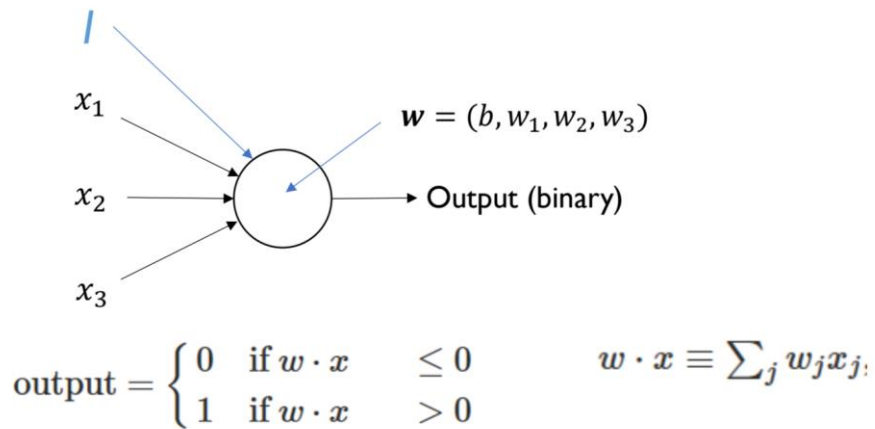
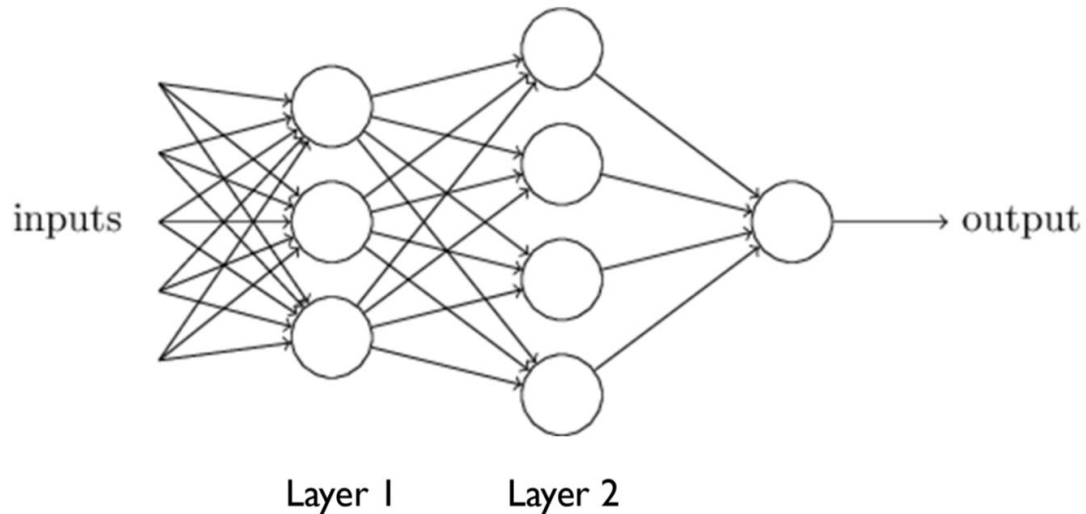


Deep Learning for Computer Vision

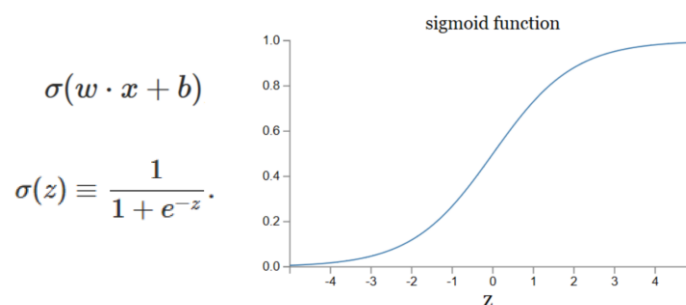
Perceptron - single neuron



Multi-Layer Perceptron

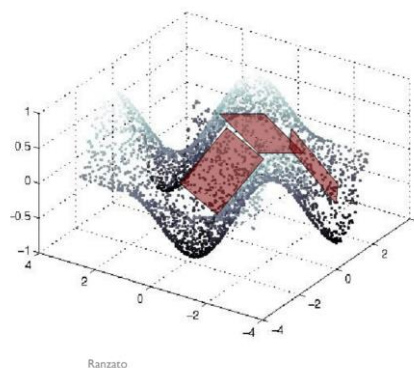
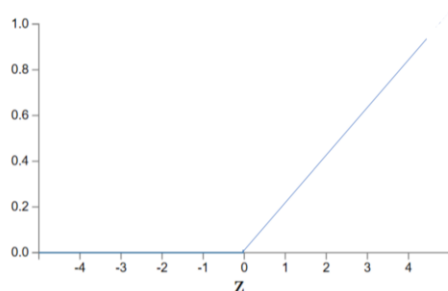


Nonlinearity – activation function

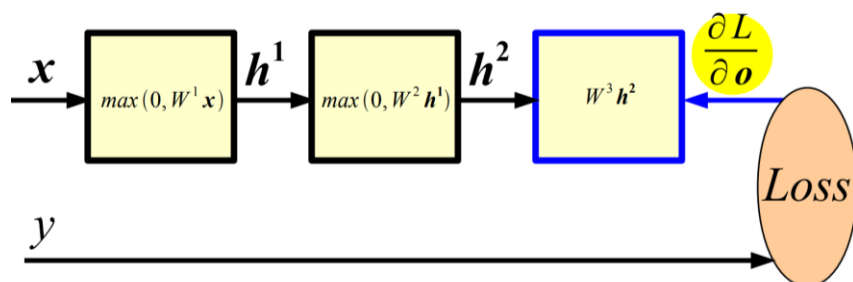


ReLU

$$f(x) = \max(0, x)$$



Learning – Backpropagation Algorithm



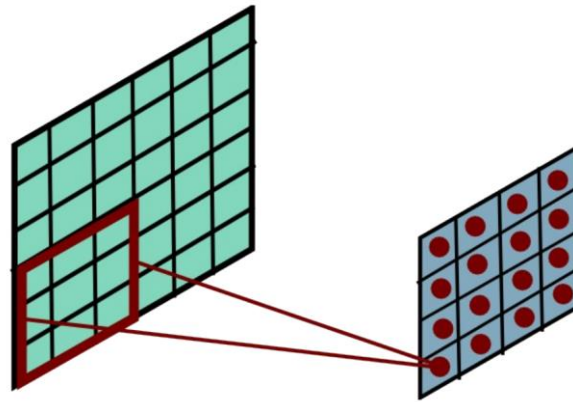
Given $\partial L / \partial \mathbf{o}$ and assuming we can easily compute the Jacobian of each module, we have:

$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial W^3}$$

$$\frac{\partial L}{\partial \mathbf{h}^2} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial \mathbf{h}^2}$$

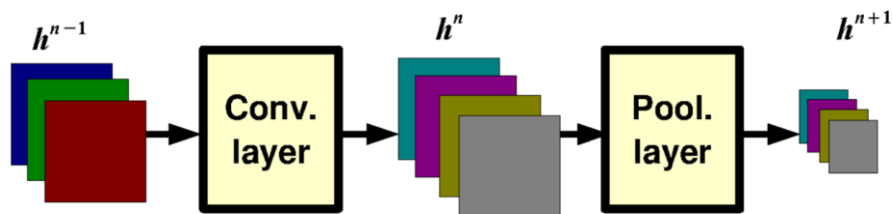
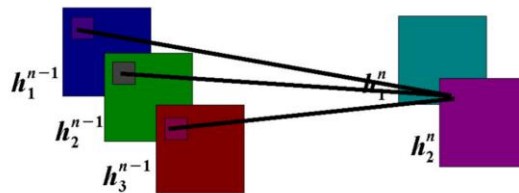
$$\frac{\partial L}{\partial W^3} = (p(c|\mathbf{x}) - \mathbf{y}) \mathbf{h}^{2T} \quad \frac{\partial L}{\partial \mathbf{h}^2} = W^{3T} (p(c|\mathbf{x}) - \mathbf{y})$$

Convolutional Neural Nets (CNNs)

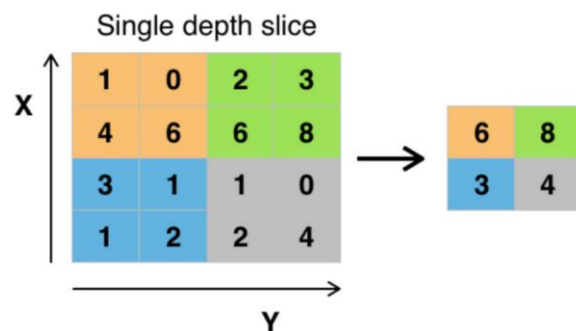


$$h_j^n = \max(0, \sum_{k=1}^K h_k^{n-1} * w_{kj}^n)$$

output feature map input feature map kernel

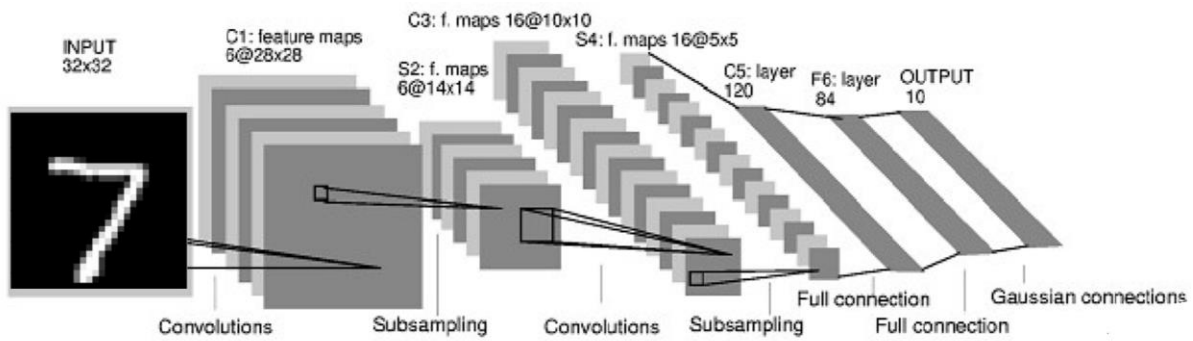


Pooling Layers - downsampling

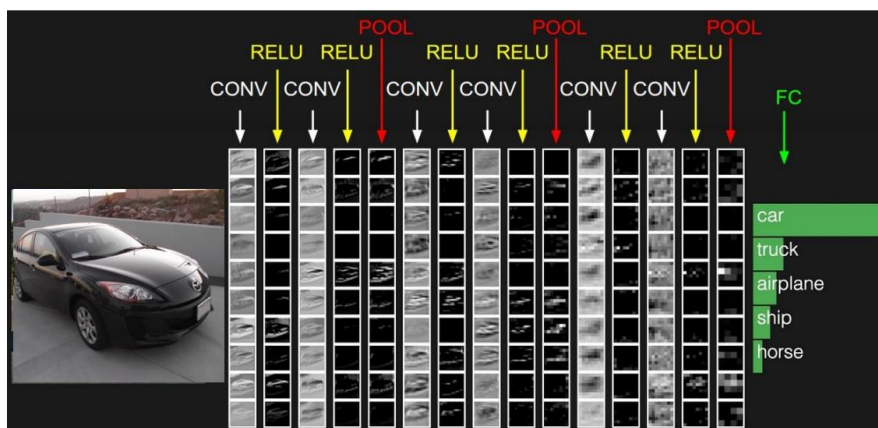


(max pooling)

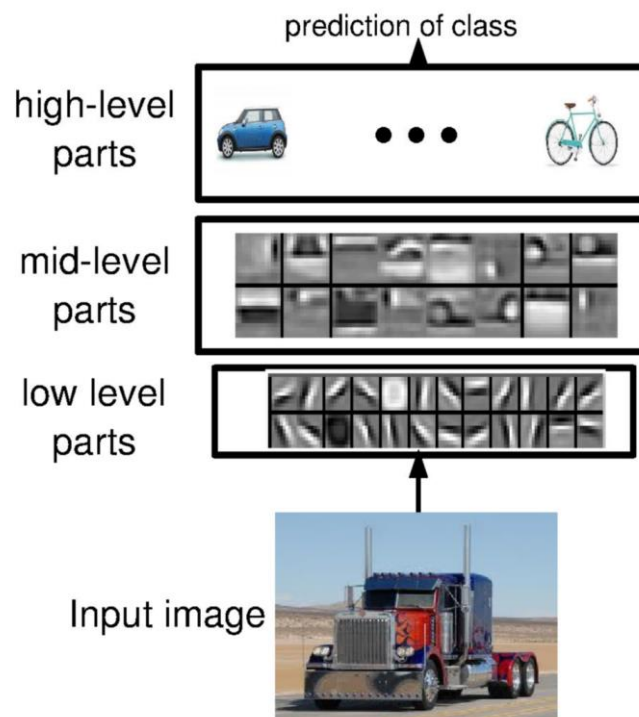
Yann LeCun's MNIST Architecture



Layering

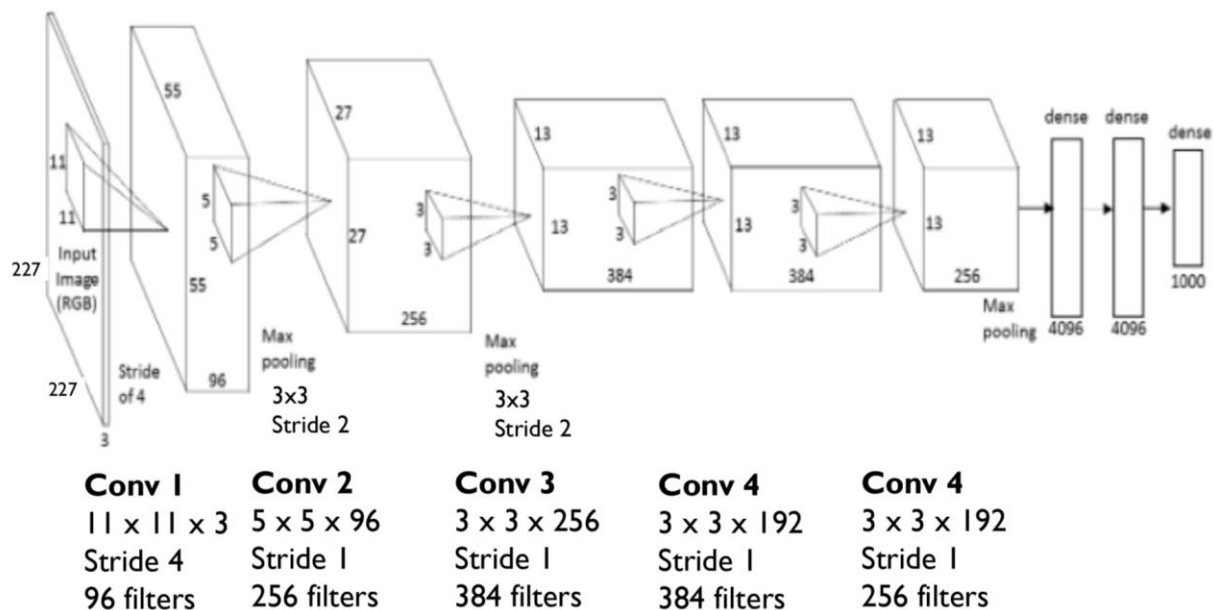


Hierarchical Feature Learning



AlexNet (Krizhevsky et al. 2012)

params	AlexNet	FLOPs
4M	FC 1000	4M
16M	FC 4096 / ReLU	16M
37M	FC 4096 / ReLU	37M
	Max Pool 3x3s2	
442K	Conv 3x3s1, 256 / ReLU	74M
1.3M	Conv 3x3s1, 384 / ReLU	112M
884K	Conv 3x3s1, 384 / ReLU	149M
	Max Pool 3x3s2	
	Local Response Norm	
307K	Conv 5x5s1, 256 / ReLU	223M
	Max Pool 3x3s2	
	Local Response Norm	
35K	Conv 11x11s4, 96 / ReLU	105M

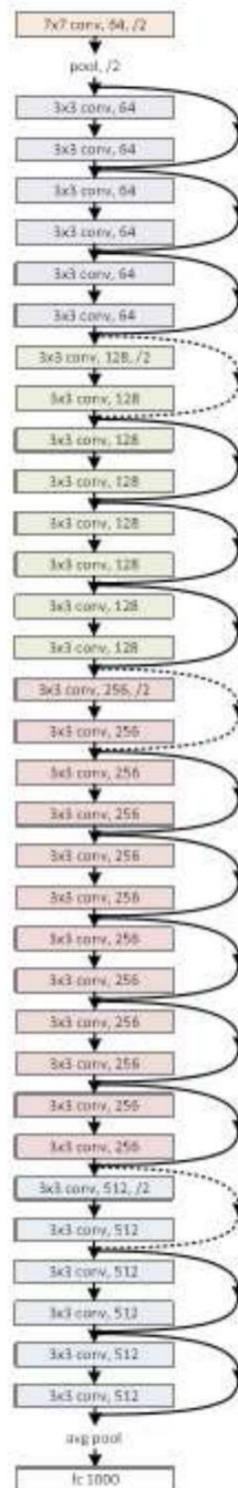


ResNet (He et al., 2015)

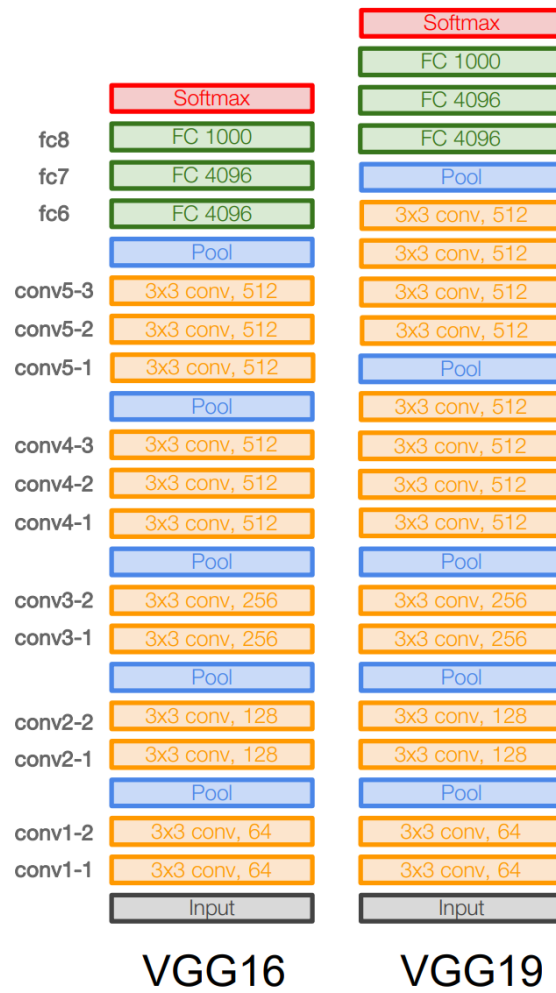
plain net



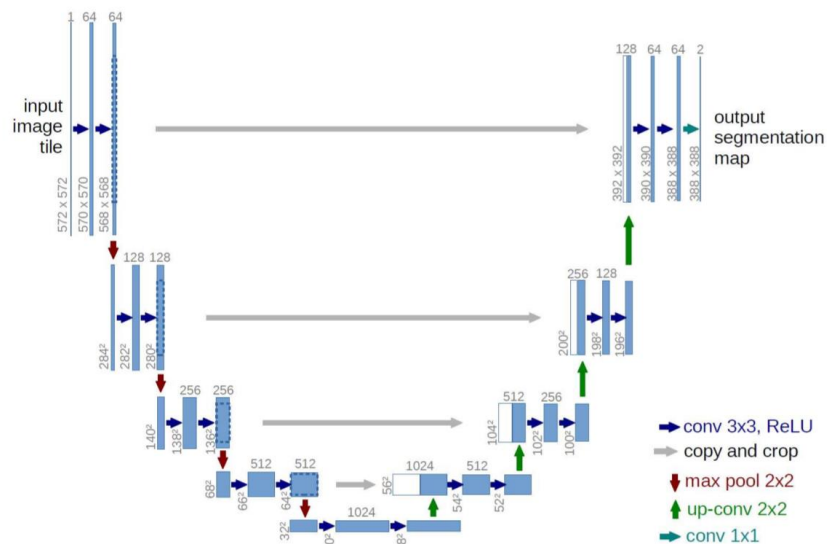
ResNet



VGG Net (Simonyan and Zisserman, 2014)

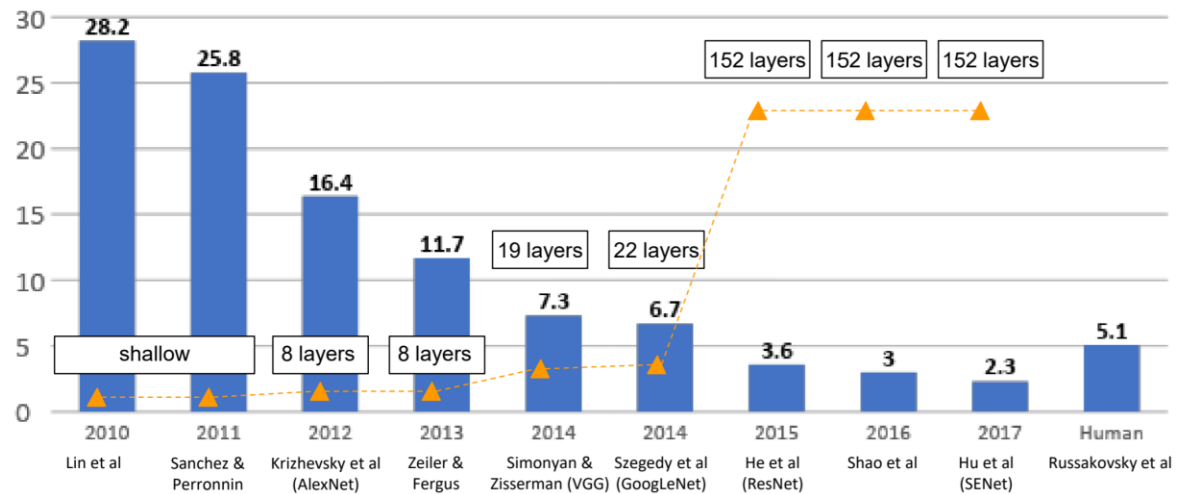


UNets

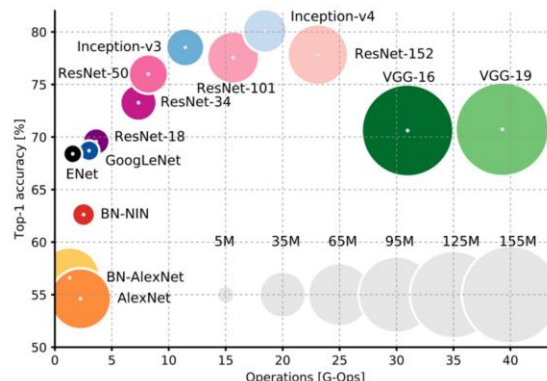
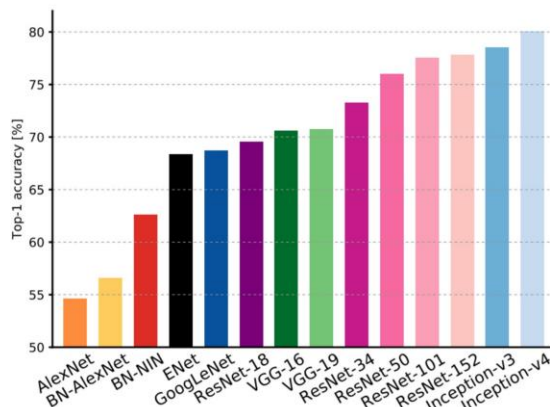


Performance and Model Complexity

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Comparing complexity...



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017.

ImageNet Dataset

ImageNet Large Scale Visual Recognition Challenge (ILSVRC), a benchmark in image classification and object detection

1-K has 1000 object classes, 1,281,167 training, 50,000 validation, 100,000 test images.

22-K has 21841 classes, 14,197,122 images

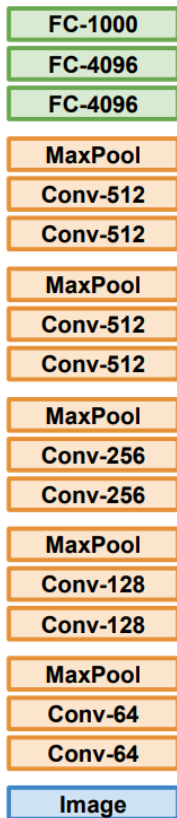


CIFAR-10

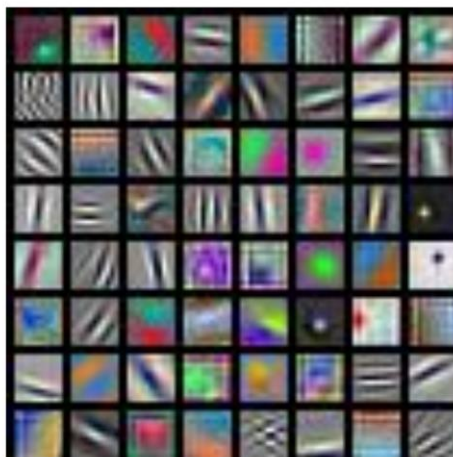
60,000 32x32 color images, 10 classes

Transfer Learning

1. Train on Imagenet



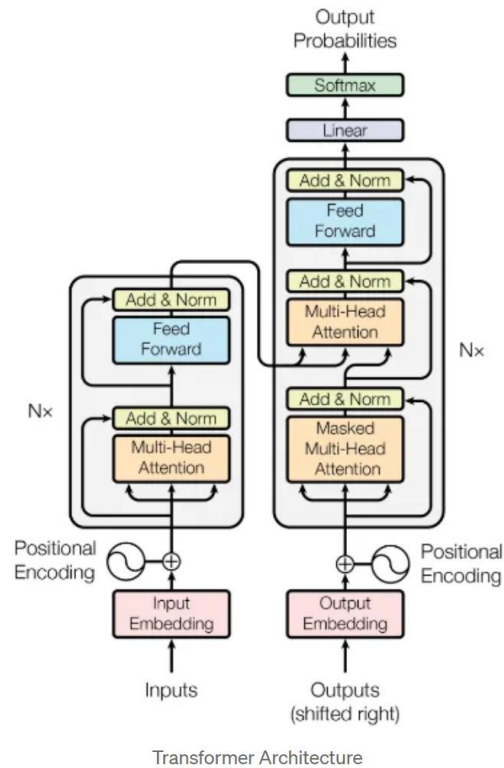
2. Small Dataset (C classes)



AlexNet:
64 x 3 x 11 x 11

(learned features, first stage)

The Transformer Architecture



Attention Is All You Need, Vaswani et al. (2017)

ViTs – Vision Transformers

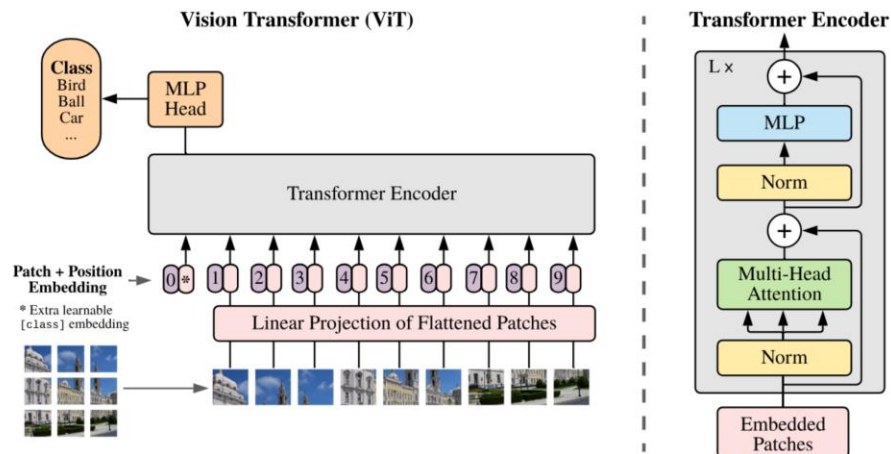
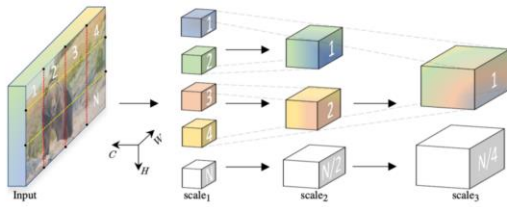
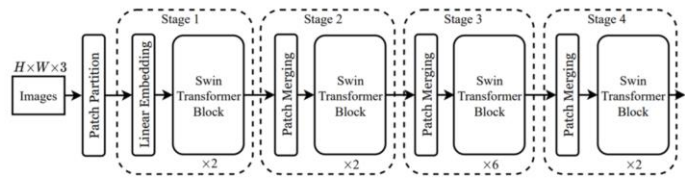


Figure from:
Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ArXiv 2020

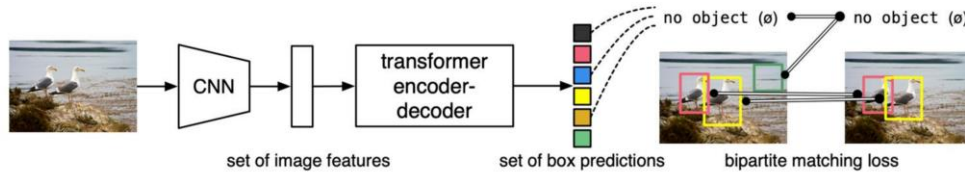
Vision Transformers



Fan et al, "Multiscale Vision Transformers", ICCV 2021



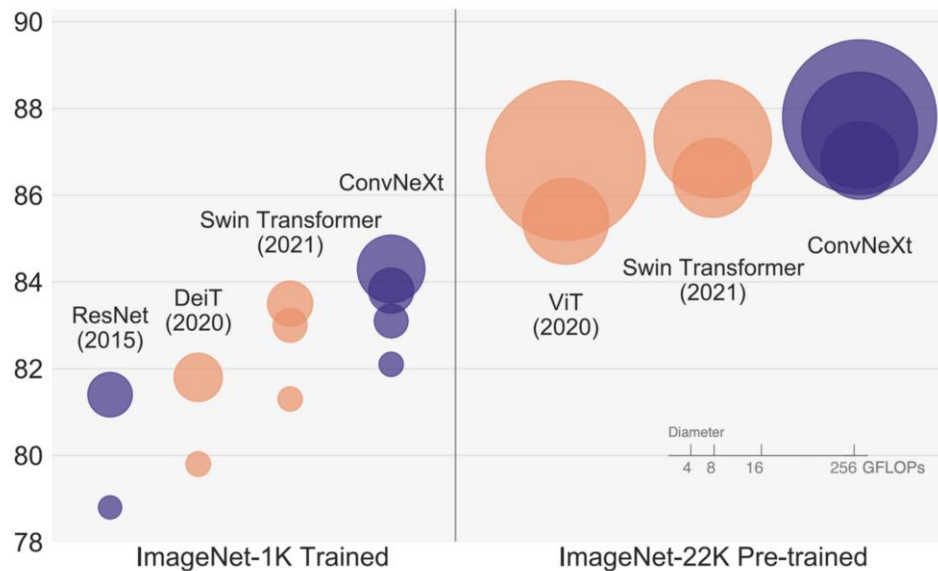
Liu et al, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows", CVPR 2021



Carion et al, "End-to-End Object Detection with Transformers", ECCV 2020

Performance

ImageNet-1K Acc.

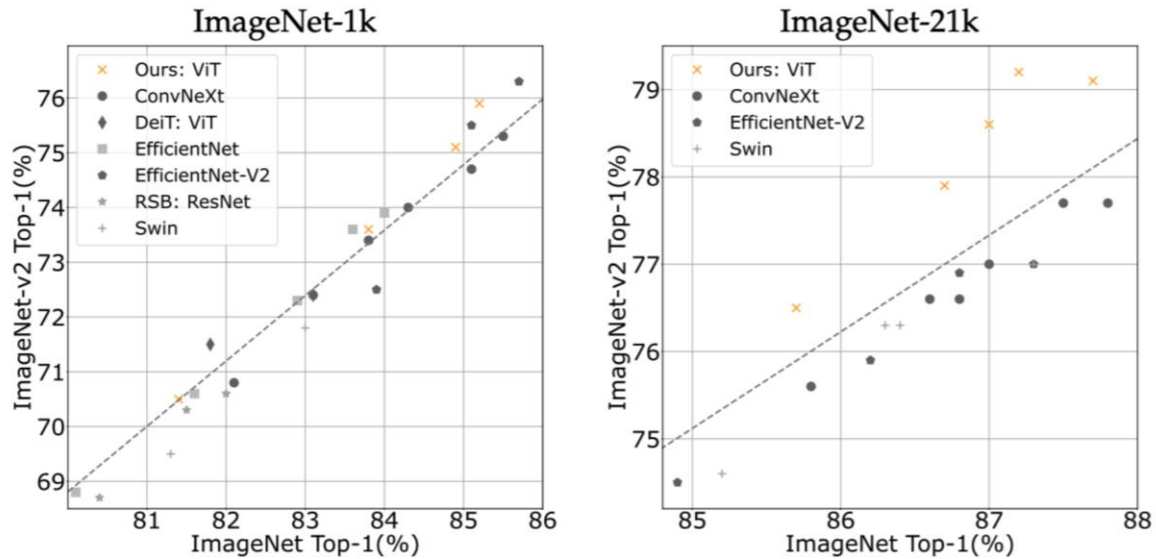


ConvNets strike back!

A ConvNet for the 2020s. Liu et al. CVPR 2022

DeiT III: Revenge of the ViT

Hugo Touvron^{*,†} Matthieu Cord[†] Hervé Jégou^{*}



Resources

Wikipedia

Deep Learning for Vision Lecture, Marc'Aurelio Ranzato – DeepMind

Deep Learning for CV Course Notes, Fei-Fei Li & Ehsan Adeli — Stanford Univ.

Attention Is All You Need, Vaswani et al., Advances in Neural Information Processing Systems (2017).

Computer Vision Course Notes, Srinath Sridhar – Brown Univ.

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) website.

Imagenet classification with deep convolutional neural networks, A Krizhevsky, I Sutskever, GE Hinton, Advances in neural information processing systems (2012).

Deep Residual Learning for Image Recognition, Kaiming He Xiangyu Zhang Shaoqing Ren Jian Sun, IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2015).