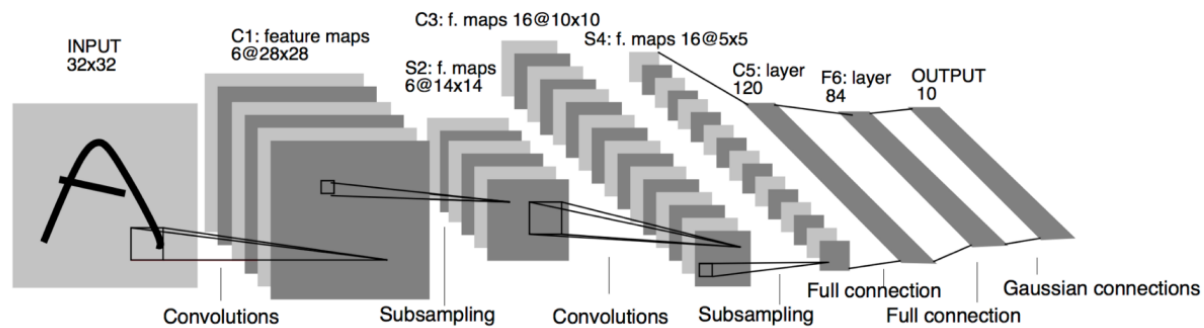# Deep Learning – some popular architectures and history

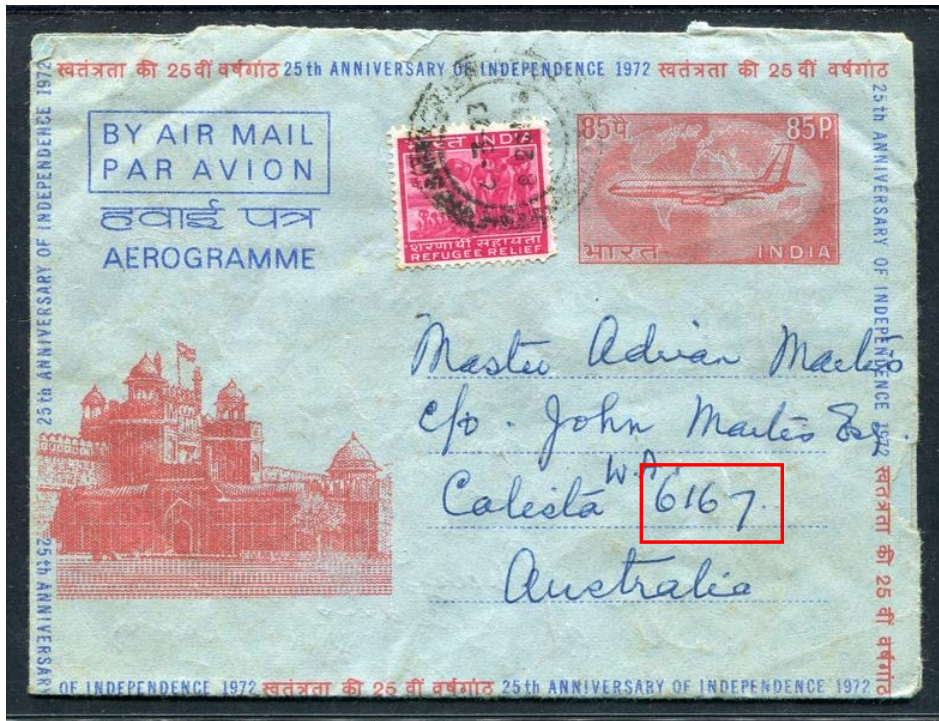Bernhard Kainz

# LeNet-5

**Gradient-Based Learning Applied to Document Recognition**

YANN LECUN, MEMBER, IEEE, LÉON BOTTOU, YOSHUA BENGIO, AND PATRICK HAFFNER



LeCun et al. 1998

# Handwritten digit recognition

# MNIST

- Cantered and scaled
- 50.000 training samples
- 10.000 test samples
- 28 x 28 images
- 10 classes

# Demo from 1995

# LeNet-5



This is expensive if you have many outputs, here only 10

```python
import torch
import torch.nn as nn
import torch.nn.functional as F


class Net(nn.Module):

    def __init__(self):
        super(Net, self).__init__()
        # 1 input image channel, 6 output channels, 3x3 square convolution
        # kernel
        self.conv1 = nn.Conv2d(1, 6, 3)
        self.conv2 = nn.Conv2d(6, 16, 3)
        # an affine operation: y = Wx + b
        self.fc1 = nn.Linear(16 * 6 * 6, 120)  # 6*6 from image dimension
        self.fc2 = nn.Linear(120, 84)
        self.fc3 = nn.Linear(84, 10)

    def forward(self, x):
        # Max pooling over a (2, 2) window
        x = F.max_pool2d(F.relu(self.conv1(x)), (2, 2))
        # If the size is a square you can only specify a single number
        x = F.max_pool2d(F.relu(self.conv2(x)), 2)
        x = x.view(-1, self.num_flat_features(x))
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return x

    def num_flat_features(self, x):
        size = x.size()[1:]  # all dimensions except the batch dimension
        num_features = 1
        for s in size:
            num_features *= s
        return num_features


net = Net()
print(net)
```
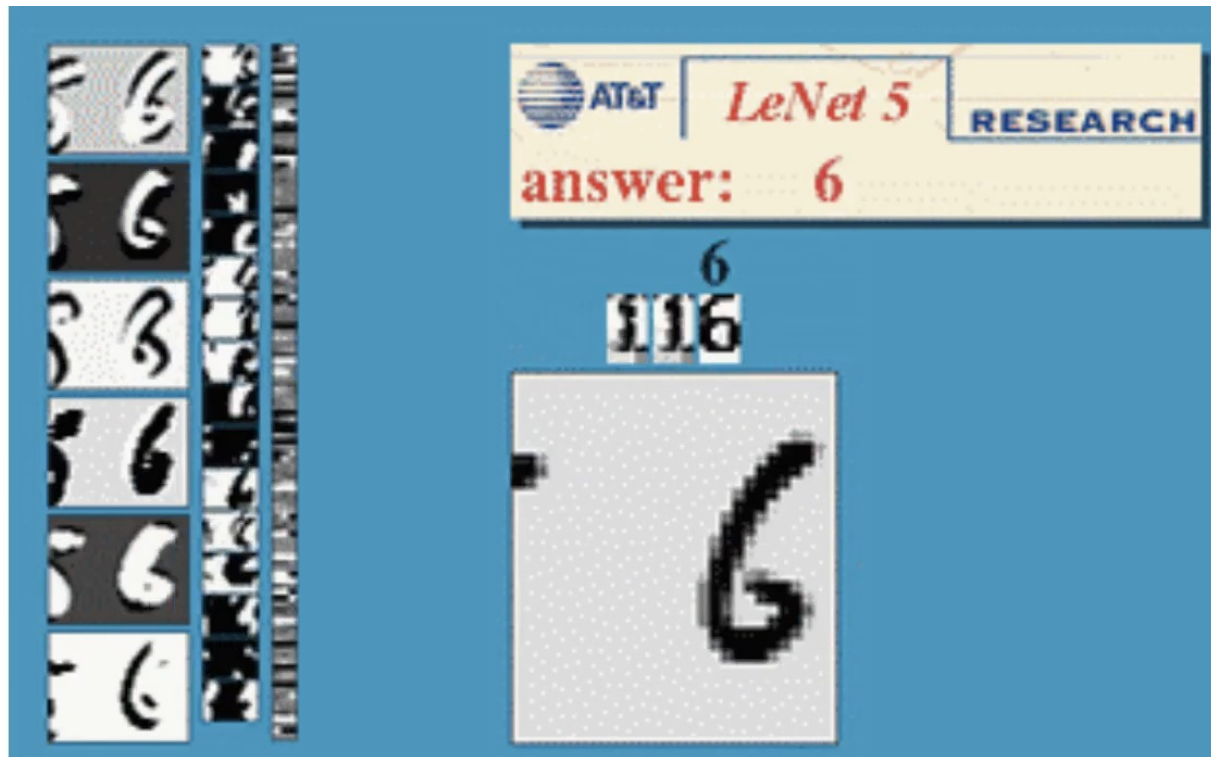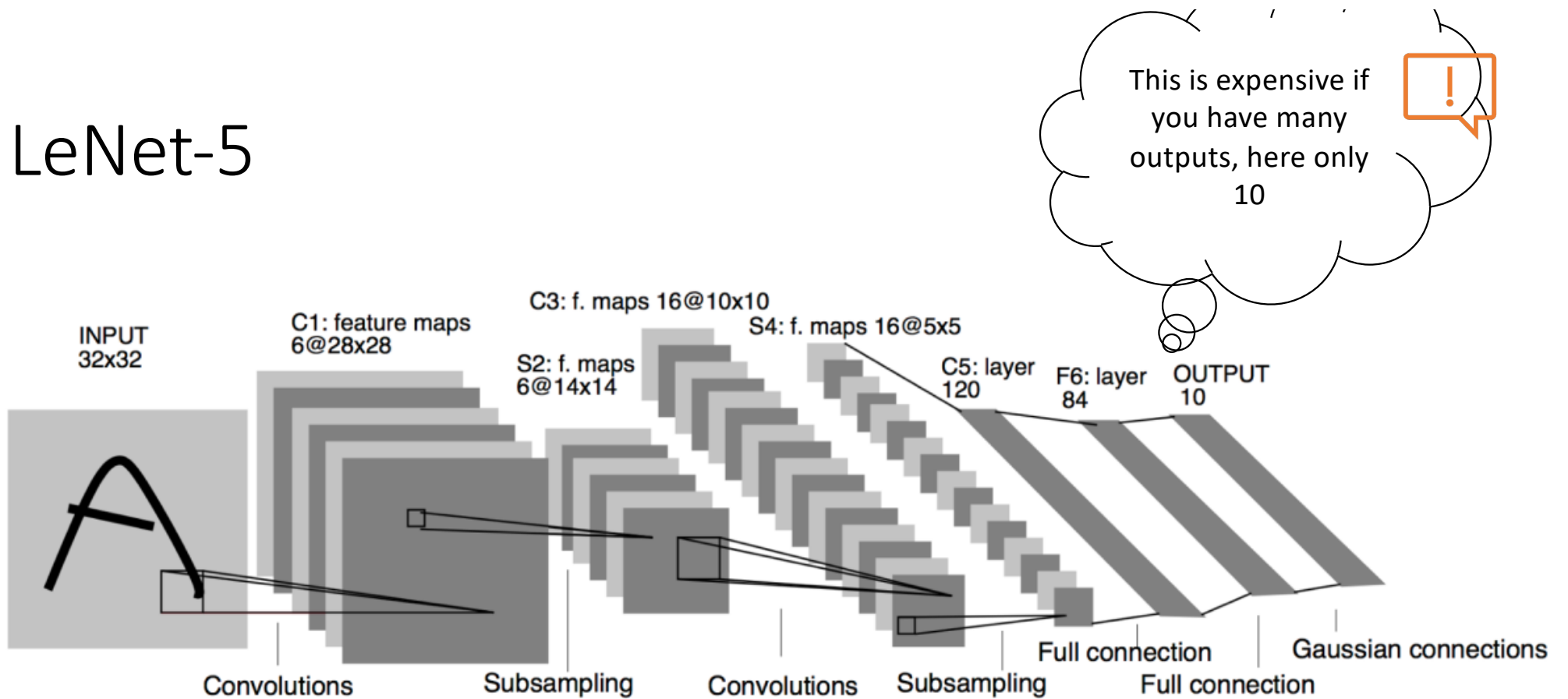


INPUT 32x32 · C1: feature maps 6@28x28 · S2: f. maps 6@14x14 · C3: f. maps 16@10x10 · S4: f. maps 16@5x5 · C5: layer 120 · F6: layer 84 · OUTPUT 10

Convolutions · Subsampling · Convolutions · Subsampling · Full connection · Full connection · Gaussian connections

https://pytorch.org/tutorials/beginner/blitz/neural_networks_tutorial.html

– Bernhard Kainz

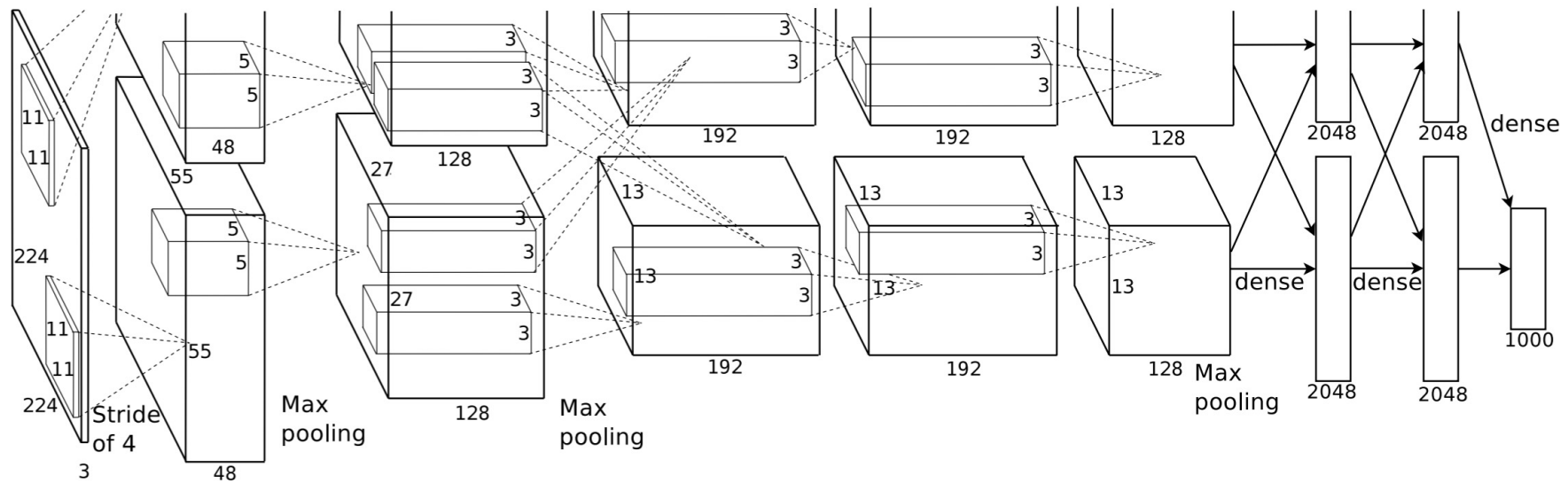# AlexNet

# AlexNet



Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.
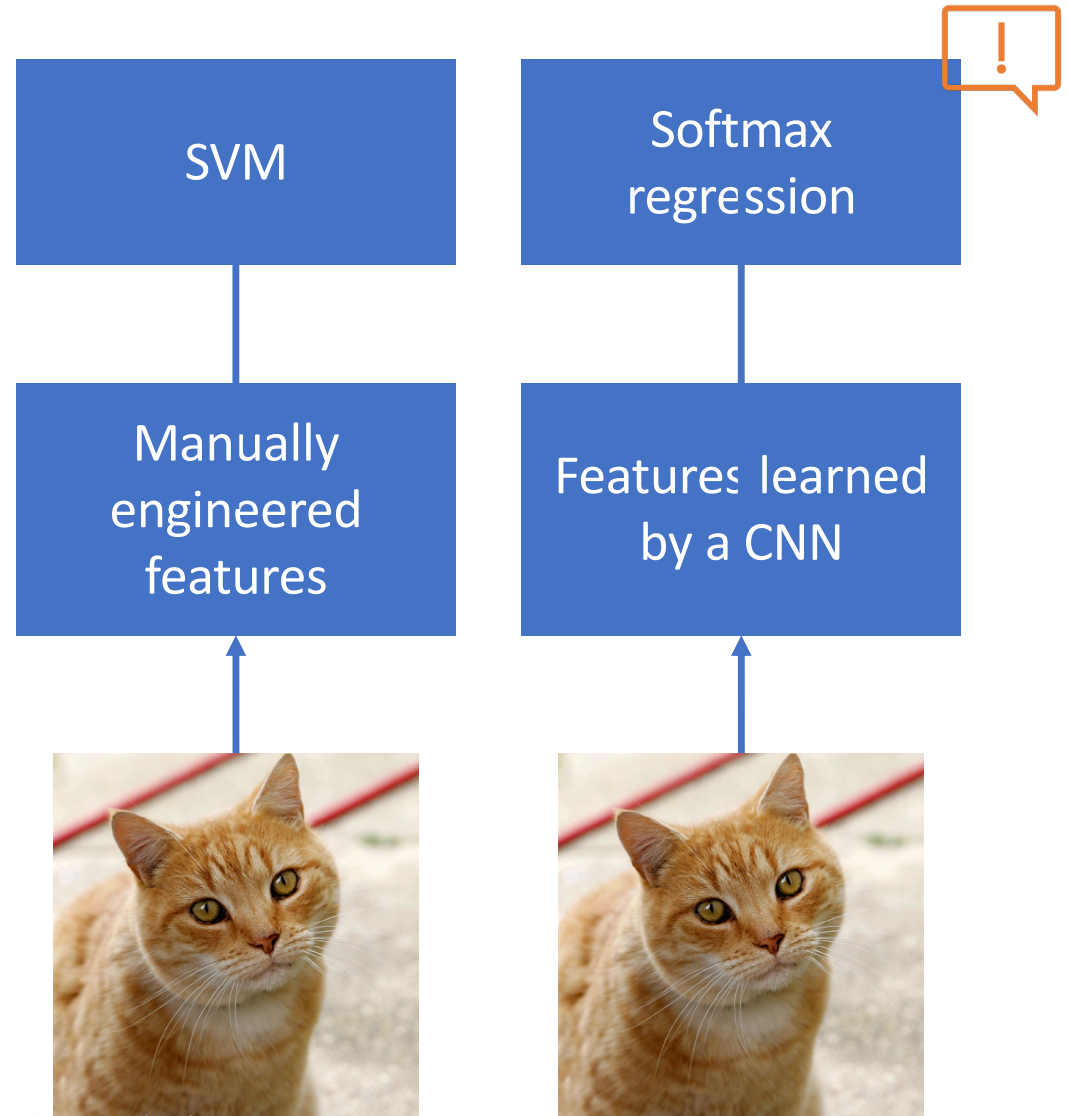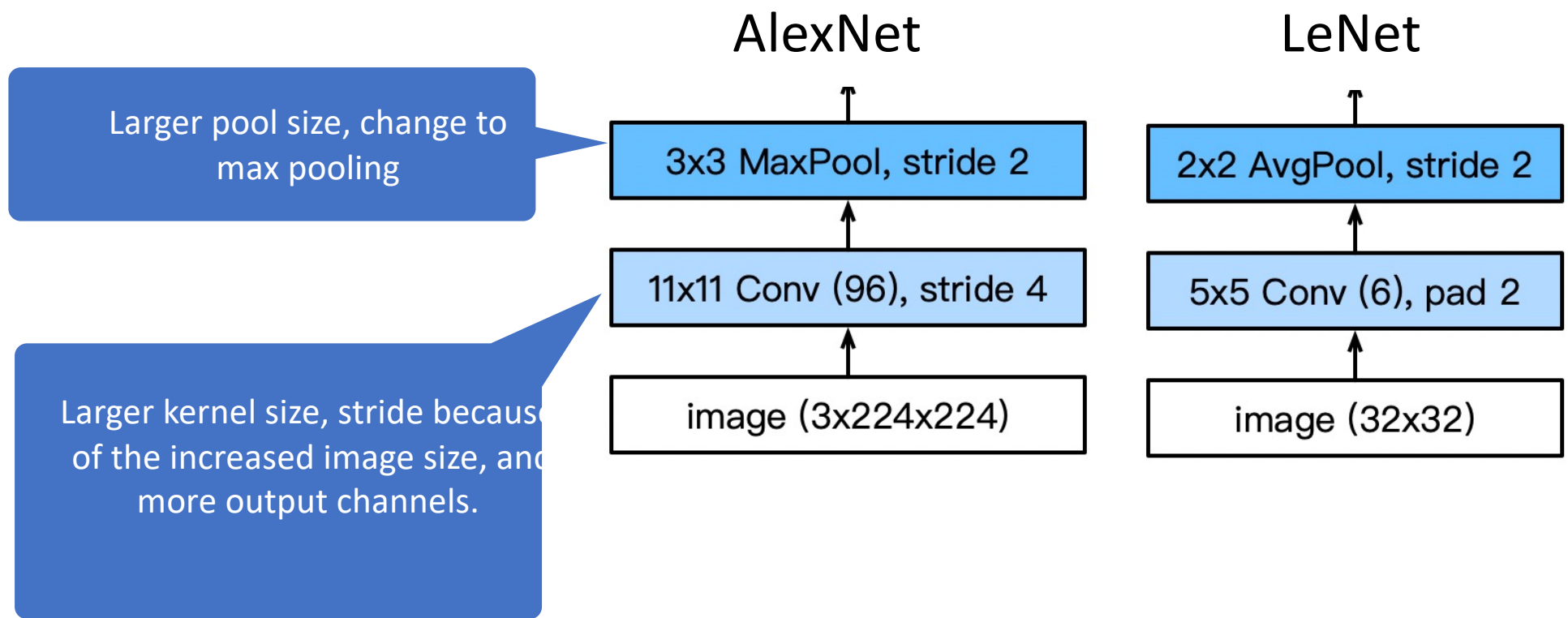
# ImageNet (2010)



| Images | Color images with nature objects | Gray image for hand-written digits |
|---|---|---|
| Size | 469 x 387 | 28 x 28 |
| # examples | 1.2 M | 60 K |
| # classes | 1,000 | 10 |

# AlexNet

- AlexNet won ImageNet competition in 2012

- Deeper and bigger LeNet

- Key modifications:
  - Add a dropout layer after two hidden dense layers
    (better robustness / regularization)
  - Change activation function from sigmoid to ReLu
    (no more vanishing gradient)
  - MaxPooling
  - Heavy data augmentation
  - Model ensembling
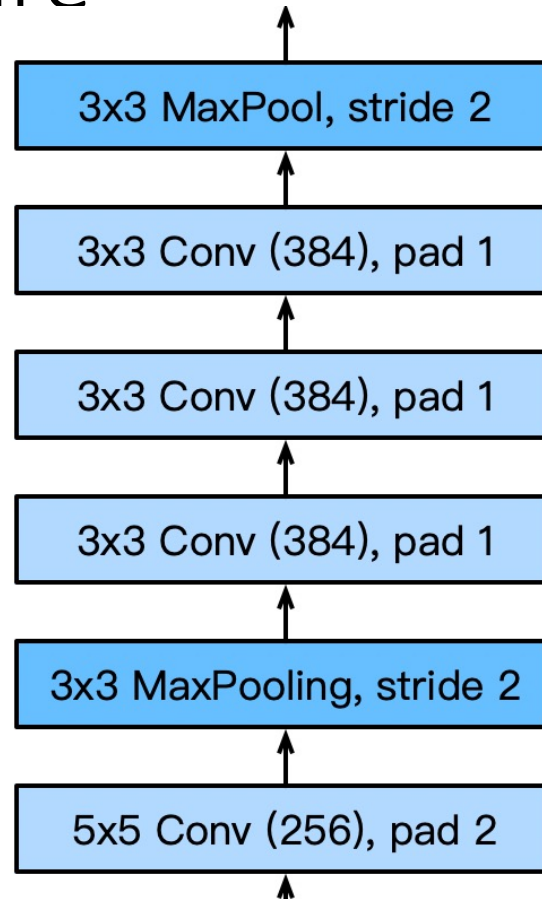
- Paradigm shift for computer vision

SVM

Manually engineered features

Softmax regression

Features learned by a CNN

Deep Learning – Bernhard Kainz

Slide adopted from Alex Smola

# AlexNet Architecture

AlexNet                    LeNet

Larger pool size, change to max pooling

| 3x3 MaxPool, stride 2 | | 2x2 AvgPool, stride 2 |

| 11x11 Conv (96), stride 4 | | 5x5 Conv (6), pad 2 |

| image (3x224x224) | | image (32x32) |

Larger kernel size, stride because of the increased image size, and more output channels.

Slide adopted from Alex Smola

Deep Learning – Bernhard Kainz

# AlexNet Architecture

## AlexNet



3x3 MaxPool, stride 2

3x3 Conv (384), pad 1

3x3 Conv (384), pad 1

**3 additional convolutional layers**

3x3 Conv (384), pad 1

3x3 MaxPooling, stride 2

**More output channels.**

5x5 Conv (256), pad 2

## LeNet

2x2 AvgPool, stride 2

5x5 Conv (16)

# AlexNet Architecture

AlexNet

LeNet

1000 classes output → Dense (1000)  | Dense (10)

Dense (4096) | Dense (84)

Increase hidden size from 120 to 4096 → Dense (4096) | Dense (120)

# Complexity

| | #parameters | | FLOP | |
|---|---|---|---|---|
| | **AlexNet** | **LeNet** | **AlexNet** | **LeNet** |
| **Conv1** | 35K | 150 | 101M | 1.2M |
| **Conv2** | 614K | 2.4K | 415M | 2.4M |
| **Conv3-5** | 3M | | 445M | |
| **Dense1** | 26M | 0.48M | 26M | 0.48M |
| **Dense2** | 16M | 0.1M | 16M | 0.1M |
| **Total** | 46M | 0.6M | 1G | 4M |
| **Increase** | 11x | 1x | 250x | 1x |



Dense (1000)

Dense (4096)

Dense (4096)

Max Pooling

3x3 Conv (384)

3x3 Conv (384)

3x3 Conv (384)

Max Pooling

5x5 Conv (256)

Max Pooling

11x11 Conv (96), stride 4

image (224x224)

Slide adopted from Alex Smola

# demo



Silicon Valley: Season 4 Episode 4: Not Hotdog (HBO)
https://www.youtube.com/watch?v=pqTntG1RXSY

```python
14    class AlexNet(nn.Module):
15
16        def __init__(self, num_classes=1000):
17            super(AlexNet, self).__init__()
18            self.features = nn.Sequential(
19                nn.Conv2d(3, 64, kernel_size=11, stride=4, padding=2),
20                nn.ReLU(inplace=True),
21                nn.MaxPool2d(kernel_size=3, stride=2),
22                nn.Conv2d(64, 192, kernel_size=5, padding=2),
23                nn.ReLU(inplace=True),
24                nn.MaxPool2d(kernel_size=3, stride=2),
25                nn.Conv2d(192, 384, kernel_size=3, padding=1),
26                nn.ReLU(inplace=True),
27                nn.Conv2d(384, 256, kernel_size=3, padding=1),
28                nn.ReLU(inplace=True),
29                nn.Conv2d(256, 256, kernel_size=3, padding=1),
30                nn.ReLU(inplace=True),
31                nn.MaxPool2d(kernel_size=3, stride=2),
32            )
33            self.avgpool = nn.AdaptiveAvgPool2d((6, 6))
34            self.classifier = nn.Sequential(
35                nn.Dropout(),
36                nn.Linear(256 * 6 * 6, 4096),
37                nn.ReLU(inplace=True),
38                nn.Dropout(),
39                nn.Linear(4096, 4096),
40                nn.ReLU(inplace=True),
41                nn.Linear(4096, num_classes),
42            )
43
44        def forward(self, x):
45            x = self.features(x)
46            x = self.avgpool(x)
47            x = torch.flatten(x, 1)
48            x = self.classifier(x)
49            return x
50
```
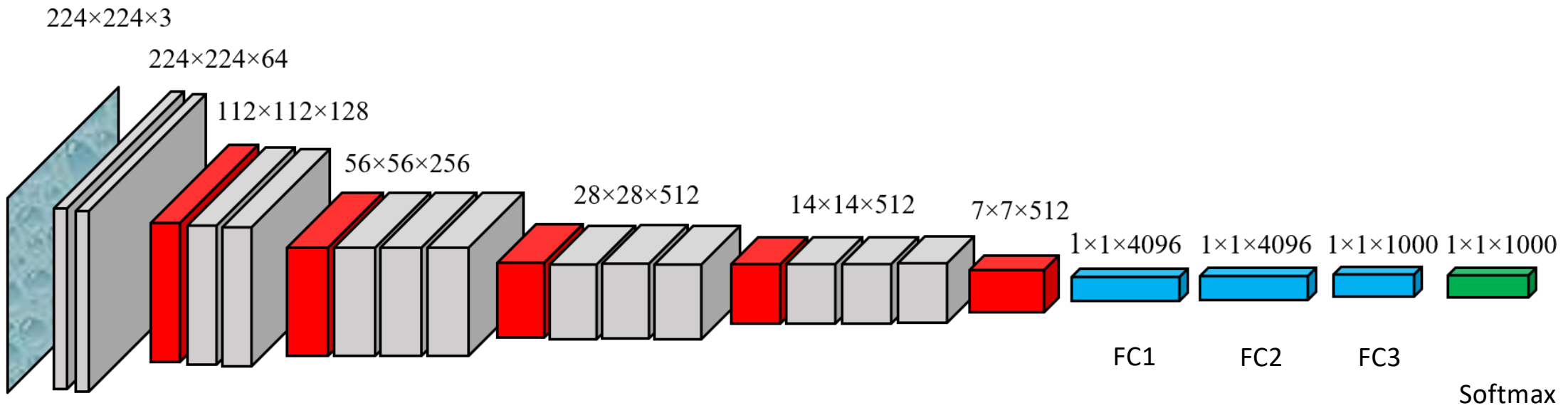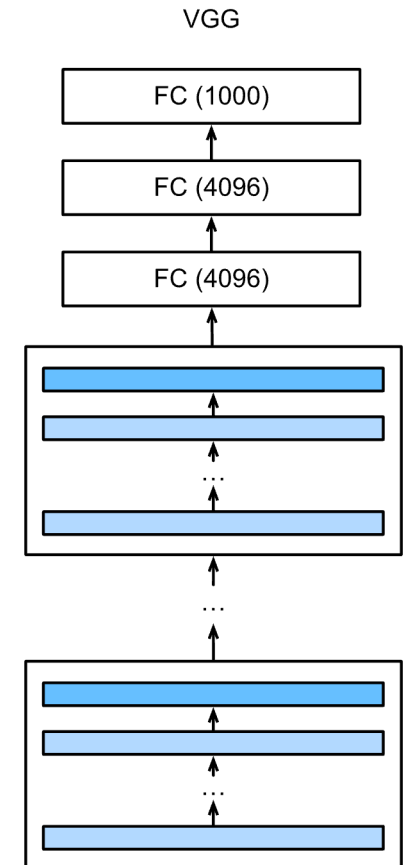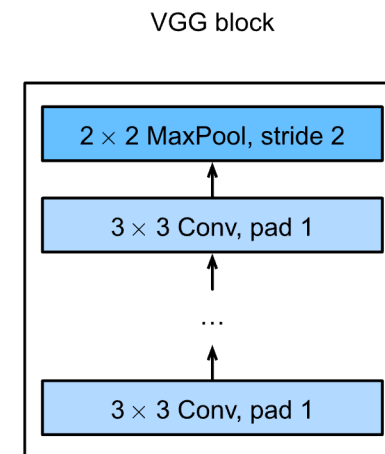


https://github.com/pytorch/vision/blob/master/torchvision/models/alexnet.py

Deep Learning – Bernhard Kainz

# VGG

Lecture inspired by Alex Smola with add-ons

# VGG



224×224×3

224×224×64

112×112×128

56×56×256

28×28×512

14×14×512

7×7×512

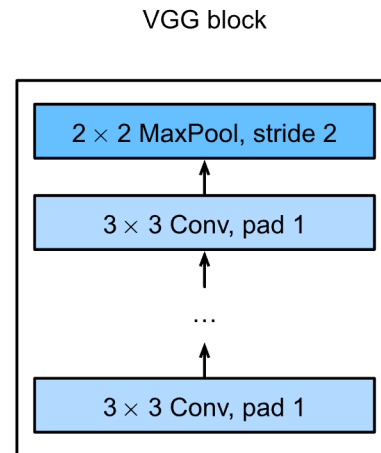1×1×4096  1×1×4096  1×1×1000  1×1×1000

FC1  FC2  FC3

Softmax

- AlexNet = bigger than LeNet

- Bigger = better?

- Options
  - **More** dense layers (too expensive)
  - **More** convolutions
  - Group into **blocks**

Deep Learning – Bernhard Kainz    http://d2l.ai/chapter_convolutional-modern/vgg.html

# VGG blocks

VGG block



- Deeper vs. wider?
  - 13x13?
  - 5x5?
  - 3x3?
  - Deep and narrow = better
- VGG block
  - 3x3 convolutions (pad 1)
    (n layers, m channels)
  - 2x2 max-pooling
    (stride 2)

## VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE-SCALE IMAGE RECOGNITION

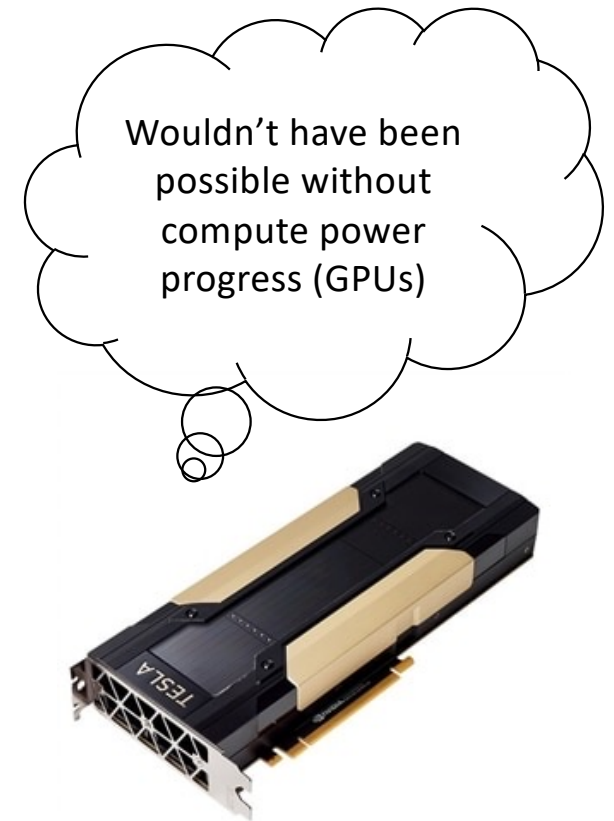**Karen Simonyan*** **& Andrew Zisserman+**
Visual Geometry Group, Department of Engineering Science, University of Oxford
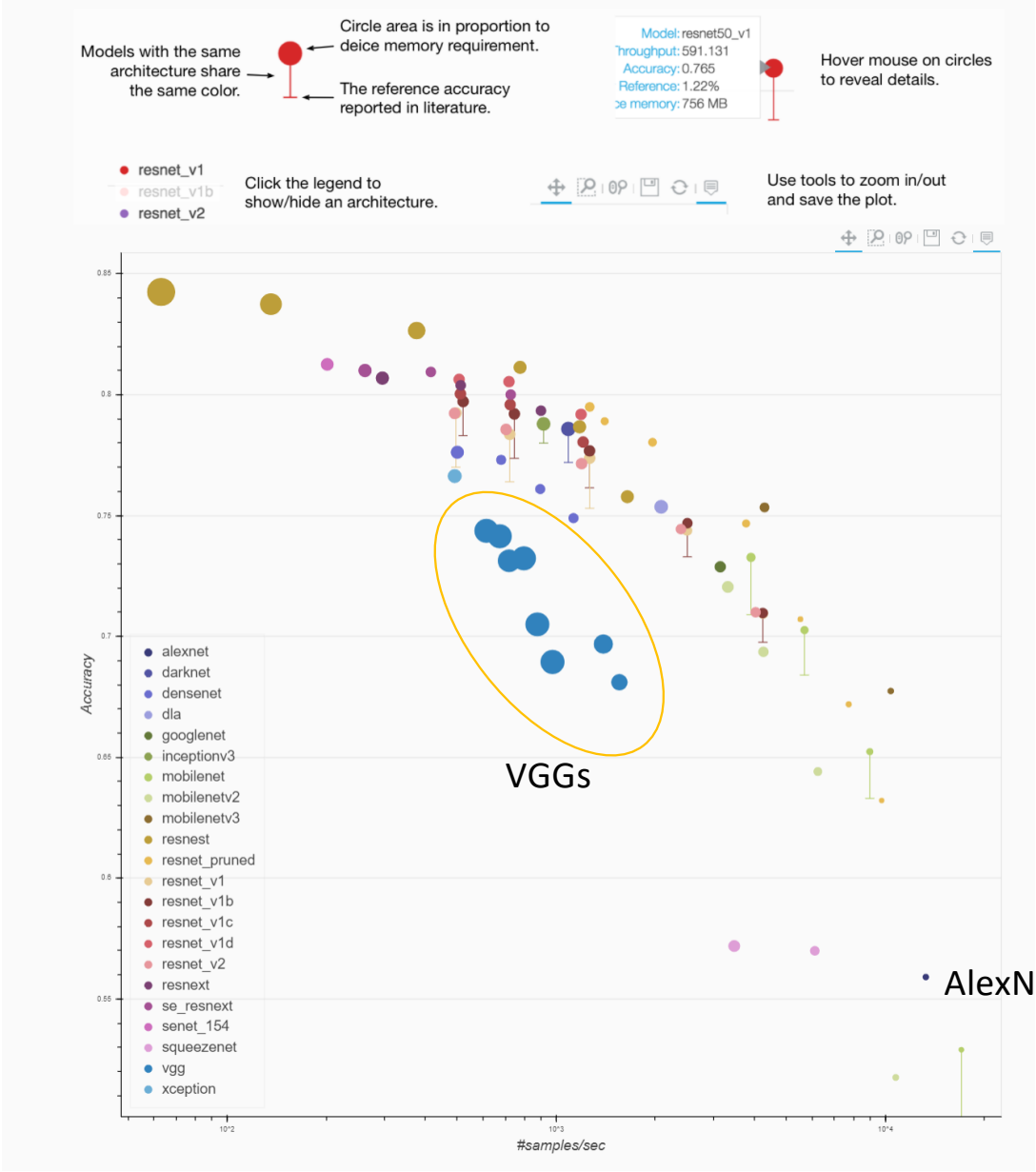{karen,az}@robots.ox.ac.uk

### ABSTRACT

In this work we investigate the effect of the convolutional network depth on its accuracy in the large-scale image recognition setting. Our main contribution is a thorough evaluation of networks of increasing depth using an architecture with very small ($3 \times 3$) convolution filters, which shows that a significant improvement on the prior-art configurations can be achieved by pushing the depth to 16–19 weight layers. These findings were the basis of our ImageNet Challenge 2014 submission, where our team secured the first and the second places in the localisation and classification tracks respectively. We also show that our representations generalise well to other datasets, where they achieve state-of-the-art results. We have made our two best-performing ConvNet models publicly available to facilitate further research on the use of deep visual representations in computer vision.

# progress

- LeNet (1995)
  - 2 convolution + pooling layers
  - 2 hidden dense layers
- AlexNet
  - Bigger and deeper LeNet
  - ReLu, Dropout, preprocessing
- VGG
  - Bigger and deeper AlexNet (repeated VGG blocks)

Wouldn't have been possible without compute power progress (GPUs)

Models with the same architecture share the same color.

Circle area is in proportion to deice memory requirement.

The reference accuracy reported in literature.

Model: resnet50_v1
Throughput: 591.131
Accuracy: 0.765
Reference: 1.22%
ce memory: 756 MB

Hover mouse on circles to reveal details.

● resnet_v1
● resnet_v1b
● resnet_v2

Click the legend to show/hide an architecture.

Use tools to zoom in/out and save the plot.

● alexnet
● darknet
● densenet
● dla
● googlenet
● inceptionv3
● mobilenet
● mobilenetv2
● mobilenetv3
● resnest
● resnet_pruned
● resnet_v1
● resnet_v1b
● resnet_v1c
● resnet_v1d
● resnet_v2
● resnext
● se_resnext
● senet_154
● squeezenet
● vgg
● xception

VGGs

AlexNet

Accuracy

#samples/sec

https://cv.gluon.ai/model_zoo/classification.html

```python
148
149    def vgg16(pretrained: bool = False, progress: bool = True, **kwargs: Any) -> VGG:
150        r"""VGG 16-layer model (configuration "D")
151        `"Very Deep Convolutional Networks For Large-Scale Image Recognition" <https://arxiv.org/pdf/1409.1556.pdf>`_
152
153        Args:
154            pretrained (bool): If True, returns a model pre-trained on ImageNet
155            progress (bool): If True, displays a progress bar of the download to stderr
156        """
157        return _vgg('vgg16', 'D', False, pretrained, progress, **kwargs)


93    def _vgg(arch: str, cfg: str, batch_norm: bool, pretrained: bool, progress: bool, **kwargs: Any) -> VGG:
94
95        if pretrained:
96            kwargs['init_weights'] = False
97        model = VGG(make_layers(cfgs[cfg], batch_norm=batch_norm), **kwargs)
98        if pretrained:
99            state_dict = load_state_dict_from_url(model_urls[arch],
100                                                  progress=progress)
101            model.load_state_dict(state_dict)
102        return model


    def make_layers(cfg: List[Union[str, int]], batch_norm: bool = False) -> nn.Sequential:
        layers: List[nn.Module] = []
        in_channels = 3
        for v in cfg:
            if v == 'M':
                layers += [nn.MaxPool2d(kernel_size=2, stride=2)]
            else:
                v = cast(int, v)
                conv2d = nn.Conv2d(in_channels, v, kernel_size=3, padding=1)
                if batch_norm:
                    layers += [conv2d, nn.BatchNorm2d(v), nn.ReLU(inplace=True)]
                else:
                    layers += [conv2d, nn.ReLU(inplace=True)]
                in_channels = v
        return nn.Sequential(*layers)
```

```python
25    class VGG(nn.Module):
26
27        def __init__(
28            self,
29            features: nn.Module,
30            num_classes: int = 1000,
31            init_weights: bool = True
32        ) -> None:
33            super(VGG, self).__init__()
34            self.features = features
35            self.avgpool = nn.AdaptiveAvgPool2d((7,
36            self.classifier = nn.Sequential(
37                nn.Linear(512 * 7 * 7, 4096),
38                nn.ReLU(True),
39                nn.Dropout(),
40                nn.Linear(4096, 4096),
41                nn.ReLU(True),
42                nn.Dropout(),
43                nn.Linear(4096, num_classes),
44            )
45            if init_weights:
46                self._initialize_weights()
47
48        def forward(self, x: torch.Tensor) -> tor
49            x = self.features(x)
50            x = self.avgpool(x)
51            x = torch.flatten(x, 1)
52            x = self.classifier(x)
53            return x
```
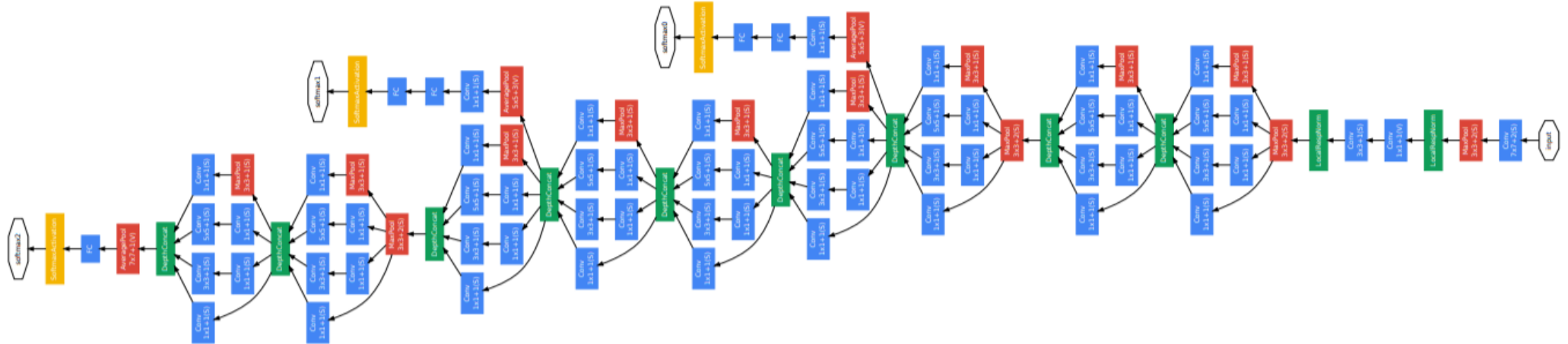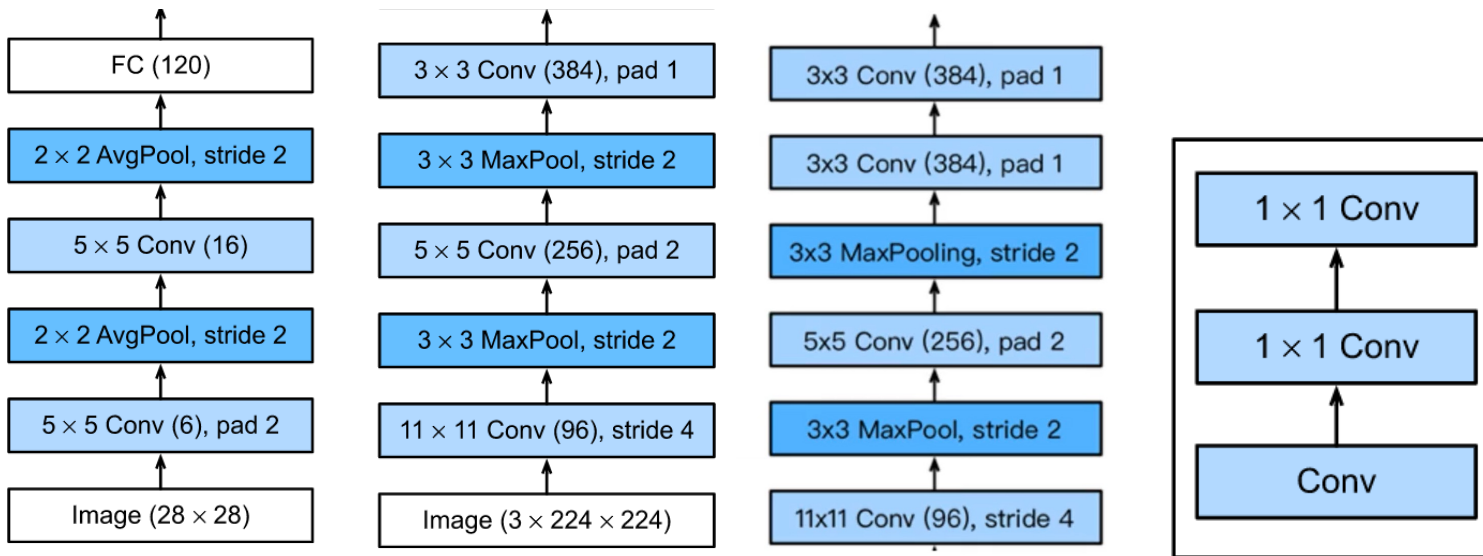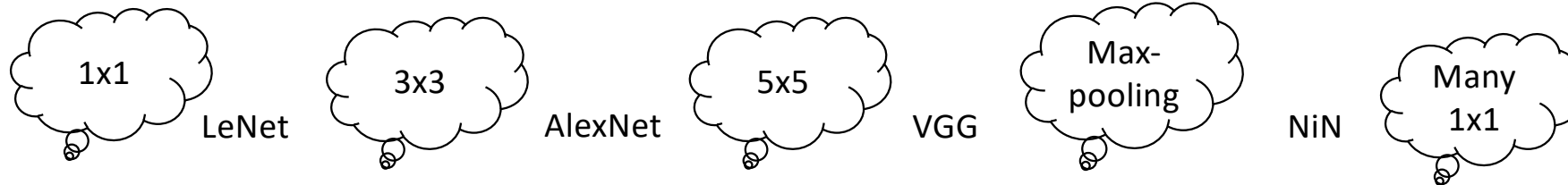


Bernhard Kainz

https://github.com/pytorch/vision/blob/master/torchvision/models/vgg.py

# Inception (GoogLeNet)



Deep Learning – Bernhard Kainz

https://arxiv.org/abs/1409.4842

# Which convolution is the best!?

1x1    3x3    5x5    Max-pooling    Many 1x1

LeNet    AlexNet    VGG    NiN



**LeNet**
- FC (120)
- 2 × 2 AvgPool, stride 2
- 5 × 5 Conv (16)
- 2 × 2 AvgPool, stride 2
- 5 × 5 Conv (6), pad 2
- Image (28 × 28)

**AlexNet**
- 3 × 3 Conv (384), pad 1
- 3 × 3 MaxPool, stride 2
- 5 × 5 Conv (256), pad 2
- 3 × 3 MaxPool, stride 2
- 11 × 11 Conv (96), stride 4
- Image (3 × 224 × 224)

**VGG**
- 3x3 Conv (384), pad 1
- 3x3 Conv (384), pad 1
- 3x3 MaxPooling, stride 2
- 5x5 Conv (256), pad 2
- 3x3 MaxPool, stride 2
- 11x11 Conv (96), stride 4

**NiN**
- 1 × 1 Conv
- 1 × 1 Conv
- Conv

Deep Learning – Bernhard Kainz

# Inception block

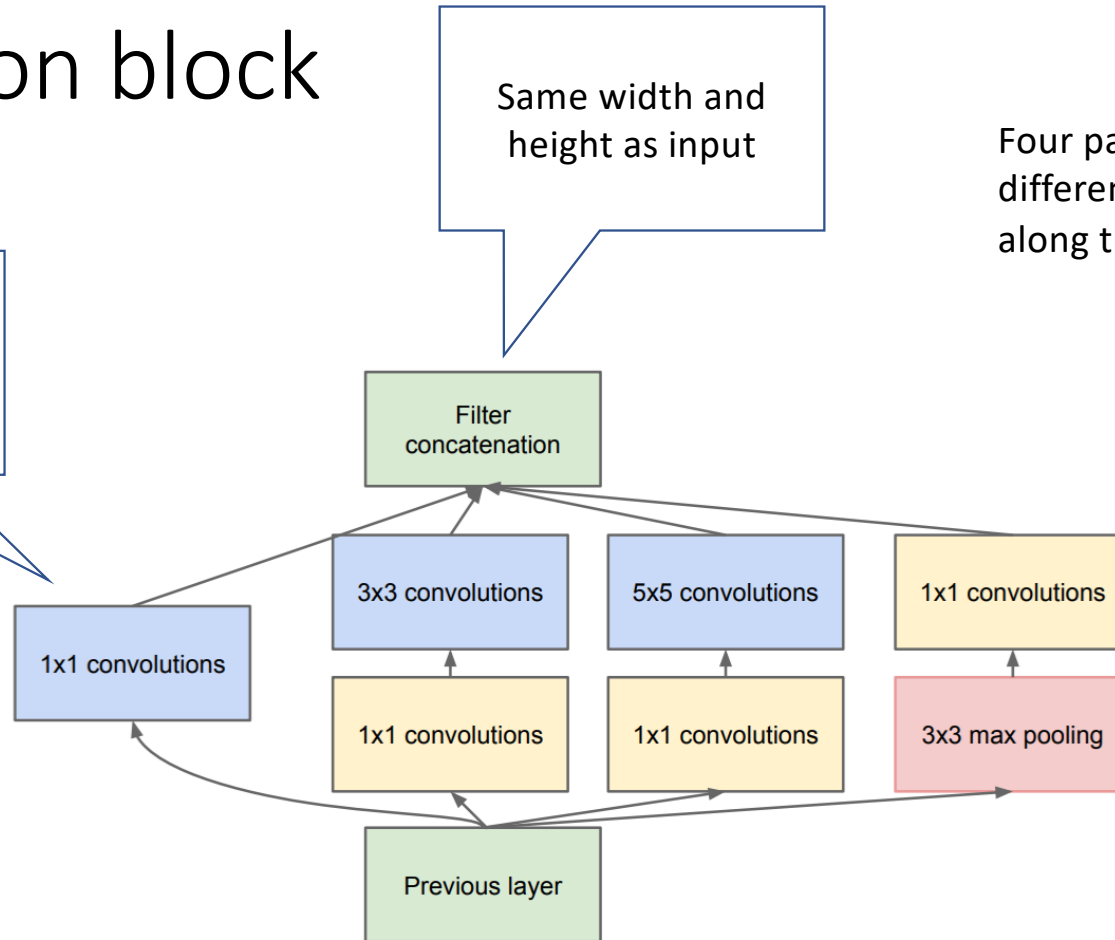https://arxiv.org/abs/1409.4842

# Inception block

Same width and height as input

Four paths extract information from different aspects, then concatenate along the output channel
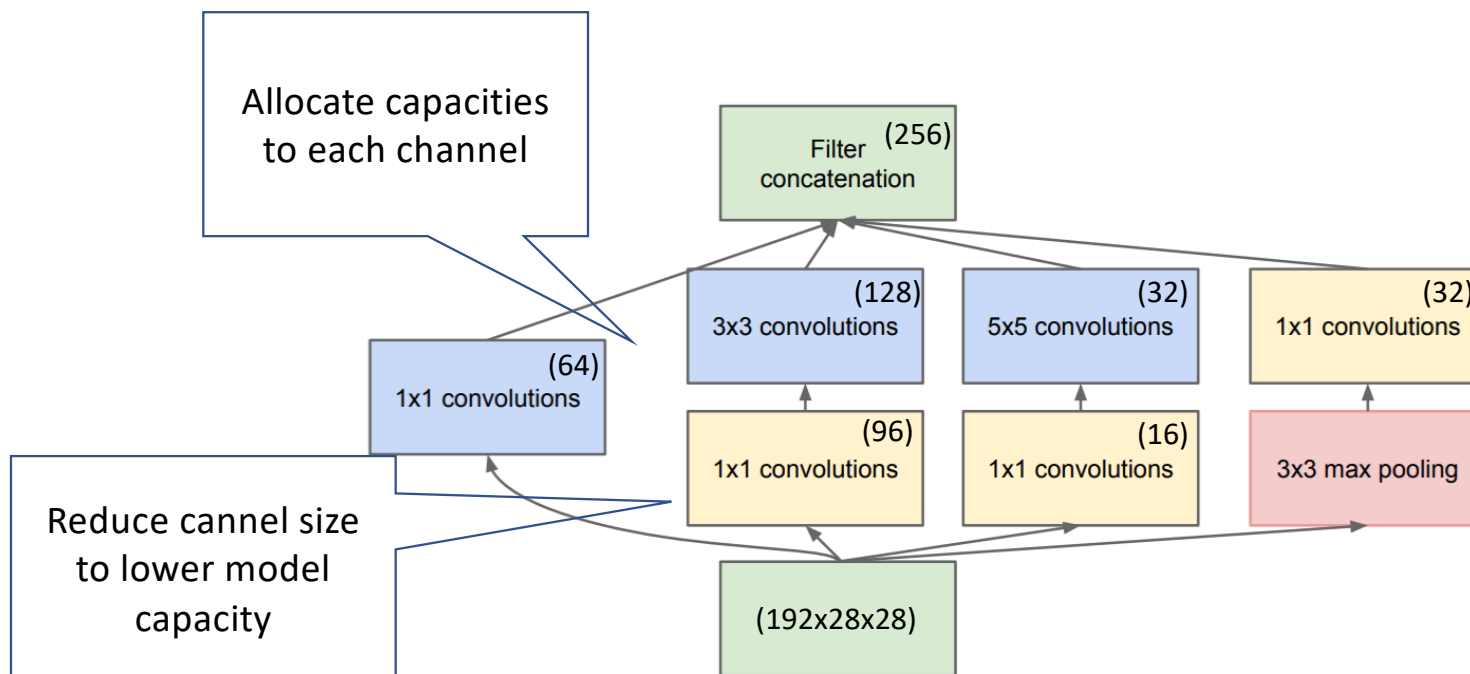
Extract with different size convolutions

Extract spatial information with pooling



Filter concatenation

3x3 convolutions

5x5 convolutions

1x1 convolutions

1x1 convolutions

1x1 convolutions

1x1 convolutions

3x3 max pooling

Previous layer

https://arxiv.org/abs/1409.4842

# Inception block

The first inception block has channel sizes specified

Allocate capacities to each channel

Reduce cannel size to lower model capacity



Filter concatenation (256)

(128) 3x3 convolutions

(32) 5x5 convolutions

(32) 1x1 convolutions

(64) 1x1 convolutions

(96) 1x1 convolutions

(16) 1x1 convolutions

3x3 max pooling

(192x28x28)

Deep Learning – Bernhard Kainz

https://arxiv.org/abs/1409.4842

# Inception blocks

- Inception blocks have fewer parameters and less computation complexity than single 3x3 or 5x5 convolution layers

- They are a mix of different functions, which makes them a powerful function class

- Computing and memory wise they are efficient (good generalisation)

| | #parameters | FLOPS |
|---|---|---|
| Inception | 0.16 M | 128 M |
| 3x3 Conv | 0.44 M | 346 M |
| 5x5 Conv | 1.22 M | 963 M |

As: replace all conv block with 3x3 or 5x5 in Inception

# Less operations?

- $k^2 \times \boxed{c_{in}} \times c_{out} \times m_h \times \boxed{\times m_w}$

  fixed            fixed

- $c_{in} \times m_h \times m_w \times [ \sum_{paths\ j} k_j^2 \times c_{out,j} ]$

allocating compute
to different channels
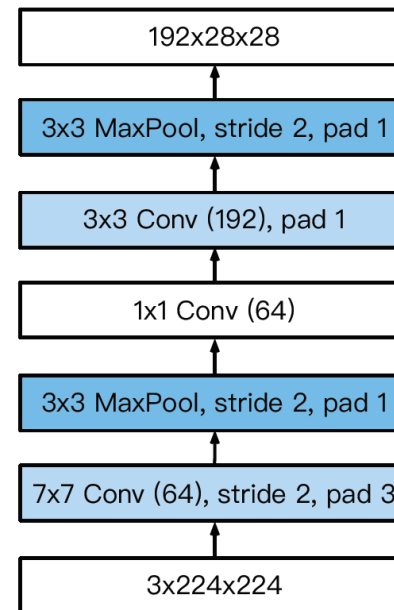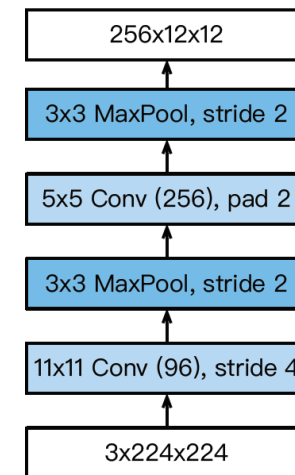= better computing

# Inception

- 5 stages with 9 inception blocks



https://d2l.ai/

Deep Learning – Bernhard Kainz

# Stage 1 and 2

- Smaller kernel size and output channels because of more layers

GoogLeNet

| 192x28x28 |
|---|
| 3x3 MaxPool, stride 2, pad 1 |
| 3x3 Conv (192), pad 1 |
| 1x1 Conv (64) |
| 3x3 MaxPool, stride 2, pad 1 |
| 7x7 Conv (64), stride 2, pad 3 |
| 3x224x224 |

AlexNet

| 256x12x12 |
|---|
| 3x3 MaxPool, stride 2 |
| 5x5 Conv (256), pad 2 |
| 3x3 MaxPool, stride 2 |
| 11x11 Conv (96), stride 4 |
| 3x224x224 |

https://d2l.ai/

# Stage 3



Channel allocation is different

Increases output channel

480x14x14

3x3 MaxPool, stride 2, pad 2

Concatenation (480)

3x3 Conv (192)    5x5 Conv (96)    1x1 Conv (64)

1x1 Conv (128)    1x1 Conv (128)    1x1 Conv (32)    3x3 MaxPool

Input (256)

Concatenation (256)

3x3 Conv (128)    5x5 Conv (32)    1x1 Conv (32)

1x1 Conv (64)    1x1 Conv (96)    1x1 Conv (16)    3x3 MaxPool

192x28x28

# Stage 4 & 5



832x7x7

3x3 MaxPool, stride 2, pad 2

(832)

Increases output channels

(512)

(512)

(512)

Increases output channels

(512)

480x14x14

1024x1x1

1024 dimensional feature to output layer

GlobalAvgPool

(1024)

Increase for output channel

(832)

832x7x7
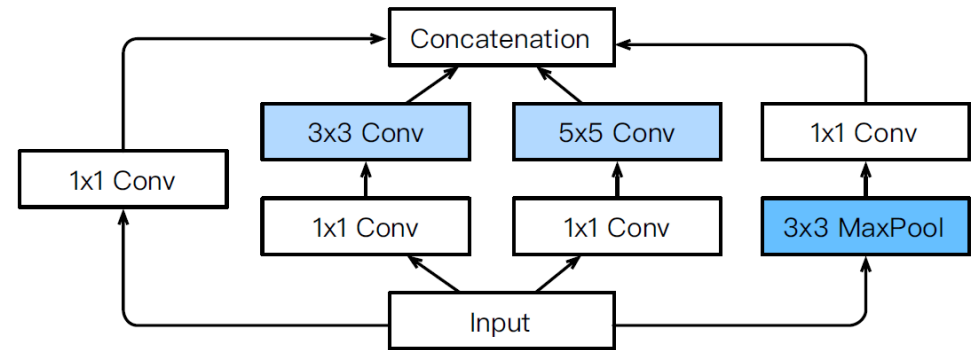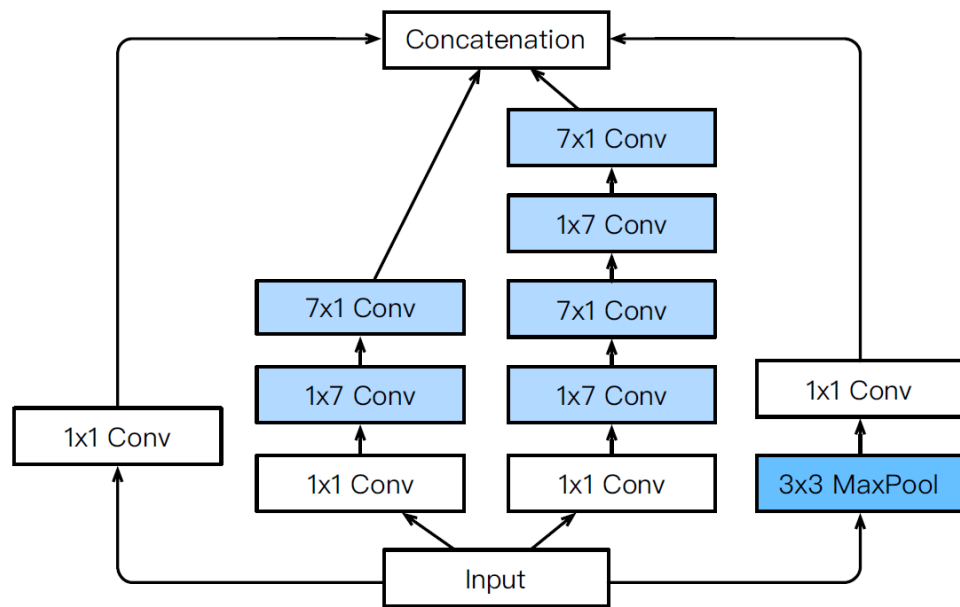
https://d2l.ai/

# Flavours of Inception Networks

- Inception-BN (v2) – added batch normalisation
- Inception-V3 – Modified the inception block
  - Replace 5x5 by multiple 3x3 convolutions
  - Replace 5x5 by 1x7 and 7x1 convolutions
  - Replace 3x3 by 1x3 and 3x1 convolutions
  - Generally deeper stack
- Inception-V4 – adds residual connections

# Inception V3 block for stage 3



https://d2l.ai/

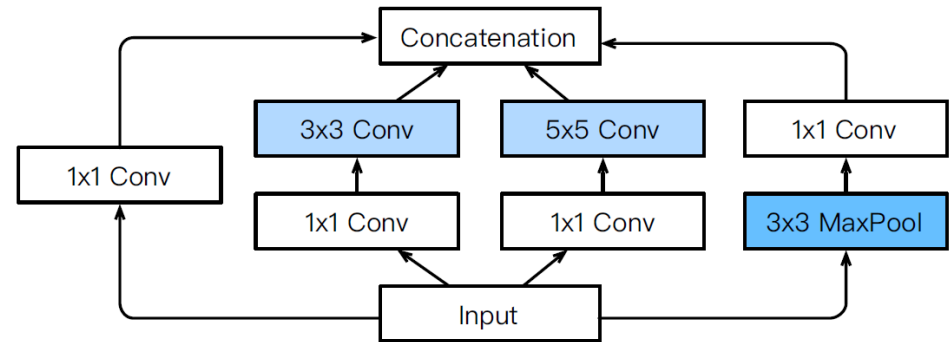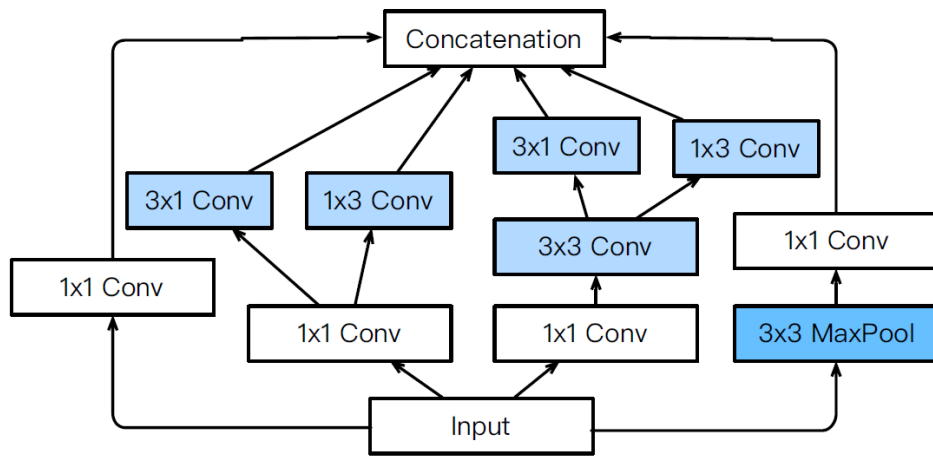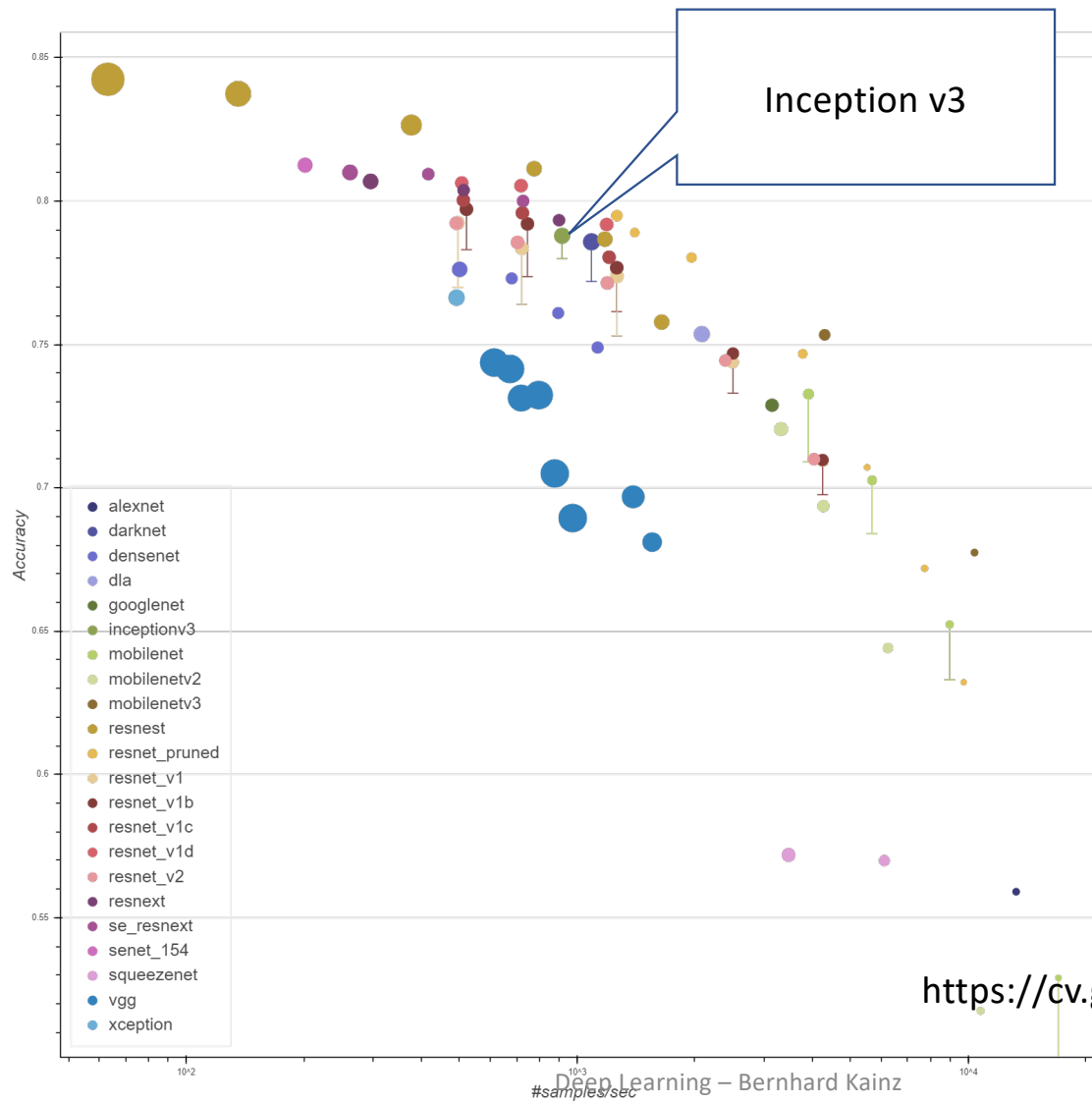Deep Learning – Bernhard Kainz

# Inception V3 block for stage 4

7 x 1

1 x 7

https://d2l.ai/

# Inception V3 block for stage 5



https://d2l.ai/

Deep Learning – Bernhard Kainz
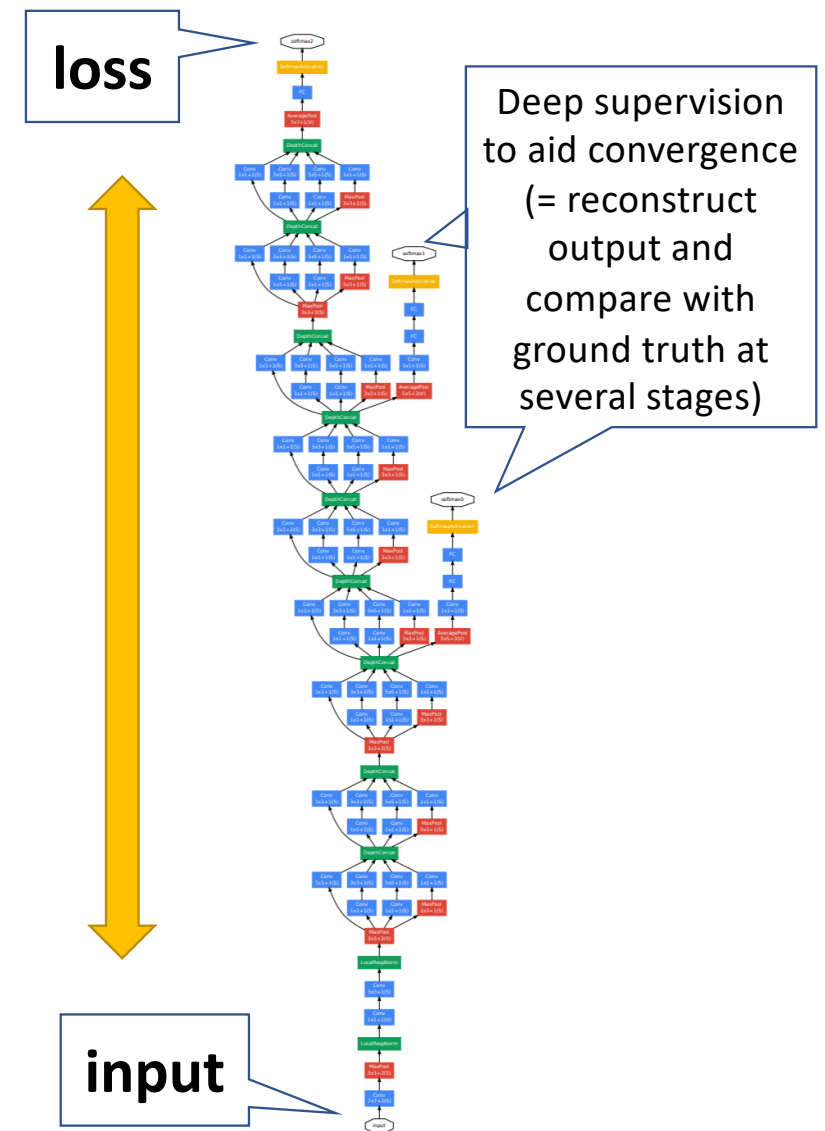
Inception v3

Deep Learning – Bernhard Kainz

# What do we learn from that?

- Dense layers are computationally and memory intensive. Real-world problems with big input tensors and many classes will prohibit their use.

- Again: 1x1 convolutions act like a multi-layer perceptron per pixel.

- Scientists are humans and need a while to understand the power of new approaches. Eventually they do but a lot of vanity is involved in the process.

- If not sure, just take all options and let the optimization decide or even learn this through trial and error (genetic algorithm, AmoebaNet)

# BatchNorm

# Batch Normalization

- Loss is calculated at last layer
  - Last layers learn quickly

- Data input is at first layer
  - First layers change - everything changes
  - Last layers need to relearn many times
  - Slow convergence

- This is like covariate shift...
  Can we avoid changing last layers while
  learning first layers?



loss

input

Deep supervision
to aid convergence
(= reconstruct
output and
compare with
ground truth at
several stages)

Deep Learning – Bernhard Kainz

# Batch Normalization

- Can we avoid changing last layers while learning first layers?
- Fix mean and variance

$$\mu_B = \frac{1}{|B|}\sum_{i \in B} x_i \ and \ \sigma_B^2 = \frac{1}{|B|}\sum_{i \in B}(x_i - \mu_B)^2 + \varepsilon$$

and adjust it separately

$$x_{i+1} = \gamma \frac{x_i - \mu_B}{\sigma_B} + \beta$$

mean

variance

# Batch Normalization

- Doesn't really reduce covariate shift (Lipton et al. 2018) https://arxiv.org/abs/1805.10694
- Regularization by noise injection

$$x_i = \gamma \frac{x_i - \hat{\mu}_B}{\hat{\sigma}_B} + \beta$$

**Random offset**

**learned**

**learned**

**Random scale**

  - Random shift per mini batch
  - Random scale per mini batch

- No need to add dropout (both are capacity control)

- Ideal mini batch size: 64-256

# Batch Normalization

- Dense layer: One normalization for all

- Convolutional layer: One normalization per channel

- Compute **new mean and variance** for every minibatch
  - Acts as regularisation
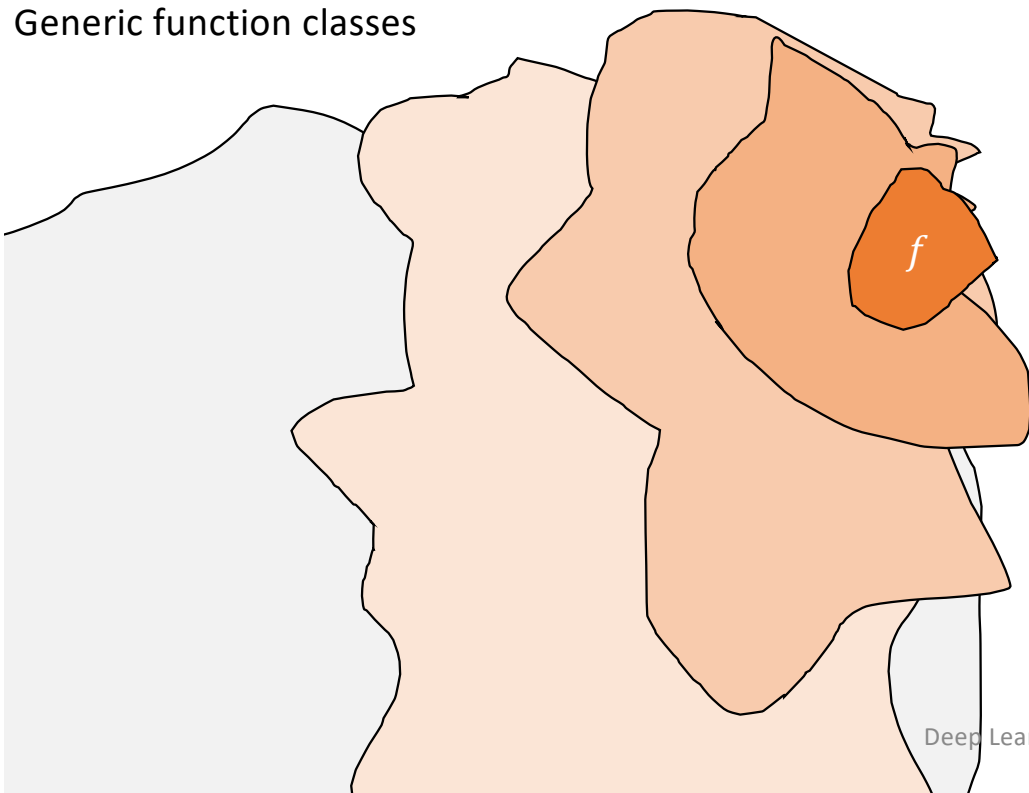  - Be careful when scaling up to multi-GPU training



https://xkcd.com/

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.

# ResNet

# Does adding layers improve accuracy?

The true solution is here ●

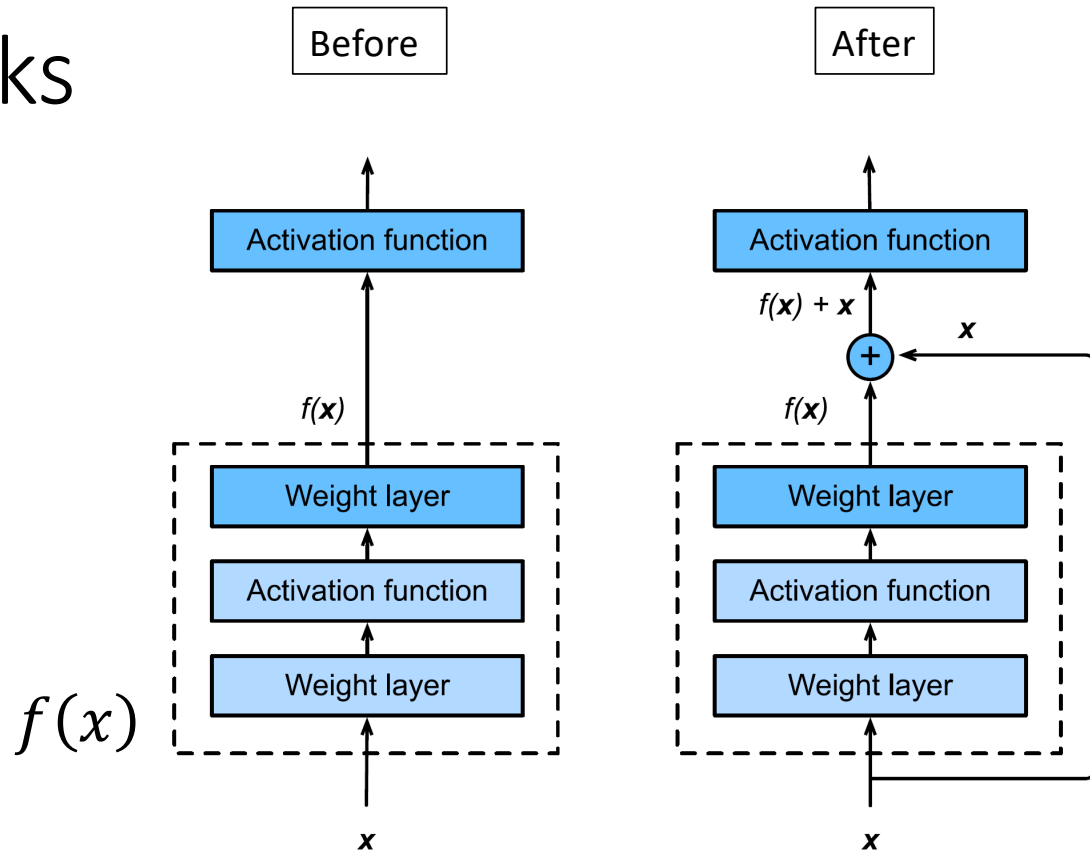Generic function classes

The true solution is here ●

Nested function classes

# Residual Networks

- Adding a layer changes function class
- We want to add to the function class
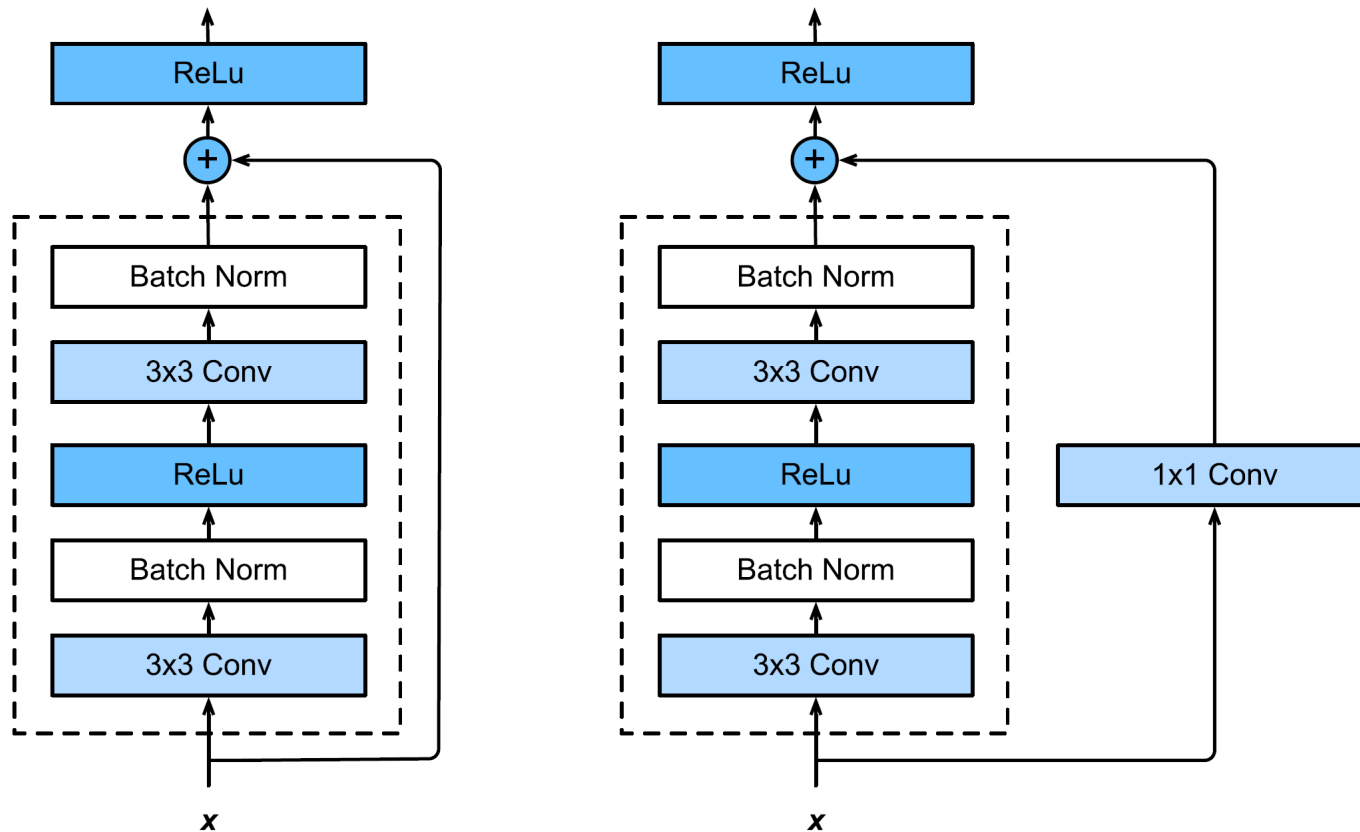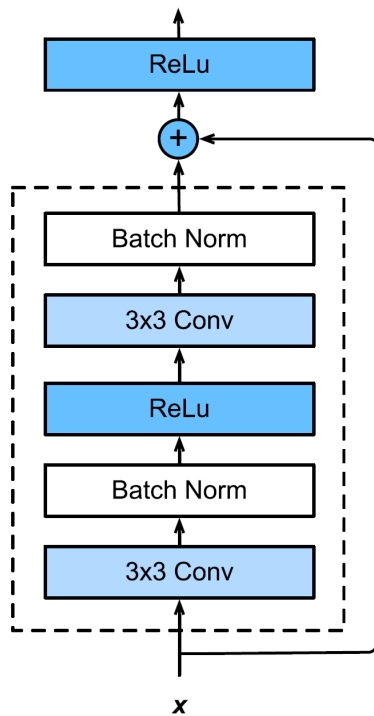- 'Taylor expansion' style parametrization

$f(x)$

Before

After



He et al. 2015 https://arxiv.org/abs/1512.03385

https://d2l.ai/

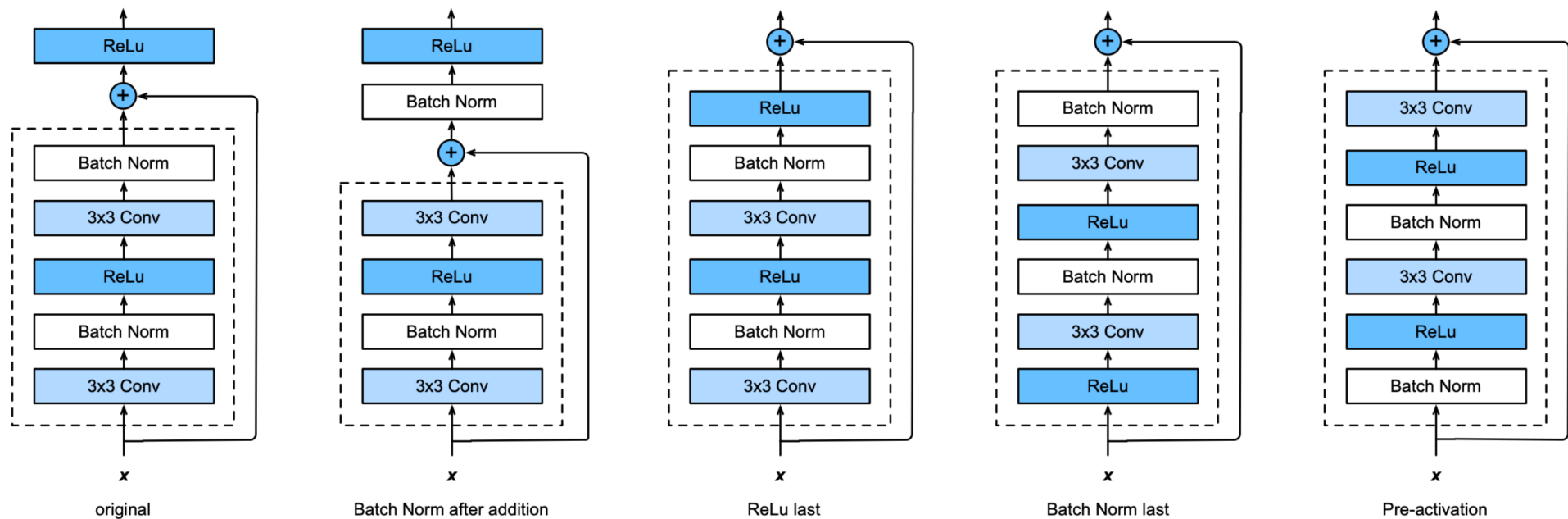# ResNet Block

https://d2l.ai/

# ResNet Block



```
def forward(self, X):
    Y = self.bn1(self.conv1(X))
    Y = nd.relu(Y)
    Y = self.bn2(self.conv2(Y))
    if self.conv3:
        X = self.conv3(X)
    return nd.relu(Y + X)
```
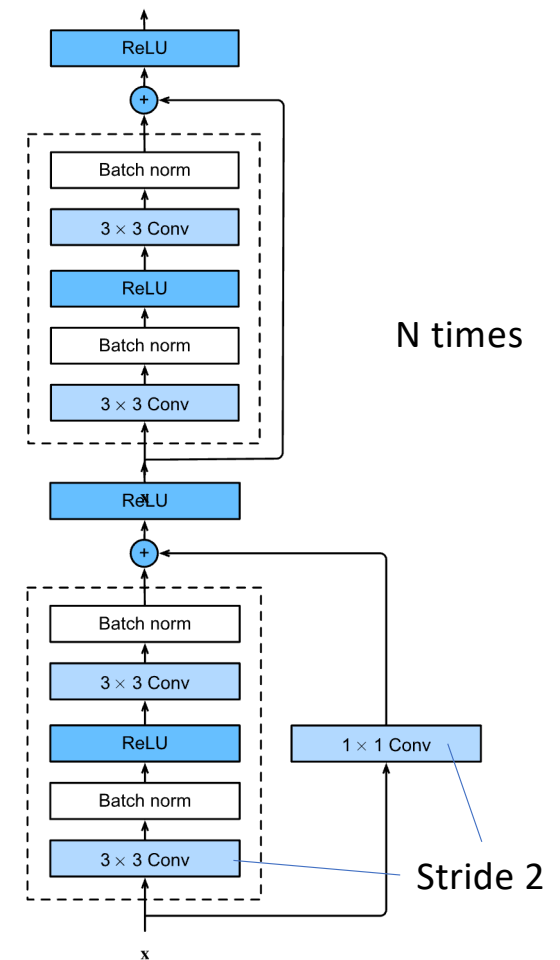
# ResNet block flavours



Trial and error of every permutation

https://d2l.ai/

# ResNet Module

- Downsample per module (stride=2)

- Enforce some nontrivial nonlinearity per module (via 1x1 convolution)

- Stack up in blocks
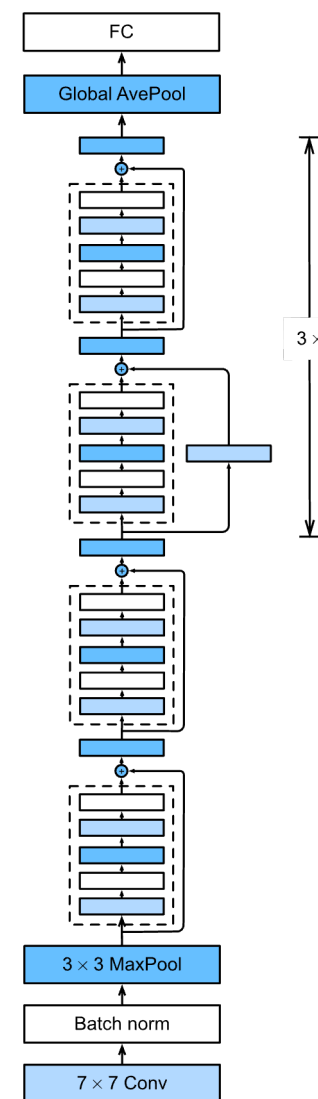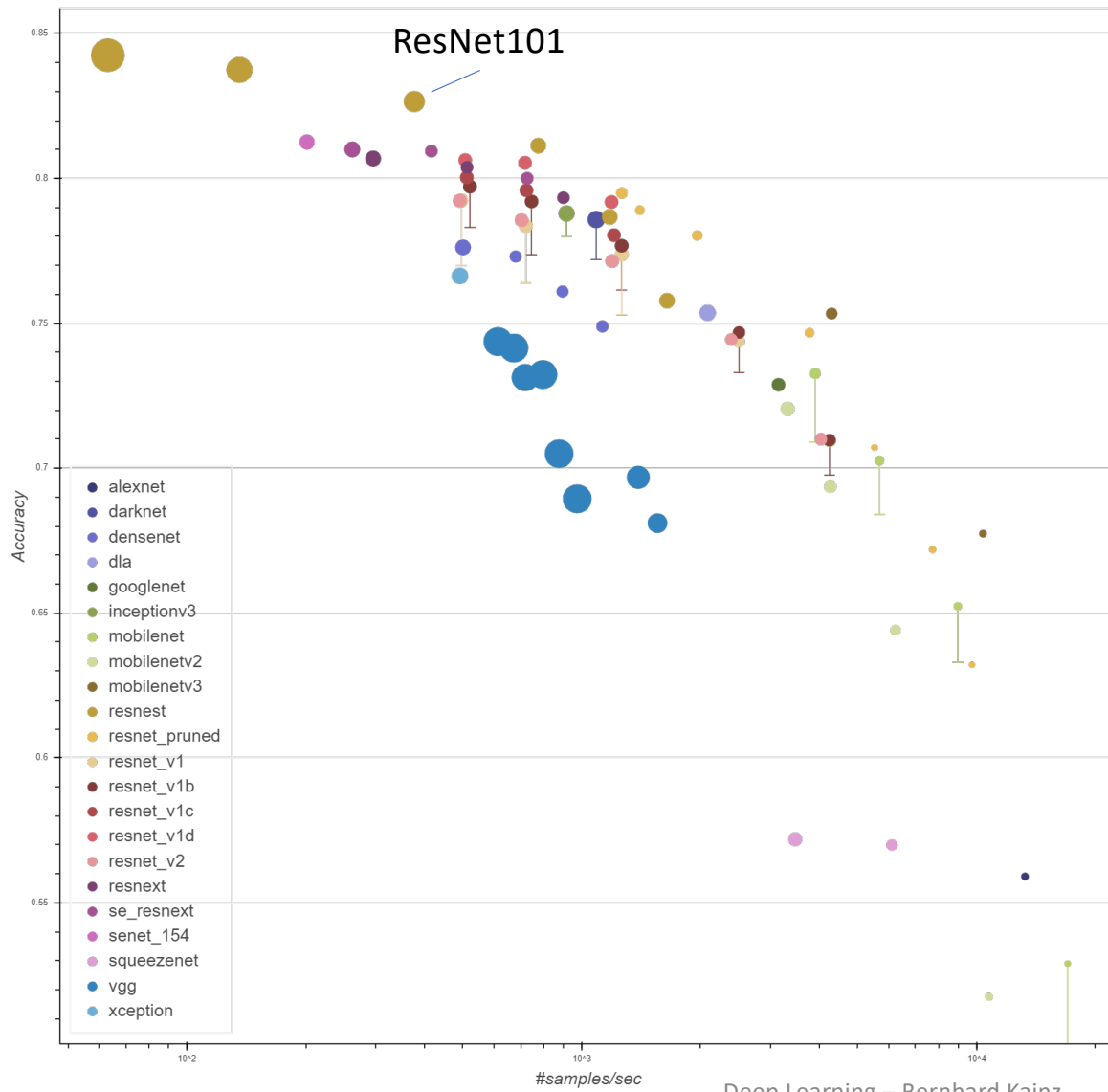


N times

Stride 2

https://d2l.ai/

# ResNet

Same block structure as e.g. VGG or GoogleNet

- Residual connection to add to expressiveness
- Pooling/stride for dimensionality reduction
- Batch Normalization for capacity control

- Trainable at scale
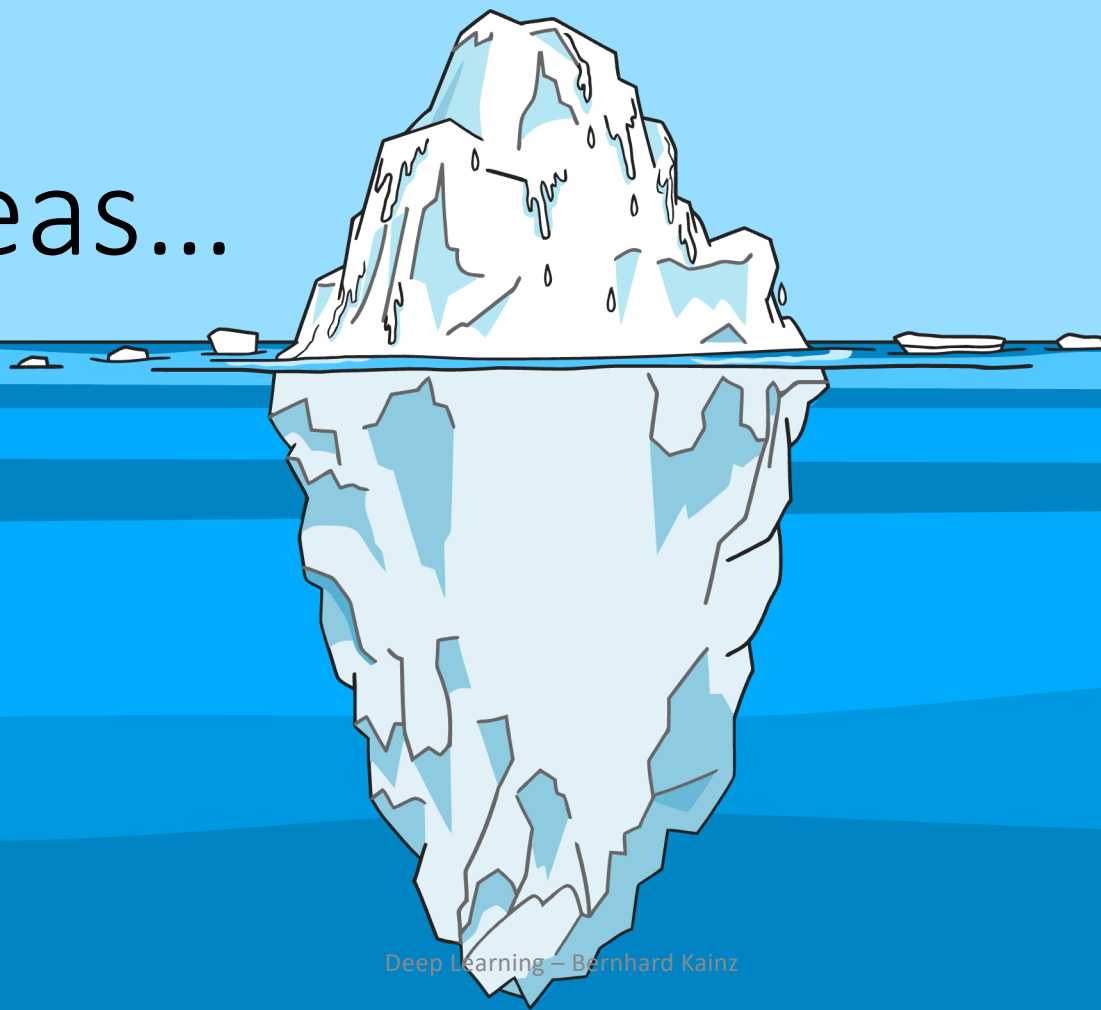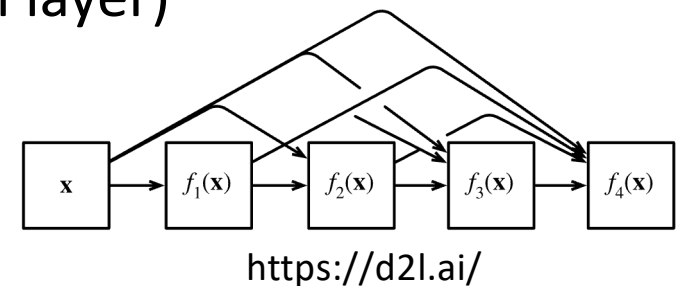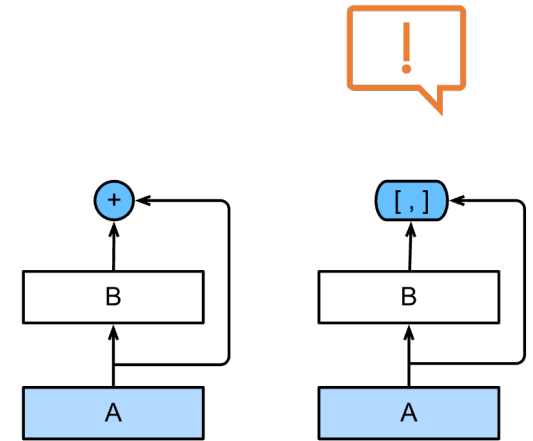- Variant name depends on how many blocks (18 layers = ResNet-18 -> )

ResNet101

https://cv.gluon.ai/model_zoo/classification.html

Deep Learning – Bernhard Kainz

More ideas...

Deep Learning – Bernhard Kainz

# DenseNet

- Huang et al., 2016 https://arxiv.org/abs/1608.06993
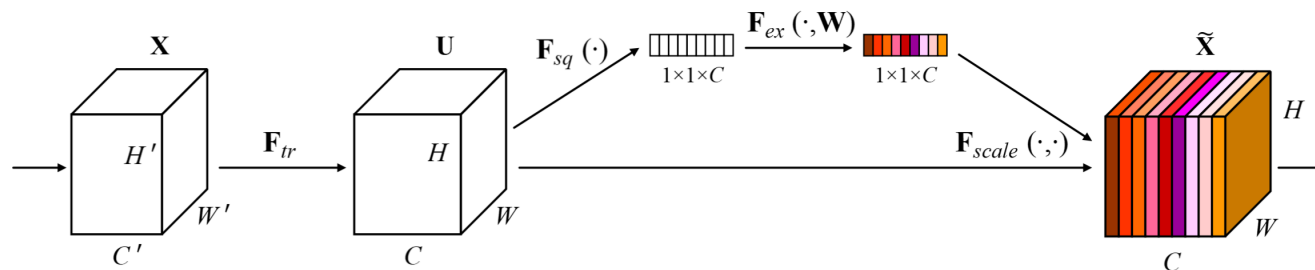
- ResNet combines $x$ and $f(x)$

- DenseNet uses higher order 'Taylor series' expansion

$$x_{i+1} = [x_i, f_i(x_i)]$$
$$x_1 = x$$
$$x_2 = [x, f_1(x)]$$
$$x_3 = [x, f_1(x), f_2([x, f_1(x)])]$$

- Occasionally need to reduce resolution (transition layer)

https://d2l.ai/

# Squeeze-Excite Net

- Hu et al., 2017 https://arxiv.org/abs/1709.01507



- Learn global weighting function per channel
- Allows for fast information transfer between pixels in different locations of the image

# Things to explore

- AutoML (find best model architecture automatically Google Cloud AutoML)

- Hypernetworks (a network that proposes the weights for another network), also neural processes

- Networks with memory, e.g. kanerva machine

- Almost no new basic architectures accepted nowadays (see https://nips.cc/virtual/2020/public/cal_main.html NeurIPS 2020 programme, focuses on meta findings)

- Attention! (second part of the course)

# Summary

- **Inception**
  - Inhomogeneous mix of convolutions (varying depth)
  - Batch norm regularization

- **ResNet**
  - Taylor expansion of functions
  - ResNext decomposes convolutions

- **Model Zoo**
  - DenseNet, ShuffleNet, Separable Convolutions, …



Deep Learning – Bernhard Kainz

https://in.pinterest.com/pin/556124253981912489/
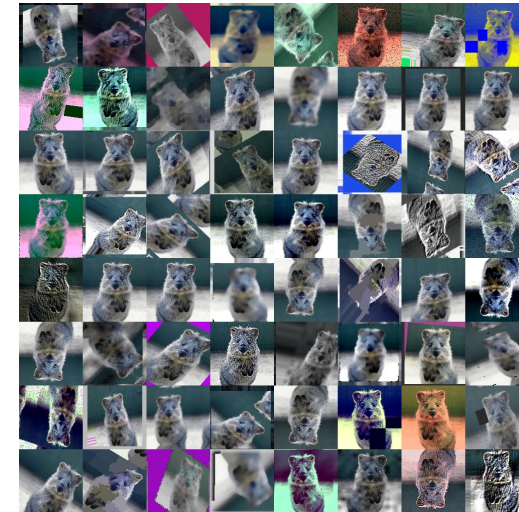
# What do we learn from that

- Deeper is not necessarily better if the function space is not regularised

- ResNet is the workhorse of Deep learning (for now. Do you have a better idea that hasn't been tried yet? Let me know but look on arXiv first!)

- Lot's of variations have been proposed but it often boils down to how you train a network and for what purpose.

# Data Augmentation

# Input augmentation

- Artificially inflate training data size through applying expected transformations during training
- https://github.com/aleju/imgaug
- https://pytorch.org/docs/stable/torchvision/transforms.html
- Excellent regularizer against overfitting



Deep Learning – Bernhard Kainz

# Transformations

- Random
  - flipping
  - scaling
  - rotations
  - intensity/contrast variations
  - cropping/padding
  - noise
  - affine transformations
  - perspective transformations