# Model Generalization

Quang-Vinh Dinh
Ph.D. in Computer Science

*Year 2020*

# Model Training

input (224,224,3) → Block-1 (112,112,64) → Block-2 (56,56,256) → Block-3 (28,28,512) → Block-4 (14,14,512) → Block-5 (7,7,512) → Dense Block → output

Legend:
- (3x3) Convolution padding='same' stride=1 + ReLU
- (2x2) max pooling
- Dense Layer + ReLU (4096)
- Dense Layer + Softmax (1000)

**Timeline:**

LeNet 1994 — AlexNet 2012 — Network-in-Network 2013 — VGG 2014 — Inception 2014 — U-Net 2015 — ResNet 2015 — SqueezeNet 2016 — DenseNet 2016 — MobileNets 2017 — EfficientNet 2019

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

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

#channels, wider, deeper, layer_i, resolution HxW, higher resolution

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

1

# Image Data

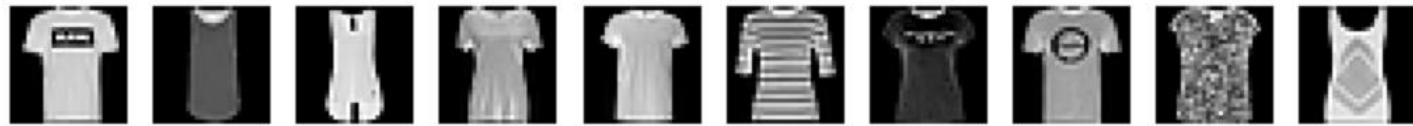**Fashion-MNIST dataset**

Grayscale images

Resolution=28x28

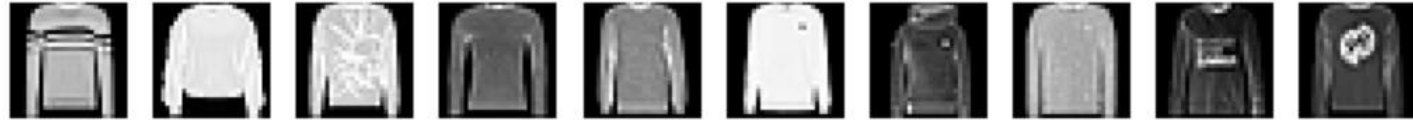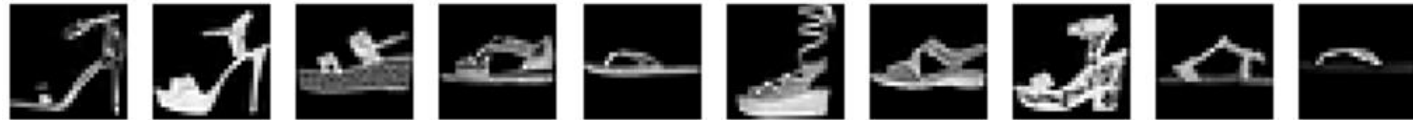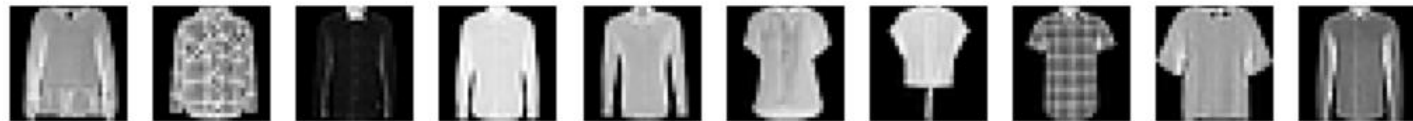Training set: 60000 samples

Testing set: 10000 samples
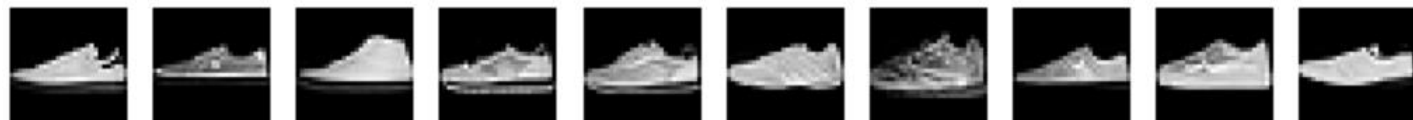
T-shirt

Trouser
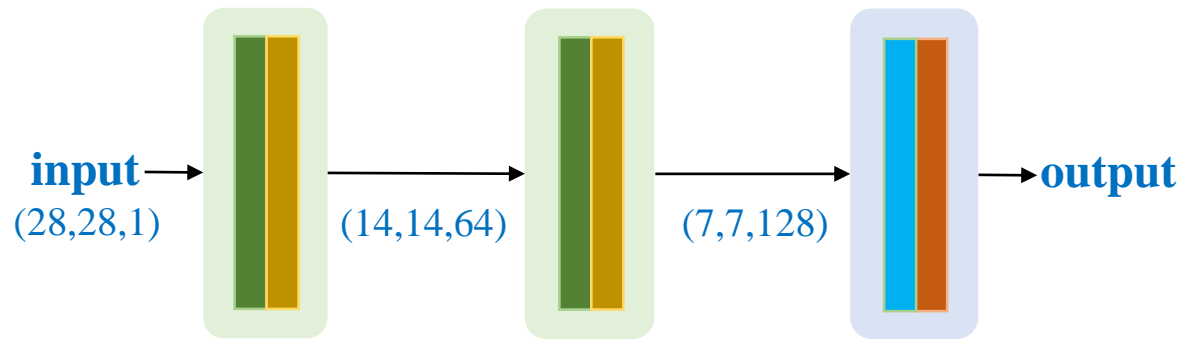
Pullover

Dress

Coat

Sandal

Shirt

Sneaker

Bag

Ankle Boot

# Network Training

## ❖ Fashion-MNIST dataset



input
(28,28,1)          (14,14,64)          (7,7,128)          output

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

Flatten

(2x2) max pooling

Dense Layer-10
+ Softmax

(7,7,128)

6272



3

# 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**



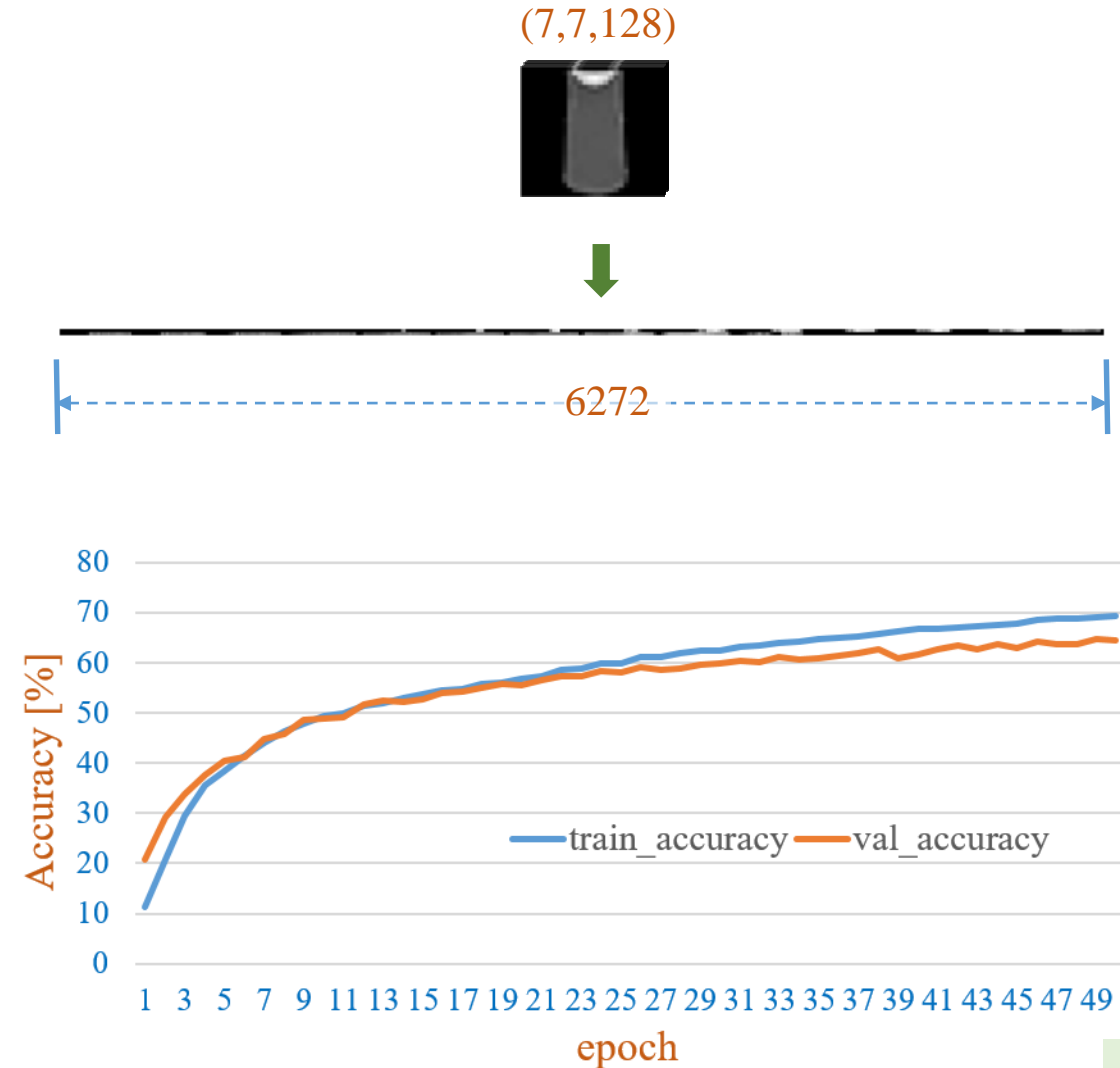**input** → (32,32,3) → (16,16,64) → (8,8,128) → **output**

(7,7,128)

6272

| | |
|---|---|
| (3x3) Convolution padding='same' stride=1 + Sigmoid | Flatten |
| (2x2) max pooling | Dense Layer-10 + Softmax |



Accuracy: 0.6930 - Val_accuracy: 0.6459

5

# Network Training

❖ **Cifar-10 dataset:**

    ❖ **Keep adding more layers**

**The network does not learn**

input (32,32,3) → (16,16,64) → (8,8,128) → (4,4,256) → output



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

■ Flatten

■ Dense Layer-512 + Sigmoid

■ (2x2) max pooling

■ Dense Layer-10 + Softmax

Accuracy: 0.0998 - Val_accuracy: 0.1000

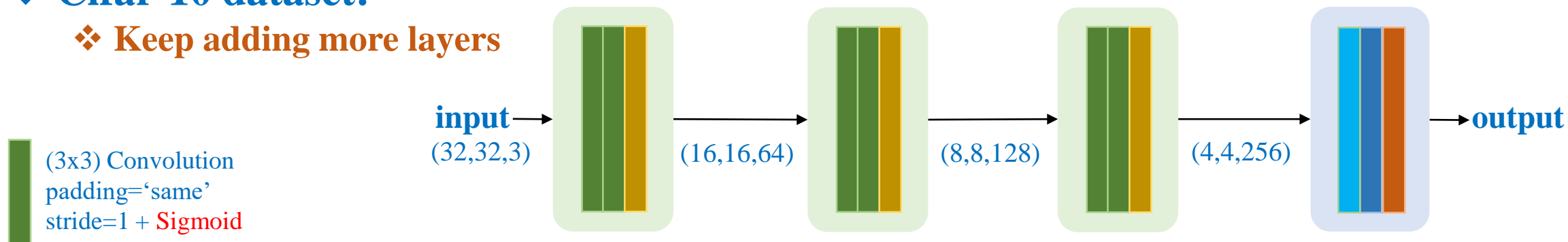# Network Training

❖ **Cifar-10 dataset:**

❖ **Keep adding more layers**



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

Dense Layer-512
+ Sigmoid

input → (32,32,3) → (16,16,64) → (8,8,128) → (4,4,256) → output

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

**Values are too small**

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')
```

8

# Network Training

❖ **Cifar-10 dataset:**
  ❖ **Use ReLU**

**Training Accuracy reaches up to 99%**

input → (16,16,64) → (8,8,128) → (4,4,256) → output
(32,32,3)



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

■ Flatten

■ Dense Layer-512 + ReLU

■ (2x2) max pooling

■ Dense Layer-10 + Softmax

Adding more layers; Hope reach to 100%

9

# Network Training

## Use ReLU and add more layers

input → (16,16,64) → (8,8,128) → (4,4,256) → (2,2,512) → output

(32,32,3)



Legend:

- (3x3) Convolution padding='same' stride=1 + ReLU
- (2x2) max pooling
- Flatten
- Dense Layer-512 + ReLU
- Dense Layer-10 + Softmax

Network does not learn again

# Network Training

❖ **Summary of the current network**



**Cifar-10 Dataset**

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

**Network Construction** (Convs , ReLU, max pooling, Dense layers)

**Parameter Initialization** (Glorot uniform)

**Training** (Adam and cross-entropy loss)

Network does not learn

# Network Training

❖ **Solution 1: Observation**



**MNIST Dataset** → **Data Normalization** (scale to [0,1)) → **Network Construction** (Convs , ReLU, max pooling, Dense layers) → **Parameter Initialization** (Glorot uniform) → **Training** (Adam and cross-entropy loss)

The current network performs excellently for MNIST dataset

# Network Training

❖ **Solution 1: Idea**

**Current Network**

Train (failed)

**Cifar-10 Dataset**

$>$

complex

**Current Network**

Train (ok)

**MNIST Dataset**

How to reduce the complexity
of the Cifar-10 dataset

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

0                                    255

0     1

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

$X =$

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

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

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

13

# 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

$$Y = aX + b$$

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



$$Y = aX + b$$

$$\bar{Y} = \frac{Y - \mu_Y}{\sigma_Y} = \frac{(aX + b) - \frac{1}{n}\sum_i(aX_i + b)}{\sqrt{\frac{1}{n}\sum_i\left((aX_i + b) - \frac{1}{n}\sum_i(aX_i + b)\right)^2}}$$

$$= \frac{aX - \frac{1}{n}\sum_i aX_i}{\sqrt{\frac{1}{n}\sum_i\left(aX_i - \frac{1}{n}\sum_j aX_j\right)^2}}$$

$$= \frac{X - \frac{1}{n}\sum_i X_i}{\sqrt{\frac{1}{n}\sum_i\left(X_i - \frac{1}{n}\sum_j X_j\right)^2}} = \frac{X - \mu_X}{\sqrt{\frac{1}{n}\sum_i(X_i - \mu_X)^2}} = \bar{X}$$

14

# Network Training

❖ **Solution 2**

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

**MNIST Dataset** → **Data Normalization** (convert to 0-mean and 1-deviation) → **Network Construction** (Convs , ReLU, max pooling, Dense layers) → **Parameter Initialization** (Glorot uniform) → **Training** (Adam and cross-entropy loss)

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

**Batch Normalization**

# Network Training

❖ **Solution 2: Batch normalization**



**Batch Normalization**

zeros (always)

bias

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

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

$m$ is mini-batch size

Compute mean and variance

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

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**



bias

**Batch Normalization**

zeros (always)

**What if**

$$\gamma = \sqrt{\sigma^2 + \epsilon} \text{ and } \beta = \mu$$

Input data for a node in batch normalization layer

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

$m$ is mini-batch size

Compute mean and variance

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

Normalize $X_i$

$$\hat{X}_i = \frac{X_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$
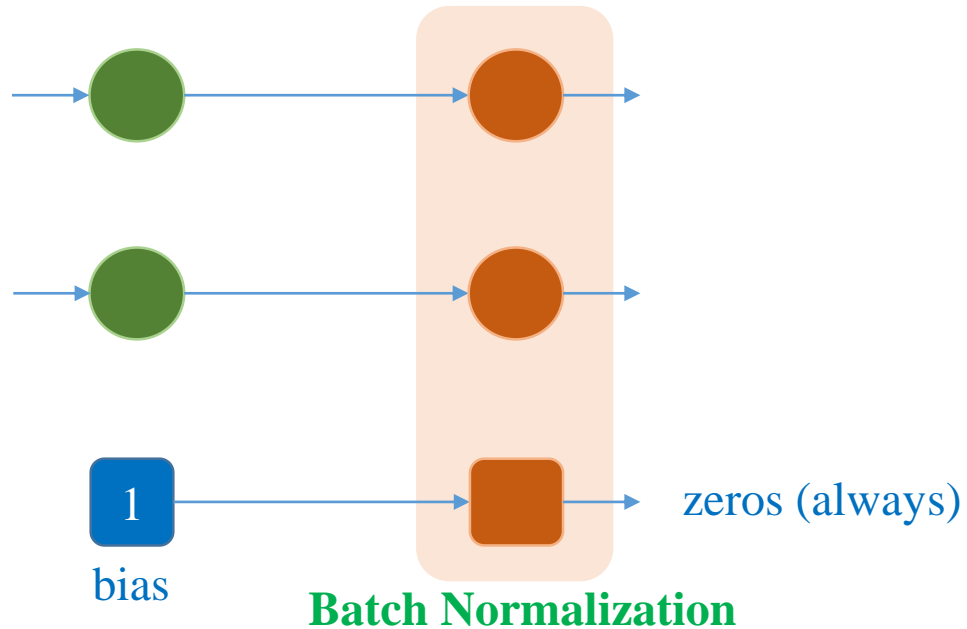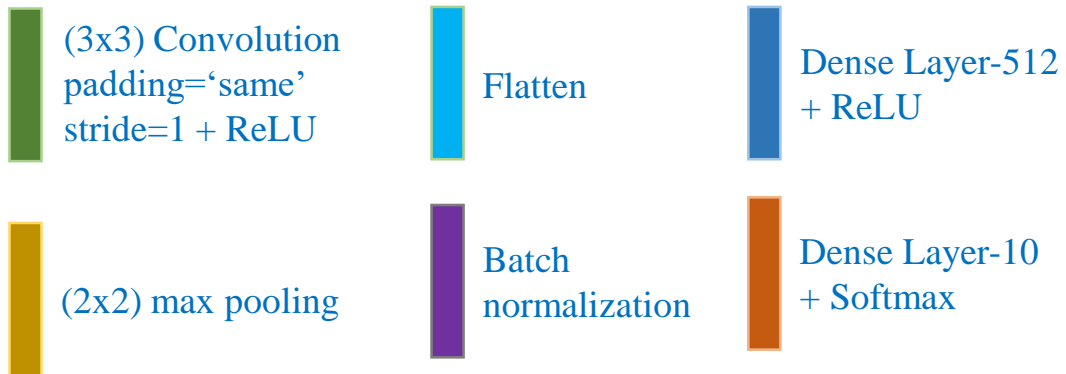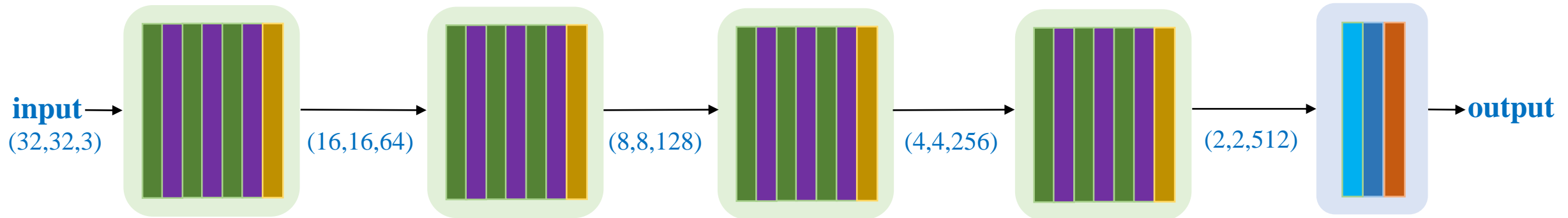
$\epsilon$ is a very small value

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
(32,32,3)

(16,16,64)

(8,8,128)

(4,4,256)

(2,2,512)

output

Legend:
- (3x3) Convolution padding='same' stride=1 + ReLU
- Flatten
- Dense Layer-512 + ReLU
- (2x2) max pooling
- Batch normalization
- Dense Layer-10 + Softmax

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

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

model.add(keras.layers.Conv2D(64, (3, 3),
                              strides=1, padding='same',
                              activation = 'relu'))
model.add(tf.keras.layers.BatchNormalization())
model.add(keras.layers.Conv2D(64, (3, 3),
                              strides=1, padding='same',
                              activation = 'relu'))
model.add(tf.keras.layers.BatchNormalizationCon())
model.add(keras.layers.Conv2D(64, (3, 3),
                              strides=1, padding='same',
                              activation='relu'))
model.add(tf.keras.layers.BatchNormalization())
model.add(keras.layers.MaxPooling2D(2))
```

18

# Network Training

❖ **Solution 2: Batch normalization**

Speed up training

Reduce the dependence on initial weights

Model Generalization

Input data for a node in batch normalization layer

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

$m$ is mini-batch size

Compute mean and variance

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

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**

**Backward**

Input data for a node in batch normalization layer

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

$m$ is mini-batch size

Compute mean and variance

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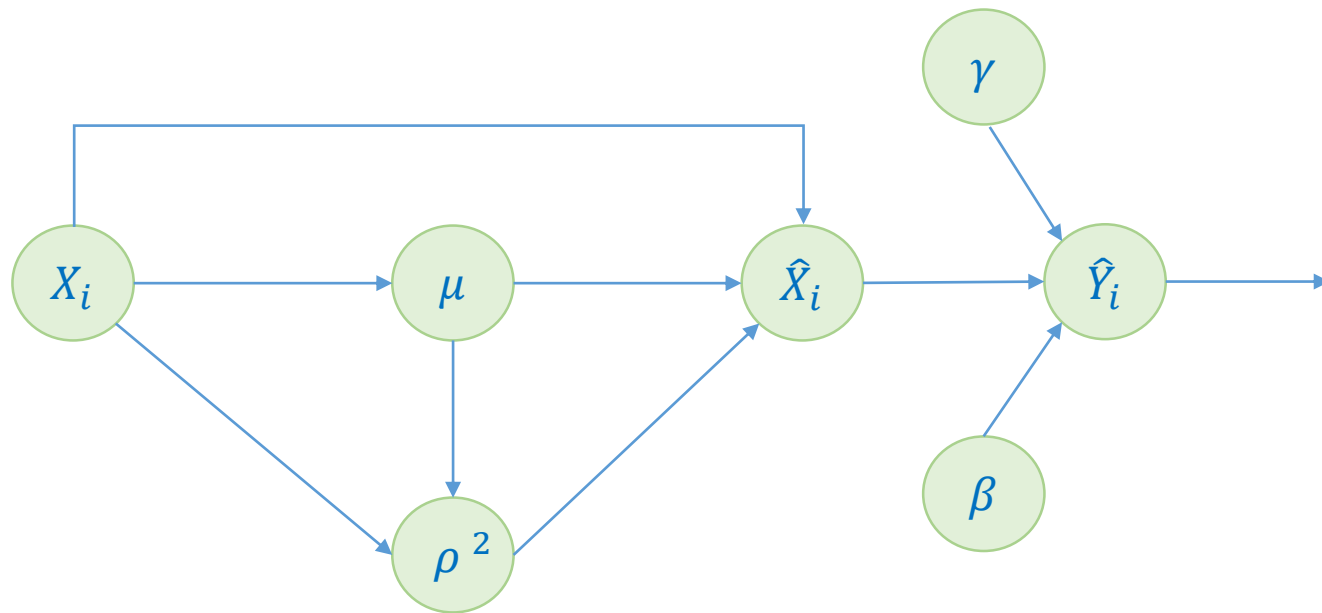
Normalize $X_i$

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

$\epsilon$ 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

20

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

$$\hat{X}_i = \frac{X_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \qquad\qquad Y_i = \gamma\hat{X}_i + \beta$$
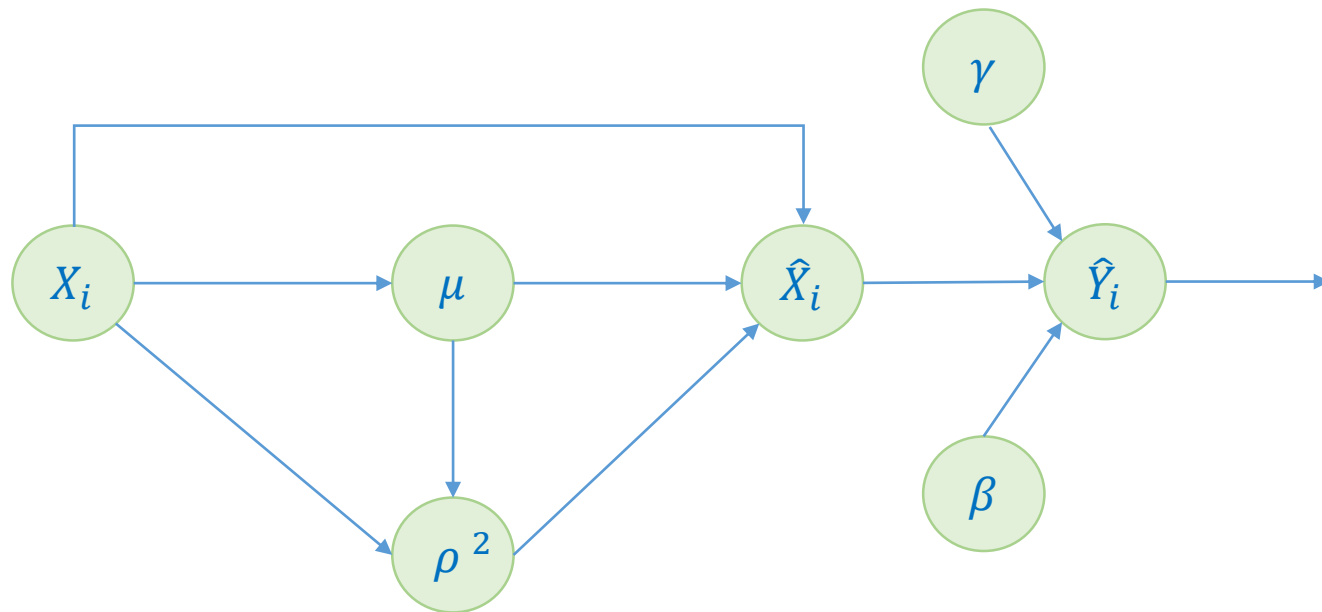
$$\mu = \frac{1}{m}\sum_{i=1}^{m} X_i \qquad \sigma^2 = \frac{1}{m}\sum_{i=1}^{m}(X_i - \mu)^2$$
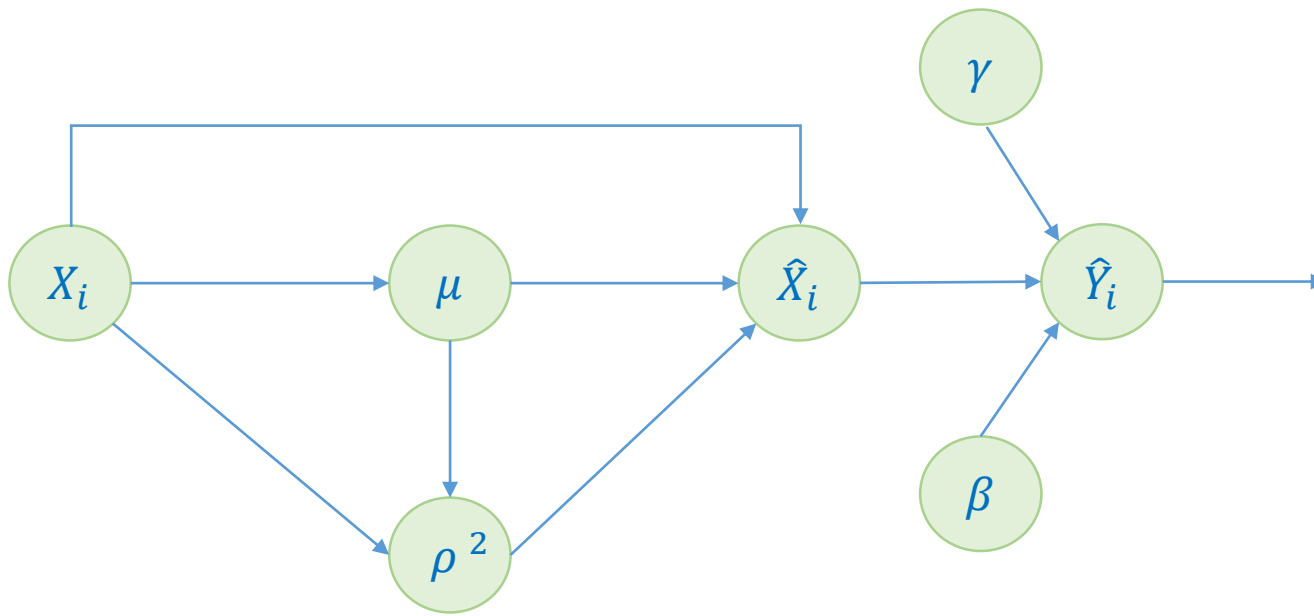
$$\hat{X}_i = \frac{X_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \qquad Y_i = \gamma\hat{X}_i + \beta$$

$$\frac{\partial L}{\partial \gamma} = \sum_{i=1}^{m} \frac{\partial L}{\partial Y_i}\hat{X}_i$$

$$\frac{\partial L}{\partial \beta} = \sum_{i=1}^{m} \frac{\partial L}{\partial Y_i}$$

$$\frac{\partial L}{\partial \hat{X}_i} = \frac{\partial L}{\partial Y_i}\gamma$$

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

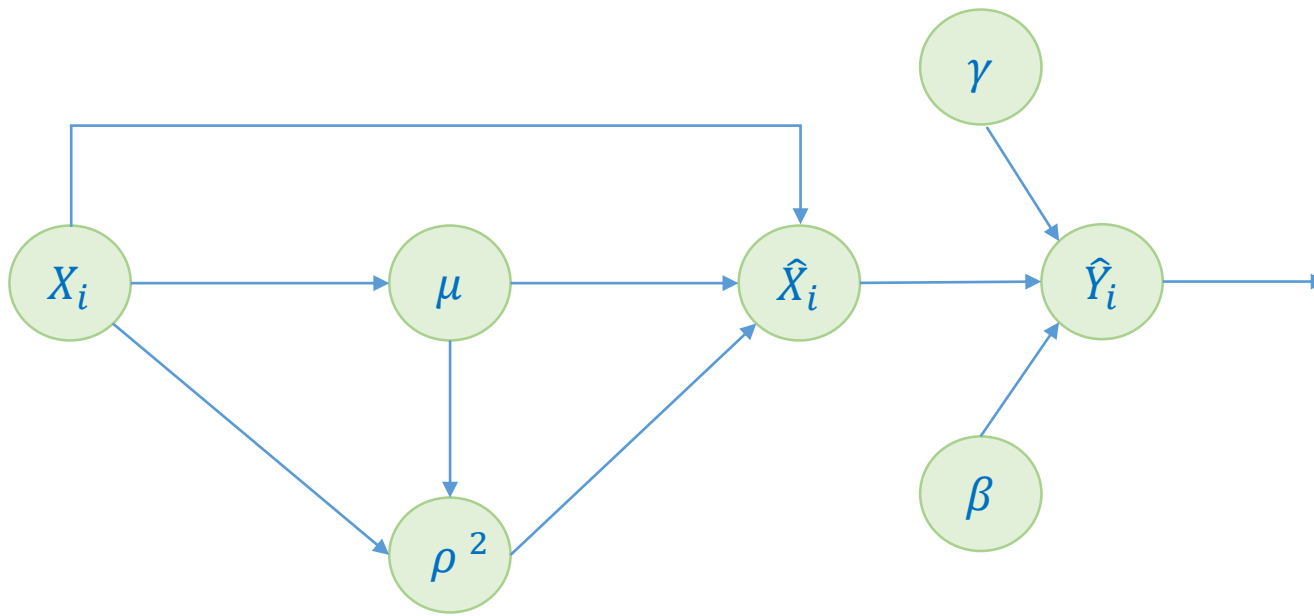$$\hat{X}_i = \frac{X_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \qquad Y_i = \gamma \hat{X}_i + \beta$$

$$\frac{\partial L}{\partial \gamma} = \sum_{i=1}^{m} \frac{\partial L}{\partial Y_i} \hat{X}_i \qquad \frac{\partial L}{\partial \beta} = \sum_{i=1}^{m} \frac{\partial L}{\partial Y_i} \qquad \frac{\partial L}{\partial \hat{X}_i} = \frac{\partial L}{\partial Y_i} \gamma$$

$$\frac{\partial L}{\partial \sigma^2} = \sum_{i=1}^{m} \frac{\partial L}{\partial \hat{X}_i} \frac{\partial \hat{X}_i}{\partial \sigma^2} = \sum_{i=1}^{m} \frac{\partial L}{\partial \hat{X}_i}(X_i - \mu)\frac{-1}{2}(\sigma^2 + \epsilon)^{\frac{-3}{2}}$$

$$\frac{\partial L}{\partial \mu} = \sum_{i=1}^{m} \frac{\partial L}{\partial \hat{X}_i}\frac{-1}{\sqrt{\sigma^2 + \epsilon}} + \frac{\partial L}{\partial \sigma^2}\frac{1}{m}\sum_{i=1}^{m} 2(X_i - \mu)$$

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

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

$$\frac{\partial L}{\partial \gamma} = \sum_{i=1}^{m} \frac{\partial L}{\partial Y_i}\hat{X}_i$$

$$\frac{\partial L}{\partial \beta} = \sum_{i=1}^{m} \frac{\partial L}{\partial Y_i}$$

$$\frac{\partial L}{\partial \hat{X}_i} = \frac{\partial L}{\partial Y_i}\gamma$$

$$\frac{\partial L}{\partial \mu} = \sum_{i=1}^{m} \frac{\partial L}{\partial \hat{X}_i}\frac{-1}{\sqrt{\sigma^2 + \epsilon}} + \frac{\partial L}{\partial \sigma^2}\frac{1}{m}\sum_{i=1}^{m} 2(X_i - \mu)$$

$$\frac{\partial L}{\partial \sigma^2} = \sum_{i=1}^{m} \frac{\partial L}{\partial \hat{X}_i}\frac{\partial \hat{X}_i}{\partial \sigma^2} = \sum_{i=1}^{m} \frac{\partial L}{\partial \hat{X}_i}(X_i - \mu)\frac{-1}{2}(\sigma^2 + \epsilon)^{\frac{-3}{2}}$$

$$\frac{\partial L}{\partial X_i} = \frac{\partial L}{\partial \hat{X}_i}\frac{\partial \hat{X}_i}{\partial X_i} + \frac{\partial L}{\partial \mu}\frac{\partial \mu}{\partial X_i} + \frac{\partial L}{\partial \sigma^2}\frac{\partial \sigma^2}{\partial X_i}$$
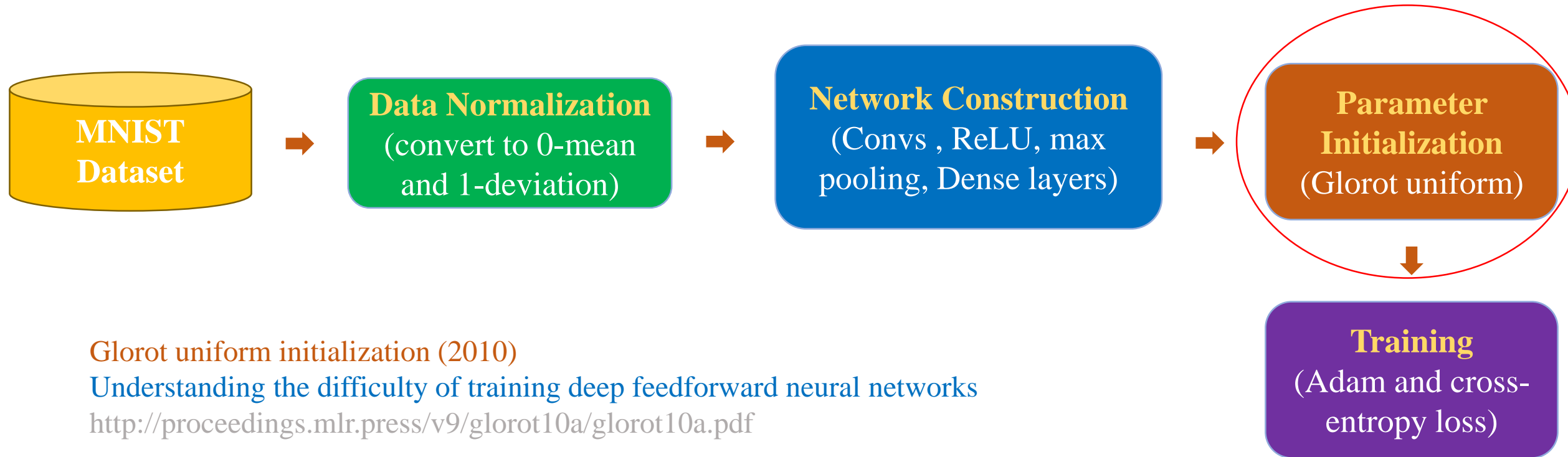
$$\frac{\partial \hat{X}_i}{\partial X_i} = \frac{1}{\sqrt{\sigma^2 + \epsilon}}$$

$$\frac{\partial \mu}{\partial X_i} = \frac{1}{m}$$

$$\frac{\partial \sigma^2}{\partial X_i} = \frac{2(X_i - \mu)}{m}$$

# Network Training

❖ **Solution 3: Use more robust initialization**



MNIST Dataset → Data Normalization (convert to 0-mean and 1-deviation) → Network Construction (Convs , ReLU, max pooling, Dense layers) → Parameter Initialization (Glorot uniform) → Training (Adam and cross-entropy loss)

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**

He initialization (2015)

Adapt to ReLU activation

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

$$E(XY) = E(X)E(Y)$$
$$var(XY) = var(X)var(Y) + var(X)\big(E(Y)\big)^2 + var(y)\big(E(X)\big)^2$$

# Network Training

❖ **Solution 4: Skip connection**



| | |
|---|---|
| ■ (green) | (3x3) Convolution padding='same' stride=1 + ReLU |
| ■ (orange) | (2x2) max pooling |
| ■ (yellow-green) | (3x3) Convolution padding='same' stride=2 + ReLU |
| ■ (light blue) | Flatten |
| ■ (blue) | Dense Layer-512 + ReLU |
| ■ (dark orange) | Dense Layer-10 + Softmax |

# Network Training

❖ **Solution 4: Skip connection**

**Backward**

# Model Generalization

# Model Generalization

**Cifar-10 Dataset**

input (32,32,3) → (16,16,32) → (8,8,64) → (4,4,128) → (2,2,256) → output

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

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

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

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

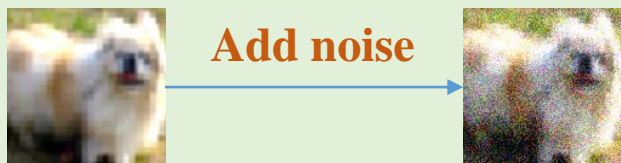(3x3) Convolution padding='same' stride=1 + ReLU

(2x2) max pooling

Flatten

Dense Layer-10 + Softmax

Dense Layer-512 + ReLU

Aim to reduce this gap



train_accuracy    val_accuracy

31

# Model Generalization

❖ **Trick 1: 'Learn hard, ' – randomly add noise to training data**



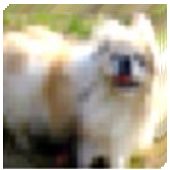Add noise

**In Keras**

```
1  if tf.random.uniform(()) > 0.5:
2          noise = tf.random.normal((32, 32, 3))/100.0
3          image = image+noise
4
5      return image, label
```



epoch

**val_accuracy increases from ~80.2% to ~80.9%**

# Model Generalization

❖ **Trick 2: Batch normalization**

mini-batch 1          mini-batch 2

$$(\mu_1, \sigma_1) \neq (\mu_2, \sigma_2)$$

very
likely

**Add noise to the output of BN layers**

Input data for a node in batch normalization layer

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

$m$ is mini-batch size

Compute mean and variance

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

Normalize $X_i$

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

$\epsilon$ is a very small value

Scale and shift $\hat{X}_i$

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

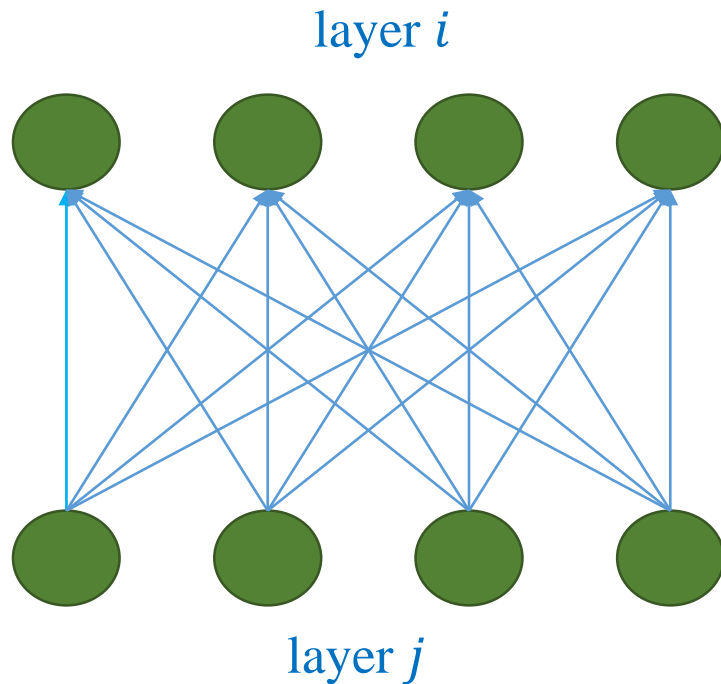$\gamma$ and $\beta$ are two learning parameters

# Model Generalization

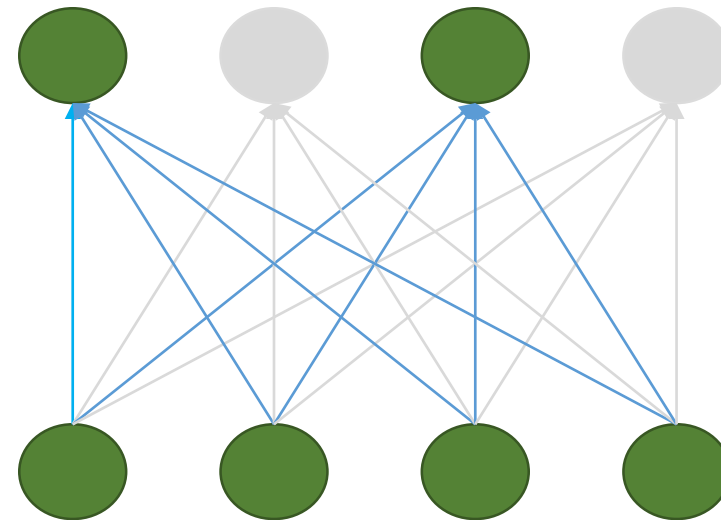❖ **Trick 2: Batch normalization**



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

Flatten

(2x2) max pooling

Dense Layer-512 + ReLU

Batch normalization

Dense Layer-10 + Softmax

**val_accuracy increases from ~80.9% to ~82%**

# Model Generalization

❖ **Trick 3: Dropout**

layer $i$



layer $j$

Apply dropout 50% to layer $i$



~50% nodes randomly selected in the $i^{th}$ layer are set to zeros (kind of noise adding)
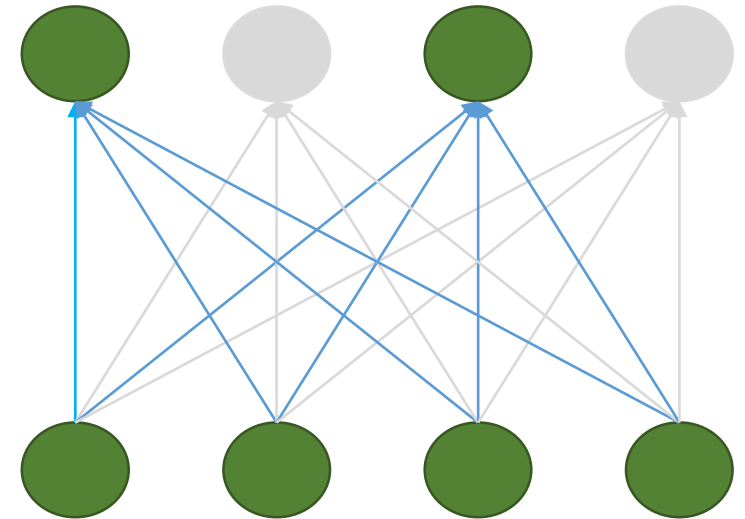
# Model Generalization

❖ **Trick 3: Dropout**

$$f(k, h) = \begin{cases} p & if \quad k = 1 \\ 1 - p & if \quad k = 0 \end{cases}$$

$$a = D \odot \sigma(Z)$$

https://deepnotes.io/dropout

$$\frac{\partial C}{\partial X} = \frac{\partial C}{\partial y} \times \frac{\partial y}{\partial X}$$

$$= \frac{\partial C}{\partial y} \times \frac{\partial \begin{cases} X_{ij} & if \quad D_{ij} = 1 \\ 0 & if \quad D_{ij} = 0 \end{cases}}{\partial X}$$

$$= \frac{\partial C}{\partial y} \times \begin{cases} 1 & if \quad D_{ij} = 1 \\ 0 & if \quad D_{ij} = 0 \end{cases}$$

$$= \frac{\partial C}{\partial y} \times D$$

Apply dropout 50% to layer $i$



~50% nodes randomly selected in the $i^{th}$ layer are set to zeros
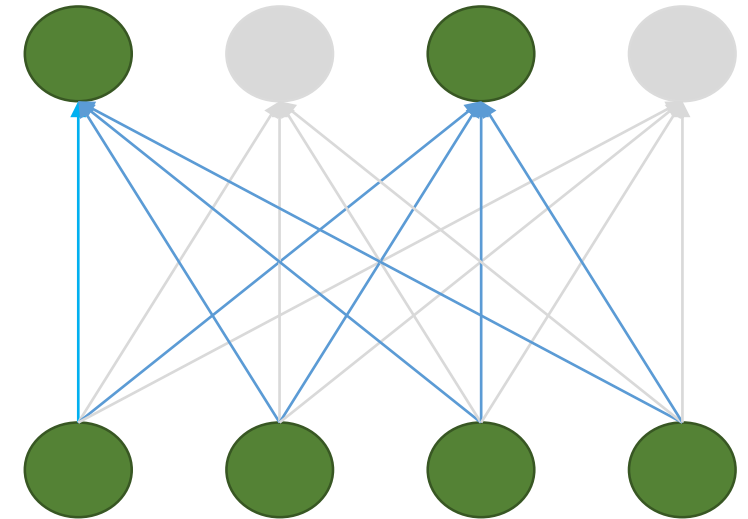
# Model Generalization

❖ **Trick 3: Dropout**

```python
class Dropout():

    def __init__(self,prob=0.5):
        self.prob = prob
        self.params = []


    def forward(self,X):
        self.mask = np.random.binomial(1,self.prob,size=X.shape) / self.prob
        out = X * self.mask
        return out.reshape(X.shape)


    def backward(self,dout):
        dX = dout * self.mask
        return dX,[]
```
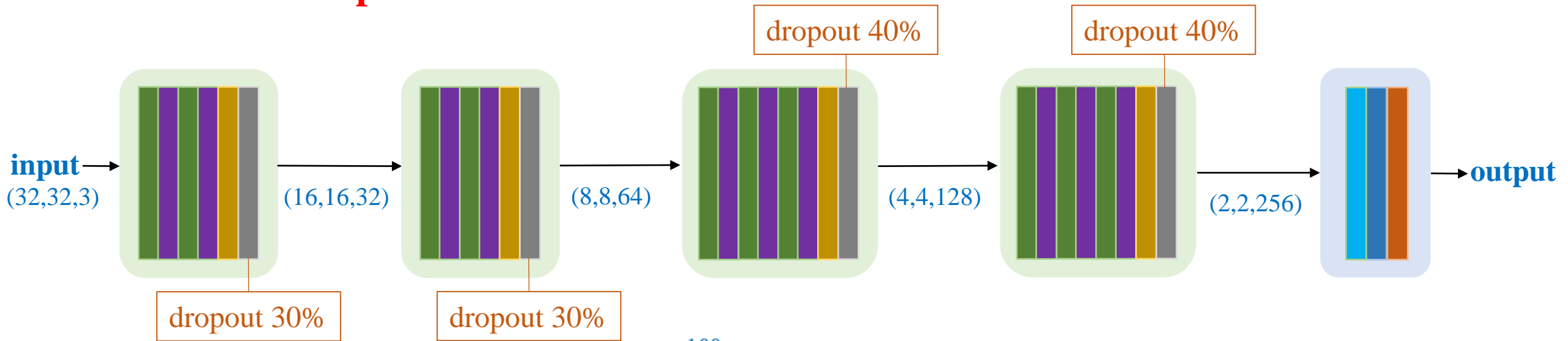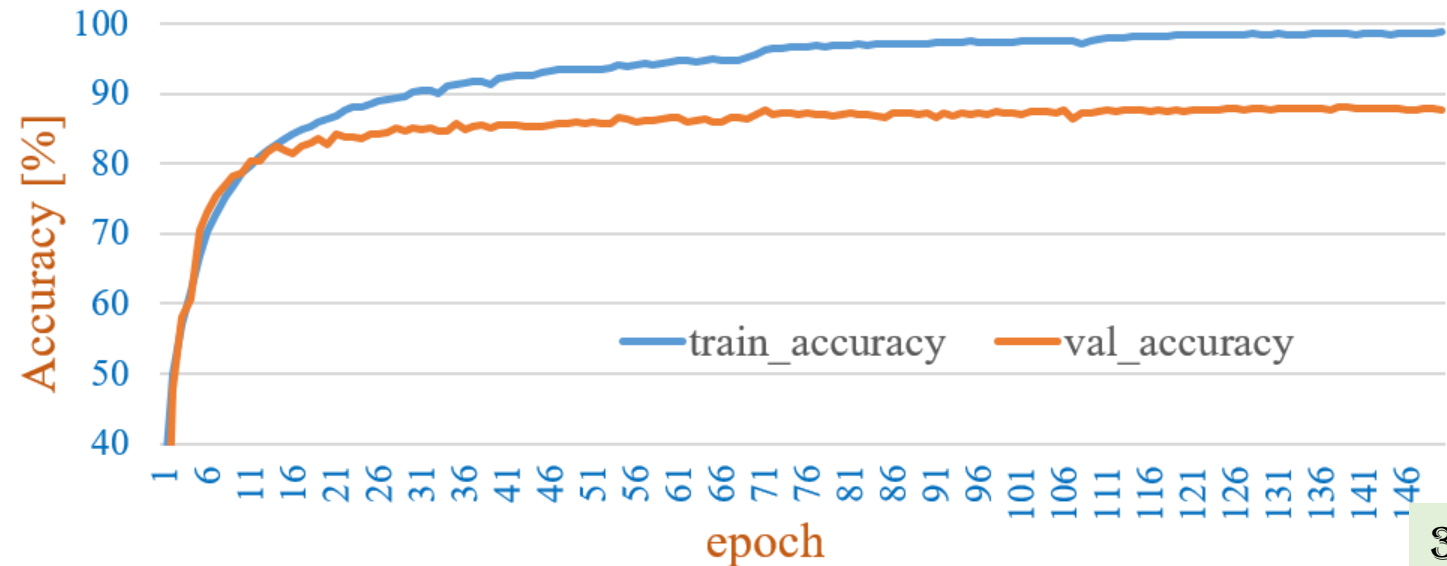
Apply dropout 50% to layer $i$



~50% nodes randomly selected in the $i^{th}$ layer are set to zeros

https://deepnotes.io/dropout

# Model Generalization

❖ **Trick 3: Dropout**



**val_accuracy increases from ~82% to ~87.9%**

# Model Generalization

❖ **Trick 4: Kernel regularization**

$$L = crossentropy + \lambda_1 \|W\| + \lambda_2 \|W\|^2$$

$L_1$ regularization    $L_2$ regularization
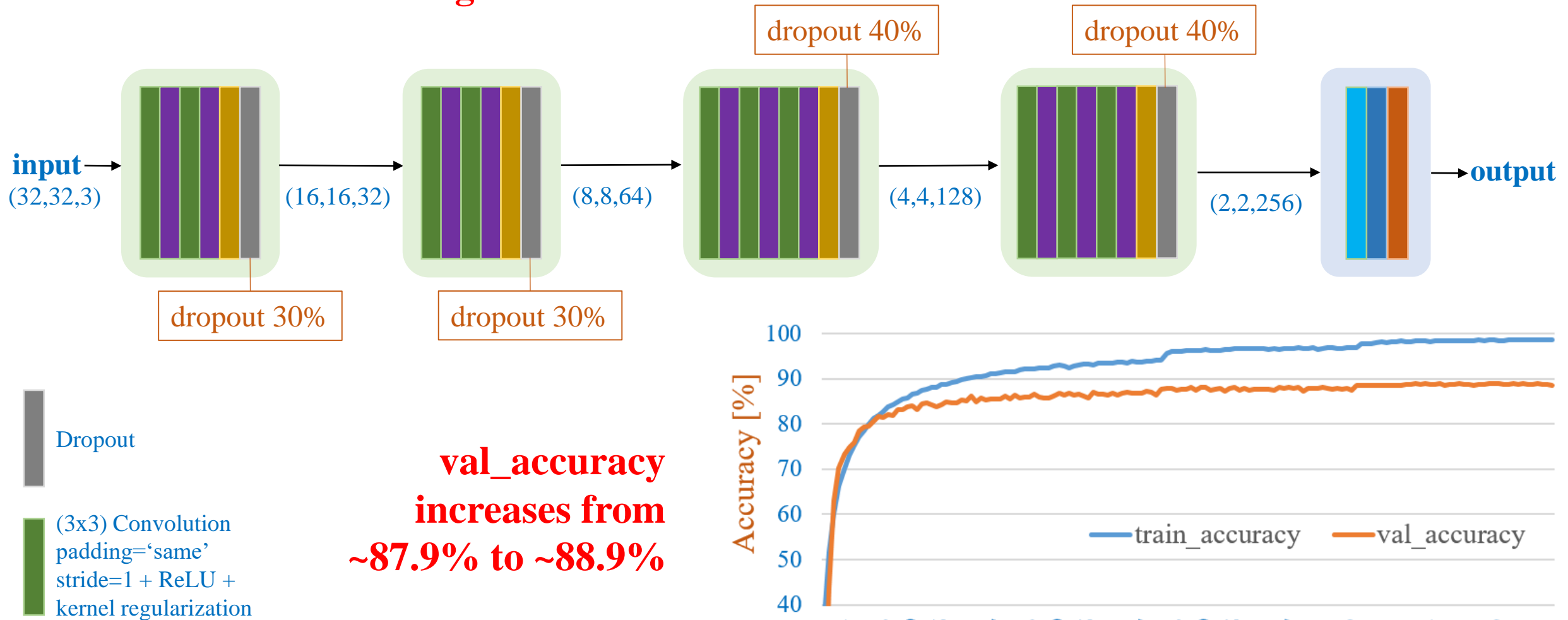
Prevent network from
focusing on specific features

Smaller weights
→ simpler models

In keras

```
Conv2D(32, (3,3), padding='same', activation='relu',
        kernel_regularizer=regularizers.l1_l2(decay1,decay2))
```
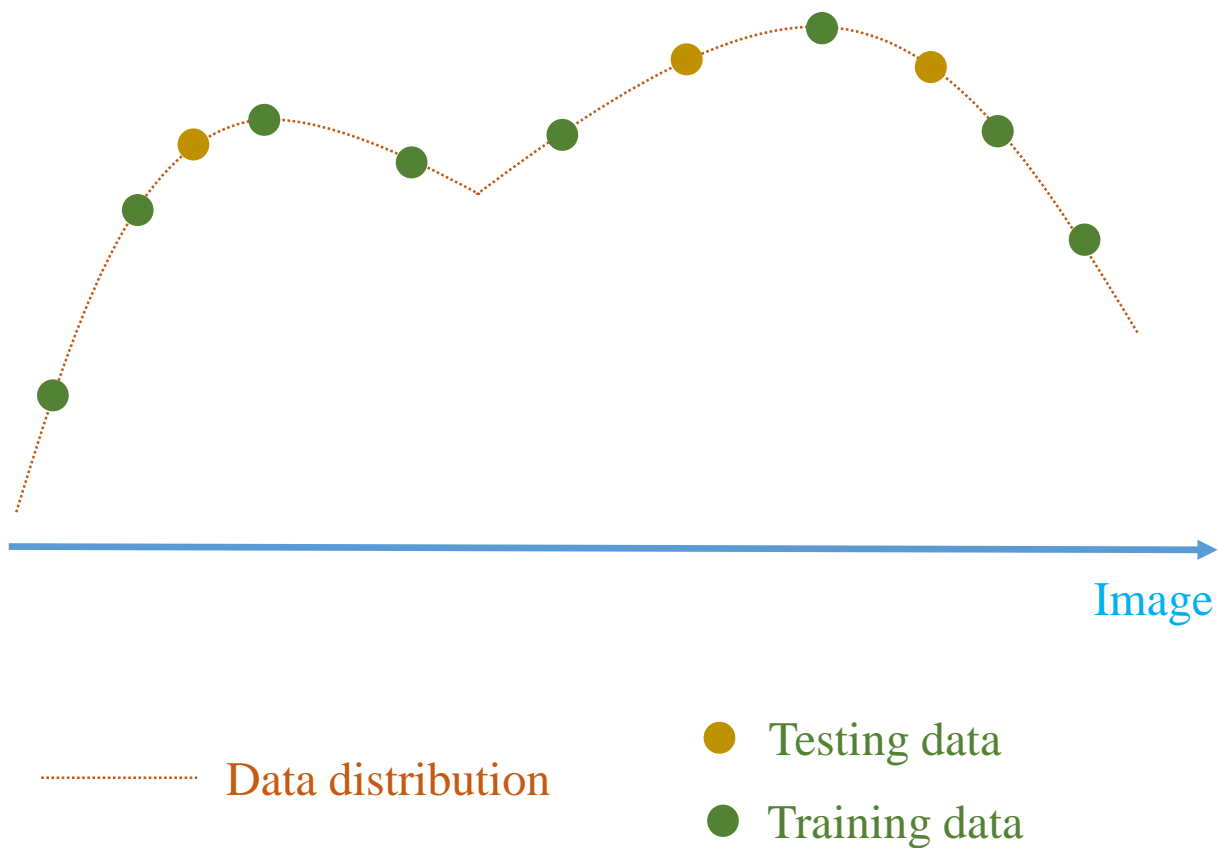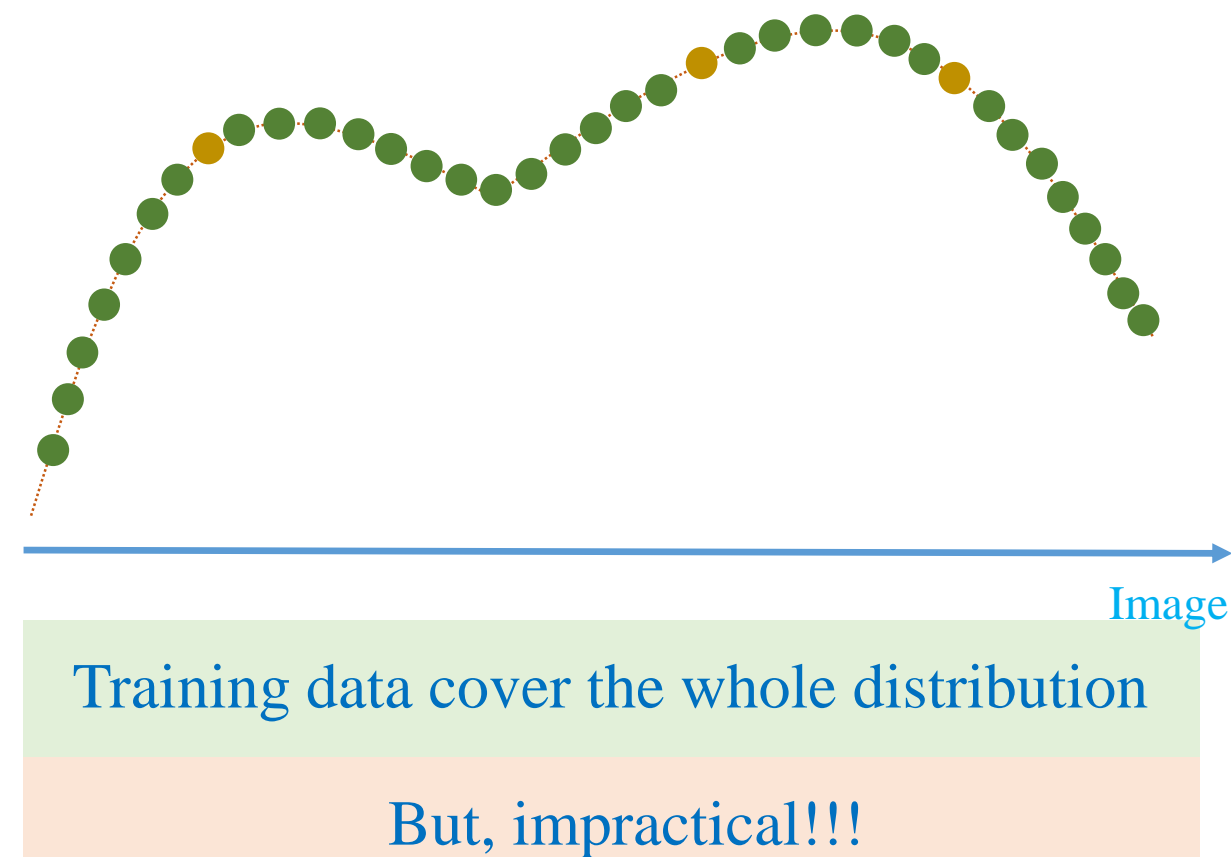
# Model Generalization

❖ **Trick 5: Data augmentation**

A normal case

A perfect case: Have unlimited training



Image

Image

- ⬤ (yellow) Testing data
- ⬤ (green) Training data
- ⋯⋯⋯ Data distribution

Training data cover the whole distribution

But, impractical!!!

# Model Generalization

❖ **Trick 5: Data augmentation**



Increase data by altering the training data

horizontal flip

rotate

crop and resize

Image

Data distribution

● Testing data

● Training data

● Augmented training data

# Model Generalization

❖ **Trick 5: Data augmentation**

**Horizontal flip**



**Horizontal flip + crop-and-resize**



**val_accuracy reaches to ~90.7%**

**val_accuracy reaches to ~91.2%**

# Model Generalization

❖ **Trick 6: Instance normalization**
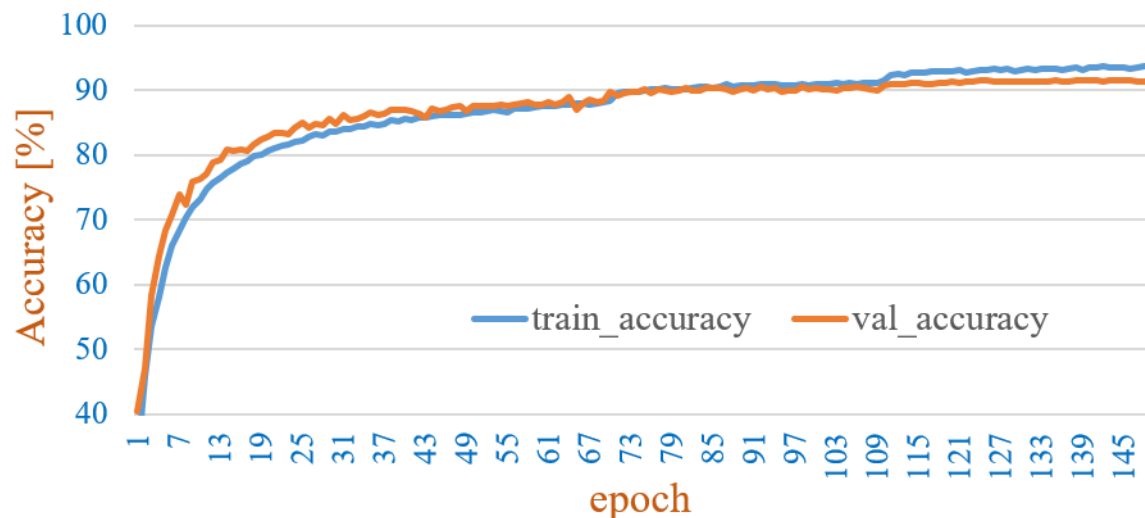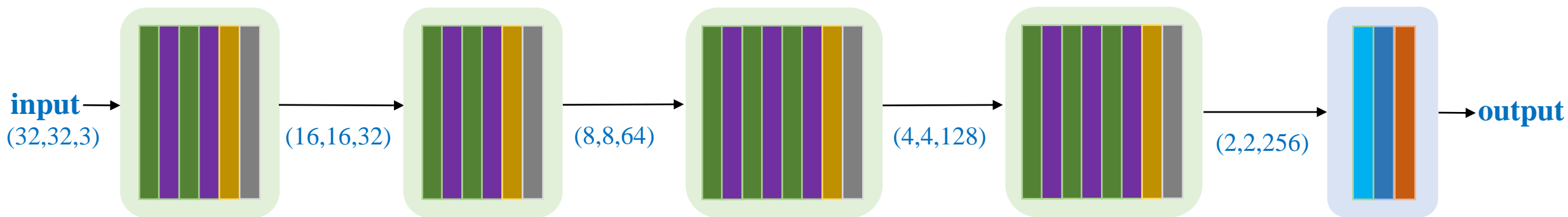


https://arxiv.org/pdf/1803.08494.pdf

"applying IN which does not only reduce the difference caused by domain changes, but also the illumination variation in single spectral images"

AFD-Net Aggregated Feature Difference Learning for Cross-Spectral Image Patch Matching (ICCV, 2019)
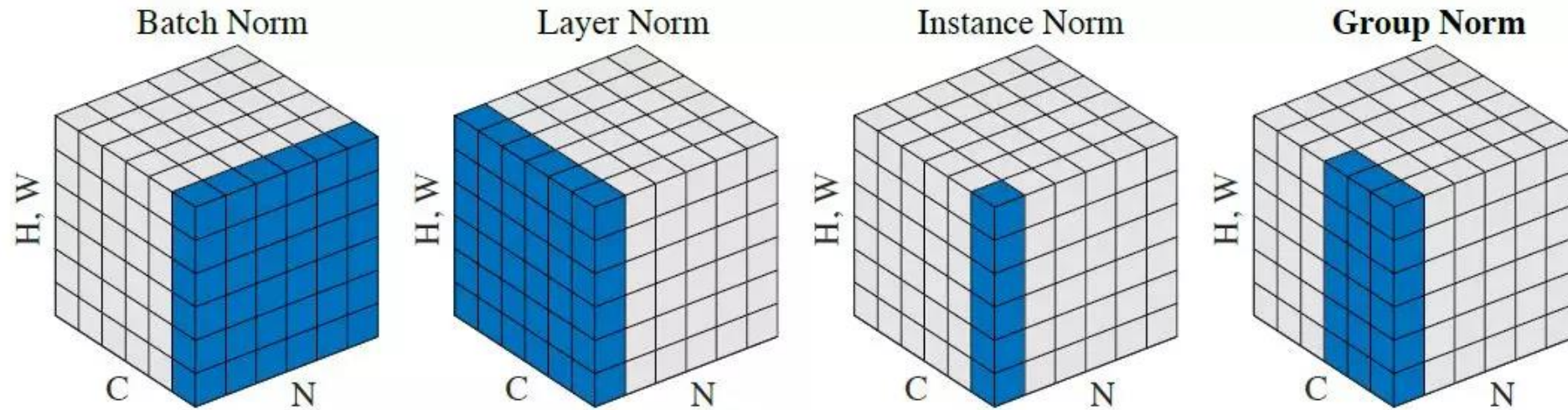
# Model Generalization

❖ **Trick 6: Instance normalization**



input
(32,32,3)

(16,16,32)

(8,8,64)

(4,4,128)

(2,2,256)

output

**Dropout**

**Instance Normalization +
Batch normalization**

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

**val_accuracy reaches to ~91.6%**

# Model Generalization

❖ **Trick 6: More about normalization**



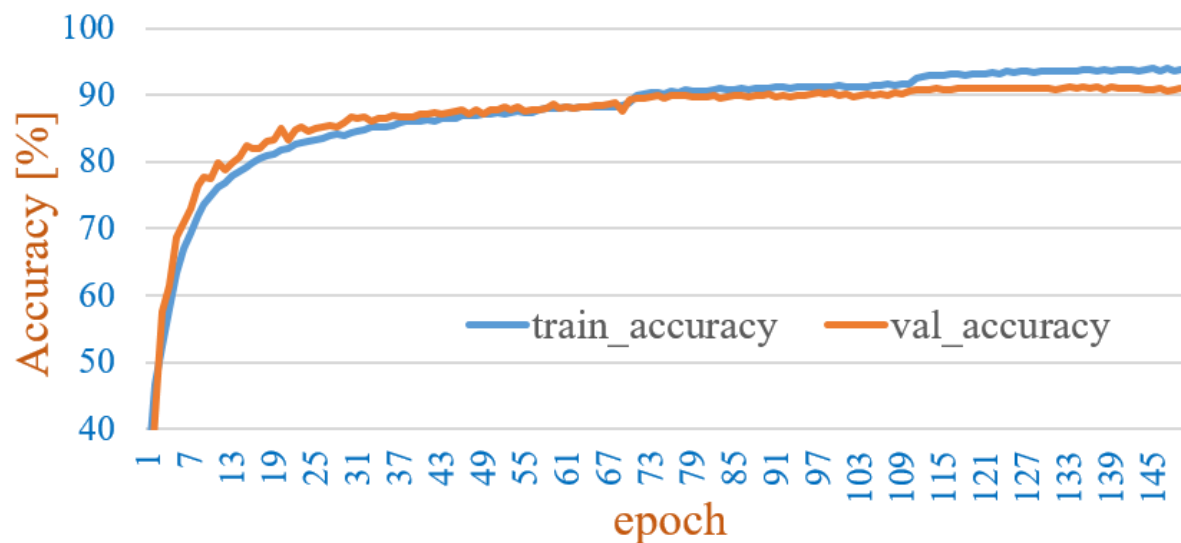https://arxiv.org/pdf/1803.08494.pdf

# Model Generalization

❖ **Summary**

**Horizontal flip + crop-and-resize**



val_accuracy reaches to ~91.6%

train_accuracy reaches to ~93.7%
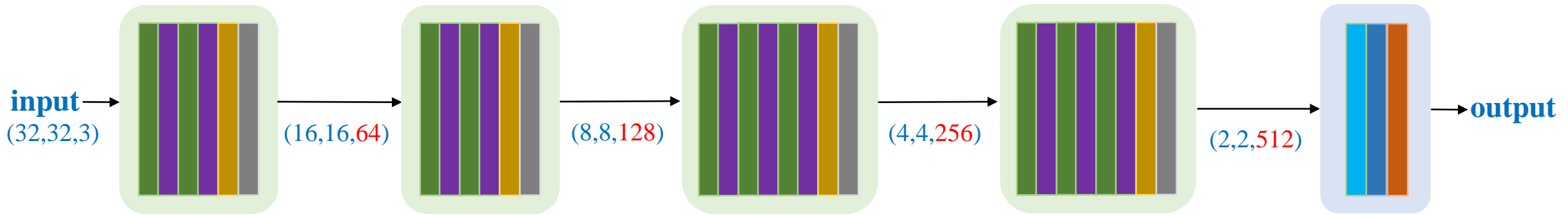
Batch normalization

Dropout

Kernel regularization

Data augmentation

**Idea: try to increase train_accuracy, expect val_accuracy increases too**

➡ **Increase model capacity**

# Model Generalization

❖ **Increase model capacity**



input
(32,32,3) → (16,16,64) → (8,8,128) → (4,4,256) → (2,2,512) → output

**val_accuracy reaches to ~93.6%**

**train_accuracy reaches to ~97.9%**