# MLDL-I
# Machine Learning and Deep Learning - I

# ImageNet

- 1.2 M Training Data*

- 50K Validation Data*

- 100K Test Data*

- 1000 Classes*

*indicates the 2010 Competition

# Transfer Learning



```python
inp = Input(shape=input_shape)
base_model = tf.keras.applications.VGG16(weights='imagenet',
                                         include_top=False, # drop classifer head
                                         input_shape=input_shape)

x = base_model(inp)
x = Flatten()(x)
x = Dense(128, activation='relu')(x)
out = Dense(num_classes, activation='sigmoid')(x)
```

# CNN Architectures

# AlexNet



- Used ReLU
- Around 60 M Param
- Used 2 GTX 580
- (VRAM 6 GB Total)
- Overlapping Pooling
- Used Dropout

**2012** ImageNet Classification with Deep Convolutional Neural Networks

Link

# VGG



- Simplified Architecture
- 3x3 Conv ,2x2 MaxPool
- Resolution down, Channel up

Link

Convolutional Kernel

Convolutional Kernel

# VGG

Very deep convolutional networks for large-scale image recognition

Link

such layers have a $7 \times 7$ effective receptive field. So what have we gained by using, for instance, a stack of three $3 \times 3$ conv. layers instead of a single $7 \times 7$ layer? First, we incorporate three non-linear rectification layers instead of a single one, which makes the decision function more discriminative. Second, we decrease the number of parameters: assuming that both the input and the output of a three-layer $3 \times 3$ convolution stack has $C$ channels, the stack is parametrised by $3\left(3^2 C^2\right) = 27C^2$ weights; at the same time, a single $7 \times 7$ conv. layer would require $7^2 C^2 = 49C^2$ parameters, i.e. $81\%$ more. This can be seen as imposing a regularisation on the $7 \times 7$ conv. filters, forcing them to have a decomposition through the $3 \times 3$ filters (with non-linearity injected in between).
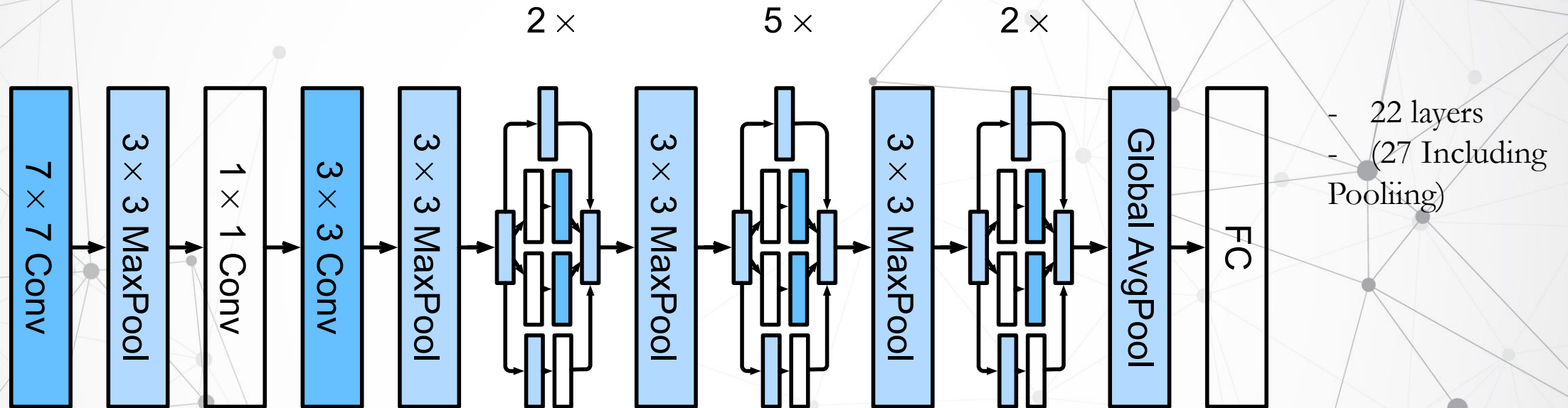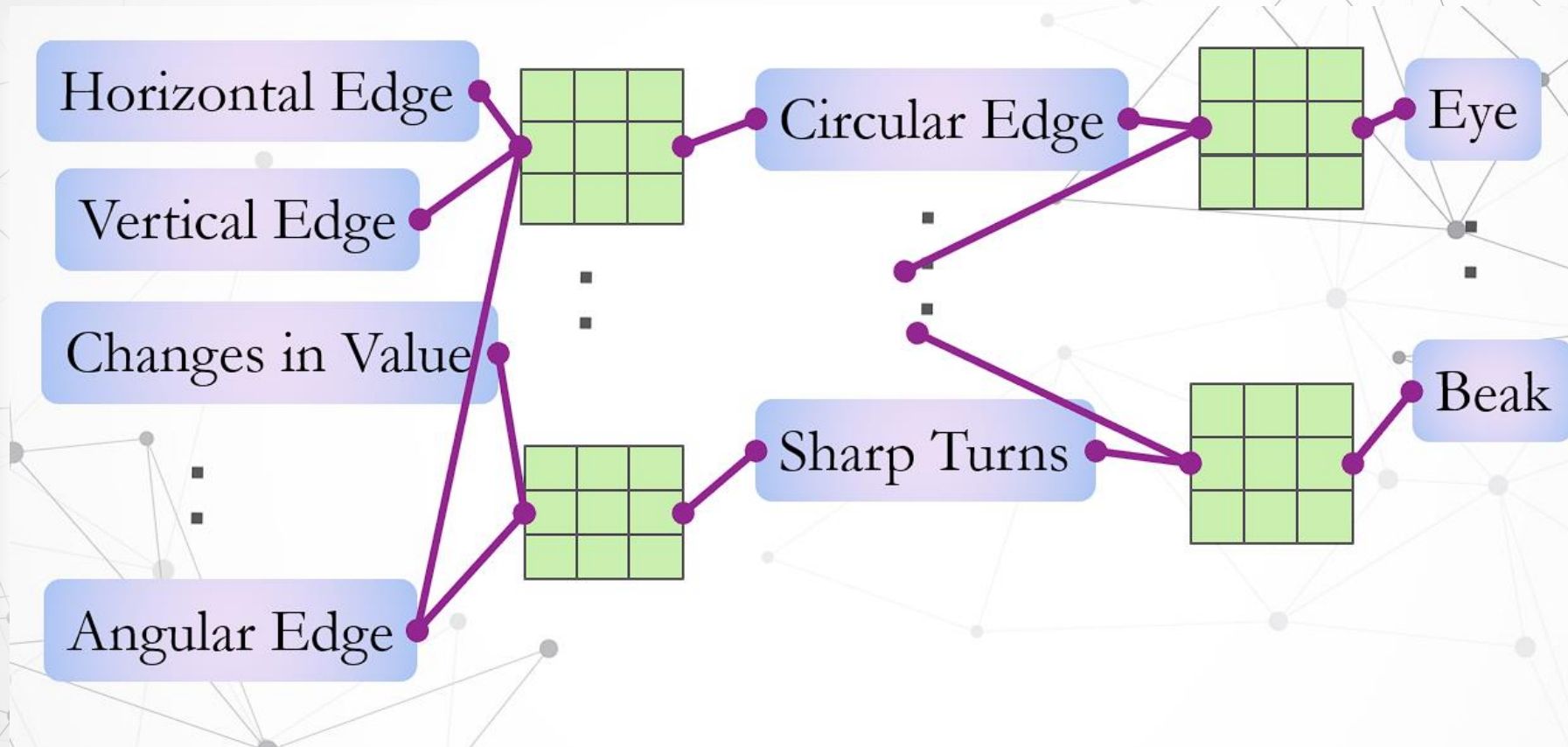
Second, we observe that the classification error decreases with the increased ConvNet depth: from 11 layers in A to 19 layers in E. Notably, in spite of the same depth, the configuration C (which contains three $1 \times 1$ conv. layers), performs worse than the configuration D, which uses $3 \times 3$ conv. layers throughout the network. This indicates that while the additional non-linearity does help (C is better than B), it is also important to capture spatial context by using conv. filters with non-trivial receptive fields (D is better than C). The error rate of our architecture saturates when the depth reaches 19 layers, but even deeper models might be beneficial for larger datasets. We also compared the net B with a shallow net with five $5 \times 5$ conv. layers, which was derived from B by replacing each pair of $3 \times 3$ conv. layers with a single $5 \times 5$ conv. layer (which has the same receptive field as explained in Sect. 2.3). The top-1 error of the shallow net was measured to be $7\%$ higher than that of B (on a center crop), which confirms that a deep net with small filters outperforms a shallow net with larger filters.

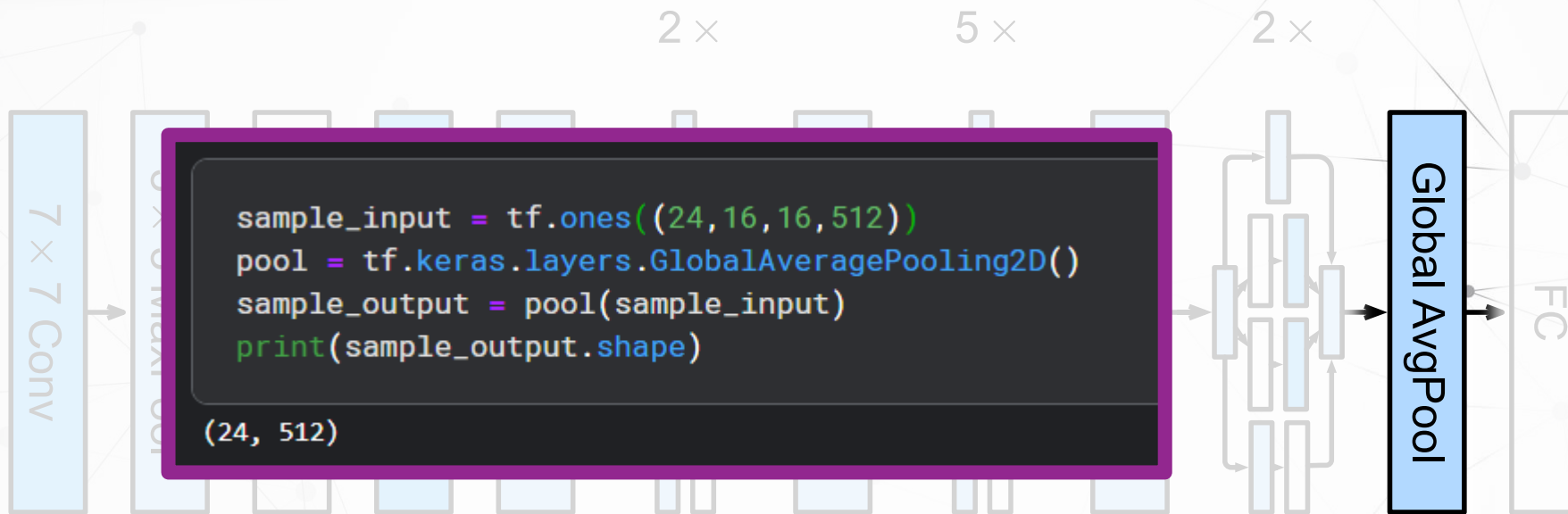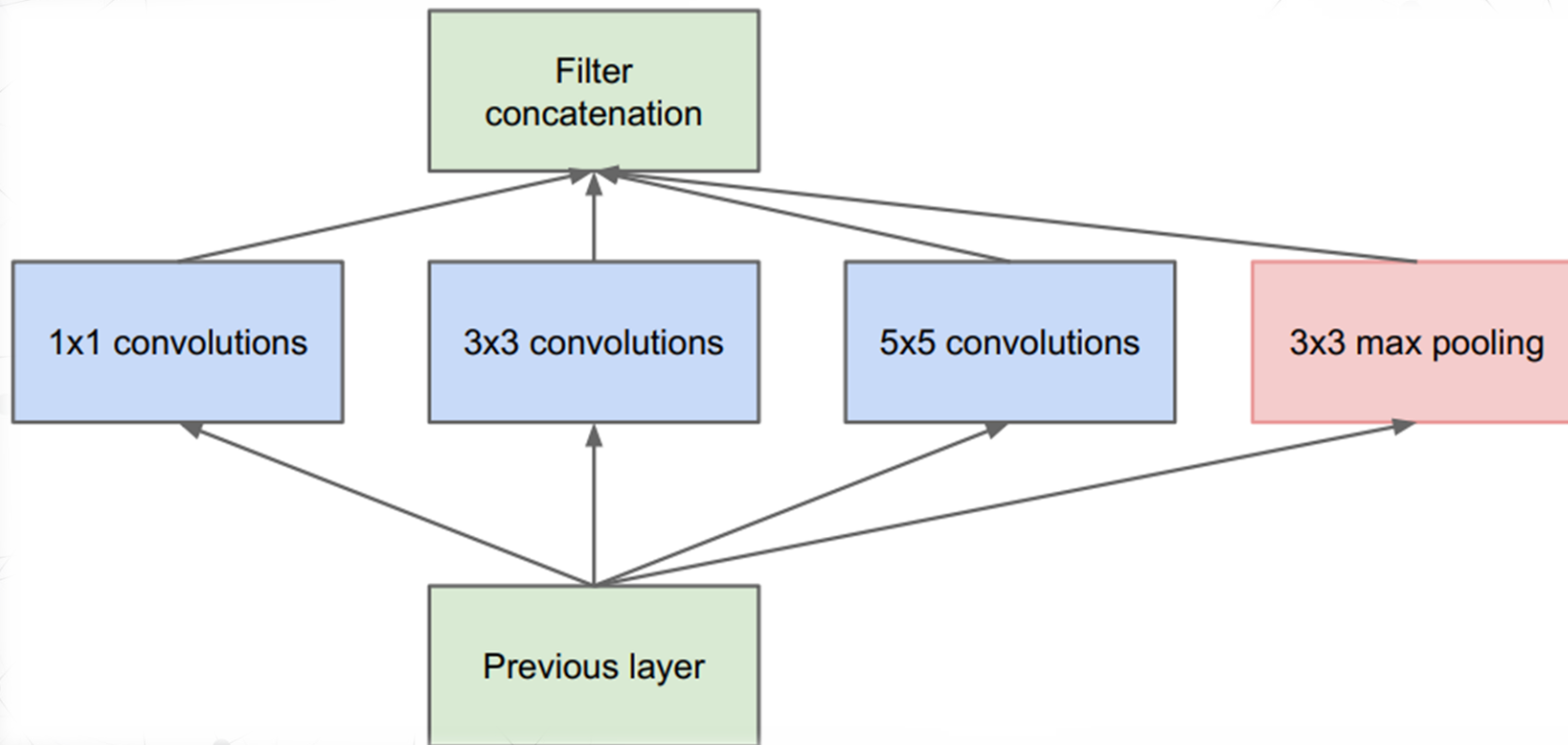2015   Very deep convolutional networks for large-scale image recognition

# GoogLeNet



- 22 layers
- (27 Including Pooling)

2015 Going Deeper with Convolutions

Link

# Global Average Pooling

2015 — Going Deeper with Convolutions

# Global Average Pooling



```python
sample_input = tf.ones((24,16,16,512))
pool = tf.keras.layers.GlobalAveragePooling2D()
sample_output = pool(sample_input)
print(sample_output.shape)
```

```
(24, 512)
```

Global AvgPool

FC

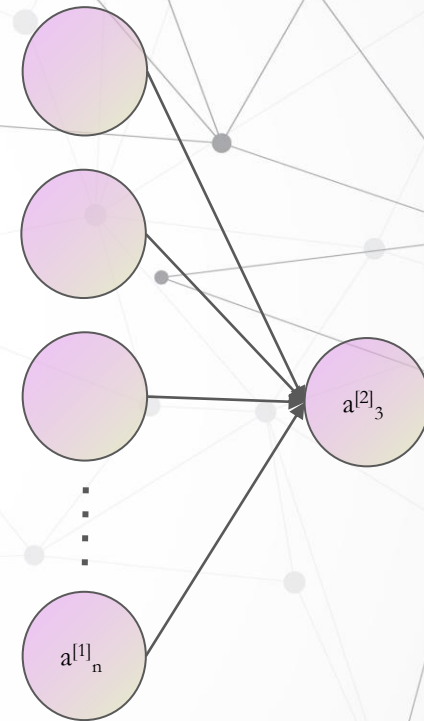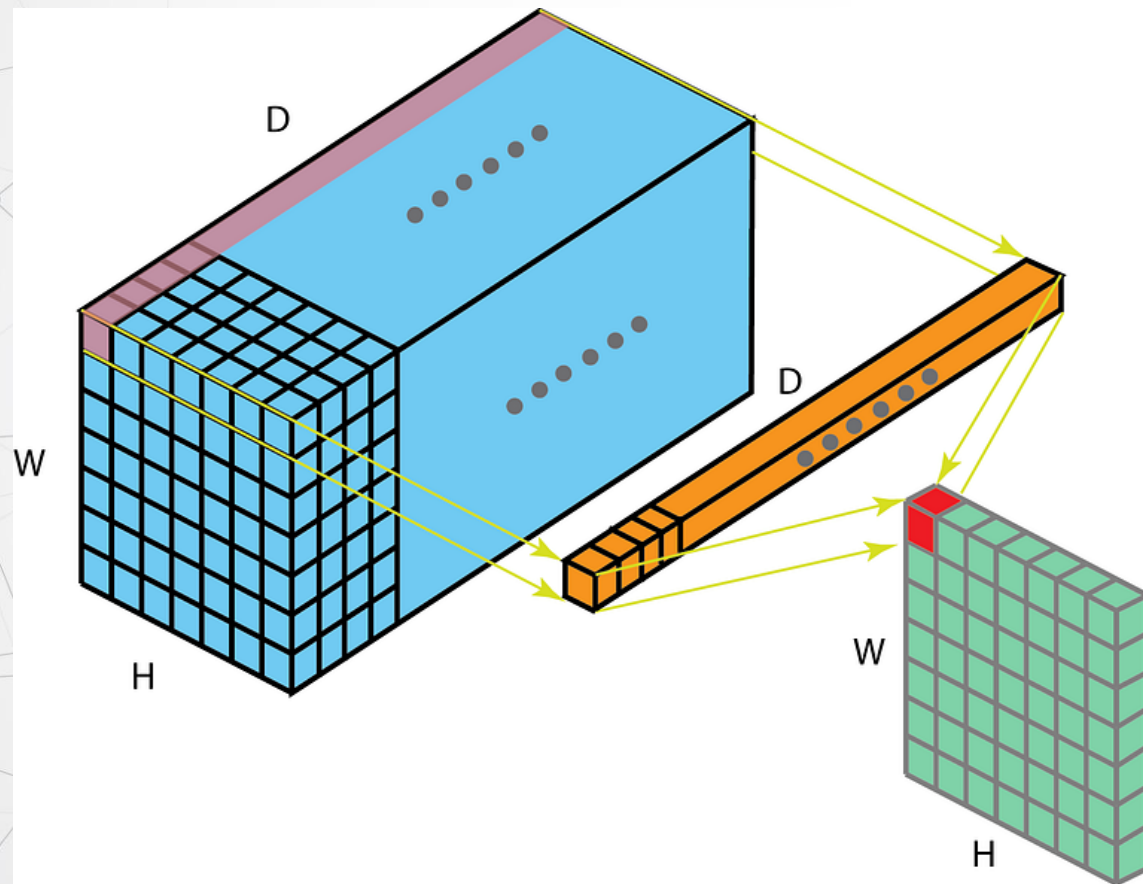7 × 7 Conv

2 ×          5 ×          2 ×

a major effect. We found that a move from fully connected layers to average pooling improved the top-1 accuracy by about 0.6%, however the use of dropout remained essential even after removing the fully connected layers.
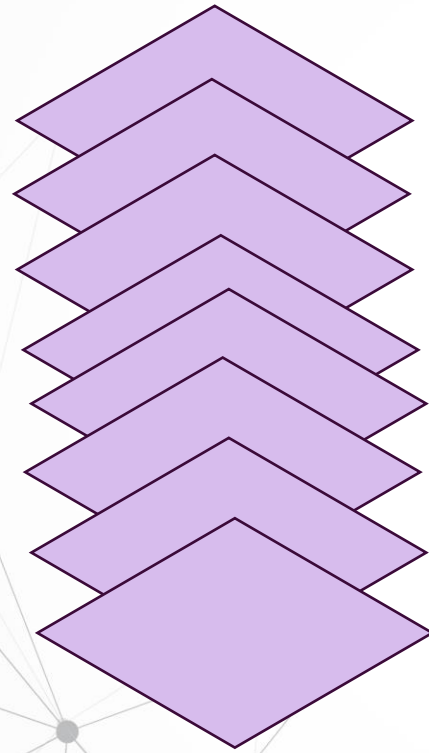
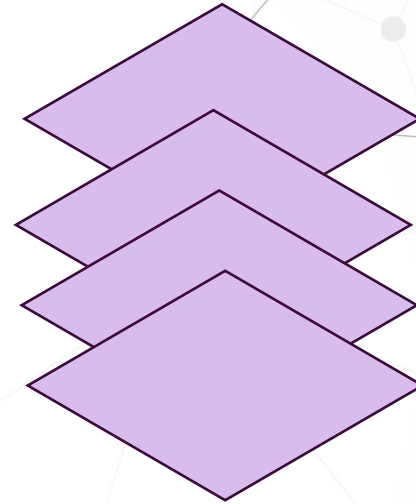2015   Going Deeper with Convolutions

# Inception

# 1x1 Convolutions

# 1x1 Convolutions

254 x 254 x 256

1x1 Conv2D (64)
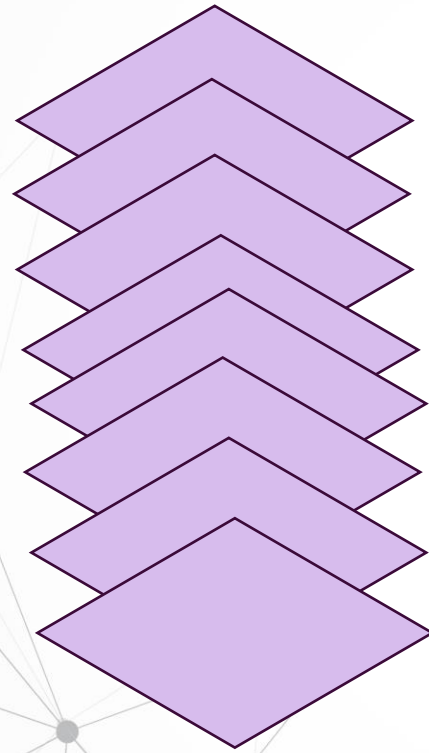
254 x 254 x 64

(256x1)x64 = 16,384

# 1x1 Convolutions

Artificial intelligence (AI) is the intelligence of machines or softwares, as opposed to the intelligence of human beings or animals. AI applications include advanced web search engines (e.g., Google Search), recommendation s[...] Netflix), understanding human speech (such as S[...] Waymo), generative or creative tools (ChatGPT a[...] level in strategic games (such as chess and Go).

Artificial intelligence was founded as an academi[...] has experienced several waves of optimism, follo[...] funding (known as an "AI winter"),followed by r[...] funding. AI research has tried and discarded man[...] the brain, modeling human problem solving, for[...] imitating animal behavior. In the first decades of the 21st century, highly mathematical and statistical machine learning has dominated the field, and this technique has proved highly successful, helping to solve many challenging problems throughout industry and academia.
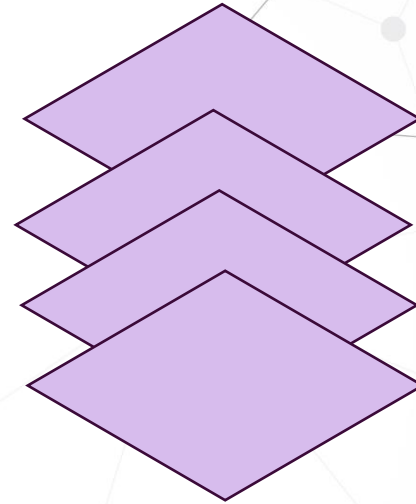
AI is the intelligence of machines/software, used in various applications like search engines, recommendation systems, speech recognition, self-driving cars, creative tools, and strategic game-playing. It has experienced cycles of optimism, disappointment, and renewed funding since its inception in 1956. Different approaches have been explored, with machine learning being dominant in recent years, solving complex problems in various domains
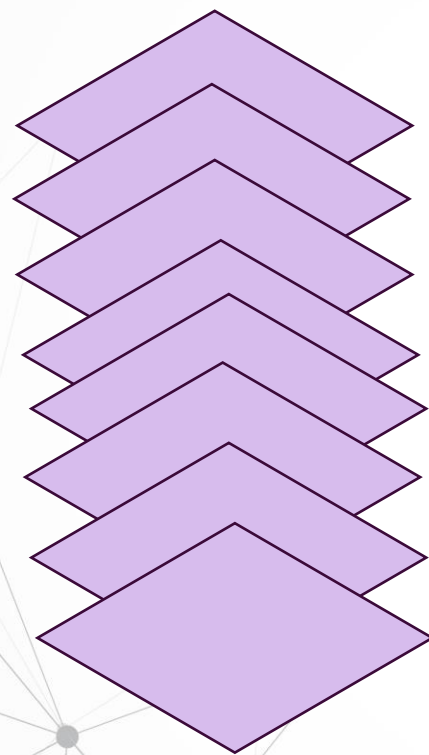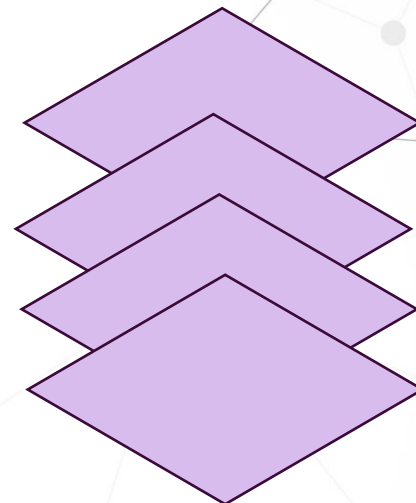
# 1x1 Convolutions



1x1
Conv2D
(64)

254 x 254 x 64

254 x 254 x 256

$(256 \times 1) \times 64 = 16{,}384$

254 x 254 x 256
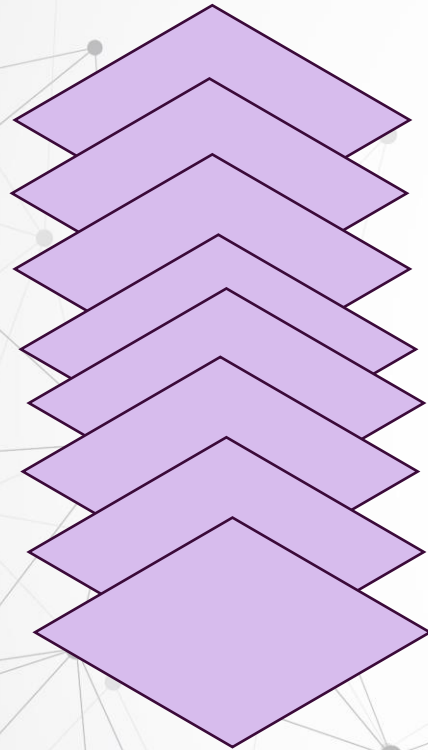
5x5 Conv2D (64)

254 x 254 x 32
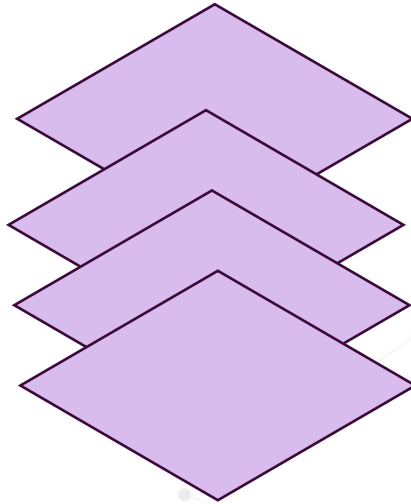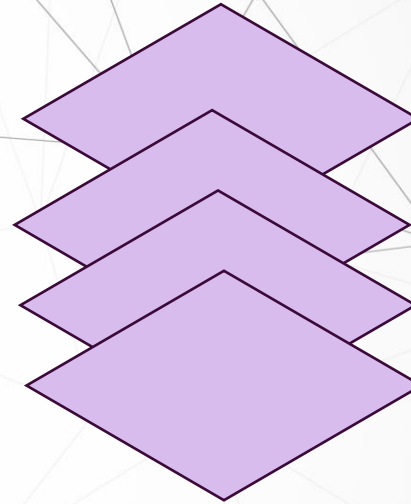
(5x5x256) x 64 = 409,600

# Bottleneck



254 x 254 x 256

1x1
Conv2D
(64)

254 x 254 x 64

5x5
Conv2D
(64)

254 x 254 x 64
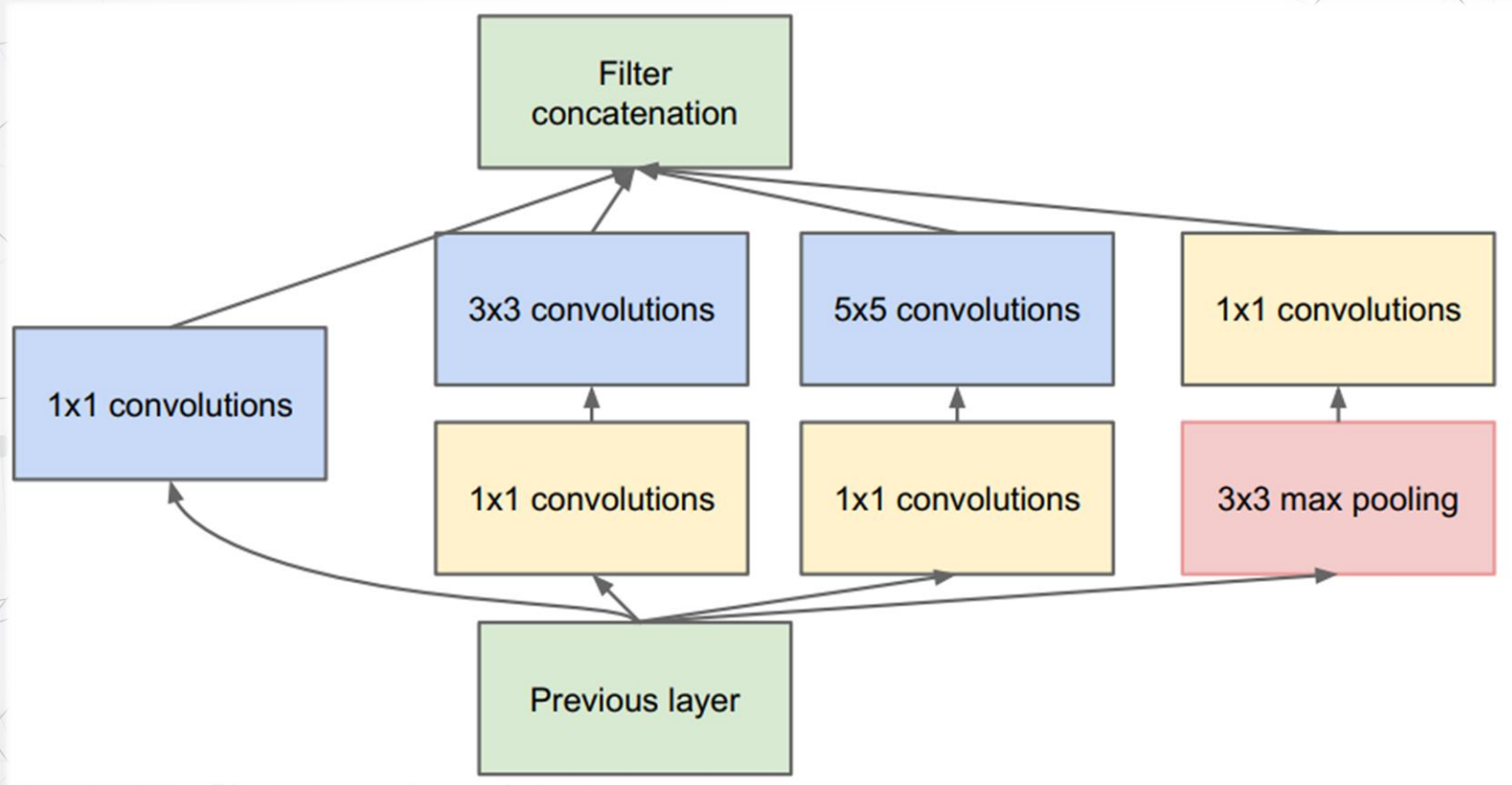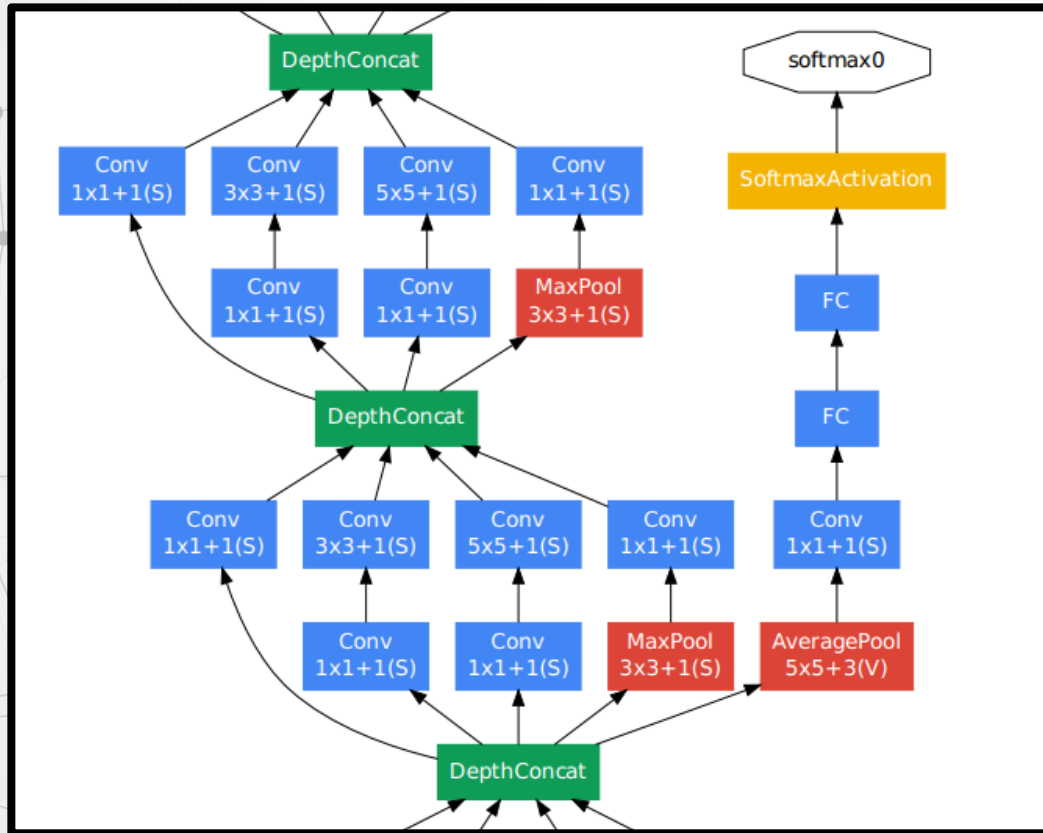
(256x1)x64 + (5x5x64) x 64 = 118,784

(5x5x256) x 64 = 409,600

# GoogLeNet

Going Deeper with Convolutions

Link

# GoogLeNet

The diagram on the left shows the Inception module with:
- DepthConcat (top)
- Conv 1x1+1(S), Conv 3x3+1(S), Conv 5x5+1(S), Conv 1x1+1(S)
- Conv 1x1+1(S), Conv 1x1+1(S), MaxPool 3x3+1(S)
- DepthConcat
- Conv 1x1+1(S), Conv 3x3+1(S), Conv 5x5+1(S), Conv 1x1+1(S), Conv 1x1+1(S)
- Conv 1x1+1(S), Conv 1x1+1(S), MaxPool 3x3+1(S), AveragePool 5x5+3(V)
- DepthConcat
- softmax0, SoftmaxActivation, FC, FC

5 ×   2 ×

3 × 3 MaxPool   Global AvgPool   FC

Given relatively large depth of the network, the ability to propagate gradients back through all the layers in an effective manner was a concern. The strong performance of shallower networks on this task suggests that the features produced by the layers in the middle of the network should be very discriminative. By adding auxiliary classifiers connected to these intermediate layers, discrimination in the lower stages in the classifier was expected. This was thought to combat the vanishing gradient problem while

Link

# ResNet



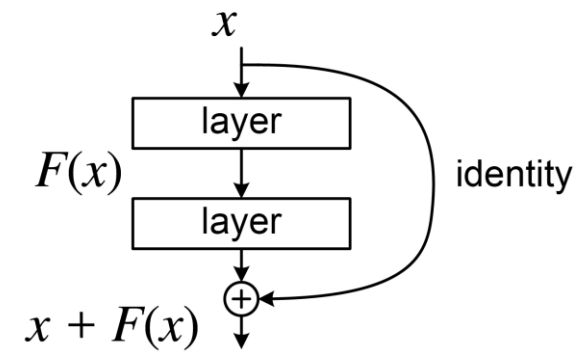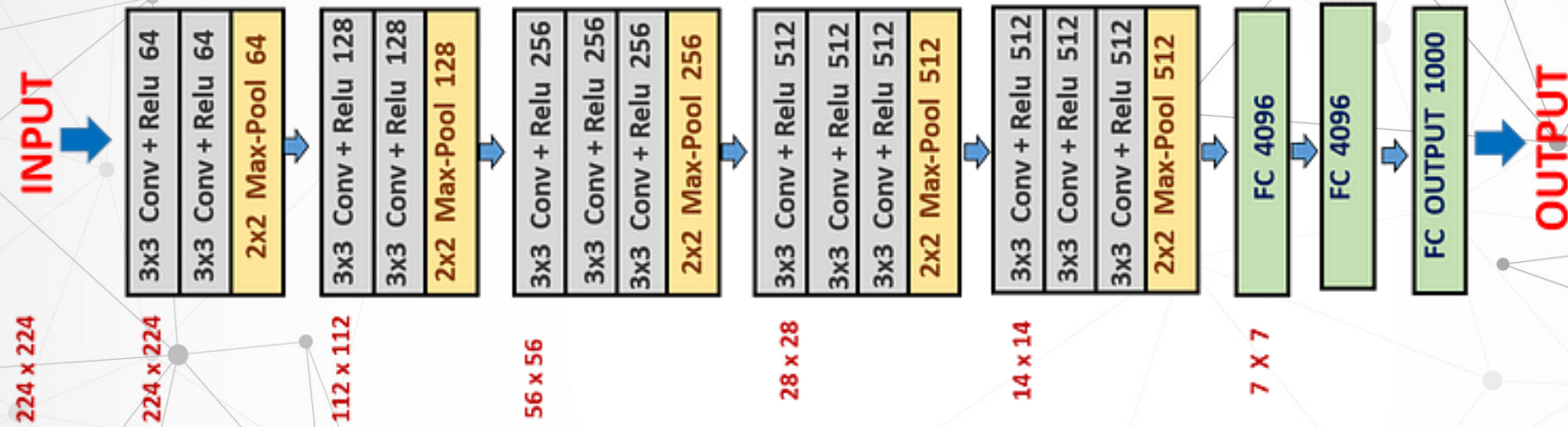Deep Residual Learning for Image Recognition

When deeper networks are able to start converging, a *degradation* problem has been exposed: with the network depth increasing, accuracy gets saturated (which might be unsurprising) and then degrades rapidly. Unexpectedly, such degradation is *not caused by overfitting*, and adding more layers to a suitably deep model leads to *higher training error*, as reported in [10, 41] and thoroughly verified by our experiments. Fig. 1 shows a typical example.
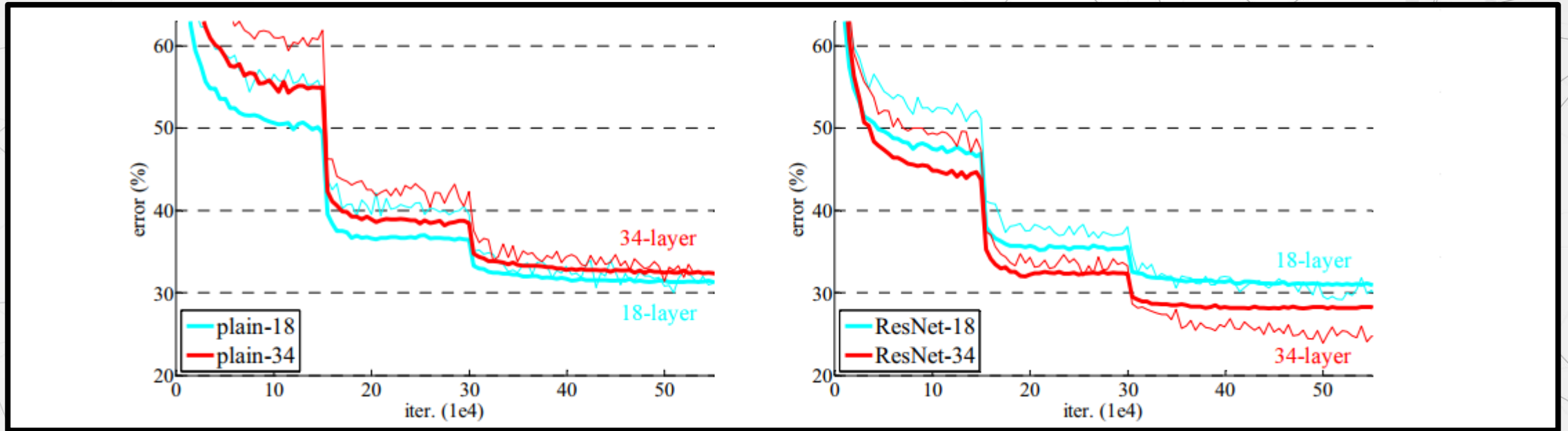
# ResNet

# ResNet
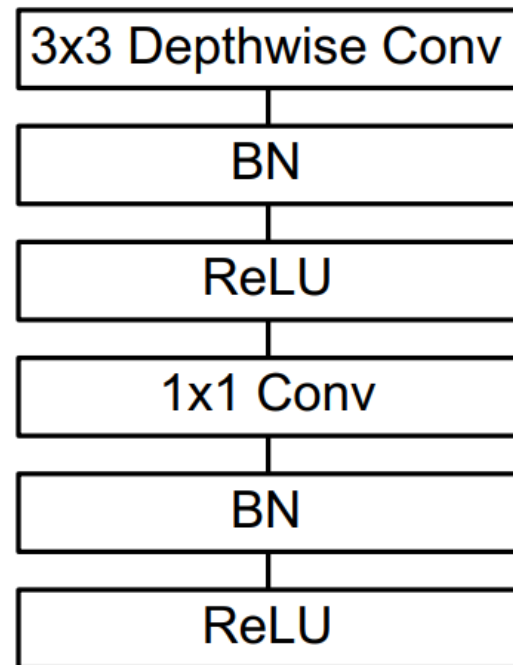


$F(x)$

identity

$x + F(x)$

# ResNet



Deep Residual Learning for Image Recognition

# MobileNet

| 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 |

```python
sample_input = tf.ones((24,16,16,512))
layer = tf.keras.layers.DepthwiseConv2D(
        kernel_size = (3,3),
        strides=(1, 1),
        padding="same")
sample_output = layer(sample_input)
print(sample_output.shape)
```
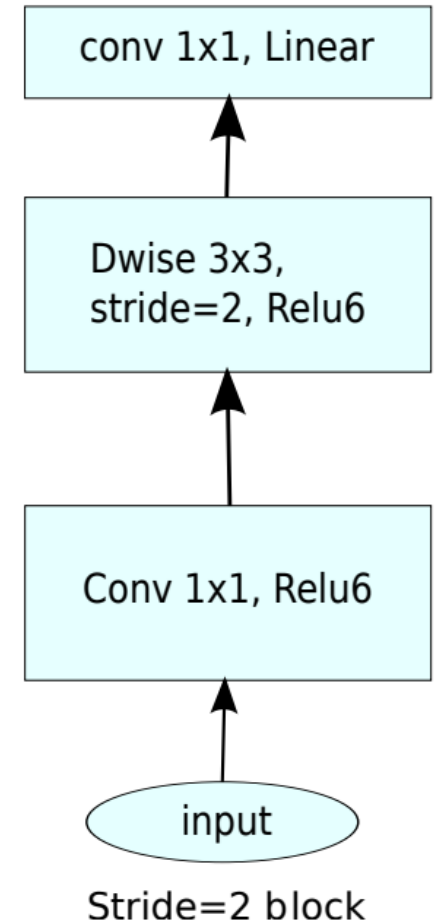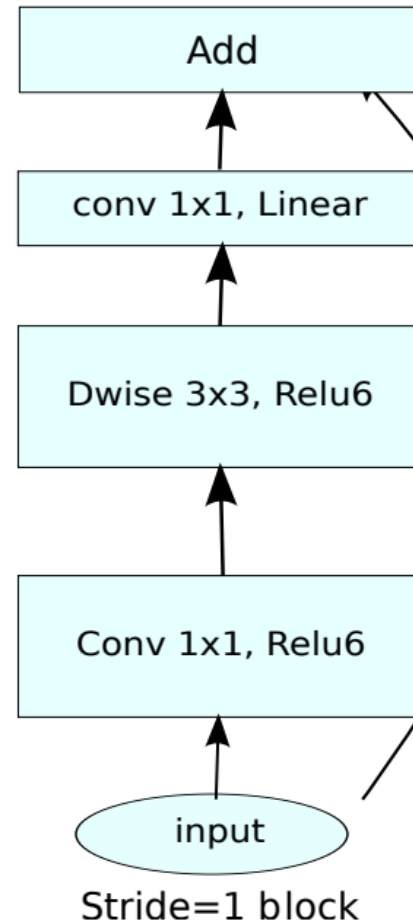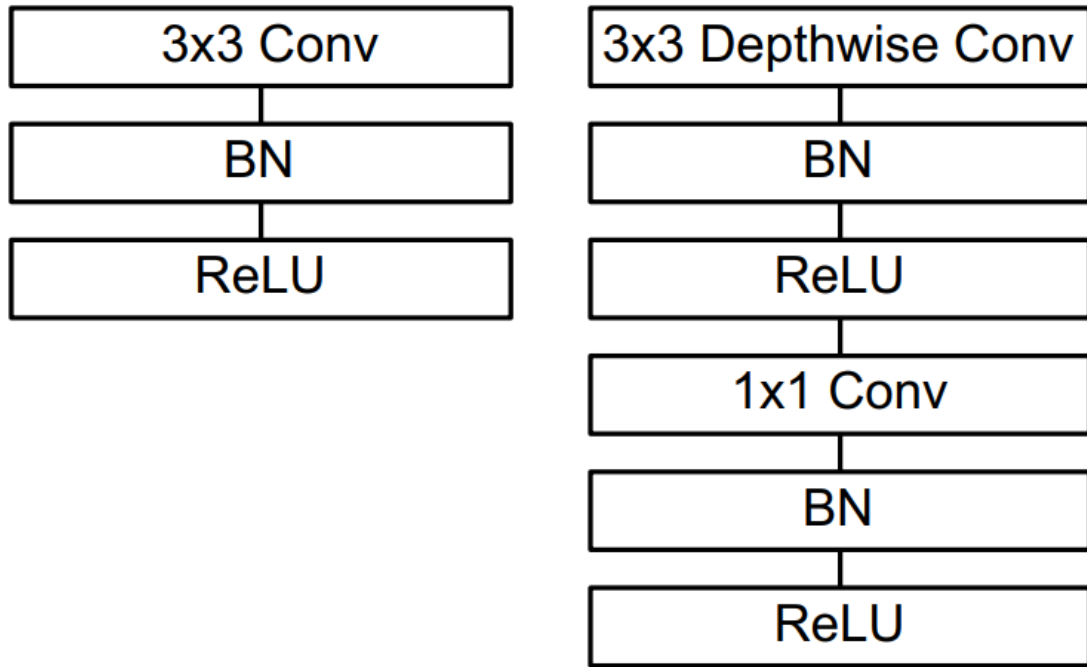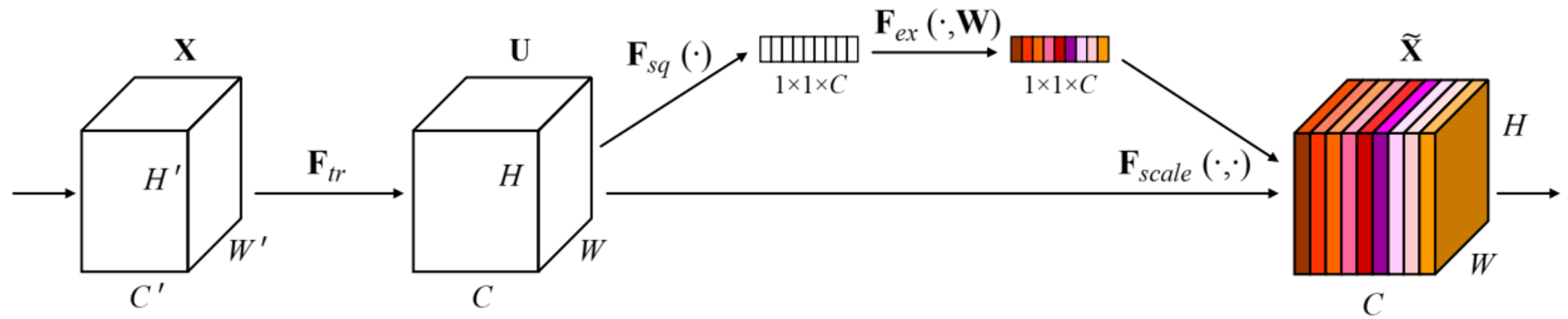
```
(24, 16, 16, 512)
```

```python
sample_input = tf.ones((24,16,16,512))
layer1 = tf.keras.layers.DepthwiseConv2D(
        kernel_size = (3,3),
        strides=(1, 1),
        padding="same")
interm_output = layer1(sample_input)
print("Interm Output shape: ",interm_output.shape)
layer2 = tf.keras.layers.Conv2D(filters = 512,
                                kernel_size = (1,1),
                                padding='same')
final_output = layer2(interm_output)
print("Final Output shape: ",final_output.shape)
```

```
Interm Output shape:  (24, 16, 16, 512)
Final Output shape:  (24, 16, 16, 512)
```

Link

# MobileNetV2



3x3 Conv
BN
ReLU

3x3 Depthwise Conv
BN
ReLU
1x1 Conv
BN
ReLU

Add
conv 1x1, Linear
Dwise 3x3, Relu6
Conv 1x1, Relu6
input

Stride=1 block

conv 1x1, Linear
Dwise 3x3, stride=2, Relu6
Conv 1x1, Relu6
input

Stride=2 block

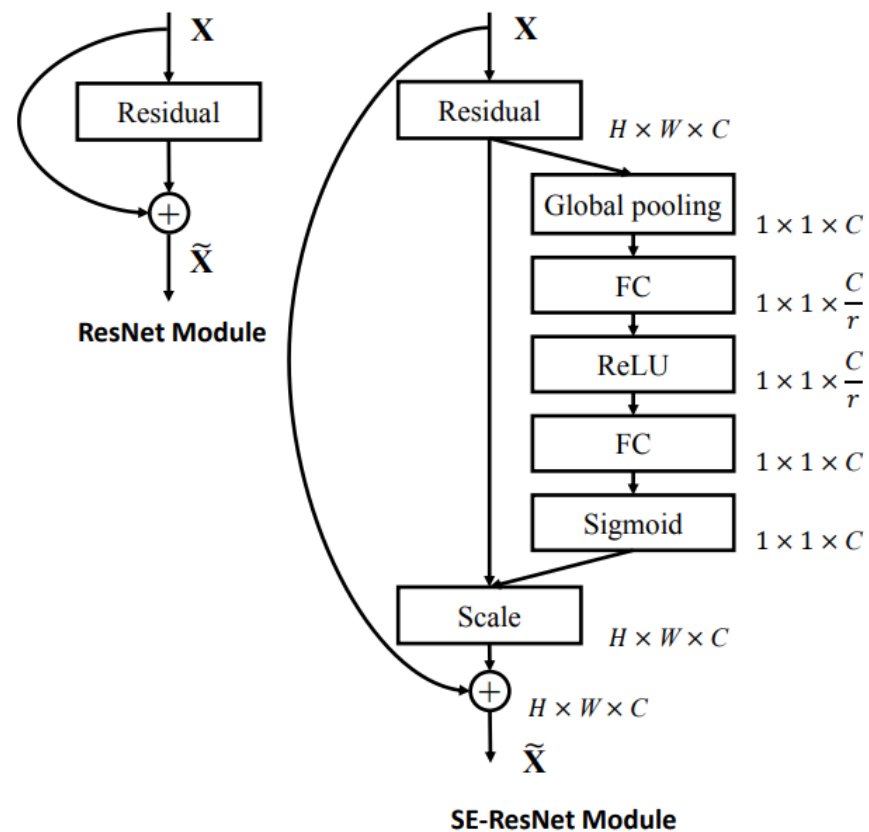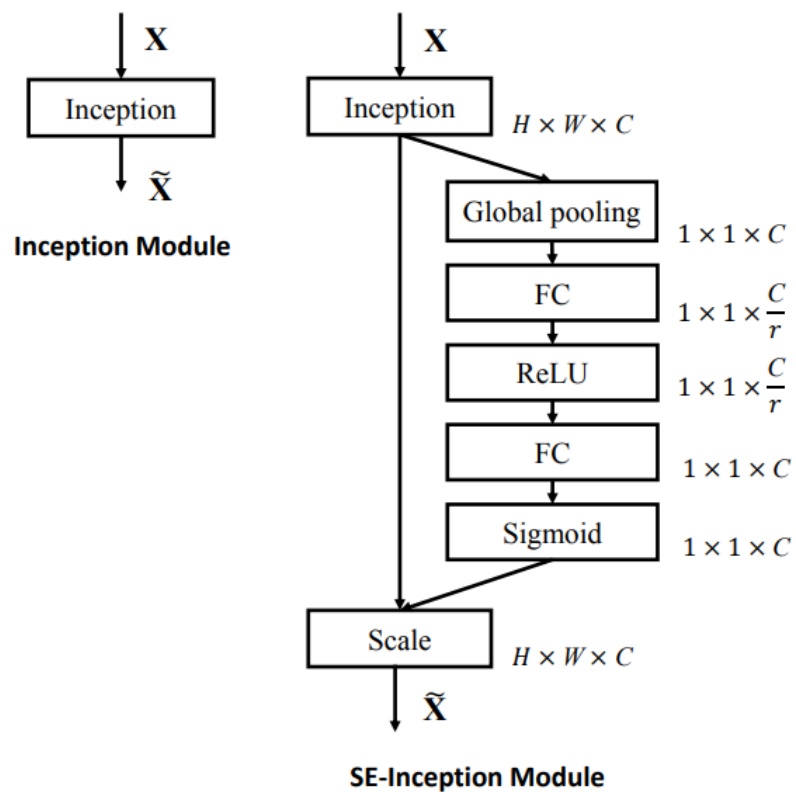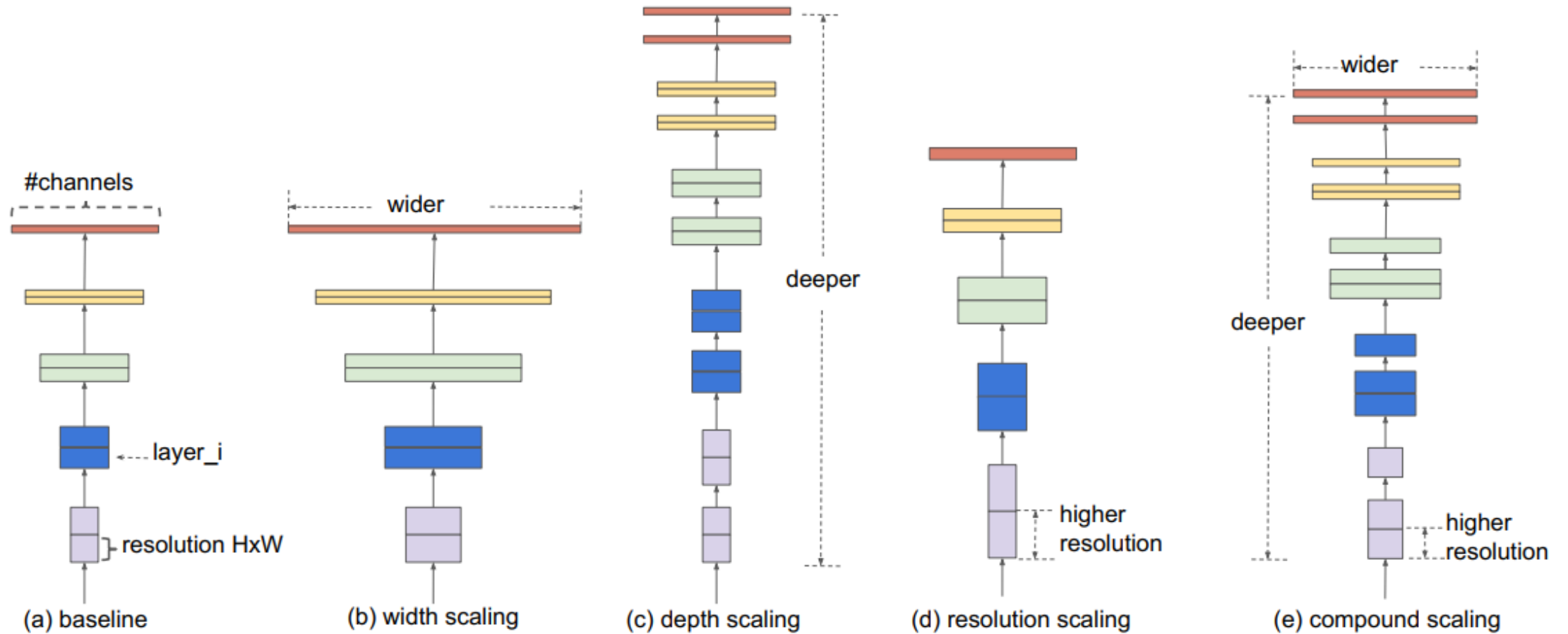MobileNetV2: Inverted Residuals and Linear Bottlenecks

# SENet



```python
sample_input = tf.ones((24,16,16,512))
se = tf.keras.layers.GlobalAveragePooling2D()(sample_input)
se = tf.keras.layers.Dense(512, activation='sigmoid', use_bias=False)(se)
sample_output = tf.keras.layers.multiply([sample_input, se])
print(sample_output.shape)
```

(24, 16, 16, 512)

Squeeze-and-Excitation Networks

# SENet

# EfficientNet



(a) baseline    (b) width scaling    (c) depth scaling    (d) resolution scaling    (e) compound scaling

# EfficientNet

In this paper, we propose a new **compound scaling method**, which use a compound coefficient $\phi$ to uniformly scales network width, depth, and resolution in a principled way:

$$\text{depth: } d = \alpha^\phi$$
$$\text{width: } w = \beta^\phi$$
$$\text{resolution: } r = \gamma^\phi \tag{3}$$
$$\text{s.t. } \alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$
$$\alpha \geq 1, \beta \geq 1, \gamma \geq 1$$

- STEP 1: we first fix $\phi = 1$, assuming twice more resources available, and do a small grid search of $\alpha, \beta, \gamma$ based on Equation 2 and 3. In particular, we find the best values for EfficientNet-B0 are $\alpha = 1.2, \beta = 1.1, \gamma = 1.15$, under constraint of $\alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$.
- STEP 2: we then fix $\alpha, \beta, \gamma$ as constants and scale up baseline network with different $\phi$ using Equation 3, to obtain EfficientNet-B1 to B7 (Details in Table 2).
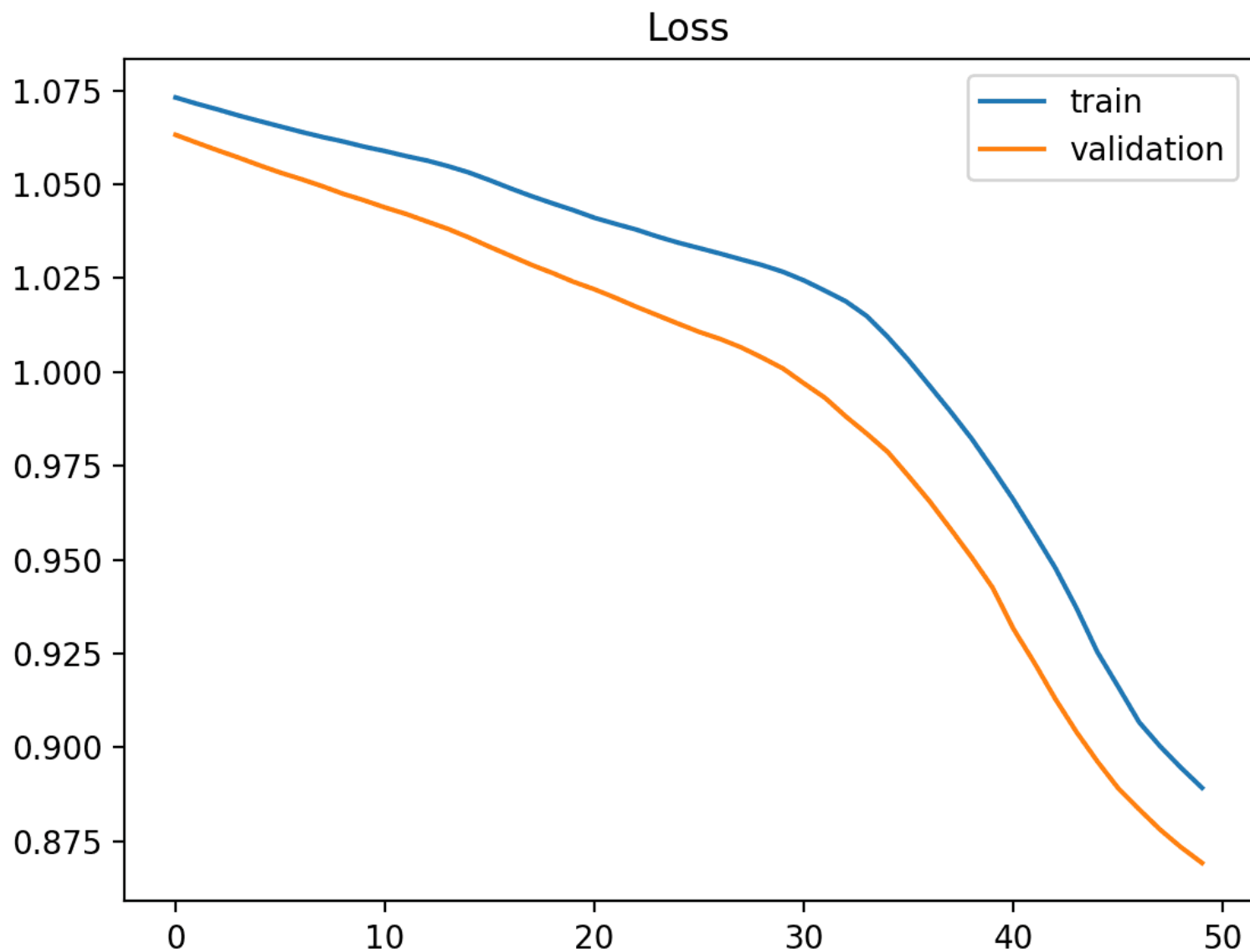
EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks
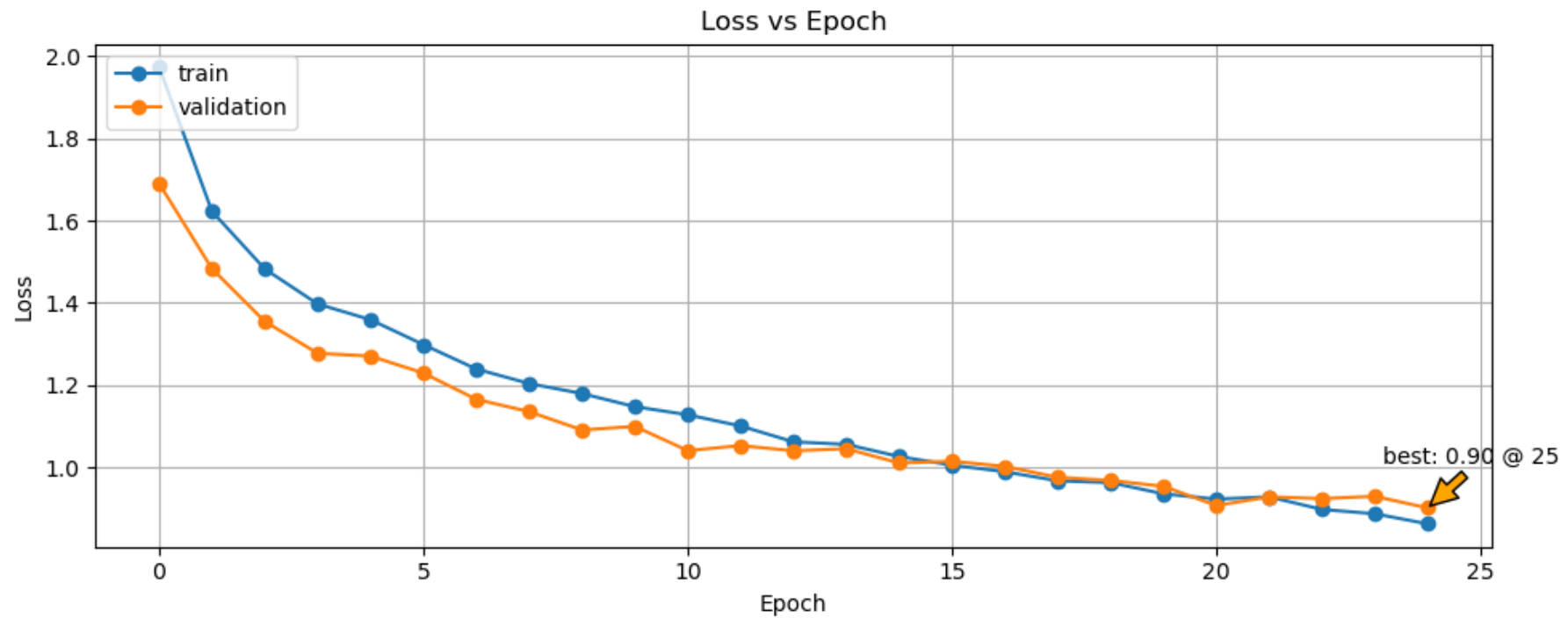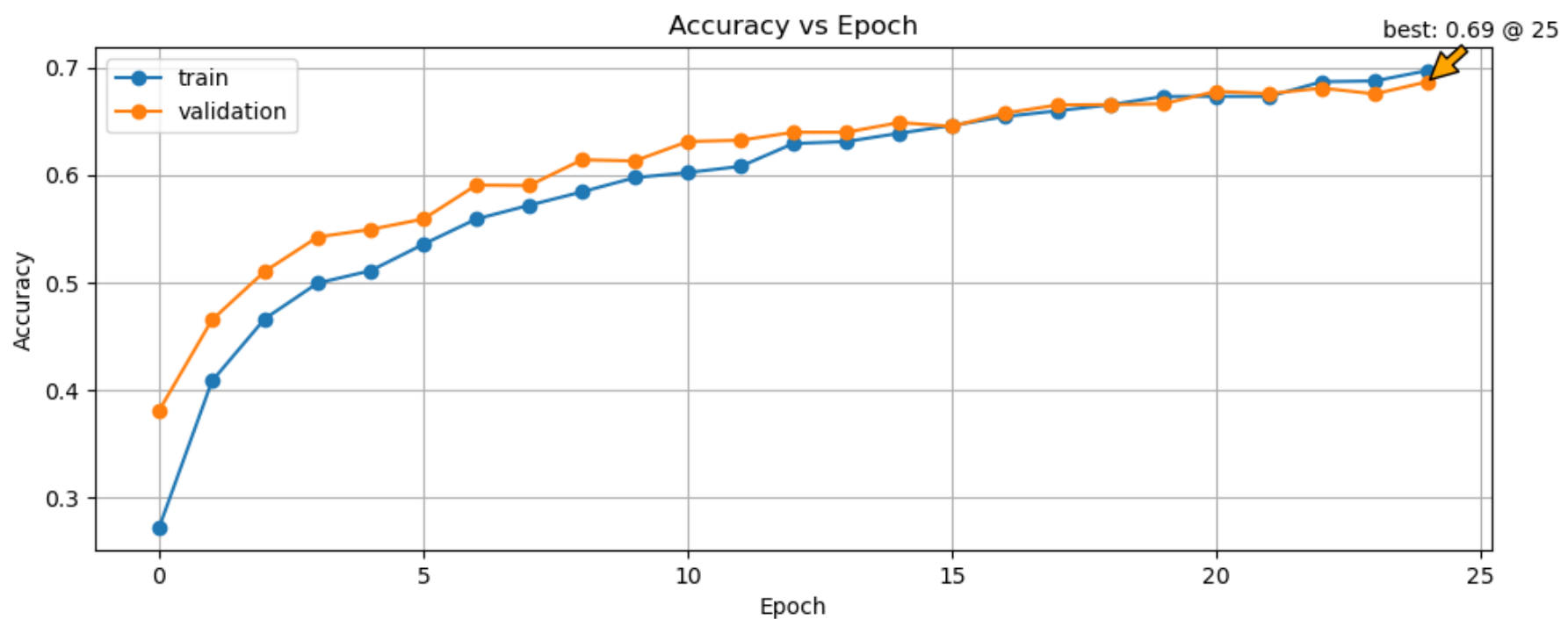
# Convolutional Kernel Shapes

# Underfitting vs Overfitting

# Underfitting vs Overfitting
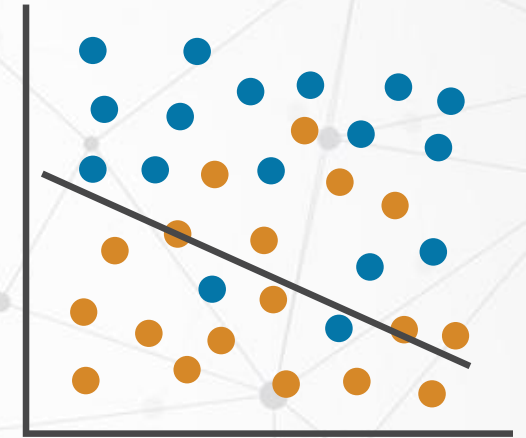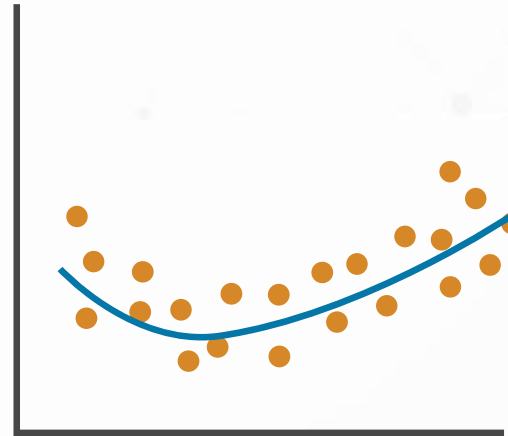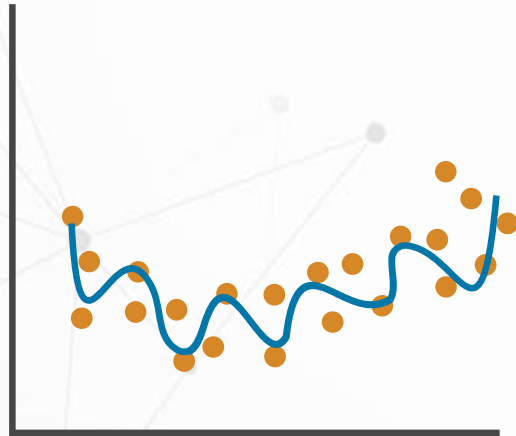


Loss

# Underfitting vs Overfitting

Right Fit

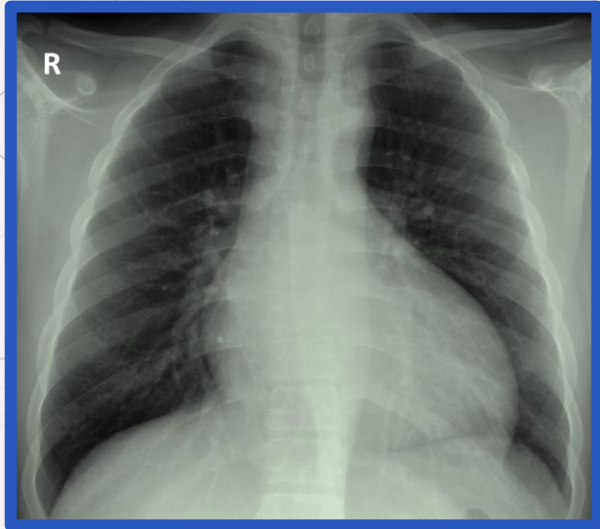Underfitting

Classification

Regression

# How to address Underfitting?

- Increase Model Depth

- Increase Training Epoch/ Increase Learning Rate

- Reduce Data

# How to address Overfitting?

- Increase Training Data

  - Use augmentation

- Reduce Model Depth / Change Architecture
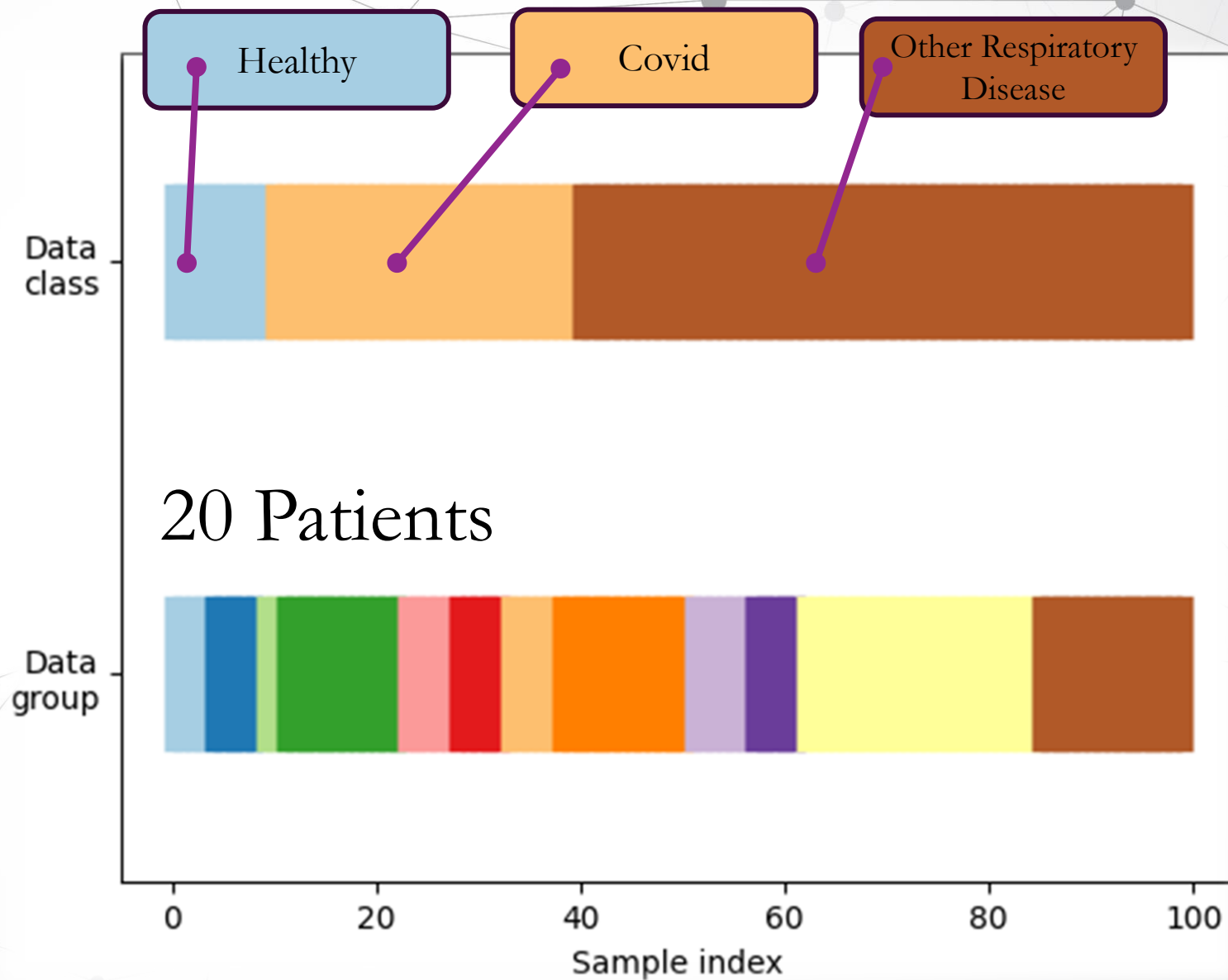
- Use Regularization / Dropout

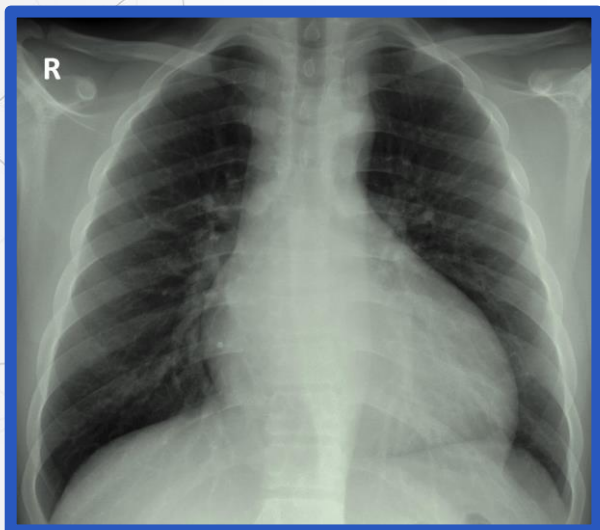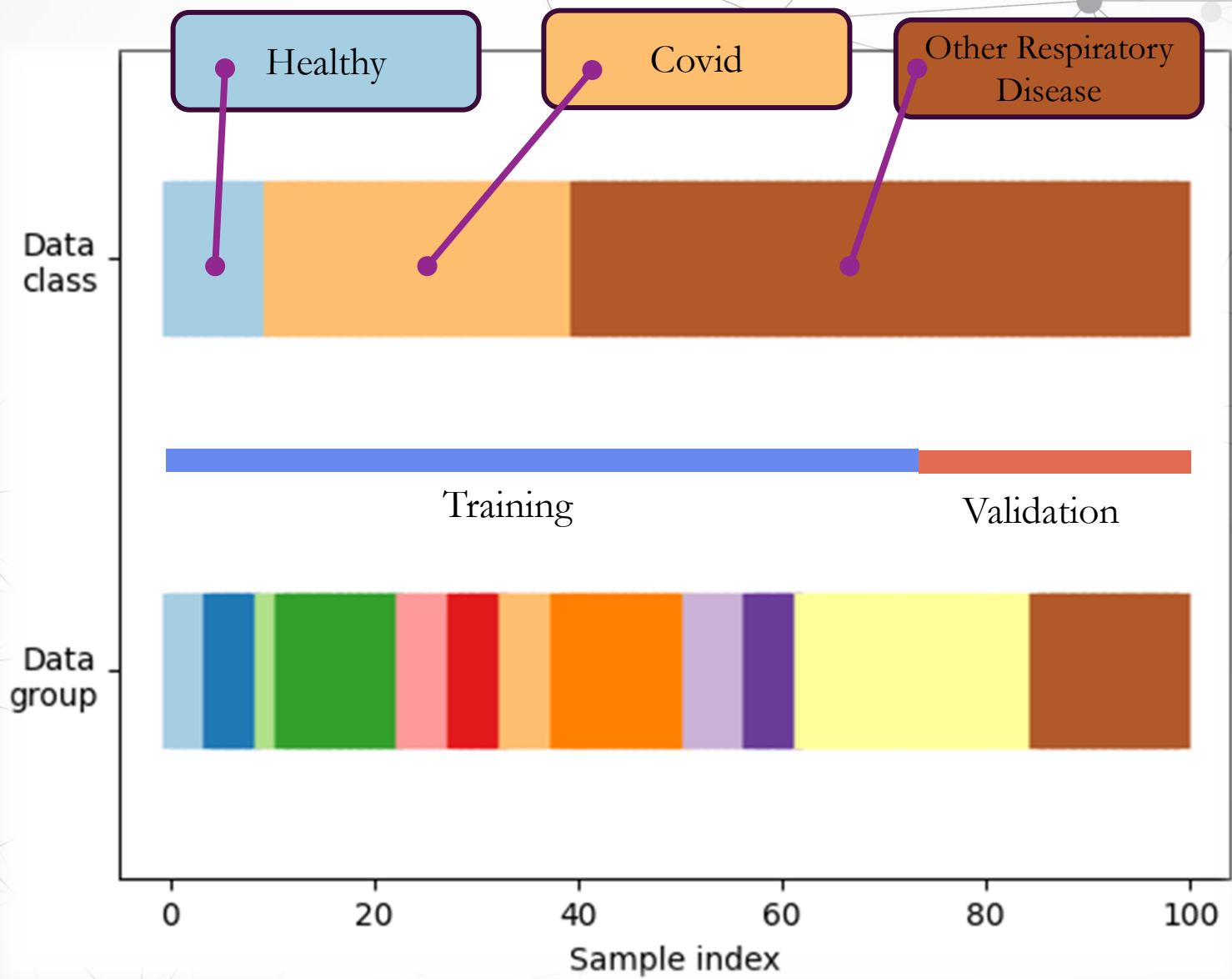*A different distribution of groups compared to the datasets is used for easier visualization

StratifiedGroupKFold

# Up Next: Regularization

# Regularization

$$\mathcal{L}_1 = \sum_{i=1}^{N} |w_i|$$

$$Total\ Cost\ Function = Loss + \mathcal{L}_1$$

$$\mathcal{L}_1 = \sum_{i=1}^{N} |w_i|^2$$

$$Total\ Cost\ Function = Loss + \mathcal{L}_1$$

# Optimizers

```python
tf..keras.optimizers.SGD(learning_rate = 0.01)
```

```python
tf.keras.optimizers.SGD(learning_rate = 0.01,
                        momentum=0.9)
```
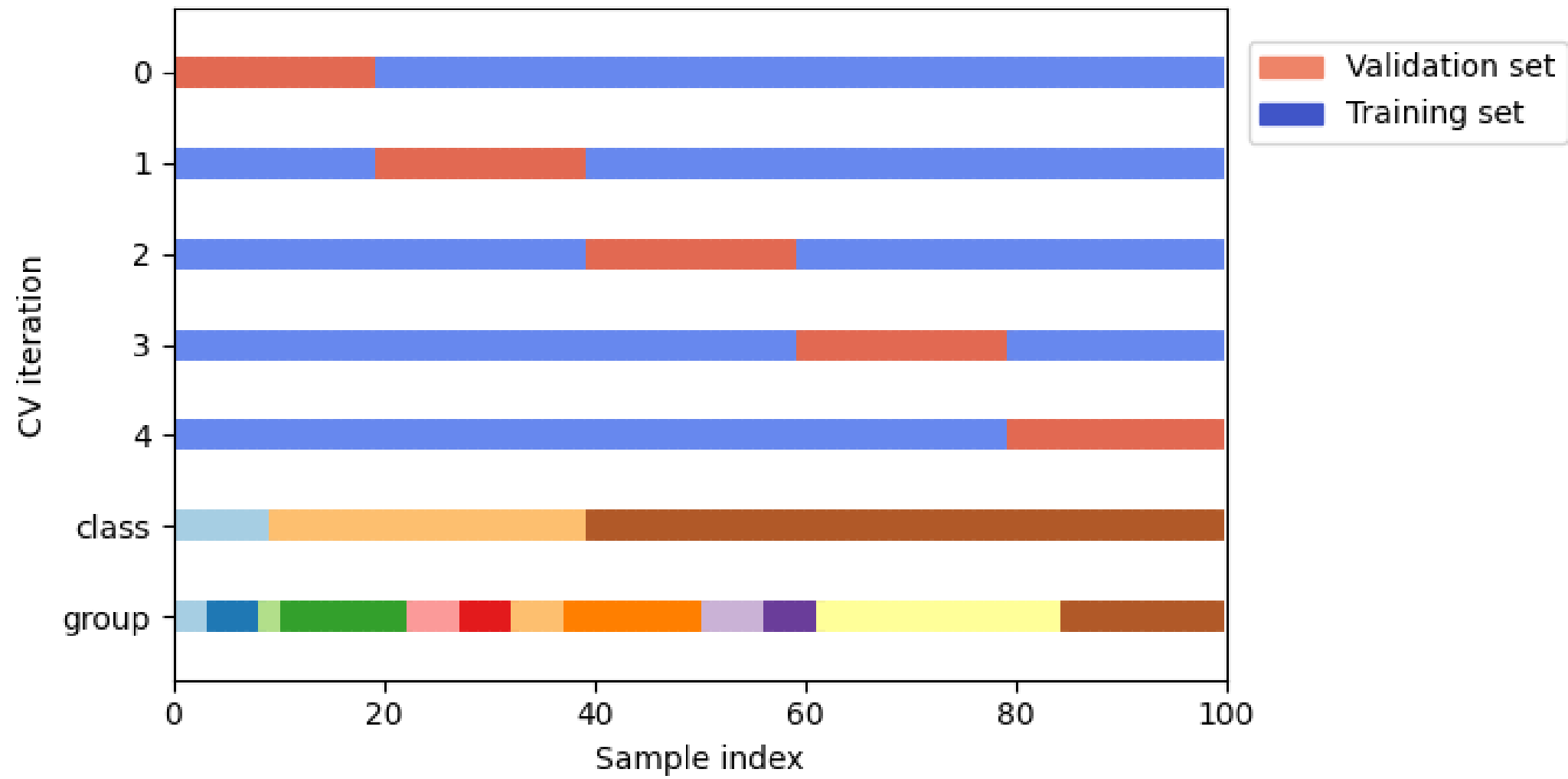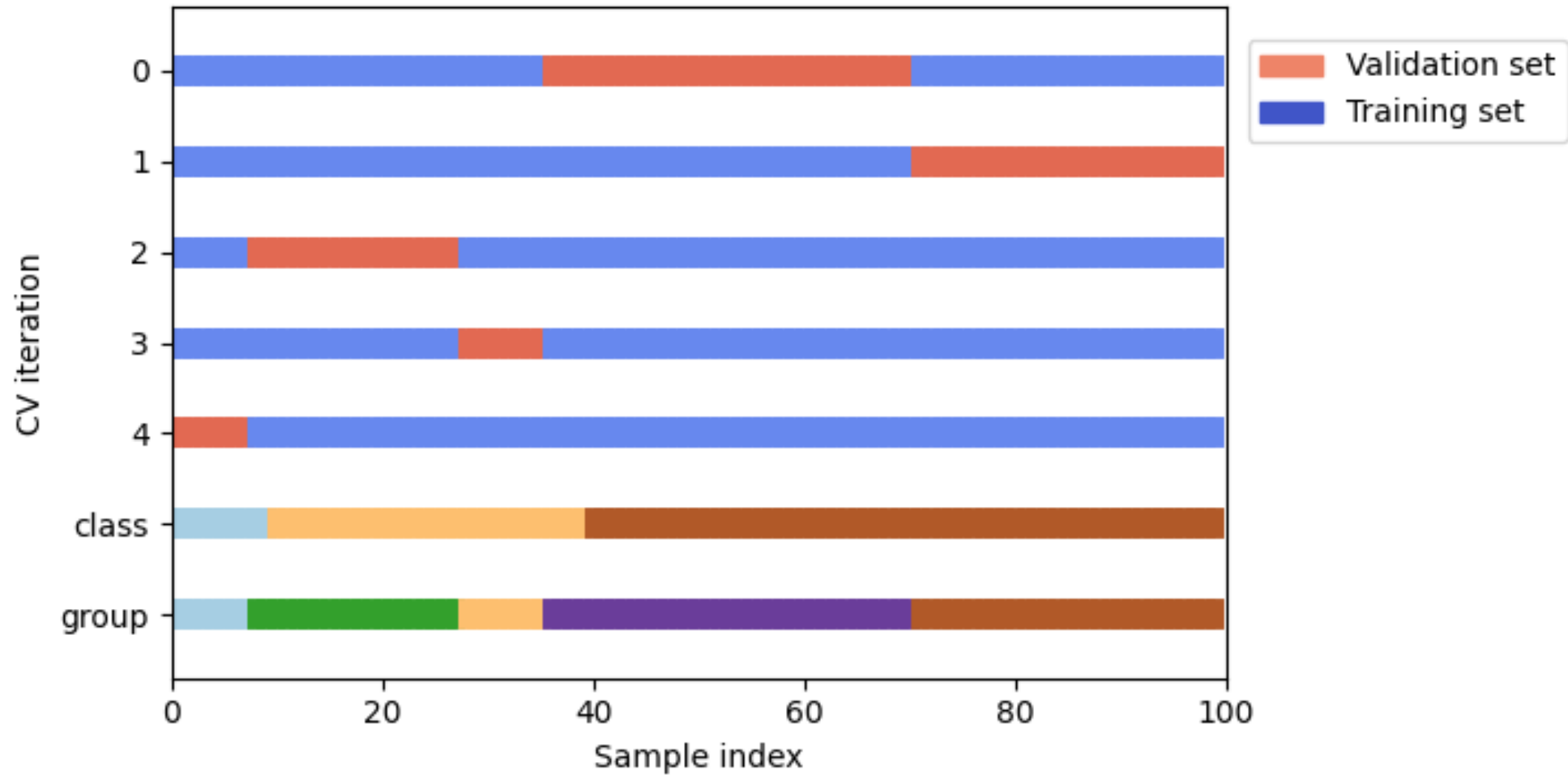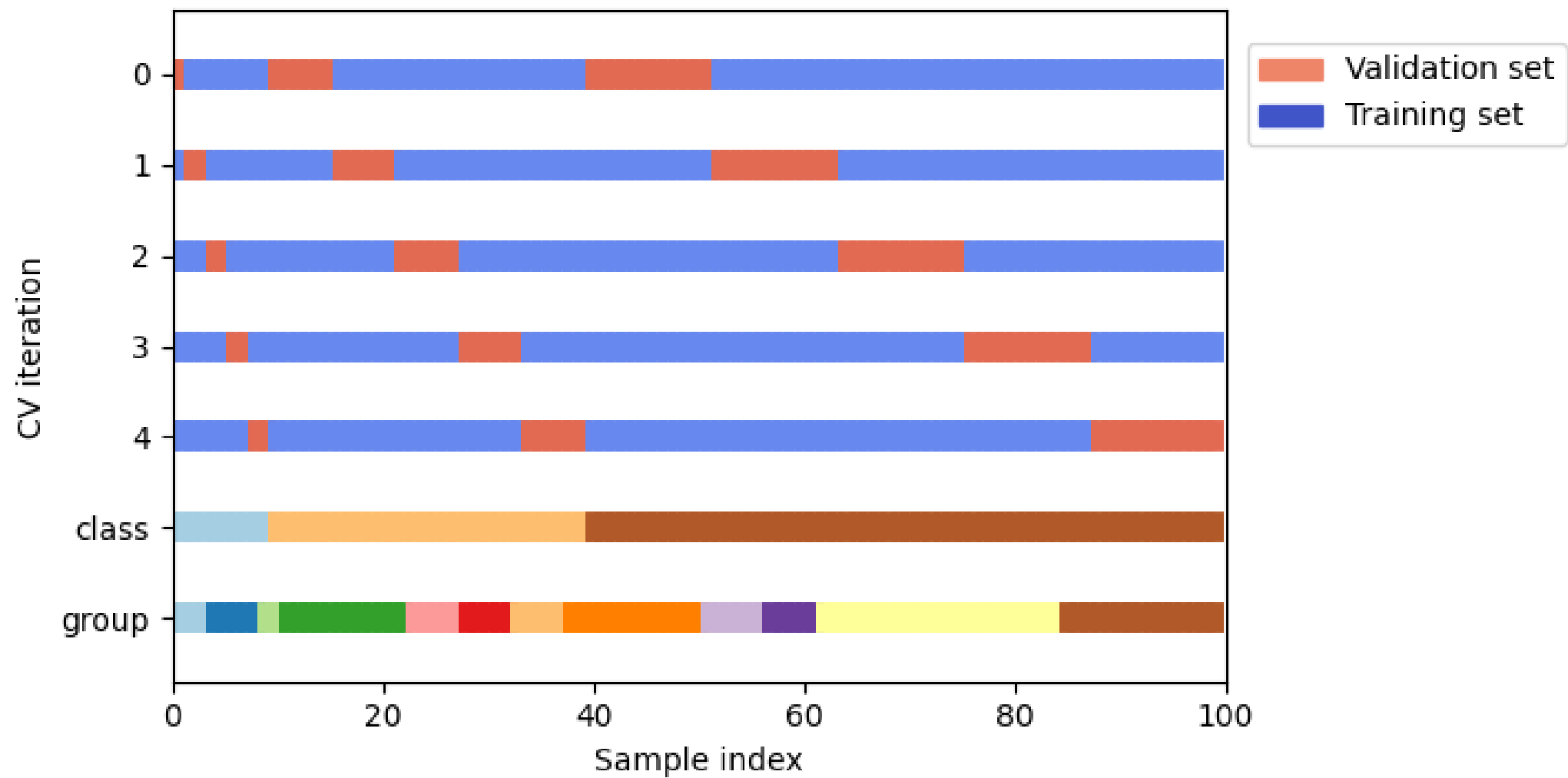
```python
tf.keras.optimizers.RMSprop(learning_rate = 0.01)
```

```python
model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=lr),
              loss='binary_crossentropy',
              metrics=[tf.keras.metrics.BinaryAccuracy(name='accuracy')])
```

# Optimizers

```
tf..keras.optimizers.SGD(learning_rate = 0.01)
```

**Calculate Prediction**
$$\hat{y} = w^T x + b$$

**Estimate Error**
$$J(w, b) = \frac{1}{m} \sum_{i=1}^{m} (\hat{y}_i - y_{i_{actual}})^2$$
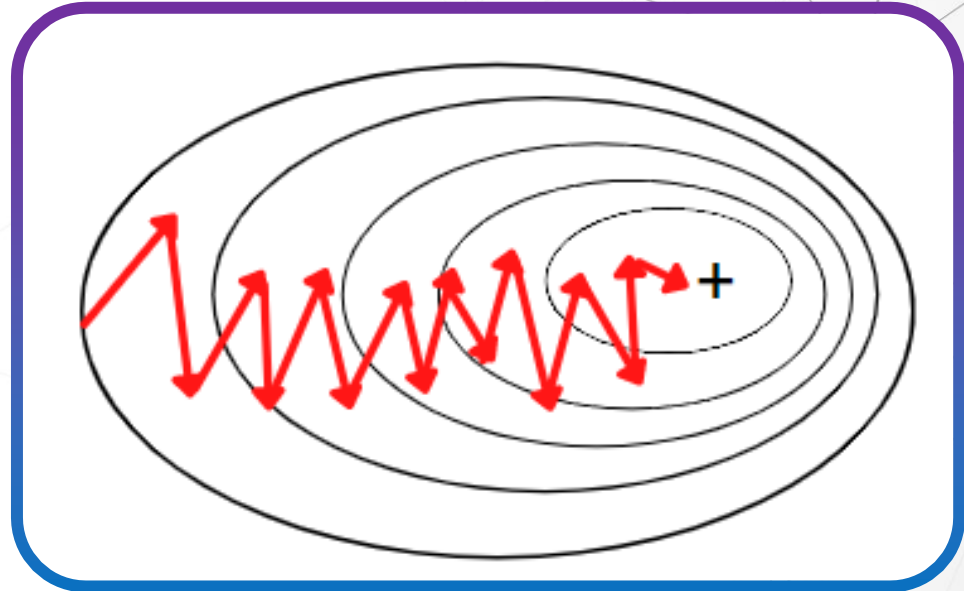
**Update Parameters**
$$w = w - \alpha \frac{\partial J}{\partial w} \qquad b = b - \alpha \frac{\partial J}{\partial b}$$

# Optimizers

```
tf.keras.optimizers.SGD(learning_rate = 0.01,
                        momentum=0.9)
```

**Calculate Prediction**
$$\hat{y} = w^T x + b$$

**Estimate Error**
$$J(w, b) = \frac{1}{m} \sum_{i=1}^{m} (\hat{y}_i - y_{i_{actual}})^2$$

**Update Parameters**
$$w = w - \alpha V_{dw} \quad b = b - \alpha V_{db}$$

$$V_{dw} = \beta V_{dw} + (1 - \beta) V_{dw}$$

$$V_{db} = \beta V_{db} + (1 - \beta) V_{db}$$

Which amounts to $\frac{1}{1-\beta}$ $values$

# Optimizers

```
tf.keras.optimizers.SGD(learning_rate = 0.01,
                        momentum=0.9)
```
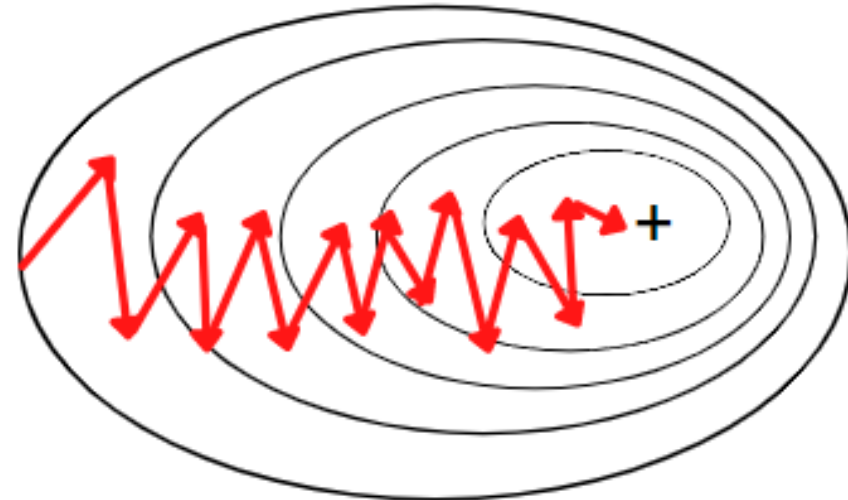
**Calculate Prediction**
$$\hat{y} = w^T x + b$$

**Estimate Error**
$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} (\hat{y}_i - y_{i_{actual}})^2$$

**Update Parameters**
$$w = w - \alpha V_{dw} \qquad b = b - \alpha V_{db}$$

# Optimizers

```
tf.keras.optimizers.RMSprop(learning_rate = 0.01)
```

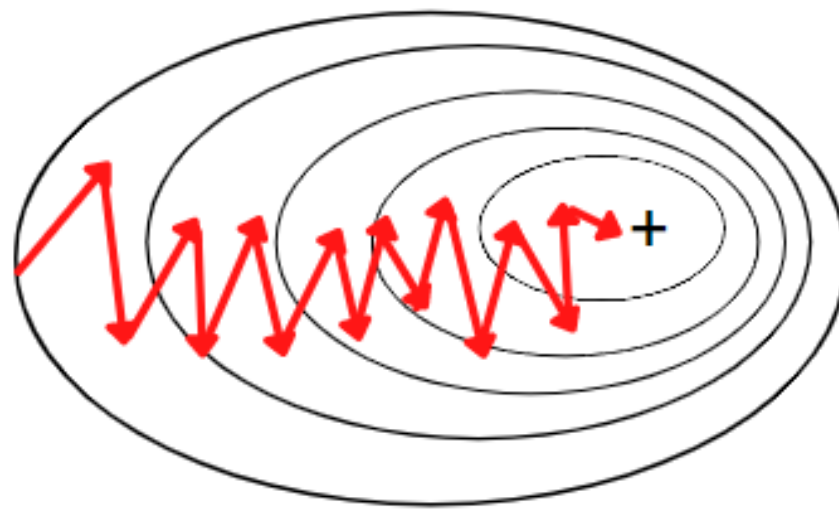**Calculate Prediction**
$$\hat{y} = w^T x + b$$

**Estimate Error**
$$J(w, b) = \frac{1}{m}\sum_{i=1}^{m}(\hat{y}_i - y_{i_{actual}})^2$$

**Update Parameters**
$$w = w - \alpha\frac{dw}{\sqrt{S_{dw}}} \quad b = b - \alpha\frac{db}{\sqrt{S_{db}}}$$

$$S_{dw} = \beta S_{dw} + (1- \beta)\, dw^2$$
$$S_{db} = \beta S_{db} + (1- \beta)\, db^2$$

# Optimizers

```python
tf..keras.optimizers.SGD(learning_rate = 0.01)
```

```python
tf.keras.optimizers.SGD(learning_rate = 0.01,
                        momentum=0.9)
```

```python
tf.keras.optimizers.RMSprop(learning_rate = 0.01)
```

```python
model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=lr),
              loss='binary_crossentropy',
              metrics=[tf.keras.metrics.BinaryAccuracy(name='accuracy')])
```

# Any Question?