

FACULTY OF MATHEMATICS,
PHYSICS AND INFORMATICS

Comenius University
Bratislava

Neural Networks for Computer Vision

Lecture 8: CNN Architectures

Ing. Viktor Kocur, PhD., RNDr. Zuzana Černeková, PhD.

8.11.2022

Contents



- Recap - CNNs
- Architectures
 - ▶ AlexNet
 - ▶ VGG
 - ▶ Inception
 - ▶ ResNets
 - ▶ Efficient Architectures
 - ▶ ConvNext
- Neural Architecture Search
- Design Spaces
- Transfer Learning

Acknowledgment



Some of the slides are directly adopted from slides for CS231n¹ course at Stanford University!

¹Fei-Fei Li, Ranjay Krishna, and Danfei Xu. *Stanford CS231n lecture slides*. <http://cs231n.stanford.edu/slides/>

Recap - CNNs



Main building blocks of a CNN

- Convolutional layers
- Pooling layers
- Fully-connected layers
- Activations

Recap - CNNs

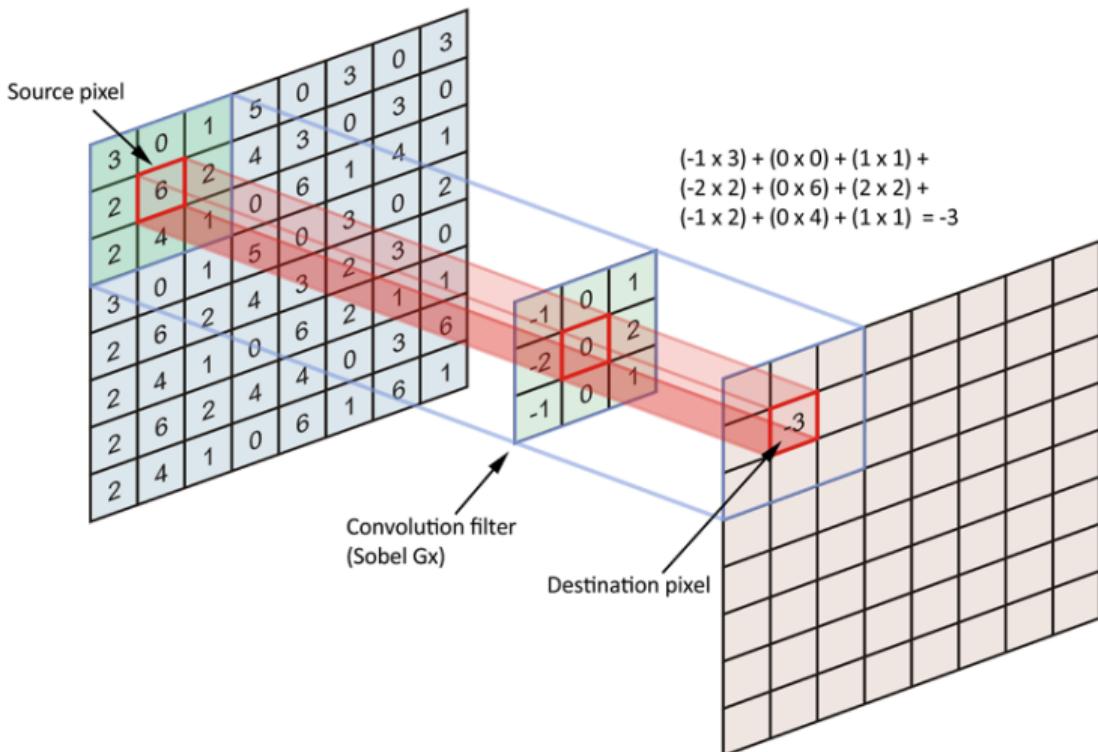


Main building blocks of a CNN

- Convolutional layers
- Pooling layers
- Fully-connected layers
- Activations

There are many potential combinations! We will now discuss historically most important architectures and what we can learn from them!

Recap - Convolution



Recap - Convolution

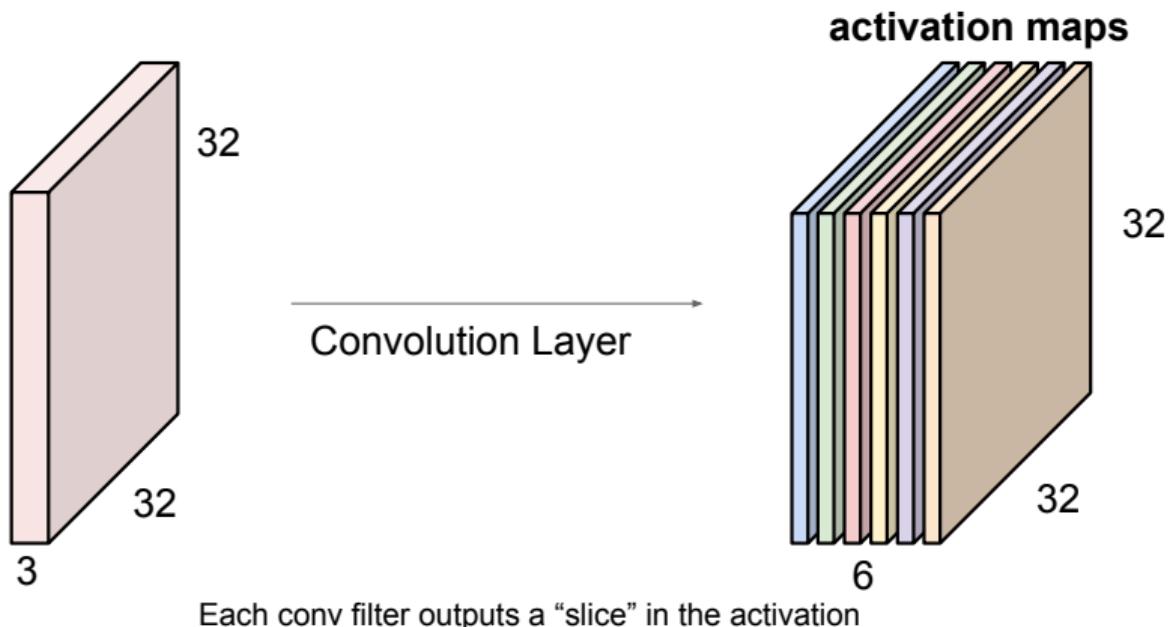


For one dimension of the input image (height or width independently) after applying convolution we get the following formula for the dimension of the output:

$$N_{out} = \frac{N_{in} - F + 2P}{S} + 1, \quad (1)$$

where N_{out} is the output size, N_{in} is the input size, F is the size of the kernel, P is the padding and S is the stride.

Recap - Convolution



Recap - Parameters



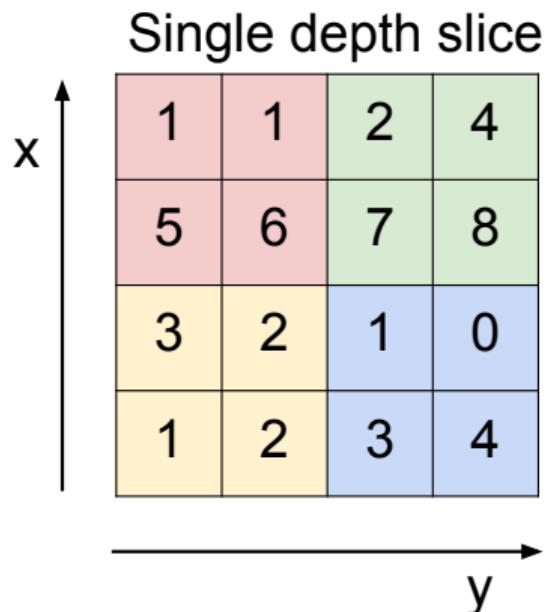
For a convolution with C_{in} input channels, C_{out} output channels, kernel size $K_h \times K_w$ the number of parameters is:

$$P_{conv} = (C_{in} \cdot K_w \cdot K_h + 1) \cdot C_{out} \quad (2)$$

For a fully-connected layer with C_{in} input channels/neurons and C_{out} output channels is:

$$P_{fc} = (C_{in} + 1) * C_{out} \quad (3)$$

Recap - Pooling

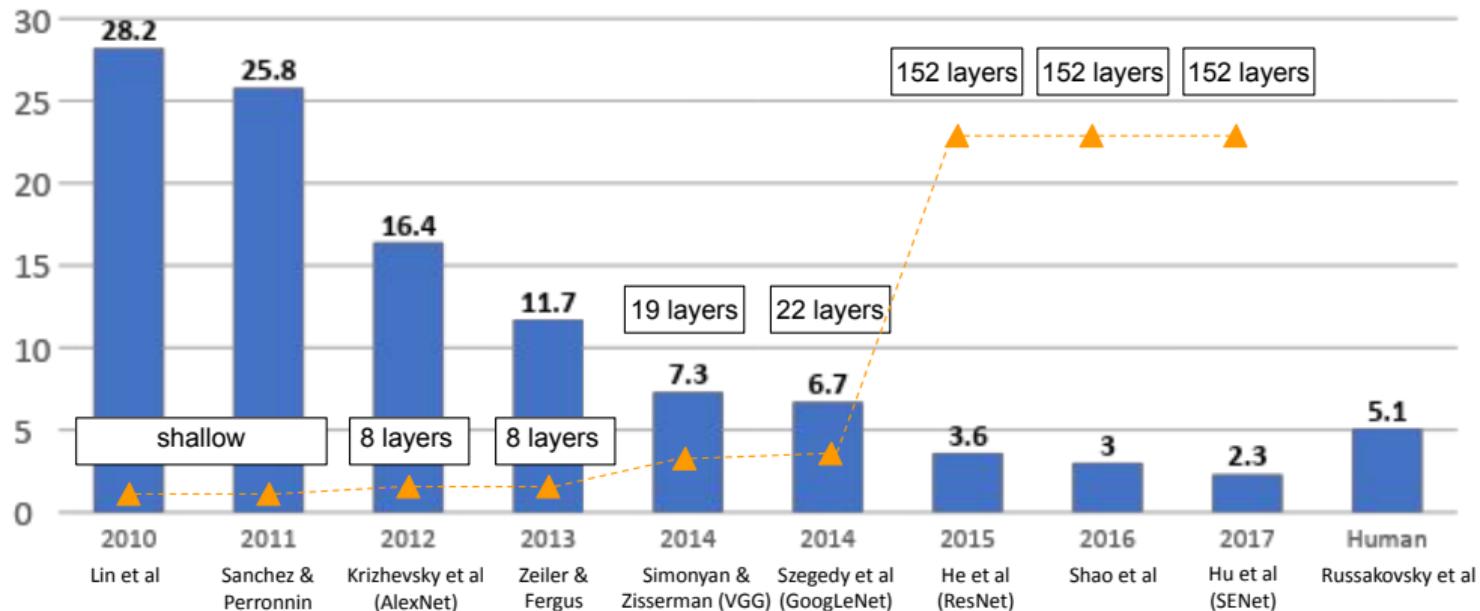


max pool with 2x2 filters
and stride 2

The result of applying a 2x2 max pooling filter with stride 2 to the input grid. The output is a 2x2 grid where each cell contains the maximum value from its 2x2 receptive field in the input. The top-left cell (6) is pink, top-right (8) is light green, bottom-left (3) is light yellow, and bottom-right (4) is light blue.

6	8
3	4

ILSVRC Winners





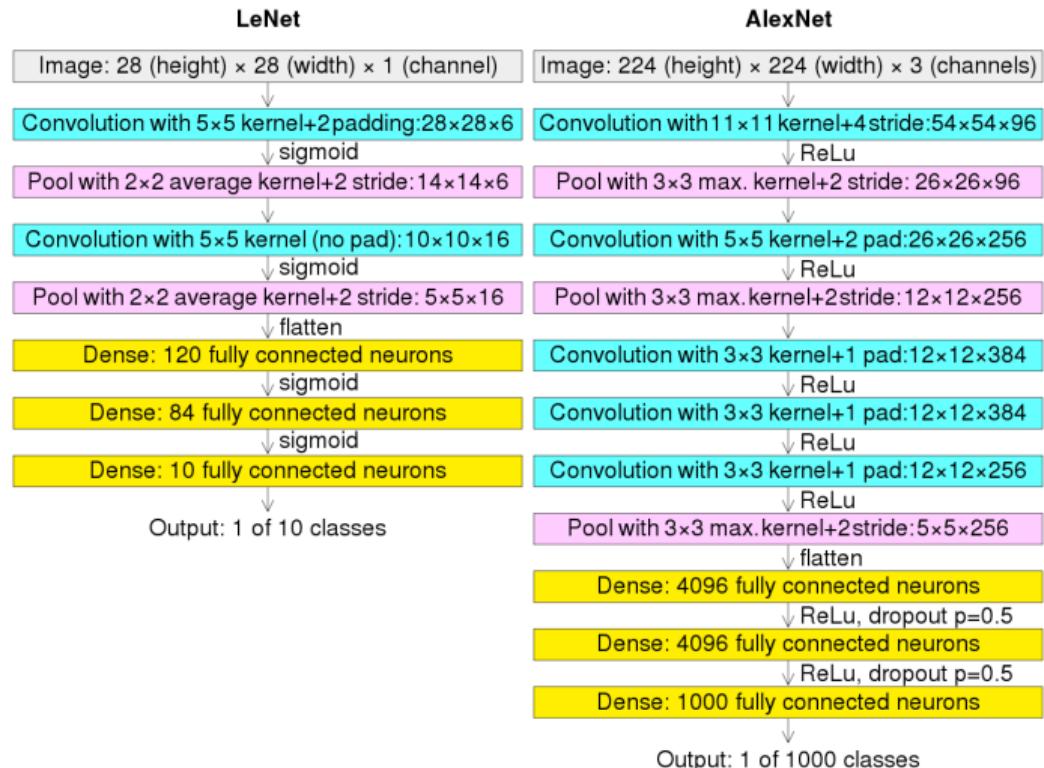
AlexNet² won the ILSVRC 2012 and started CNN revolution!

- Based on earlier work - LeNet³
- Relied on heavy data augmentation
- Training on two GPUs
- Used ReLU activations
- Batch size: 128
- SGD momentum + manual learning rate changes
- L2 regularization
- Ensemble of 7 CNNs to reduce error

²Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. "Imagenet classification with deep convolutional neural networks." In: *Advances in neural information processing systems* 25 (2012), pp. 1097–1105

³Yann LeCun et al. "Backpropagation applied to handwritten zip code recognition." In: *Neural computation* 1.4 (1989), pp. 541–551

AlexNet vs LeNet





ZFNet⁴ won the challenge next year with some modifications:

- 7×7 stride 2 conv instead of 11×11 stride 4 in first layer
- More channels

⁴ Matthew D Zeiler and Rob Fergus. "Visualizing and understanding convolutional networks." In: *European conference on computer vision*. Springer. 2014, pp. 818–833

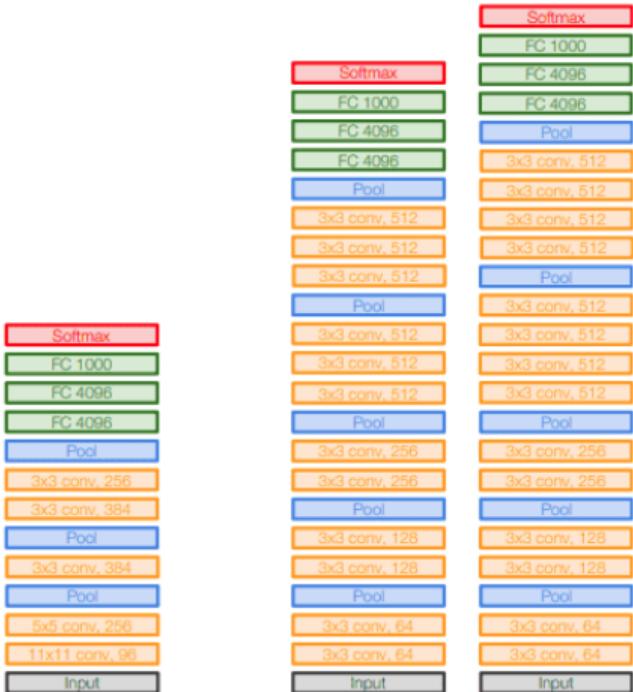


The next significant advance were VGG⁵ networks. The main differences were:

- Using only 3×3 convolutions
- Deeper networks are better!
- Trained early layers first without later ones then joined them together.

⁵Karen Simonyan and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition." In: *arXiv preprint arXiv:1409.1556* (2014)

VGG

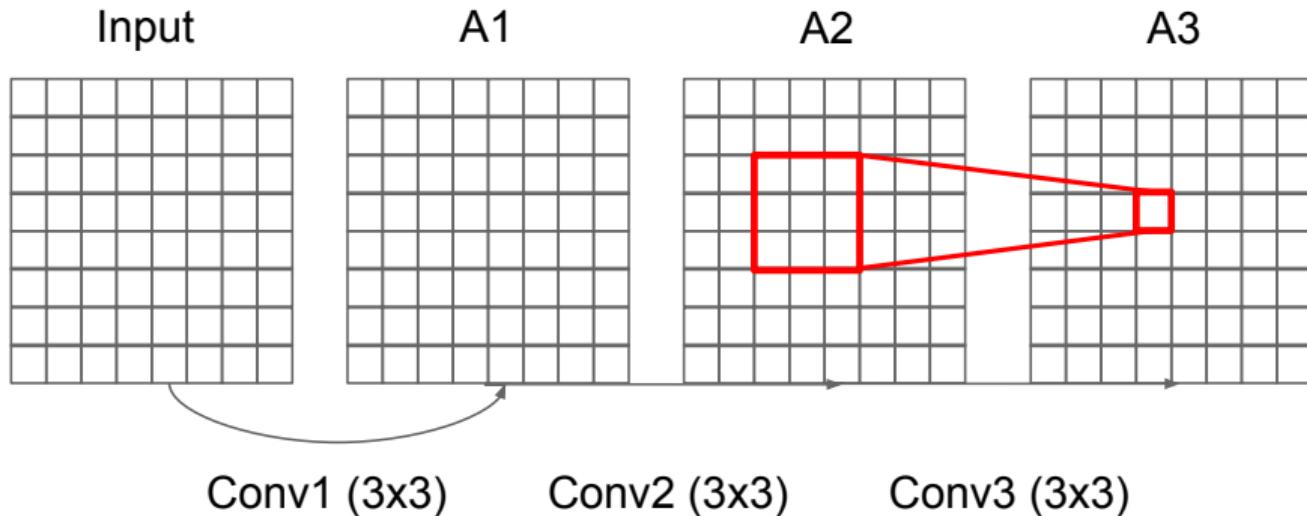


AlexNet

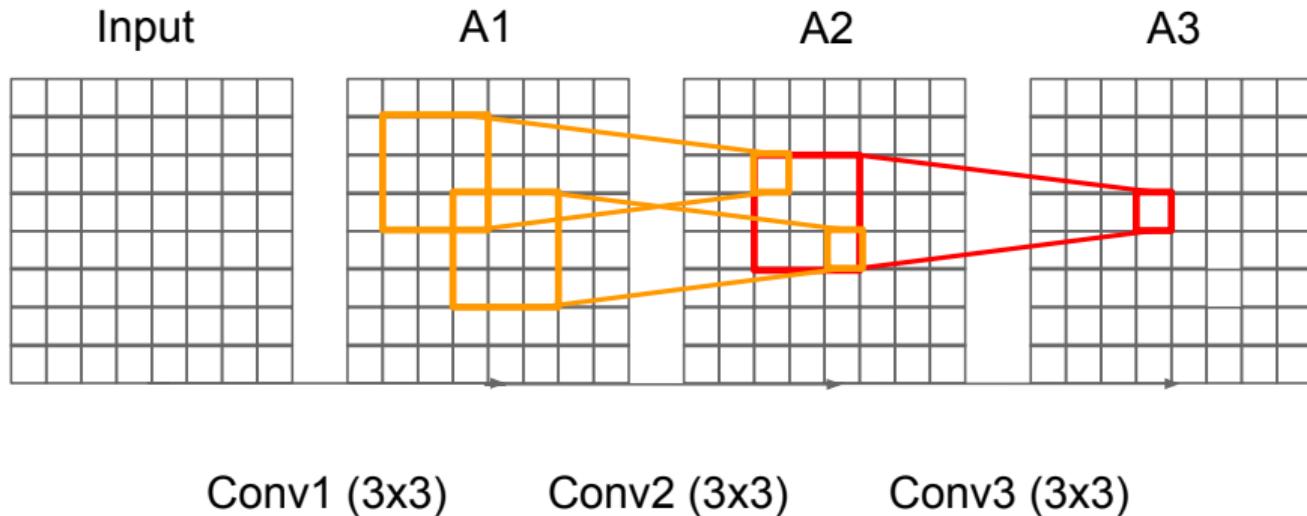
VGG16

VGG19

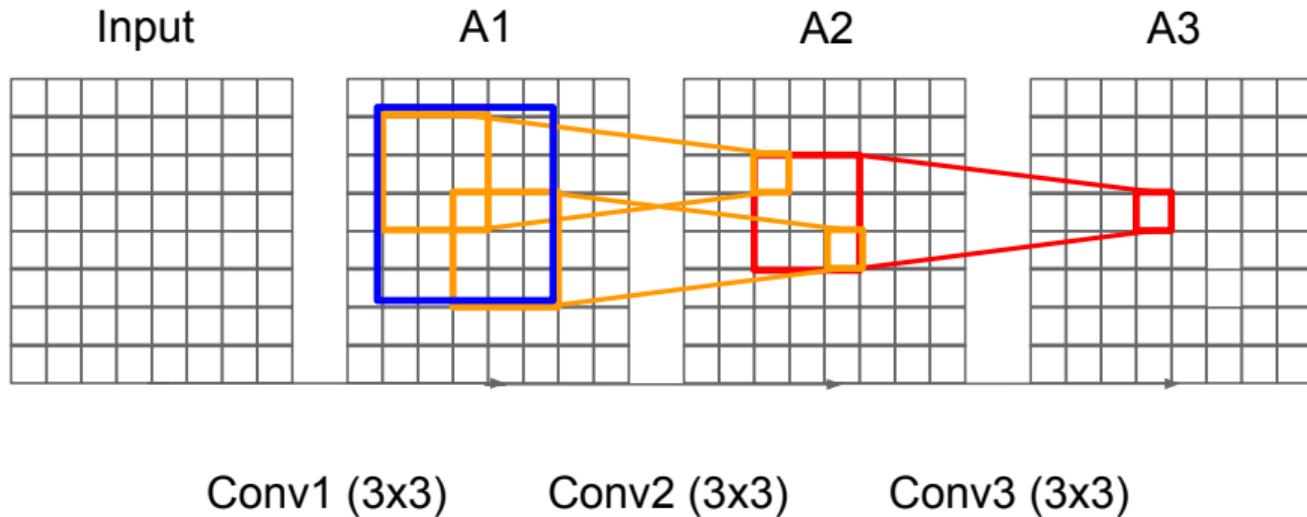
Stacking convolutions with smaller kernels



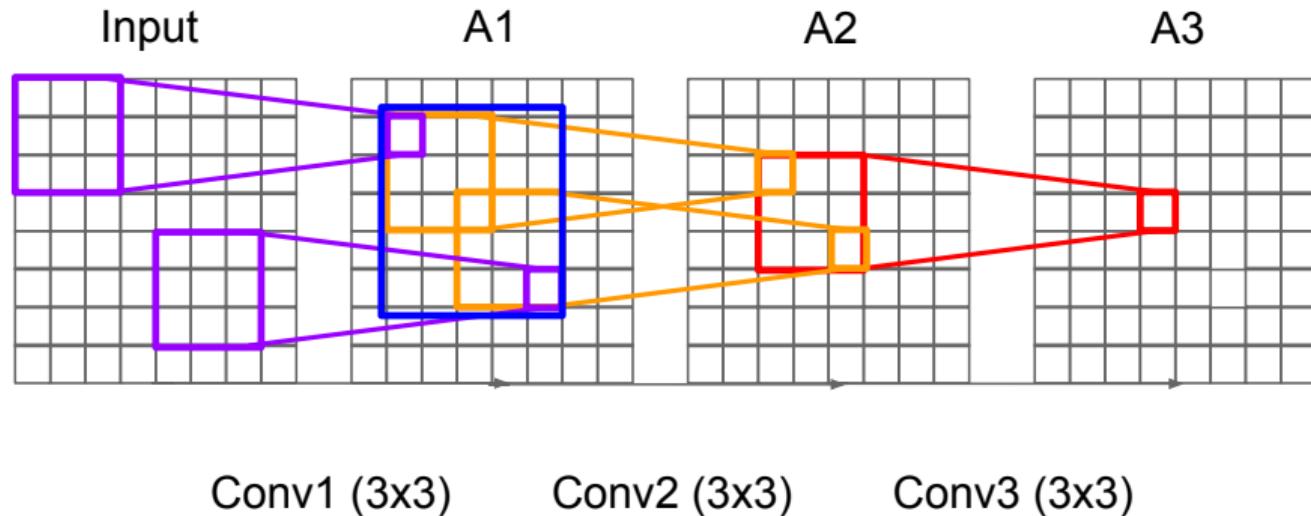
Stacking convolutions with smaller kernels



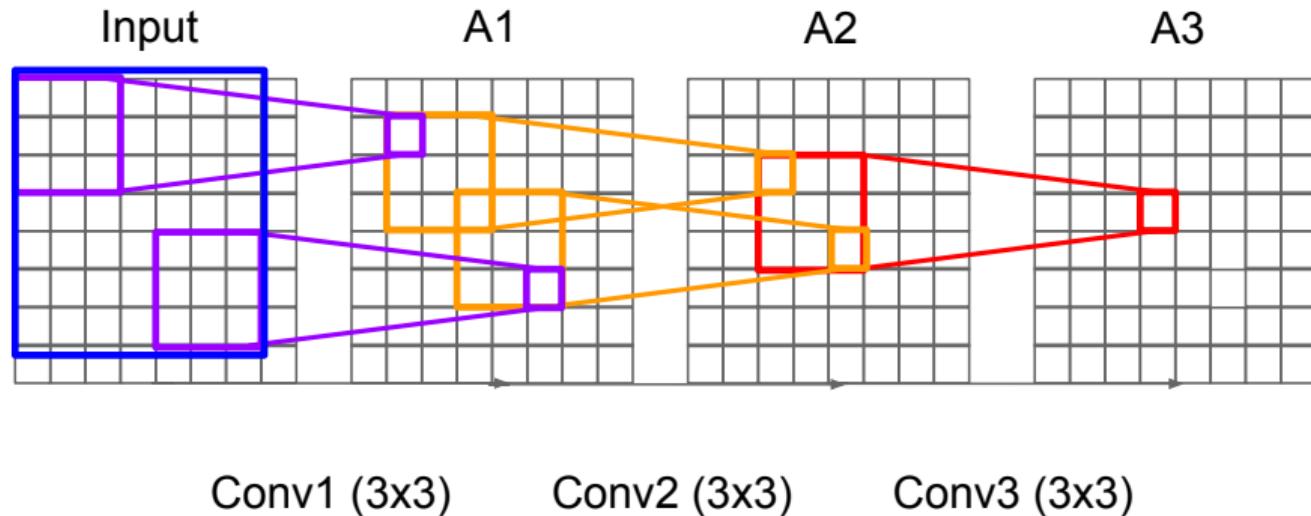
Stacking convolutions with smaller kernels



Stacking convolutions with smaller kernels



Stacking convolutions with smaller kernels



VGG parameters and memory



INPUT: [224x224x3] **memory:** 224*224*3=150K **params:** 0 **(not counting biases)**

CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2\text{M}$ params: $(3 \times 3 \times 3) \times 64 = 1.728$

CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2M$ params: $(3 \times 3 \times 64) \times 64 = 36,864$

POOL2: [112x112x64] memory: 112*112*64=800K params: 0

CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6M$ params: $(3 \times 3 \times 64) \times 128 = 73.728$

CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: $(3 \times 3 \times 128) \times 128 = 147,456$

POOL2: [56x56x128] memory: $56 \times 56 \times 128 = 400\text{K}$ params: 0

CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800\text{K}$ params: $(3 \times 3 \times 128) \times 256 = 294.912$

CONV3-256: [56x56x256] **memory:** $56 \times 56 \times 256 = 800\text{K}$ **params:** $(3 \times 3 \times 256) \times 256 = 589,824$

CONV3-256: [56x56x256] **memory:** $56 \times 56 \times 256 = 800\text{K}$ **params:** $(3 \times 3 \times 256) \times 256 = 589,824$

POOL2: [28x28x256] memory: $28 \times 28 \times 256 = 200K$ params: 0

CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400\text{K}$ params: $(3 \times 3 \times 256) \times 512 = 1,179,648$

CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [28x28x512] **memory:** $28*28*512=400K$ **params:** $(3*3*512)*512 = 2,359,296$

POOL2: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: 0

CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [14x14x512] **memory:** $14 \times 14 \times 512 = 100K$ **params:** $(3 \times 3 \times 512) \times 512 = 2,359,296$

POOL2: [7x7x512] memory: $7 \times 7 \times 512 = 25K$ params: 0

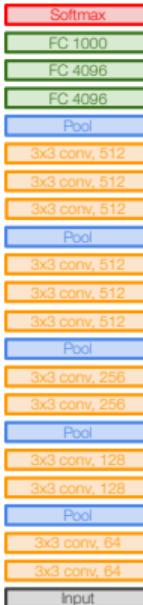
FC: [1x1x4096] memory: 4096 params: $7 \times 7 \times 512 \times 4096 = 102,760,448$

FC: [1x1x4096] memory: 4096 params: $4096 \times 4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params: $4096 \times 1000 = 4,096,000$

TOTAL memory: $24M * 4$ bytes $\approx 96MB$ / image (for a forward pass)

TOTAL params: 138M parameters

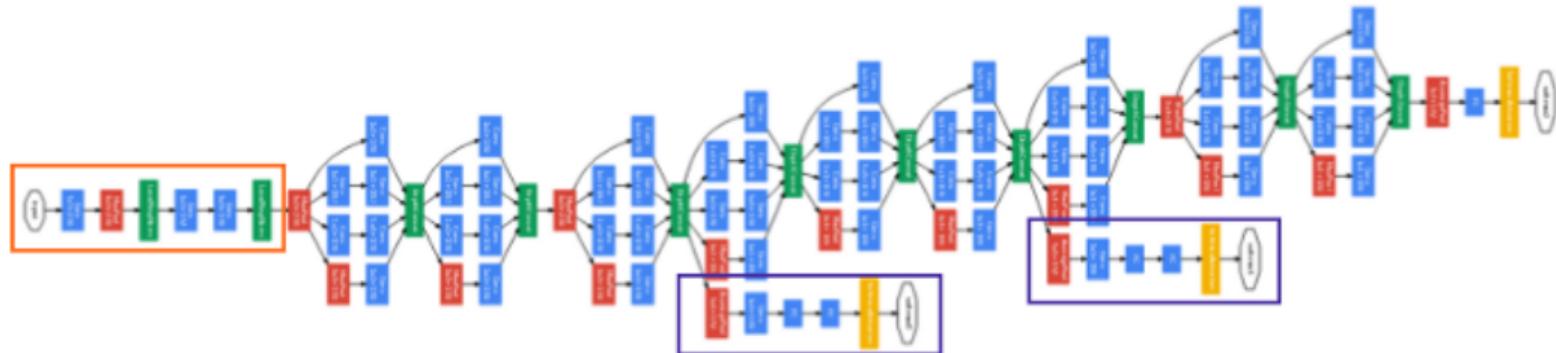


VGG16

Inception v1

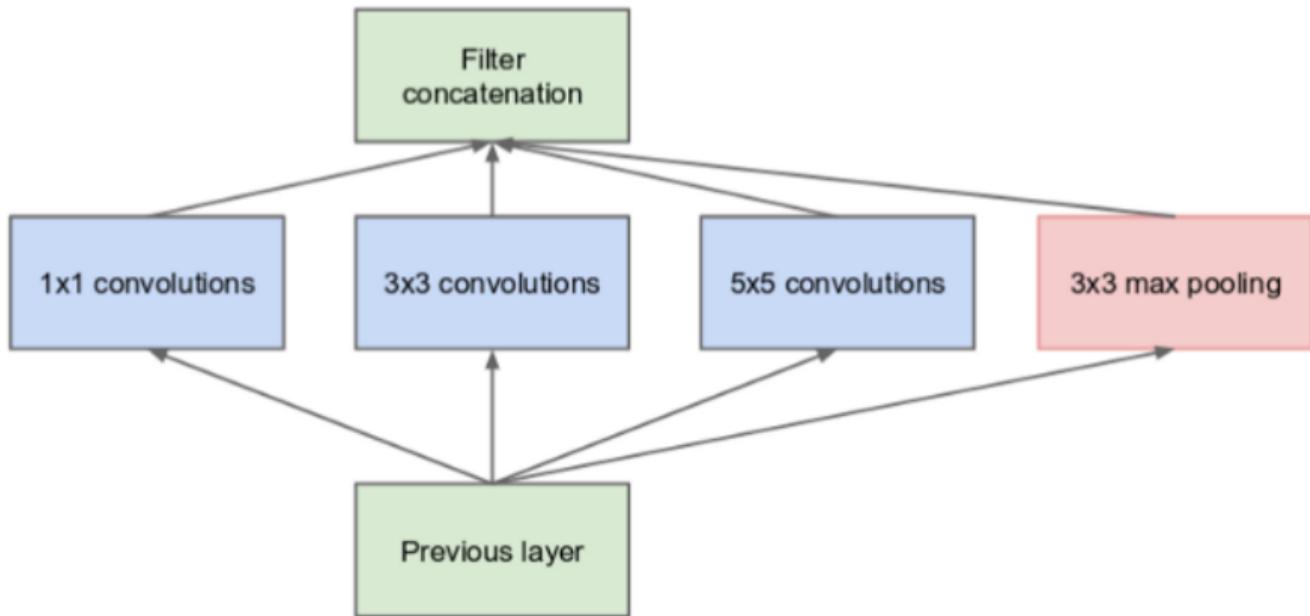


The winner of ILSVRC in 2014 was Googlenet, a.k.a Inception V1⁶

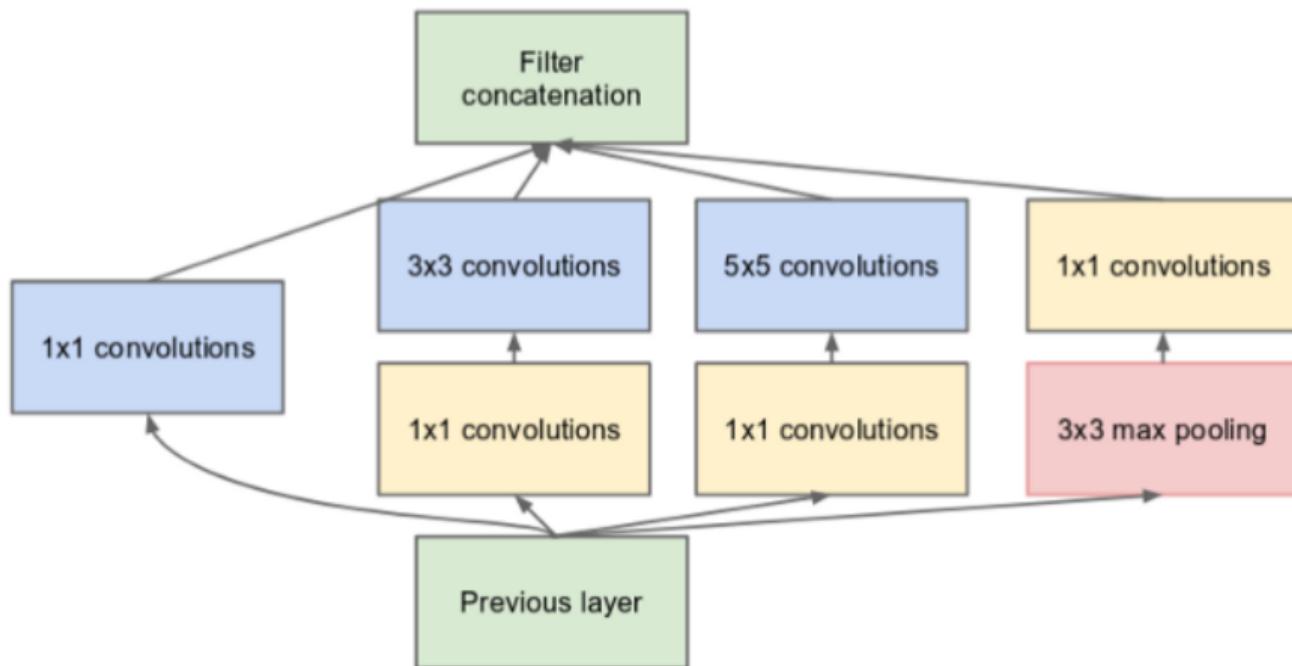


⁶Christian Szegedy et al. "Going deeper with convolutions." In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 1–9

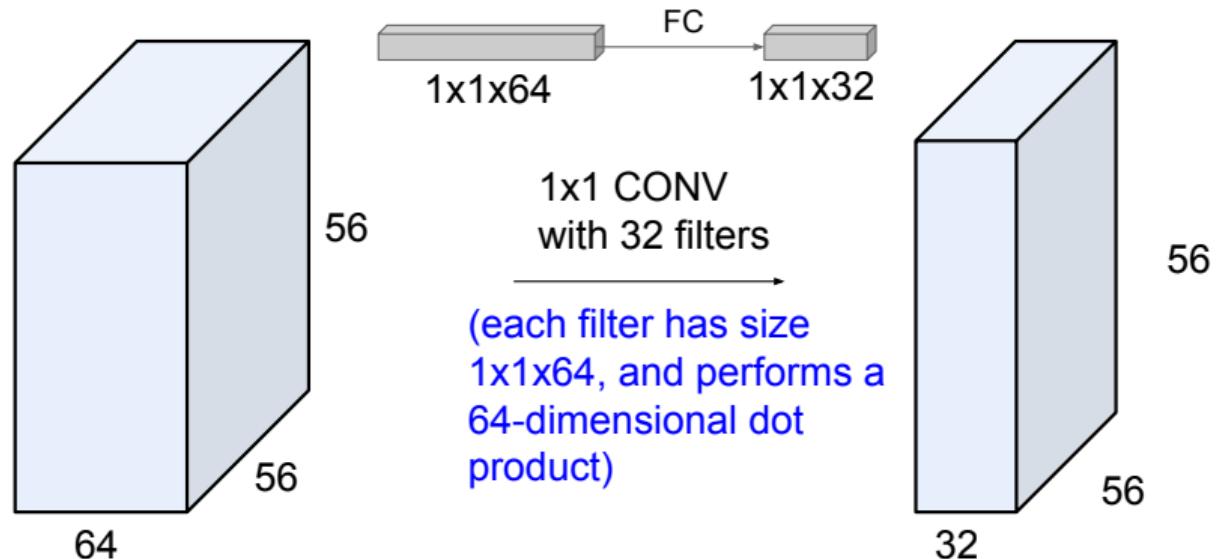
Inception blocks



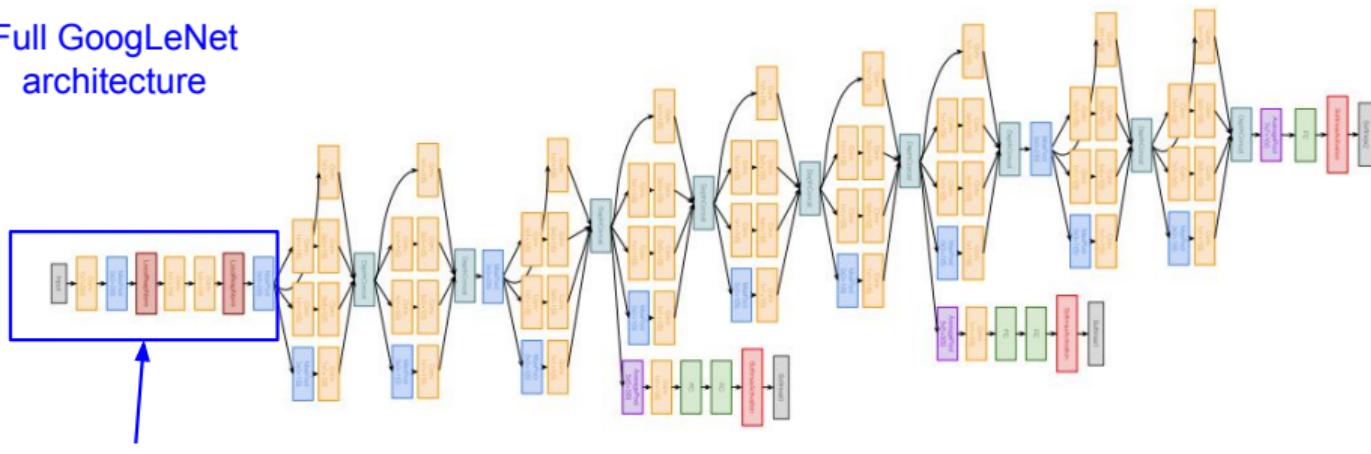
Inception blocks



1×1 convolution

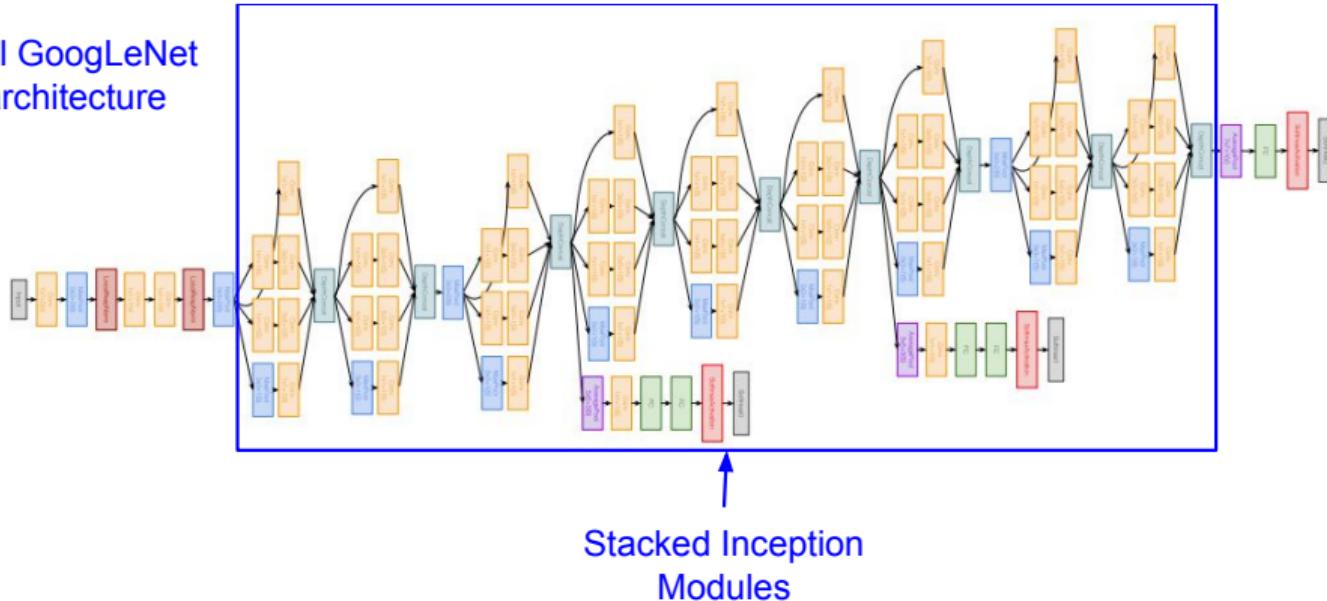


Full GoogLeNet
architecture

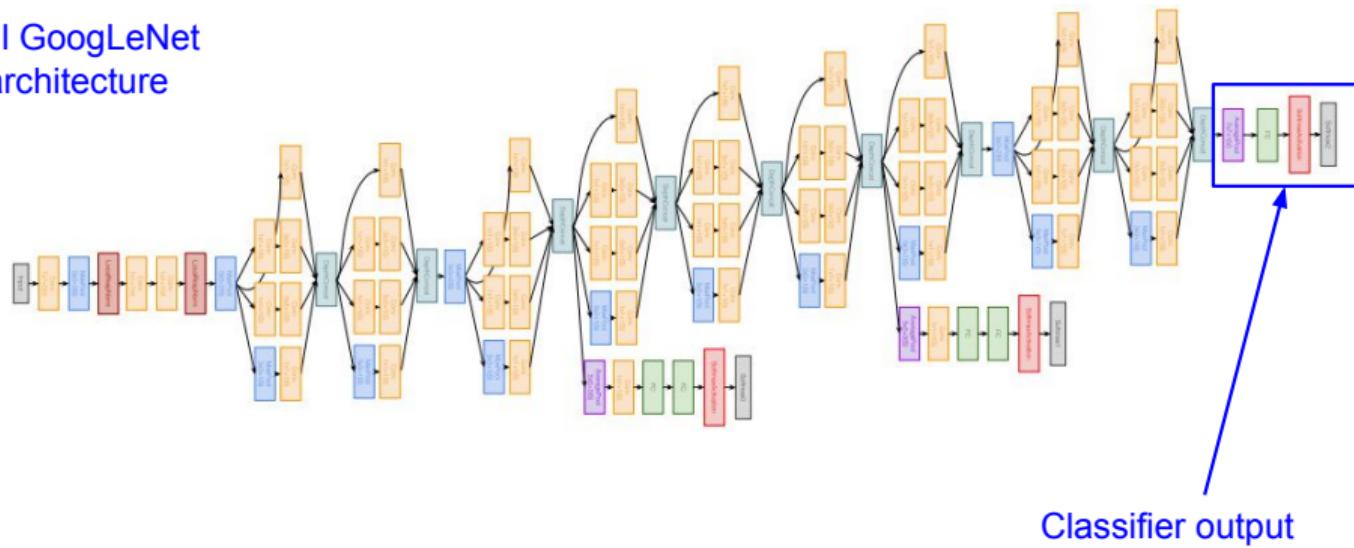


Stem Network:
Conv-Pool-
2x Conv-Pool

Full GoogLeNet
architecture

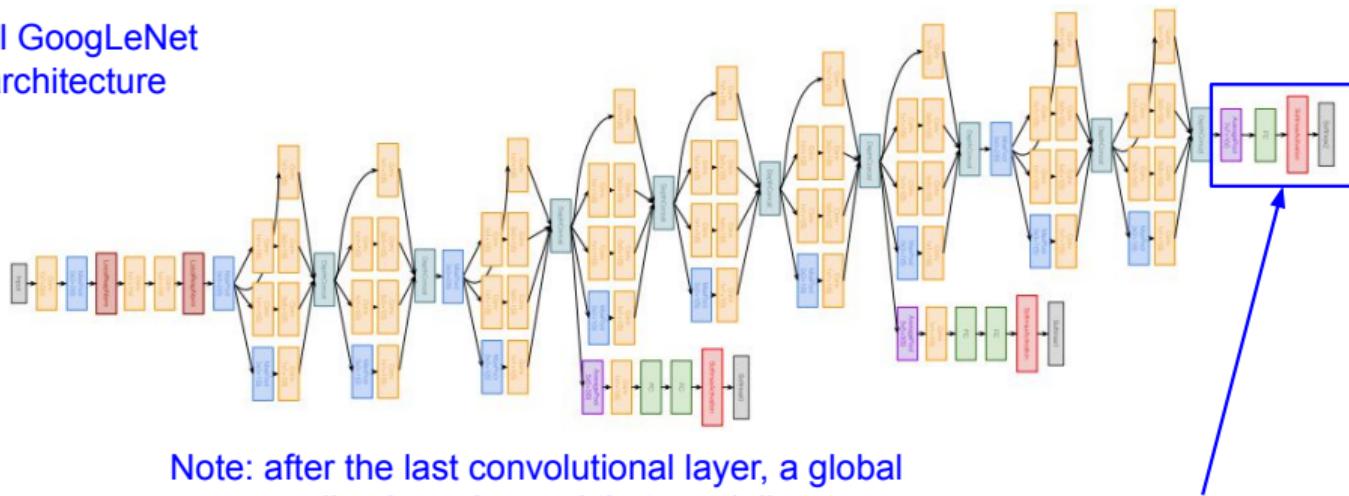


Full GoogLeNet architecture



Classifier output

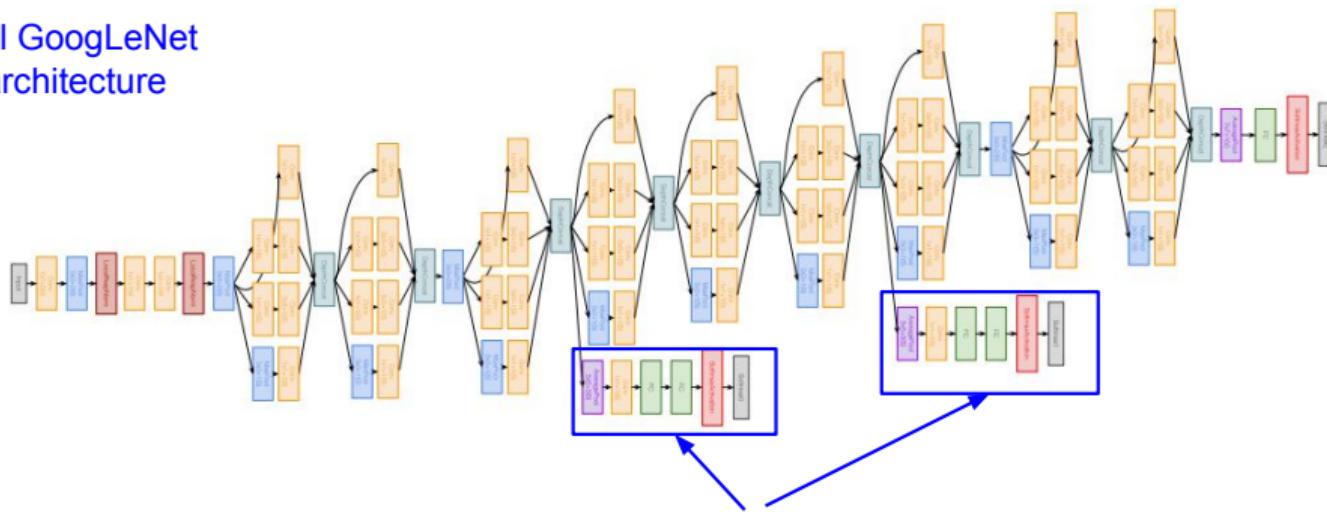
Full GoogLeNet architecture



Note: after the last convolutional layer, a global average pooling layer is used that spatially averages across each feature map, before final FC layer. No longer multiple expensive FC layers!

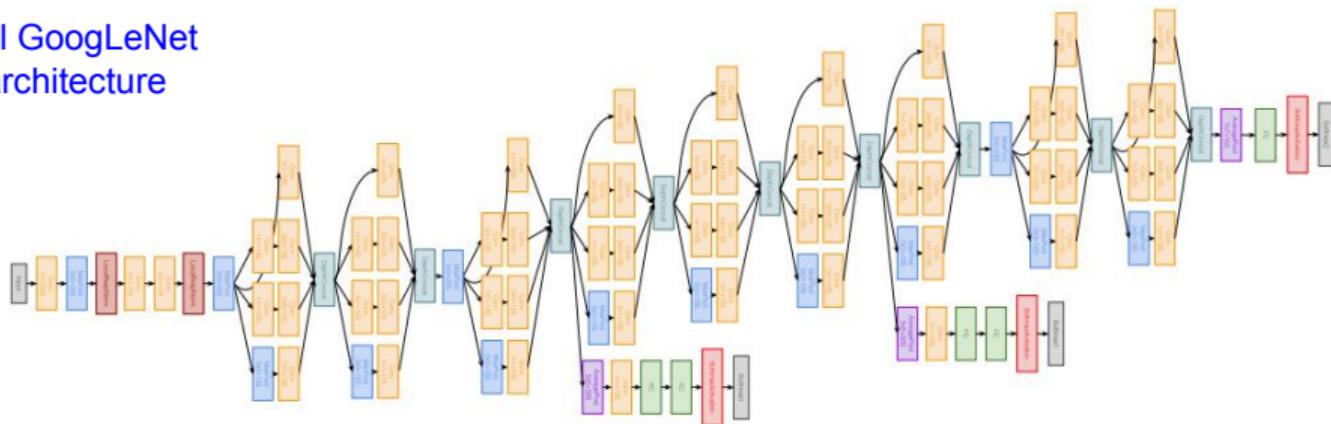
Classifier output

Full GoogLeNet architecture



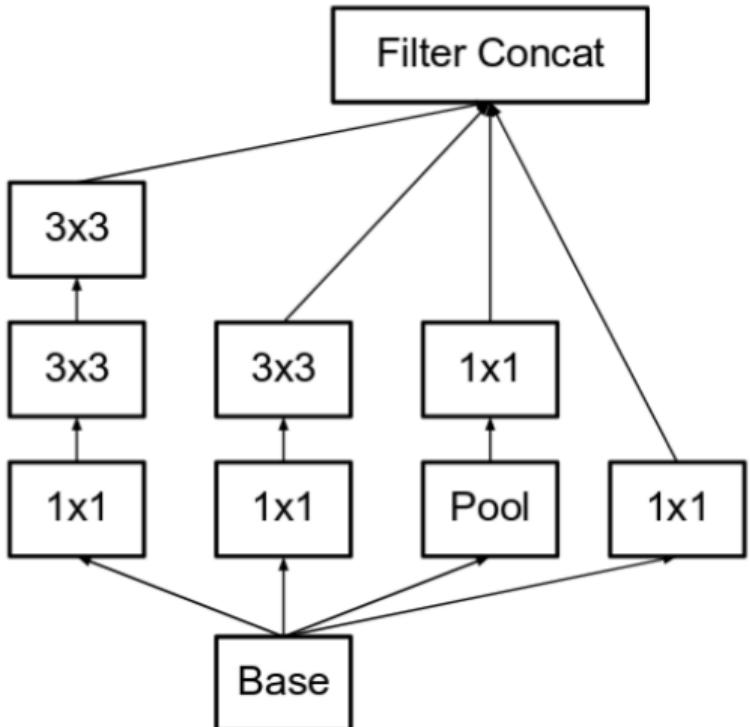
Auxiliary classification outputs to inject additional gradient at lower layers
(AvgPool-1x1Conv-FC-FC-Softmax)

Full GoogLeNet architecture

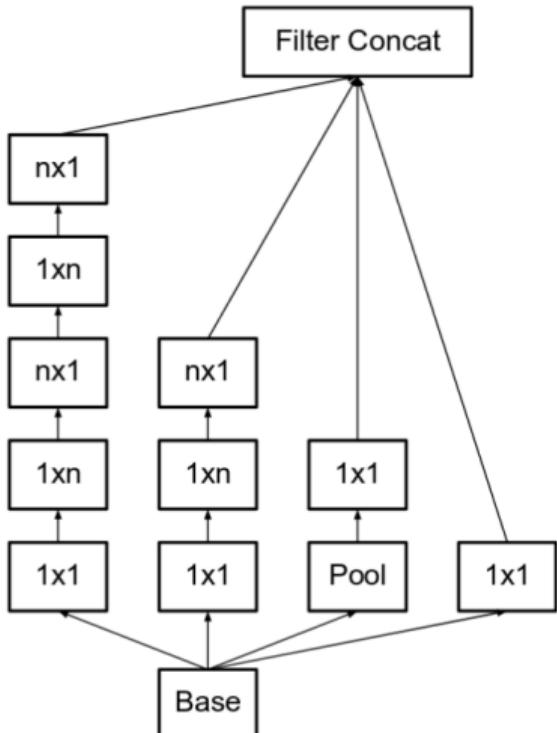


22 total layers with weights
(parallel layers count as 1 layer => 2 layers per Inception module. Don't count auxiliary output layers)

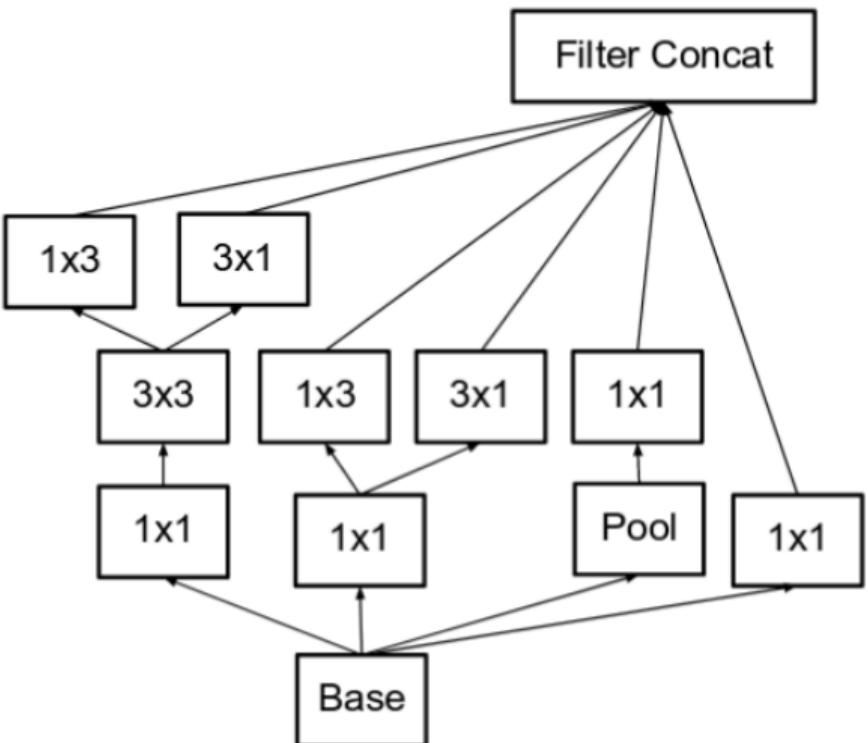
Inception v2 blocks - I



Inception v2 blocks - II



Inception v2 blocks - III



Vanishing/exploding gradients



The problem of vanishing/exploding gradients prevents training naive deep networks. So far we saw two solutions:

- VGG - train the early layers by themselves, add more later
- Inception - use auxillary training layers

Training deep networks

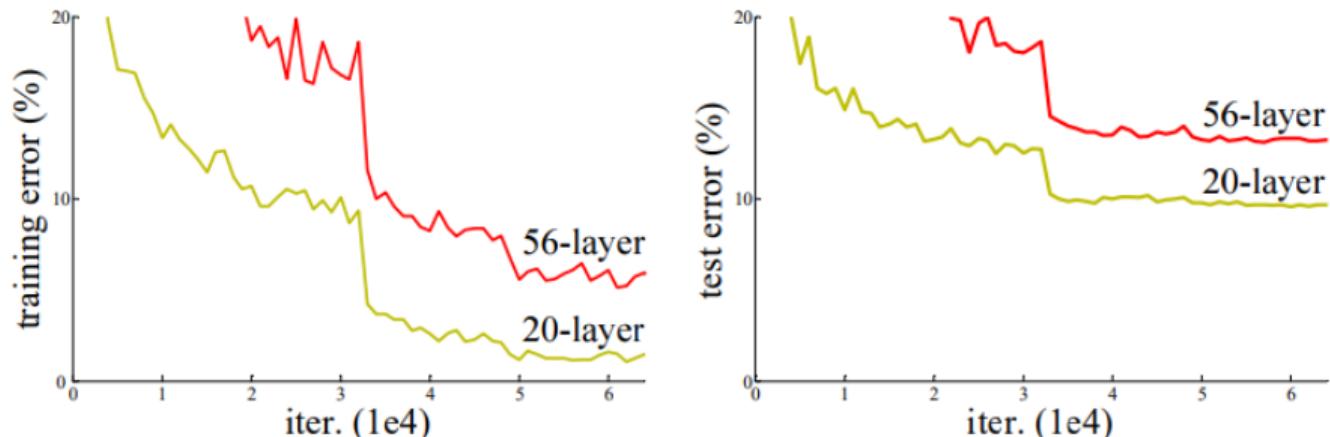
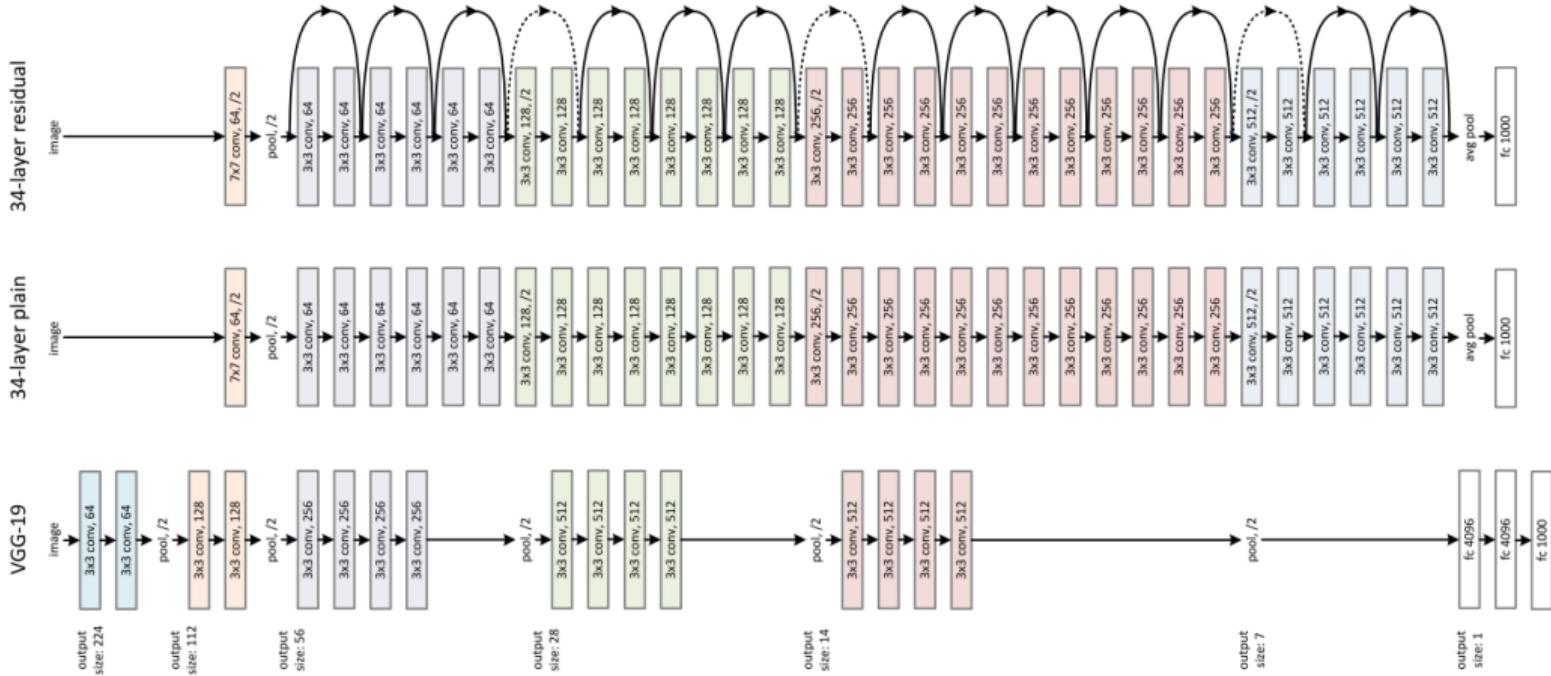
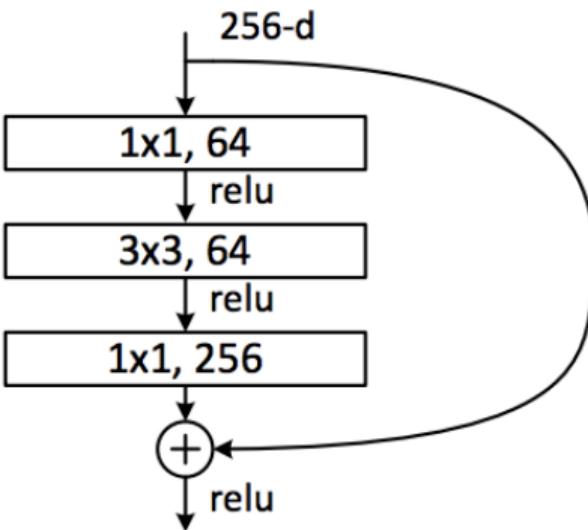
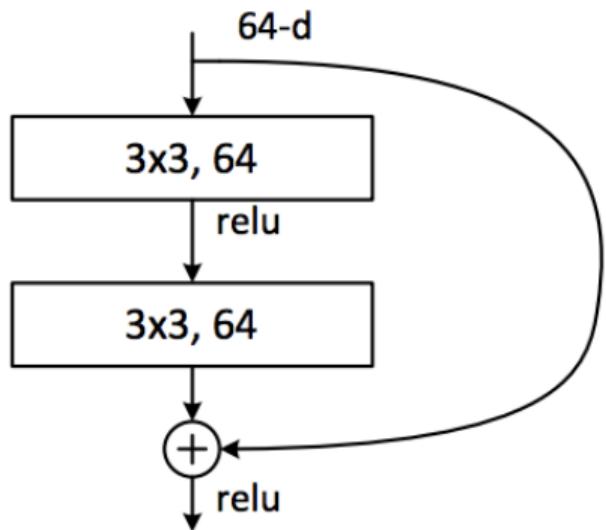


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

Introducing residual/skip connections



Residual connection



ResNets



Residual networks paper showed that a network of arbitrary depths can be trained! (With diminishing results)

⁷ Kaiming He et al. "Deep residual learning for image recognition." In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778



Residual networks paper showed that a network of arbitrary depths can be trained! (With diminishing results)

Even 6 years later ResNets are still a good choice for a backbone architecture especially when prototyping!

⁷ Kaiming He et al. "Deep residual learning for image recognition." In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778



Residual networks paper showed that a network of arbitrary depths can be trained! (With diminishing results)

Even 6 years later ResNets are still a good choice for a backbone architecture especially when prototyping!

The original paper⁷ is one of the most cited papers in CV.

⁷ Kaiming He et al. "Deep residual learning for image recognition." In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778

Resnet variants and successors



- Inception-ResNet⁸ - Combination of Inception style blocks and residual connections
- ResNet v2⁹ - Switch around the order of operations within a block
- Wide ResNets¹⁰ - Having more channels is better than going deeper
- ResNext¹¹ - Multiple pathways, similar to Inception
- Squeeze-and-ExcitationNetwork¹² - Added recalibration layer - Won ILSVRC 2017

⁸Christian Szegedy et al. "Inception-v4, inception-resnet and the impact of residual connections on learning." In: *Thirty-first AAAI conference on artificial intelligence*. 2017

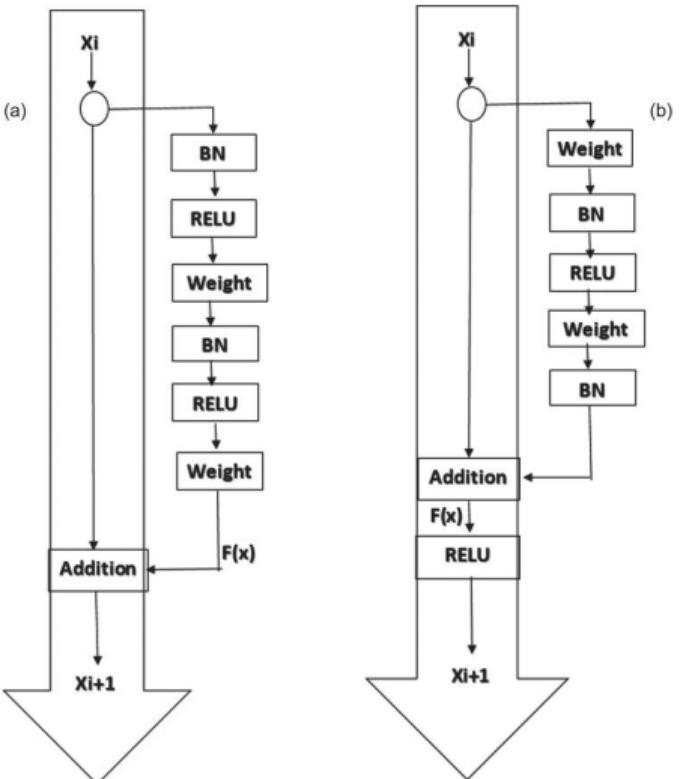
⁹Kaiming He et al. "Identity mappings in deep residual networks." In: *European conference on computer vision*. Springer. 2016, pp. 630–645

¹⁰Sergey Zagoruyko and Nikos Komodakis. "Wide residual networks." In: *arXiv preprint arXiv:1605.07146* (2016)

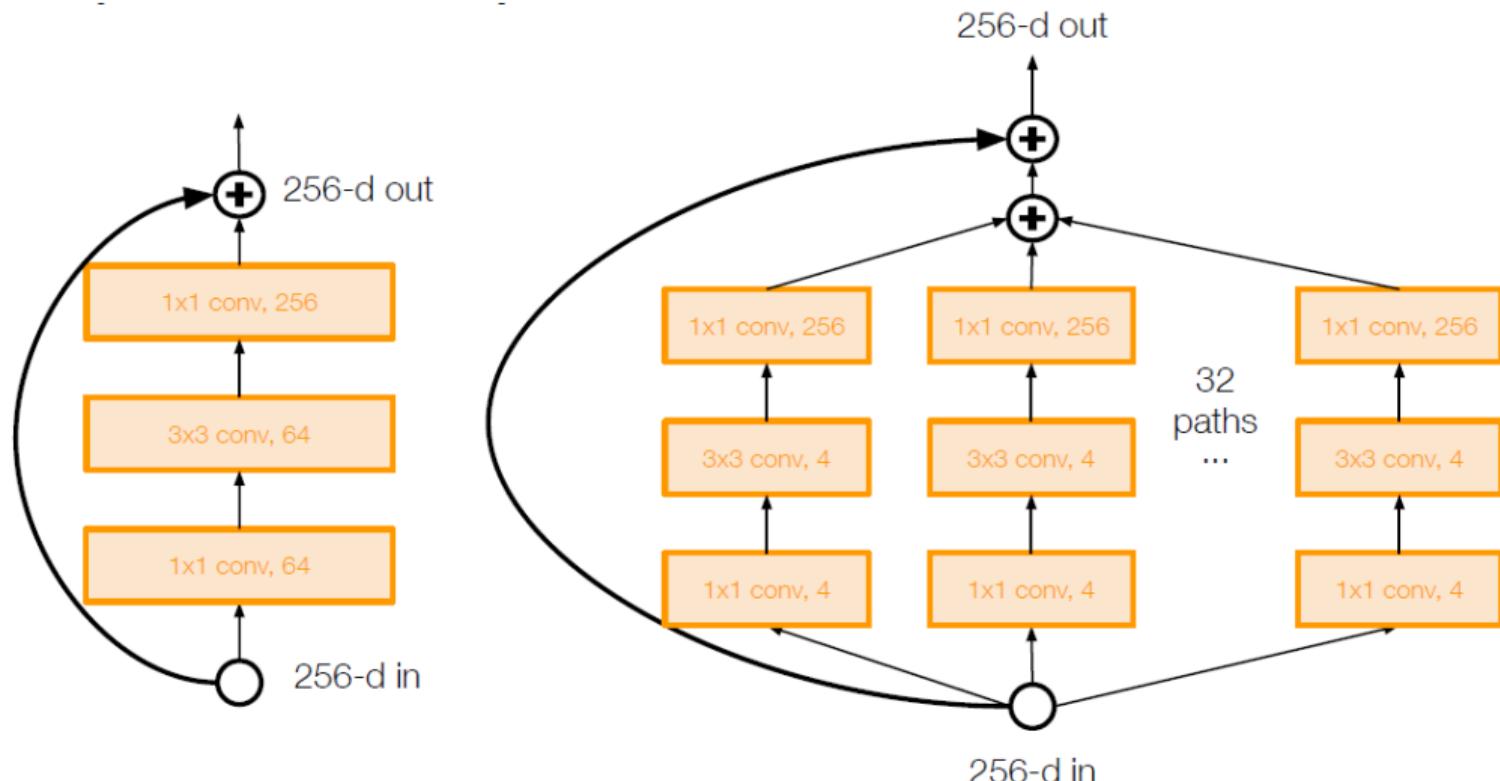
¹¹Saining Xie et al. "Aggregated residual transformations for deep neural networks." In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017, pp. 1492–1500

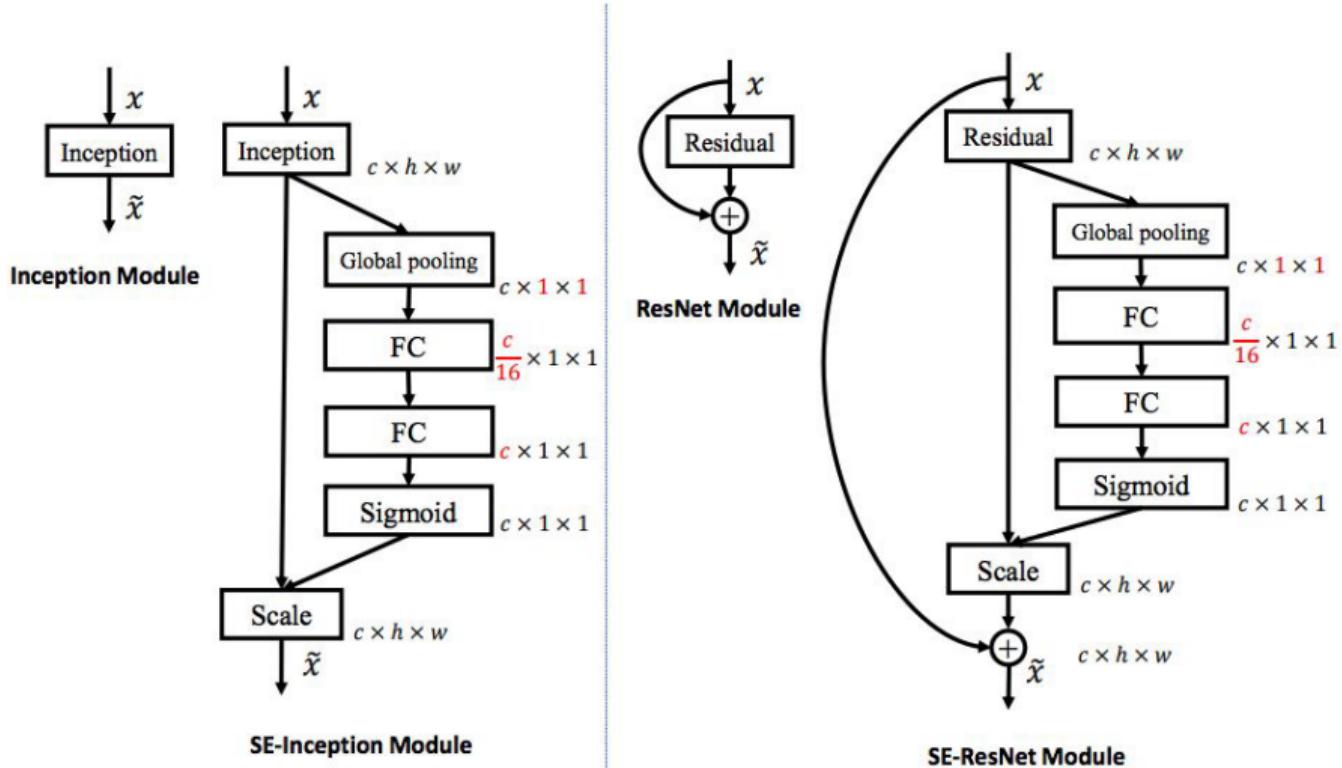
¹²Jie Hu, Li Shen, and Gang Sun. "Squeeze-and-excitation networks." In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018, pp. 7132–7141

ResNet v2

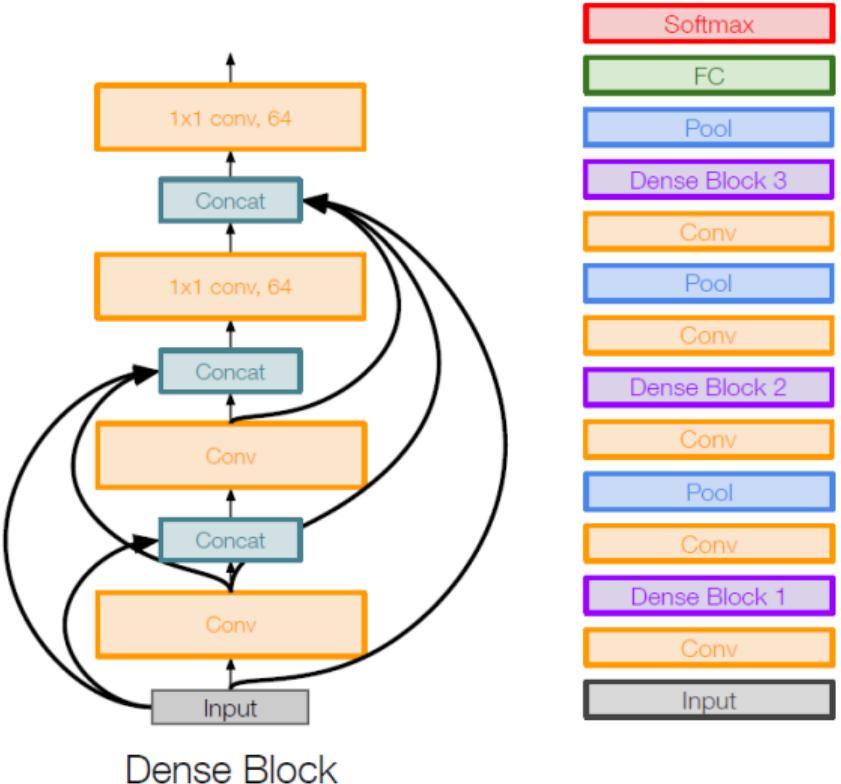


ResNext





DenseNet



Neural Architecture Search



So far the networks were usually designed by hand. It is possible to treat this as an optimization problem in the space of architectures. This is however a difficult task a very influential paper¹³ used some tricks to get it working:

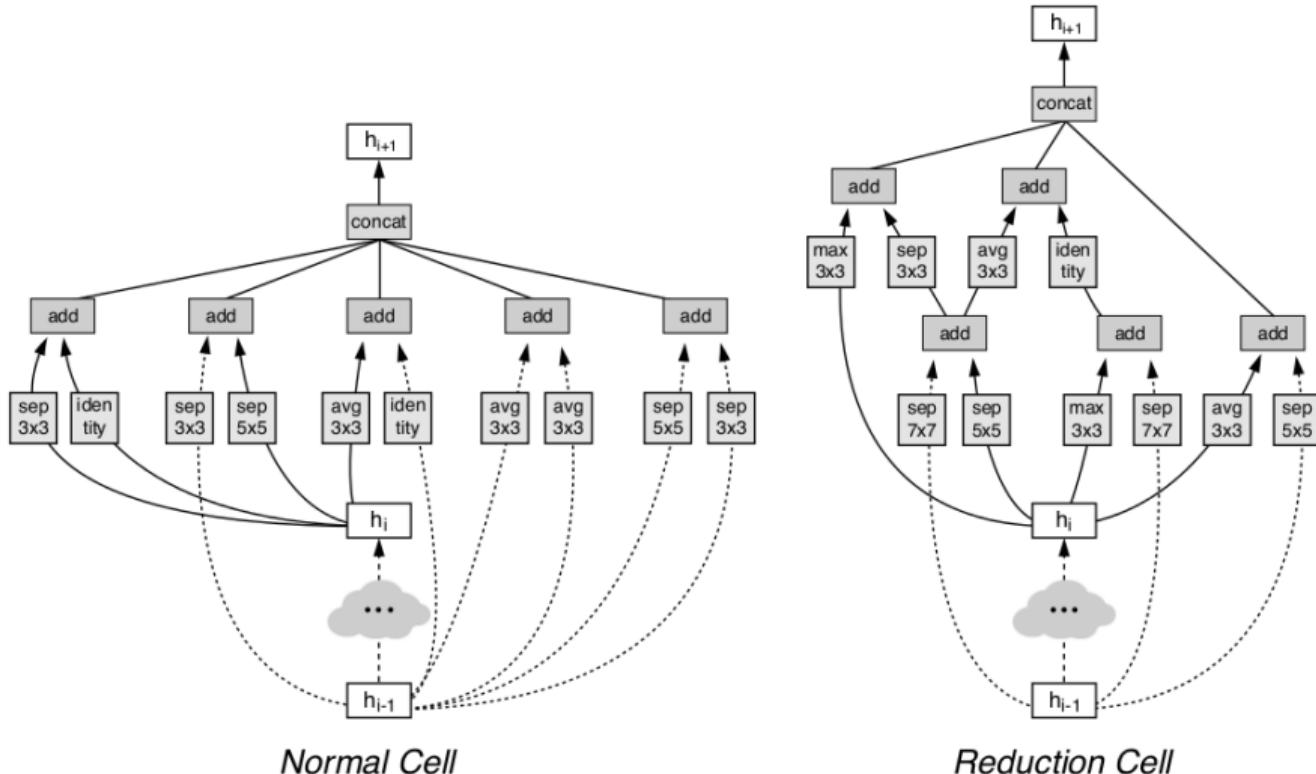
- Consider a network composed of two types of blocks and only optimize the block structure
- Optimize on smaller dataset and transfer to larger one later by stacking multiple blocks

Since then there have been some other approaches to NAS such as ENAS¹⁴ where the optimization starts with a large computational graph and tries to find the optimal subgraph.

¹³ Barret Zoph et al. "Learning transferable architectures for scalable image recognition." In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018, pp. 8697–8710

¹⁴ Hieu Pham et al. "Efficient neural architecture search via parameters sharing." In: *International Conference on Machine Learning*. PMLR. 2018, pp. 4095–4104

NASNet



Efficient architectures



So far the architectures were mostly designed to optimize for accuracy. However it is possible to optimize for efficiency as well. There are multiple approaches.

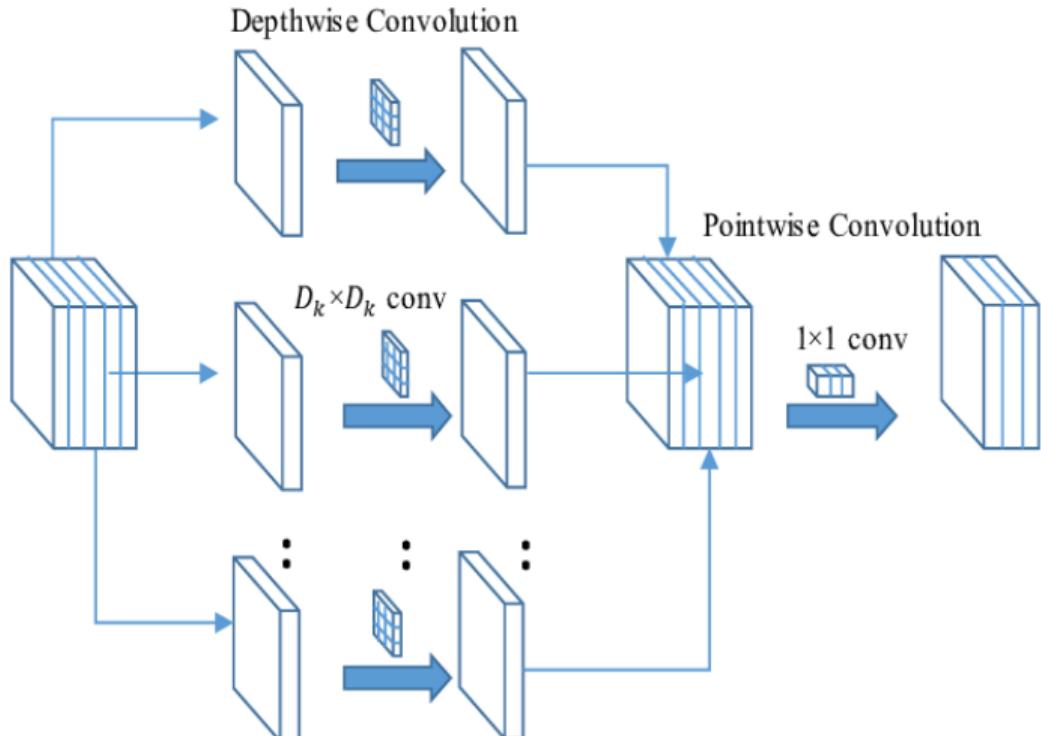
- Depthwise Separable Convolutions
 - ▶ Introduced in MobileNet¹⁵
 - ▶ Also used in Xception¹⁶ - Inception-like architecture
- ShuffleNet¹⁷ - Group pointwise convolutions with group shuffling
- NAS - MNASNet¹⁸ optimized for accuracy + computational speed on mobile devices

¹⁵ Andrew G Howard et al. "Mobilennets: Efficient convolutional neural networks for mobile vision applications." In: *arXiv preprint arXiv:1704.04861* (2017)

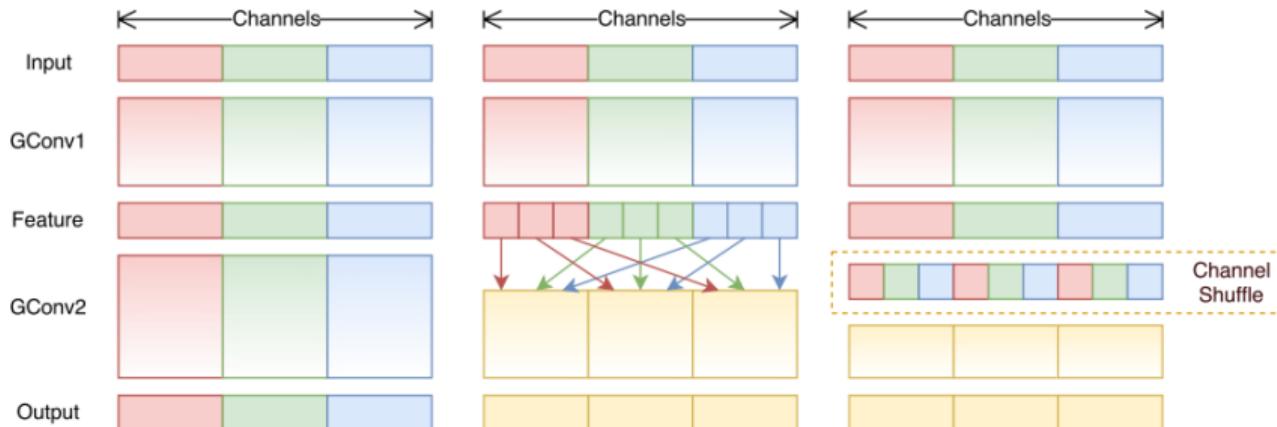
¹⁶ François Chollet. "Xception: Deep learning with depthwise separable convolutions." In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017, pp. 1251–1258

¹⁷ Xiangyu Zhang et al. "Shufflenet: An extremely efficient convolutional neural network for mobile devices." In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018, pp. 6848–6856

¹⁸ Mingxing Tan et al. "Mnasnet: Platform-aware neural architecture search for mobile." In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2019, pp. 2820–2828



ShuffleNet





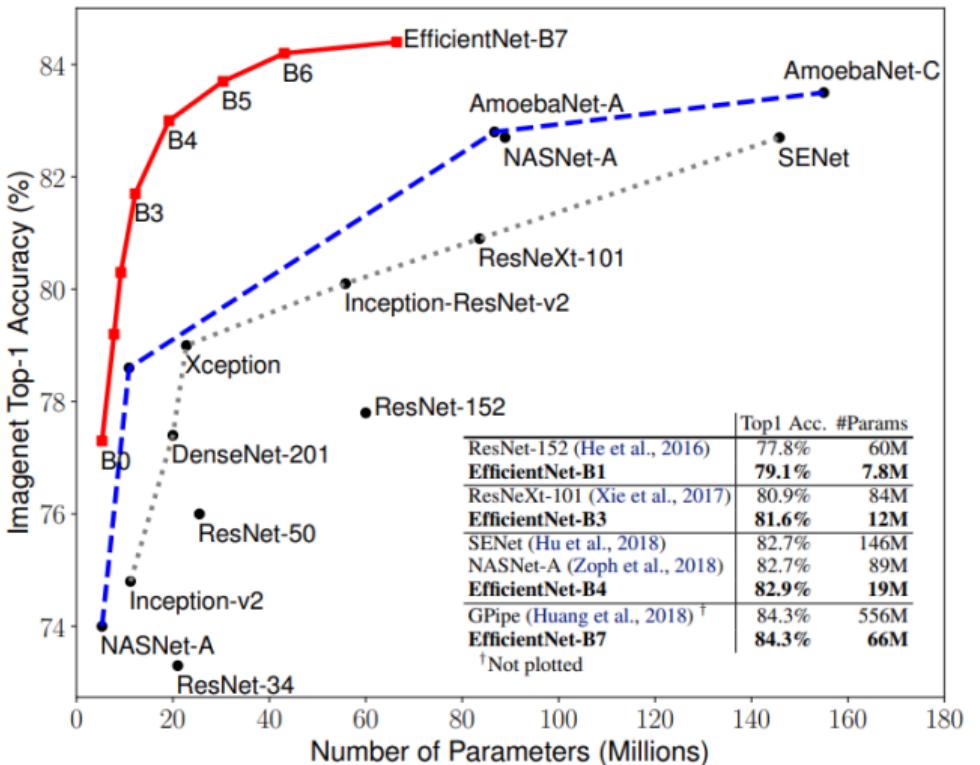
Most networks can be scaled roughly in three dimensions:

- Depth - number of layers
- Width - number of feature channels
- Input size - dimensions of input image

In the EfficientNet¹⁹ paper authors propose a mechanism that scales the network along the three dimensions to ensure better performance. They also use NAS to find optimal architecture and propose several models EfficientNet-B0 to EfficientNet-B7.

¹⁹ Mingxing Tan and Quoc Le. "Efficientnet: Rethinking model scaling for convolutional neural networks." In: *International Conference on Machine Learning*. PMLR. 2019, pp. 6105–6114

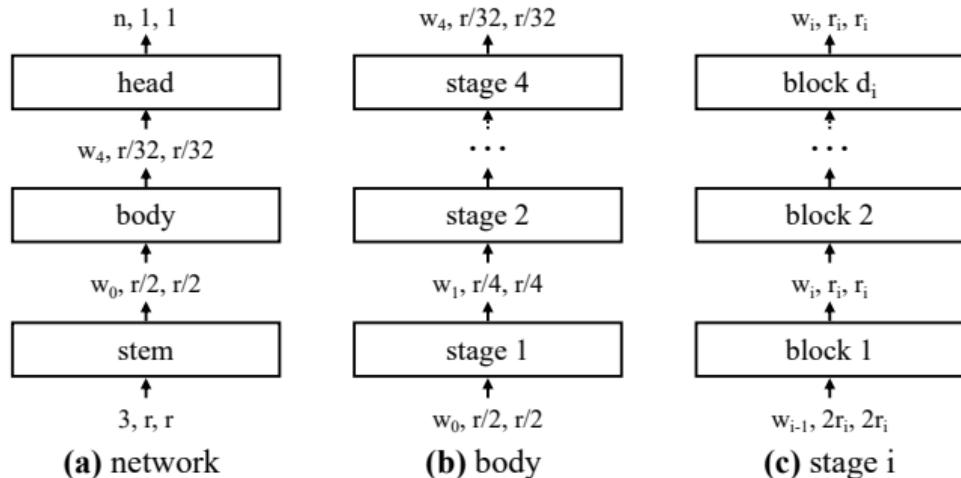
EfficientNets



Network Design Spaces

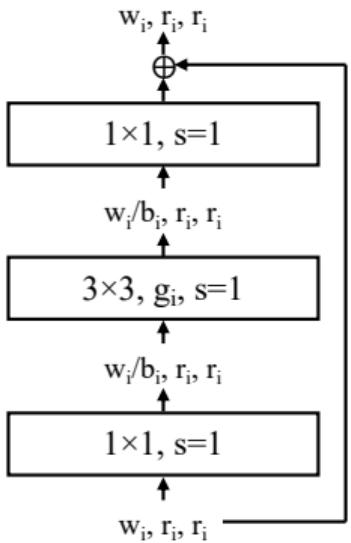


To investigate good design principles for the neural network architecture it is possible to construct a network design space²⁰ (a set of possible architectures) and manually finding the principles that improve accuracy for the whole set.

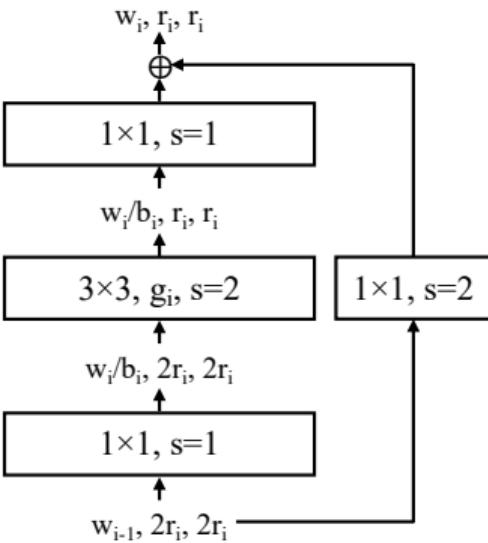


²⁰Ilija Radosavovic et al. "Designing network design spaces." In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020, pp. 10428–10436

Reducing the design space



(a) X block, $s=1$



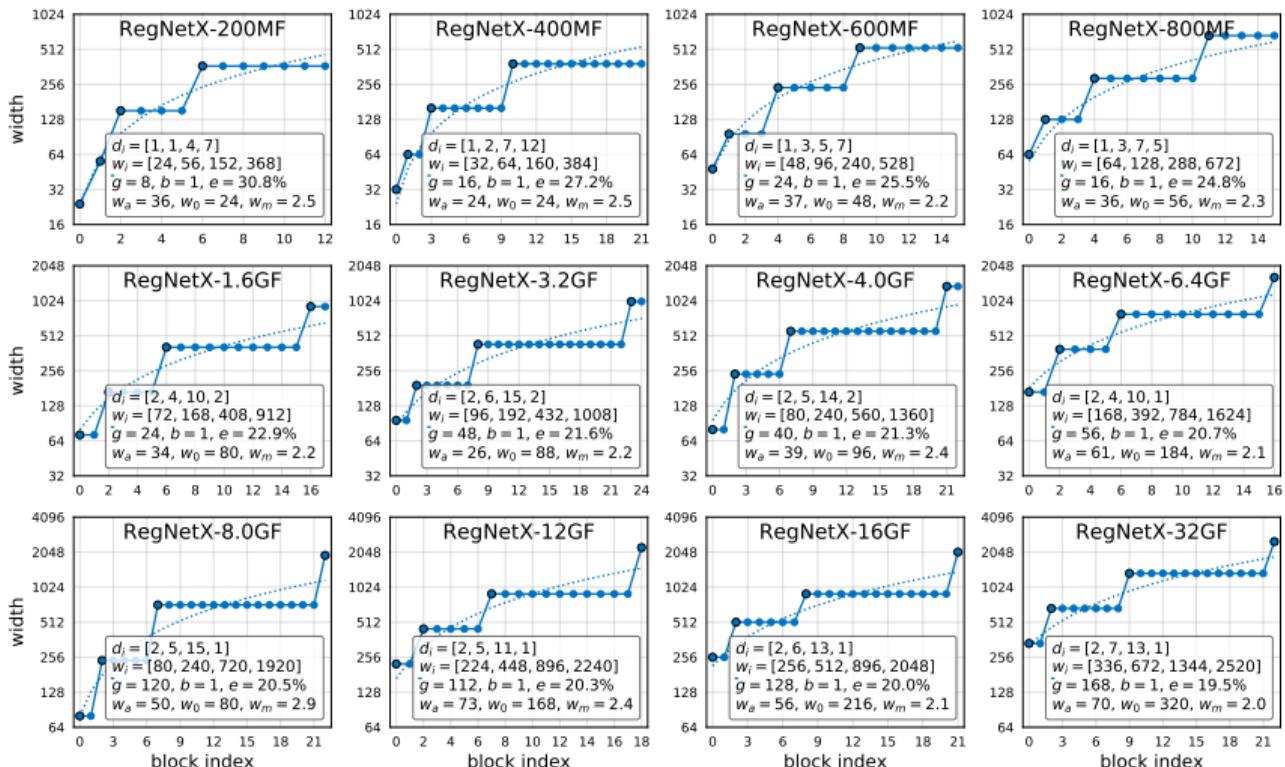
(b) X block, $s=2$

The design space is iteratively reduced using multiple rules (for example shared groups in blocks can be the same for every block).

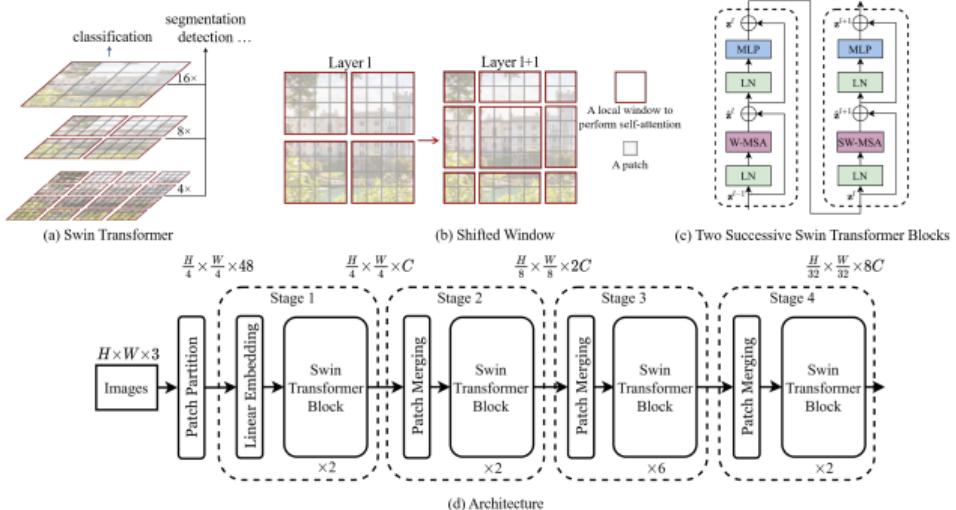


	restriction	dim.	combinations	total
AnyNet A	none	16	$(16 \cdot 128 \cdot 3 \cdot 6)^4$	$\sim 1.8 \cdot 10^{18}$
AnyNet B	$+ b_{i+1} = b_i$	13	$(16 \cdot 128 \cdot 6)^4 \cdot 3$	$\sim 6.8 \cdot 10^{16}$
AnyNet C	$+ g_{i+1} = g_i$	10	$(16 \cdot 128)^4 \cdot 3 \cdot 6$	$\sim 3.2 \cdot 10^{14}$
AnyNet D	$+ w_{i+1} \geq w_i$	10	$(16 \cdot 128)^4 \cdot 3 \cdot 6 / (4!)$	$\sim 1.3 \cdot 10^{13}$
AnyNet E	$+ d_{i+1} \geq d_i$	10	$(16 \cdot 128)^4 \cdot 3 \cdot 6 / (4!)^2$	$\sim 5.5 \cdot 10^{11}$
RegNet	quantized linear	6	$\sim 64^4 \cdot 6 \cdot 3$	$\sim 3.0 \cdot 10^8$

Best RegNetX models



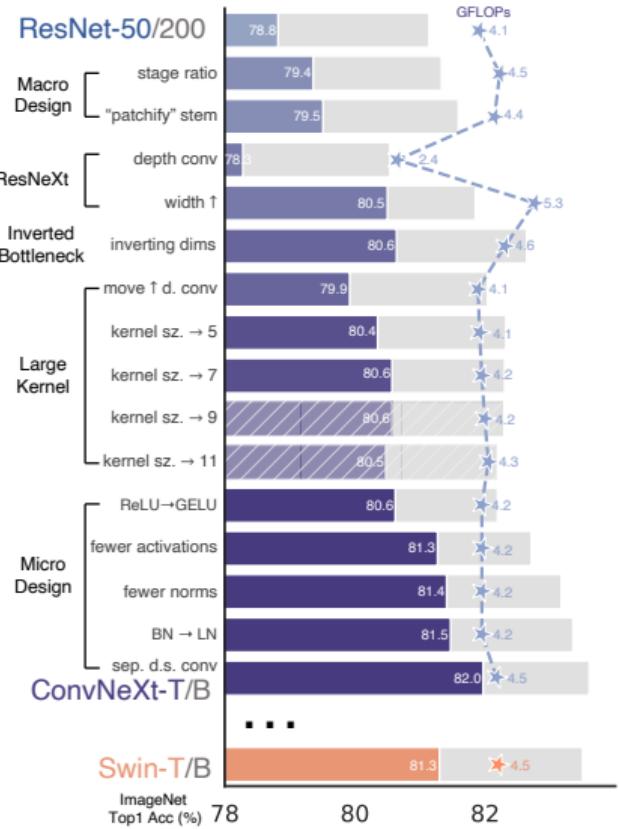
Transformers



In 2021-2022 there was a trend to move away from CNNs and instead use Transformers even for CV Tasks! In the end it turns out that Transformers are probably not better than CNNs. We will discuss this in greater detail next week.



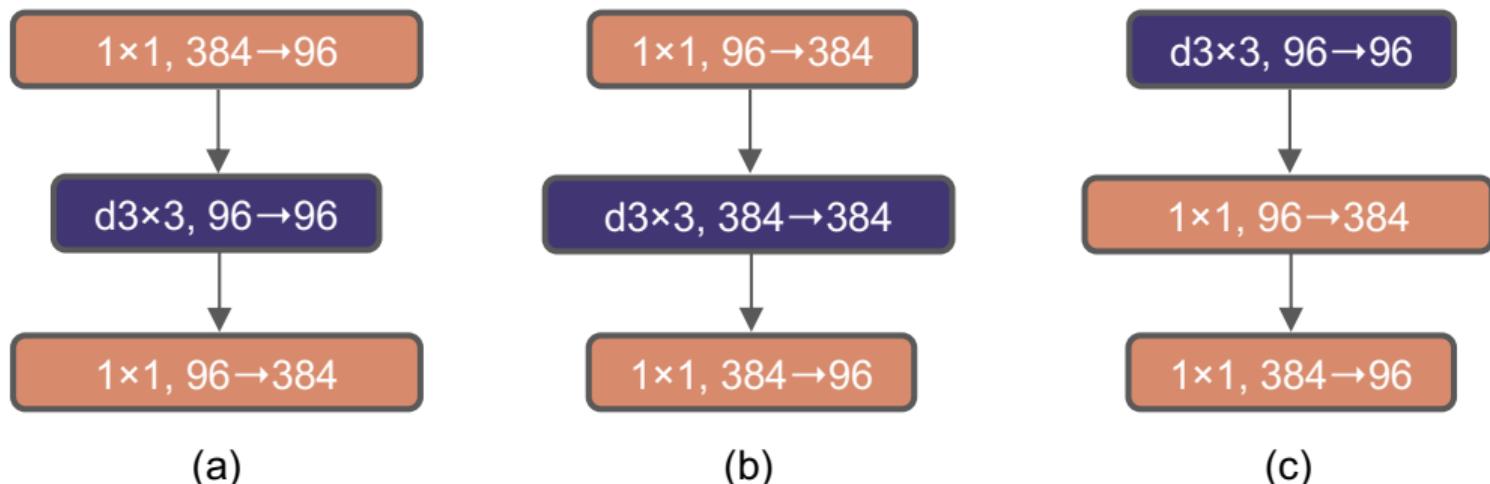
ConvNext



ConvNeXts²¹ were designed as a response to the rise of Transformers in CV. With several changes to the architecture of ResNets they were able to achieve SOTA results.

²¹Zhuang Liu et al. "A convnet for the 2020s." In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 11976–11986

ConvNext - inversion

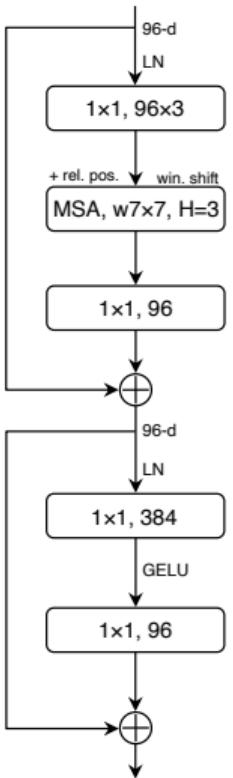


Block modifications and resulted specifications. (a) is a ResNeXt block; in (b) we create an inverted bottleneck block and in (c) the position of the spatial depthwise conv layer is moved up.

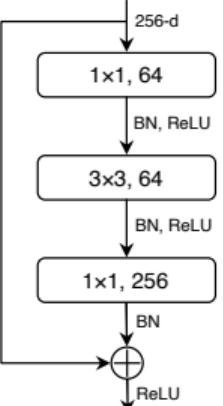
ConvNext - block comparison



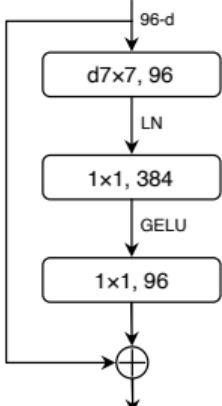
Swin Transformer Block



ResNet Block



ConvNeXt Block



Transfer Learning



Training deep neural nets requires lots of annotated data! If our dataset is small we might have a problem.

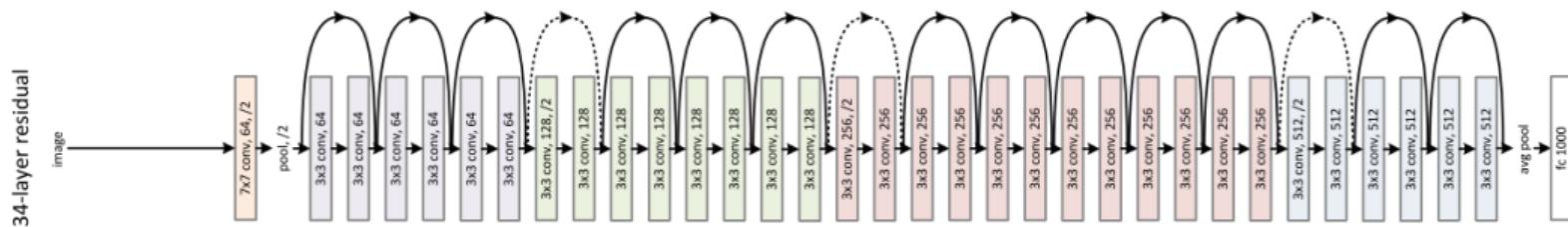
Transfer Learning



Training deep neural nets requires lots of annotated data! If our dataset is small we might have a problem.

Solution: Pretrain the network on bigger dataset (of similar data) and fine-tune on our dataset.

Transfer Learning - Freezing Layers



If your own dataset is very small you can freeze the convolutional layers and keep train only the final FC layer.

Transfer Learning as Initialization



In practice we use pre-trained weights even if our dataset is large! This removes issues with parameter initialization. It often works even if data is from different domain!

Most common architectures can be initialized in both TF and Pytorch with a single line of code. Check the model ZOO's:

- <https://pytorch.org/vision/stable/models.html>
- <https://keras.io/api/applications/>

How to choose the best architecture



Some tips for selecting an architecture for your project:

- Do not create your own - use existing architectures
- If possible use models directly from frameworks model zoo
- If you need/want to make modifications - find code on GitHub
- ResNets are usually a good baseline - easy to compare to other approaches
- Other than that go for what achieves best ImageNet acc for your computational constraints
- Do not use VGG or AlexNet
- Be mindful of the image input size, network depth (next slide)

Selecting model



- For early prototypes use small models (ResNet18) with small input size
- If everything works try a larger model and input size
- Image input size should not be larger than receptive field of network
- You can sometimes increase/decrease receptive field by modifying some elements of the network (e.g. stem in ResNets)
- If you operate under some computational-constraints (e.g. realtime video processing on edge device) keep that in mind and go for smaller nets with smaller input sizes!

ResNet - receptive fields



layer	resnet18	resnet34	resnet50	resnet101
conv1	7	7	7	7
maxpool	11	11	11	11
layer1	43	59	35	35
layer2	99	179	91	91
layer3	211	547	267	811
layer4	435	899	427	971

Due to use of different blocks ResNet50 has smaller receptive field than ResNet18!