

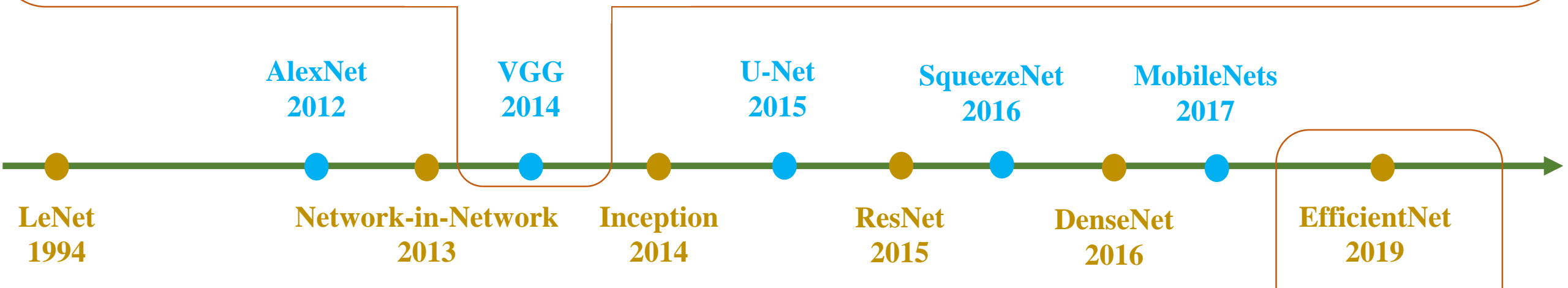
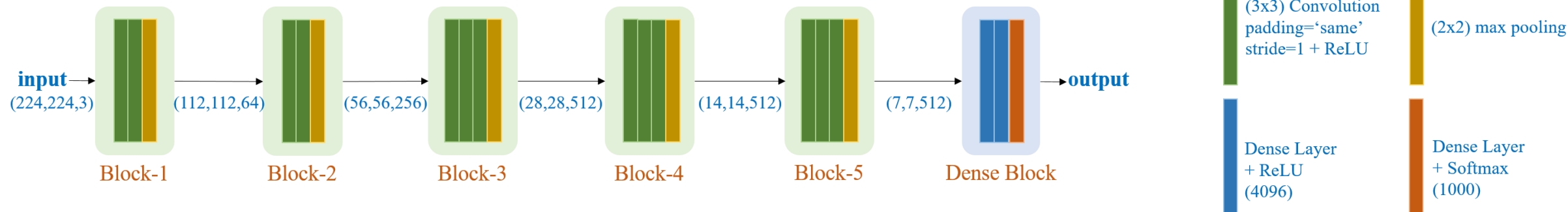
# CNN Training

## Problem-Solving Approach (Draft)

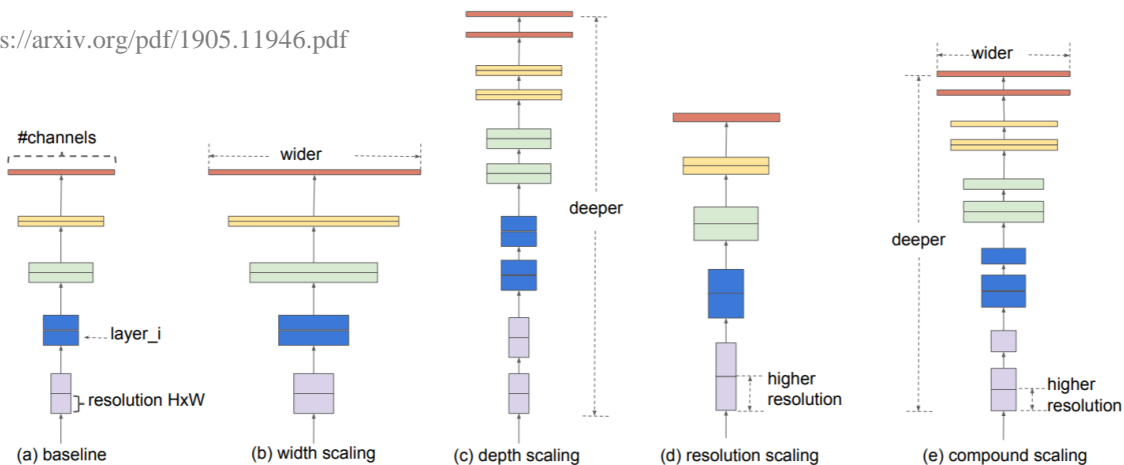
Quang-Vinh Dinh  
Ph.D. in Computer Science

# Outline

- **Introduction to Numpy**
- **Numpy Array Indexing**
- **Numpy Array Operations**
- **Broadcasting**
- **Data Processing**



<https://arxiv.org/pdf/1905.11946.pdf>



**Quoc V. Le**

Research Scientist, Google Brain  
 Verified email at stanford.edu - [Homepage](#)  
 Machine Learning Artificial Intelligence

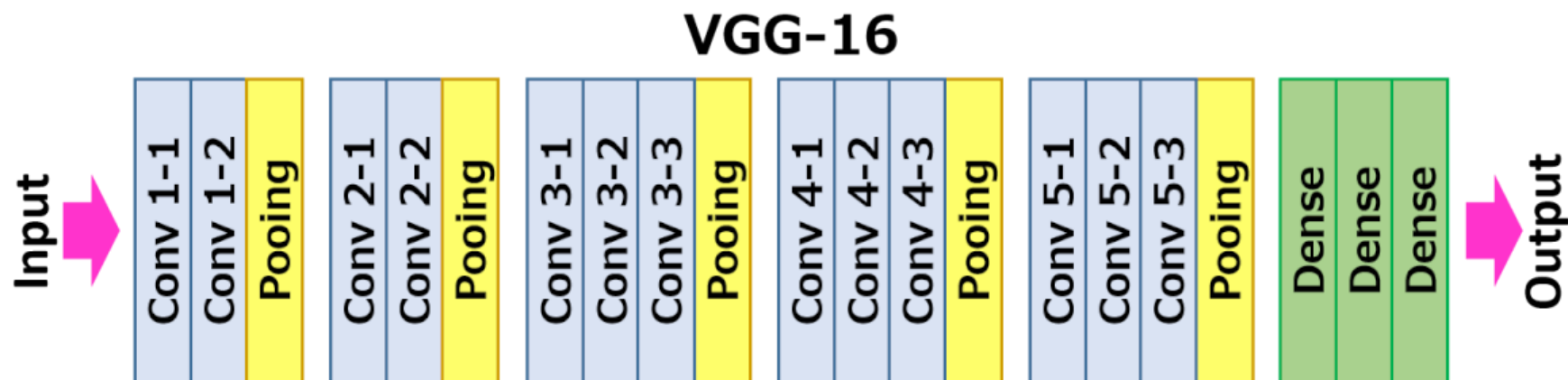
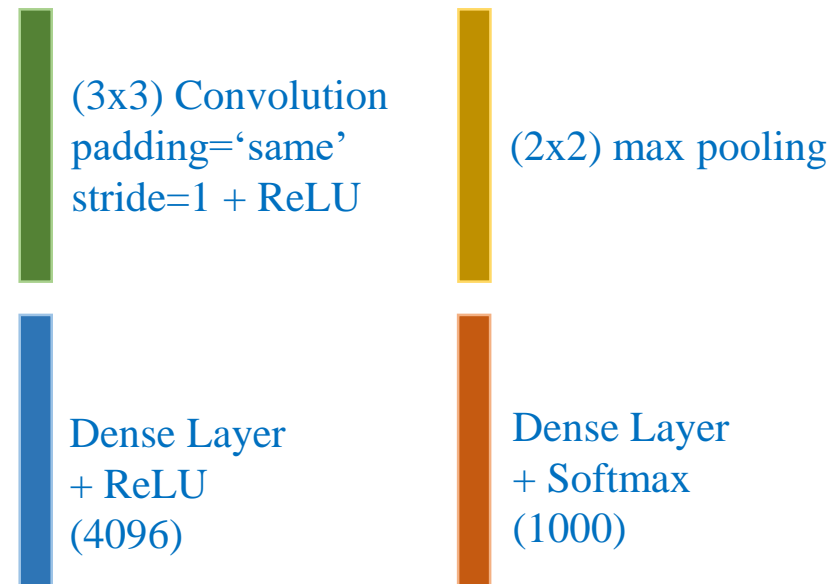
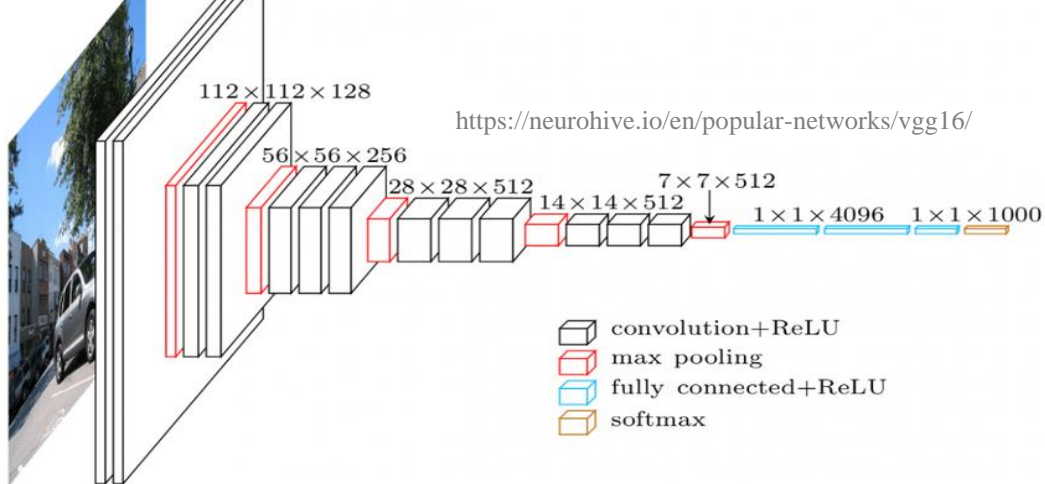
Sequence to sequence learning with neural networks

I Sutskever, O Vinyals, QV Le  
 Advances in neural information processing systems, 3104-3112

Efficientnet: Rethinking model scaling for convolutional neural networks

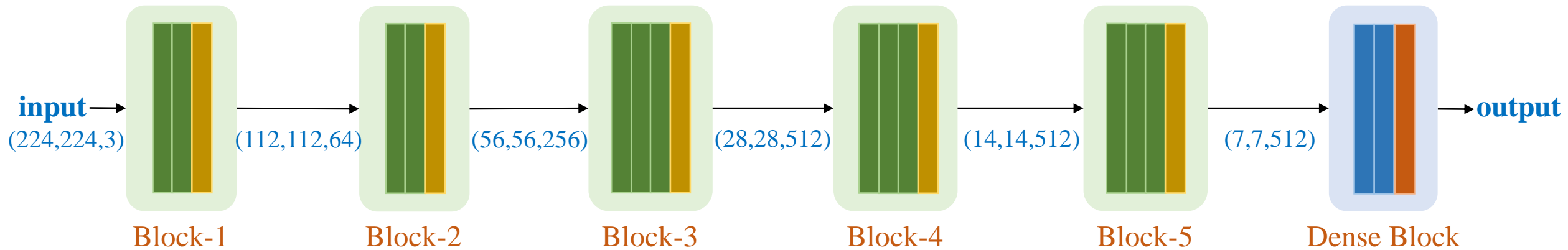
M Tan, QV Le  
 arXiv preprint arXiv:1905.11946

# VGG16



# CNN Architectures

## ❖ VGG16 for ImageNet



(3x3) Convolution  
padding='same'  
stride=1 + ReLU



(2x2) max pooling



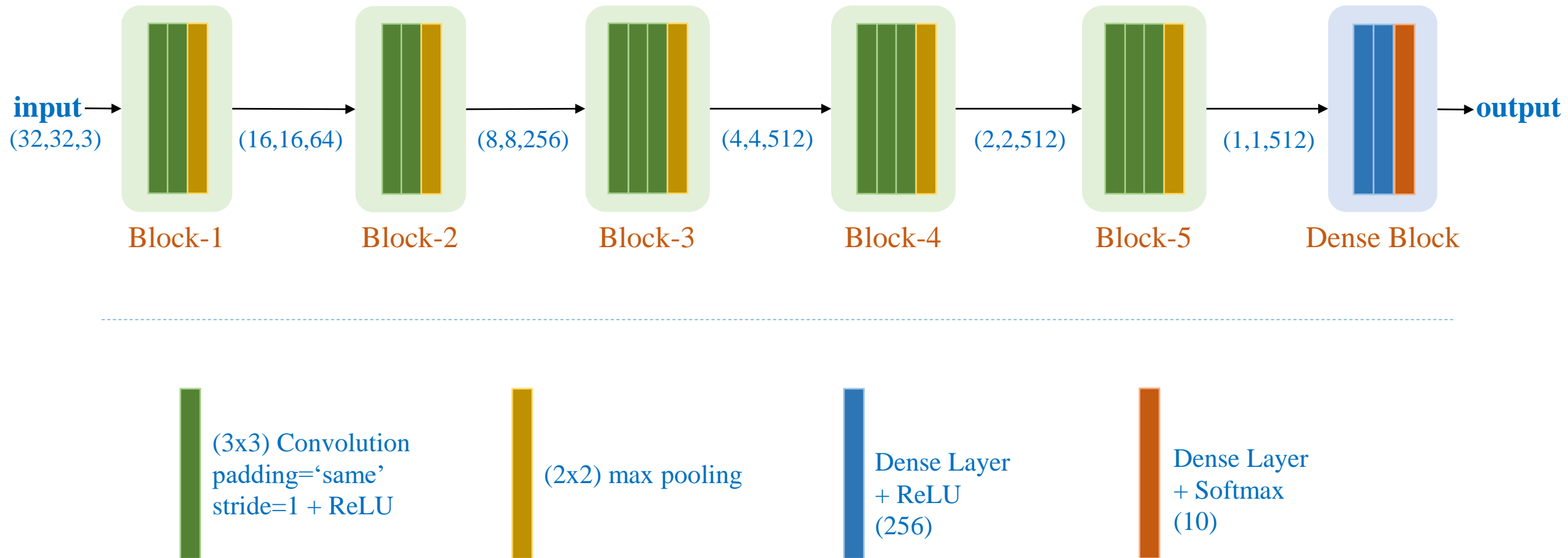
Dense Layer  
+ ReLU  
(4096)



Dense Layer  
+ Softmax  
(1000)

# CNN Architectures

## ❖ VGG16-like for Cifar-10



# VGG16 in Keras

Layer (type)	Output Shape	Param #
=====		
input_1 (InputLayer)	[(None, 244, 244, 3)]	0
block1_conv1 (Conv2D)	(None, 244, 244, 64)	1792
block1_conv2 (Conv2D)	(None, 244, 244, 64)	36928
block1_pool (MaxPooling2D)	(None, 122, 122, 64)	0
block2_conv1 (Conv2D)	(None, 122, 122, 128)	73856
block2_conv2 (Conv2D)	(None, 122, 122, 128)	147584
block2_pool (MaxPooling2D)	(None, 61, 61, 128)	0
block3_conv1 (Conv2D)	(None, 61, 61, 256)	295168
block3_conv2 (Conv2D)	(None, 61, 61, 256)	590080
block3_conv3 (Conv2D)	(None, 61, 61, 256)	590080
block3_pool (MaxPooling2D)	(None, 30, 30, 256)	0
block4_conv1 (Conv2D)	(None, 30, 30, 512)	1180160
block4_conv2 (Conv2D)	(None, 30, 30, 512)	2359808
block4_conv3 (Conv2D)	(None, 30, 30, 512)	2359808
block4_pool (MaxPooling2D)	(None, 15, 15, 512)	0

```

1 import tensorflow as tf
2 from tensorflow import keras
3
4 VGG16 = keras.applications.VGG16
5 model = model = VGG16(input_shape=(244,244,3),
6                           classes=10,
7                           weights=None)
8 model.summary()

```

block5_conv1 (Conv2D)	(None, 15, 15, 512)	2359808
block5_conv2 (Conv2D)	(None, 15, 15, 512)	2359808
block5_conv3 (Conv2D)	(None, 15, 15, 512)	2359808
block5_pool (MaxPooling2D)	(None, 7, 7, 512)	0
flatten (Flatten)	(None, 25088)	0
fc1 (Dense)	(None, 4096)	102764544
fc2 (Dense)	(None, 4096)	16781312
predictions (Dense)	(None, 10)	40970
=====		
Total params: 134,301,514		
Trainable params: 134,301,514		
Non-trainable params: 0		

```
# model
model = keras.models.Sequential()
model.add(tf.keras.Input(shape=(32, 32, 3)))

num_of_blocks = 5
convs = [2, 2, 3, 3, 3]
filters = [64, 256, 512, 512, 512]

# 5 blocks
for block in range(num_of_blocks):
    for conv in range(convs[block]):
        model.add(keras.layers.Conv2D(filters[block], 3,
                                       padding='same',
                                       activation='relu'))
    model.add(keras.layers.MaxPooling2D(2))

# Dense blocks
model.add(keras.layers.Flatten())
model.add(keras.layers.Dense(256, activation='relu'))
model.add(keras.layers.Dense(256, activation='relu'))
model.add(keras.layers.Dense(10, activation='softmax'))
```

Model: "sequential"

Layer (type)	Output Shape	Param #
=====		
conv2d (Conv2D)	(None, 32, 32, 64)	1792
conv2d_1 (Conv2D)	(None, 32, 32, 64)	36928
max_pooling2d (MaxPooling2D)	(None, 16, 16, 64)	0
conv2d_2 (Conv2D)	(None, 16, 16, 256)	147712
conv2d_3 (Conv2D)	(None, 16, 16, 256)	590080
max_pooling2d_1 (MaxPooling2D)	(None, 8, 8, 256)	0

# VGG16 in Keras

conv2d_4 (Conv2D)	(None, 8, 8, 512)	1180160
conv2d_5 (Conv2D)	(None, 8, 8, 512)	2359808
conv2d_6 (Conv2D)	(None, 8, 8, 512)	2359808
max_pooling2d_2 (MaxPooling2D)	(None, 4, 4, 512)	0
conv2d_7 (Conv2D)	(None, 4, 4, 512)	2359808
conv2d_8 (Conv2D)	(None, 4, 4, 512)	2359808
conv2d_9 (Conv2D)	(None, 4, 4, 512)	2359808
max_pooling2d_3 (MaxPooling2D)	(None, 2, 2, 512)	0
conv2d_10 (Conv2D)	(None, 2, 2, 512)	2359808
conv2d_11 (Conv2D)	(None, 2, 2, 512)	2359808
conv2d_12 (Conv2D)	(None, 2, 2, 512)	2359808
max_pooling2d_4 (MaxPooling2D)	(None, 1, 1, 512)	0
flatten (Flatten)	(None, 512)	0
dense (Dense)	(None, 256)	131328
dense_1 (Dense)	(None, 256)	65792
dense_2 (Dense)	(None, 10)	2570
=====		
Total params: 21,034,826		
Trainable params: 21,034,826		
Non-trainable params: 0		



# Dg

---



# Image Data

## Fashion-MNIST dataset

Grayscale images

Resolution=28x28

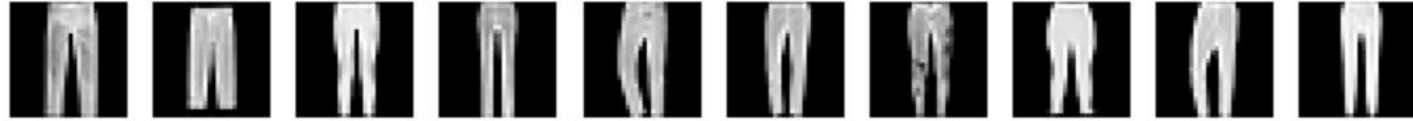
Training set: 60000 samples

Testing set: 10000 samples

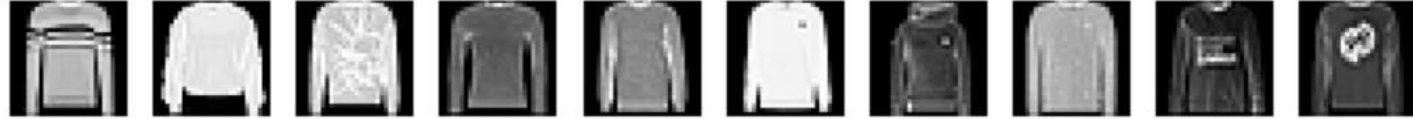
T-shirt



Trouser



Pullover



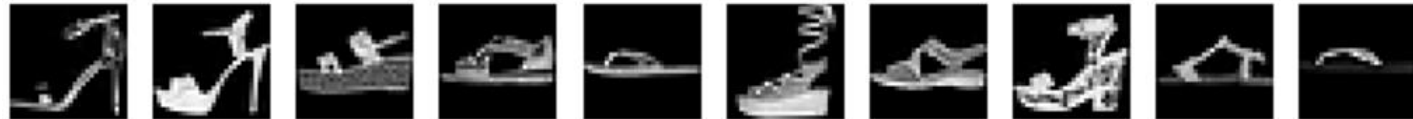
Dress



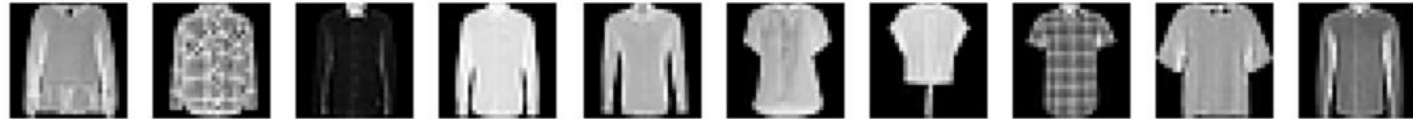
Coat



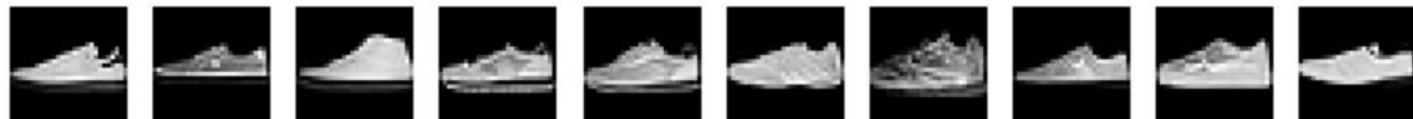
Sandal



Shirt



Sneaker



Bag

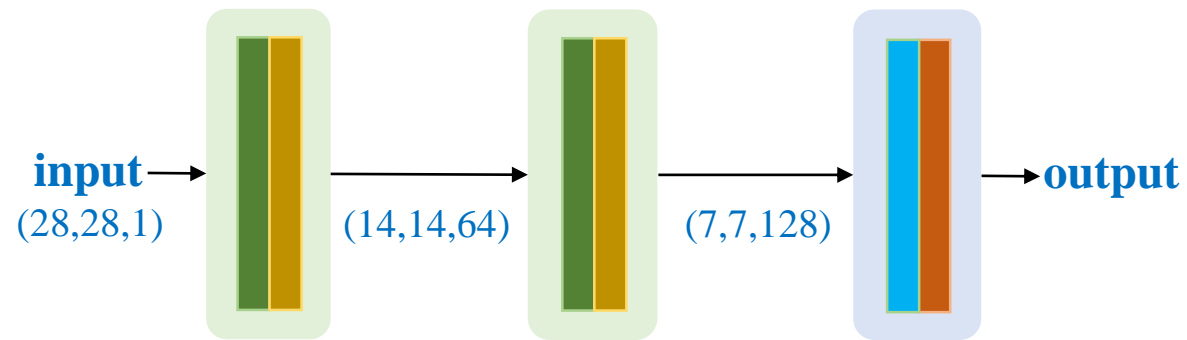



Ankle  
Boot




# Network Training

## ❖ Fashion-MNIST dataset

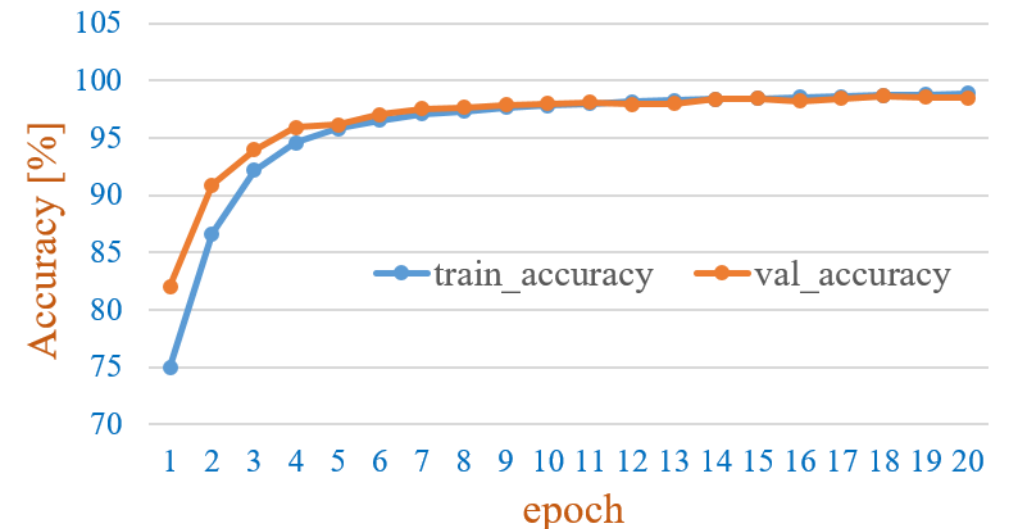
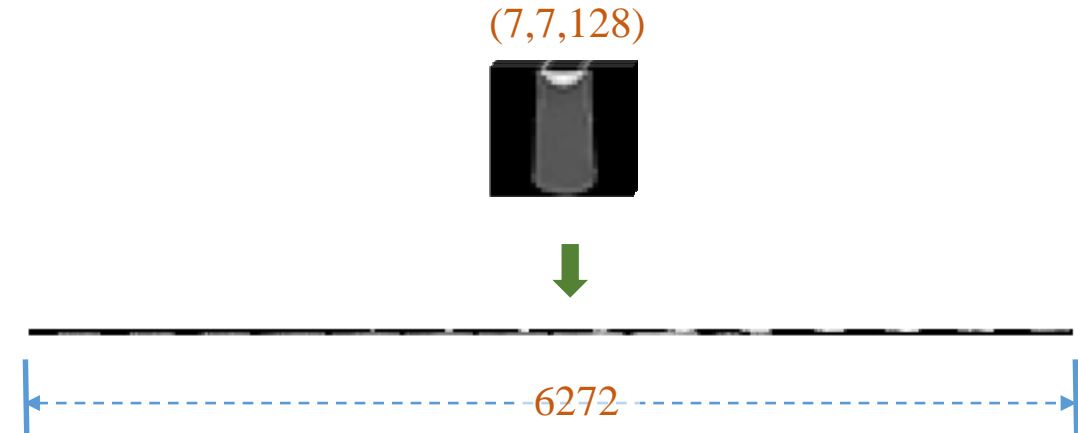


 (3x3) Convolution  
padding='same'  
stride=1 + Sigmoid

 (2x2) max pooling

 Flatten

 Dense Layer-10  
+ Softmax



# Network Training

## ❖ Fashion-MNIST dataset

### X-data format

(batch, height, width, channel)

### Data normalization [0,1]

(3x3) Convolution with 64 filters,  
stride=1, padding='same'  
+ Sigmoid activation  
+ glorot\_uniform initialization

Adam optimizer and Cross-entropy loss

```
1 import tensorflow as tf
2 import tensorflow.keras as keras
3
4 # data preparation
5 mnist = tf.keras.datasets.mnist
6 (x_train, y_train), (x_test, y_test) = mnist.load_data()
7
8 x_train = x_train.reshape((60000, 28, 28, 1))
9 x_test = x_test.reshape((10000, 28, 28, 1))
10
11 # normalize
12 x_train, x_test = x_train / 255.0, x_test / 255.0
13
14 # model
15 model = keras.models.Sequential()
16 model.add(tf.keras.Input(shape=(28, 28, 1)))
17
18 model.add(keras.layers.Conv2D(64, 3, padding='same',
19                                activation='sigmoid'))
20 model.add(keras.layers.MaxPooling2D(2))
21
22 model.add(keras.layers.Conv2D(128, 3, padding='same',
23                                activation='sigmoid'))
24 model.add(keras.layers.MaxPooling2D(2))
25
26 # flatten
27 model.add(keras.layers.Flatten())
28 model.add(keras.layers.Dense(10, activation='softmax'))
29 model.summary()
30
31 # training
32 model.compile(optimizer='adam', metrics=['accuracy'],
33               loss='sparse_categorical_crossentropy')
34 history = model.fit(x_train, y_train, batch_size=256,
35                     validation_data=(x_test, y_test),
36                                     epochs=20, verbose=1)
```



# Network Training

**Cifar-10 dataset**  
(more complex dataset)

Color images

Resolution=32x32

Training set: 50000 samples

Testing set: 10000 samples

airplane



automobile



bird



cat



deer



dog



frog



horse



ship

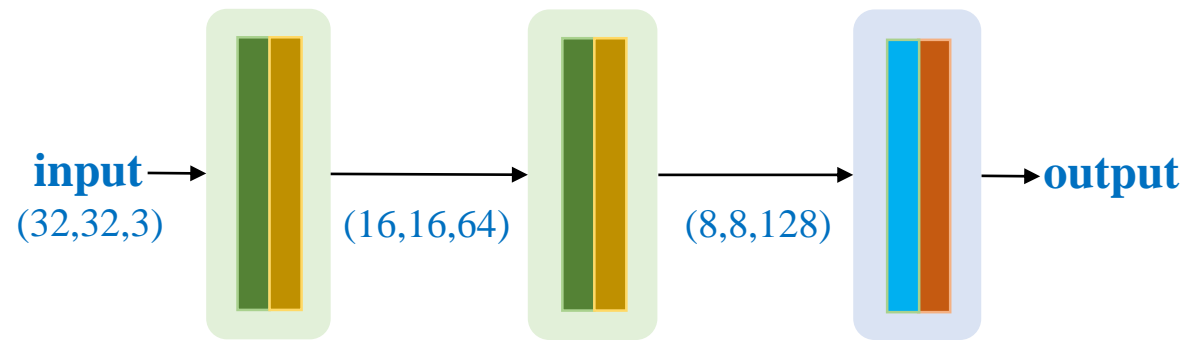



truck




# Network Training


## ❖ Cifar-10 dataset

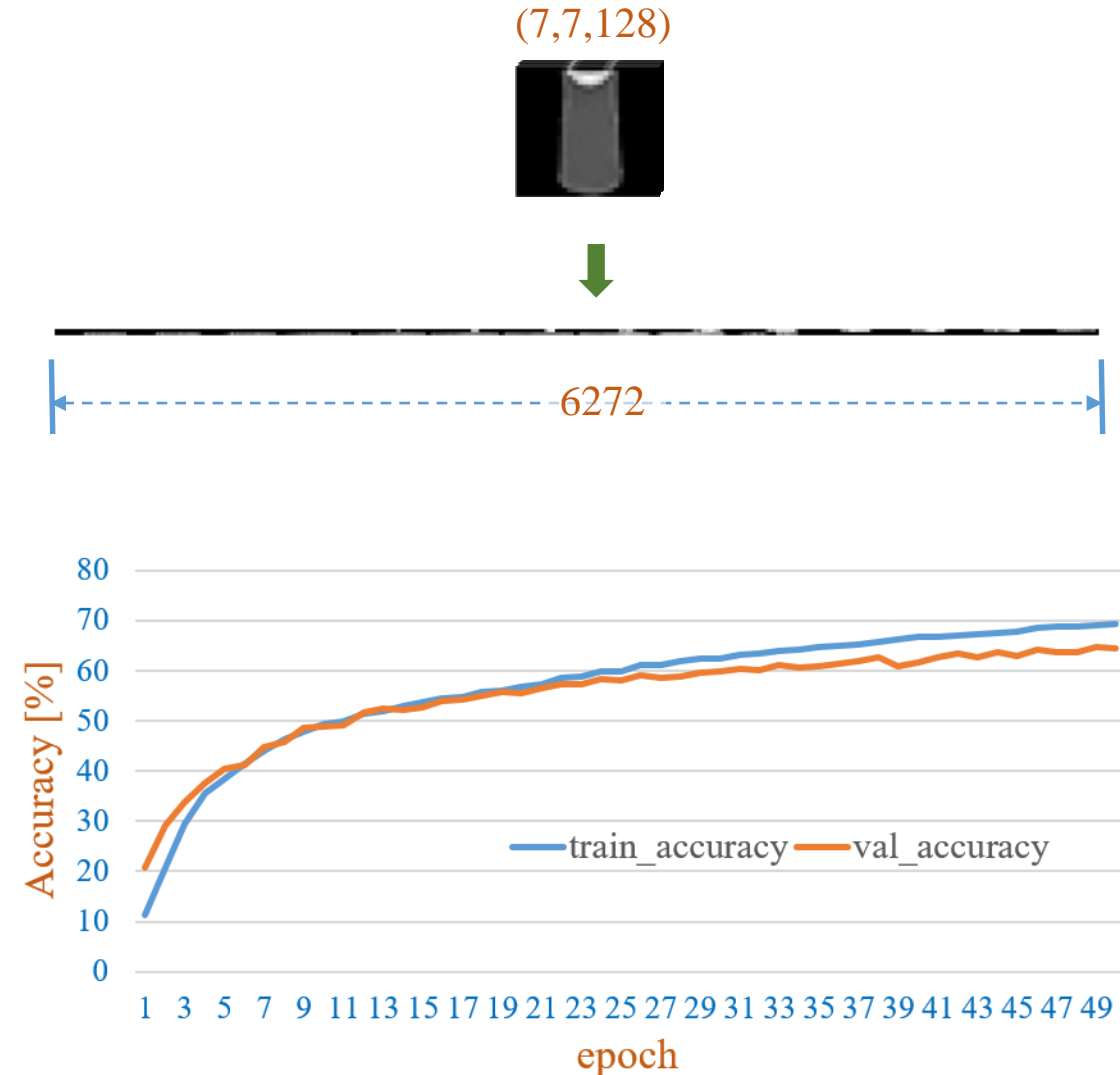


 (3x3) Convolution  
padding='same'  
stride=1 + Sigmoid

 (2x2) max pooling

 Flatten

 Dense Layer-10  
+ Softmax

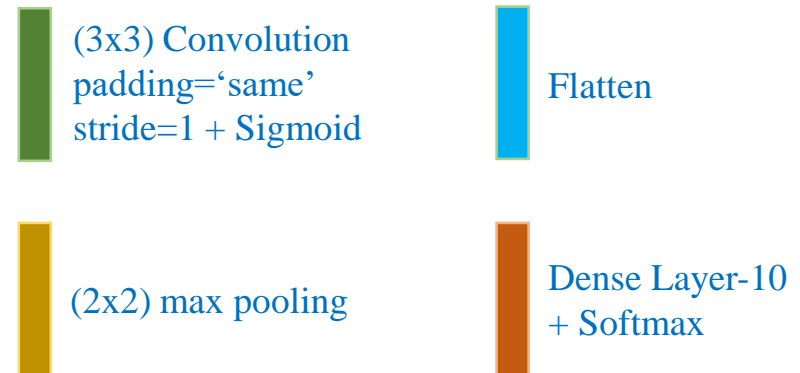
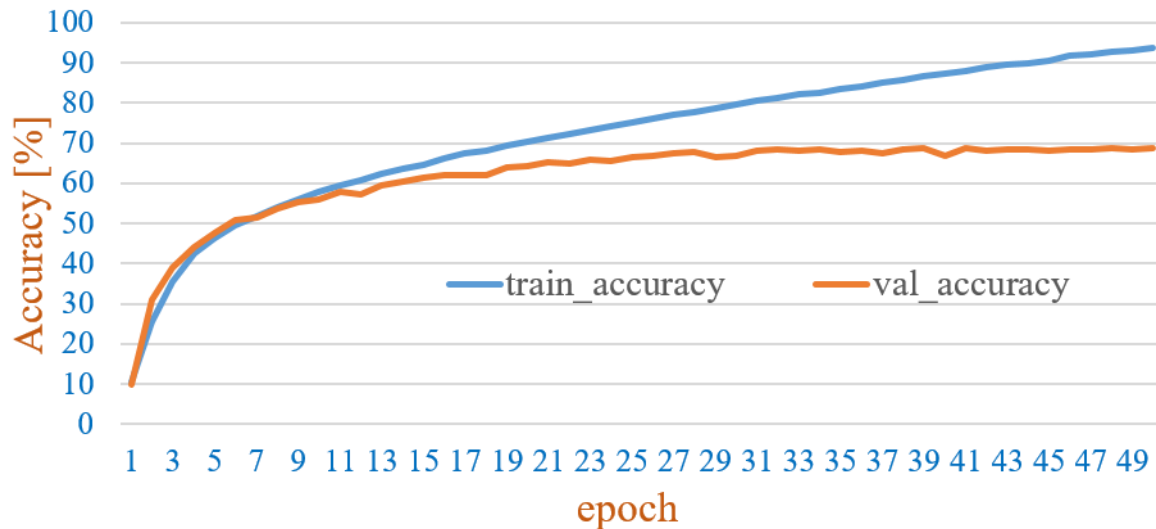
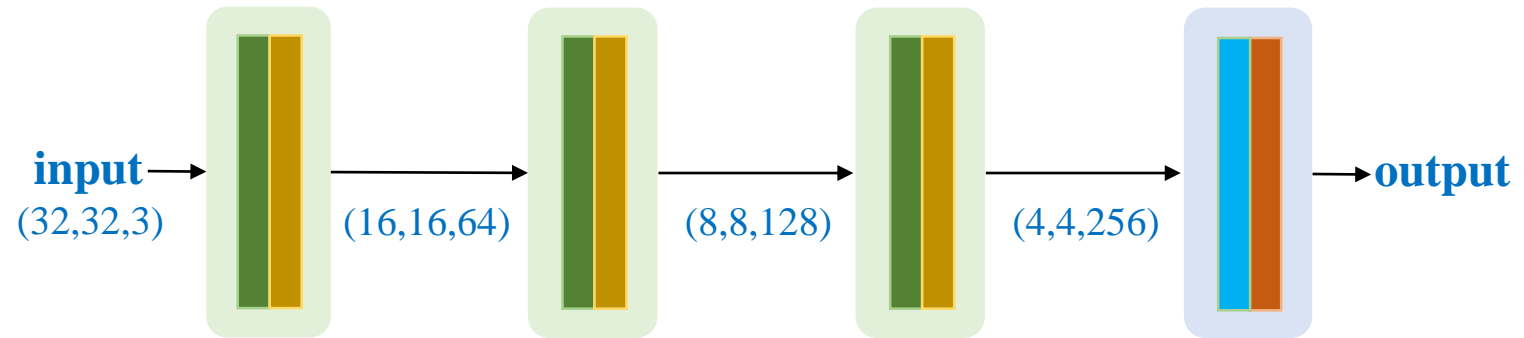


# Network Training

## ❖ Cifar-10 dataset:

### ❖ Adding more layers

**Good news: Network accuracy increases about 25%**

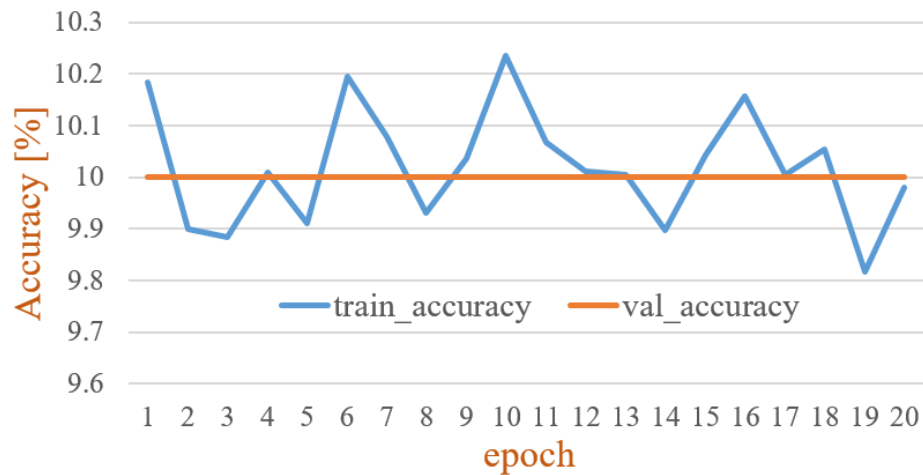
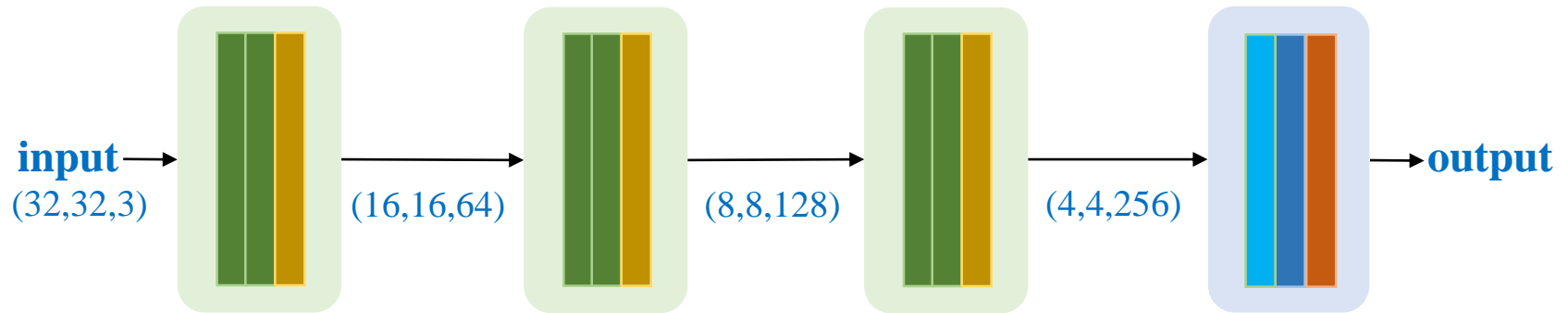


# Network Training

## ❖ Cifar-10 dataset:

### ❖ Keep adding more layers

The network does  
not learn



(3x3) Convolution  
padding='same'  
stride=1 + Sigmoid

(2x2) max pooling

Flatten

Dense Layer-10  
+ Softmax

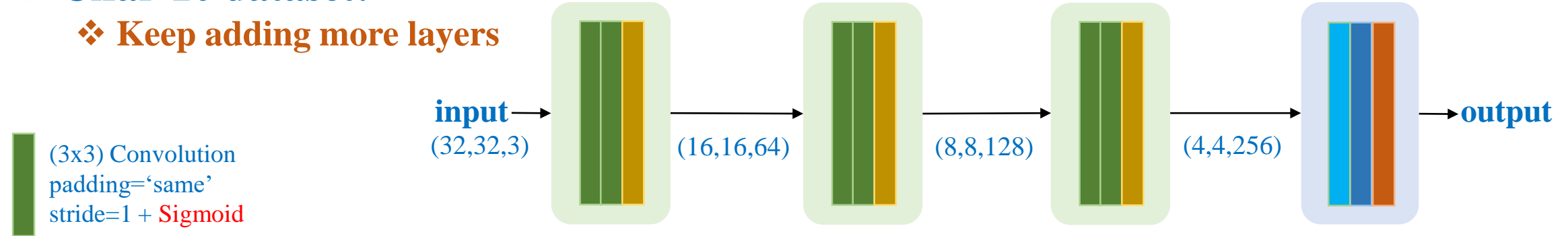
Dense Layer-512  
+ Sigmoid



# Network Training

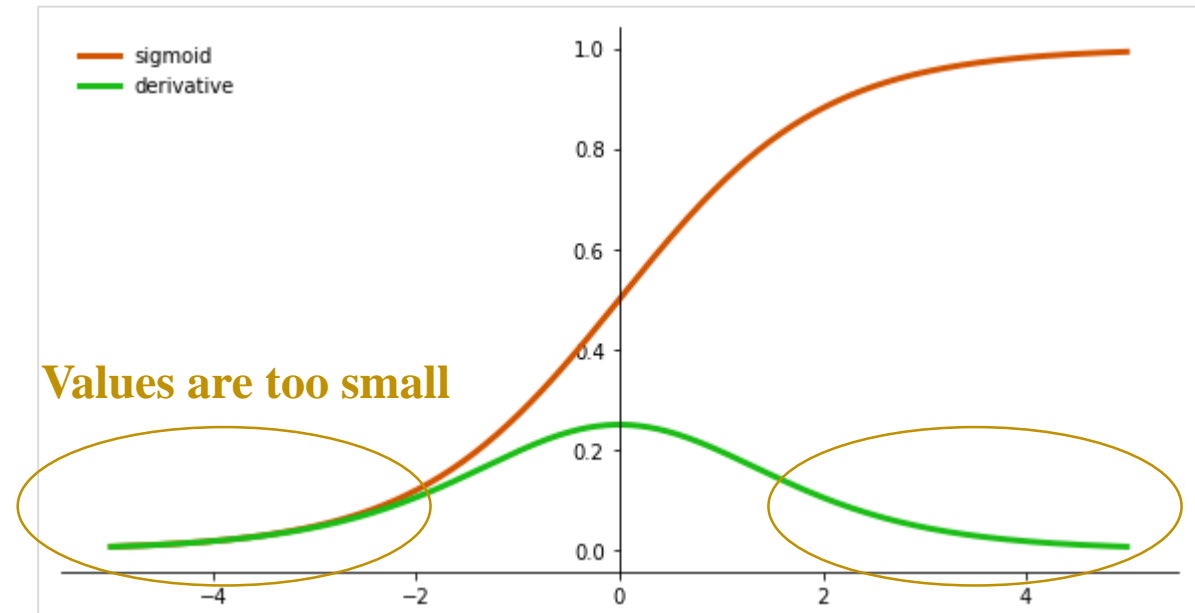
## ❖ Cifar-10 dataset:

### ❖ Keep adding more layers



$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$


Vanishing Problem



# Network Training

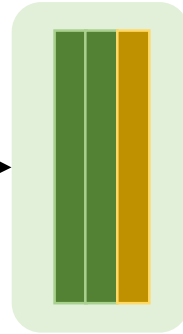
## ❖ Cifar-10 dataset:

### ❖ Keep adding more layers

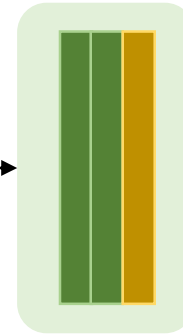
 (3x3) Convolution  
padding='same'  
stride=1 + **ReLU**

 Dense Layer-512  
+ **ReLU**

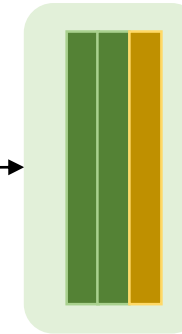
**input**  
(32,32,3)



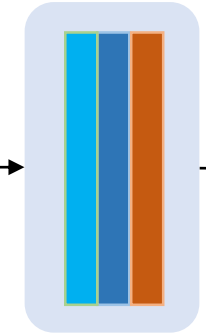
(16,16,64)



(8,8,128)



(4,4,256)



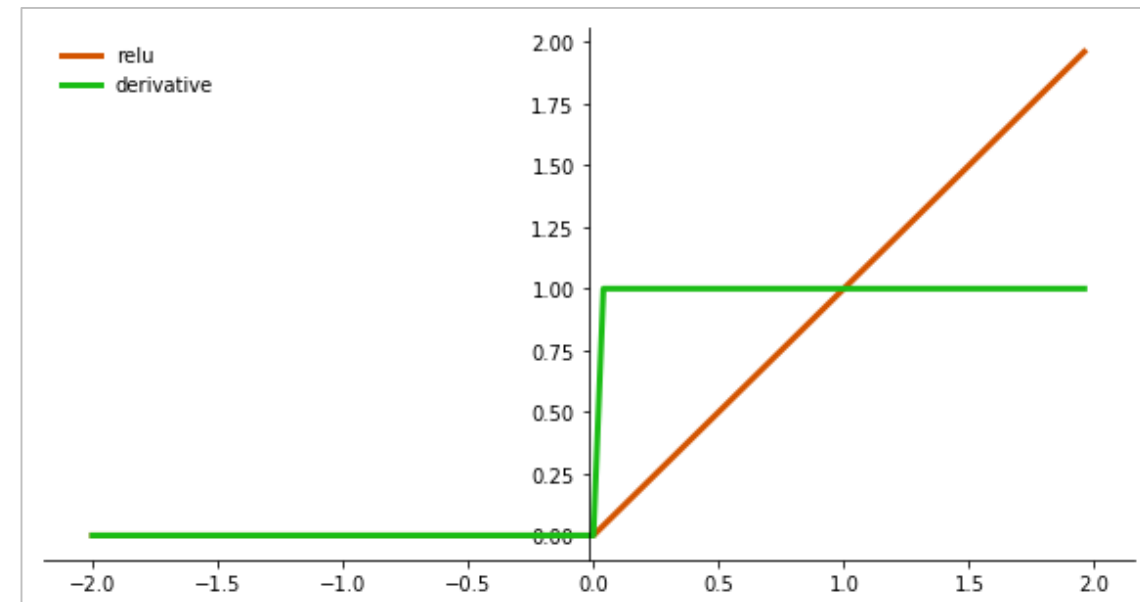
**output**

$$\text{ReLU}(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

```
Conv2D(num_filters, kernel_size, activation='sigmoid')
```



```
Conv2D(num_filters, kernel_size, activation='relu')
```

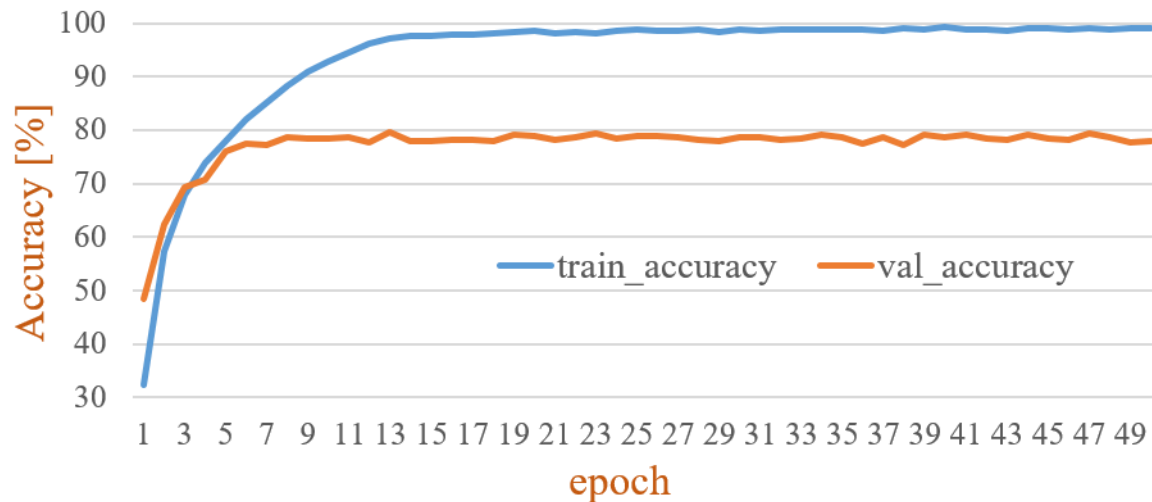
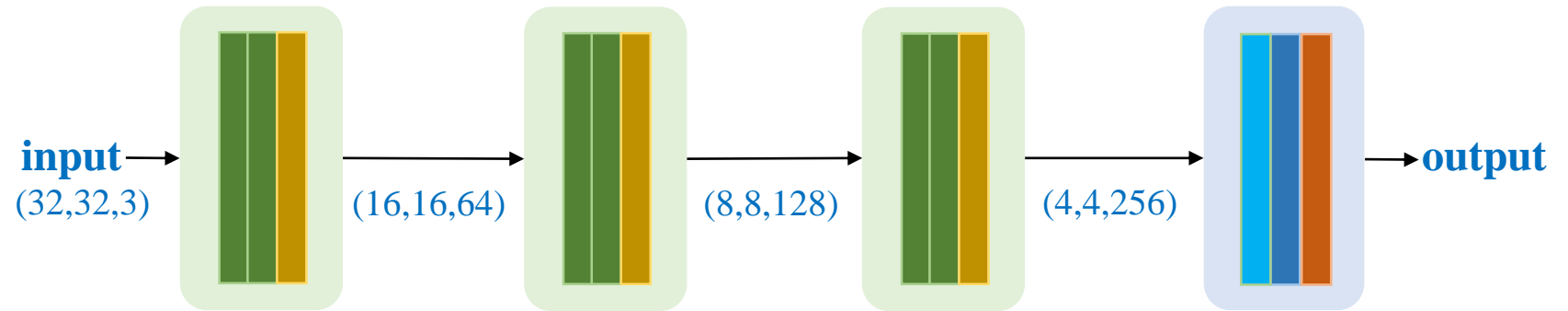


# Network Training

❖ **Cifar-10 dataset:**

❖ **Use ReLU**

**Training Accuracy  
reaches up to 99%**



(3x3) Convolution  
padding='same'  
stride=1 + ReLU

(2x2) max pooling

Flatten

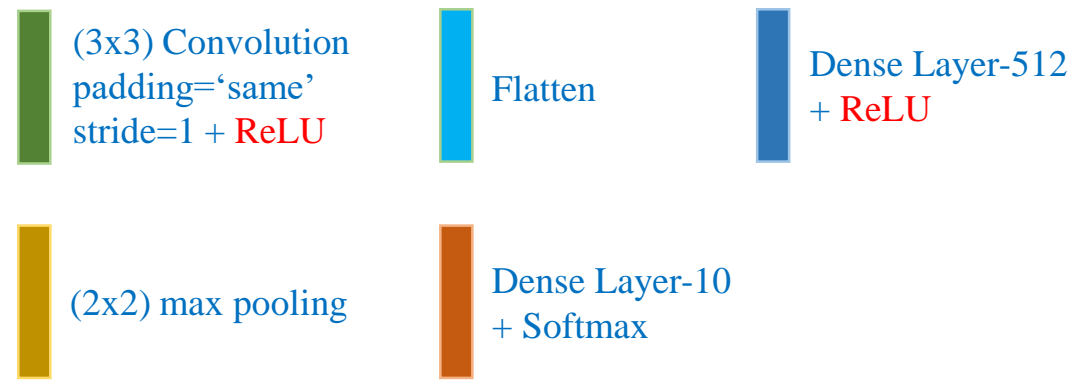
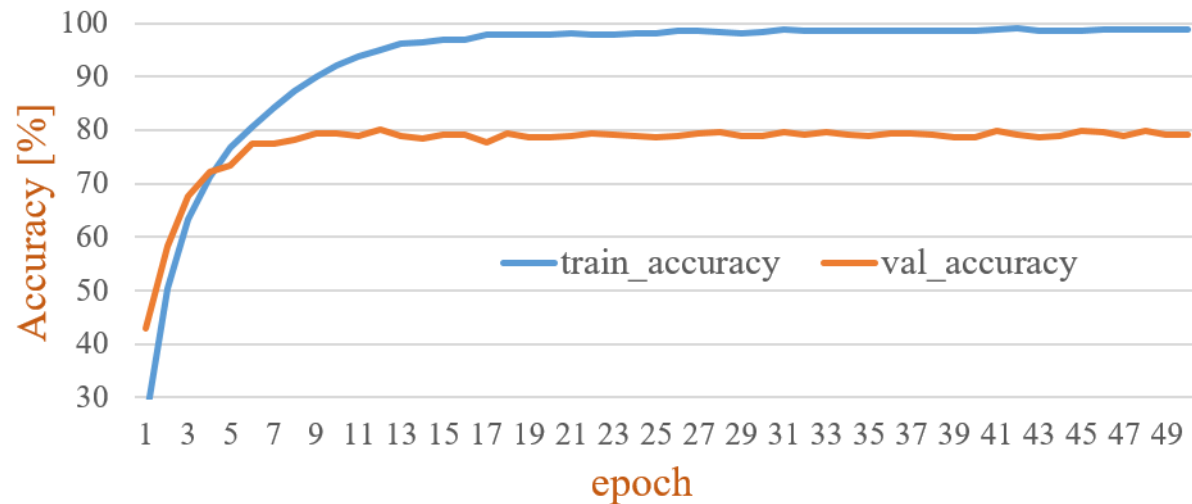
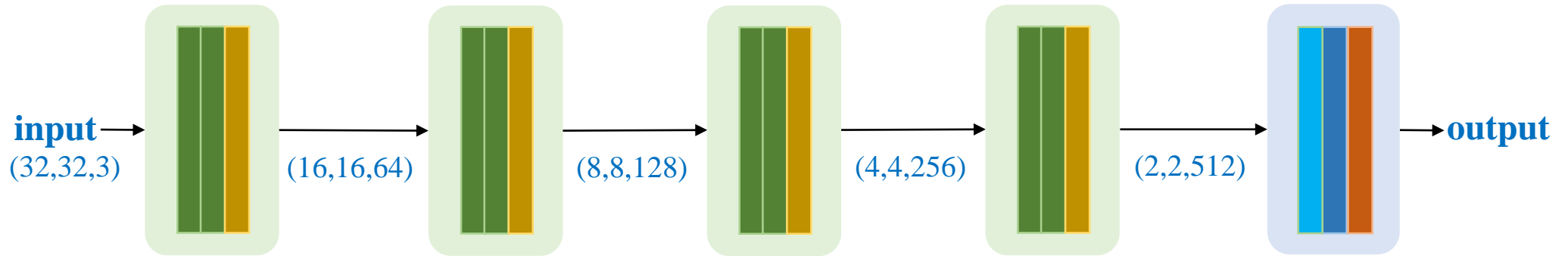
Dense Layer-10  
+ Softmax

Dense Layer-512  
+ ReLU

Adding more layers; Hope reach to 100%

# Network Training

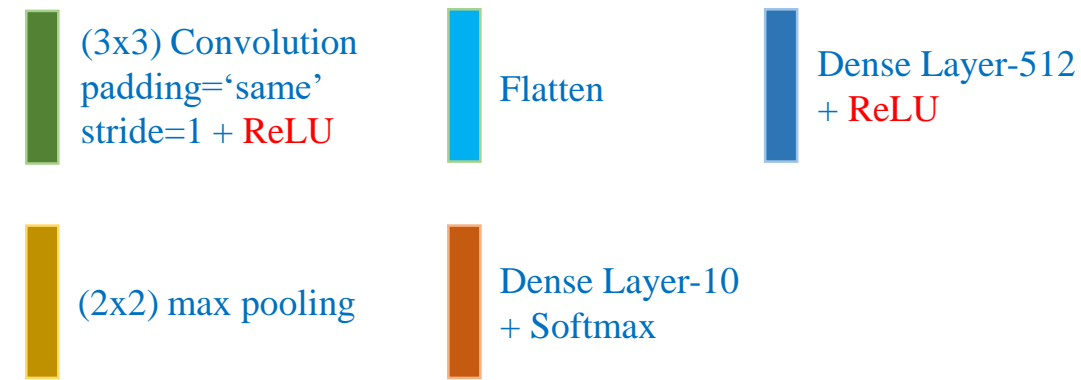
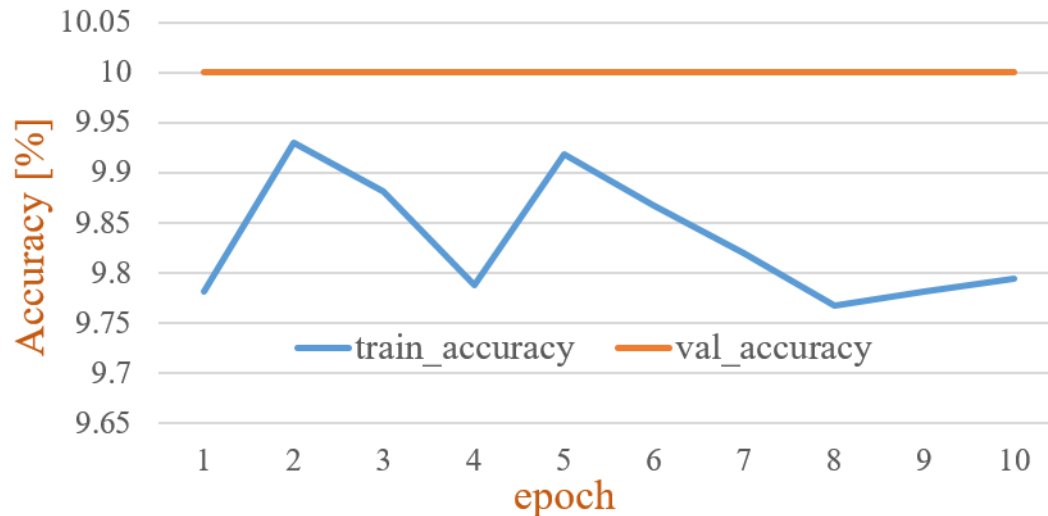
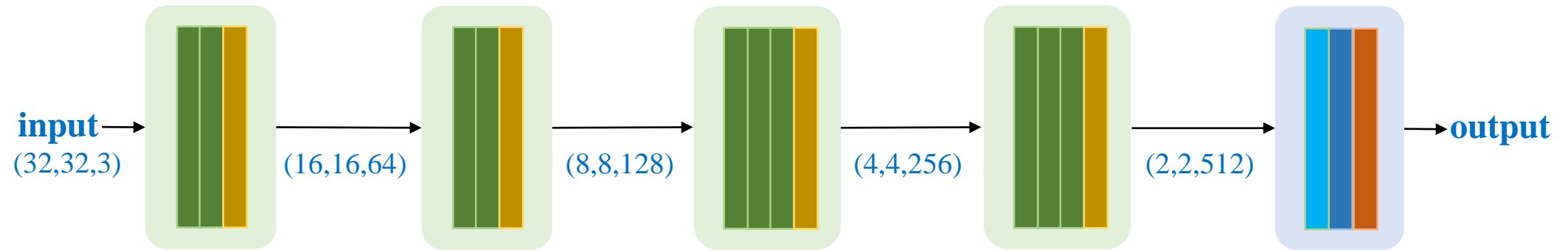
Use ReLU and add more layers



Same performance; Keep adding more layers

# Network Training

Use ReLU and add more layers



Network does not learn again

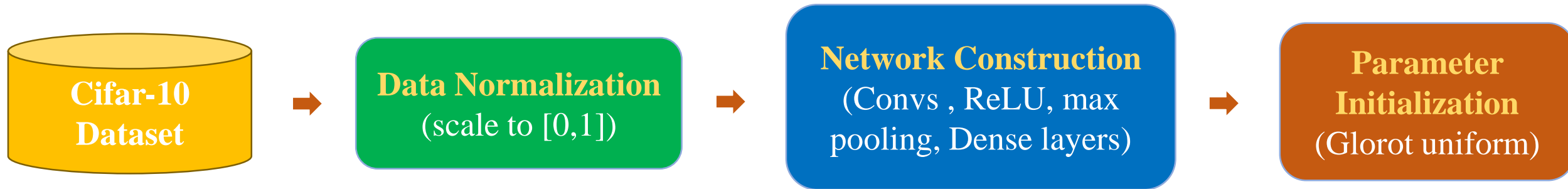
# Dg

---

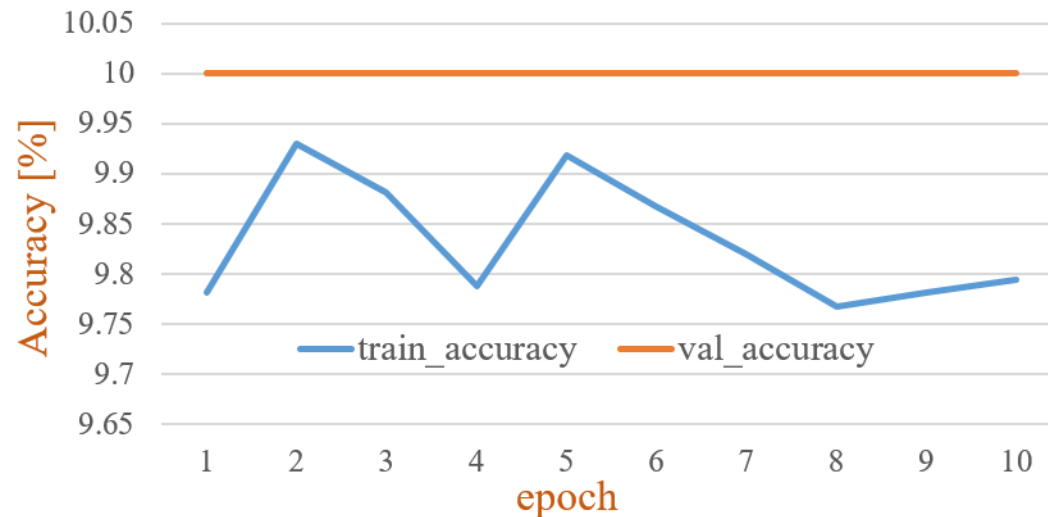


# Network Training

## ❖ Summary of the current network

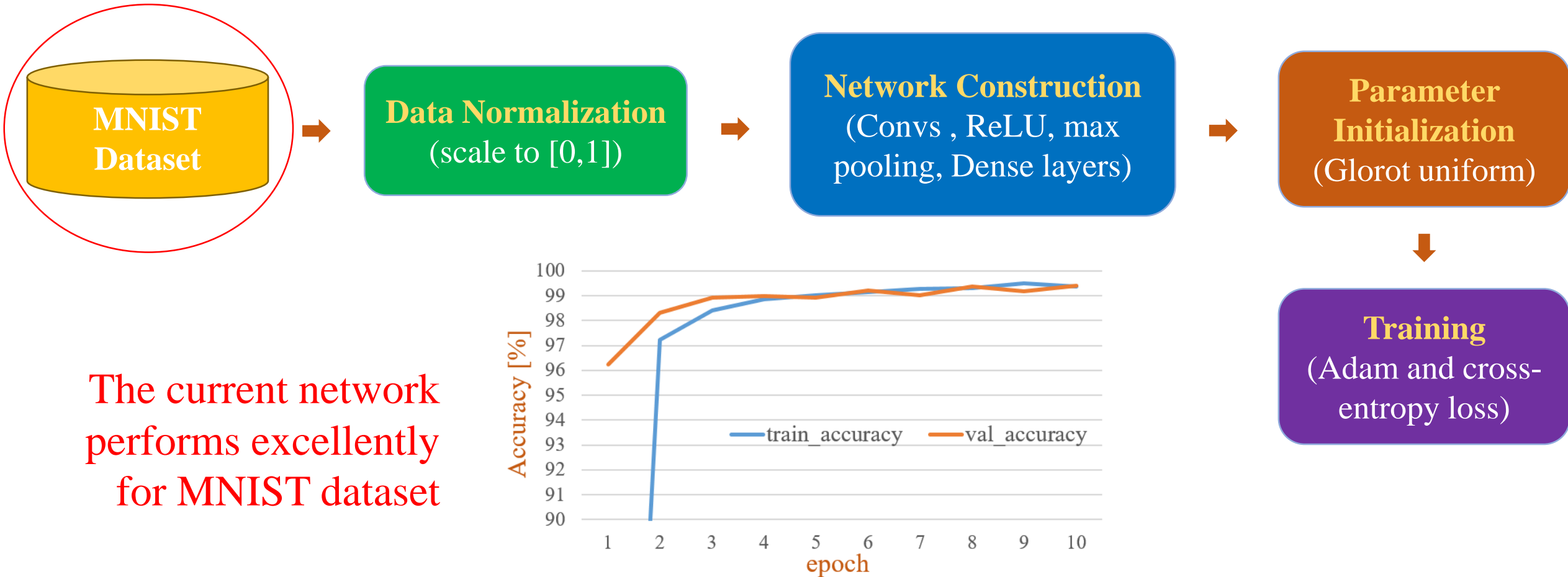


Network does  
not learn



# Network Training

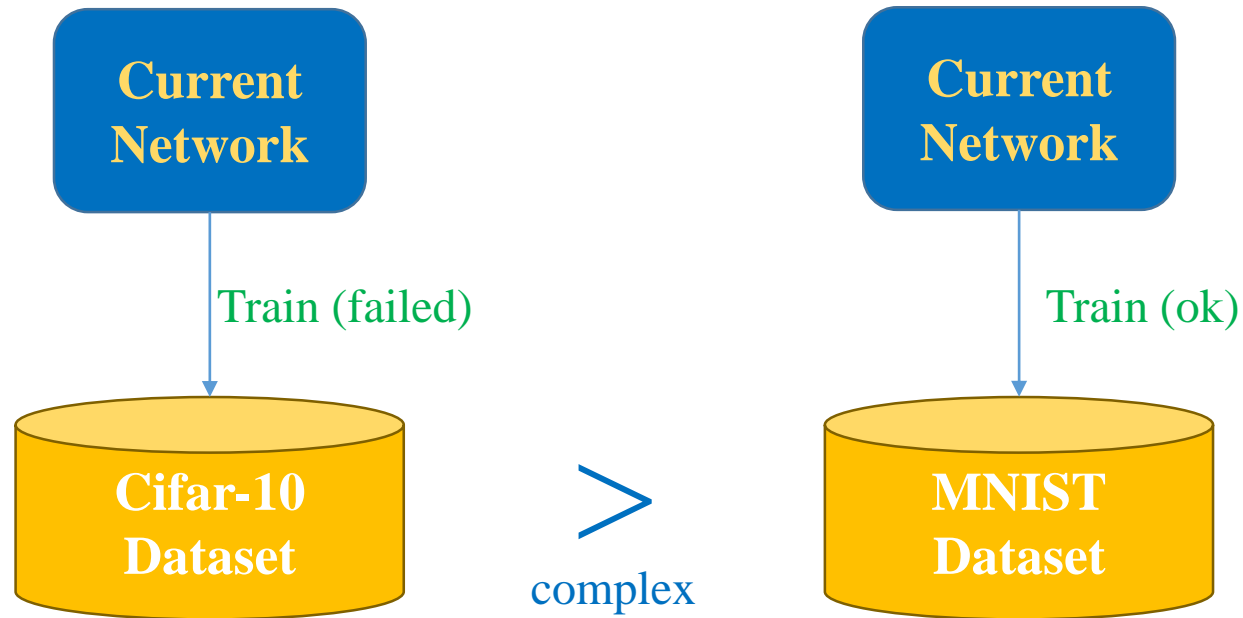
## ❖ Solution 1: Observation





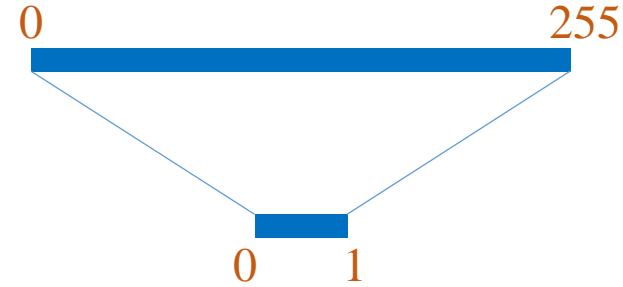
# Network Training

## ❖ Solution 1: Idea



How to reduce the complexity of the Cifar-10 dataset

**Data Normalization**  
(scale to [0,1])



**Data Normalization**  
(convert to 0-mean and 1-deviation)

$X =$



$$X = \frac{X - \mu}{\sigma}$$

$$\mu = \frac{1}{n} \sum_i X_i$$

$$\sigma = \sqrt{\frac{1}{n} \sum_i (X_i - \mu)^2}$$

# Network Training

## ❖ Solution 1: Idea

$$\bar{X} = \frac{X - \mu}{\sigma}$$

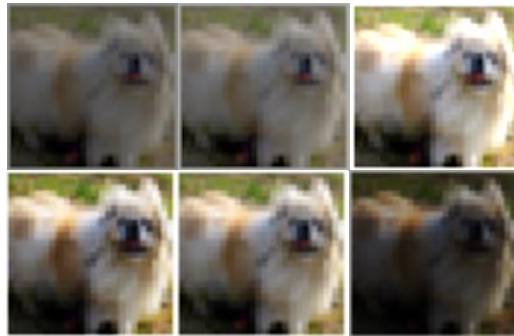
$$\mu = \frac{1}{n} \sum_i X_i$$

$$\sigma = \sqrt{\frac{1}{n} \sum_i (X_i - \mu)^2}$$

This normalization helps network to be invariant to linear transformation

$$Y = aX + b$$

$$\bar{Y} = \frac{Y - \mu_Y}{\sigma_Y} = \bar{X}$$



$$Y = aX + b$$

$$\begin{aligned} \bar{Y} &= \frac{Y - \mu_Y}{\sigma_Y} = \frac{(aX + b) - \frac{1}{n} \sum_i (aX_i + b)}{\sqrt{\frac{1}{n} \sum_i \left( (aX_i + b) - \frac{1}{n} \sum_j (aX_j + b) \right)^2}} \\ &= \frac{aX - \frac{1}{n} \sum_i aX_i}{\sqrt{\frac{1}{n} \sum_i \left( aX_i - \frac{1}{n} \sum_j aX_j \right)^2}} \\ &= \frac{X - \frac{1}{n} \sum_i X_i}{\sqrt{\frac{1}{n} \sum_i \left( X_i - \frac{1}{n} \sum_j X_j \right)^2}} = \frac{X - \mu_X}{\sqrt{\frac{1}{n} \sum_i (X_i - \mu_X)^2}} = \bar{X} \end{aligned}$$

# Network Training

## ❖ Solution 1: 0-mean and unit-deviation normalization

**Data Normalization**  
(convert to 0-mean  
and 1-deviation)

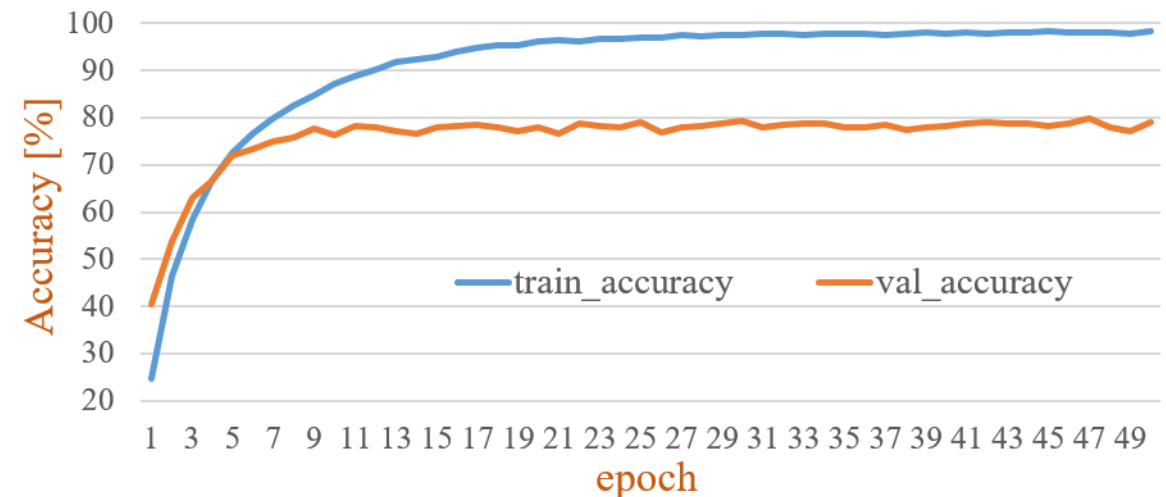
$$X = \frac{X - \mu_d}{\sigma_d}$$

$\mu_d$  is the mean of dataset

$\sigma_d$  is the deviation for the whole dataset

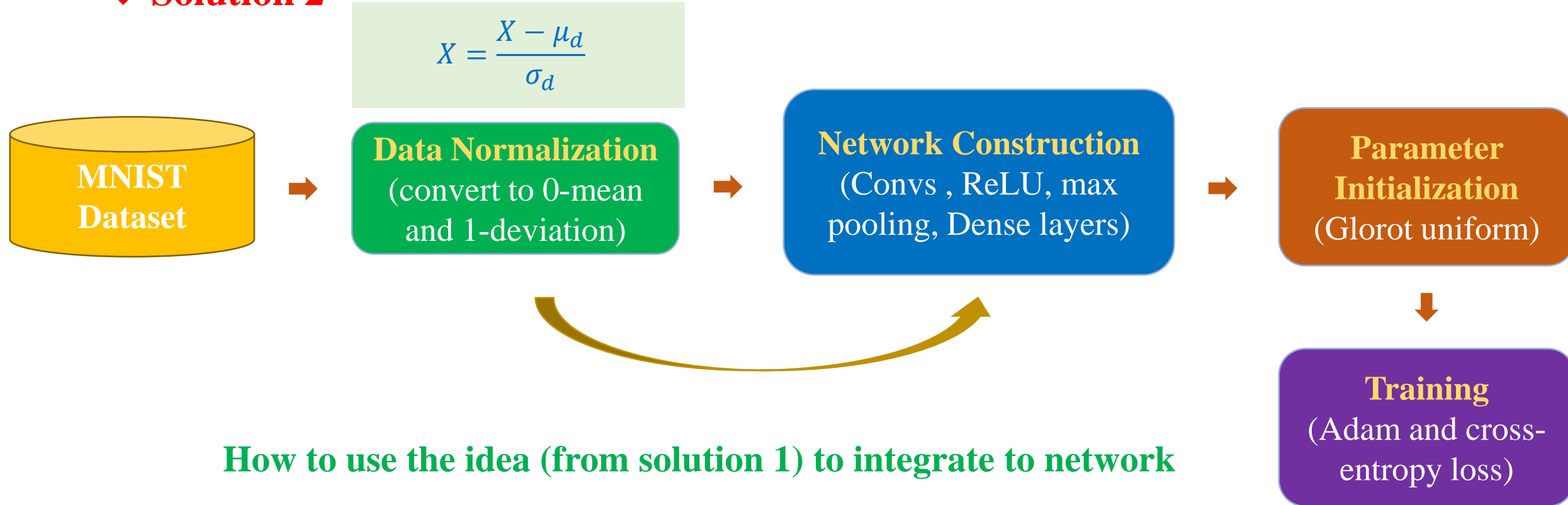
Normalize each channel separately

```
1 # normalize
2 mean = np.array([[[[125.3, 122.9, 113.8]]]])
3 std  = np.array([[[[62.9, 62.1, 66.7]]]])
4
5 x_train = (x_train - mean) / std
6 x_test  = (x_test - mean) / std
```



# Network Training

## ❖ Solution 2

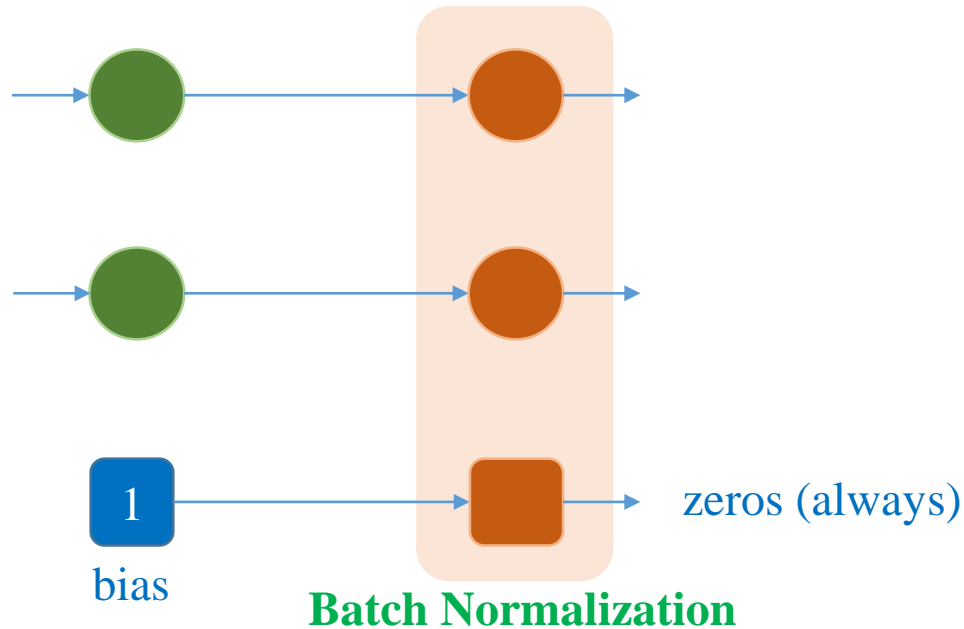


How to use the idea (from solution 1) to integrate to network

**Batch Normalization**

# Network Training

## ❖ Solution 2: Batch normalization



Do not need bias when using BN

$\mu$  and  $\sigma$  are updated in forward pass  
 $\gamma$  and  $\beta$  are updated in backward pass

Input data for a node in batch normalization layer

$$X = \{X_1, \dots, X_m\}$$

$m$  is mini-batch size

Compute mean and variance

$$\mu = \frac{1}{m} \sum_{i=1}^m X_i \quad \sigma^2 = \frac{1}{m} \sum_{i=1}^m (X_i - \mu)^2$$

Normalize  $X_i$

$$\hat{X}_i = \frac{X_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

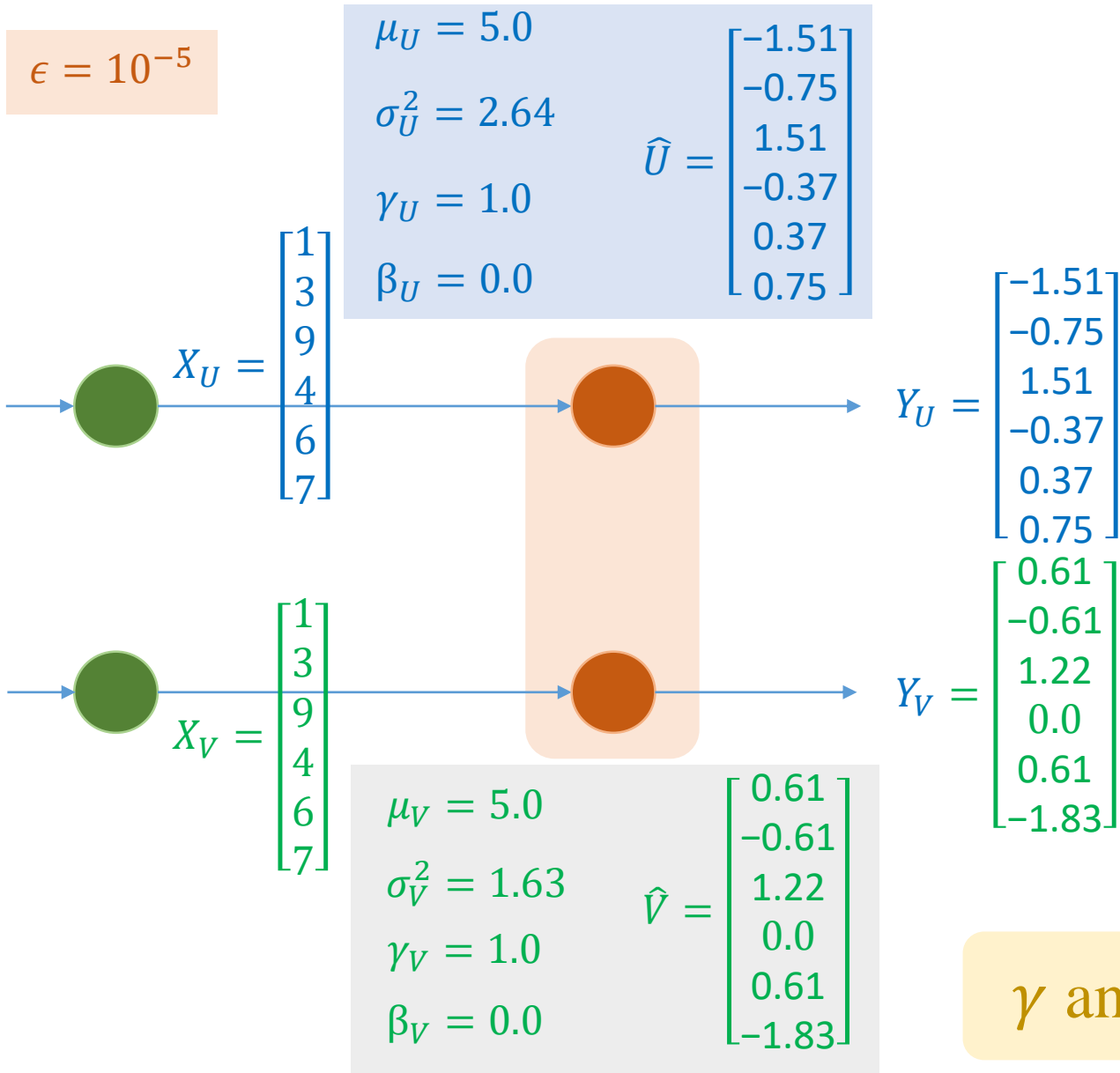
$\epsilon$  is a very small value

Scale and shift  $\hat{X}_i$

$$Y_i = \gamma \hat{X}_i + \beta$$

$\gamma$  and  $\beta$  are two learning parameters

# Network Training



## Solution 2: Batch normalization

Input data for a node in batch normalization layer

$$X = \{X_1, \dots, X_m\}$$

$m$  is mini-batch size

Compute mean and variance

$$\mu = \frac{1}{m} \sum_{i=1}^m X_i \quad \sigma^2 = \frac{1}{m} \sum_{i=1}^m (X_i - \mu)^2$$

Normalize  $X_i$

$$\hat{X}_i = \frac{X_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$\epsilon$  is a very small value

Scale and shift  $\hat{X}_i$

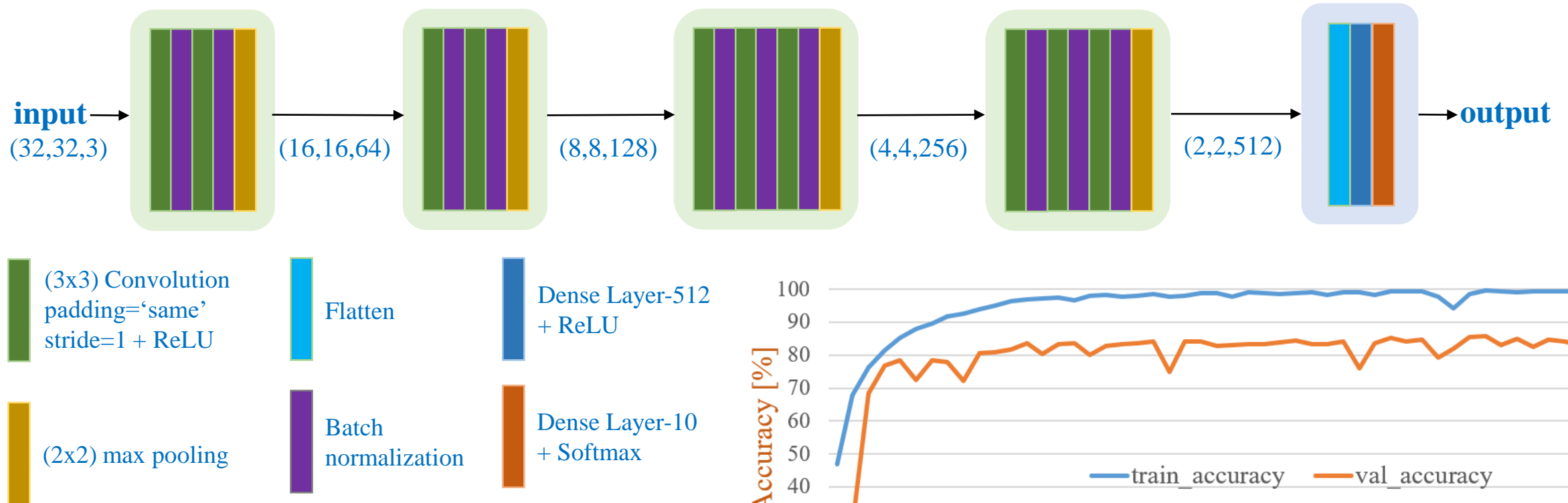
$$Y_i = \gamma \hat{X}_i + \beta$$

$\gamma$  and  $\beta$  are two learning parameters

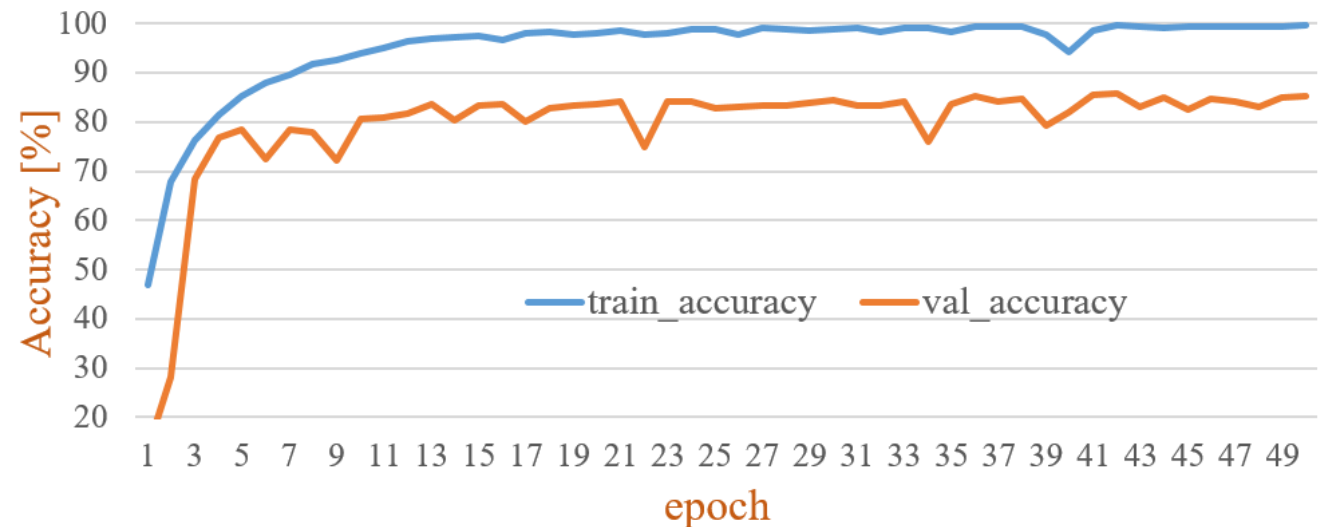
$\gamma$  and  $\beta$  are updated in training process

# Network Training

## ❖ Solution 2: Batch normalization



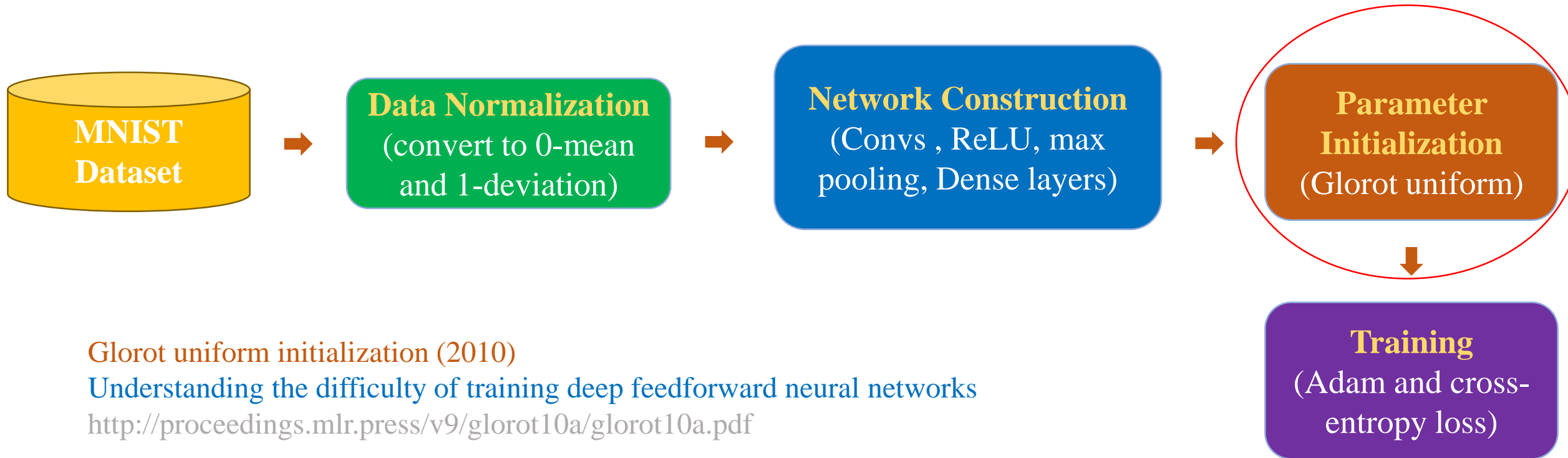
```
model.add(Conv2D(num_filter, kernel_size,  
                activation='relu'))  
model.add(keras.layers.BatchNormalization())
```



We could use BN before activation

# Network Training

## ❖ Solution 3: Use more robust initialization



Glorot uniform initialization (2010)

Understanding the difficulty of training deep feedforward neural networks

<http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>

He initialization (2015)

Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification

<https://arxiv.org/pdf/1502.01852.pdf>



# Network Training

## ❖ Solution 3: He Initialization

Glorot uniform initialization (2010)

$$W \sim U \left[ -\frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}} \right]$$

$n_j$  is #inputs in layer  $j$

Assuming activation functions are linear

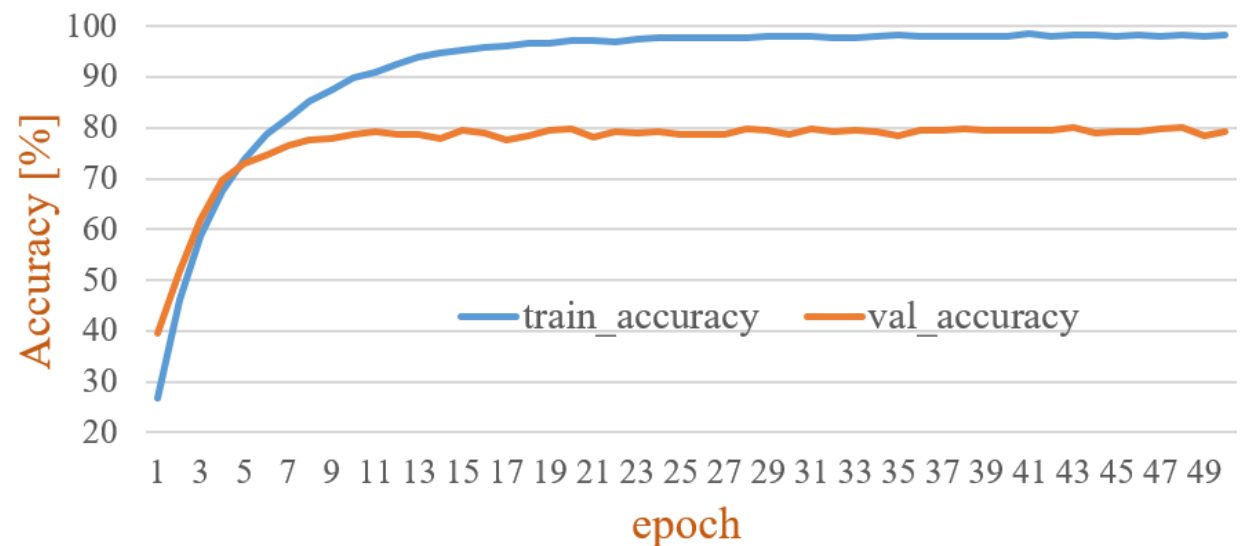
He initialization (2015)

Taking activation function into account

Adapt to ReLU activation

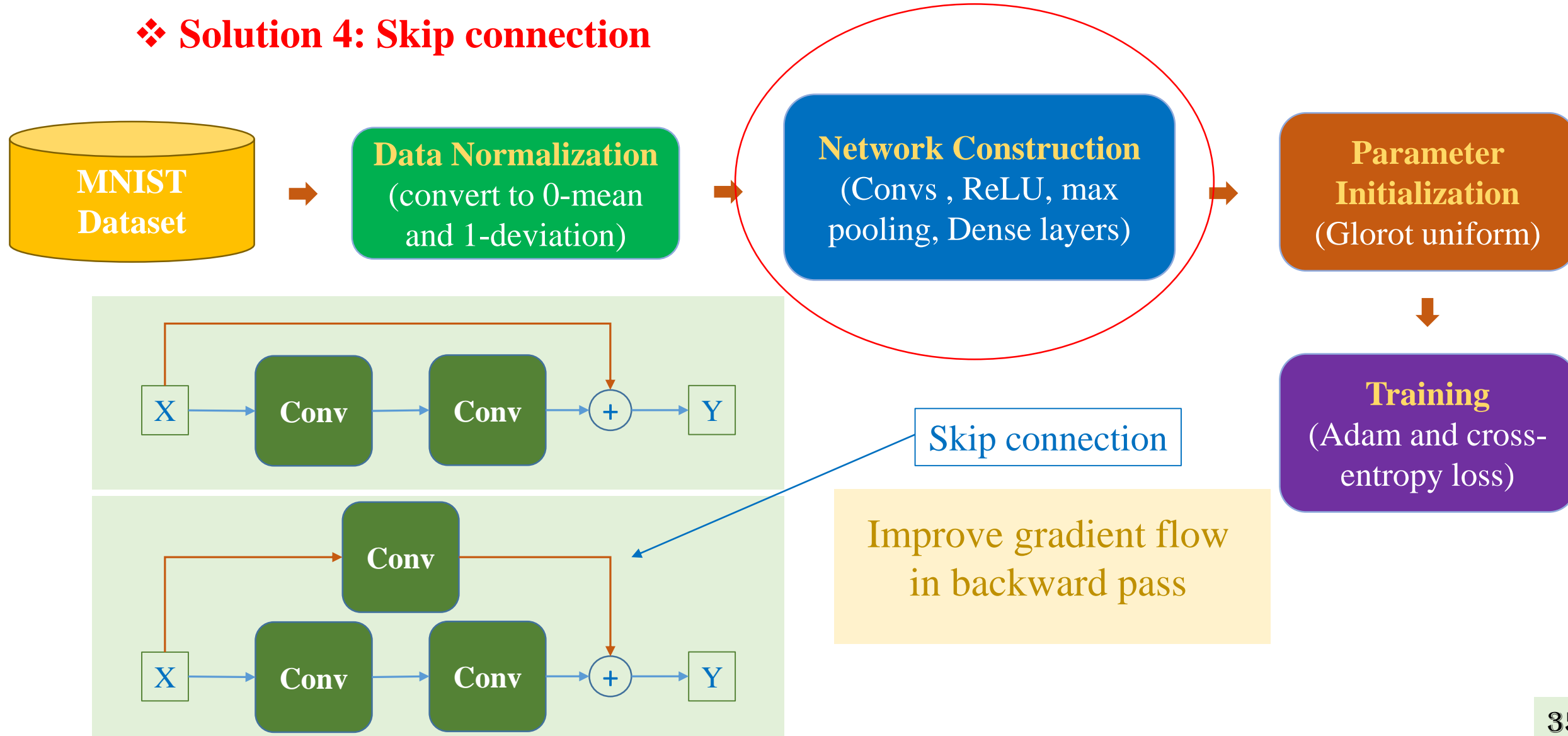
$$W \sim \mathcal{N} \left( 0, \frac{2}{n_j} \right)$$

```
initializer = keras.initializers.he_normal()  
Conv2D(num_filter, kernel_size, activation='relu',  
       kernel_initializer=initializer)
```



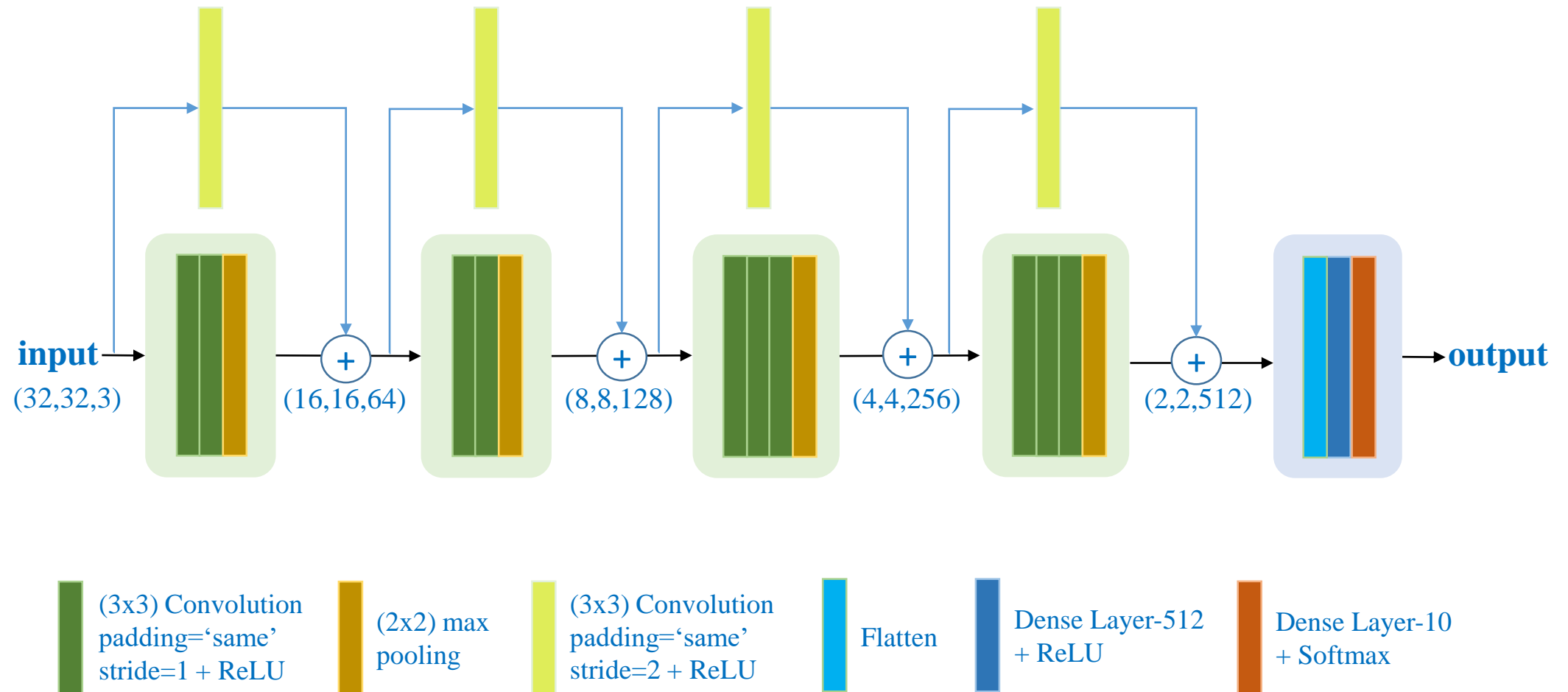
# Network Training

## ❖ Solution 4: Skip connection



# Network Training

## ❖ Solution 4: Skip connection

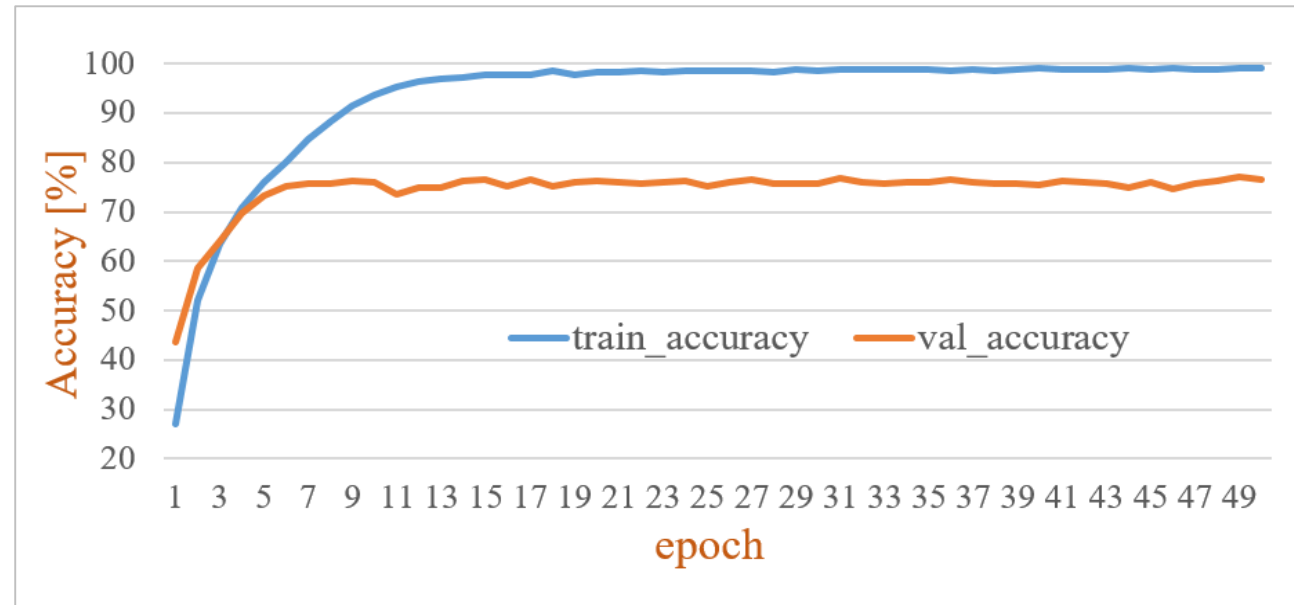
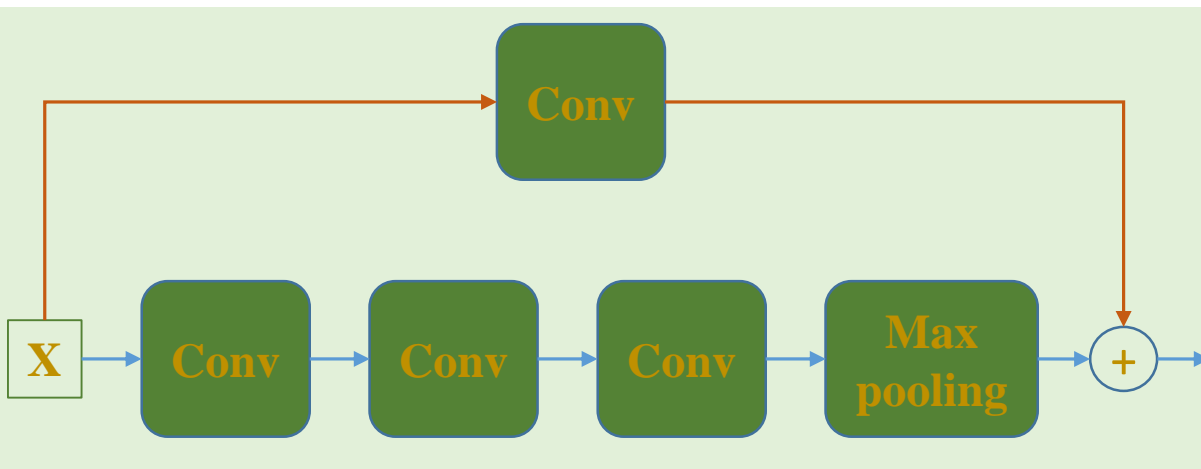


# Network Training

## ❖ Solution 4: Skip connection

```
# block 1
previous = x
x = keras.layers.Conv2D(64, (3, 3), padding='same', activation='relu')(x)
x = keras.layers.Conv2D(64, (3, 3), padding='same', activation='relu')(x)
x = keras.layers.Conv2D(64, (3, 3), padding='same', activation='relu')(x)
x = keras.layers.MaxPooling2D(2)(x)

residual = keras.layers.Conv2D(64, (3, 3), strides=2, padding='same', activation='relu')(previous)
x = keras.layers.add([x, residual])
```



There are several variants that use fully skip connection, concatenation, long skip connection

# Reading

---

## Skip connection

<https://theaisummer.com/skip-connections/>

## Trying to overfit Data

<http://karpathy.github.io/2019/04/25/recipe/>

# Summary

Use max pooling before ReLU to save some computations

