



Convolutional Neural Networks

A. M. Sadeghzadeh, Ph.D.

Sharif University of Technology
Computer Engineering Department (CE)
Data and Network Security Lab (DNSL)



March 12, 2023

Most slides have been adapted from Justin Johnson, EECS 498-007, University of Michigan 2020, and Fei Fei Li, cs231n, Stanford 2017

Today's agenda

1 Recap

2 Vanishing/Exploding Gradient

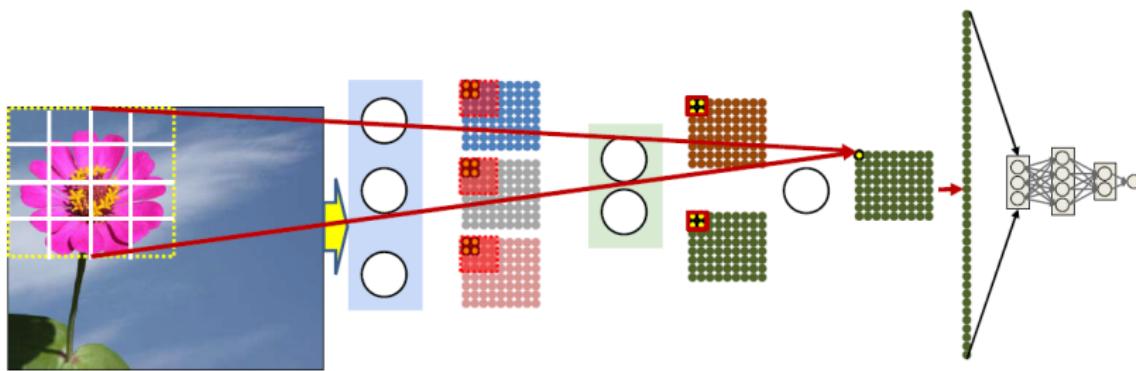
3 Batch Normalization

4 CNN Architectures

Recap

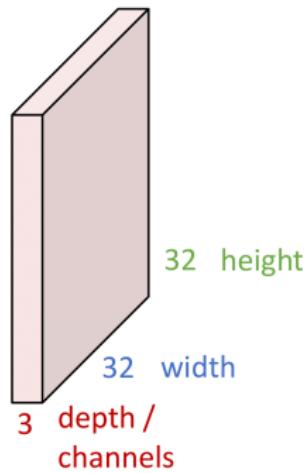
Convolutional Neural Networks (CNNs)

- Advantage of CNNs
 - Shift invariant
 - Parameter sharing
 - Sparse weights
 - Preserve spatial structures



Convolution Layer

3x32x32 image



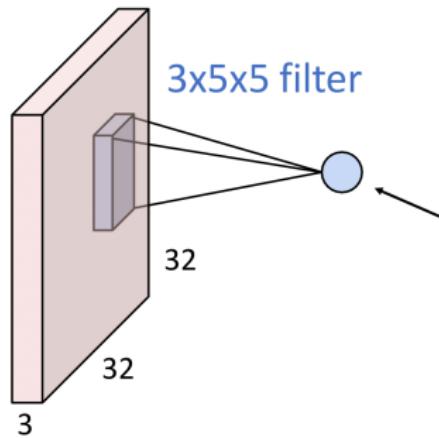
3x5x5 filter



Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Convolution Layer

3x32x32 image

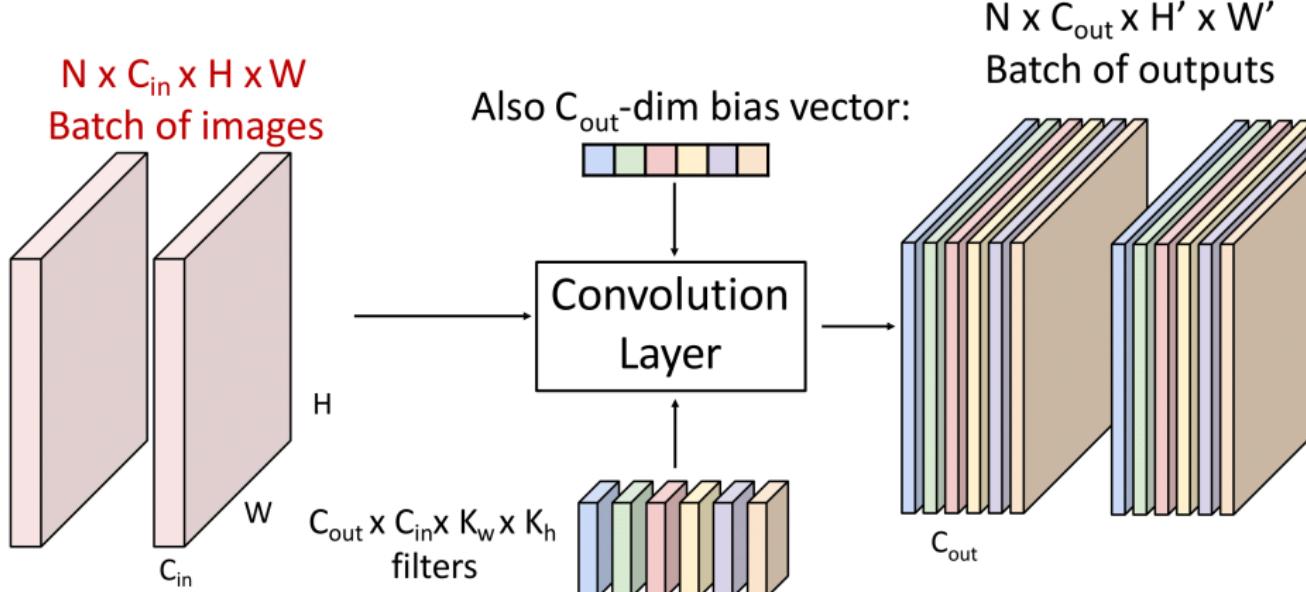


1 number:

the result of taking a dot product between the filter and a small 3x5x5 chunk of the image
(i.e. $3 \times 5 \times 5 = 75$ -dimensional dot product + bias)

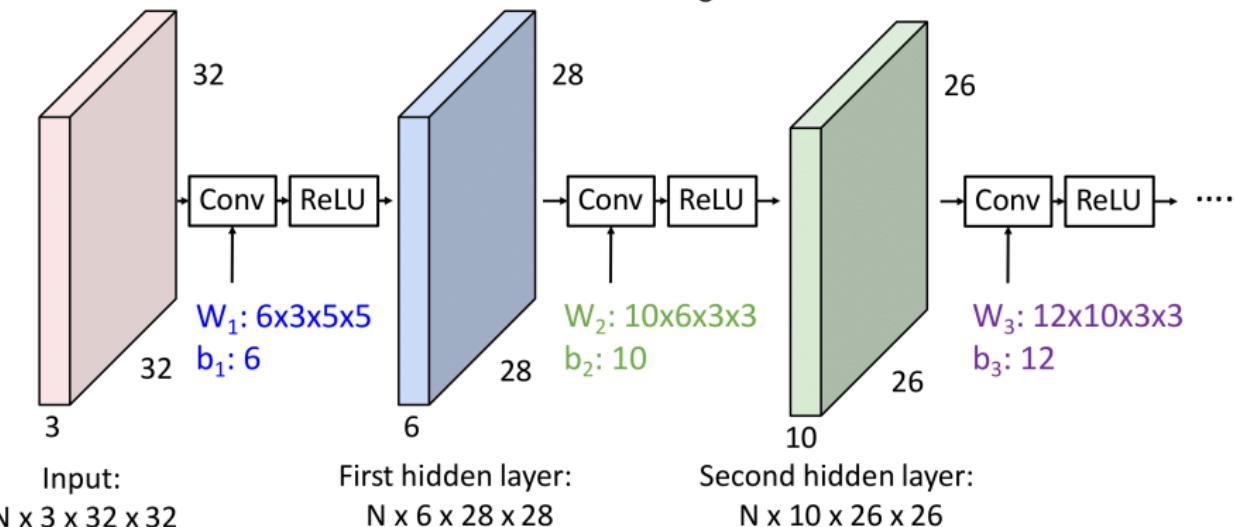
$$w^T x + b$$

Convolution Layer

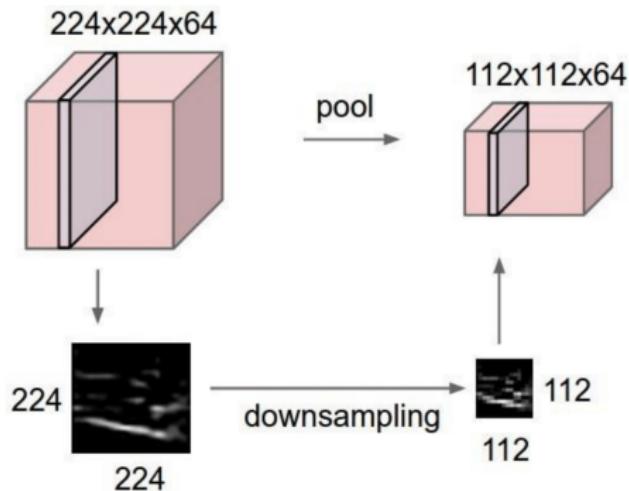


Stacking convolutions

Q: What happens if we stack two convolution layers? (Recall $y=W_2W_1x$ is a linear classifier)
A: We get another convolution!

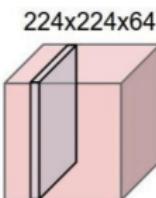
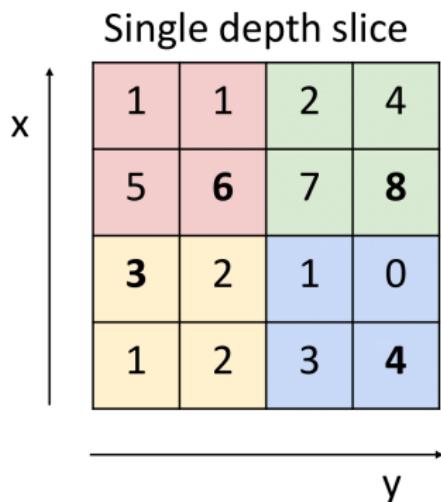


Pooling Layers: Another way to downsample



Hyperparameters:
Kernel Size
Stride
Pooling function

Max Pooling



Max pooling with 2x2 kernel size and stride 2



6	8
3	4

Introduces **invariance** to
small spatial shifts
No learnable parameters!

Pooling Summary

Input: $C \times H \times W$

Hyperparameters:

- Kernel size: K
- Stride: S
- Pooling function (max, avg)

Common settings:

max, $K = 2, S = 2$

max, $K = 3, S = 2$ (AlexNet)

Output: $C \times H' \times W'$ where

- $H' = (H - K) / S + 1$
- $W' = (W - K) / S + 1$

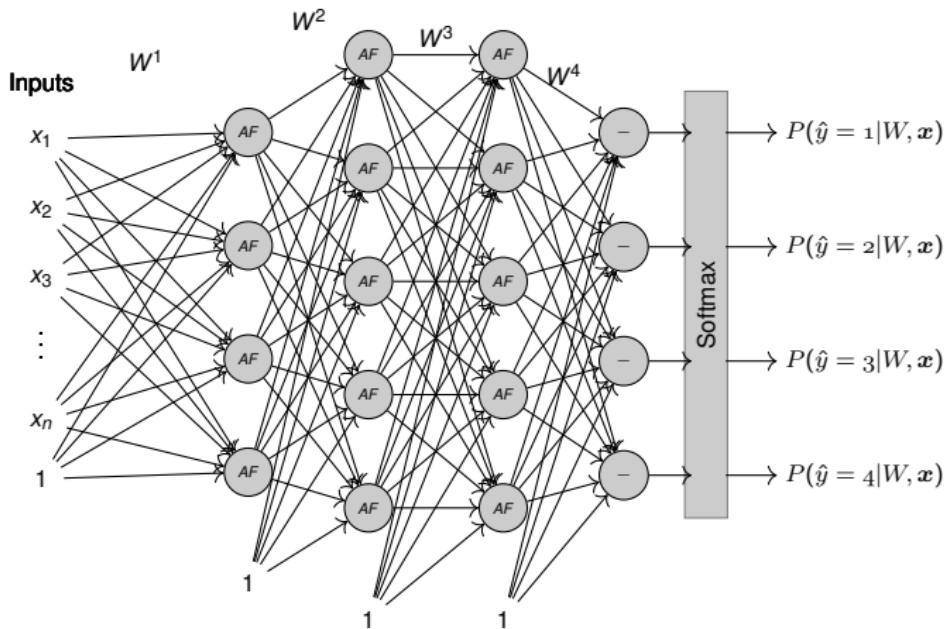
Learnable parameters: None!

Vanishing/Exploding Gradient

The vanishing gradient problem for deep networks

- A particular problem with training deep networks
 - The gradient of the error with respect to weights is unstable

Fully connected neural networks



$$\hat{y} = f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^2(W^2(f^1(W^1 \mathbf{x})))))$$

The problem with training deep networks

- A multilayer perceptron is a nested function

$$\hat{y} = f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^2(W^2(f^1(W^1 \mathbf{x}))))))$$

- W^k is the weights matrix at the k^{th} layer
- f^k is the activation function of the k^{th} layer
- The error for X can be written as

$$L = CE(f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^2(W^2(f^1(W^1 \mathbf{x})))))), y)$$

- CE is cross entropy

Training deep networks

- Let $\mathbf{z}^\ell = W^\ell \mathbf{o}^{\ell-1}$ and $\mathbf{o}^\ell = f^\ell(\mathbf{z}^\ell)$

$$\frac{\partial L}{\partial \mathbf{o}^{\ell-1}} = \text{---}$$

Training deep networks

- Let $\mathbf{z}^\ell = W^\ell \mathbf{o}^{\ell-1}$ and $\mathbf{o}^\ell = f^\ell(\mathbf{z}^\ell)$

$$\frac{\partial L}{\partial \mathbf{o}^{\ell-1}} = \dots \frac{\partial L}{\partial \mathbf{o}^\ell}$$

Training deep networks

- Let $\mathbf{z}^\ell = W^\ell \mathbf{o}^{\ell-1}$ and $\mathbf{o}^\ell = f^\ell(\mathbf{z}^\ell)$

$$\frac{\partial L}{\partial \mathbf{o}^{\ell-1}} = -\frac{\partial \mathbf{o}^\ell}{\partial \mathbf{z}^\ell} \frac{\partial L}{\partial \mathbf{o}^\ell}$$

Training deep networks

- Let $\mathbf{z}^\ell = W^\ell \mathbf{o}^{\ell-1}$ and $\mathbf{o}^\ell = f^\ell(\mathbf{z}^\ell)$

$$\frac{\partial L}{\partial \mathbf{o}^{\ell-1}} = \frac{\partial \mathbf{z}^\ell}{\partial \mathbf{o}^{\ell-1}} \frac{\partial \mathbf{o}^\ell}{\partial \mathbf{z}^\ell} \frac{\partial L}{\partial \mathbf{o}^\ell}$$

Training deep networks

- Let $\mathbf{z}^\ell = W^\ell \mathbf{o}^{\ell-1}$ and $\mathbf{o}^\ell = f^\ell(\mathbf{z}^\ell)$

$$\frac{\partial L}{\partial \mathbf{o}^{\ell-1}} = \frac{\partial \mathbf{z}^\ell}{\partial \mathbf{o}^{\ell-1}} \frac{\partial \mathbf{o}^\ell}{\partial \mathbf{z}^\ell} \frac{\partial L}{\partial \mathbf{o}^\ell}$$

$$\nabla_{\mathbf{o}^{\ell-1}} L = W^\ell \cdot \nabla_{\mathbf{z}^\ell} f^\ell \cdot \nabla_{\mathbf{o}^\ell} L$$

- Where
 - $\nabla_{\mathbf{x}} f$ is the jacobian matrix of $f(\mathbf{x})$ w.r.t \mathbf{x}

Training deep networks

- For

$$L = CE(f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^1(W^1 \mathbf{x})))), y)$$

Training deep networks

- For

$$L = CE(f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^1(W^1 \mathbf{x})))), y)$$

- We get

$$\nabla_{W^k} L = \mathbf{o}^{k-1} \cdot \nabla_{\mathbf{z}^k} f^k \cdot W^{k+1} \cdot \nabla_{\mathbf{z}^{k+1}} f^{k+1} \dots W^{N-1} \cdot \nabla_{\mathbf{z}^{N-1}} f^{N-1} \cdot W^N \cdot \nabla_{\mathbf{z}^N} f^N \cdot \nabla_{\hat{\mathbf{y}}} CE$$

$$\nabla_{W^k} L = \mathbf{o}^{k-1} \cdot \nabla_{\mathbf{z}^k} f^k \cdot \left(\prod_{i=k+1}^N W^i \cdot \nabla_{\mathbf{z}^i} f^i \right) \cdot \nabla_{\hat{\mathbf{y}}} CE$$

- Where

- W^N is the gradient of the affine transformation \mathbf{z}^n w.r.t the output of the $n - 1^{th}$ layer (\mathbf{o}^{n-1}) of the network
- $\nabla_{\mathbf{z}^n} f^n$ is gradient of activation function $f^n(\cdot)$ w.r.t to affine transformation \mathbf{z}^n in the n^{th} layer of the network
- \mathbf{o}^{n-1} is the gradient of the affine transformation \mathbf{z}^n w.r.t the weights of the n^{th} layer (W^n) of the network
- $\nabla_{\hat{\mathbf{y}}} CE$ is the gradient of cross entropy w.r.t the output of the network ($\hat{\mathbf{y}}$)

Training deep networks

- For

$$L = CE(f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^1(W^1 \mathbf{x})))), y)$$

- We get

$$\nabla_{W^k} L = \mathbf{o}^{k-1} \cdot \nabla_{\mathbf{z}^k} f^k \cdot W^{k+1} \cdot \nabla_{\mathbf{z}^{k+1}} f^{k+1} \dots W^{N-1} \cdot \nabla_{\mathbf{z}^{N-1}} f^{N-1} \cdot W^N \cdot \nabla_{\mathbf{z}^N} f^N \cdot \nabla_{\hat{y}} CE$$

$$\nabla_{W^k} L = \mathbf{o}^{k-1} \cdot \nabla_{\mathbf{z}^k} f^k \cdot \left(\prod_{i=k+1}^N W^i \cdot \nabla_{\mathbf{z}^i} f^i \right) \cdot \nabla_{\hat{y}} CE$$

- Where

- W^N is the gradient of the affine transformation \mathbf{z}^n w.r.t the output of the $n - 1^{th}$ layer (\mathbf{o}^{n-1}) of the network
- $\nabla_{\mathbf{z}^n} f^n$ is gradient of activation function $f^n(\cdot)$ w.r.t to affine transformation \mathbf{z}^n in the n^{th} layer of the network
- \mathbf{o}^{n-1} is the gradient of the affine transformation \mathbf{z}^n w.r.t the weights of the n^{th} layer (W^n) of the network
- $\nabla_{\hat{y}} CE$ is the gradient of cross entropy w.r.t the output of the network (\hat{y})

Training deep networks

- For

$$L = CE(f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^1(W^1 \mathbf{x})))), y)$$

- We get

$$\nabla_{W^k} L = \mathbf{o}^{k-1} \cdot \nabla_{\mathbf{z}^k} f^k \cdot \mathbf{W}^{k+1} \cdot \nabla_{\mathbf{z}^{k+1}} f^{k+1} \dots \mathbf{W}^{N-1} \cdot \nabla_{\mathbf{z}^{N-1}} f^{N-1} \cdot \mathbf{W}^N \cdot \nabla_{\mathbf{z}^N} f^N \cdot \nabla_{\hat{y}} CE$$

$$\nabla_{W^k} L = \mathbf{o}^{k-1} \cdot \nabla_{\mathbf{z}^k} f^k \cdot \left(\prod_{i=k+1}^N \mathbf{W}^i \cdot \nabla_{\mathbf{z}^i} f^i \right) \cdot \nabla_{\hat{y}} CE$$

- Where

- \mathbf{W}^N is the gradient of the affine transformation \mathbf{z}^n w.r.t the output of the $n - 1^{th}$ layer (\mathbf{o}^{n-1}) of the network
- $\nabla_{\mathbf{z}^n} f^n$ is gradient of activation function $f^n(\cdot)$ w.r.t to affine transformation \mathbf{z}^n in the n^{th} layer of the network
- \mathbf{o}^{n-1} is the gradient of the affine transformation \mathbf{z}^n w.r.t the weights of the n^{th} layer (W^n) of the network
- $\nabla_{\hat{y}} CE$ is the gradient of cross entropy w.r.t the output of the network (\hat{y})

Training deep networks

- For

$$L = CE(f^N(W^N f^{N-1}(W^{N-1} f^{N-2}(\dots f^1(W^1 \mathbf{x})))), y)$$

- We get

$$\nabla_{W^k} L = \sigma^{k-1} \cdot \nabla_{\mathbf{z}^k} f^k \cdot W^{k+1} \cdot \nabla_{\mathbf{z}^{k+1}} f^{k+1} \cdots W^{N-1} \cdot \nabla_{\mathbf{z}^{N-1}} f^{N-1} \cdot W^N \cdot \nabla_{\mathbf{z}^N} f^N \cdot \nabla_{\hat{y}} CE$$

$$\nabla_{W^k} L = \sigma^{k-1} \cdot \nabla_{\mathbf{z}^k} f^k \cdot \left(\prod_{i=k+1}^N W^i \cdot \nabla_{\mathbf{z}^i} f^i \right) \cdot \nabla_{\hat{y}} CE$$

- Repeated multiplication by the weights matrix and jacobian matrices of activation functions will result in exploding or vanishing gradients
 - Depends on the spectral norm (largest singular value) of jacobian matrices
- The gradients in the lower/earlier layers can explode or vanish
 - Resulting in insignificant or unstable gradient descent updates
 - Problem gets worse as network depth increases

Batch Normalization

Batch Normalization

Idea: “Normalize” the outputs of a layer so they have zero mean and unit variance

Why? Helps reduce “internal covariate shift”, improves optimization

We can normalize a batch of activations like this:

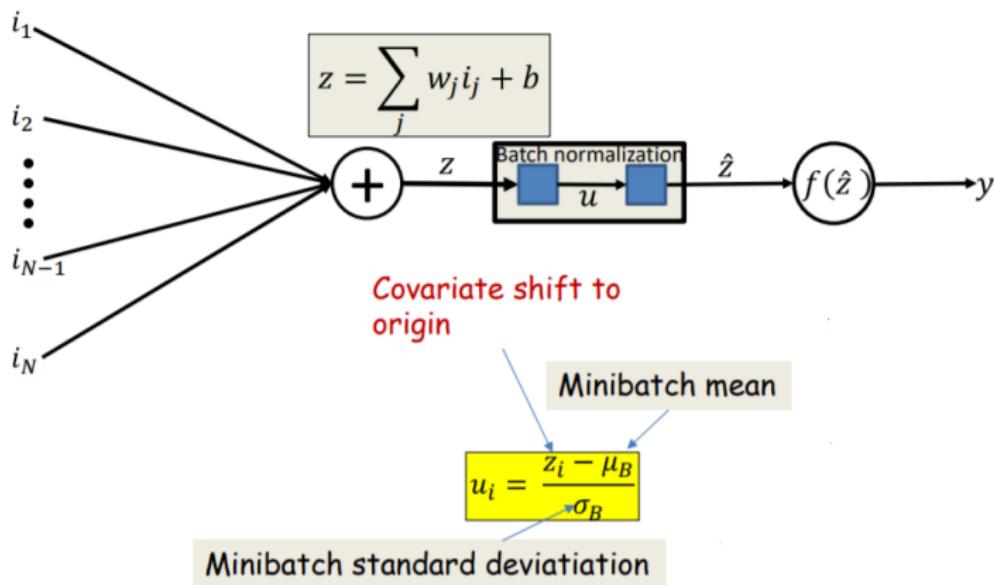
$$\hat{x}^{(k)} = \frac{x^{(k)} - \text{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

This is a **differentiable function**, so we can use it as an operator in our networks and backprop through it!

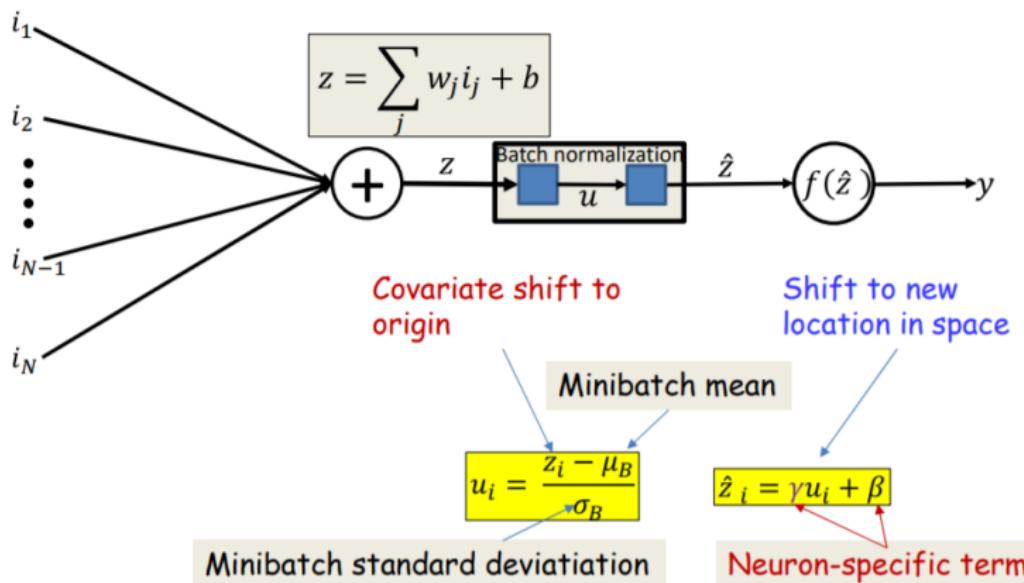
Ioffe and Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift”, ICML 2015

- "We define Internal Covariate Shift as the change in the distribution of network activations due to the change in network parameters during training."

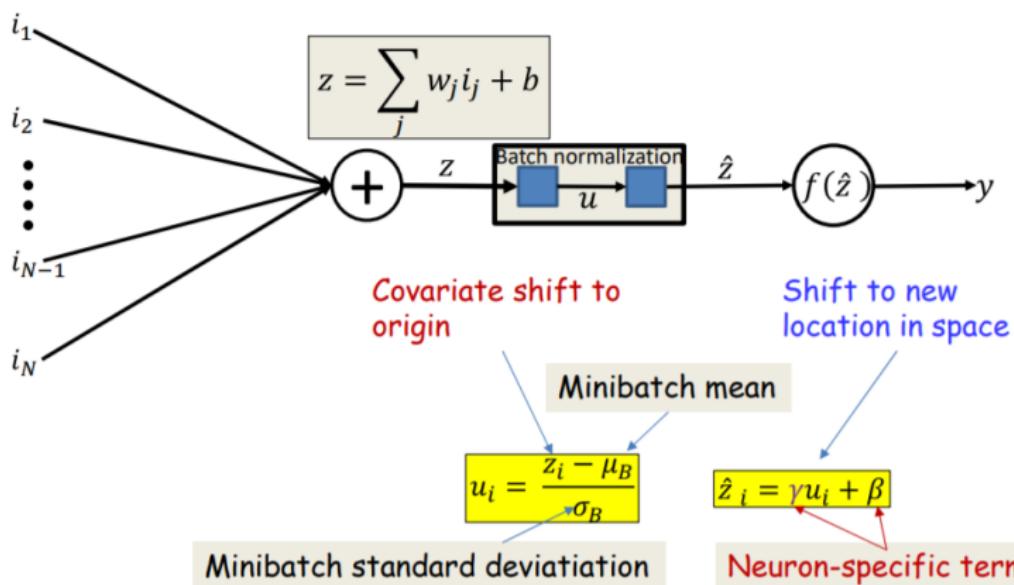
Batch Normalization



Batch Normalization



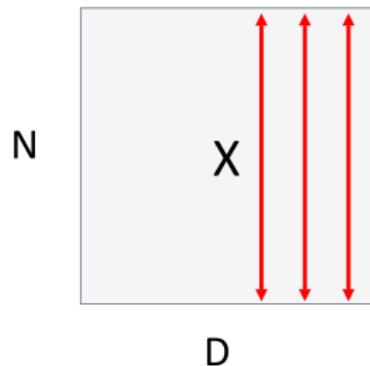
Batch Normalization



- Note that simply normalizing each input of a layer may change what the layer can represent. For instance, normalizing the inputs of a sigmoid would constrain them to the linear regime of the nonlinearity. To address this, we make sure that the transformation inserted in the network can represent the identity transform.

Batch Normalization

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 std, shape is D

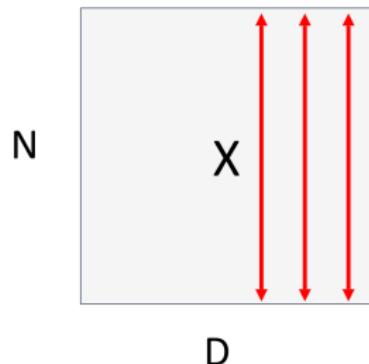
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Normalized x,
Shape is N x D

Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Batch Normalization

Input: $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} \quad \text{Per-channel mean, shape is } D$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 \quad \text{Per-channel std, shape is } D$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}} \quad \text{Normalized x, Shape is } N \times D$$

Problem: What if zero-mean, unit variance is too hard of a constraint?

Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Batch Normalization

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 std, shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$
 Normalized x,
Shape is $N \times D$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$
 Output,
Shape is $N \times D$

Learnable scale and shift parameters:

$$\gamma, \beta : D$$

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

Batch Normalization: Test-Time

Problem: Estimates depend on minibatch; can't do this at test-time!

Input: $x : N \times D$

Learnable scale and shift parameters:

$\gamma, \beta : D$

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

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} \quad \text{Per-channel mean, shape is } D$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 \quad \text{Per-channel std, shape is } D$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}} \quad \begin{matrix} \text{Normalized } x, \\ \text{Shape is } N \times D \end{matrix}$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j \quad \begin{matrix} \text{Output,} \\ \text{Shape is } N \times D \end{matrix}$$

Batch Normalization: Test-Time

Input: $x : N \times D$

$\mu_j =$ (Running) average of values seen during training Per-channel mean, shape is D

$\sigma_j^2 =$ (Running) average of values seen during training Per-channel std, shape is D

$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$ Normalized x,
Shape is $N \times D$

$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$ Output,
Shape is $N \times D$

Learnable scale and shift parameters:

$\gamma, \beta : D$

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

Batch Normalization: Test-Time

Input: $x : N \times D$

$\mu_j =$ (Running) average of values seen during training Per-channel mean, shape is D

$\sigma_j^2 =$ (Running) average of values seen during training Per-channel std, shape is D

$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$ Normalized x,
Shape is $N \times D$

$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$ Output,
Shape is $N \times D$

Learnable scale and shift parameters:

$\gamma, \beta : D$

During testing batchnorm becomes a linear operator!

Can be fused with the previous fully-connected or conv layer

Batch Normalization for ConvNets

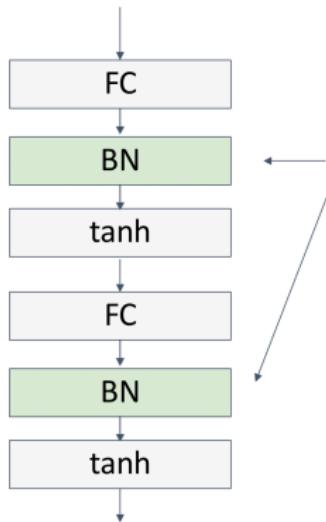
Batch Normalization for
fully-connected networks

$$\begin{aligned} \mathbf{x}: \mathbf{N} &\times \mathbf{D} \\ \text{Normalize} & \downarrow \\ \boldsymbol{\mu}, \boldsymbol{\sigma}: \mathbf{1} &\times \mathbf{D} \\ \boldsymbol{\gamma}, \boldsymbol{\beta}: \mathbf{1} &\times \mathbf{D} \\ \mathbf{y} = \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta} \end{aligned}$$

Batch Normalization for
convolutional networks
(Spatial Batchnorm, BatchNorm2D)

$$\begin{aligned} \mathbf{x}: \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W} \\ \text{Normalize} & \downarrow \quad \downarrow \quad \downarrow \\ \boldsymbol{\mu}, \boldsymbol{\sigma}: \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\ \boldsymbol{\gamma}, \boldsymbol{\beta}: \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\ \mathbf{y} = \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta} \end{aligned}$$

Batch Normalization

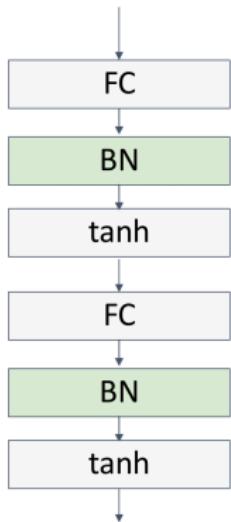


Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

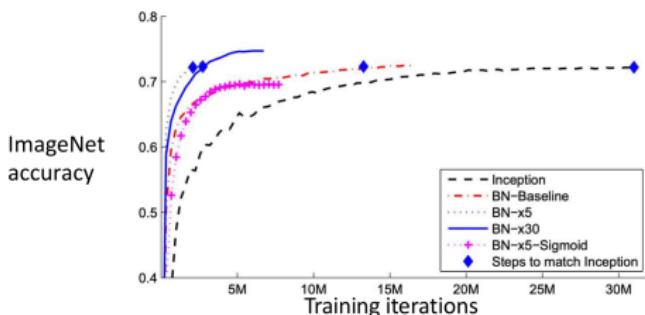
$$\hat{x}^{(k)} = \frac{x^{(k)} - \text{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Batch Normalization

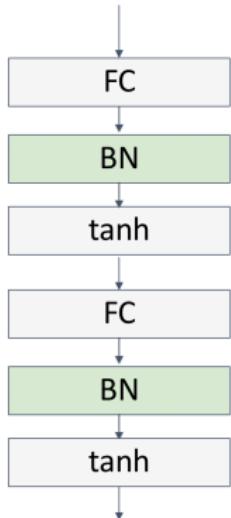


- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!



Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Batch Normalization

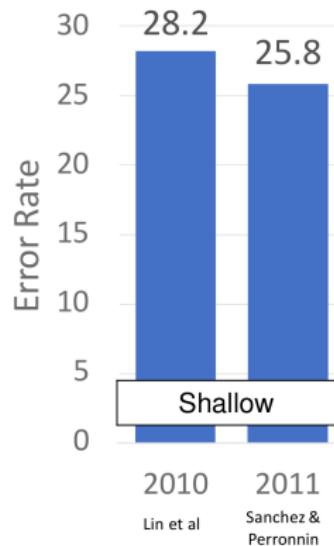


- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!
- **Not well-understood theoretically (yet)**
- Behaves differently during training and testing: this is a **very common source of bugs!**

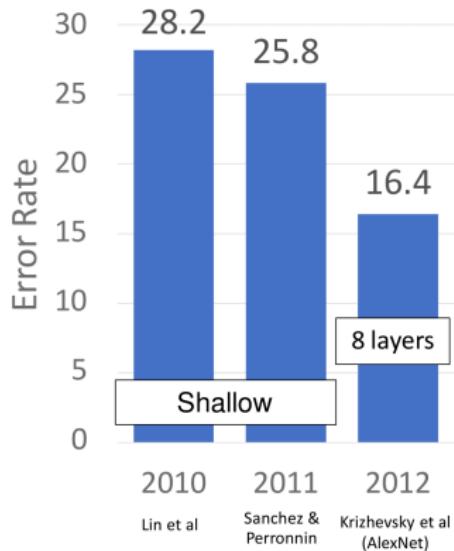
Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

CNN Architectures

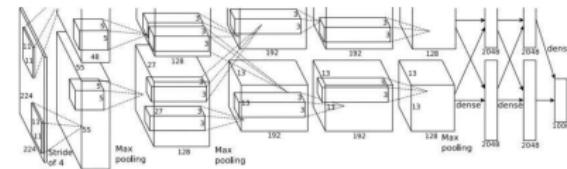
ImageNet classification challenge



ImageNet classification challenge



AlexNet



227 x 227 inputs

5 Convolutional layers

Max pooling

3 fully-connected layers

ReLU nonlinearities

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

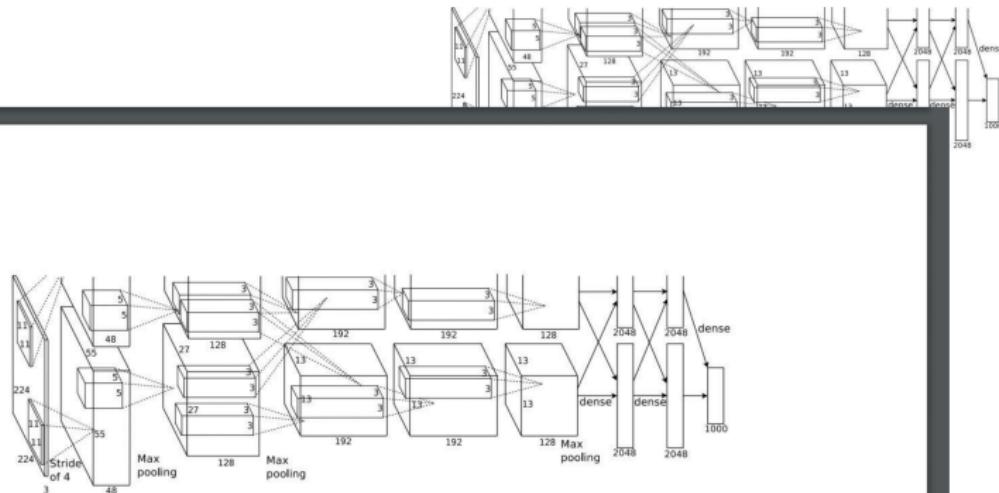
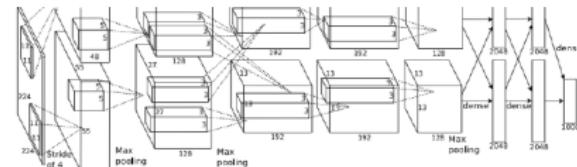


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input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

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet

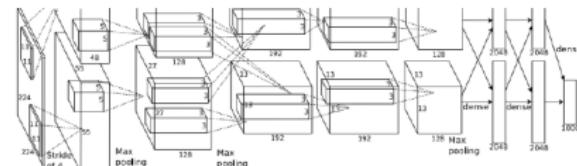


Layer	Input size		Layer					Output size	
	C	H / W	filters	kernel	stride	pad	C	H / W	
conv1	3	227	64	11	4	2	?		

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet



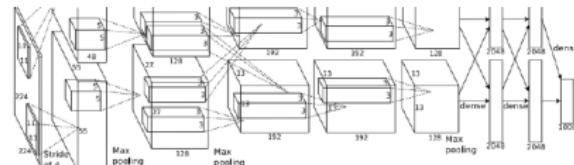
	Input size		Layer					Output size	
Layer	C	H / W	filters	kernel	stride	pad	C	H / W	
conv1	3	227	64	11	4	2	64	?	

Recall: Output channels = number of filters

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet



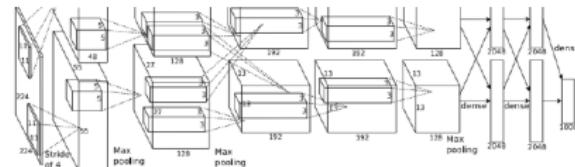
	Input size		Layer					Output size	
Layer	C	H / W	filters	kernel	stride	pad	C	H / W	
conv1	3	227	64	11	4	2	64	56	

$$\begin{aligned} \text{Recall: } W' &= (W - K + 2P) / S + 1 \\ &= 227 - 11 + 2*2) / 4 + 1 \\ &= 220/4 + 1 = 56 \end{aligned}$$

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet

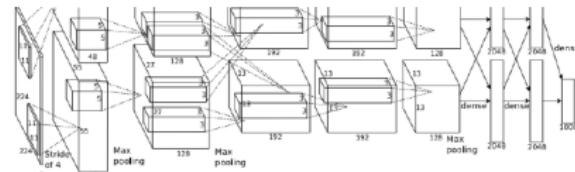


Layer	Input size		Layer					Output size			memory (KB)
	C	H / W	filters	kernel	stride	pad	C	H / W			
conv1	3	227	64	11	4	2	64	56	?		

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet



	Input size		Layer					Output size			
Layer	C	H / W	filters	kernel	stride	pad	C	H / W	memory (KB)		
conv1	3	227	64	11	4	2	64	56	784		

$$\begin{aligned} \text{Number of output elements} &= C * H' * W' \\ &= 64 * 56 * 56 = 200,704 \end{aligned}$$

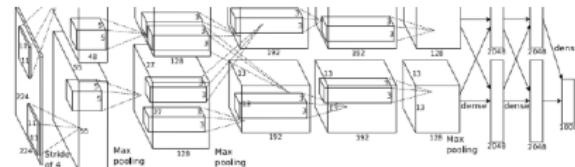
Bytes per element = 4 (for 32-bit floating point)

$$\begin{aligned} \text{KB} &= (\text{number of elements}) * (\text{bytes per elem}) / 1024 \\ &= 200704 * 4 / 1024 \\ &= \mathbf{784} \end{aligned}$$

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet

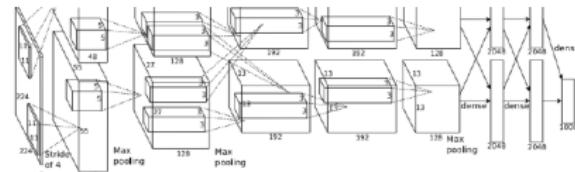


Layer	Input size		Layer					Output size			memory (KB)	params (k)
	C	H / W	filters	kernel	stride	pad	C	H / W				
conv1	3	227	64	11	4	2	64	56			784	?

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet



Layer	Input size		Layer					Output size		memory (KB)	params (k)
	C	H / W	filters	kernel	stride	pad	C	H / W			
conv1	3	227	64	11	4	2	64	56	784	23	

$$\begin{aligned}\text{Weight shape} &= C_{\text{out}} \times C_{\text{in}} \times K \times K \\ &= 64 \times 3 \times 11 \times 11\end{aligned}$$

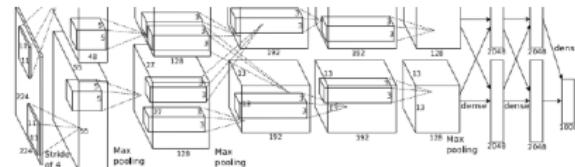
$$\text{Bias shape} = C_{\text{out}} = 64$$

$$\begin{aligned}\text{Number of weights} &= 64 * 3 * 11 * 11 + 64 \\ &= \mathbf{23,296}\end{aligned}$$

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet

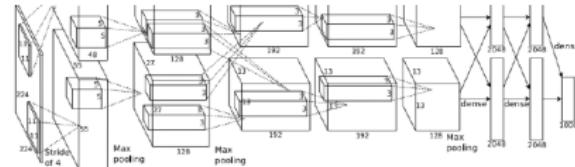


Layer	Input size		Layer					Output size			memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W					
conv1	3	227	64	11	4	2	64	56			784	23	?

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet



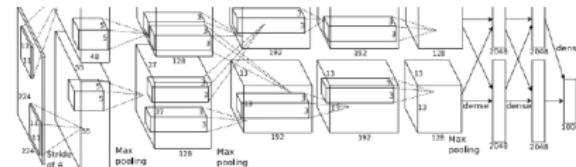
Layer	Input size		Layer					Output size		memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W				
conv1	3	227	64	11	4	2	64	56	784	23	73	

Number of floating point operations (multiply+add)
 $= (\text{number of output elements}) * (\text{ops per output elem})$
 $= (C_{\text{out}} \times H' \times W') * (C_{\text{in}} \times K \times K)$
 $= (64 * 56 * 56) * (3 * 11 * 11)$
 $= 200,704 * 363$
 $= \mathbf{72,855,552}$

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet

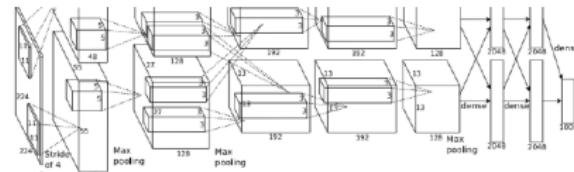


Layer	Input size		Layer				Output size		memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W			
conv1	3	227	64	11	4	2	64	56	784	23	73
pool1	64	56		3	2	0	?				

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet



Layer	Input size		Layer				Output size		memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W			
conv1	3	227	64	11	4	2	64	56	784	23	73
pool1	64	56		3	2	0	64	27			

For pooling layer:

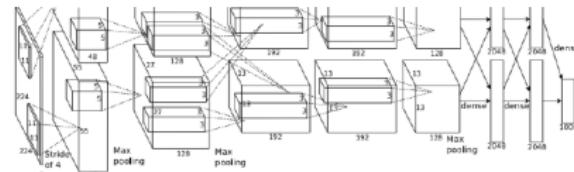
$$\# \text{output channels} = \# \text{input channels} = 64$$

$$\begin{aligned} W' &= \text{floor}((W - K) / S + 1) \\ &= \text{floor}(53 / 2 + 1) = \text{floor}(27.5) = 27 \end{aligned}$$

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet



Layer	Input size		Layer				Output size		memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W			
conv1	3	227	64	11	4	2	64	56	784	23	73
pool1	64	56		3	2	0	64	27	182	?	

$$\# \text{output elems} = C_{\text{out}} \times H' \times W'$$

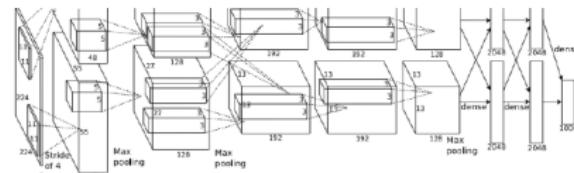
Bytes per elem = 4

$$\begin{aligned} \text{KB} &= C_{\text{out}} * H' * W' * 4 / 1024 \\ &= 64 * 27 * 27 * 4 / 1024 \\ &= \mathbf{182.25} \end{aligned}$$

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet



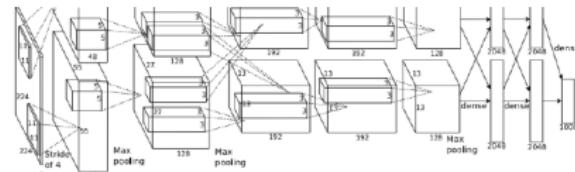
Layer	Input size		Layer				Output size		memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W			
conv1	3	227	64	11	4	2	64	56	784	23	73
pool1	64	56		3	2	0	64	27	182	0	?

Pooling layers have no learnable parameters!

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet



Layer	Input size		Layer				Output size		memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W			
conv1	3	227	64	11	4	2	64	56	784	23	73
pool1	64	56		3	2	0	64	27	182	0	0

Floating-point ops for pooling layer

$$= (\text{number of output positions}) * (\text{flops per output position})$$

$$= (C_{\text{out}} * H' * W') * (K * K)$$

$$= (64 * 27 * 27) * (3 * 3)$$

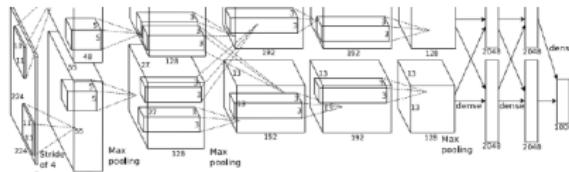
$$= 419,904$$

$$= \mathbf{0.4 \text{ MFLOP}}$$

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet



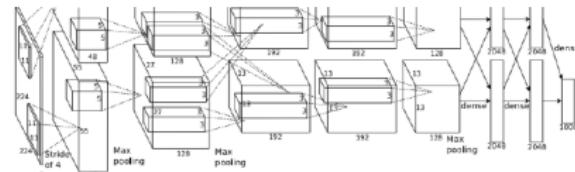
Layer	Input size		Layer					Output size		memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W				
conv1	3	227	64	11	4	2	64	56	784	23	73	
pool1	64	56		3	2	0	64	27	182	0	0	
conv2	64	27	192	5	1	2	192	27	547	307	224	
pool2	192	27		3	2	0	192	13	127	0	0	
conv3	192	13	384	3	1	1	384	13	254	664	112	
conv4	384	13	256	3	1	1	256	13	169	885	145	
conv5	256	13	256	3	1	1	256	13	169	590	100	
pool5	256	13		3	2	0	256	6	36	0	0	
flatten	256	6					9216		36	0	0	

$$\begin{aligned}\text{Flatten output size} &= C_{in} \times H \times W \\ &= 256 * 6 * 6 \\ &= \mathbf{9216}\end{aligned}$$

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

AlexNet



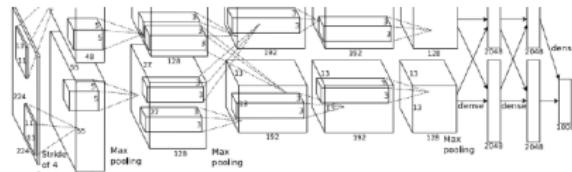
Layer	Input size		Layer					Output size		memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W				
conv1	3	227	64	11	4	2	64	56	784	23	73	
pool1	64	56		3	2	0	64	27	182	0	0	
conv2	64	27	192	5	1	2	192	27	547	307	224	
pool2	192	27		3	2	0	192	13	127	0	0	
conv3	192	13	384	3	1	1	384	13	254	664	112	
conv4	384	13	256	3	1	1	256	13	169	885	145	
conv5	256	13	256	3	1	1	256	13	169	590	100	
pool5	256	13		3	2	0	256	6	36	0	0	
flatten	256	6					9216		36	0	0	
fc6	9216		4096				4096		16	37,749	38	

$$\begin{aligned} \text{FC params} &= C_{\text{in}} * C_{\text{out}} + C_{\text{out}} \\ &= 9216 * 4096 + 4096 \\ &= 37,725,832 \end{aligned}$$

$$\begin{aligned} \text{FC flops} &= C_{\text{in}} * C_{\text{out}} \\ &= 9216 * 4096 \\ &= 37,748,736 \end{aligned}$$

AlexNet

AlexNet

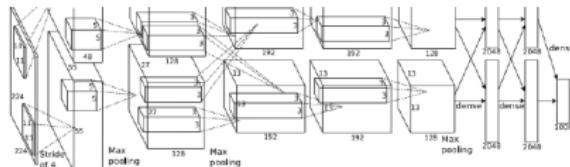


Layer	Input size		Layer					Output size		memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W				
conv1	3	227	64	11	4	2	64	56	784	23	73	
pool1	64	56		3	2	0	64	27	182	0	0	
conv2	64	27	192	5	1	2	192	27	547	307	224	
pool2	192	27		3	2	0	192	13	127	0	0	
conv3	192	13	384	3	1	1	384	13	254	664	112	
conv4	384	13	256	3	1	1	256	13	169	885	145	
conv5	256	13	256	3	1	1	256	13	169	590	100	
pool5	256	13		3	2	0	256	6	36	0	0	
flatten	256	6					9216		36	0	0	
fc6	9216		4096				4096		16	37,749	38	
fc7	4096		4096				4096		16	16,777	17	
fc8	4096		1000				1000		4	4,096	4	

AlexNet

AlexNet

How to choose this?
Trial and error =(

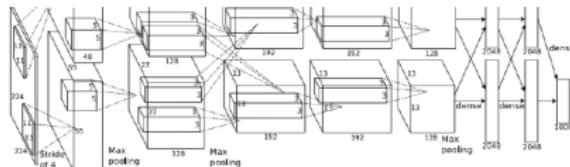


Layer	Input size		Layer				Output size		memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W			
conv1	3	227	64	11	4	2	64	56	784	23	73
pool1	64	56		3	2	0	64	27	182	0	0
conv2	64	27	192	5	1	2	192	27	547	307	224
pool2	192	27		3	2	0	192	13	127	0	0
conv3	192	13	384	3	1	1	384	13	254	664	112
conv4	384	13	256	3	1	1	256	13	169	885	145
conv5	256	13	256	3	1	1	256	13	169	590	100
pool5	256	13		3	2	0	256	6	36	0	0
flatten	256	6					9216		36	0	0
fc6	9216		4096				4096		16	37,749	38
fc7	4096		4096				4096		16	16,777	17
fc8	4096		1000				1000		4	4,096	4

AlexNet

AlexNet

Interesting trends here!



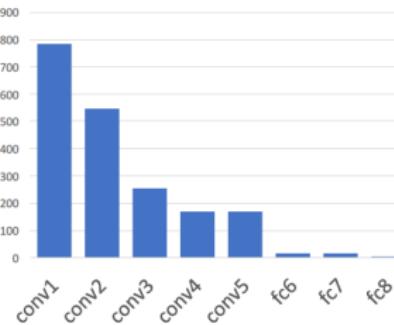
Layer	Input size		Layer					Output size		memory (KB)	params (k)	flop (M)
	C	H / W	filters	kernel	stride	pad	C	H / W				
conv1	3	227	64	11	4	2	64	56	784	23	73	
pool1	64	56		3	2	0	64	27	182	0	0	
conv2	64	27	192	5	1	2	192	27	547	307	224	
pool2	192	27		3	2	0	192	13	127	0	0	
conv3	192	13	384	3	1	1	384	13	254	664	112	
conv4	384	13	256	3	1	1	256	13	169	885	145	
conv5	256	13	256	3	1	1	256	13	169	590	100	
pool5	256	13		3	2	0	256	6	36	0	0	
flatten	256	6					9216		36	0	0	
fc6	9216		4096				4096		16	37,749	38	
fc7	4096		4096				4096		16	16,777	17	
fc8	4096		1000				1000		4	4,096	4	

AlexNet

AlexNet

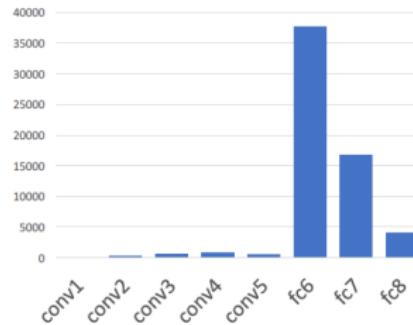
Most of the **memory usage** is in the early convolution layers

Memory (KB)



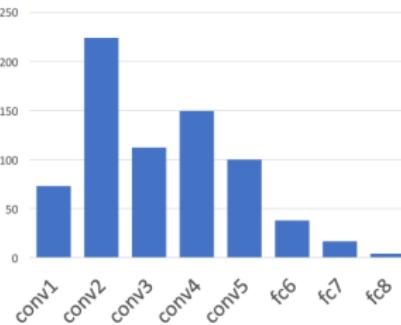
Nearly all **parameters** are in the fully-connected layers

Params (K)



Most **floating-point ops** occur in the convolution layers

MFLOP



References

- Deep Learning, Ian Goodfellow, MIT Press, Ch. 9
- Aston Zhang, Dive into Deep Learning, Ch 7 & Ch. 8