

# Convolutional Neural Networks

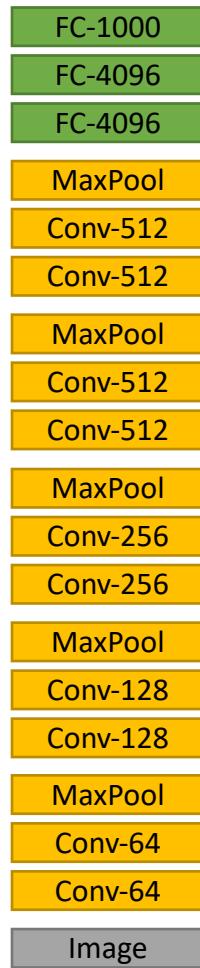
Cunjian Chen

Department of Computer Science and Engineering  
Michigan State University

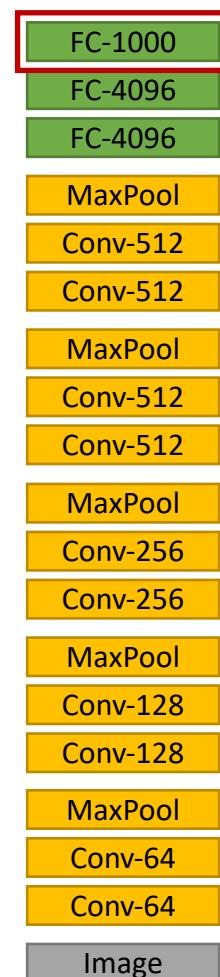
1/13/2020

# How to Train CNNs?

## 1. Train on ImageNet



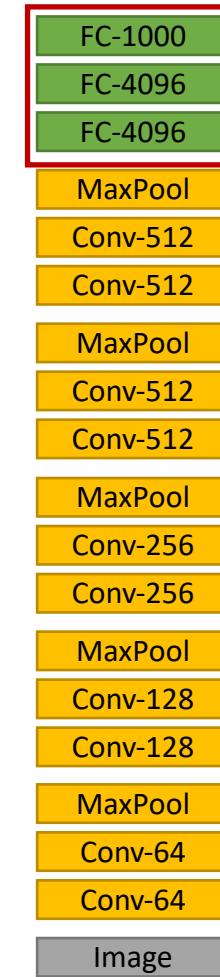
## 2. Small Dataset (C classes)



Reinitialize  
This and train

Freeze these

## 3. Bigger dataset



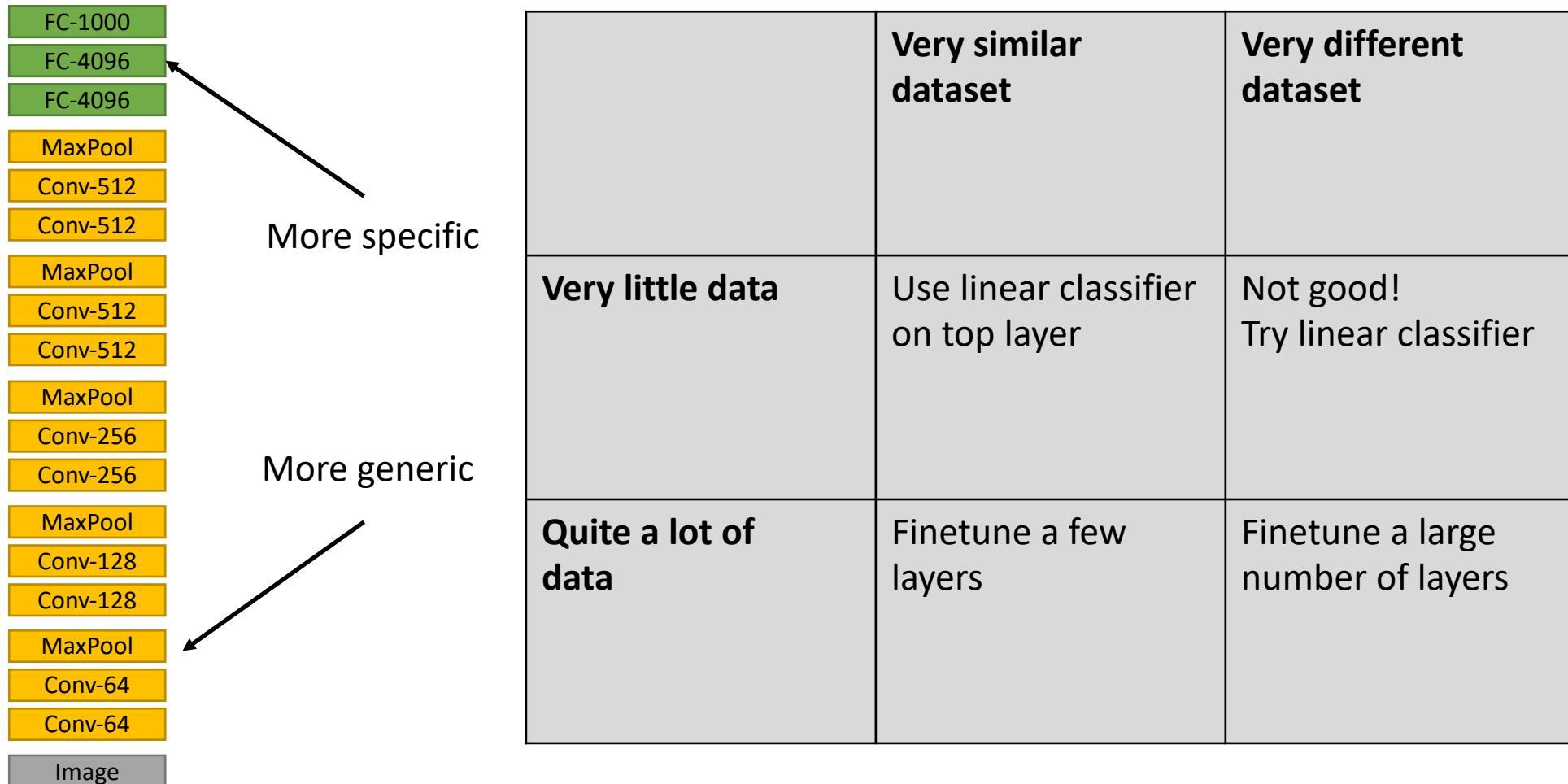
Train these

With bigger dataset,  
train more layers

Freeze these

Lower learning rate  
when finetuning;  
1/10 of original LR is  
good starting point

# How to Train CNNs?



# How to Train CNNs?

```
import torchvision.models as models
resnet18 = models.resnet18(pretrained=True)
alexnet = models.alexnet(pretrained=True)
squeezenet = models.squeeze1_0(pretrained=True)
vgg16 = models.vgg16(pretrained=True)
densenet = models.densenet161(pretrained=True)
inception = models.inception_v3(pretrained=True)
googlenet = models.googlenet(pretrained=True)
shufflenet = models.shufflenet_v2_x1_0(pretrained=True)
mobilenet = models.mobilenet_v2(pretrained=True)
resnext50_32x4d = models.resnext50_32x4d(pretrained=True)
wide_resnet50_2 = models.wide_resnet50_2(pretrained=True)
mnasnet = models.mnasnet1_0(pretrained=True)
```

## TORCHVISION.MODELS

# How to Train CNNs?

Two major transfer learning scenarios:

- **Finetuning the convnet:** Instead of random initialization, we initialize the network with a pretrained network, like the one that is trained on imagenet 1000 dataset. Rest of the training looks as usual.
- **ConvNet as fixed feature extractor:** Here, we will freeze the weights for all of the network except that of the final fully connected layer. This last fully connected layer is replaced with a new one with random weights and only this layer is trained.

# How to Train CNNs?

```
model_ft = models.resnet18(pretrained=True) → Model definition
num_ftrs = model_ft.fc.in_features
# Here the size of each output sample is set to 2.
# Alternatively, it can be generalized to nn.Linear(num_ftrs, len(class_names)).
model_ft.fc = nn.Linear(num_ftrs, 2) → Model finetuned

model_ft = model_ft.to(device)

criterion = nn.CrossEntropyLoss() → Loss function

# Observe that all parameters are being optimized
optimizer_ft = optim.SGD(model_ft.parameters(), lr=0.001, momentum=0.9) → Optimizer

# Decay LR by a factor of 0.1 every 7 epochs
exp_lr_scheduler = lr_scheduler.StepLR(optimizer_ft, step_size=7, gamma=0.1) → Learning rate policy

model_ft = train_model(model_ft, criterion, optimizer_ft, exp_lr_scheduler,
                      num_epochs=25) → Model Training
```

# How to Train CNNs?

```
model_conv = torchvision.models.resnet18(pretrained=True)
for param in model_conv.parameters():
    param.requires_grad = False freeze the parameters

# Parameters of newly constructed modules have requires_grad=True by default
num_ftrs = model_conv.fc.in_features
model_conv.fc = nn.Linear(num_ftrs, 2)

model_conv = model_conv.to(device)

criterion = nn.CrossEntropyLoss()

# Observe that only parameters of final layer are being optimized as
# opposed to before.
optimizer_conv = optim.SGD(model_conv.fc.parameters(), lr=0.001, momentum=0.9)

# Decay LR by a factor of 0.1 every 7 epochs
exp_lr_scheduler = lr_scheduler.StepLR(optimizer_conv, step_size=7, gamma=0.1)
```

# CNN Training Procedure

A typical training procedure for a neural network is as follows:

- Define the neural network that has some learnable parameters (or weights)
- Iterate over a dataset of inputs
- Process input through the network
- Compute the loss (how far is the output from being correct)
- Propagate gradients back into the network's parameters
- Update the weights of the network, typically using a simple update rule:  $\text{weight} = \text{weight} - \text{learning\_rate} * \text{gradient}$

```
learning_rate = 0.01
for f in net.parameters():
    f.data.sub_(f.grad.data * learning_rate)
```

# CNN Training Procedure

---

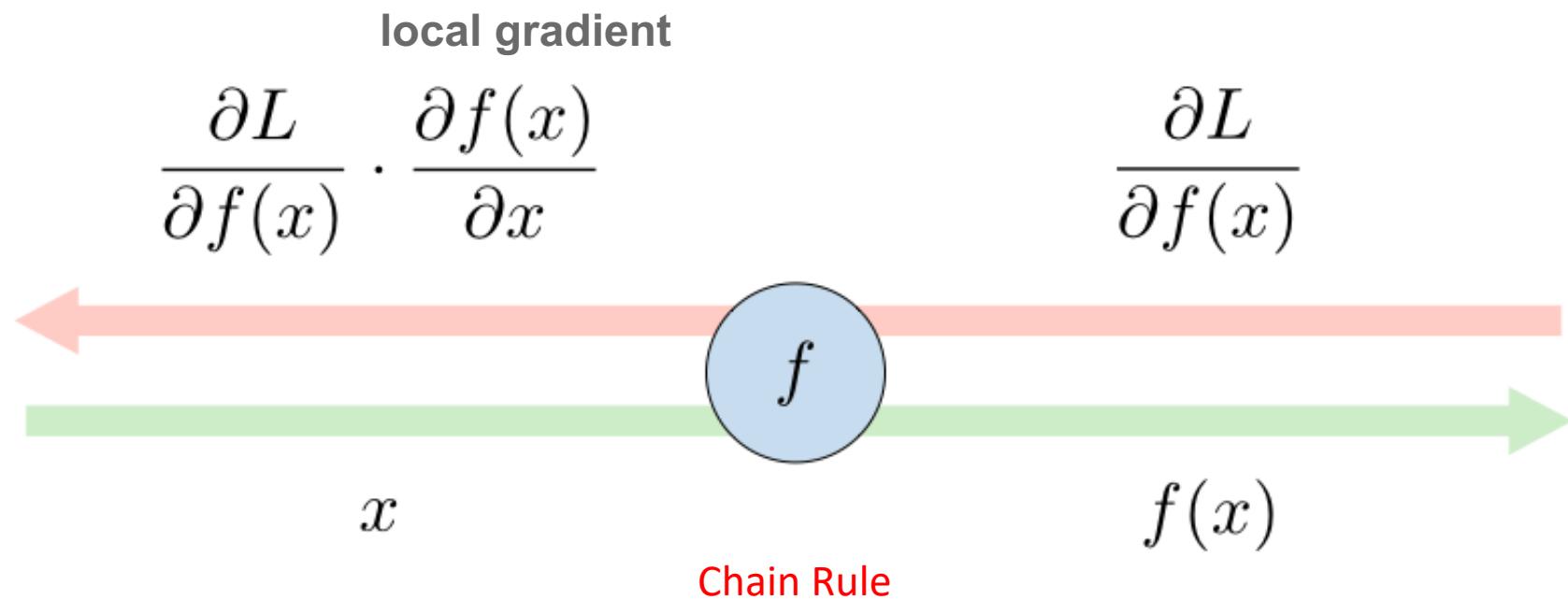
**Algorithm 1** Train a neural network with mini-batch stochastic gradient descent.

---

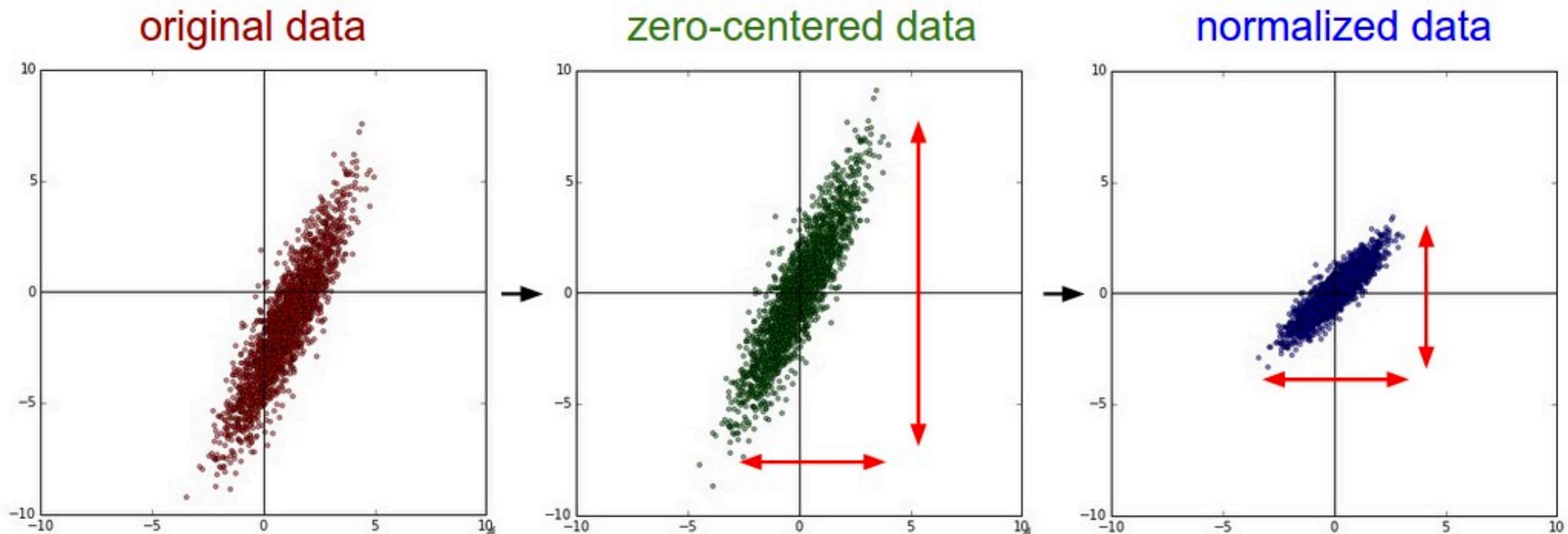
```
initialize(net)
for epoch = 1, . . . , K do
    for batch = 1, . . . , #images/b do
        images  $\leftarrow$  uniformly random sample b images
         $X, y \leftarrow$  preprocess(images)
         $z \leftarrow$  forward(net,  $X$ )
         $\ell \leftarrow$  loss( $z, y$ )
        grad  $\leftarrow$  backward( $\ell$ )
        update(net, grad)
    end for
end for
```

---

# Backpropagation



# Data Preprocessing



[0,255]: Normalize RGB channels by subtracting 123.68, 116.779, 103.939 and dividing by 58.393, 57.12, 57.375, respectively

[0,1]: Normalize RGB channels by subtracting 0.485, 0.456, 0.406 and dividing by 0.229, 0.224, 0.225, respectively

# CNN Preprocessing Pipeline

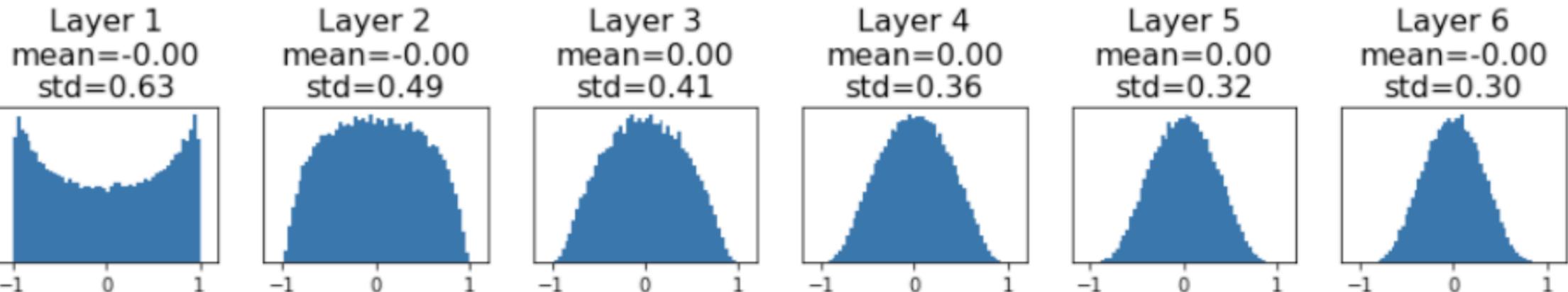
1. Randomly sample an image and decode it into 32-bit floating point raw pixel values in [0, 255].
2. Randomly crop a rectangular region whose aspect ratio is randomly sampled in [3/4, 4/3] and area randomly sampled in [8%, 100%], then resize the cropped region into a 224-by-224 square image (**RandomResizedCrop**).
3. Flip horizontally with 0.5 probability (**RandomHorizontalFlip**).
4. Scale hue, saturation, and brightness with coefficients uniformly drawn from [0.6, 1.4].
5. Add PCA noise with a coefficient sampled from a normal distribution  $N(0, 0.1)$ .
6. Normalize RGB channels by subtracting 123.68, 116.779, 103.939 and dividing by 58.393, 57.12, 57.375, respectively (**Normalize**).

# Weight Initialization

```
dims = [4096] * 7          "Xavier" initialization:  
hs = []                      std = 1/sqrt(Din)  
x = np.random.randn(16, dims[0])  
for Din, Dout in zip(dims[:-1], dims[1:]):  
    W = np.random.randn(Din, Dout) / np.sqrt(Din)  
    x = np.tanh(x.dot(W))  
    hs.append(x)
```

"Just right": Activations are nicely scaled for all layers!

For conv layers, Din is  $\text{kernel\_size}^2 * \text{input\_channels}$



# Weight Initialization

Weight Initialization: Xavier Initialization

“Xavier” initialization:  
std = 1/sqrt(Din)

**Derivation:** Variance of output = Variance of input

$$y = Wx$$

$$y_i = \sum_{j=1}^{Din} x_j w_j$$

$$\begin{aligned} \text{Var}(y_i) &= \text{Din} * \text{Var}(x_i w_i) && [\text{Assume } x, w \text{ are iid}] \\ &= \text{Din} * (\text{E}[x_i^2] \text{E}[w_i^2] - \text{E}[x_i]^2 \text{E}[w_i]^2) && [\text{Assume } x, w \text{ independent}] \\ &= \text{Din} * \text{Var}(x_i) * \text{Var}(w_i) && [\text{Assume } x, w \text{ are zero-mean}] \end{aligned}$$

If  $\text{Var}(w_i) = 1/\text{Din}$  then  $\text{Var}(y_i) = \text{Var}(x_i)$

# Weight Update (Optimizers)

## SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

```
while True:  
    dx = compute_gradient(x)  
    x -= learning_rate * dx
```

## SGD+Momentum

$$\begin{aligned} v_{t+1} &= \rho v_t + \nabla f(x_t) \\ x_{t+1} &= x_t - \alpha v_{t+1} \end{aligned}$$

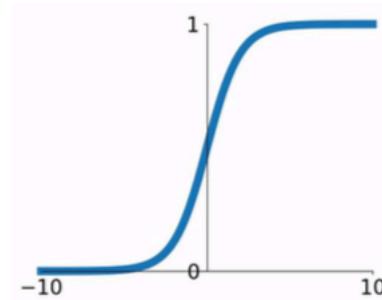
```
vx = 0  
while True:  
    dx = compute_gradient(x)  
    vx = rho * vx + dx  
    x -= learning_rate * vx
```

RMSprop, Adagrad, Adadelta, Adam, Adamax, Nadam...

# Activation Functions

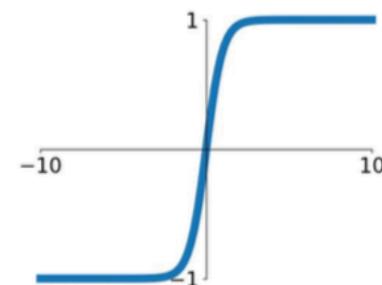
**Sigmoid**

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



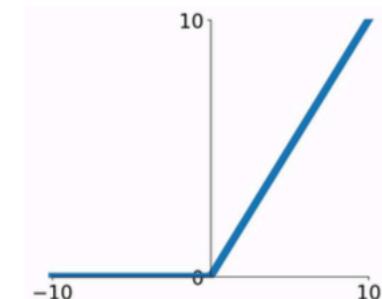
**tanh**

$$\tanh(x)$$

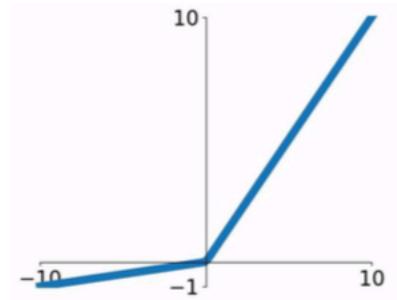


**ReLU**

$$\max(0, x)$$



**Leaky ReLU**  
 $\max(0.1x, x)$

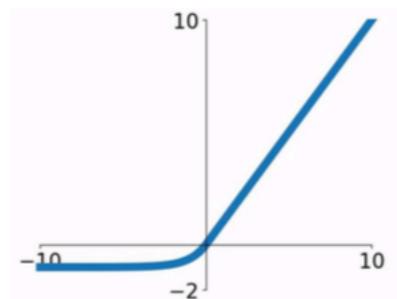


**Maxout**

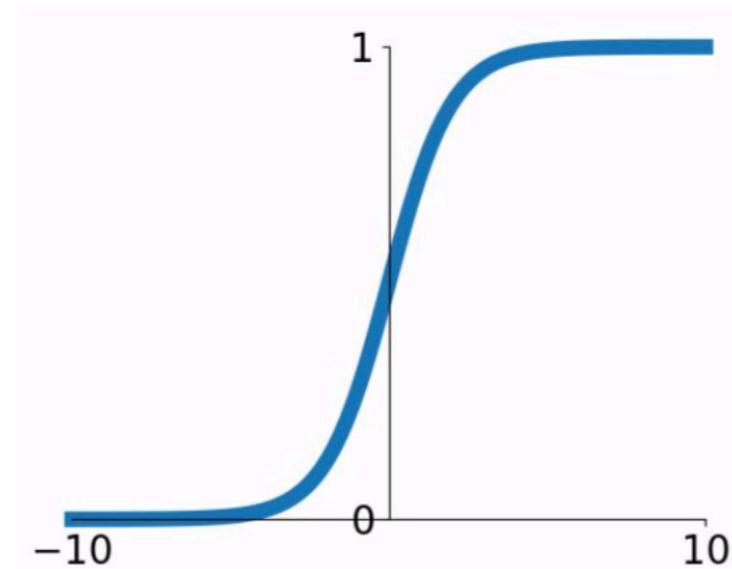
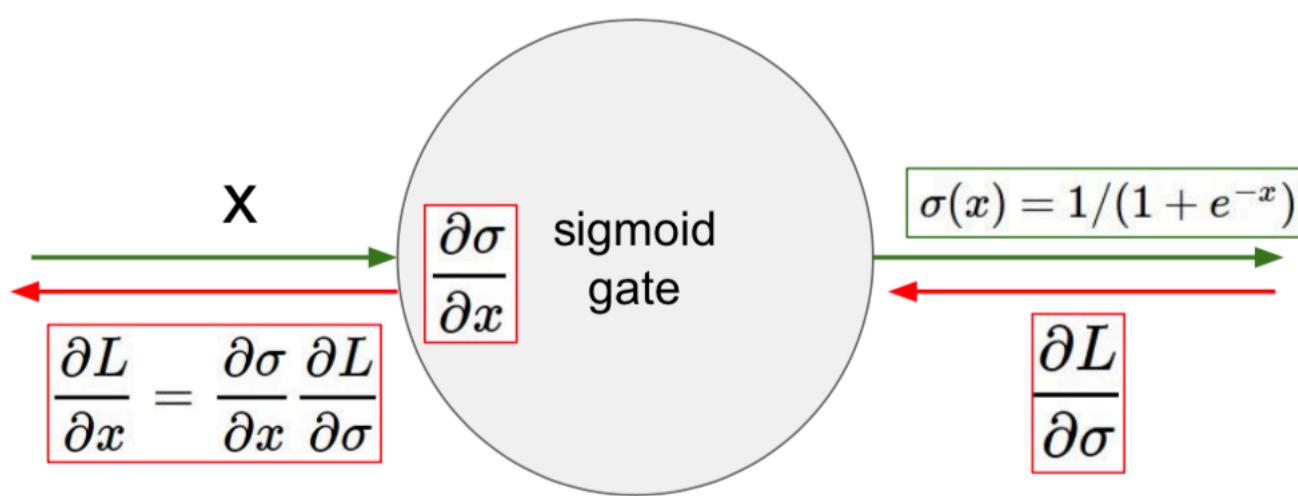
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



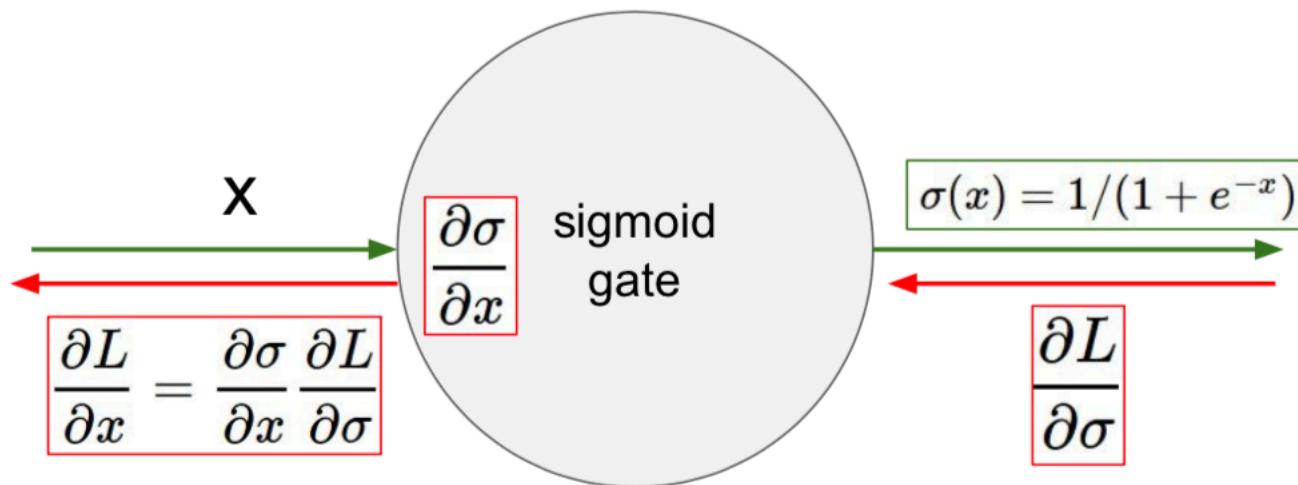
# Activation Functions



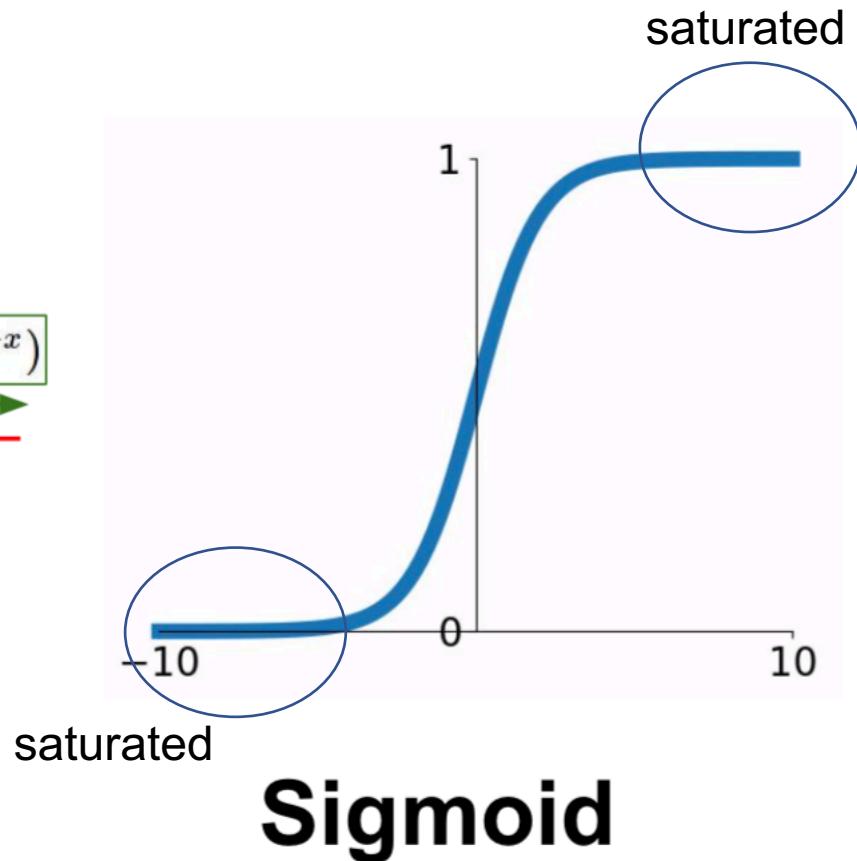
What is the gradient at  $x=-10, 0, 10$ ?

**Sigmoid**

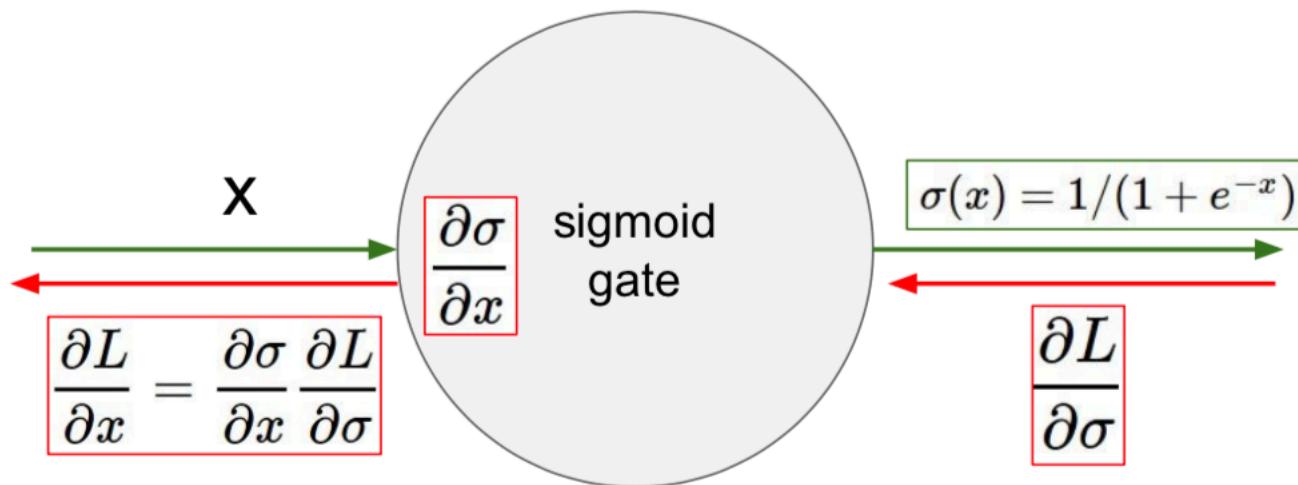
# Activation Functions



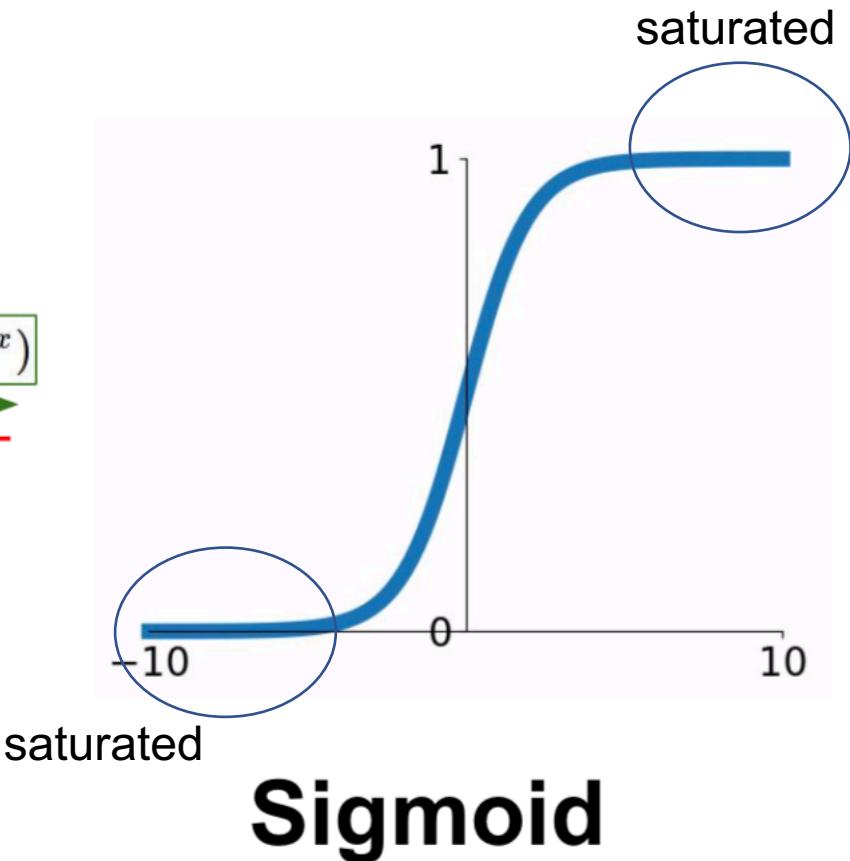
What is the gradient at  $x=-10, 0, 10$ ?



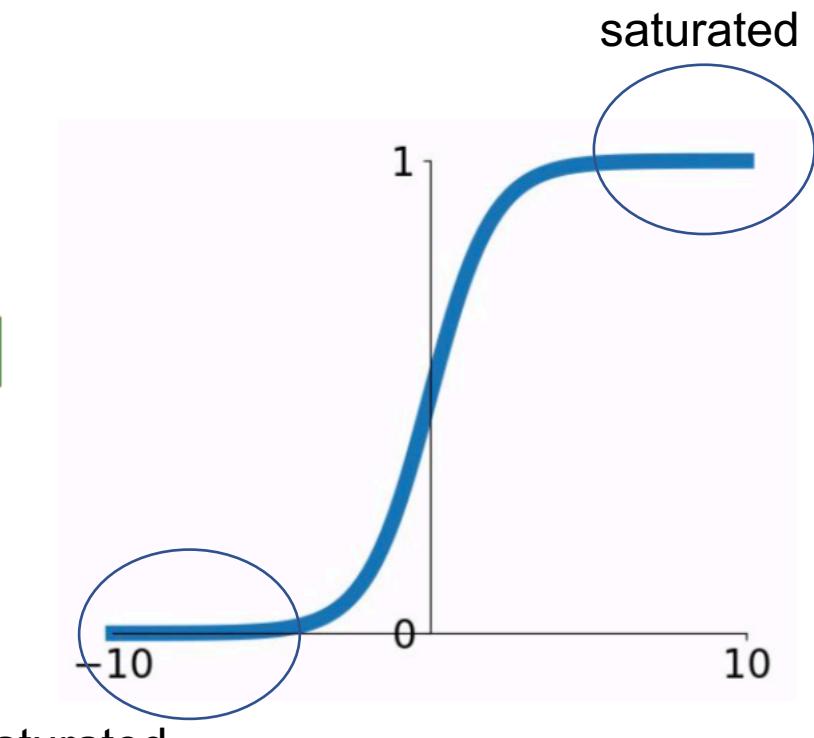
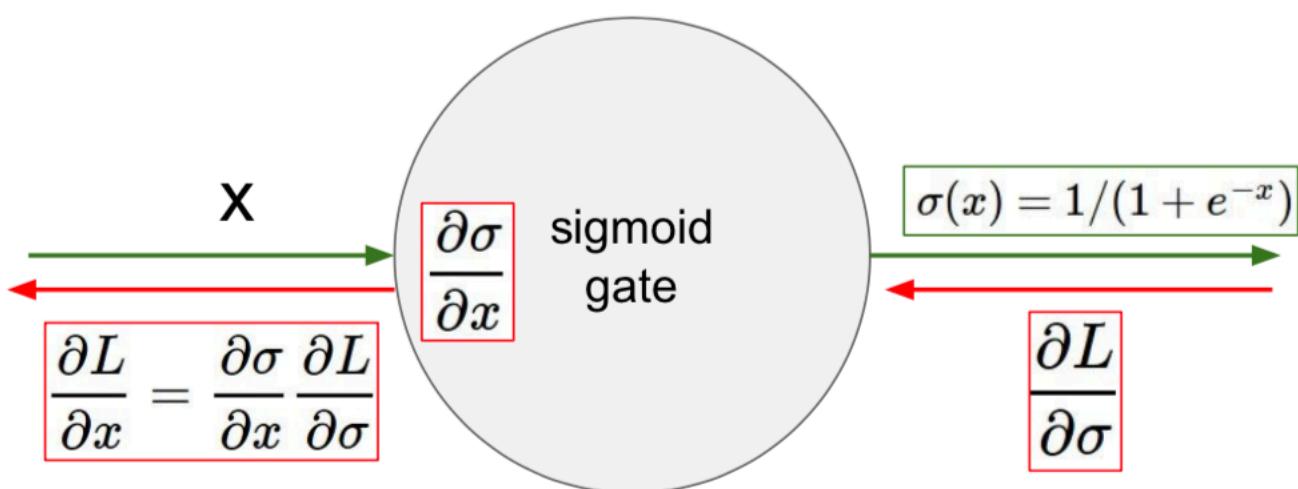
# Activation Functions



What is the other drawback?



# Activation Functions



## Sigmoid

*Sigmoid outputs are not zero-centered.*

The gradient on the weights  $w$  during backpropagation become either all positive, or all negative (depending on the gradient of the whole expression  $f$ ), given  $f = w^T x + b$  when  $x > 0$ .

(For a single element! Minibatches help)

# Activation Functions

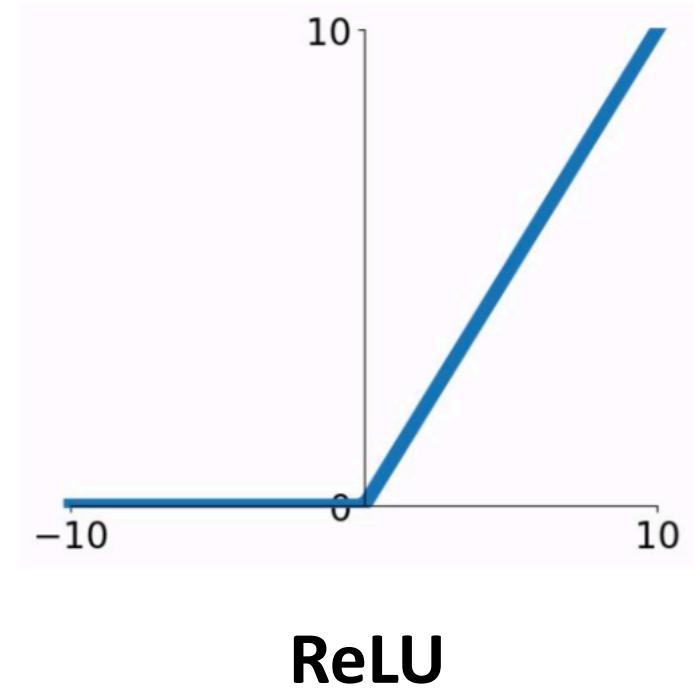
$$f(x) = \max(0, x)$$

Pros:

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

Cons:

- Not zero-centered output
- Could result in dead neurons.



# Activation Functions

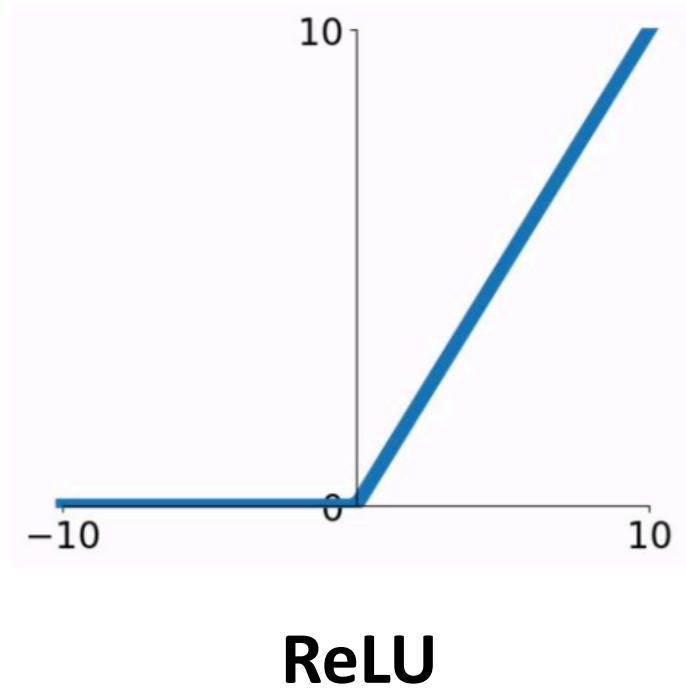
$$f(x) = \max(0, x)$$

Pros:

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

Cons:

- Not zero-centered output
- Could result in dead neurons.



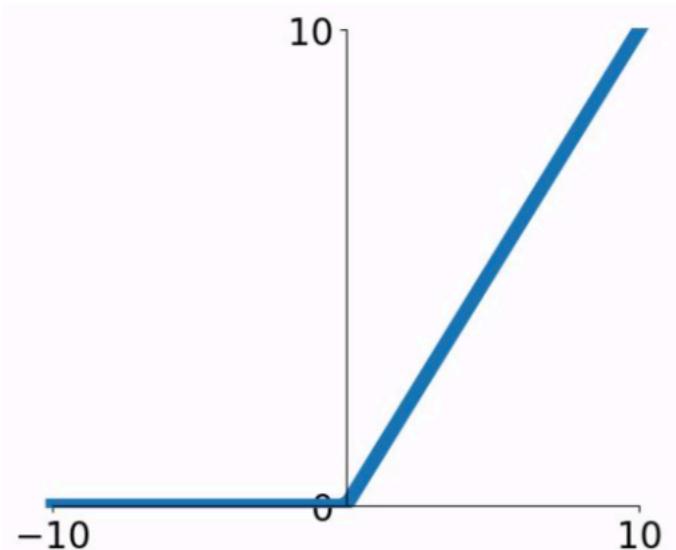
What is the gradient at x=-10, 0, 10?

# Activation Functions

$$f(x) = \begin{cases} x, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

$$f'(x) = \begin{cases} 1, & x > 0 \\ 0, & x < 0 \end{cases}$$

$$f(x) = \max(0, x)$$



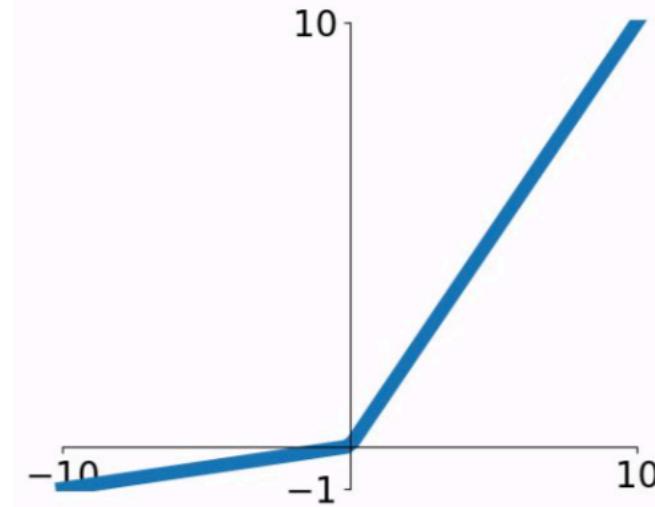
Dying ReLU problem: For activations in the region ( $x < 0$ ) of ReLU, gradient will be 0.

**ReLU**

# Activation Functions

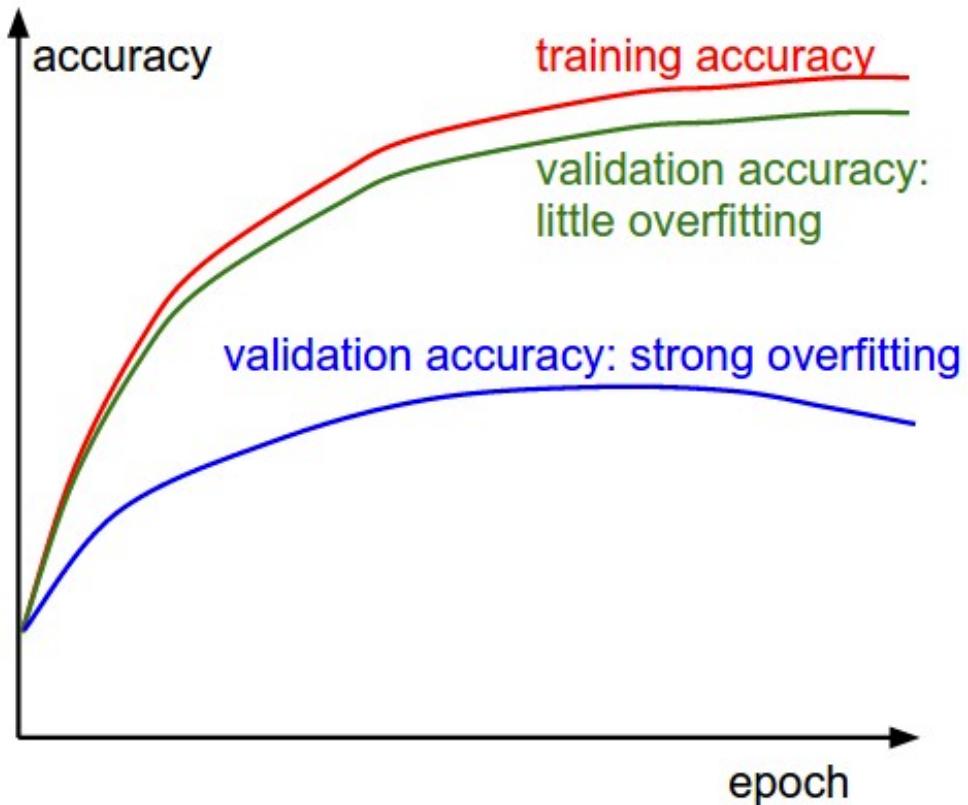
Pros:

- Does not saturate
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice  
(e.g. 6x)
- Will not die



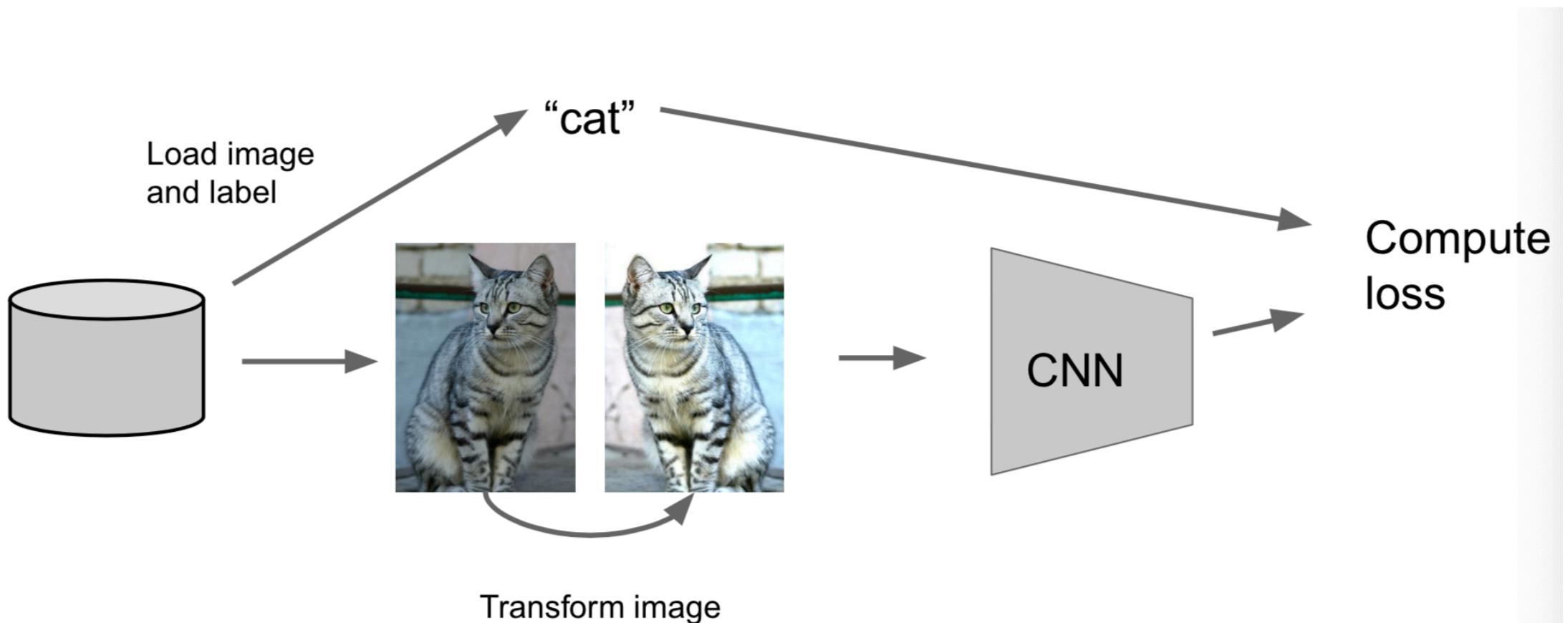
**Leaky ReLU**

# When Overfitting Happens?



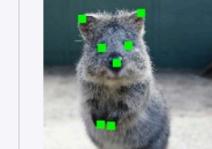
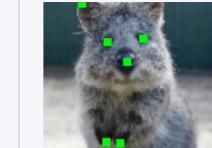
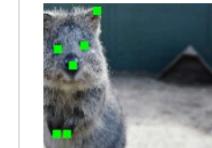
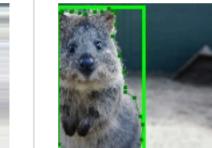
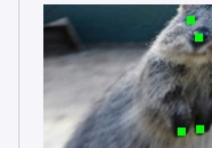
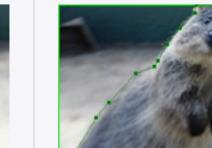
The gap between the training and validation accuracy indicates the amount of overfitting.

# Data Augmentation



Data augmentation can be performed either in real-time or off-line.

# Data Augmentation

	Image	Heatmaps	Seg. Maps	Keypoints	Bounding Boxes, Polygons
<i>Original Input</i>					
Gauss. Noise + Contrast + Sharpen					
Affine					
Crop + Pad					
Flplr + Perspective					

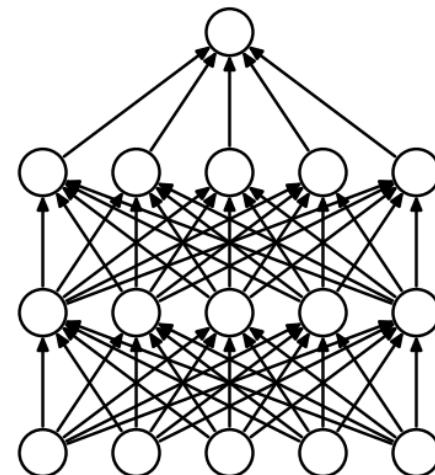
Less prone to overfitting

# Data Augmentation Implementation

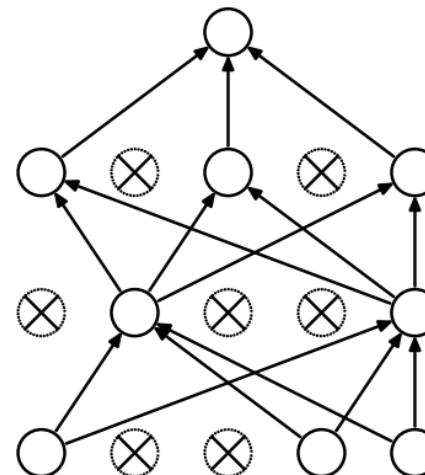
```
# Data augmentation and normalization for training
# Just normalization for validation
data_transforms = {
    'train': transforms.Compose([
        transforms.RandomResizedCrop(224),
        transforms.RandomHorizontalFlip(),
        transforms.ToTensor(),
        transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])
    ]),
    'val': transforms.Compose([
        transforms.Resize(256),
        transforms.CenterCrop(224),
        transforms.ToTensor(),
        transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])
    ]),
}
```

# Regularizations

- L2 regularization (weight decay):  $\frac{1}{2} \lambda w^2$
- L1 regularization:  $\lambda |w|$
- L1 + L2:  $\lambda_1 |w| + \lambda_2 w^2$
- Max norm:  $\|w\|_2 < c$
- Dropout
- Drop connect



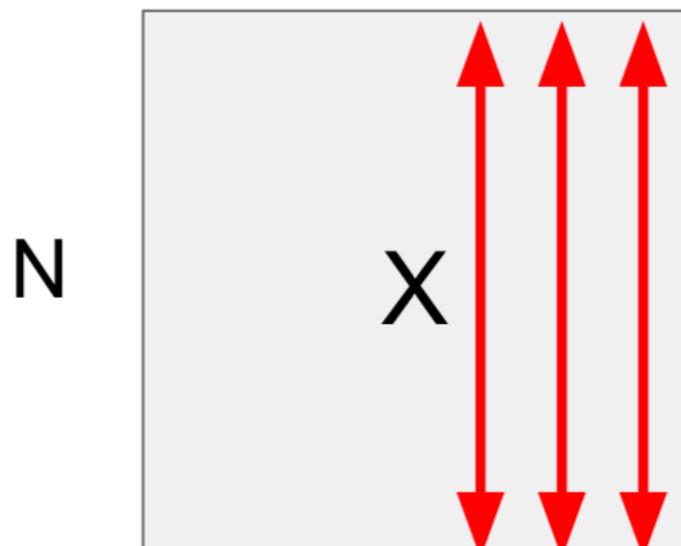
(a) Standard Neural Net



(b) After applying dropout.

# Batch Normalizations

**Input:**  $x : N \times D$



zero-mean unit-variance activations

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,  
shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,  
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Normalized x,  
Shape is  $N \times D$

D

Batch Normalization for fully-connected networks

# Batch Normalizations

**Input:**  $x : N \times D$

**Learnable scale and shift parameters:**

$\gamma, \beta : D$

Learning  $\gamma = \sigma$ ,  
 $\beta = \mu$  will recover the identity function!

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,  
shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,  
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

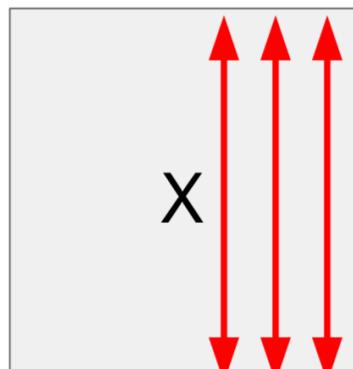
Normalized x,  
Shape is N x D

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Output,  
Shape is N x D

# Batch Normalizations

Input:  $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean,  
shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel var,  
shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Normalized x,  
Shape is  $N \times D$

```
X = torch.randn(20,100) * 5 + 10
X = Variable(X)

mu = torch.mean(X[:,1])
var_ = torch.var(X[:,1], unbiased=False)
sigma = torch.sqrt(var_ + 1e-5)
x = (X[:,1] - mu)/sigma
```

# Batch Normalizations

Batch Normalization for  
**fully-connected** networks

$\mathbf{x}$ :  $N \times D$

Normalize



$\mu, \sigma$ :  $1 \times D$

$\gamma, \beta$ :  $1 \times D$

$$y = \gamma(\mathbf{x} - \mu) / \sigma + \beta$$

Batch Normalization for  
**convolutional** networks  
(Spatial Batchnorm, BatchNorm2D)

$\mathbf{x}$ :  $N \times C \times H \times W$

Normalize

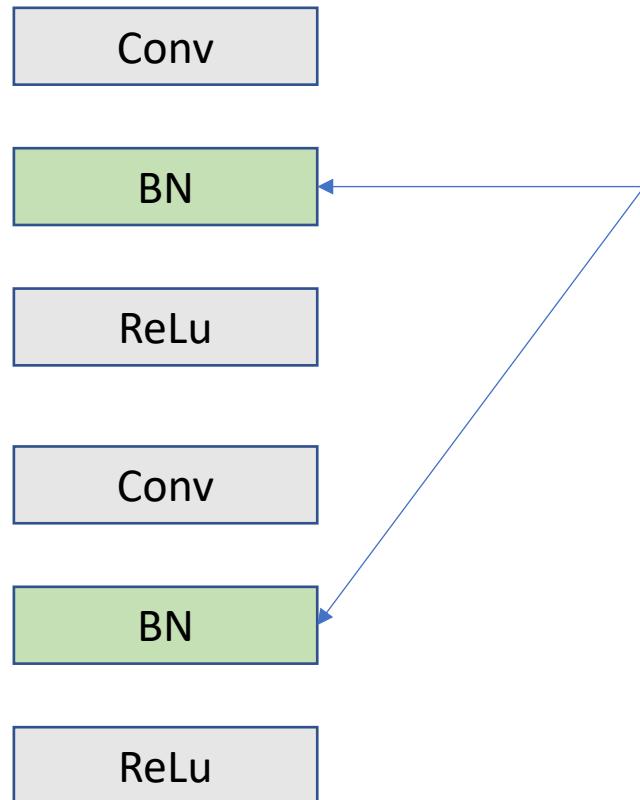


$\mu, \sigma$ :  $1 \times C \times 1 \times 1$

$\gamma, \beta$ :  $1 \times C \times 1 \times 1$

$$y = \gamma(\mathbf{x} - \mu) / \sigma + \beta$$

# Batch Normalizations



Usually insert after FC or Conv layers  
before the nonlinearity

```
out = self.conv1(x)  
out = self.bn1(out)  
out = self.relu(out)
```

```
out = self.conv2(out)  
out = self.bn2(out)  
out = self.relu(out)
```

```
out = self.conv3(out)  
out = self.bn3(out)
```

Resnet bottleneck

# Batch Normalizations

- Training: the sample mean and (uncorrected) sample variance are computed from minibatch statistics and used to normalize the incoming data. During training we also keep an exponentially decaying running mean of the mean and variance of each feature, and these averages are used to normalize data at test-time.

```
running_mean = momentum * running_mean + (1 - momentum) * sample_mean
```

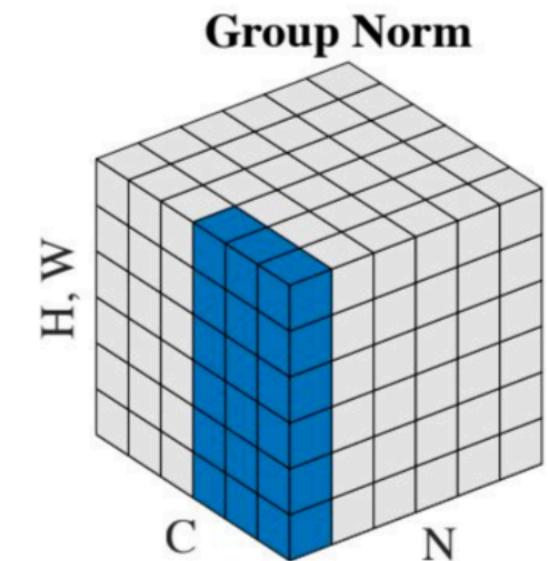
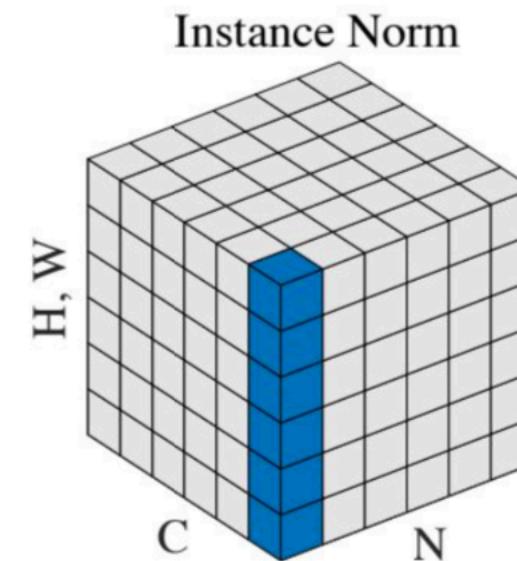
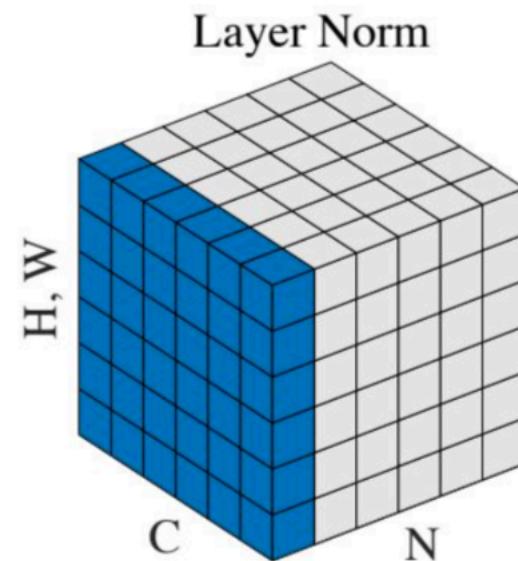
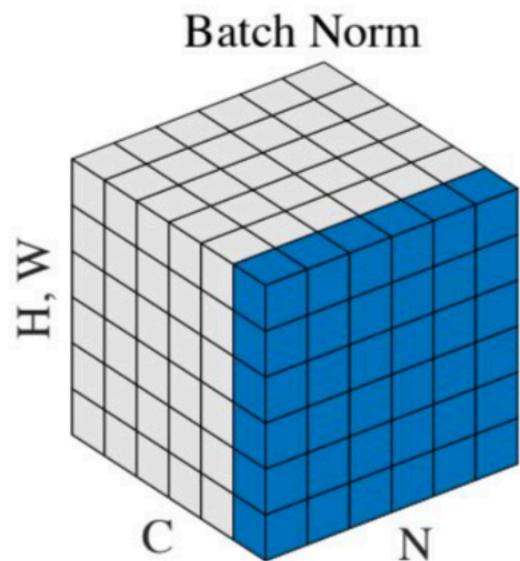
```
running_var = momentum * running_var + (1 - momentum) * sample_var
```

- Testing:

```
x_stand = (x - running_mean) / np.sqrt(running_var)
```

```
out = x_stand * gamma + beta
```

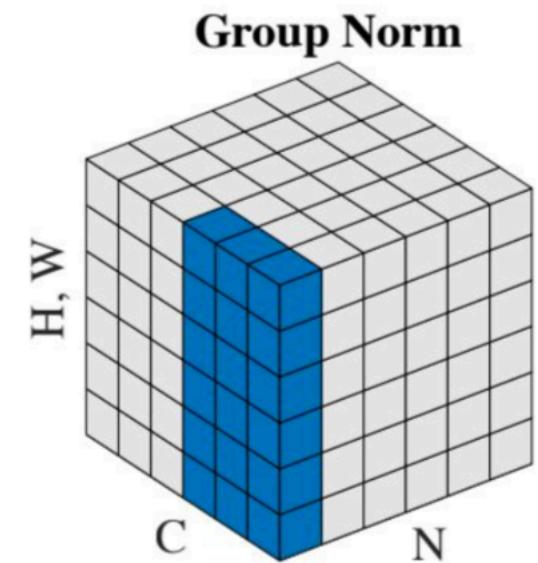
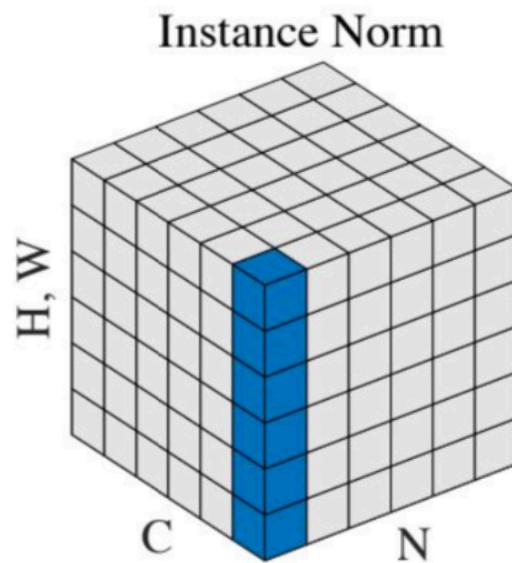
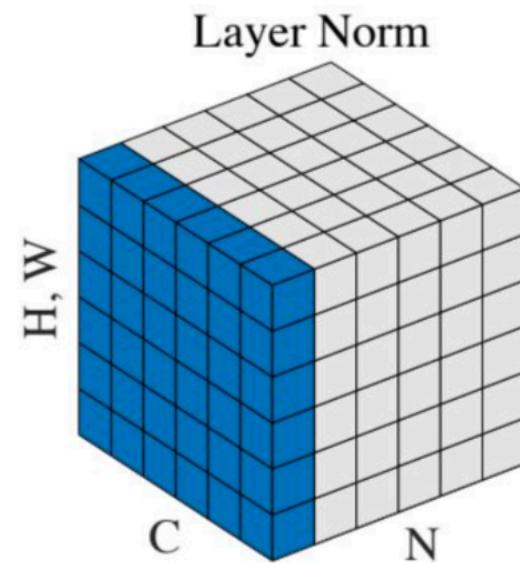
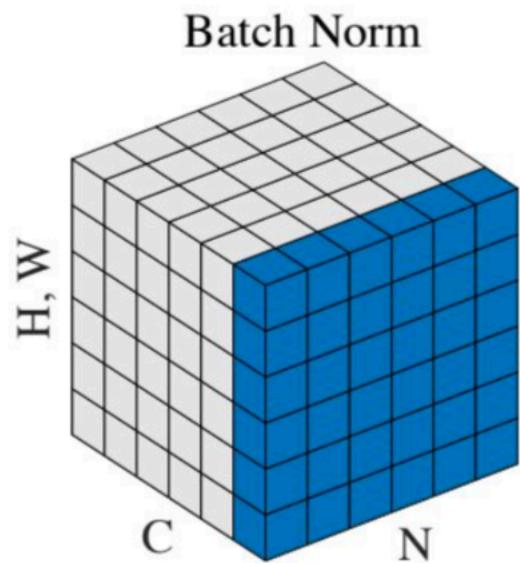
# Normalizations



Feature map tensor:  $N \times C \times H \times W$

What is shape of mean and sigma?

# Normalizations



$1 \times C \times 1 \times 1$

$N \times 1 \times 1 \times 1$

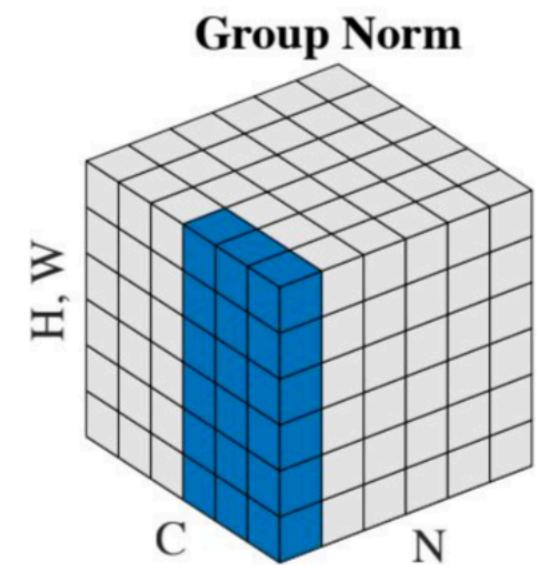
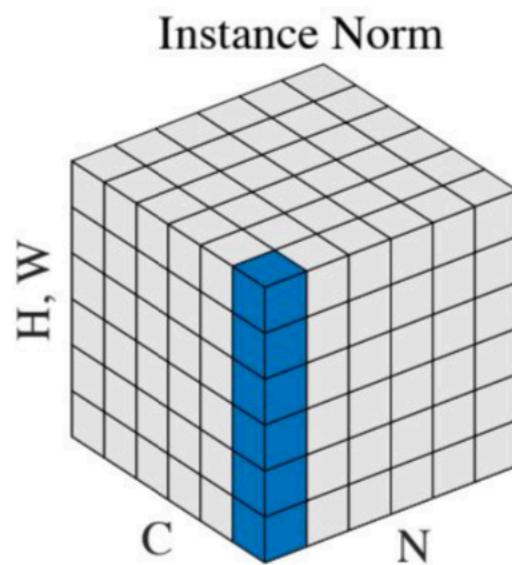
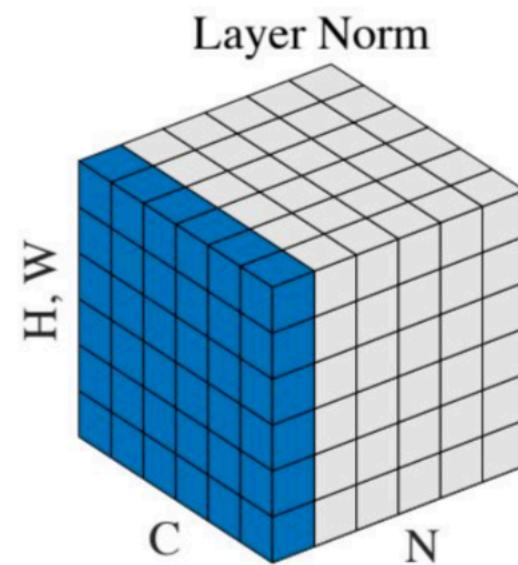
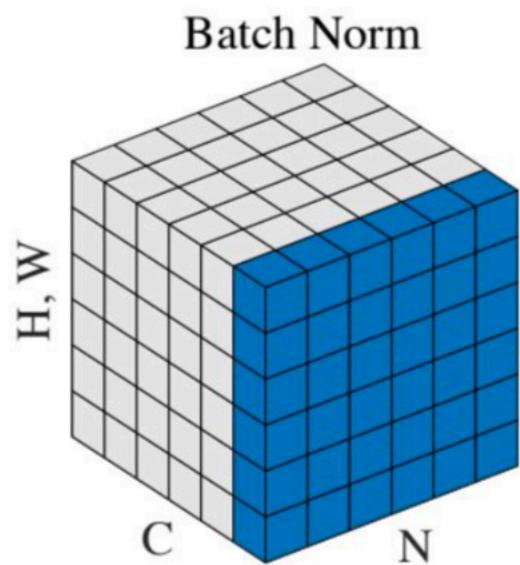
$N \times C \times 1 \times 1$

$N \times G \times 1 \times 1 \times 1$



$[N, G, C // G, H, W]$

# Normalizations



$1 \times C \times 1 \times 1$

$N \times 1 \times 1 \times 1$

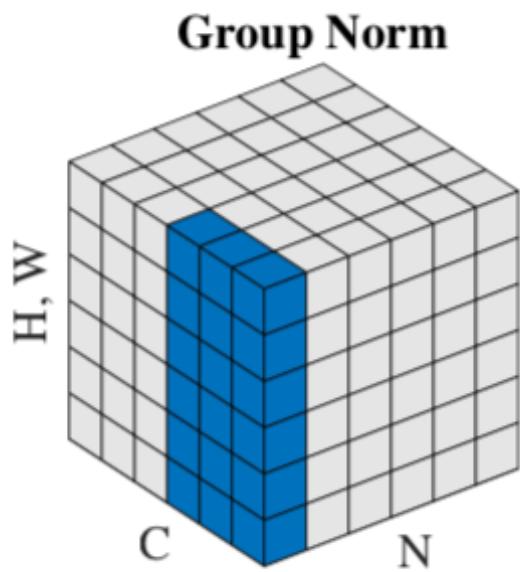
$N \times C \times 1 \times 1$

$N \times G \times 1 \times 1 \times 1$



$[N, G, C // G, H, W]$

# Normalizations



```
def GroupNorm(x, gamma, beta, G, eps=1e-5):
    # x: input features with shape [N,C,H,W]
    # gamma, beta: learnable scale and offset, with shape [1,C,1,1]
    # G: number of groups for GN

    N, C, H, W = x.shape
    x = tf.reshape(x, [N, G, C // G, H, W])

    mean, var = tf.nn.moments(x, [2, 3, 4], keep_dims=True)
    x = (x - mean) / tf.sqrt(var + eps)

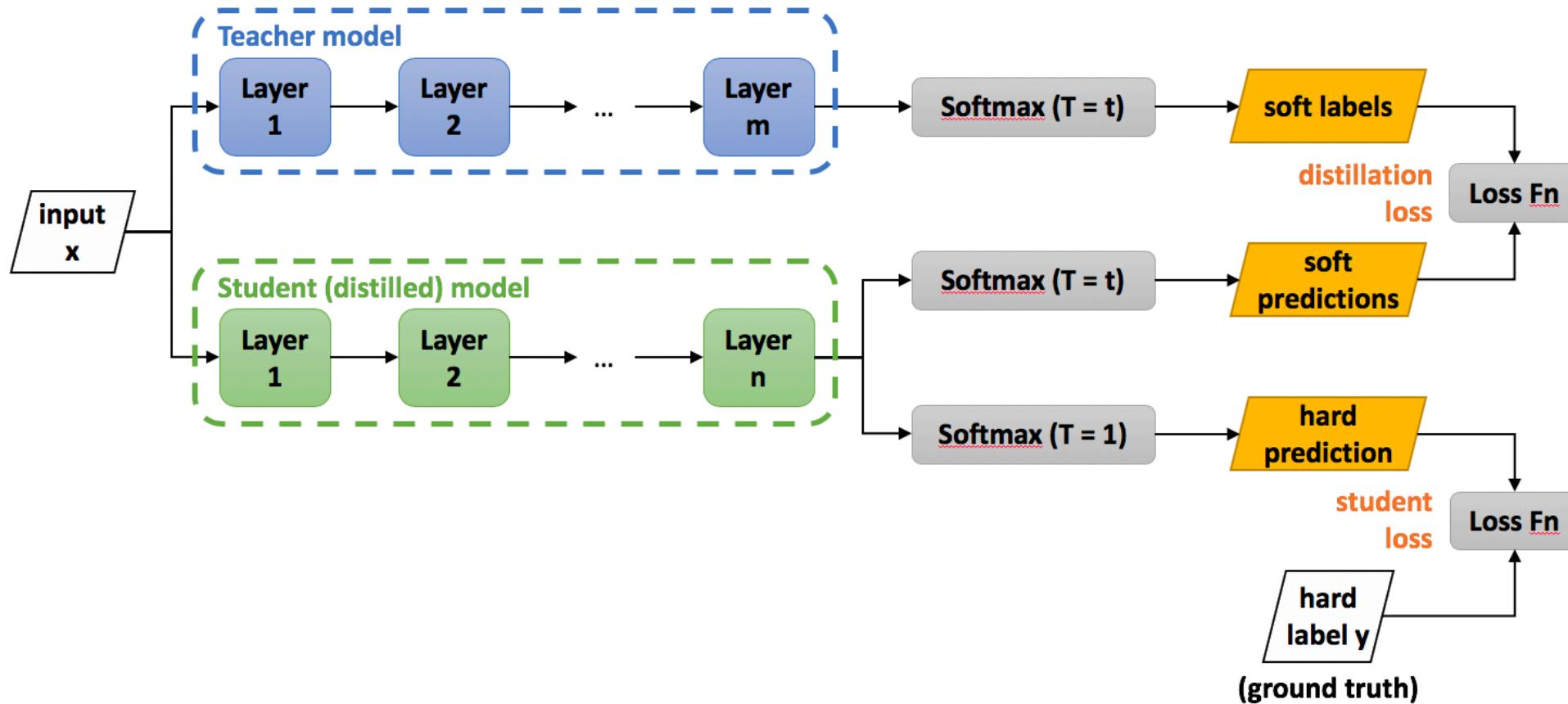
    x = tf.reshape(x, [N, C, H, W])

    return x * gamma + beta
```

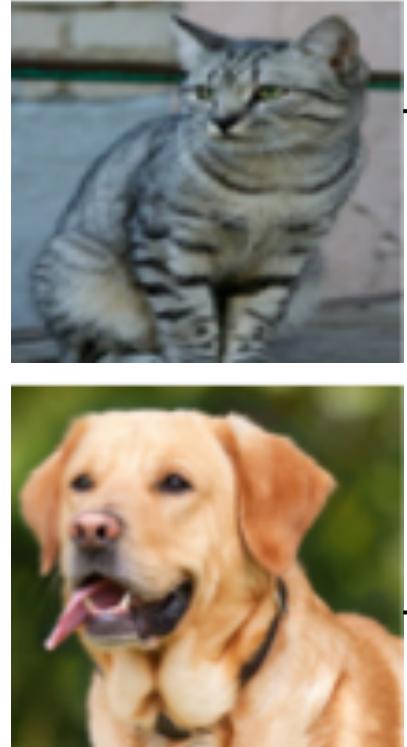
# Bag of Tricks for CNN Training

- Learning Rate Warmup + Cosine Learning Rate Decay
- Label Smoothing
- Knowledge Distillation
- Mixup Training

# Knowledge Distillation



# Mixup Training (Data Augmentation)



```
def mixup_data(x, y, alpha=1.0, use_cuda=True):  
  
    '''Compute the mixup data. Return mixed inputs, pairs of targets, and lambda'''  
    if alpha > 0.:  
        lam = np.random.beta(alpha, alpha)  
    else:  
        lam = 1.  
    batch_size = x.size()[0]  
    if use_cuda:  
        index = torch.randperm(batch_size).cuda()  
    else:  
        index = torch.randperm(batch_size)  
  
    mixed_x = lam * x + (1 - lam) * x[index,:]  
    y_a, y_b = y, y[index]  
    return mixed_x, y_a, y_b, lam  
  
def mixup_criterion(y_a, y_b, lam):  
    return lambda criterion, pred: lam * criterion(pred, y_a) + (1 - lam) * criterion(pred, y_b)
```

**Training:** Train on random blends of images  
**Testing:** Use original images

$$\tilde{x} = \lambda x_i + (1 - \lambda)x_j, \quad \text{where } x_i, x_j \text{ are raw input vectors}$$
$$\tilde{y} = \lambda y_i + (1 - \lambda)y_j, \quad \text{where } y_i, y_j \text{ are one-hot label encodings}$$

<https://github.com/hongyi-zhang/mixup>