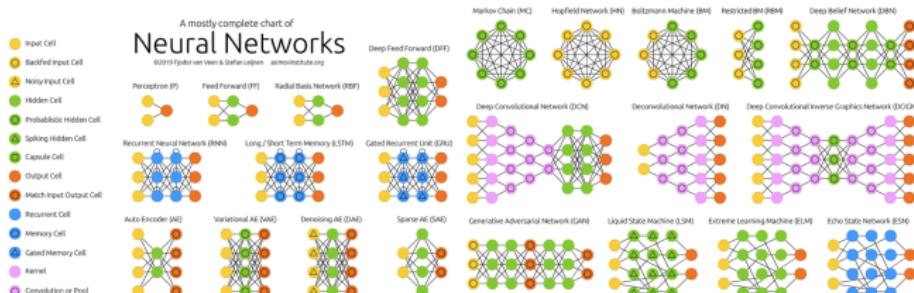


CNN ARCHITECTURES

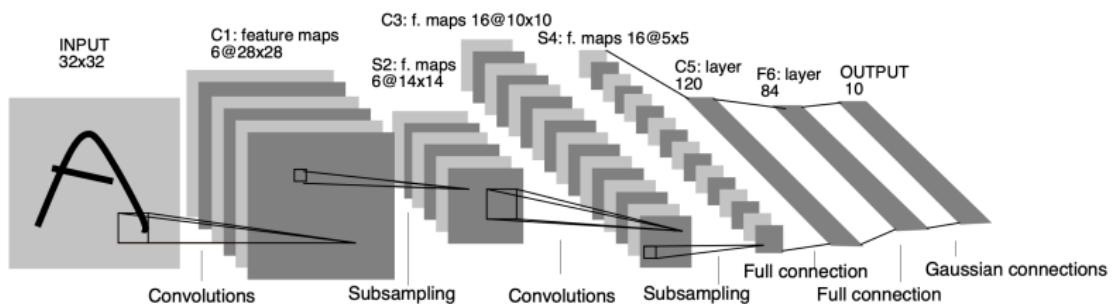
Vincent Barra
LIMOS, UMR 6158 CNRS, Université Clermont Auvergne

INTRODUCTION

- ▶ Since 2012, a HUGE collection of CNN models
- ▶ How/why the well known models have been built ?
- ▶ Finding right architecture : active area of research



LENET-5 (1998)



LENET-5 (1998)

Features

Architecture:

Conv((5× 5),6) - ReLU, Max Pooling - Conv((5× 5),16) - ReLU, Max Pooling -
FC(120)-ReLU - FC (84) - ReLU - FC(10) - Softmax

Features

One of the precursors

- ▶ First CNN to apply backpropagation
- ▶ 61 706 trainable parameters
- ▶ Number of connections limited
- ▶ Used for character recognition.



LIMOS
LABORATOIRE DINFORMATIQUE,
DE MODELISATION ET OPTIMISATION DES SYSTEMES



INTRODUCTION
○

LENET-5
○○●○○

ALEXNET
○○○○

ZEILER AND FERGUS
○○○○

VGG
○○○

INCEPTION
○○○○

RESIDUAL NETS
○○○○○○

CONCLUSION
○○

CONCLUSION
○

LENET-5 (1998)



INTERLUDE

Need for data

Can't add too much depth because no relevant training set available.... until 2012

⇒ ImageNet (10,000,000 labeled images depicting 10,000+ object categories for training).

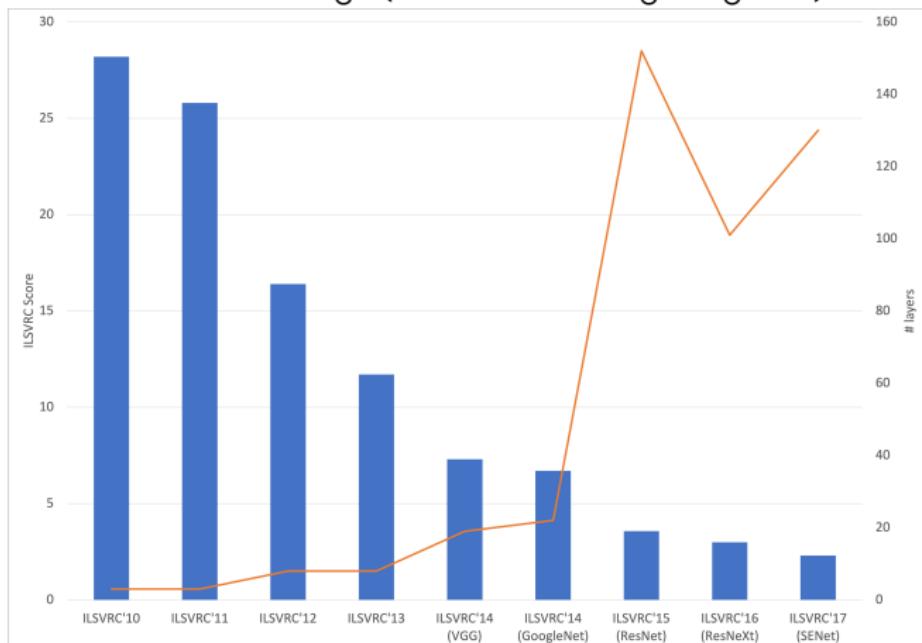
(today 14,197,122 images, 21841 synsets indexed).



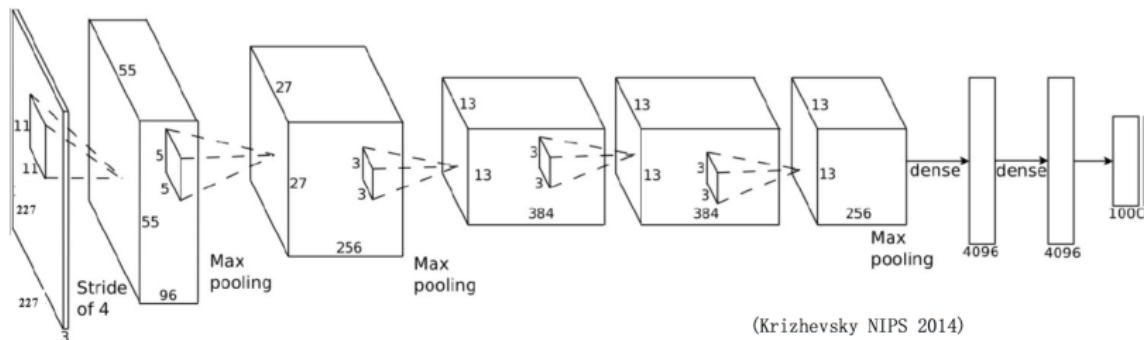
INTRODUCTION
○LENET-5
○○○○●ALEXNET
○○○○ZEILER AND FERGUS
○○○○VGG
○○○INCEPTION
○○○○RESIDUAL NETS
○○○○○CONCLUSION
○○CONCLUSION
○

INTERLUDE

ILSVRC Challenge (classification using ImageNet)



ALEXNET (2012)



(Krizhevsky NIPS 2014)

ALEXNET (2012)

Features

Architecture: (simplified)

Conv(11×11,96) - ReLU - Max Pooling - Conv(5×5,192) - ReLU - Max Pooling
- Conv(3×3,384) - ReLU - Conv(3×3,256) - ReLU - Conv(3×3,256) - ReLU -
MaxPool - Softmax

Features

- ▶ ReLU, local response normalization, data augmentation, dropout
- ▶ More than $62 \cdot 10^6$ trainable parameters
- ▶ stochastic gradient descent with momentum

INTRODUCTION
○

LENET-5
○○○○○

ALEXNET
○○●○

ZEILER AND FERGUS
○○○○

VGG
○○○

INCEPTION
○○○○

RESIDUAL NETS
○○○○○○

CONCLUSION
○○

CONCLUSION
○

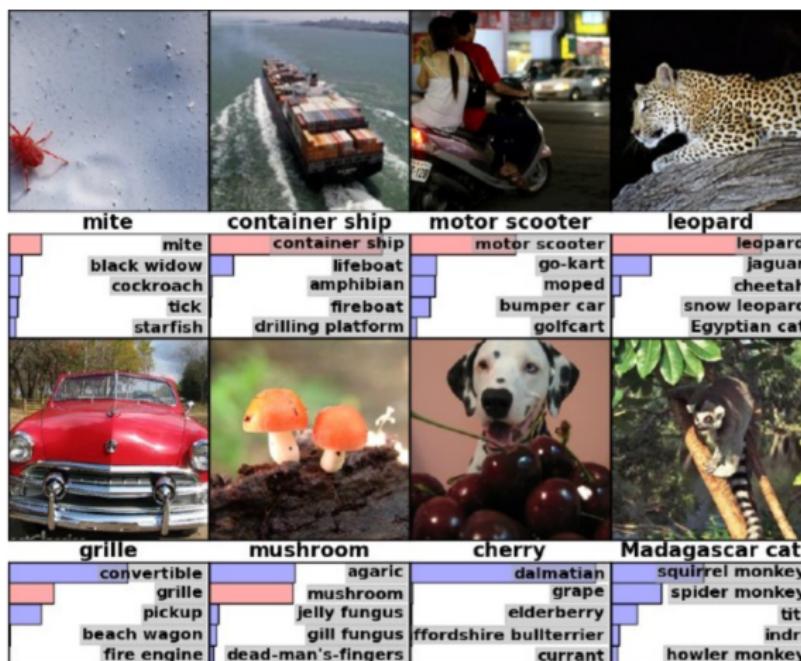
ALEXNET (2012)

Layer 1 - Learned filters

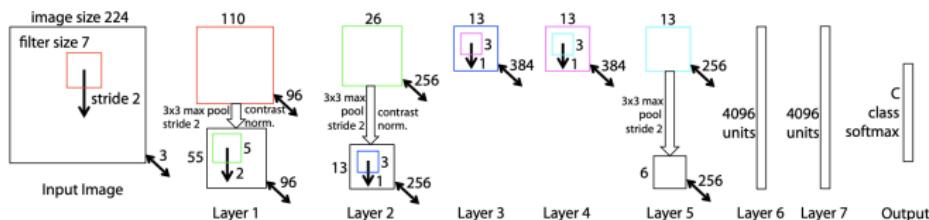


ALEXNET (2012)

AlexNet Results.



ZEILER AND FERGUS (2014)



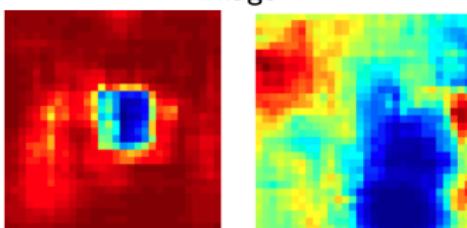
Features

- ▶ Refinement of AlexNet
- ▶ More than $107 \cdot 10^6$ trainable parameters
- ▶ Layer widths adjusted by cross-validation \Rightarrow Depth matters

ZEILER AND FERGUS (2014)



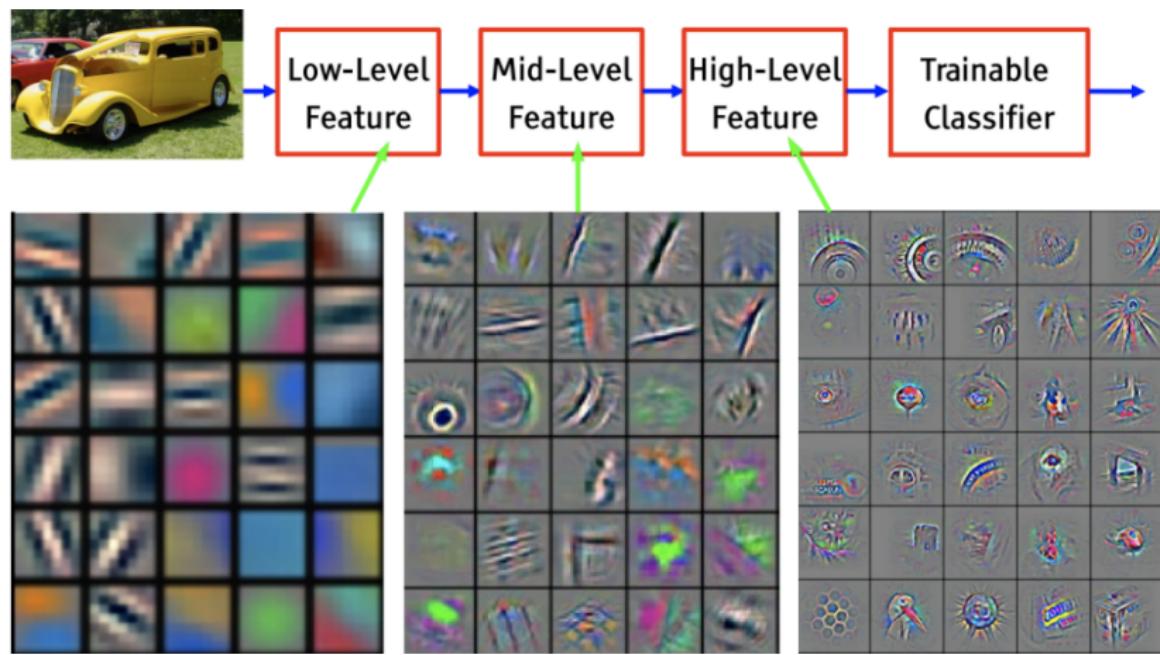
image



correct class probability

- ▶ image occluded by gray square at different positions (x, y)
- ▶ correct class probability : function of (x, y)

ZEILER AND FERGUS (2014)



Hierarchical representation

INTRODUCTION

LENET-5

ALEXNET

ZEILER AND FERGUS

VGG

INCEPTION

RESIDUAL NETS

CONCLUSION

CONCLUSION

○

○○○○○

○○○○

○○○●

○○○

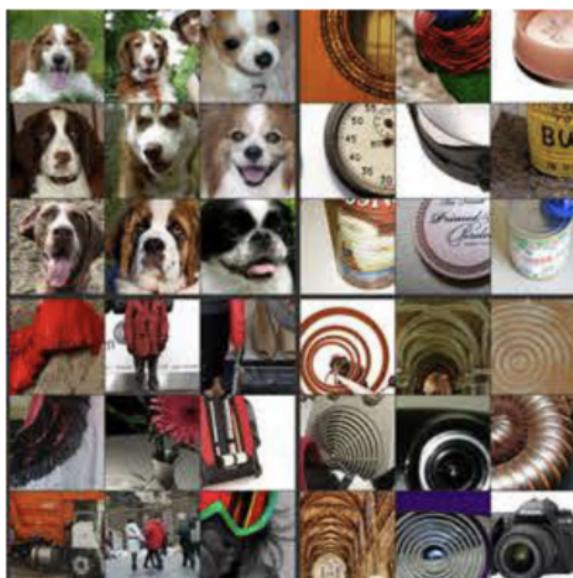
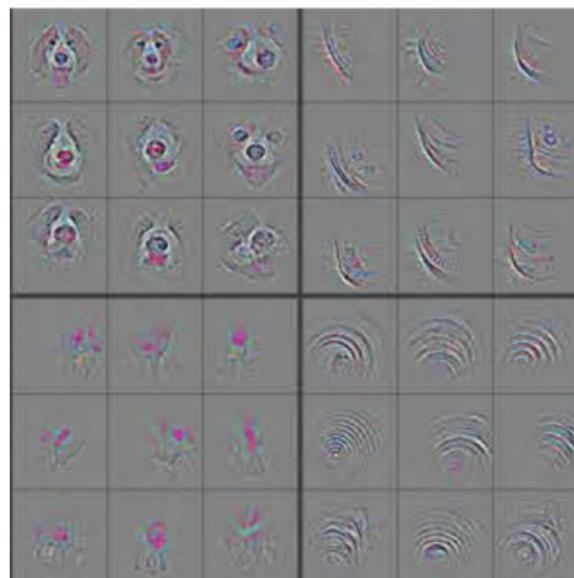
○○○○

○○○○○

○○

○

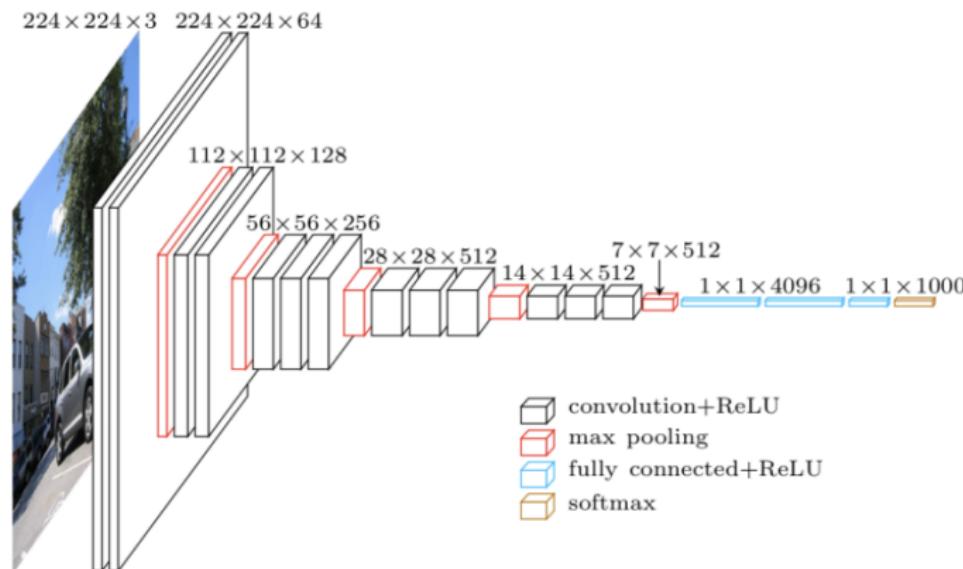
ZEILER AND FERGUS (2014)



Reconstructed patterns from top 9 activations of selected features of layer 4
and corresponding image patches

VGG (2015)

Very Deep CNN for Large-Scale Image Recognition. Family of networks.



VGG (2015)

Very Deep CNN for Large-Scale Image Recognition

A	A-LRN	B	C	D	E
11 layers	11 layers	13 layers	16 layers	16 layers	19 layers
133 M param	133 M param	133 M param	134 M param	138 M param	144M param
Input: 224 × 224 RGB image					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
max pooling					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
max pooling					
conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
max pooling					
conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
max pooling					
conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
max pooling					
FC layer 4096 neurons					
FC layer 4096 neurons					
FC layer 1000 neurons					
Softmax					

VGG (2015)

Features

- ▶ 3×3 filters: smallest possible filters to capture 4-connectivity
 - ▶ $2 \ 3 \times 3$ filters = receptive field of a 5×5 filter
 - ▶ $3 \ 3 \times 3$ filters = receptive field of a 7×7 filter
- ⇒ A large filter can be replaced by a deeper stack of successive smaller filters
- ⇒ More nonlinearities (ReLU), fewer parameters $(3 \times 3 \times K) \cdot 2 = 18K / (5 \times 5 \times K) \cdot 1 = 25K$
- ⇒ Extension to 1×1 filters+nonlinearity : increase nonlinearities without affecting receptive field sizes.

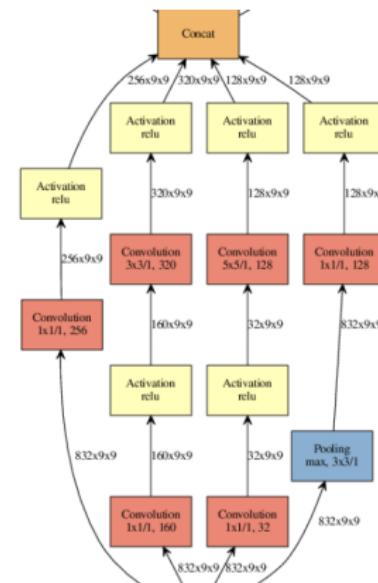
Result: 7.3% error rate in ImageNet (18.2% for AlexNet)

Although deeper, number of weights is not exploding.

GOOGLENET (2015)

Features

- ▶ Incarnation of the Inception architecture
- ▶ Multiple feedforward pathways



INCEPTION FAMILY

Idea

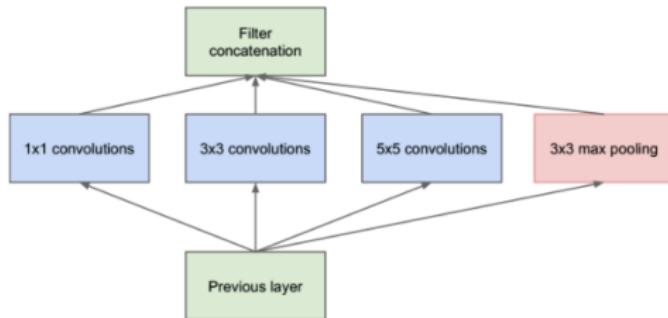
- ▶ Salient parts may have great variation in sizes
- ▶ The receptive field should vary in size
- ▶ Stacking CNN is expensive & may overfit.



INCEPTION FAMILY

Idea

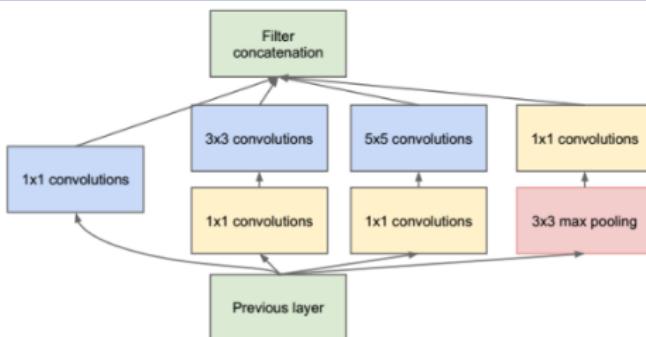
- ▶ Multiple kernel filters of different sizes ⇒ Expensive !



INCEPTION FAMILY

Idea

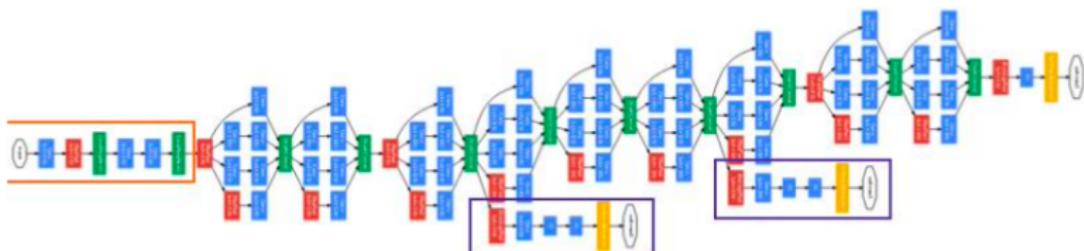
- ▶ Multiple kernel filters of different sizes \Rightarrow Expensive !
- ▶ Solution: add 1×1 convolutions



INCEPTION FAMILY

Idea

- ▶ Multiple kernel filters of different sizes \Rightarrow Expensive !
- ▶ Solution: add 1×1 convolutions
- ▶ 9 inception modules, 27 layers



INCEPTION FAMILY

Idea

- ▶ Multiple kernel filters of different sizes ⇒ Expensive !
- ▶ Solution: add 1×1 convolutions
- ▶ 9 inception modules, 27 layers

6.67% Imagenet error, compared to 18.2% of Alexnet, 25 times less parameters and faster than AlexNet

INCEPTION FAMILY

Idea

- ▶ Multiple kernel filters of different sizes ⇒ Expensive !
- ▶ Solution: add 1×1 convolutions
- ▶ 9 inception modules, 27 layers

The network was too deep :-(⇒ vanishing gradient

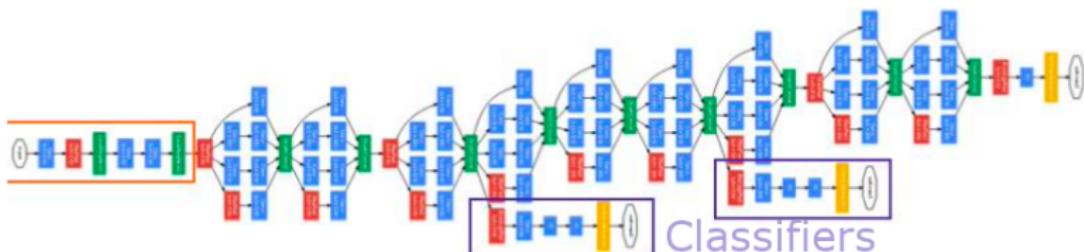
Additional ressource

See Slides "Vanishing gradient".

INCEPTION FAMILY

Idea

- ▶ Multiple kernel filters of different sizes \Rightarrow Expensive !
- ▶ Solution: add 1×1 convolutions
- ▶ 9 inception modules, 27 layers

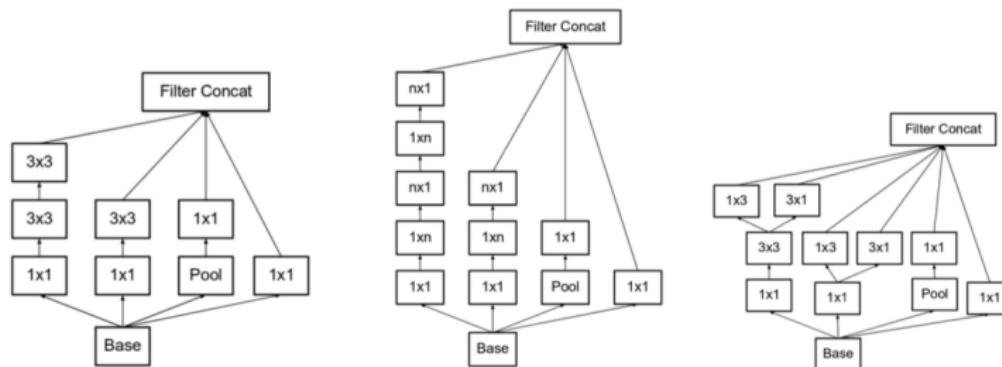


Solution: Intermediate classifiers

INCEPTION FAMILY

v2,v3,v4,v5

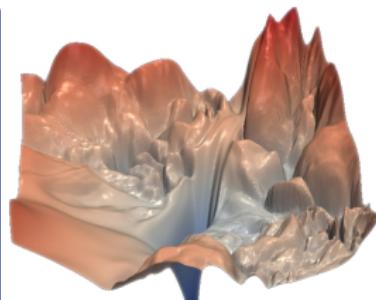
- ▶ Factorization of 5×5 in 2 3×3
- ▶ Factorization of $n \times n$ in 2 $n \times 1$ and $1 \times n$
- ▶ Optimization function (RMSProp)
- ▶ Batch normalization...



INTRODUCTION

Some facts

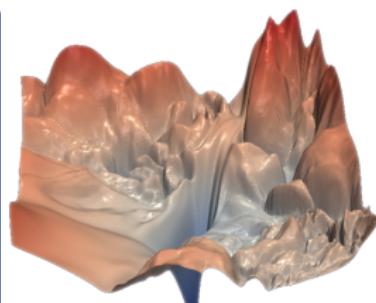
- ▶ Very deep networks stop learning after some time
- ▶ max accuracy → saturation → unlearning
- ▶ causes : not only overfitting and vanishing gradients but harder optimization



INTRODUCTION

Some facts

- ▶ Very deep networks stop learning after some time
- ▶ max accuracy → saturation → unlearning
- ▶ causes : not only overfitting and vanishing gradients but harder optimization



Some facts

- ▶ Although identity is representable, learning it proves difficult for optimization methods
- ▶ Intuition: tweak the network so it doesn't have to learn identity connections

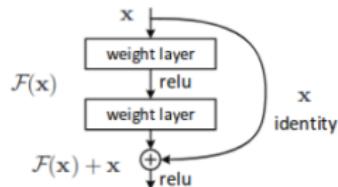
RESIDUAL BLOCK

Residual idea

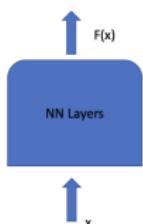
Learn residual functions rather than direct mappings.

If x is the input to the nonlinear layer, and \mathcal{F} the stacked nonlinearities

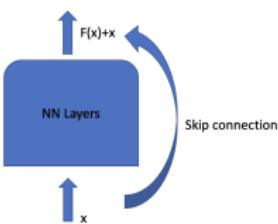
$$\mathcal{F}(x) = \mathcal{R}(x) - x \Rightarrow \mathcal{R}(x) = \mathcal{F}(x) + x$$



RESIDUAL BLOCK

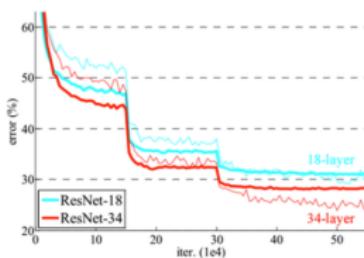
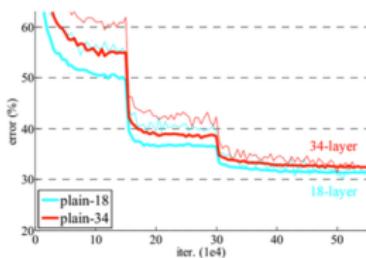


Hard to get $F(x)=x$ and
Make $y=x$ identity mapping



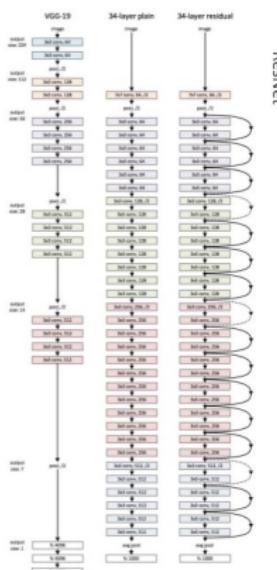
Easy to get $F(x)=0$ and
Make $y=x$ identity mapping

With residual connections, deeper nets improve their scores (here on ImageNet)



RESIDUAL BLOCK

- ▶ Very low errors on imageNet (all challenges)
- ▶ Up to 1202 layers
- ▶ Skip connections \Rightarrow faster training
- ▶ Batch norm compulsory (vanishing gradient)
- ▶ Comparison with VGG-D
parameters: 25M / 138M
- ▶ Comparison with VGG-D
computational complexity: 3.8B / 15.3B Flops.



RESIDUAL NETWORKS

Why do they work ?

- ▶ ResNets \approx implicitly ensembling shallower networks: with r residual modules there are 2^r possible subsets of modules
- ▶ Able to learn unrolled iterative refinements

RESIDUAL NETWORKS

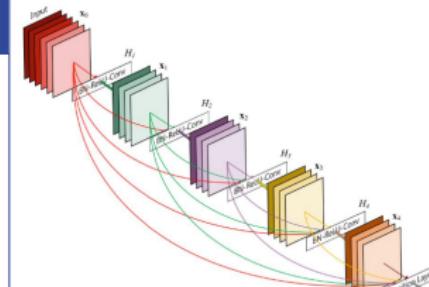
HighwayNet

≈ ResNets + gate with learnable parameters on the importance of each skip connection

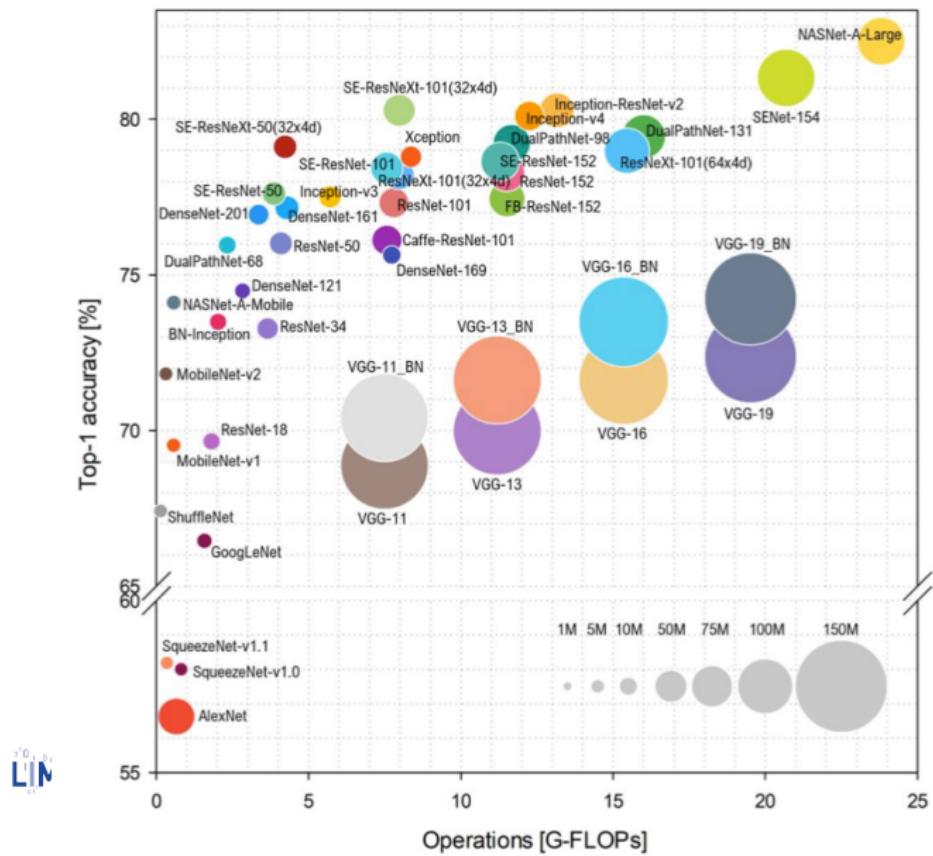
$$y = T(x, W_t) \cdot \mathcal{R}(x, W_h) + (1 - T(x, W_t)) \cdot x$$

DenseNet

- ▶ Add skip connections to multiple forward layers
- ▶ Each layer is connected to every subsequent layer
- ▶ Interest: combinaison of detected patterns



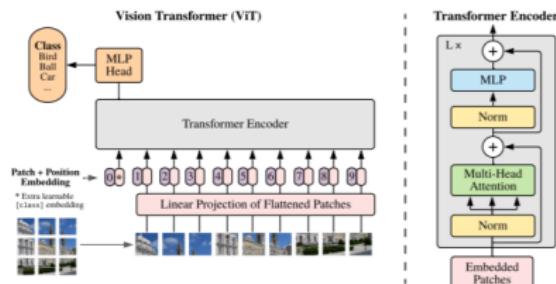
STATE-OF-THE-ART



CURRENT TRENDS (CHANGE RAPIDLY...)

Some issues

- ▶ Reduce filter sizes (except possibly at the lowest layer), factorize filters
- ▶ 1×1 convolutions to reduce and expand the number of feature maps
- ▶ Use skip connections
- ▶ Attention mechanisms
- ▶ Visual transformers
- ▶ NAS (Network Architecture Search)



STATE-OF-THE-ART

- ▶ All these models are already trained on huge databases (ImageNet, CIFAR...).
 - ▶ For a “close” problem, no need to train the CNN from scratch...
 - ▶ Subject of a further lecture (teasing :-))

See...

Lecture “Transfer Learning”.