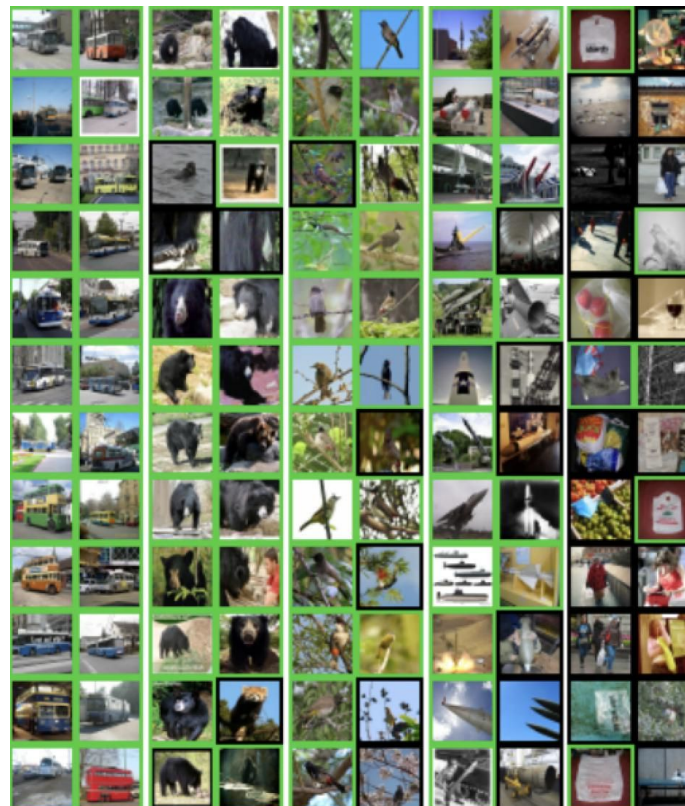# Deep Learning

Lecture 6

# Recap

- Attention
- Applications
- Types of attention
- Transformer
  - Positional encoding
  - Self-attention
  - Multi-head attention
- BERT model (MLM)
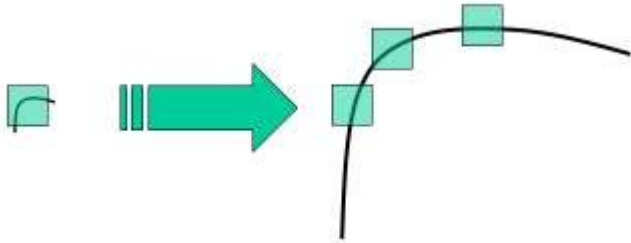
# ImageNet
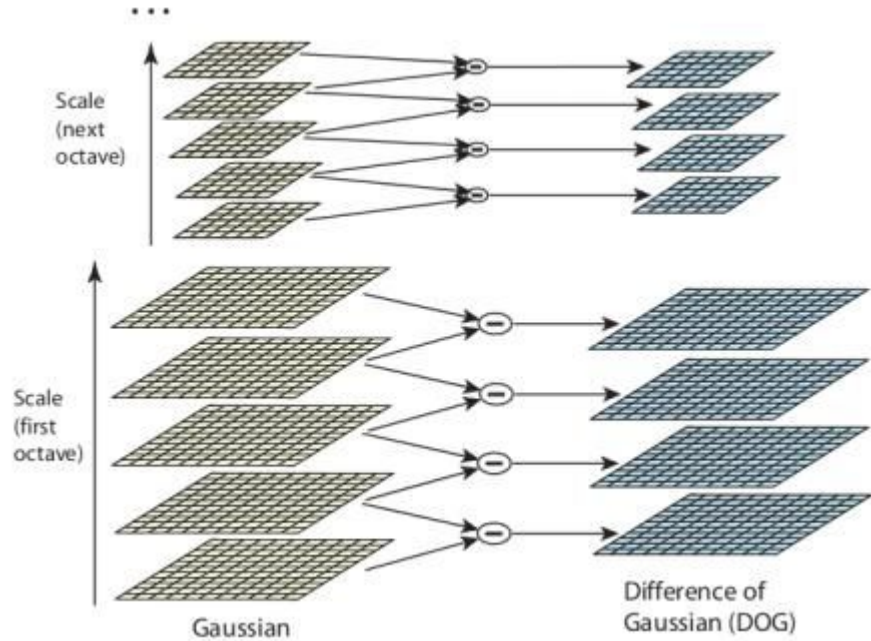
IM**A**GENET
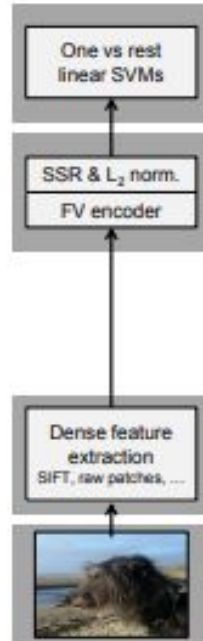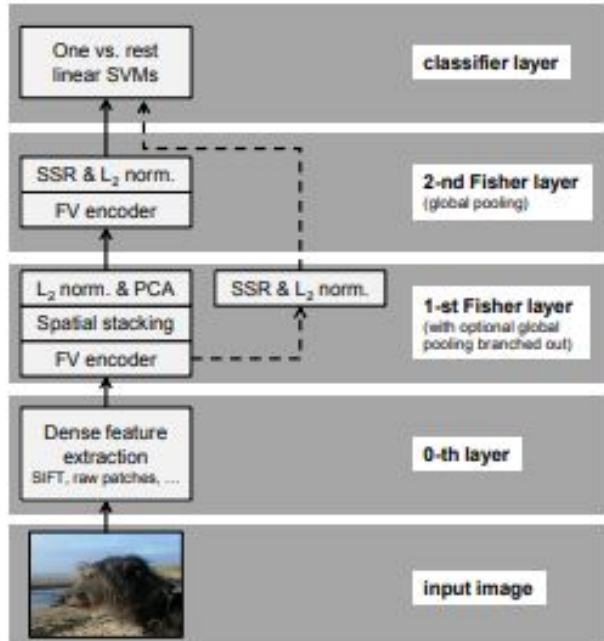


Train/Val/Test: 1.2 M / 50k / 100k images

Classes: 1000

# SIFT vectors

We want to create scale invariant feature extractor



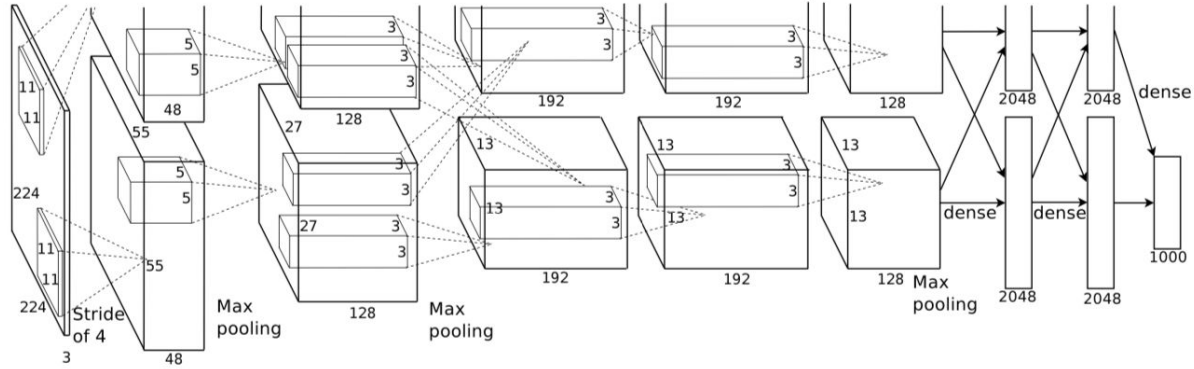Gaussian

Difference of Gaussian (DOG)

# Classical solutions



Before deep learning era solution for classification task has the following form:

- Manual feature extractor (e.g. SIFT)
- When some encoder to be able classify images of arbitrary size (Fisher vectors, codebooks)

# AlexNet



- Replaced tanh with ReLU (x6 speedup)
- Dropout + Augmentations
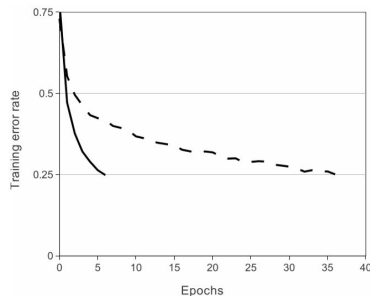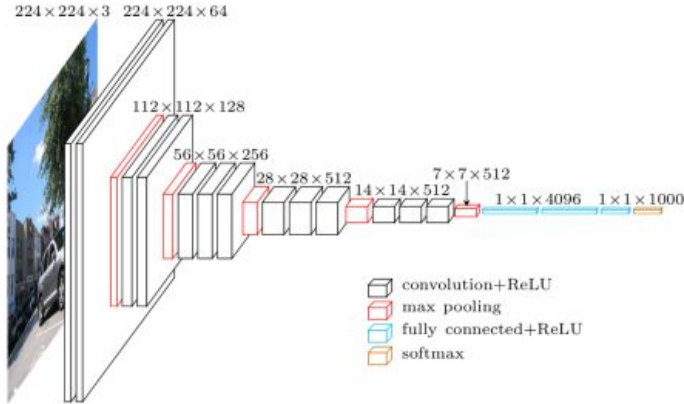- 5 conv layers (11x11,5x5,3x3,3x3,3x3)

# AlexNet



Figure 1: A four-layer convolutional neural network with ReLUs (**solid line**) reaches a 25% training error rate on CIFAR-10 six times faster than an equivalent network with tanh neurons (**dashed line**). The learning rates for each network were chosen independently to make training as fast as possible. No regularization of any kind was employed. The magnitude of the effect demonstrated here varies with network architecture, but networks with ReLUs consistently learn several times faster than equivalents with saturating neurons.
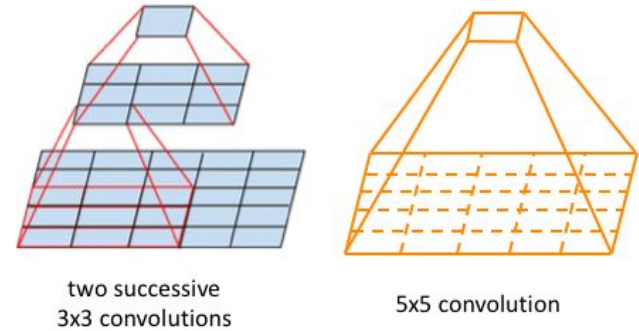
Local normalization (analogue to batch normalization)

$$b_{x,y}^i = a_{x,y}^i / \left( k + \alpha \sum_{j=\max(0,i-n/2)}^{\min(N-1,i+n/2)} (a_{x,y}^j)^2 \right)^\beta$$

# VGG



224×224×3  224×224×64
112×112×128
56×56×256
28×28×512
14×14×512
7×7×512
1×1×4096  1×1×1000

convolution+ReLU
max pooling
fully connected+ReLU
softmax

two successive
3x3 convolutions

5x5 convolution

11-19 conv. layers
All conv. filters are 3x3 size (cascade of kernels)
Stagewise training

5x5 conv is equal to two 3x3 conv
(in terms of receptive field)

# VGG

Initial architecture
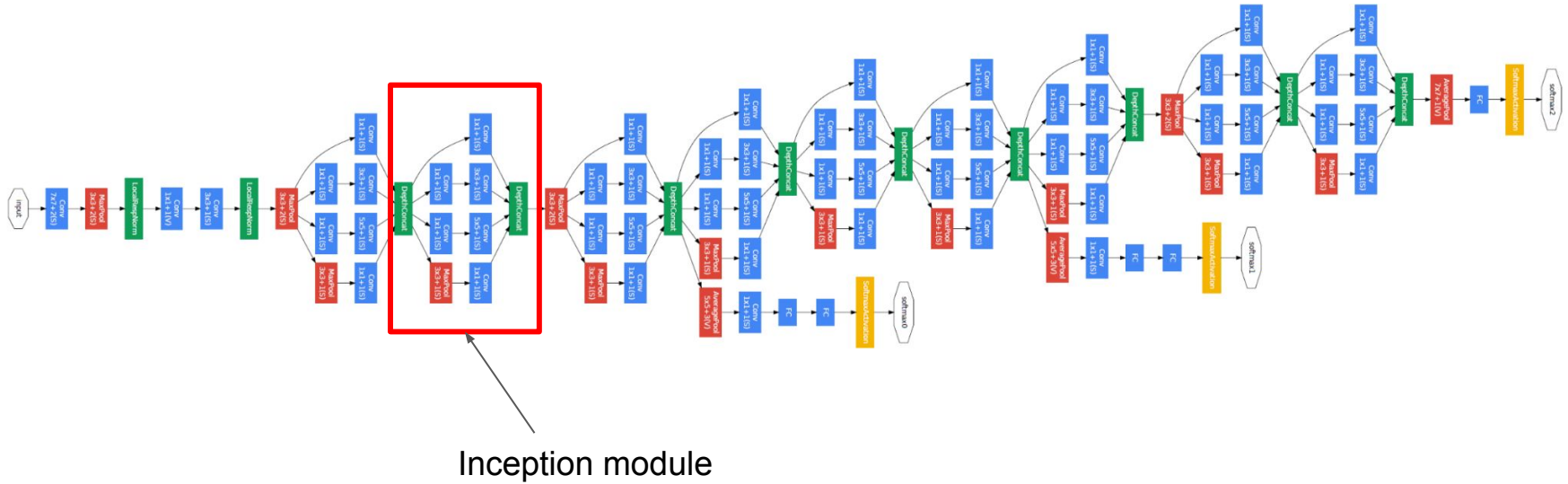
The final architecture

Table 1: **ConvNet configurations** (shown in columns). The depth of the configurations increases from the left (A) to the right (E), as more layers are added (the added layers are shown in bold). The convolutional layer parameters are denoted as "conv⟨receptive field size⟩-⟨number of channels⟩". The ReLU activation function is not shown for brevity.

| ConvNet Configuration | | | | | |
|---|---|---|---|---|---|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| input (224 × 224 RGB image) | | | | | |
| conv3-64 | conv3-64 | conv3-64 | conv3-64 | conv3-64 | conv3-64 |
| | LRN | **conv3-64** | conv3-64 | conv3-64 | conv3-64 |
| maxpool | | | | | |
| conv3-128 | conv3-128 | conv3-128 | conv3-128 | conv3-128 | conv3-128 |
| | | **conv3-128** | conv3-128 | conv3-128 | conv3-128 |
| maxpool | | | | | |
| conv3-256 | conv3-256 | conv3-256 | conv3-256 | conv3-256 | conv3-256 |
| conv3-256 | conv3-256 | conv3-256 | conv3-256 | conv3-256 | conv3-256 |
| | | | **conv1-256** | **conv3-256** | conv3-256 |
| | | | | | **conv3-256** |
| maxpool | | | | | |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
| | | | **conv1-512** | **conv3-512** | conv3-512 |
| | | | | | **conv3-512** |
| maxpool | | | | | |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
| | | | **conv1-512** | **conv3-512** | conv3-512 |
| | | | | | **conv3-512** |
| maxpool | | | | | |
| FC-4096 | | | | | |
| FC-4096 | | | | | |
| FC-1000 | | | | | |
| soft-max | | | | | |

# Inception | GoogLeNet
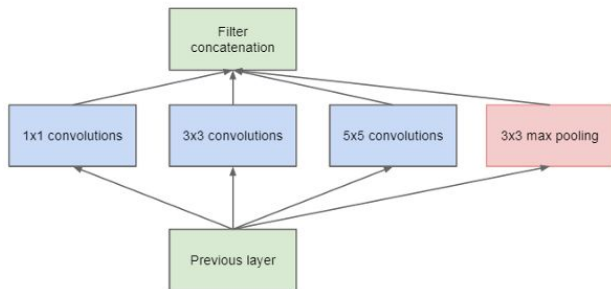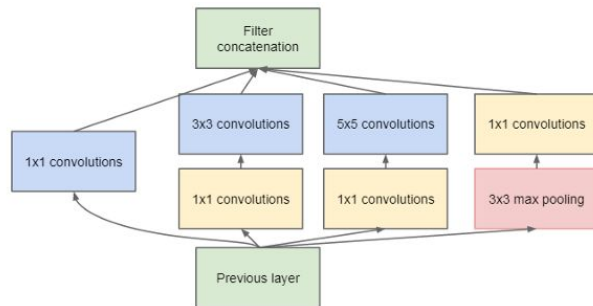


Inception module

- 22 layer
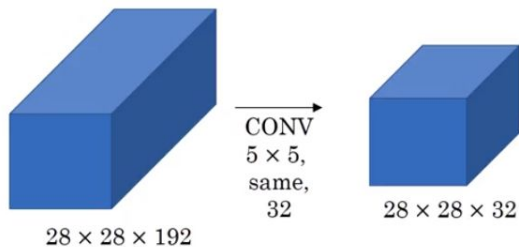- Additional outputs for classification
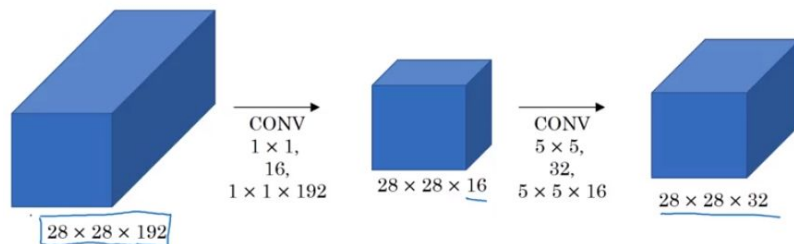
# Inception module



(a) Inception module, naïve version

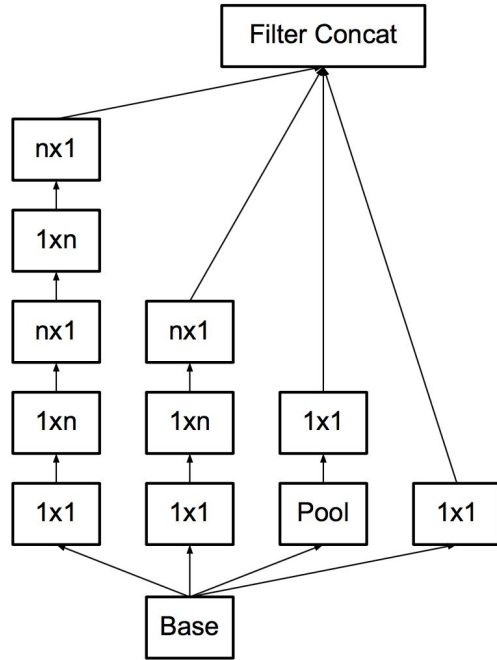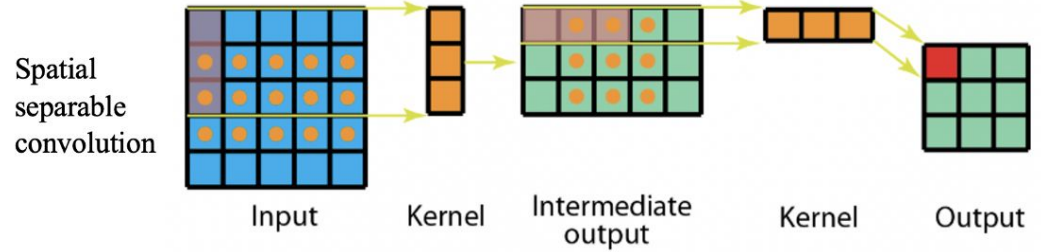(b) Inception module with dimension reductions

$\approx 120M$ calculations

$\approx 12.4M$ calculations

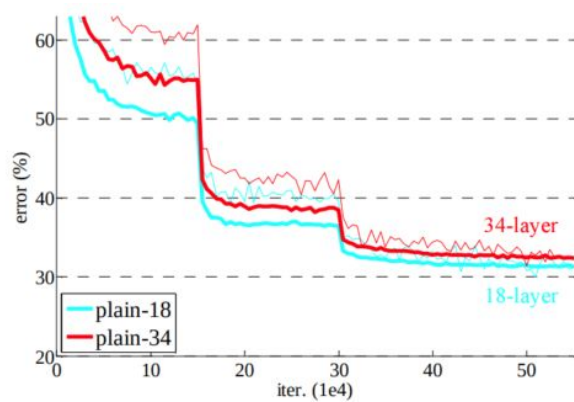ten times less calculations!

# Inception v2, v3

Filter Concat

nx1
1xn
nx1
1xn
nx1
1xn
1x1
1x1
1x1
Pool
1x1

Base

How to represent 3x3 convolution by composition of two one dimensional convolutions?

Spatial separable convolution

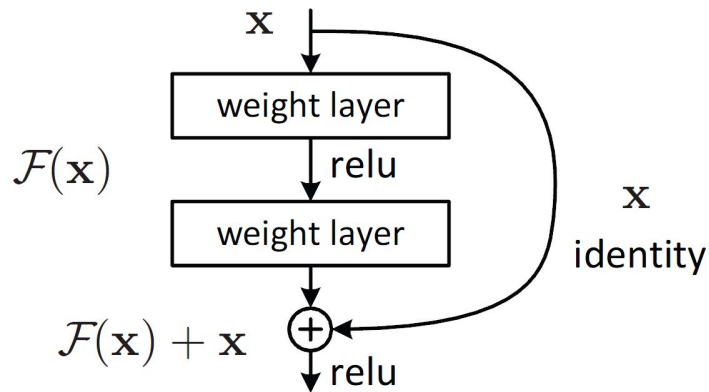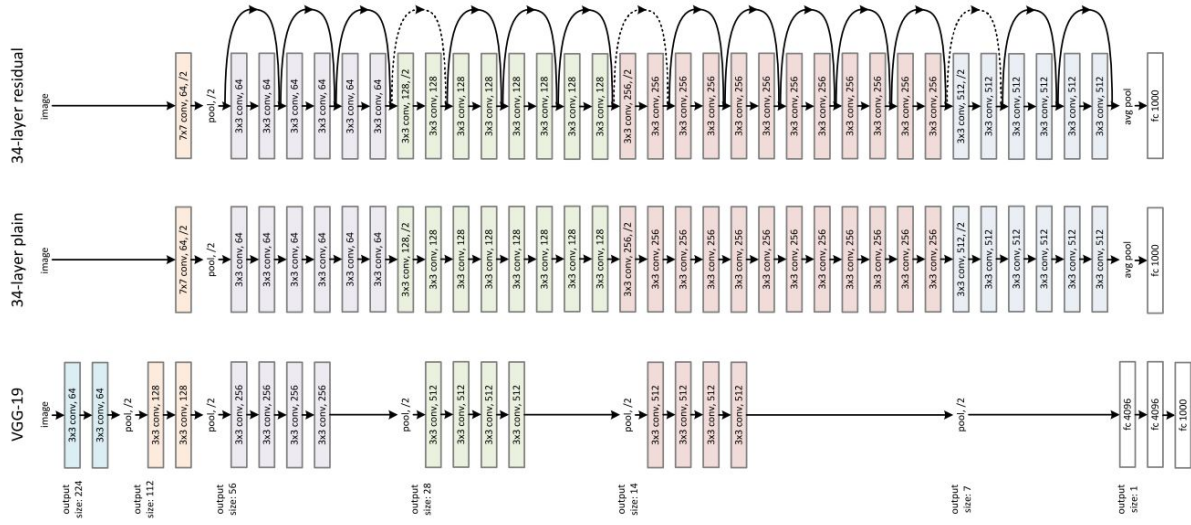Input | Kernel | Intermediate output | Kernel | Output
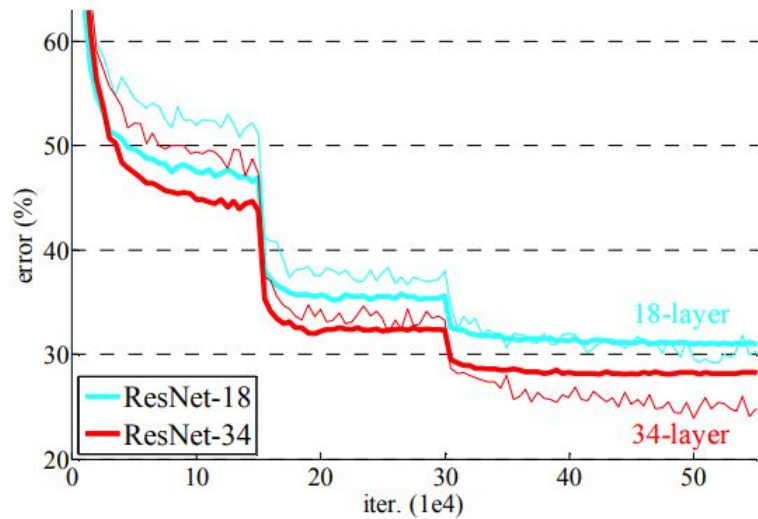
# ResNet

The result for base and deep model is the same



Creating highway to keep the gradient

# ResNet

# ResNet



Results for the deeper model is better! Success!
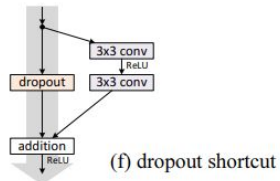
# ResNet

What is the optimal residual layer form?



(a) original
(b) constant scaling
(c) exclusive gating
(d) shortcut-only gating
(e) conv shortcut
(f) dropout shortcut

(a) original
(b) proposed

ResNet–1001, original (error: 7.61%)
ResNet–1001, proposed (error: 4.92%)

$$\mathbf{x}_{l+1} = \mathbf{x}_l + \mathcal{F}\left(\mathbf{x}_l, \mathcal{W}_l\right)$$

# Xception



Standardized Inception module form



Xception module

**Hypothesis**: cross-channel correlations and spatial correlations are sufficiently decoupled that it is preferable not to map them jointly

# Xception



Depthwise separable

1x1 convolution

How Xception architecture looks like?

Depthwise separable + 1x1 conv

# MobileConv



Convolution calculations

$$D_K \cdot D_K \cdot M \cdot N \cdot D_F \cdot D_F$$

Xception calculations

$$D_K \cdot D_K \cdot M \cdot D_F \cdot D_F + M \cdot N \cdot D_F \cdot D_F$$

$$\frac{D_K \cdot D_K \cdot M \cdot D_F \cdot D_F + M \cdot N \cdot D_F \cdot D_F}{D_K \cdot D_K \cdot M \cdot N \cdot D_F \cdot D_F}$$
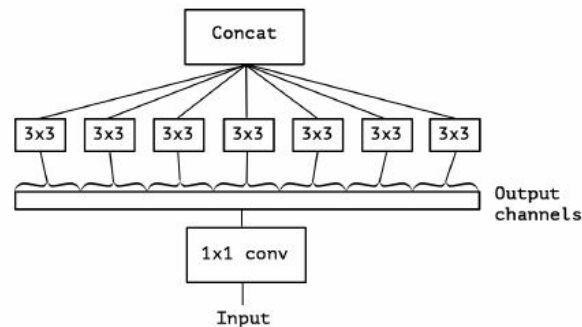
$$= \frac{1}{N} + \frac{1}{D_K^2}$$

if we take N=512 filters and kernel size D=3 -> we get approximately 9 times less calclations

# MobileNet

Table 1. MobileNet Body Architecture

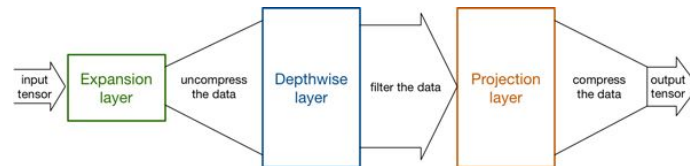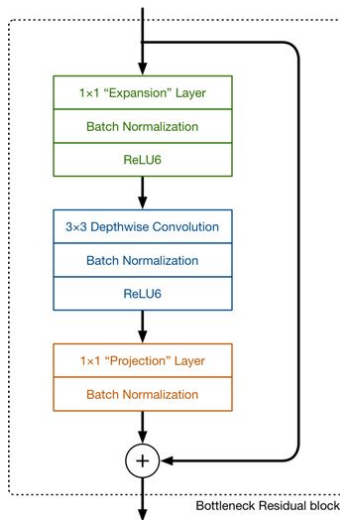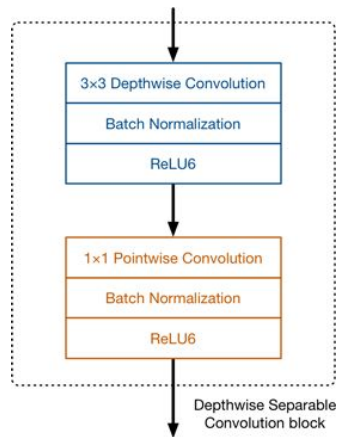| Type / Stride | Filter Shape | Input Size |
|---|---|---|
| Conv / s2 | $3 \times 3 \times 3 \times 32$ | $224 \times 224 \times 3$ |
| Conv dw / s1 | $3 \times 3 \times 32$ dw | $112 \times 112 \times 32$ |
| Conv / s1 | $1 \times 1 \times 32 \times 64$ | $112 \times 112 \times 32$ |
| Conv dw / s2 | $3 \times 3 \times 64$ 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 | $14 \times 14 \times 512$ |
| Conv / s1 | $1 \times 1 \times 512 \times 512$ | $14 \times 14 \times 512$ |
| 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 | $1 \times 1 \times 1000$ |

Table 8. MobileNet Comparison to Popular Models

| Model | ImageNet Accuracy | Million Mult-Adds | Million Parameters |
|---|---|---|---|
| 1.0 MobileNet-224 | 70.6% | 569 | 4.2 |
| GoogleNet | 69.8% | 1550 | 6.8 |
| VGG 16 | 71.5% | 15300 | 138 |

Same quality but number of calculations and number of parameters ~1.5 times less
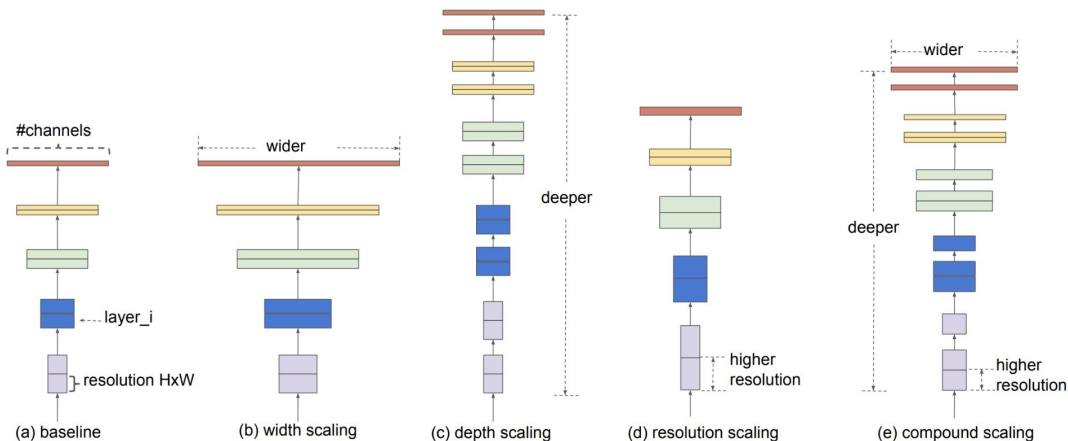
# MobileNetV2



Depthwise Separable Convolution block



Bottleneck Residual block



ReLU in low dimension can kill a lot of information



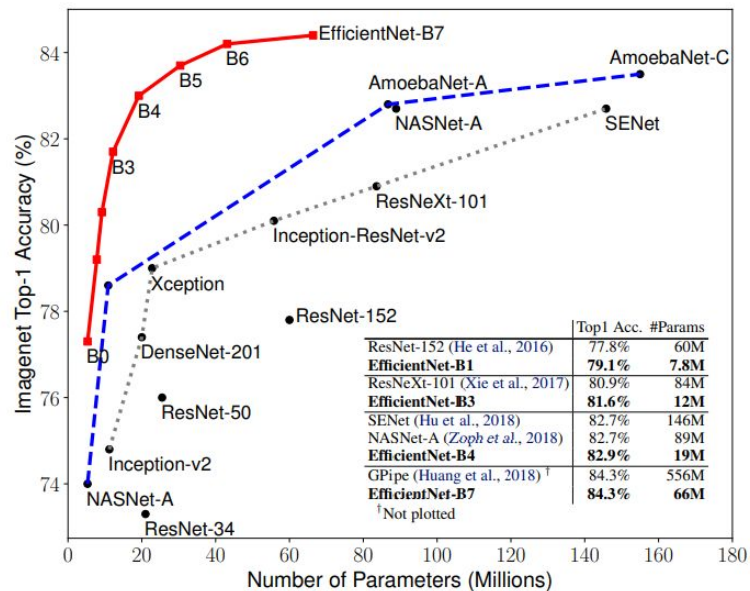Input    Output/dim=2    Output/dim=3    Output/dim=5    Output/dim=15    Output/dim=30

# EfficientNet

We can improve quality of the model by increasing resolution/depth/width of the model. What is the optimal balance between them?



$$\text{depth: } d = \alpha^\phi$$
$$\text{width: } w = \beta^\phi$$
$$\text{resolution: } r = \gamma^\phi$$
$$\text{s.t. } \alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$
$$\alpha \geq 1, \beta \geq 1, \gamma \geq 1$$

# EfficientNet



| | Top1 Acc. | #Params |
|---|---|---|
| ResNet-152 (He et al., 2016) | 77.8% | 60M |
| **EfficientNet-B1** | **79.1%** | **7.8M** |
| ResNeXt-101 (Xie et al., 2017) | 80.9% | 84M |
| **EfficientNet-B3** | **81.6%** | **12M** |
| SENet (Hu et al., 2018) | 82.7% | 146M |
| NASNet-A (Zoph et al., 2018) | 82.7% | 89M |
| **EfficientNet-B4** | **82.9%** | **19M** |
| GPipe (Huang et al., 2018) [†] | 84.3% | 556M |
| **EfficientNet-B7** | **84.3%** | **66M** |

[†]Not plotted

# EfficientNetV2

Modifications

- New NAS metric $A \cdot S^w \cdot P^v$
- FusedMBConv in early stages
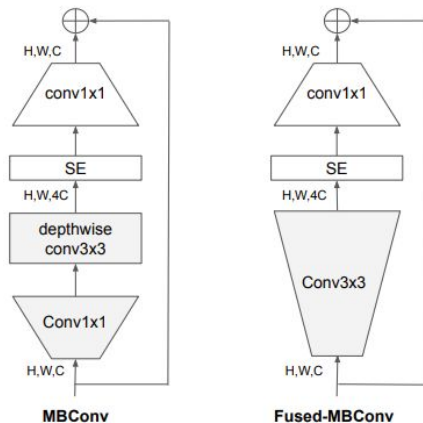- Progressive training



Figure 2. Structure of MBConv and Fused-MBConv.

# ViT

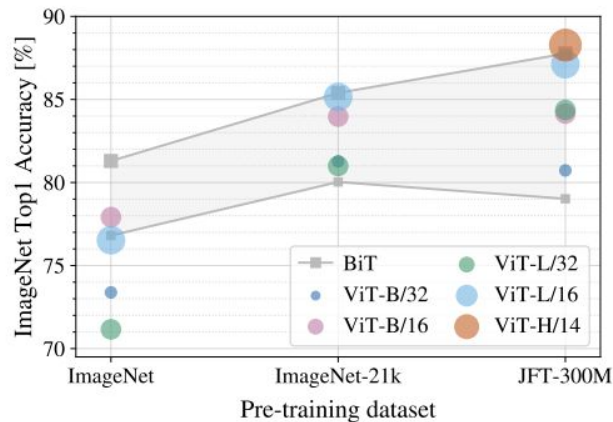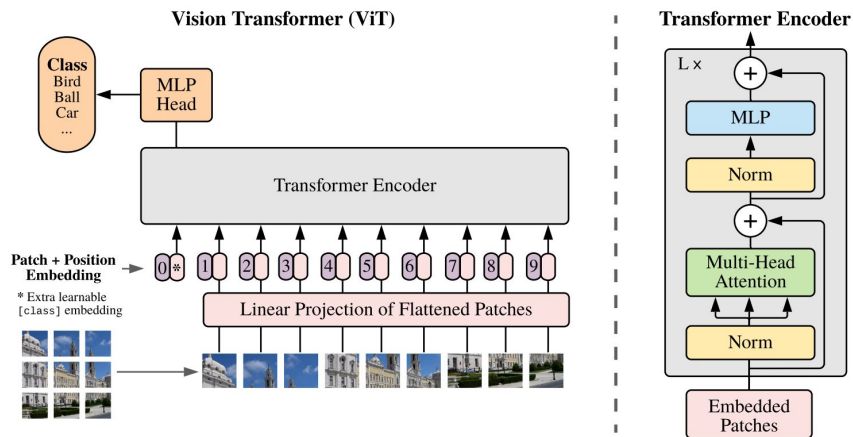How to adapt transformer for computer vision?



Figure 3: Transfer to ImageNet. While large ViT models perform worse than BiT ResNets (shaded area) when pre-trained on small datasets, they shine when pre-trained on larger datasets. Similarly, larger ViT variants overtake smaller ones as the dataset grows.

# Recap

- AlexNet
- VGG
- Inception
- ResNet
- Xception
- MobileNet
- EfficientNet
- ViT