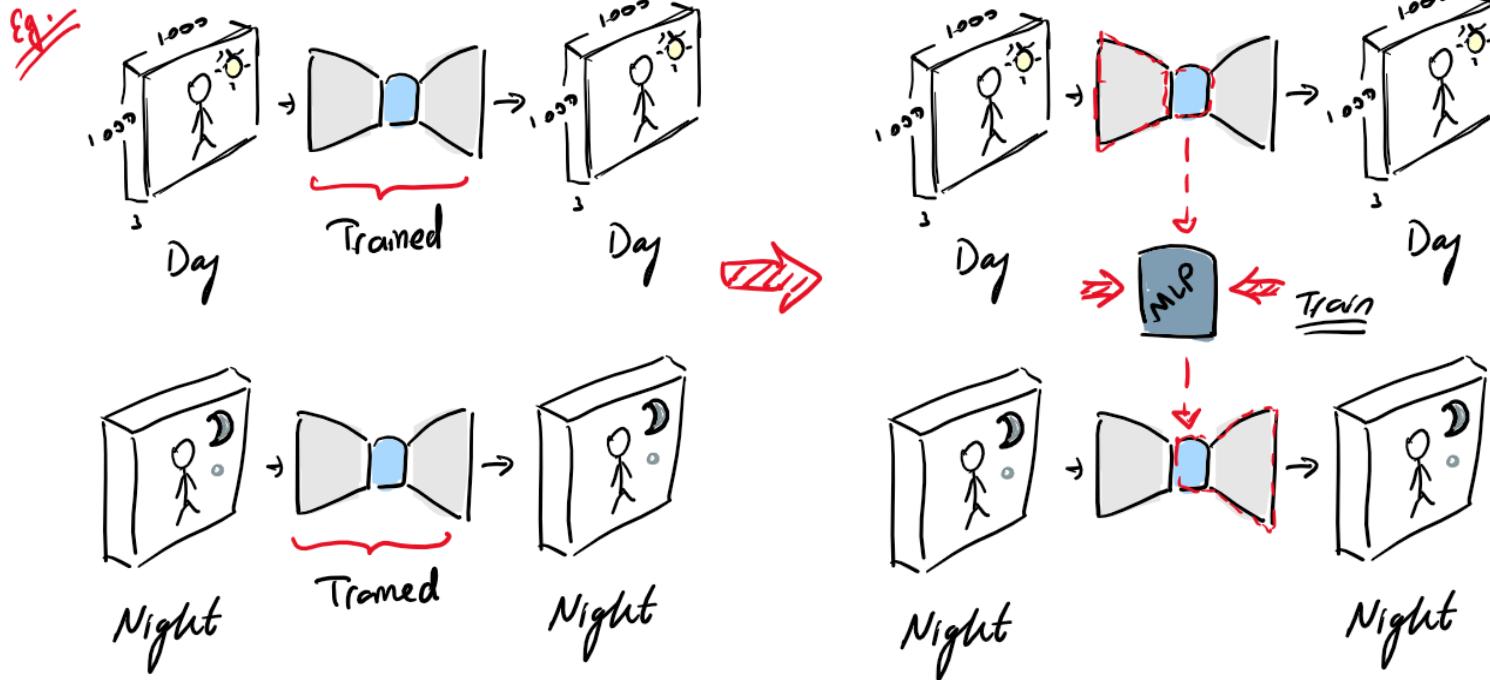
D How this is useful?

A useful Tool: Conv. Auto-encoders



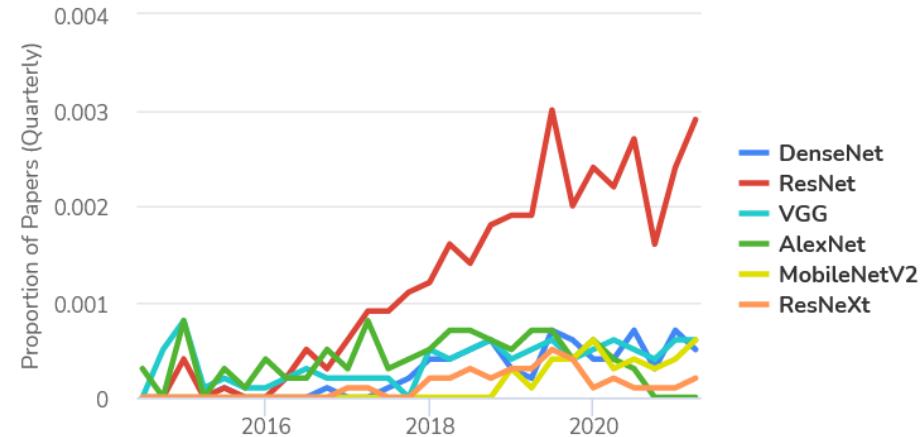
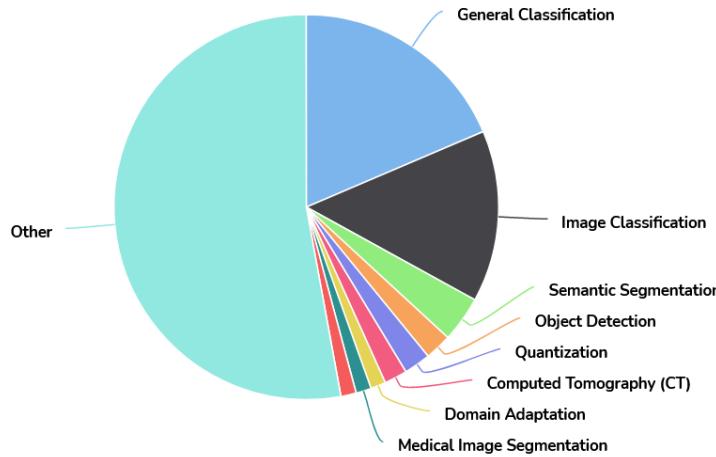
- (*) Image translation; (*) CNN pretraining
 - Day → Night
 - B&W → Color
 - Noisy → Clean
 - ...

A Useful Tool: Conv. Auto-encoders

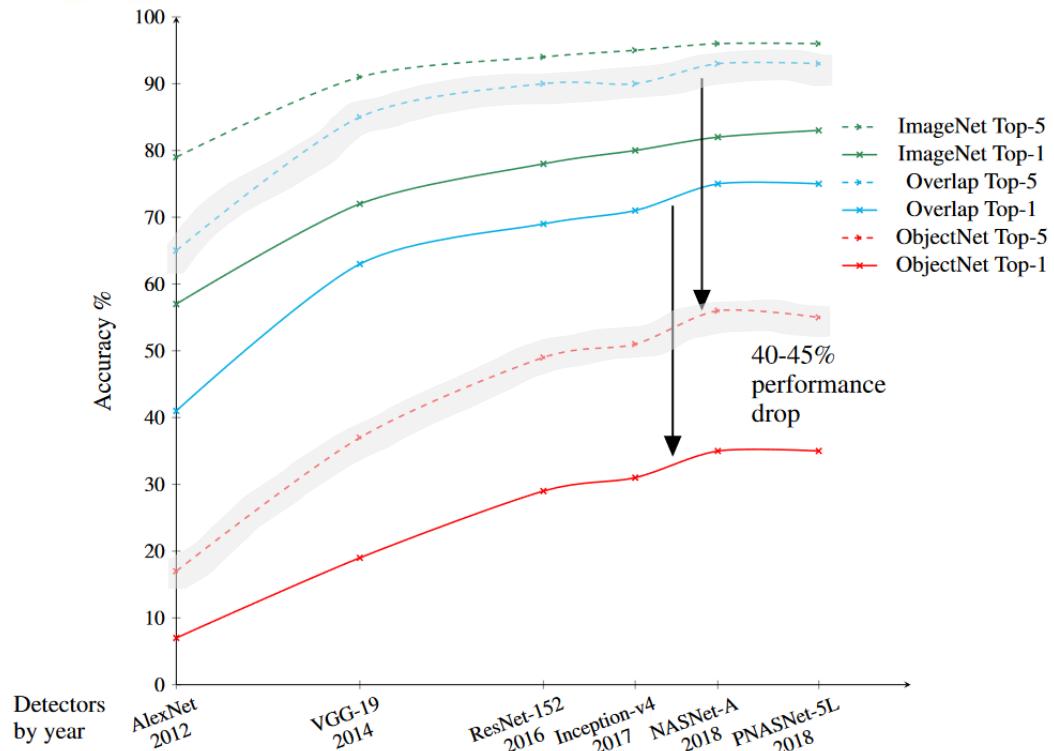


Too many options ...

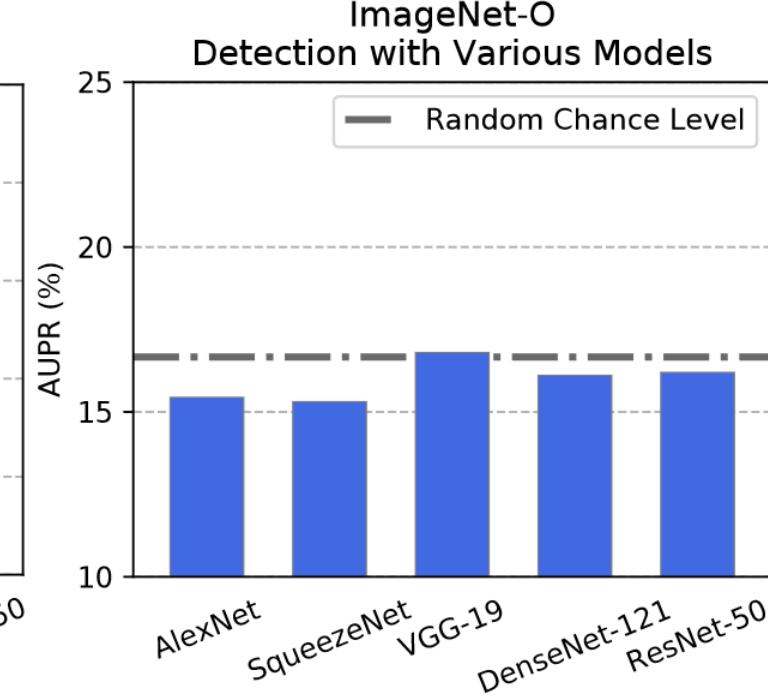
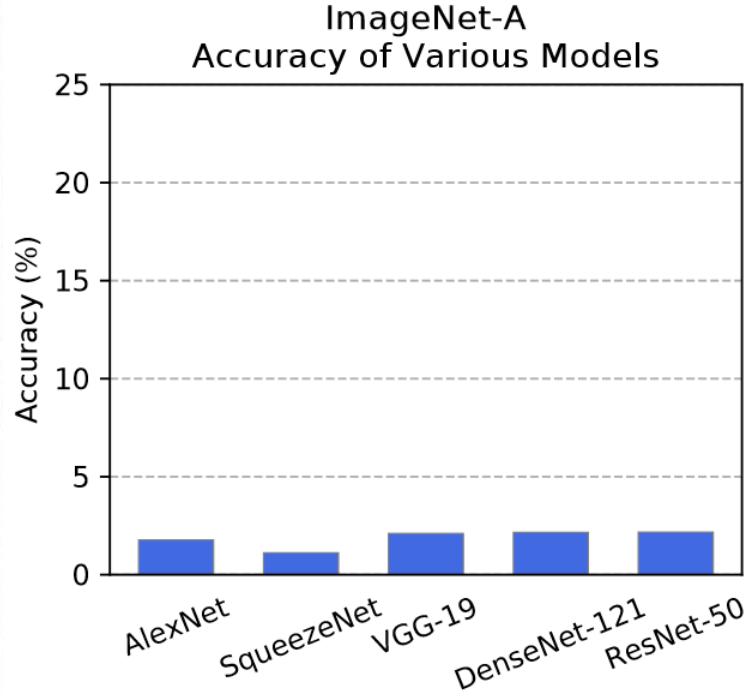
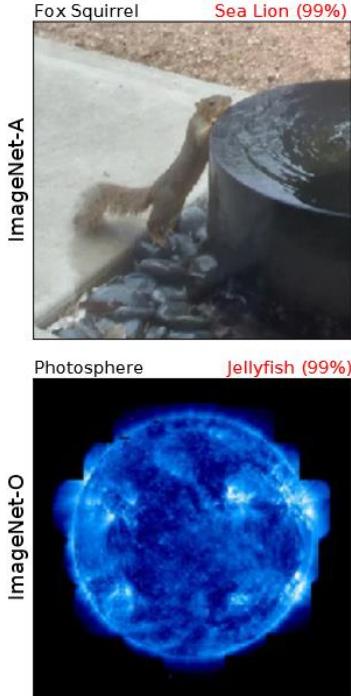
- * Models including CNN ~ 100 accessible
- * Utilized in various fields



Shallow & Brittle: no hierarchical perception

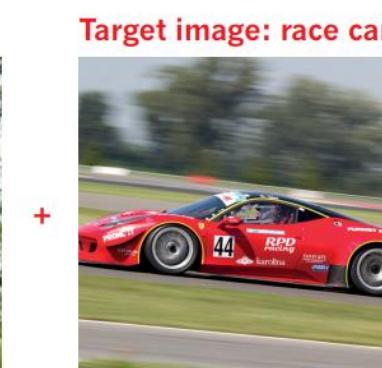
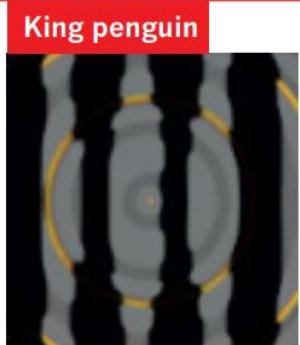
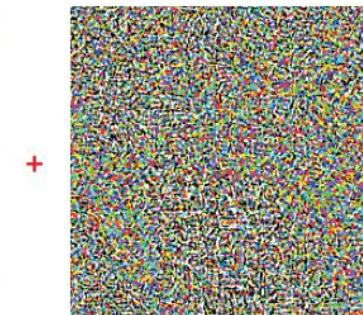
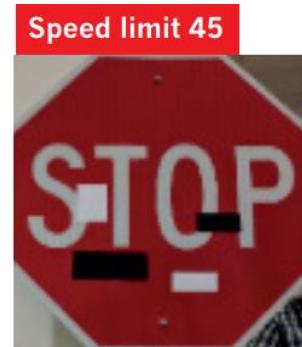


Shallow & Brittle: no hierarchical perception



<https://arxiv.org/pdf/1907.07174.pdf>

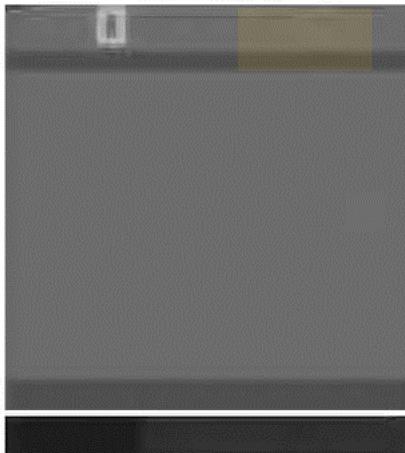
Lack of hierarchical perception: easy to trick



Lack of hierarchical perception: easy to trick

Test-Time Execution

raw input



output action distribution

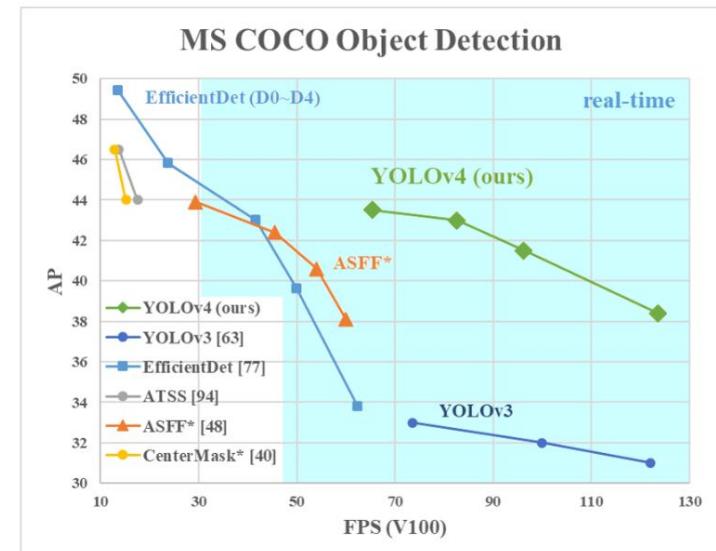
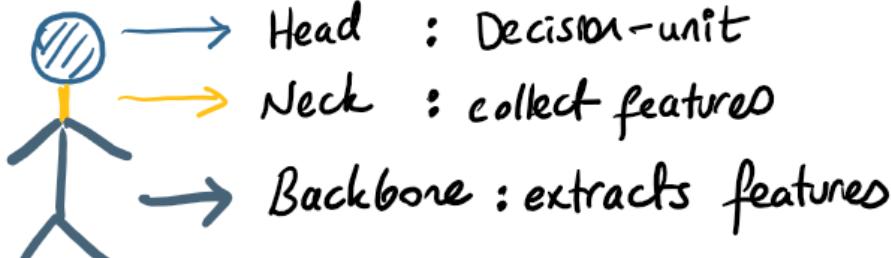
<https://arxiv.org/abs/1702.02284>

Case study: Object Detection:

* {Classify + Location + many objects}

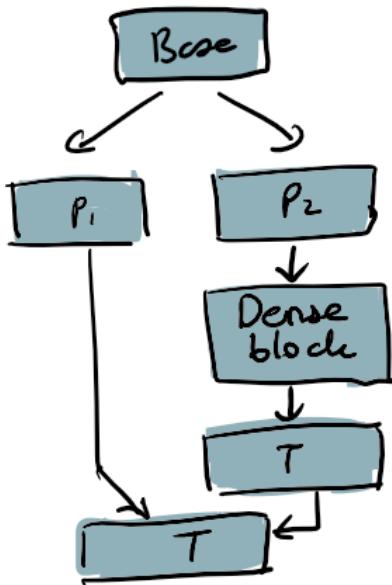
Model anatomy reflects this aim

* YOLO := "You only look once"



Case study: Object Detection:

Backbone: Dense block + Cross-Stage-Partial connections



+ Darknet 53

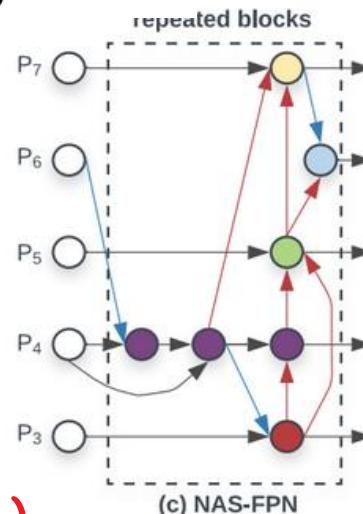
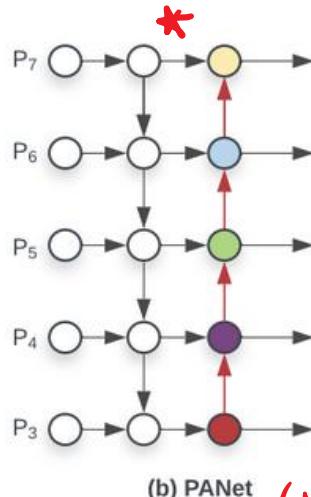
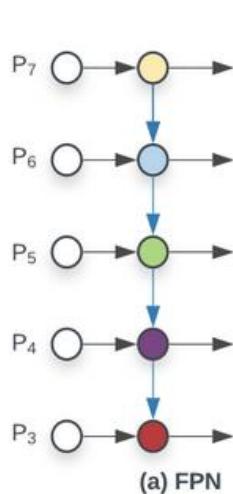
higher acc. > Resnet

Type	Filters	Size	Output
Convolutional	32	3 x 3	256 x 256
Convolutional	64	3 x 3 / 2	128 x 128
Convolutional	32	1 x 1	
1x Convolutional	64	3 x 3	
Residual			128 x 128
Convolutional	128	3 x 3 / 2	64 x 64
Convolutional	64	1 x 1	
2x Convolutional	128	3 x 3	
Residual			64 x 64
Convolutional	256	3 x 3 / 2	32 x 32
Convolutional	128	1 x 1	
8x Convolutional	256	3 x 3	
Residual			32 x 32
Convolutional	512	3 x 3 / 2	16 x 16
Convolutional	256	1 x 1	
Convolutional	512	3 x 3	
Residual			16 x 16
Convolutional	1024	3 x 3 / 2	8 x 8
Convolutional	512	1 x 1	
Convolutional	1024	3 x 3	
Residual			8 x 8
Avgpool		Global	
Connected		1000	
Softmax			

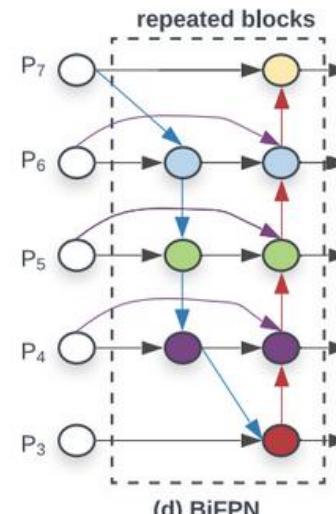
Case study: Object Detection:

Neck

- * Combination of features managed here



$P \rightarrow$ feature layer



Eff. Net

Case study: Object Detection:

Head

- * YOLOv4 ← YOLOv3
- * anchor-based detection
- * 3 levels of granularity

Bag of Freebies

- * Drop Block regularization
- * Class label smoothing
- * CutMix & Mosaic data aug.

Bag of Specials

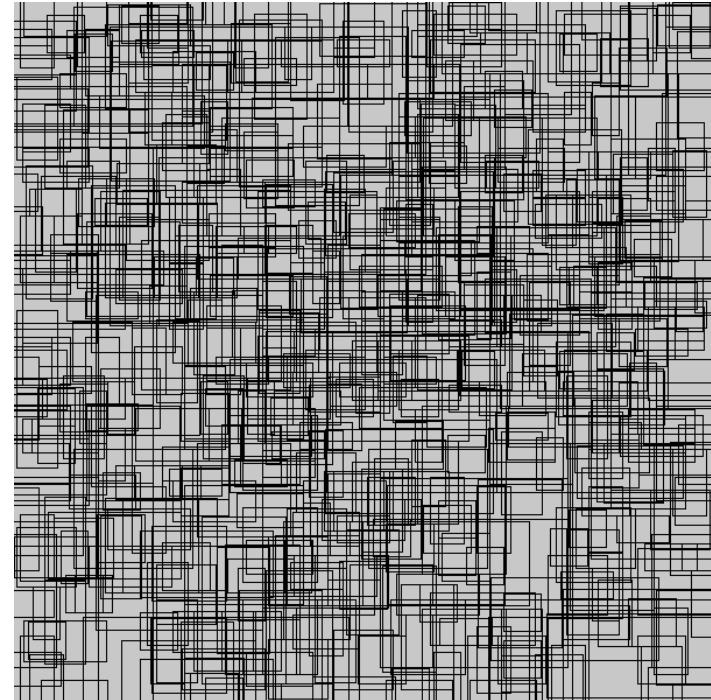
- * Mish activation
- * Cross Stage Partial connections
- * Multi-input-weighted residual connections

Case study: Object Detection:

Anchor Boxes

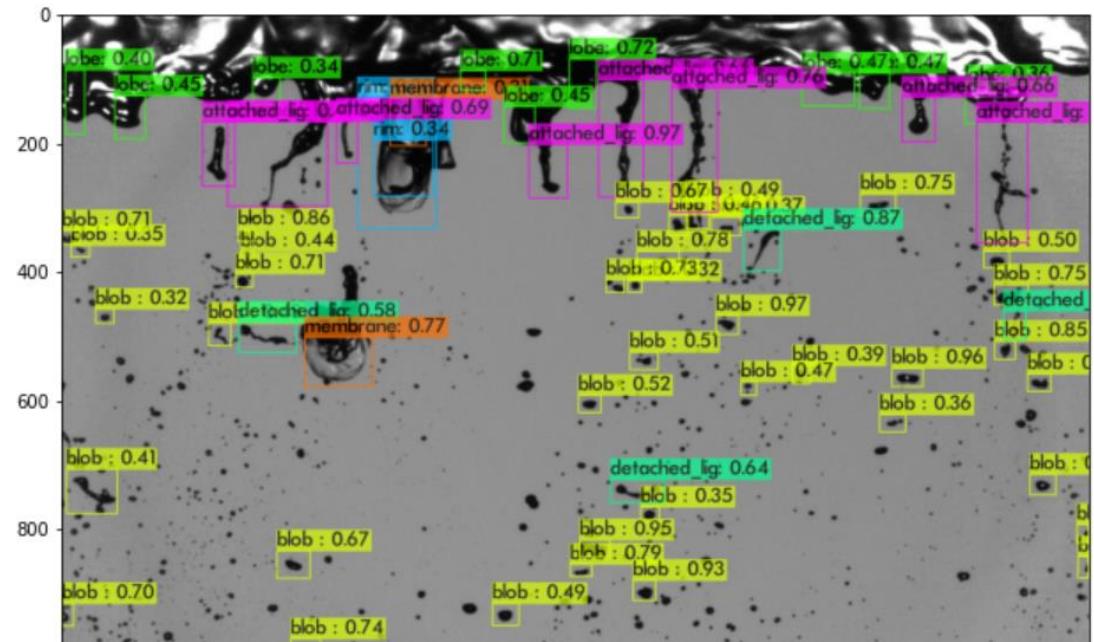
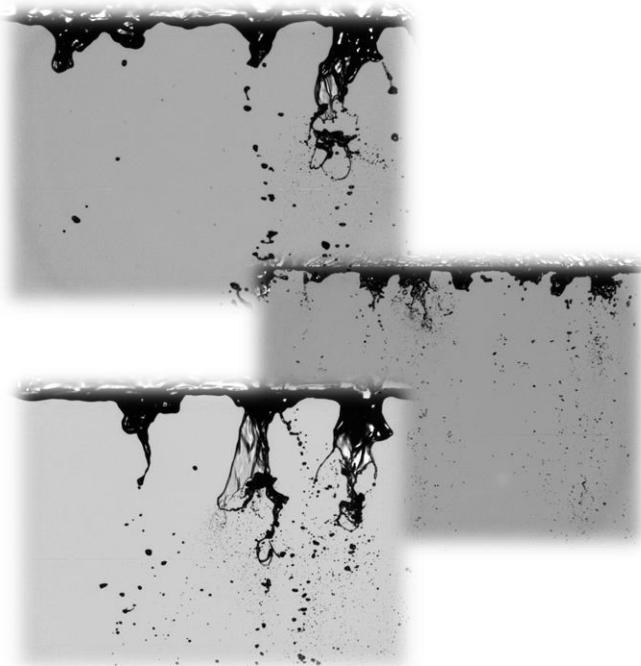
- (i) Create many boxes for each predictor
- (ii) For each box, calculate which objects bounding box has the highest overlap/non-overlap ratio. (IoU)
- (iii) If highest (IoU) > 50% \Rightarrow "Detect object,"
- (iv) ~40% - 50% \Rightarrow ambiguous \Rightarrow "not learn from this case,"
- (v) < 40% \Rightarrow No object here 

? Box dimensions \Rightarrow YOLO; K-means clustering on training data



Case study: Object Detection:

Get the pictures of each part of the video as frame per seconds, and use it for training our model.





colab

Computer Vision

544 methods • 43959 papers with code

Image Feature Extractors



[▶ See all 39 methods](#)

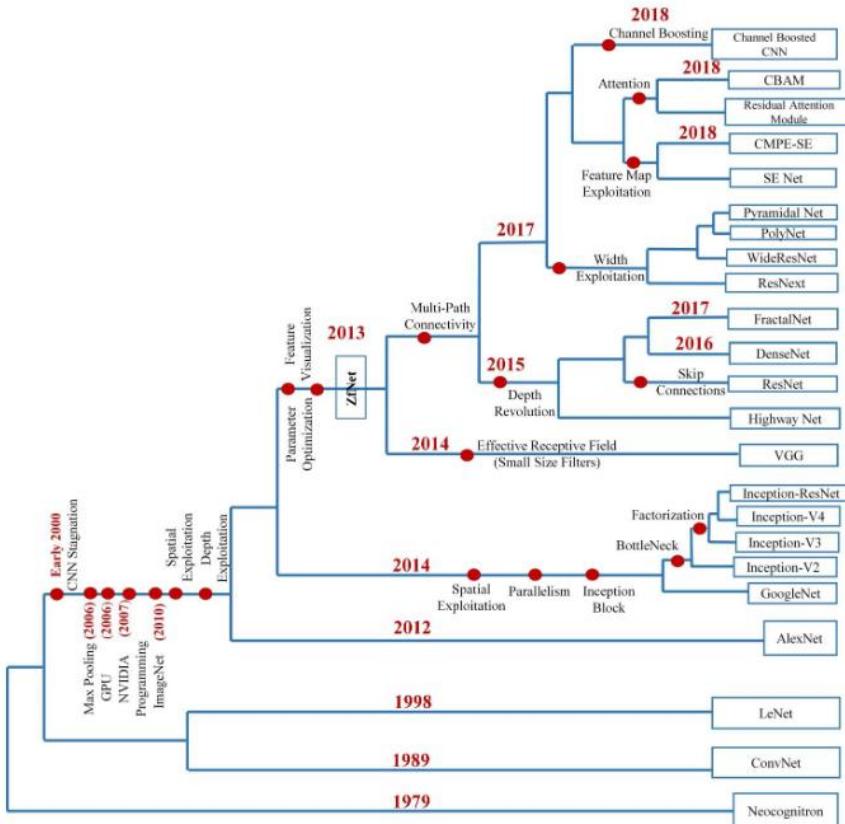
Convolutions





Convolutional Neural Networks

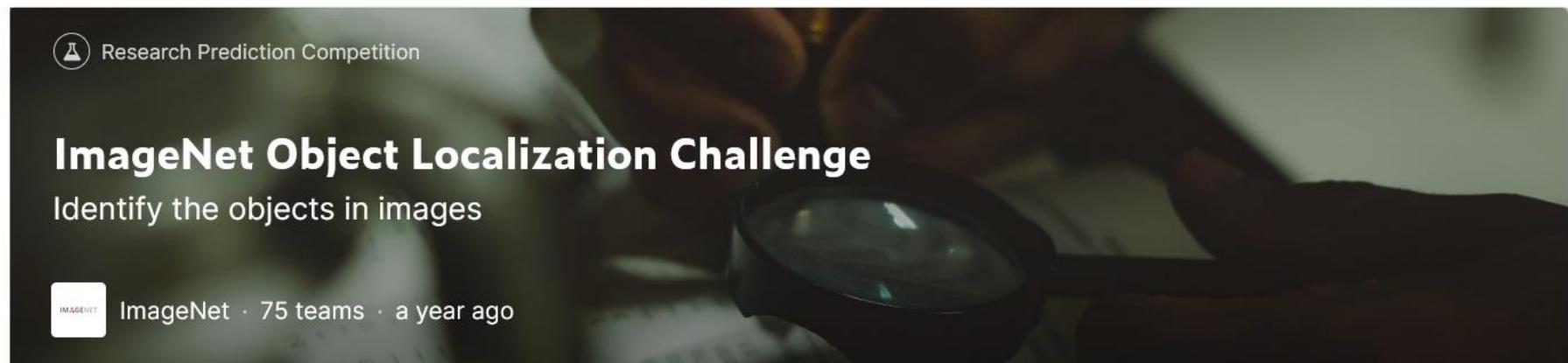
METHOD	YEAR	PAPERS
 ResNet	2015	1139
 VGG	2014	293
 AlexNet	2012	278
 DenseNet	2016	248
 MobileNetV2	2018	144
 ResNeXt	2016	107
 GoogLeNet	2014	106
 EfficientNet	2019	101



Artificial Intelligence Review

A Survey of the Recent Architectures of Deep Convolutional Neural Networks

ImageNet Challenge :



Research Prediction Competition

ImageNet Object Localization Challenge

Identify the objects in images

IMAGENET ImageNet · 75 teams · a year ago

Overview Data Code Discussion [Leaderboard](#) Rules [Join Competition](#)