

# Deep Learning

## Lecture 8

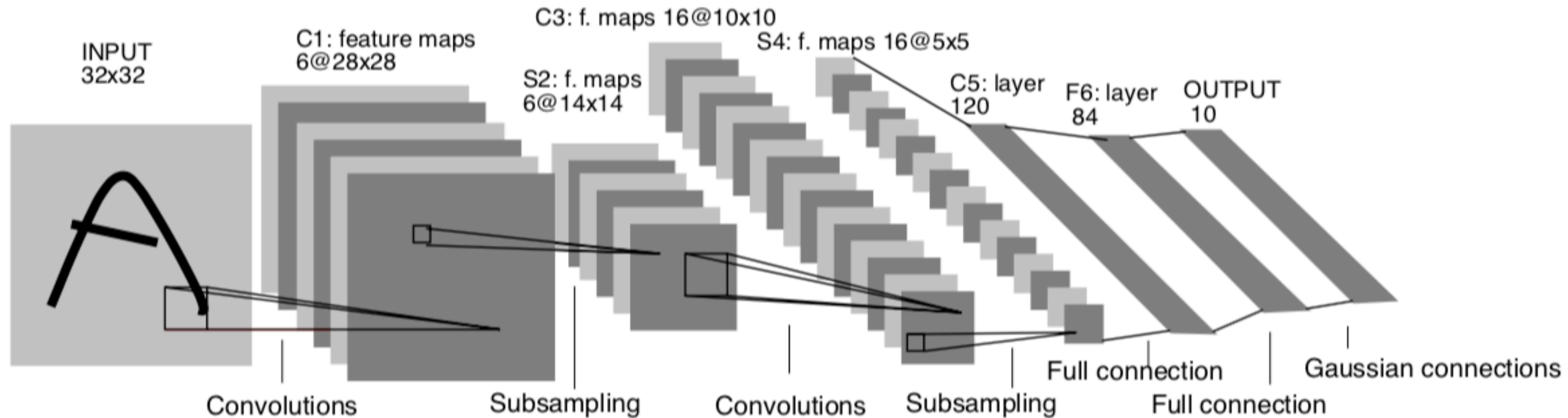
# From LeNet to ResNet

A Brief History of CNNs

Part 1

# LeNet-5

[LeCun et al., 1998]



Conv filters were 5x5, stride 1

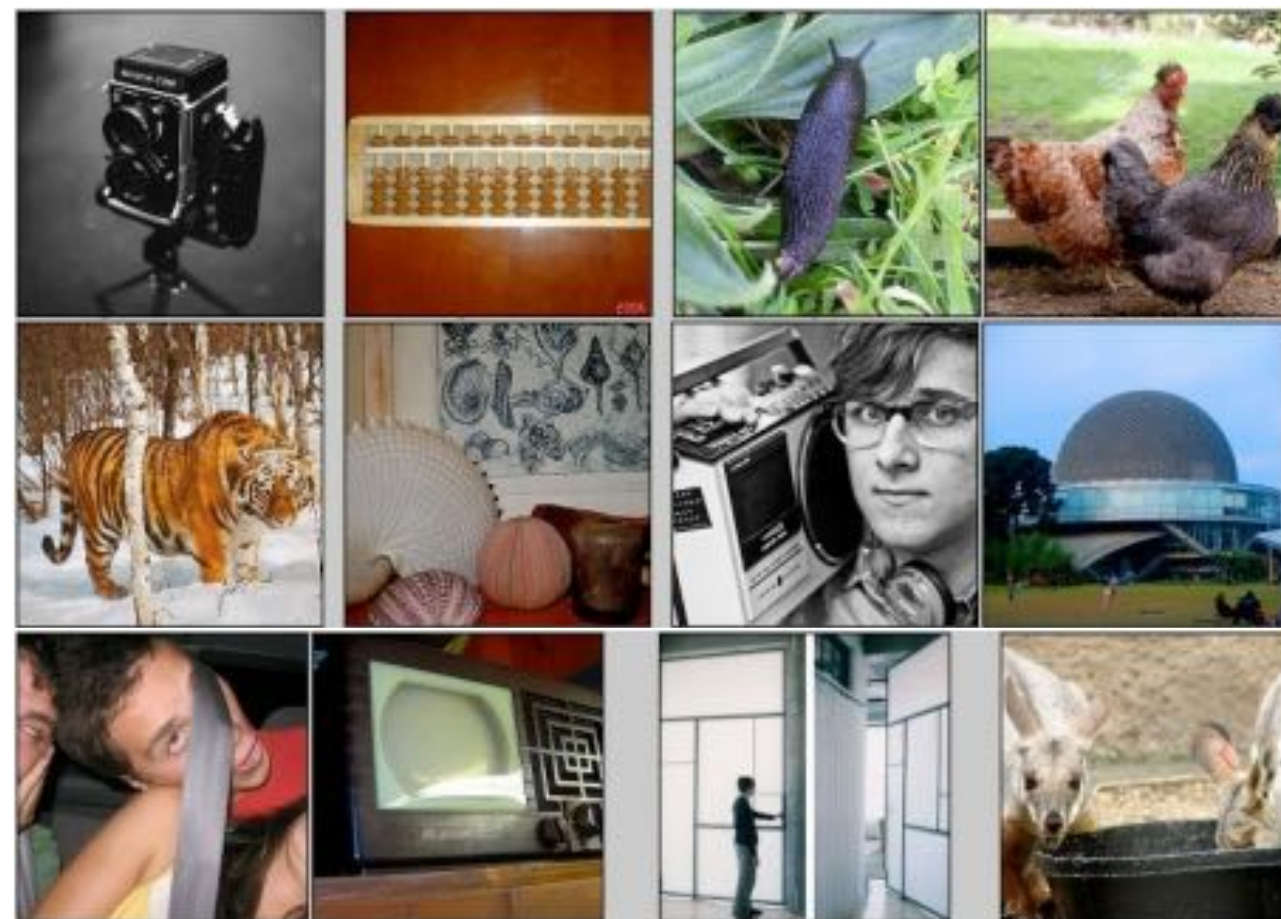
Subsampling (Pooling) layers were 2x2, stride 2

Architecture is CONV-POOL-CONV-POOL-FC-FC

Fun fact: MNIST (modified NIST) hand-written digits dataset introduced in this paper!

# Fast forward to the arrival of big visual data ...

IMAGENET

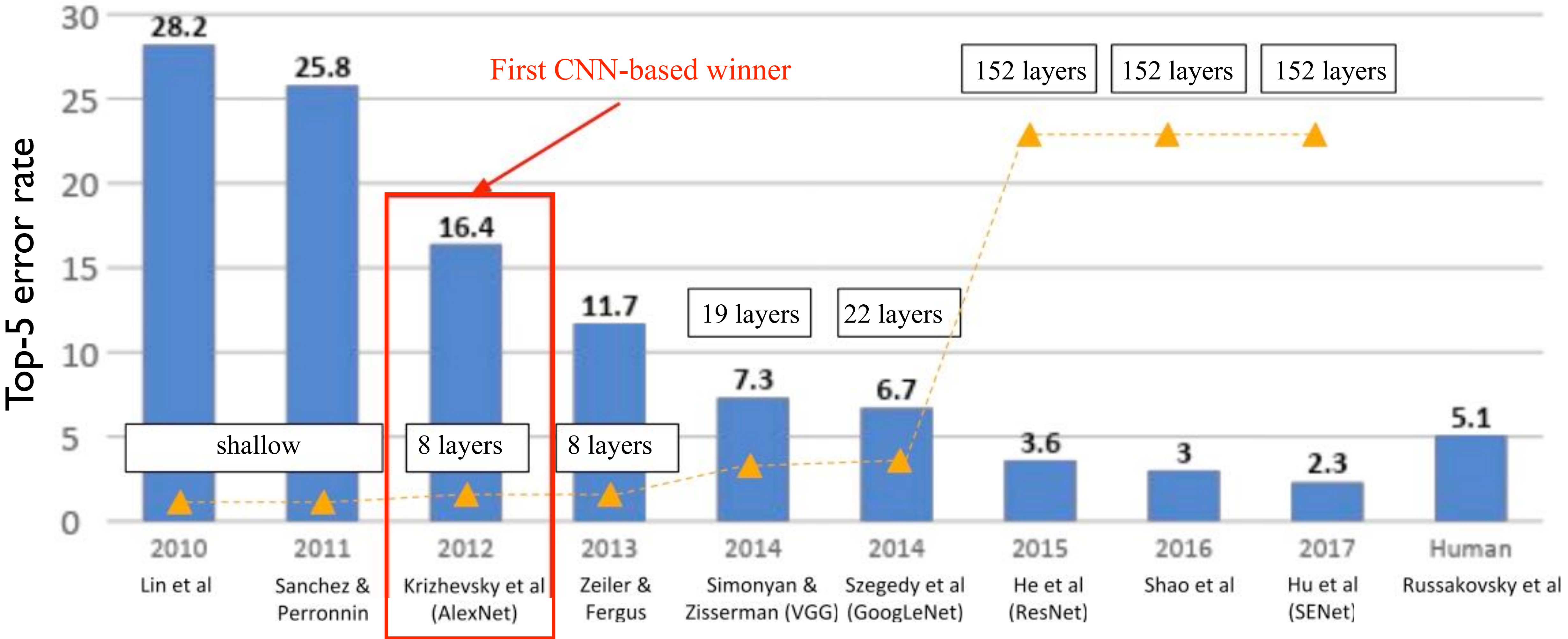


- ~14 million labeled images, 20k classes
- Images gathered from Internet
- Human labels via Amazon Mechanical Turk
- ImageNet Large-Scale Visual Recognition Challenge (ILSVRC):  
1.2 million training images, 1000 classes

[www.image-net.org/challenges/LSVRC/](http://www.image-net.org/challenges/LSVRC/)



# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



# Case Study: AlexNet

*[Krizhevsky et al. 2012]*

## Architecture:

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

CONV5

Max POOL3

FC6

FC7

FC8

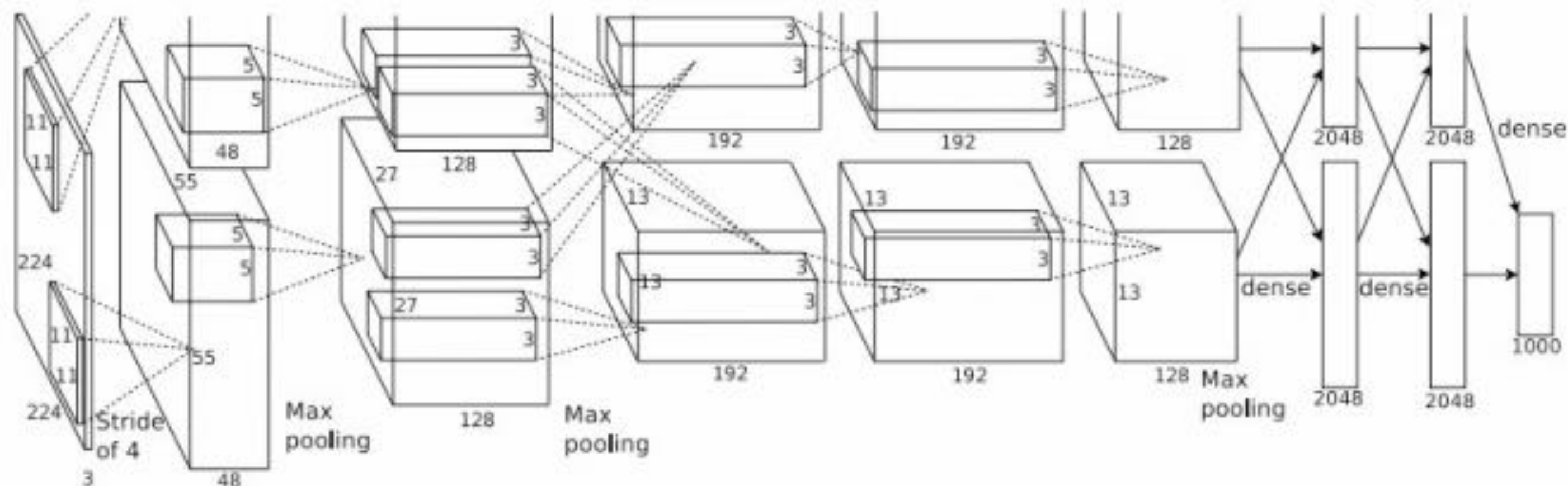
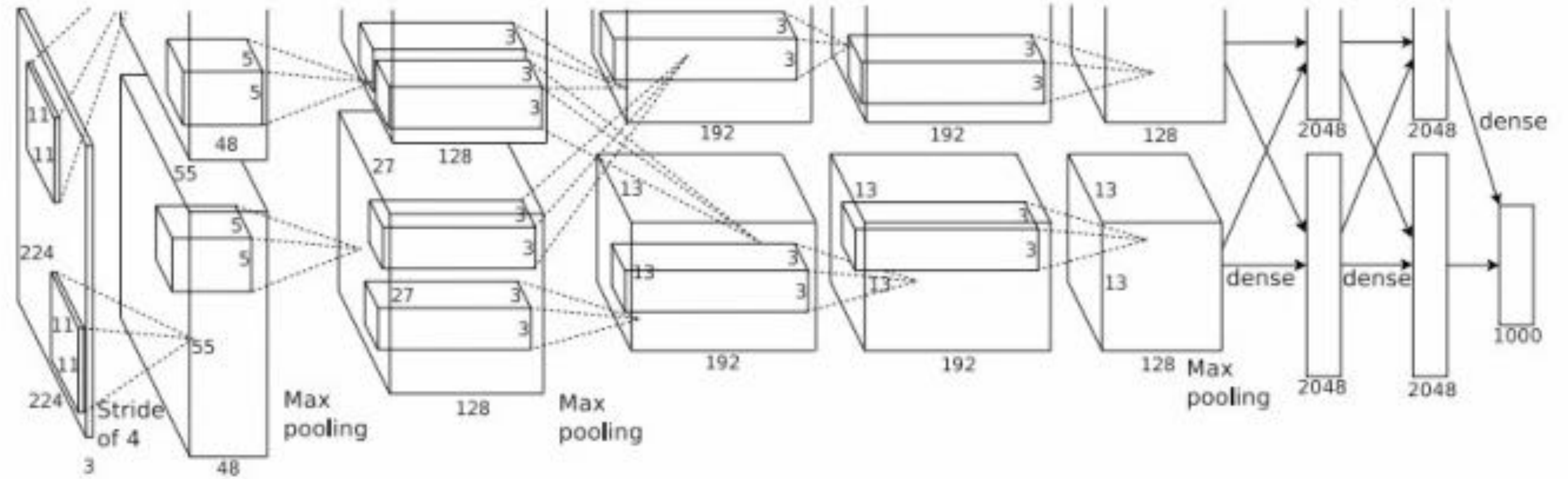


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

**First layer (CONV1):** 96 11x11 filters applied at stride 4

 $\Rightarrow$ 

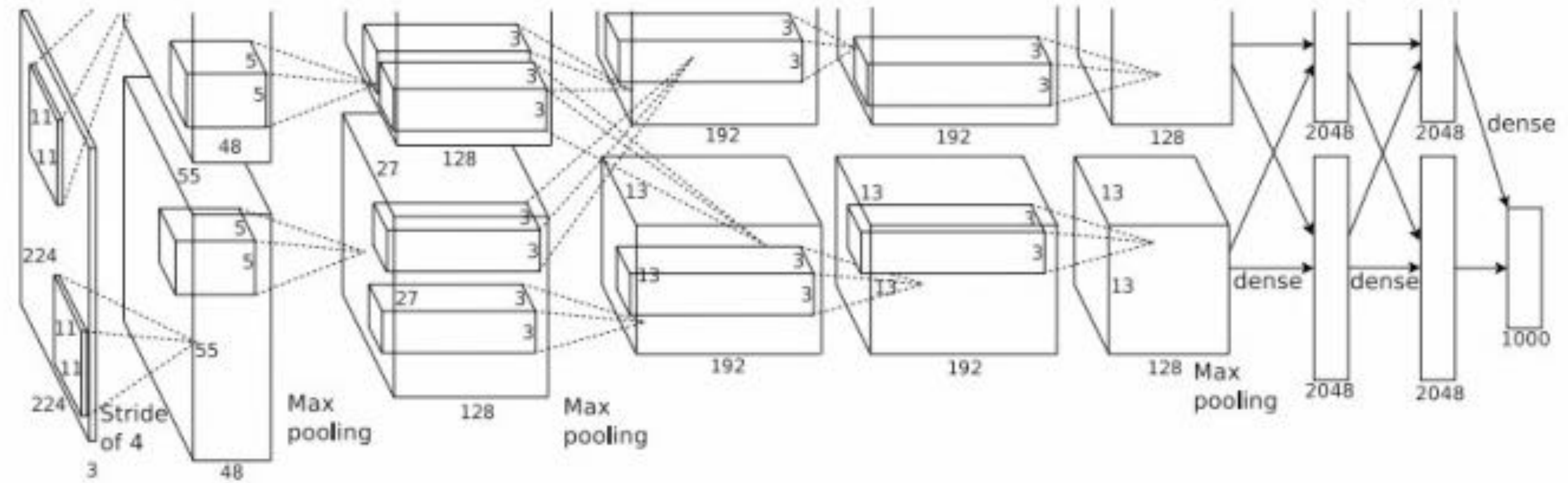
Q: what is the output volume size? Hint:  $(227-11)/4+1 = 55$

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



# Case Study: AlexNet

*[Krizhevsky et al. 2012]*



Input: 227x227x3 images

**First layer (CONV1):** 96 11x11 filters applied at stride 4

=>

Output volume [**55x55x96**]

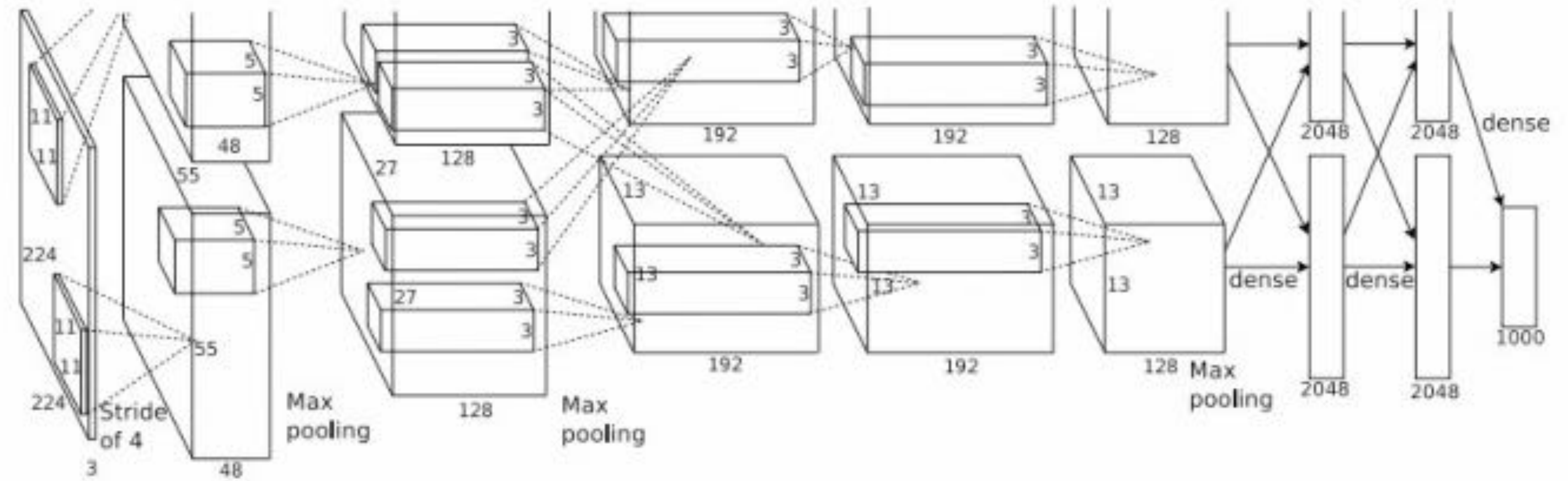
Q: What is the total number of parameters in this layer?

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



# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

**First layer (CONV1):** 96 11x11 filters applied at stride 4

 $\Rightarrow$ 

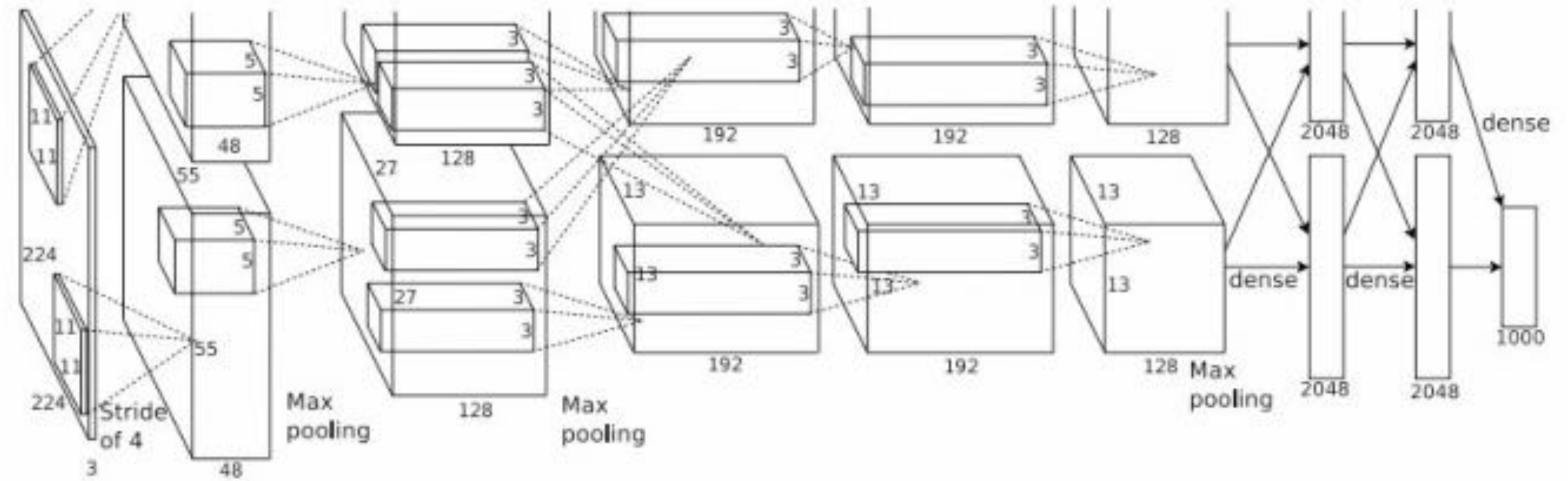
Output volume [55x55x96]

Parameters:  $(11 \times 11 \times 3) \times 96 \sim 34,848$ 

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

# Case Study: AlexNet

*[Krizhevsky et al. 2012]*



Input: 227x227x3 images

After CONV1: 55x55x96

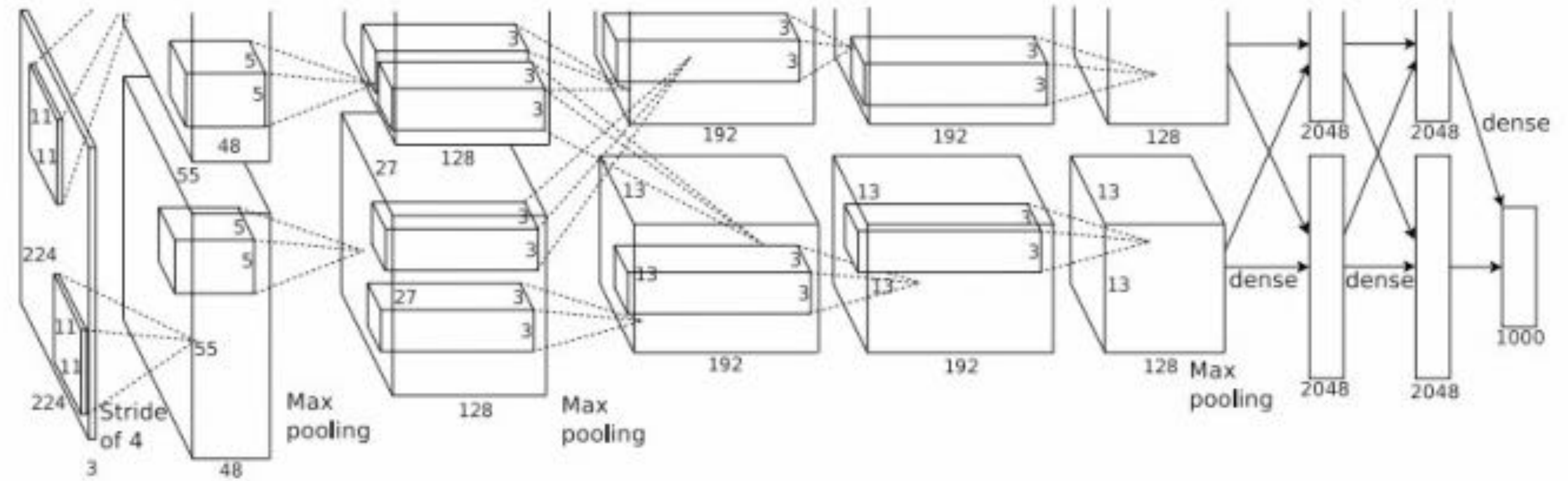
**Second layer (POOL1):** 3x3 filters applied at stride 2

Q: what is the output volume size? Hint:  $(55-3)/2+1 = 27$

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

# Case Study: AlexNet

*[Krizhevsky et al. 2012]*



Input: 227x227x3 images

After CONV1: 55x55x96

**Second layer (POOL1):** 3x3 filters applied at stride 2

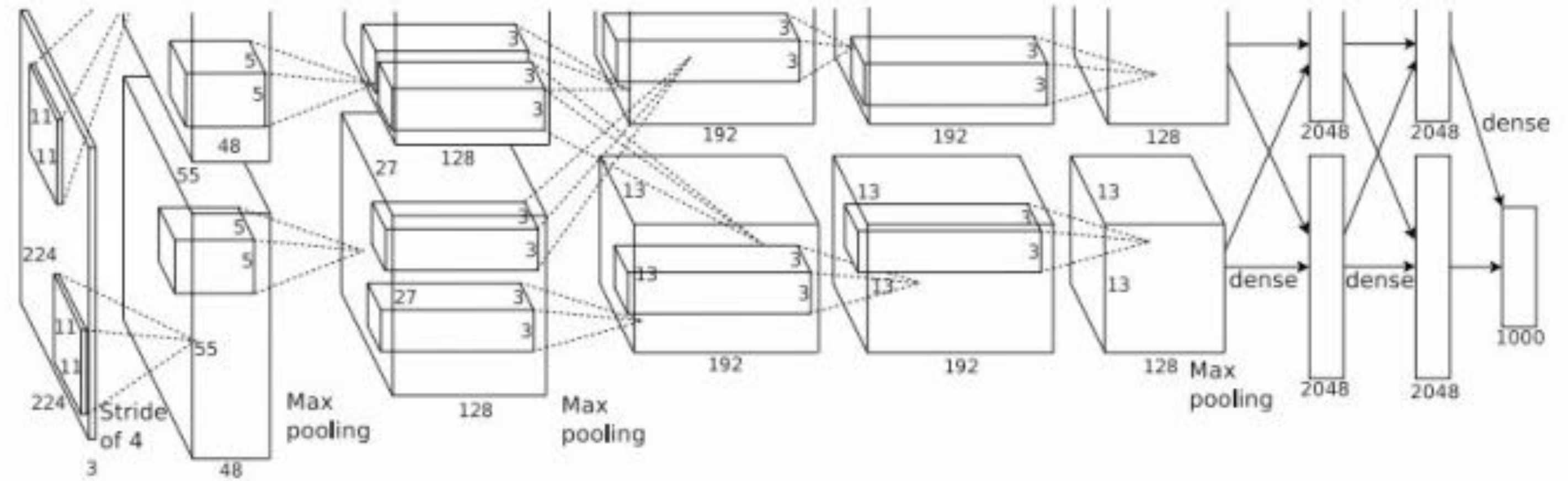
Output volume: 27x27x96

Q: what is the number of parameters in this layer?

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

# Case Study: AlexNet

*[Krizhevsky et al. 2012]*



Input: 227x227x3 images

After CONV1: 55x55x96

**Second layer (POOL1):** 3x3 filters applied at stride 2

Output volume: 27x27x96

Parameters: 0!

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



# Case Study: AlexNet

[Krizhevsky et al. 2012]

Input: 227x227x3 images

After CONV1: 55x55x96

After POOL1: 27x27x96

• • •

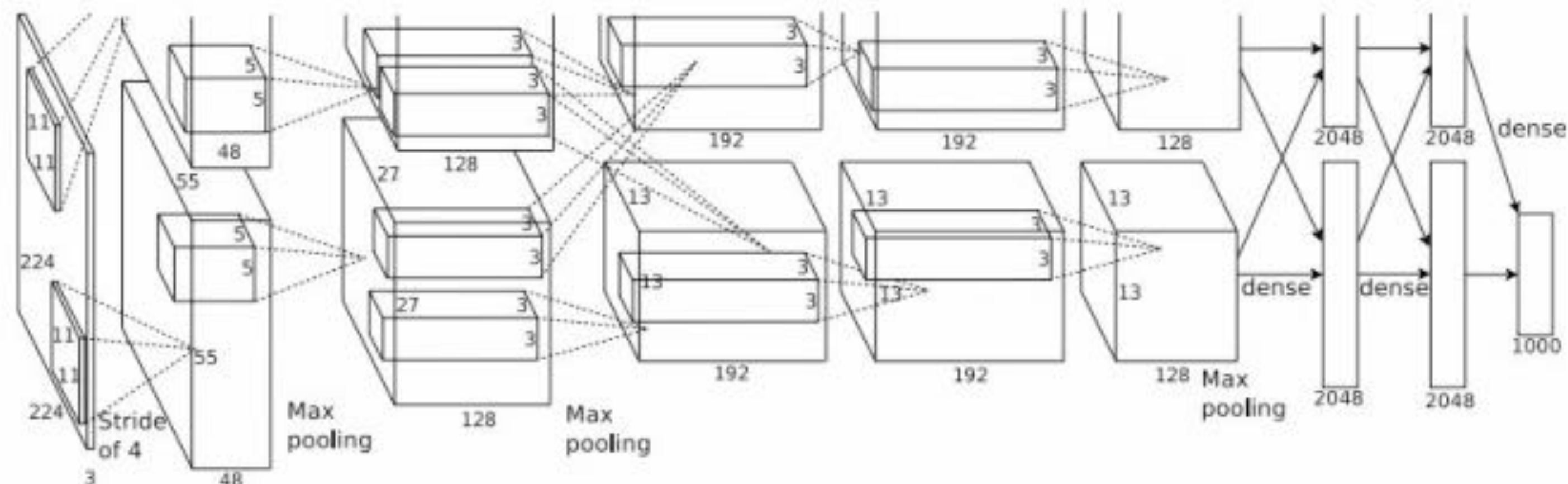


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

# Case Study: AlexNet

*[Krizhevsky et al. 2012]*

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)

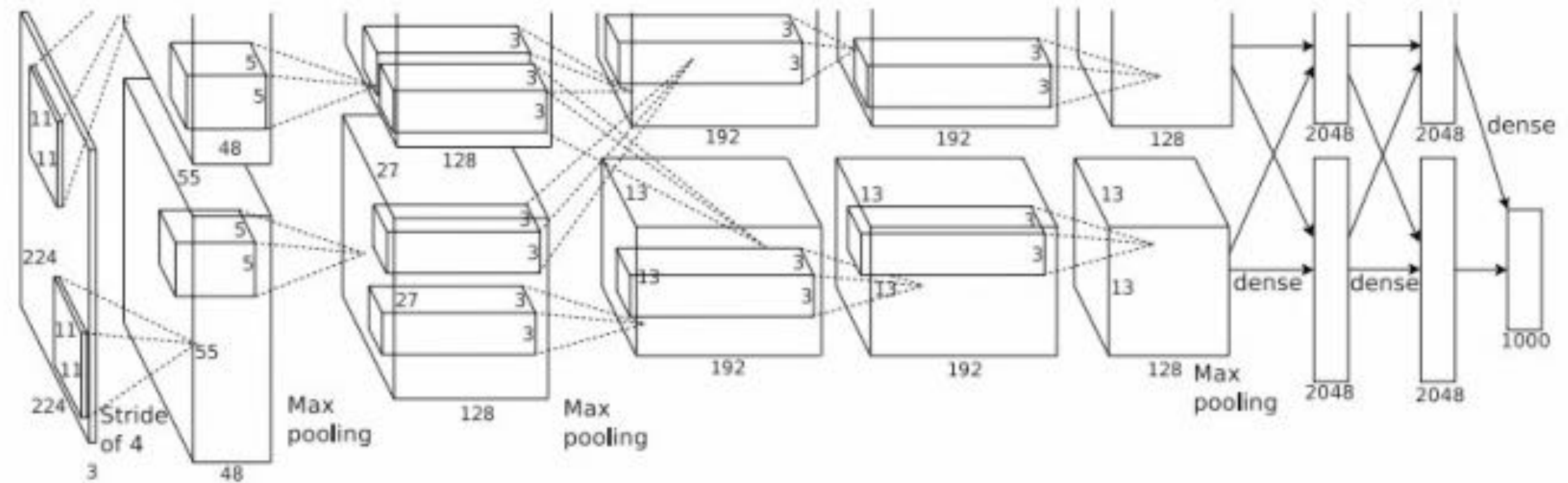


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

# Case Study: AlexNet

*[Krizhevsky et al. 2012]*

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

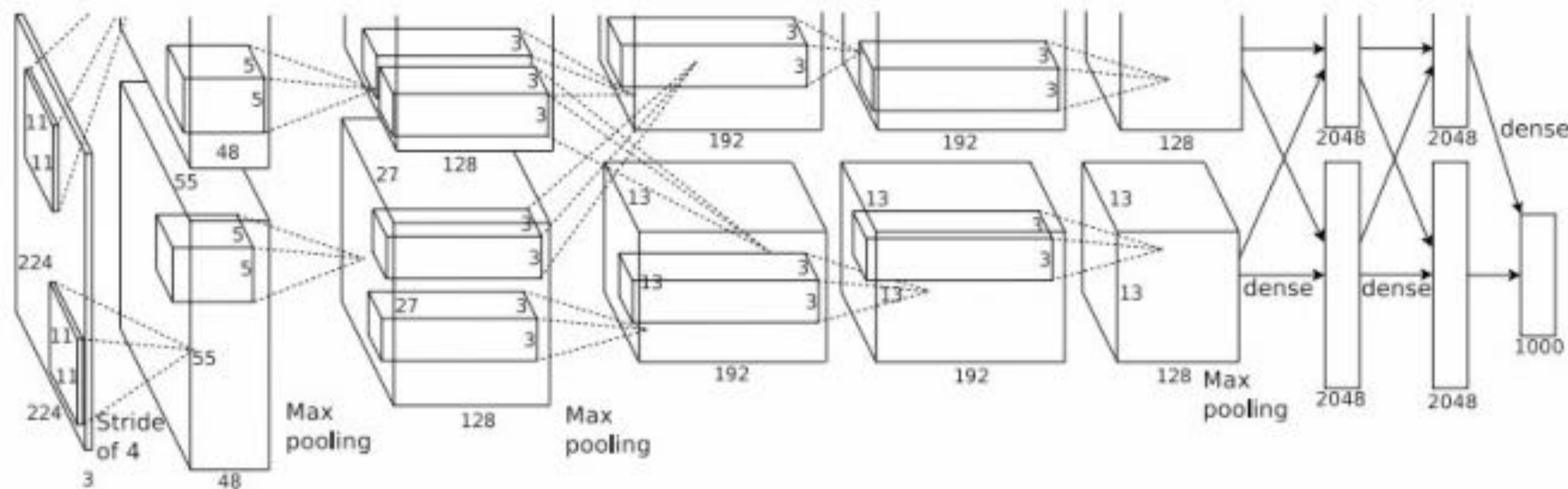
[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



## Details/Retrospectives:

- first use of ReLU
- used Local Response Normalization layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4

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



# Case Study: AlexNet

*[Krizhevsky et al. 2012]*

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

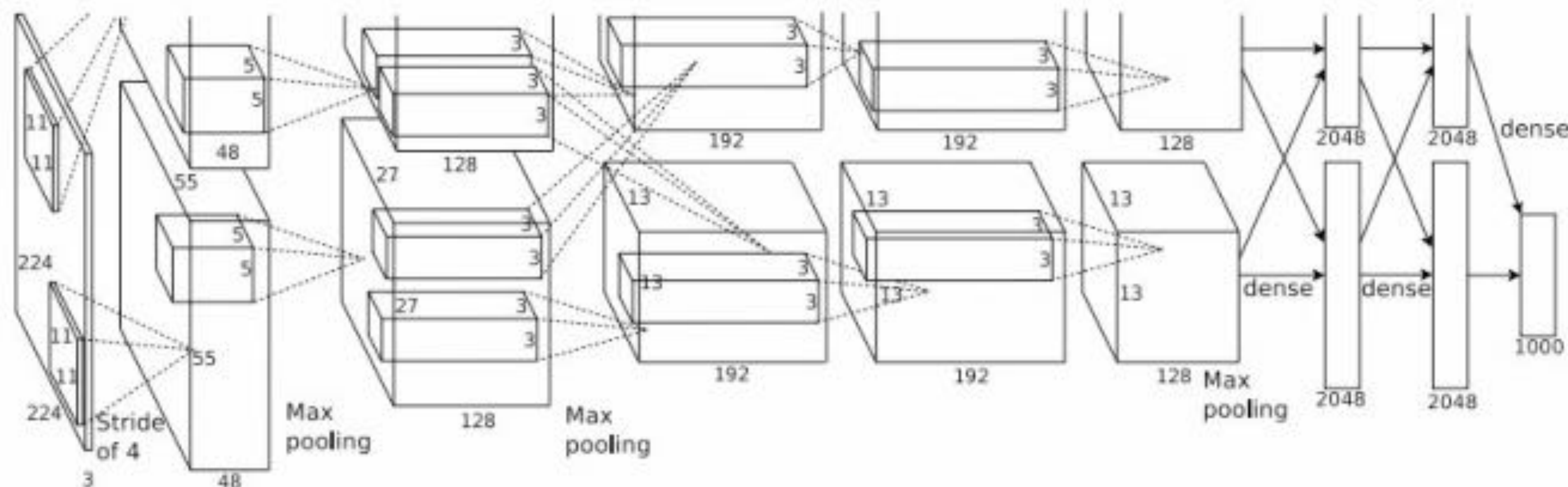
[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



## Details/Retrospectives:

- **first use of ReLU**
- used Local Response Normalization layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4

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



# Case Study: AlexNet

*[Krizhevsky et al. 2012]*

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

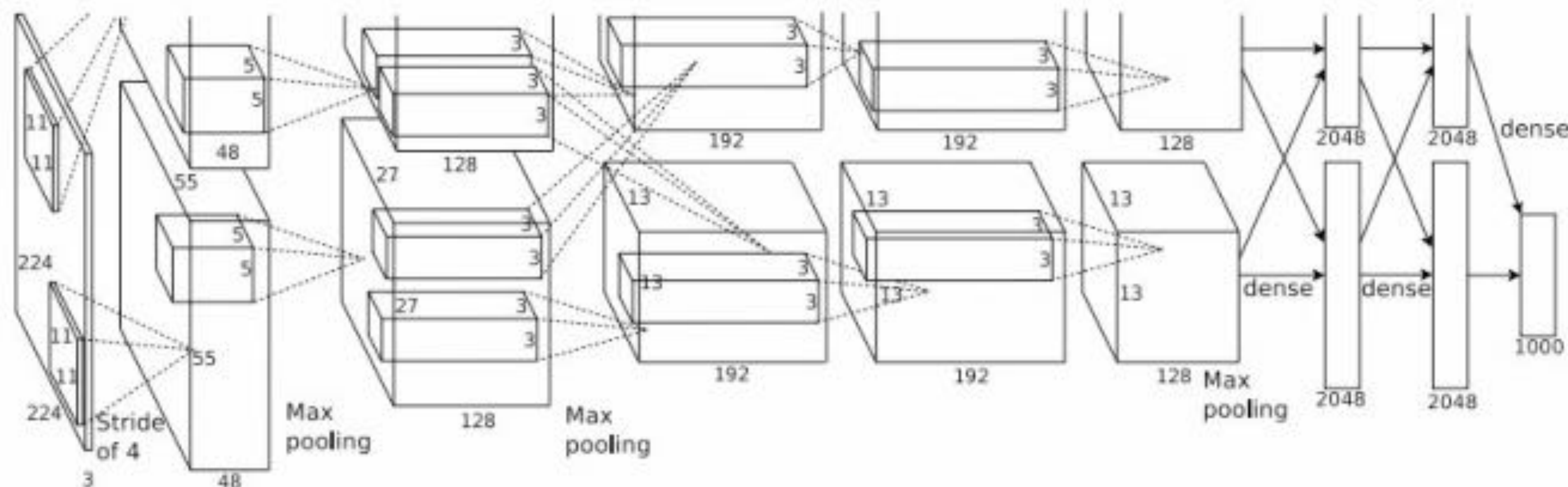
[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



## Details/Retrospectives:

- first use of ReLU
- **used Local Response Normalization layers (not common anymore)**
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4

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

# Case Study: AlexNet

*[Krizhevsky et al. 2012]*

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

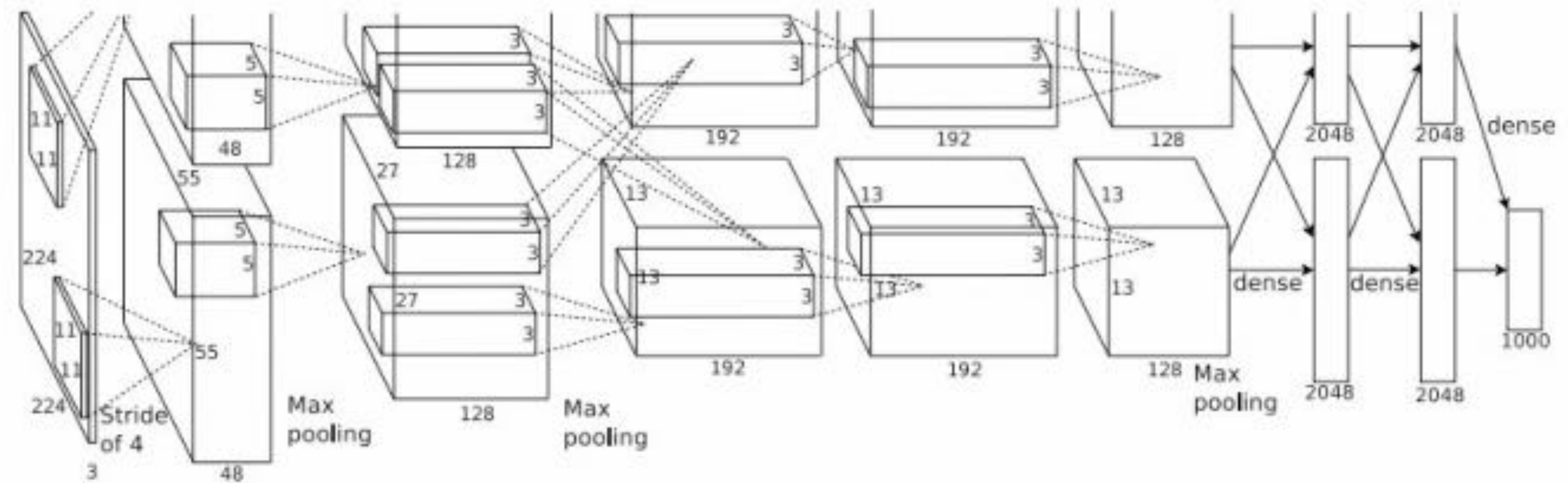
[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



## Details/Retrospectives:

- first use of ReLU
- used Local Response Normalization layers (not common anymore)
- **heavy data augmentation**
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4

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



# Regularization: Add term to loss

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$$

In common use:

L2 regularization

L1 regularization

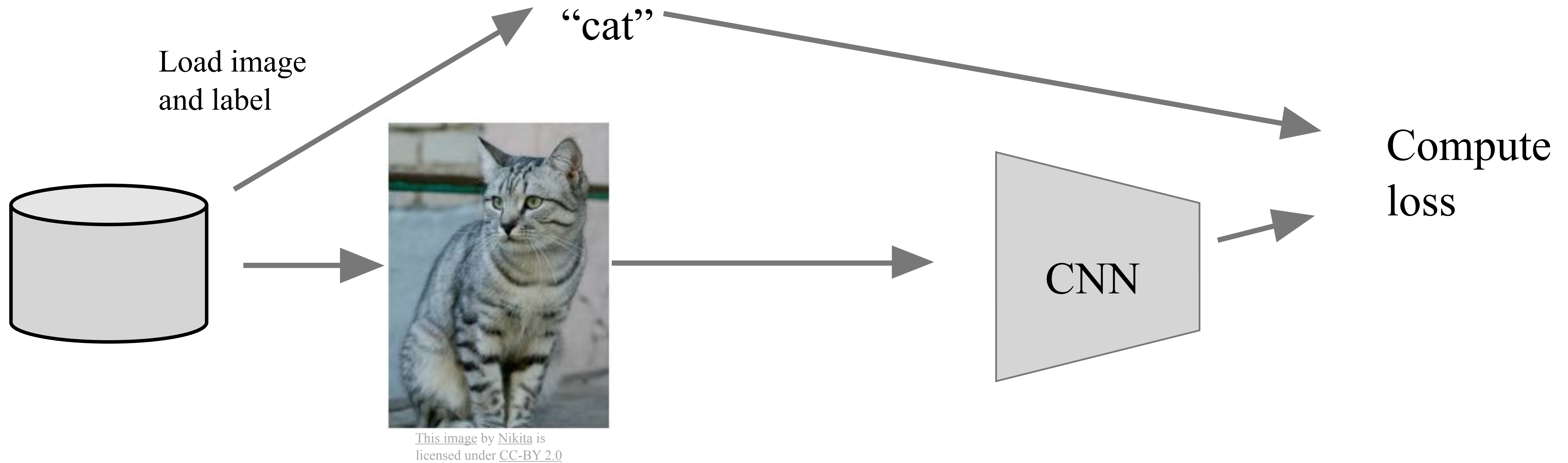
Elastic net (L1 + L2)

$$R(W) = \sum_k \sum_l W_{k,l}^2 \quad (\text{Weight decay})$$

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

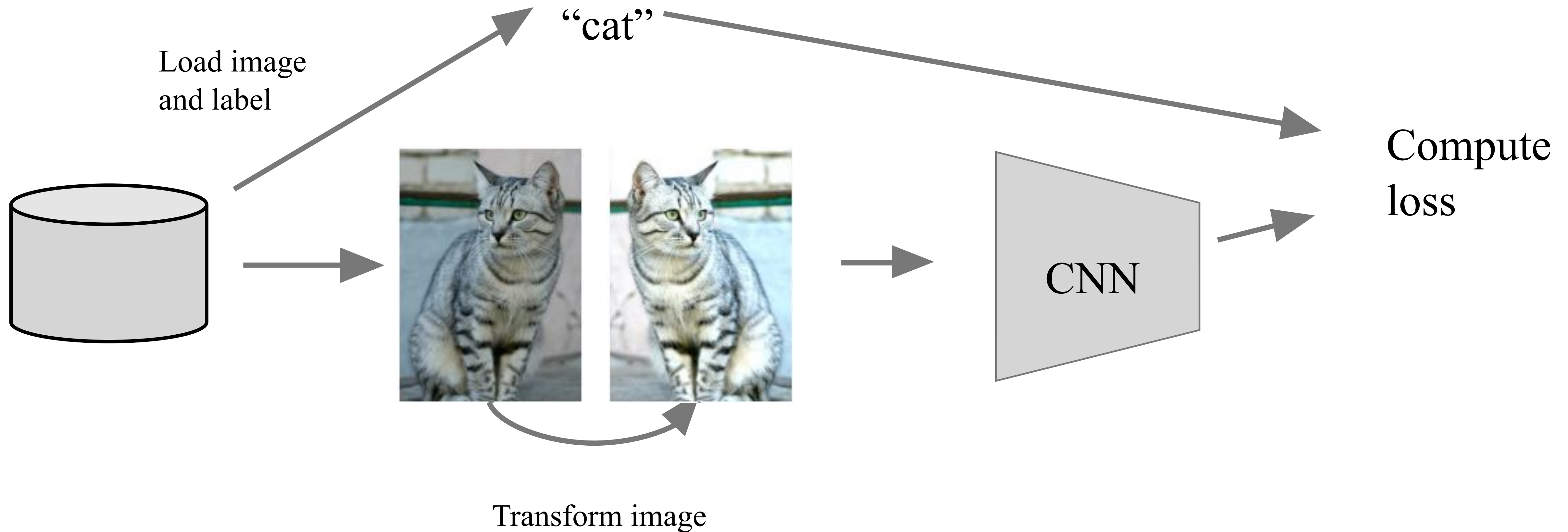
$$R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

# Regularization: Data Augmentation





# Regularization: Data Augmentation



# Data Augmentation

## Horizontal Flips



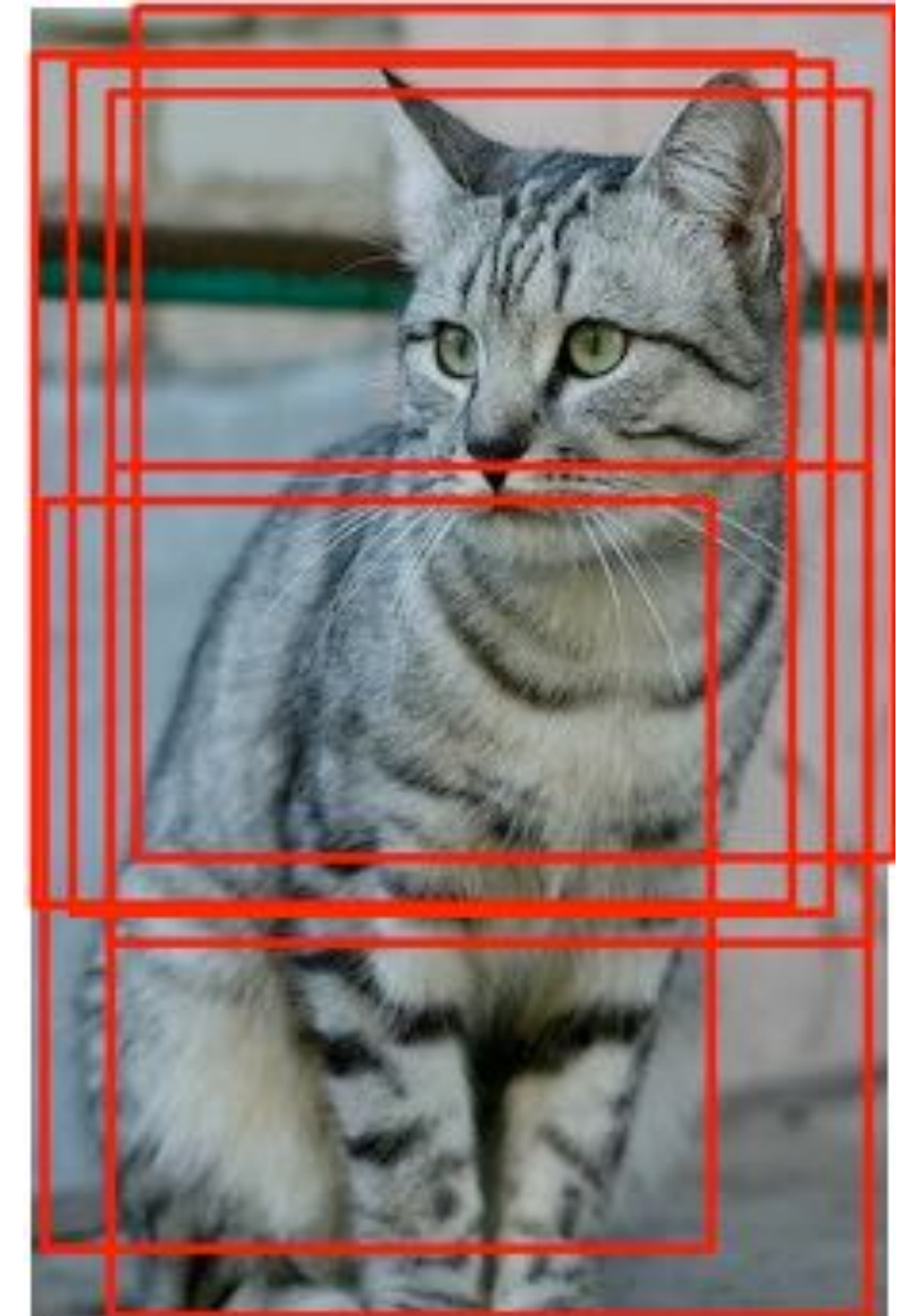


# Data Augmentation

## Random crops and scales

**Training** : sample random crops

1. Resize training image, short side = 256
2. Sample random 224 x 224 patch



# Data Augmentation

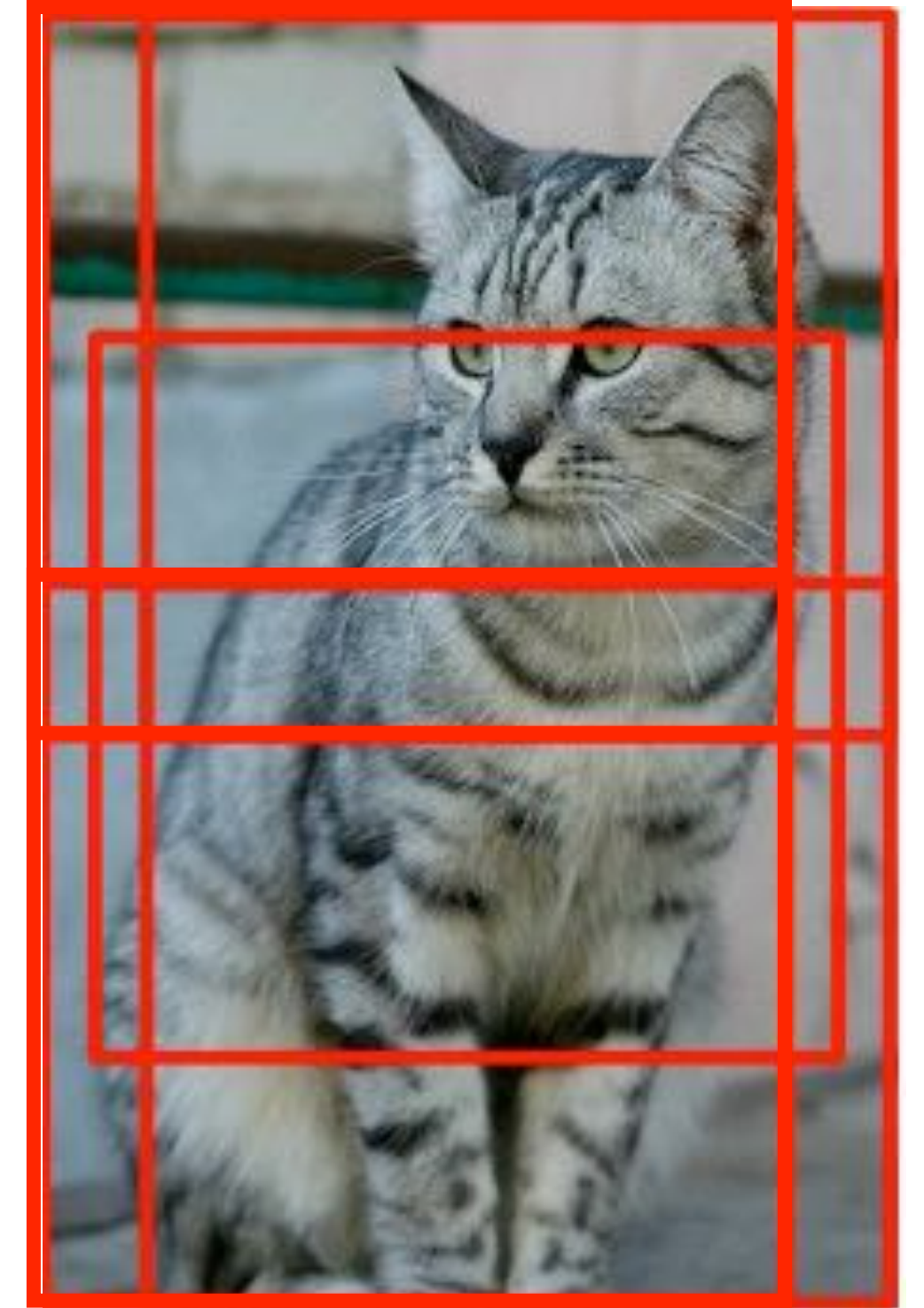
## Random crops and scales

### **Training** : sample random crops

1. Resize training image, short side = 256
2. Sample random 224 x 224 patch

### **Testing** : average a fixed set of crops

1. Extract 5 224x224 patches (corners and center) and horizontal reflections (10 patches total)
2. Average softmax of 10 patches and predict

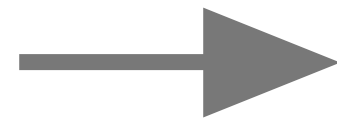




# Data Augmentation

## Color Jitter

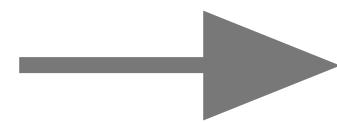
Simple: Randomize  
contrast and brightness



# Data Augmentation

## Color Jitter

Simple: Randomize  
contrast and brightness



## More Complex :

1. Apply PCA to all [R, G, B] pixels in training set
2. Sample a “color offset” along principal component directions
3. Add offset to all pixels of a training image

(As seen in *[Krizhevsky et al. 2012]*, ResNet, etc)

# Data Augmentation

Get creative for your problem!

Random mix/combinations of :

- translation
- rotation
- stretching
- shearing,
- lens distortions, ... (go crazy)



# Case Study: AlexNet

*[Krizhevsky et al. 2012]*

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

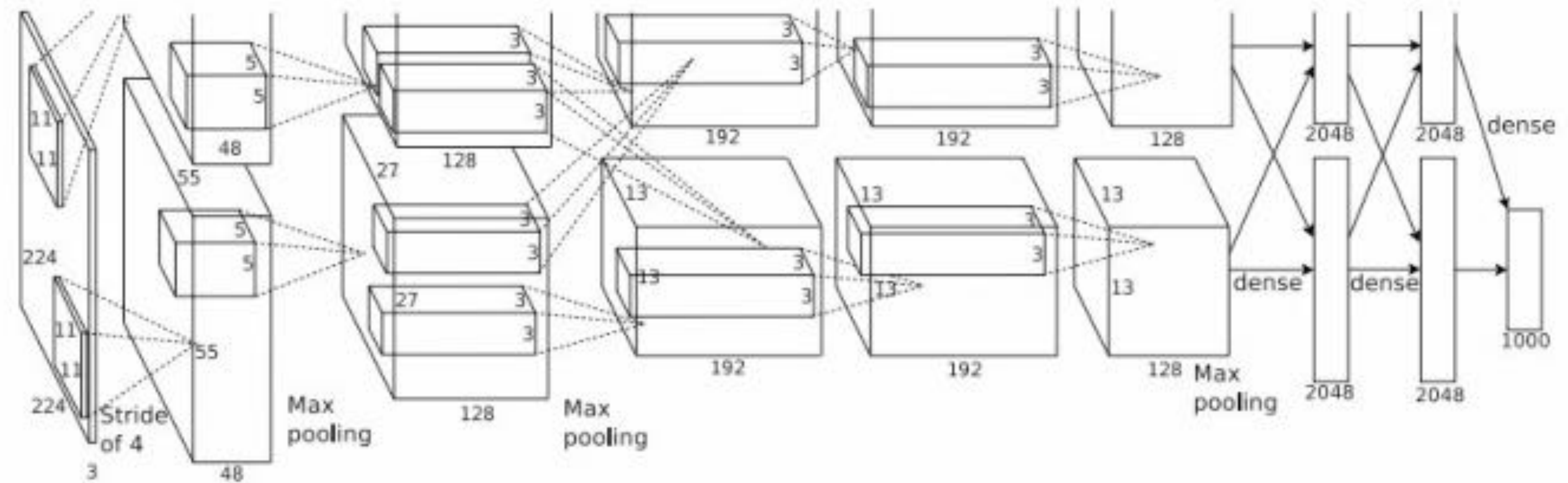
[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



## Details/Retrospectives:

- first use of ReLU
- used Local Response Normalization layers (not common anymore)
- heavy data augmentation
- **dropout 0.5**
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4

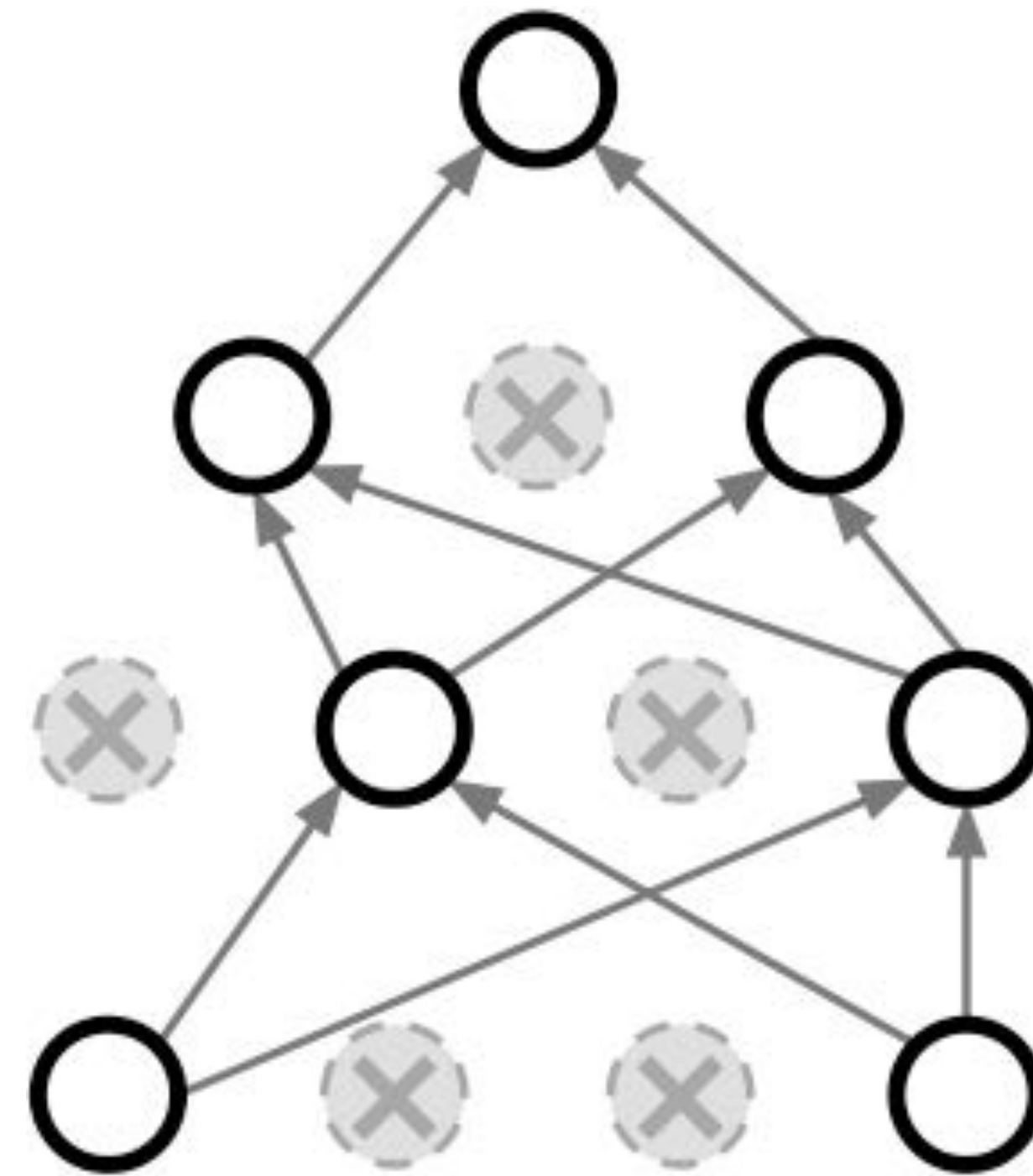
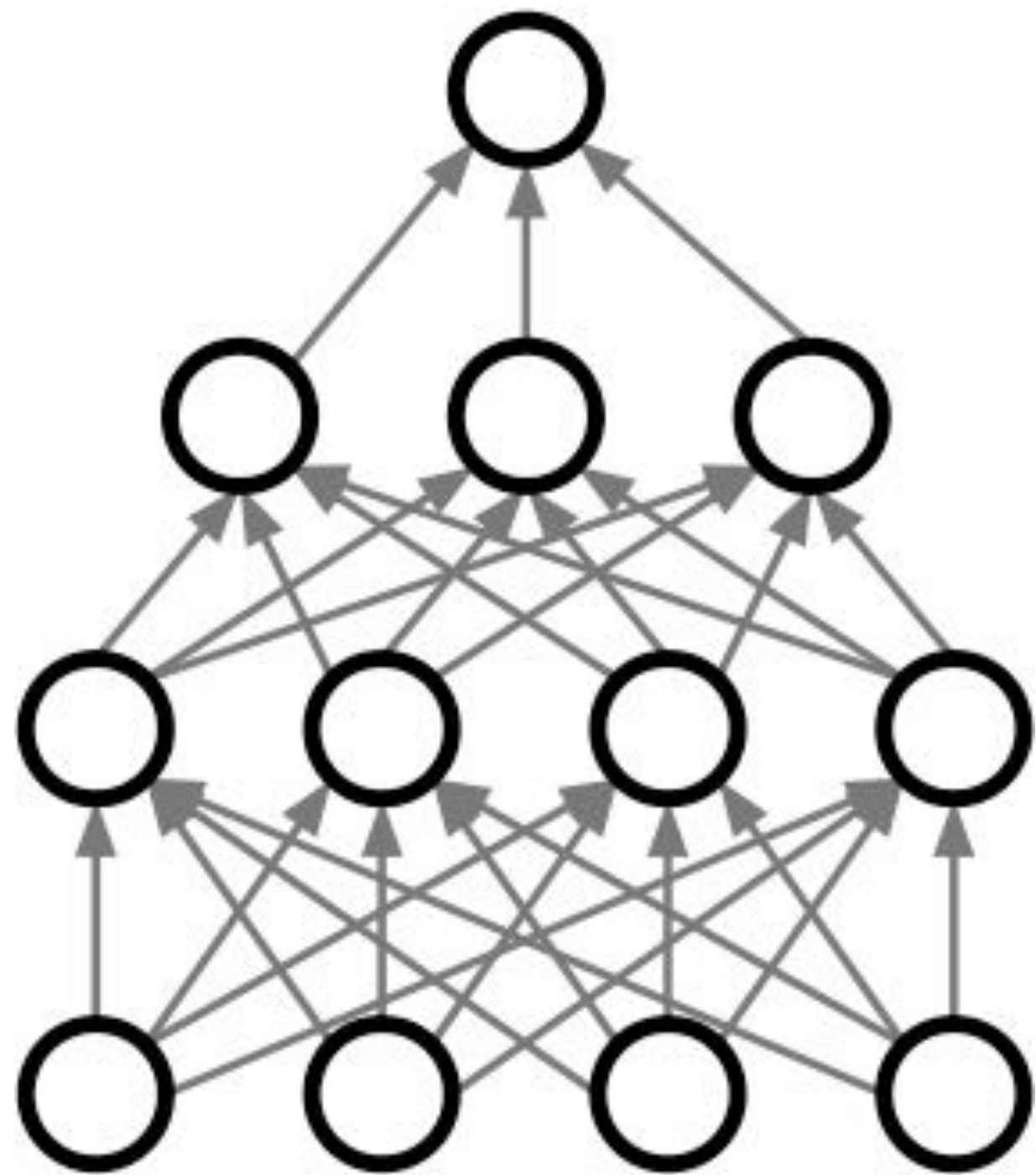
Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.



# Regularization: Dropout

In each forward pass, randomly set some neurons to zero

Probability of dropping is a hyperparameter; 0.5 is common



Srivastava et al, "Dropout: A simple way to prevent neural networks from overfitting", JMLR 2014

# Regularization: Dropout

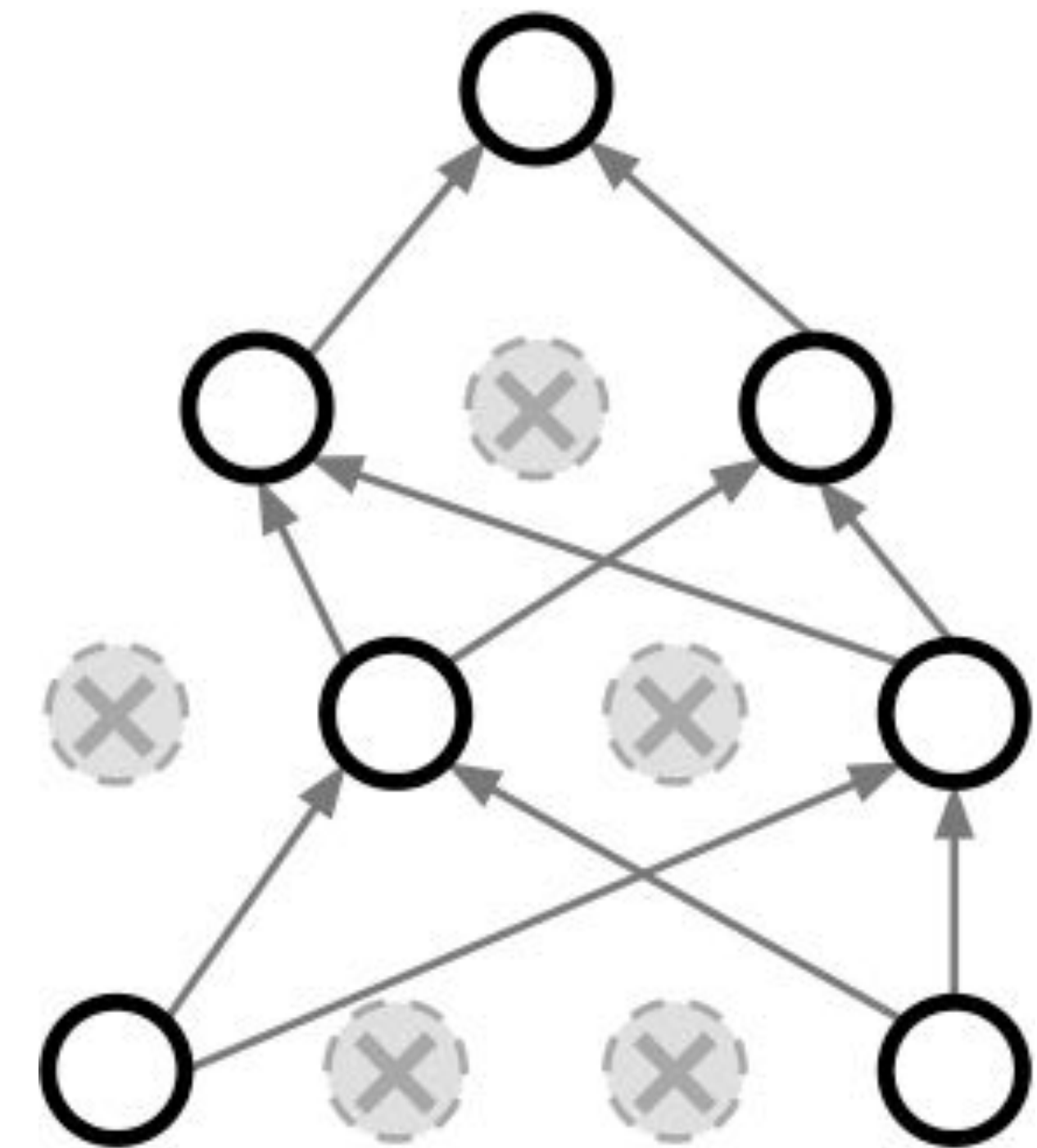
Example forward pass with a 3-layer network using dropout

```
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    """ X contains the data """

    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

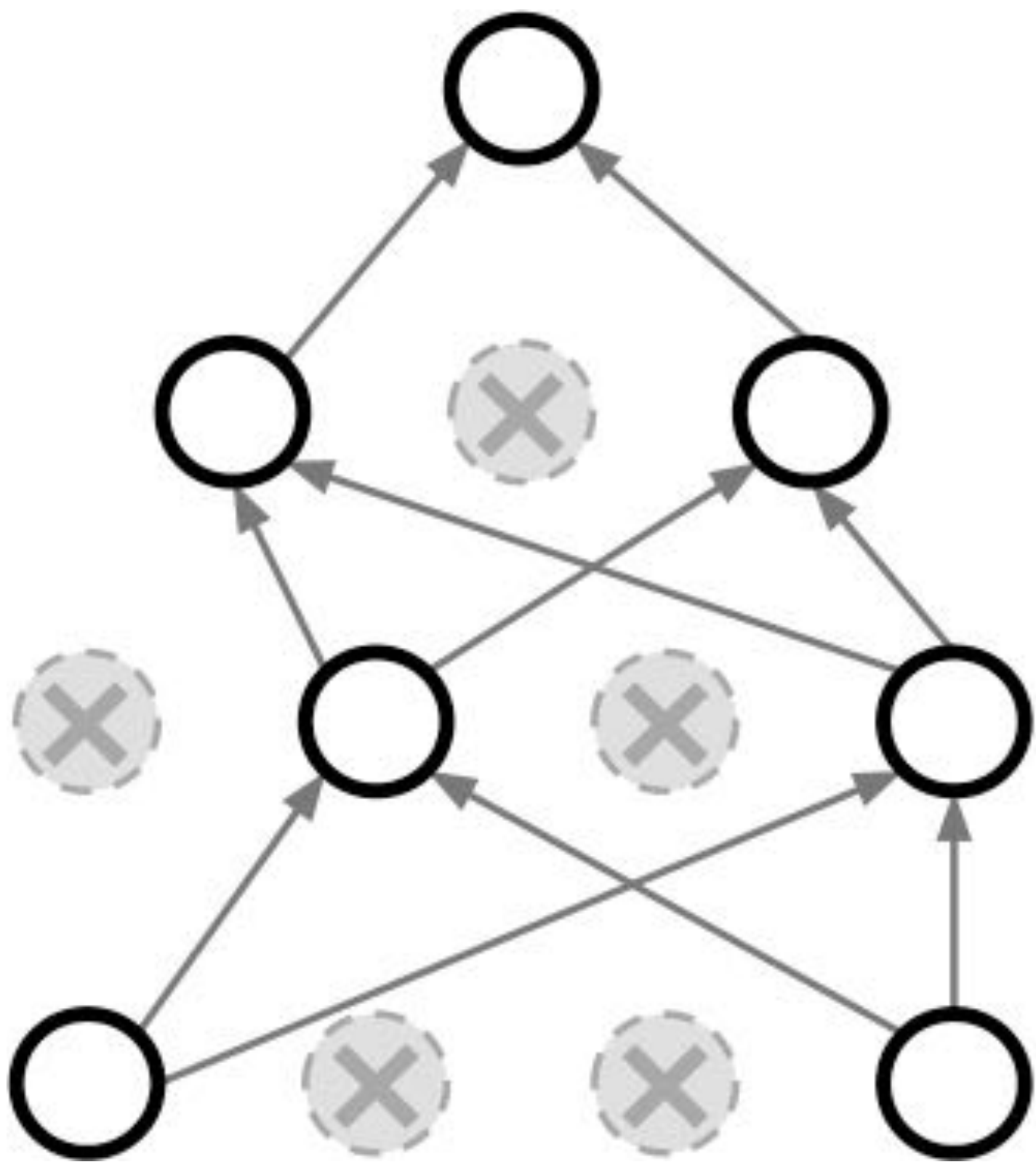
    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)
```



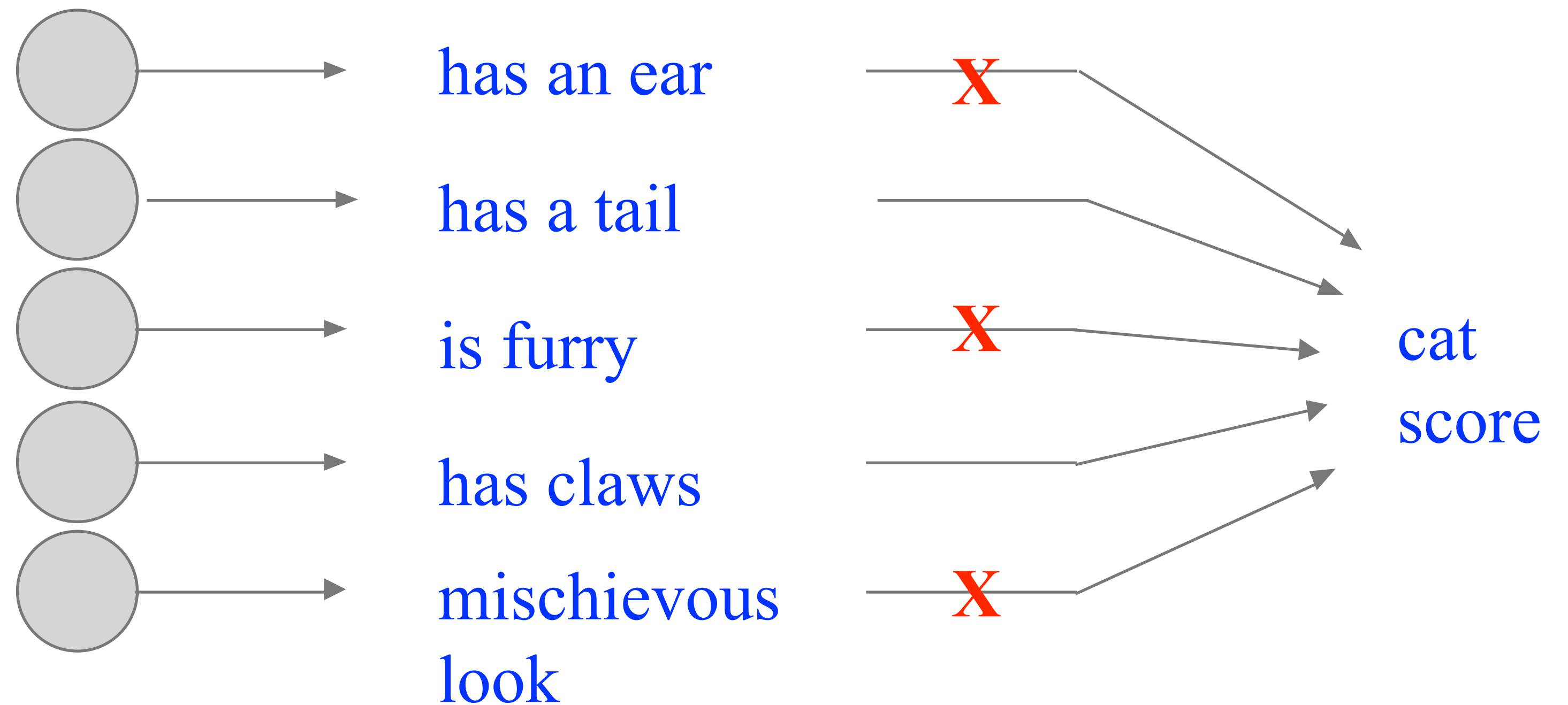


# Regularization: Dropout

How can this possibly be a good idea?



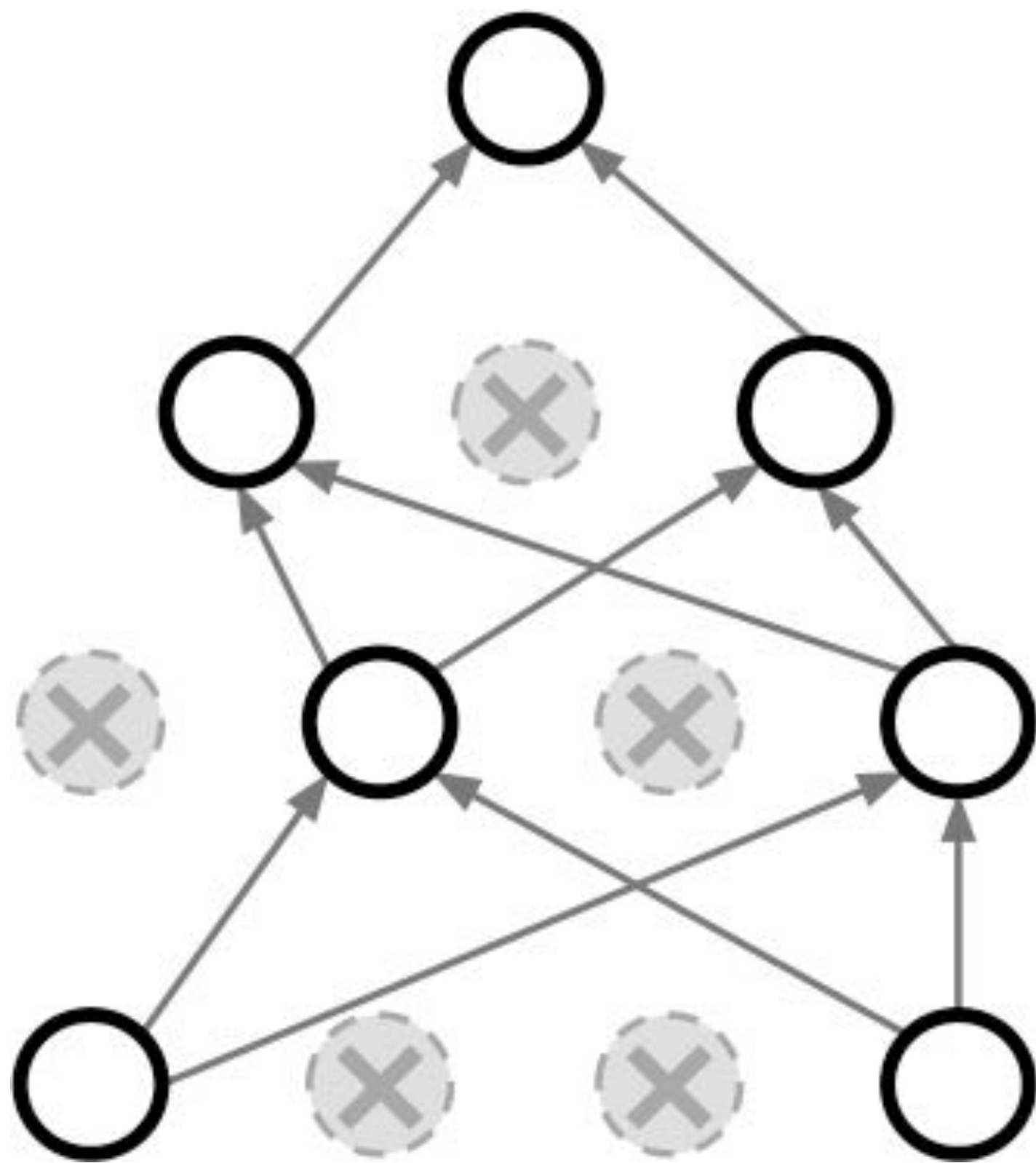
Forces the network to have a redundant representation;  
Prevents co-adaptation of features





# Regularization: Dropout

How can this possibly be a good idea?



Another interpretation:

Dropout is training a large **ensemble** of models (that share parameters).

Each binary mask is one model

An FC layer with 4096 units has  $2^{4096} \sim 10^{1233}$  possible masks!

Only  $\sim 10^{82}$  atoms in the universe...

# Dropout: Test time

Dropout makes our output random!

$$\begin{array}{c} \text{Output} \\ \text{(label)} \end{array} \quad \begin{array}{c} \text{Input} \\ \text{(image)} \end{array} \quad \begin{array}{c} \text{Random} \\ \text{mask} \end{array}$$
$$\boxed{y} = f_W(\boxed{x}, \boxed{z})$$

Want to “average out” the randomness at test-time

$$y = f(x) = E_z[f(x, z)] = \int p(z) f(x, z) dz$$

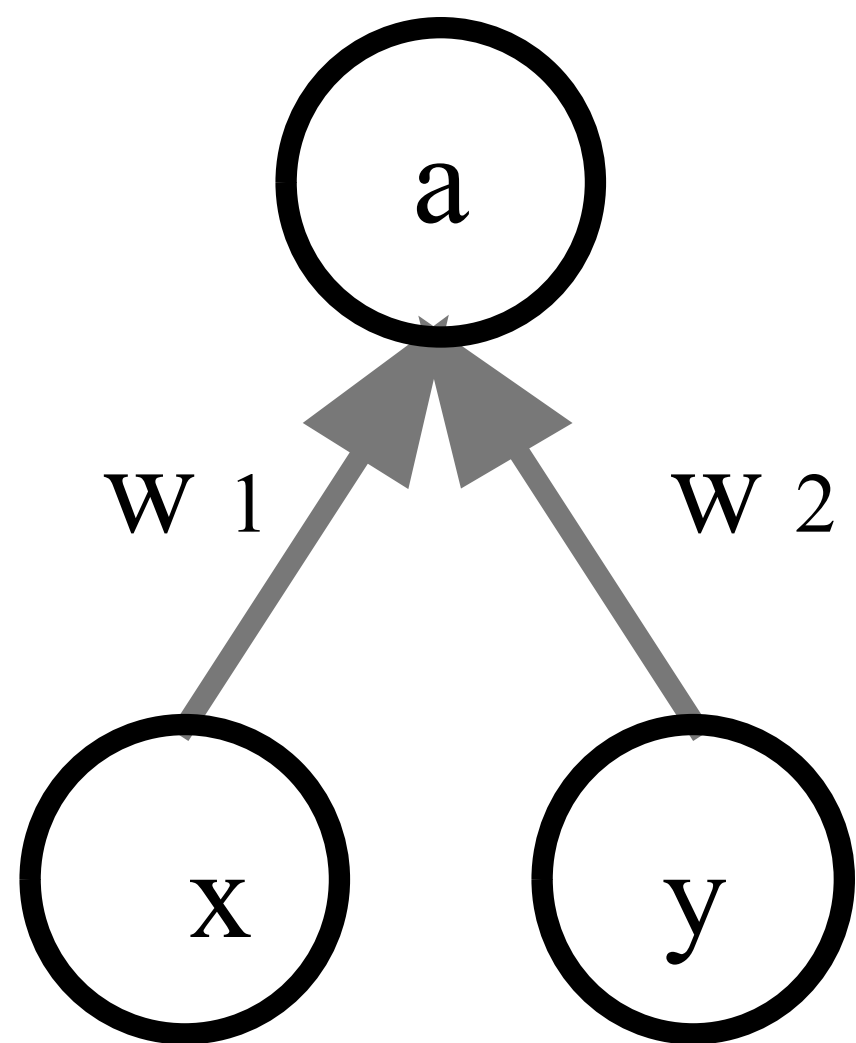
But this integral seems hard ...

# Dropout: Test time

Want to approximate the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.





# Dropout: Test time

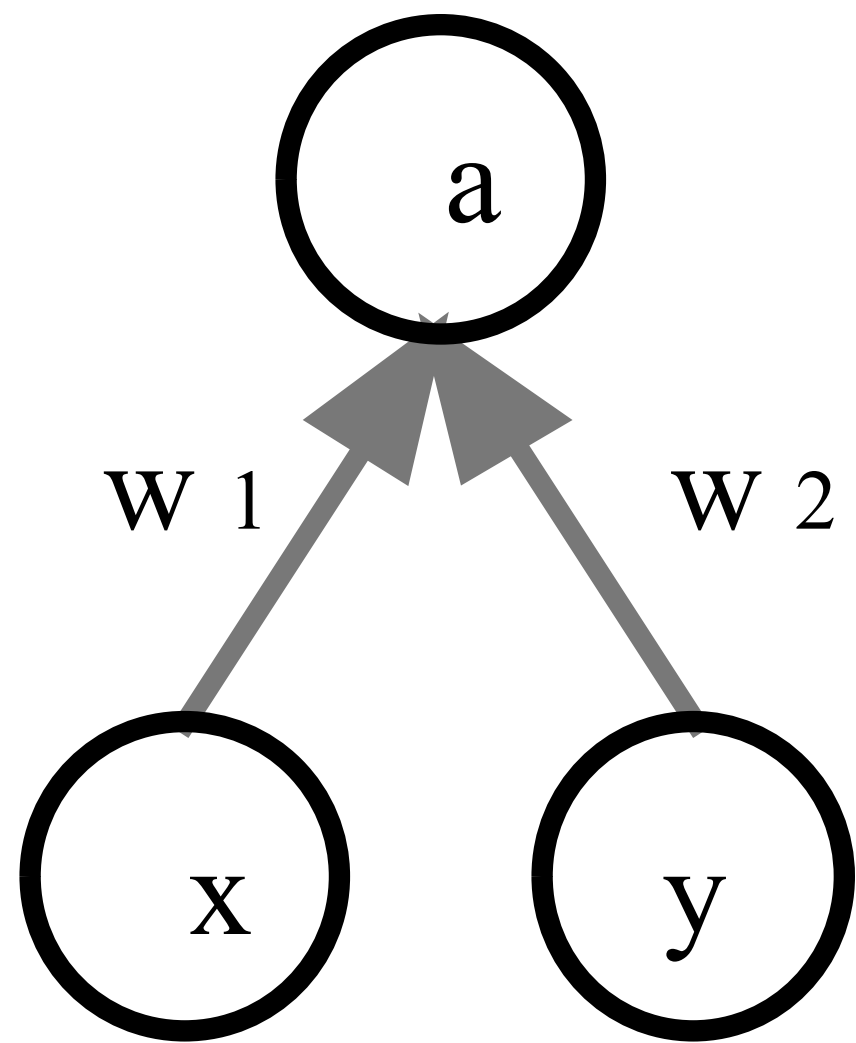
Want to approximate the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.

At test time we have:

$$E[a] = w_1x + w_2y$$



# Dropout: Test time

Want to approximate the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

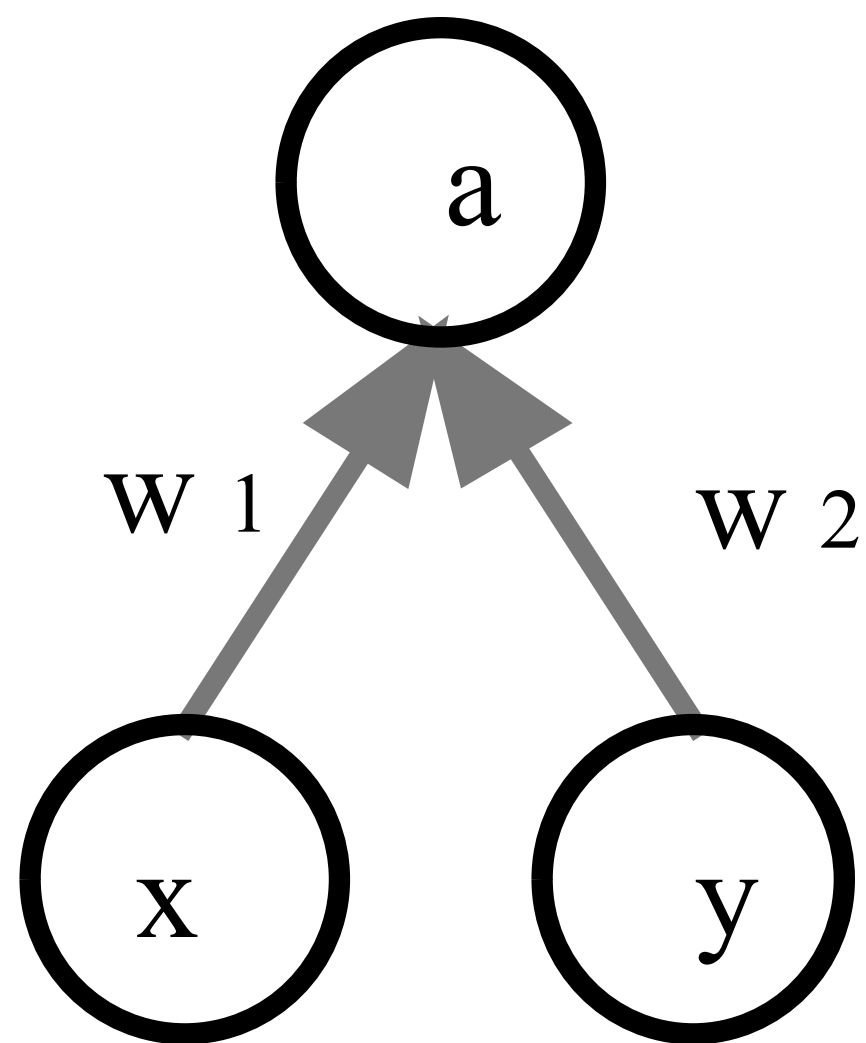
Consider a single neuron.

At test time we have:

$$E[a] = w_1x + w_2y$$

During training we have:

$$\begin{aligned} E[a] &= \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) \\ &\quad + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) \\ &= \frac{1}{2}(w_1x + w_2y) \end{aligned}$$



# Dropout: Test time

Want to approximate the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

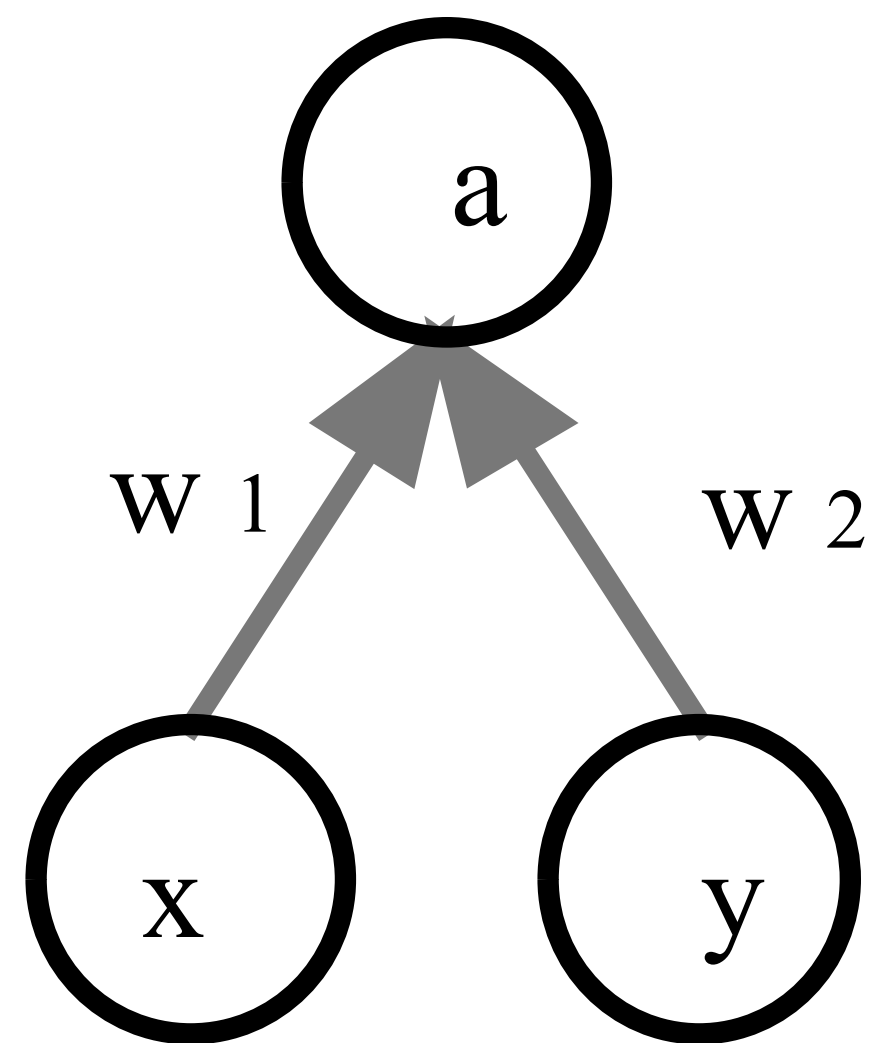
Consider a single neuron.

At test time we have:

$$E[a] = w_1x + w_2y$$

During training we have:

$$\begin{aligned} E[a] &= \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) \\ &\quad + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) \\ &= \frac{1}{2}(w_1x + w_2y) \end{aligned}$$



At test time, **multiply**  
by dropout probability



# Dropout: Test time

```
def predict(X):  
    # ensembled forward pass  
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations  
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations  
    out = np.dot(W3, H2) + b3
```

At test time all neurons are active always

=> We must scale the activations so that for each neuron:

output at test time = expected output at training time

# Dropout Summary

```
""" Vanilla Dropout: Not recommended implementation (see notes below) """

p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    """ X contains the data """

    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)

def predict(X):
    # ensembled forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations
    out = np.dot(W3, H2) + b3
```

drop in forward pass

scale at test time



# More common: “Inverted dropout”

```
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = (np.random.rand(*H1.shape) < p) / p # first dropout mask. Notice /p!
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = (np.random.rand(*H2.shape) < p) / p # second dropout mask. Notice /p!
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)

def predict(X):
    # ensembled forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) # no scaling necessary
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    out = np.dot(W3, H2) + b3
```

test time is unchanged!





# Regularization: A common pattern

**Training :** Add some kind of randomness

$$y = f_W(x, z)$$

**Testing:** Average out randomness (sometimes approximate)

$$y = f(x) = E_z[f(x, z)] = \int p(z) f(x, z) dz$$

# Case Study: AlexNet

*[Krizhevsky et al. 2012]*

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

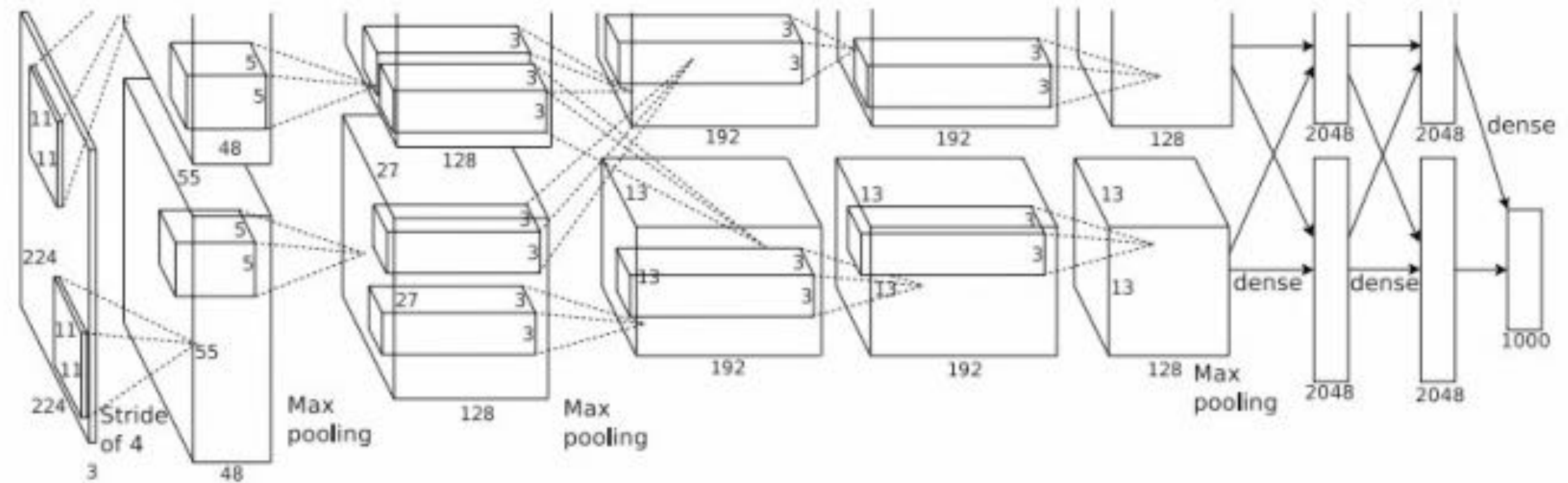
[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



## Details/Retrospectives:

- first use of ReLU
- used Local Response Normalization layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- **SGD Momentum 0.9**
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4

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

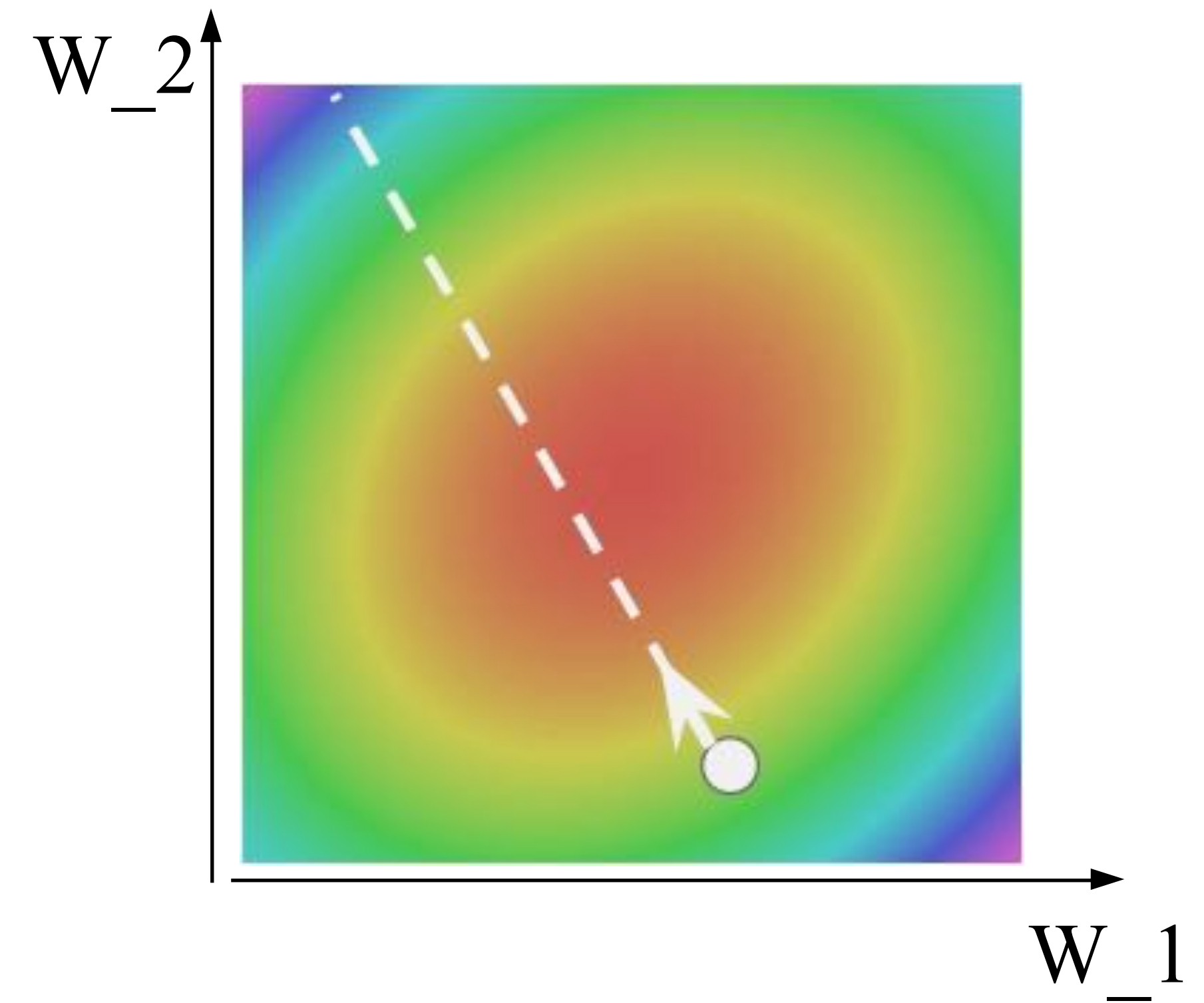
# Optimization

```
# Vanilla Gradient Descent
```

```
while True:
```

```
    weights_grad = evaluate_gradient(loss_fun, data, weights)
```

```
    weights += - step_size * weights_grad # perform parameter update
```

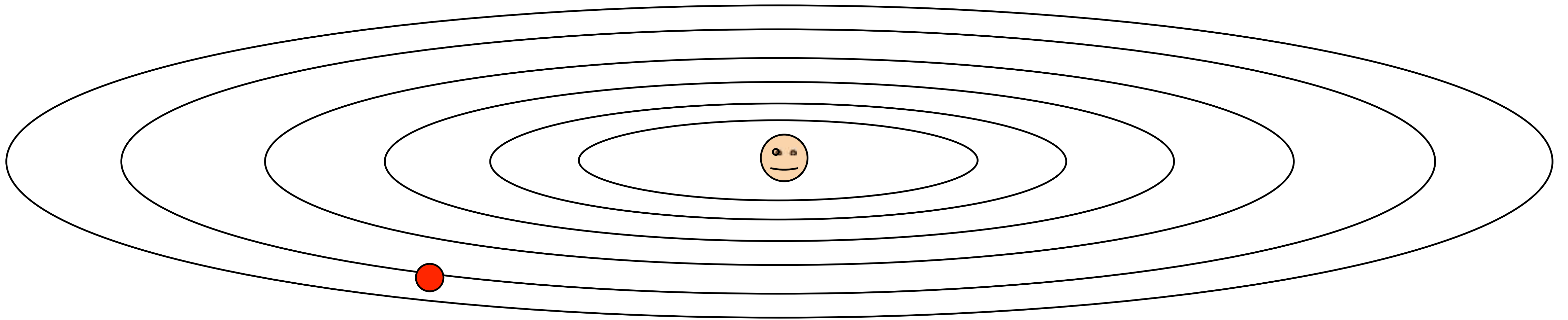




# Optimization: Problems with SGD

What if loss changes quickly in one direction and slowly in another?

What does gradient descent do?



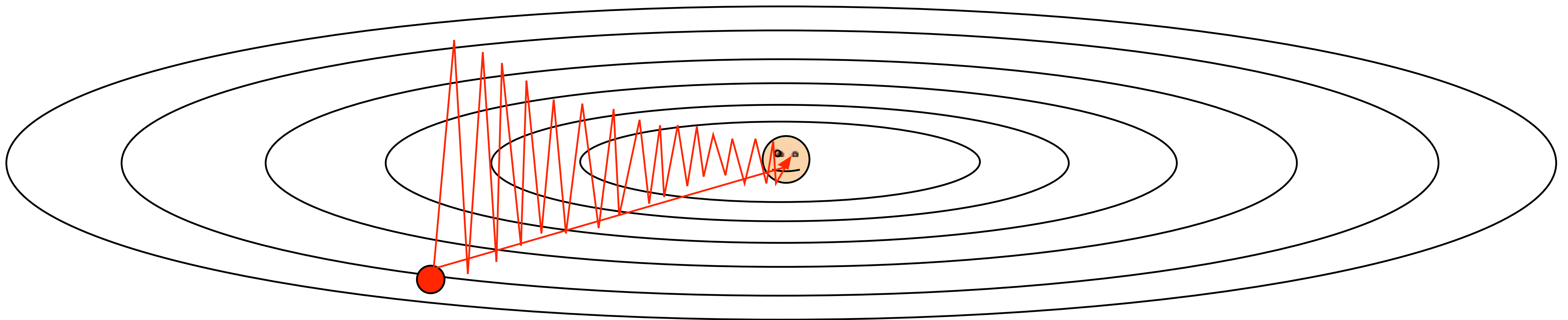
Loss function has high **condition number** : ratio of largest to smallest singular value of the Hessian matrix is large

# Optimization: Problems with SGD

What if loss changes quickly in one direction and slowly in another?

What does gradient descent do?

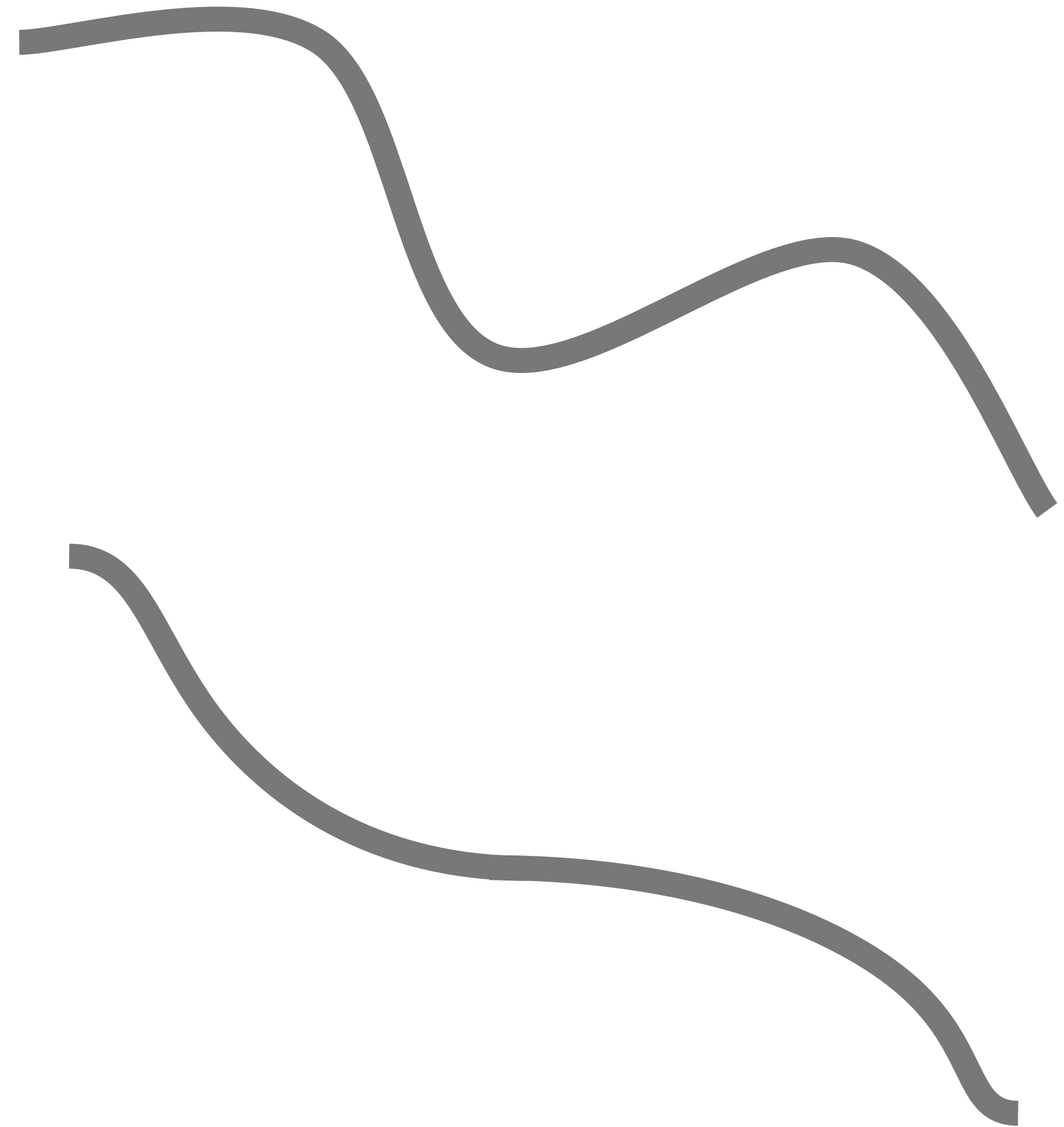
Very slow progress along shallow dimension, jitter along steep direction



Loss function has high **condition number** : ratio of largest to smallest singular value of the Hessian matrix is large

# Optimization: Problems with SGD

What if the loss function has a **local minima** or **saddle point** ?

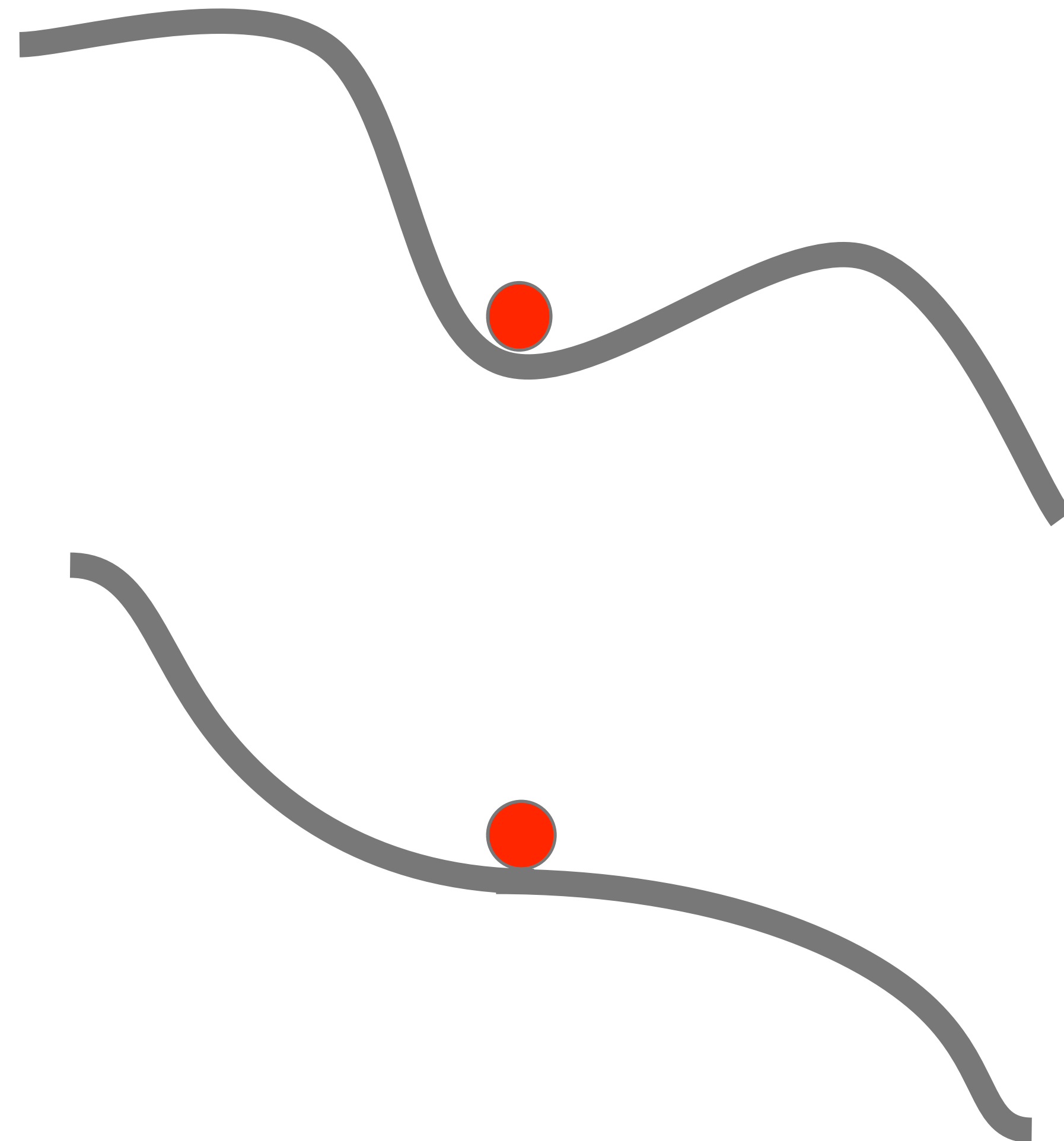




# Optimization: Problems with SGD

What if the loss function has a **local minima** or **saddle point** ?

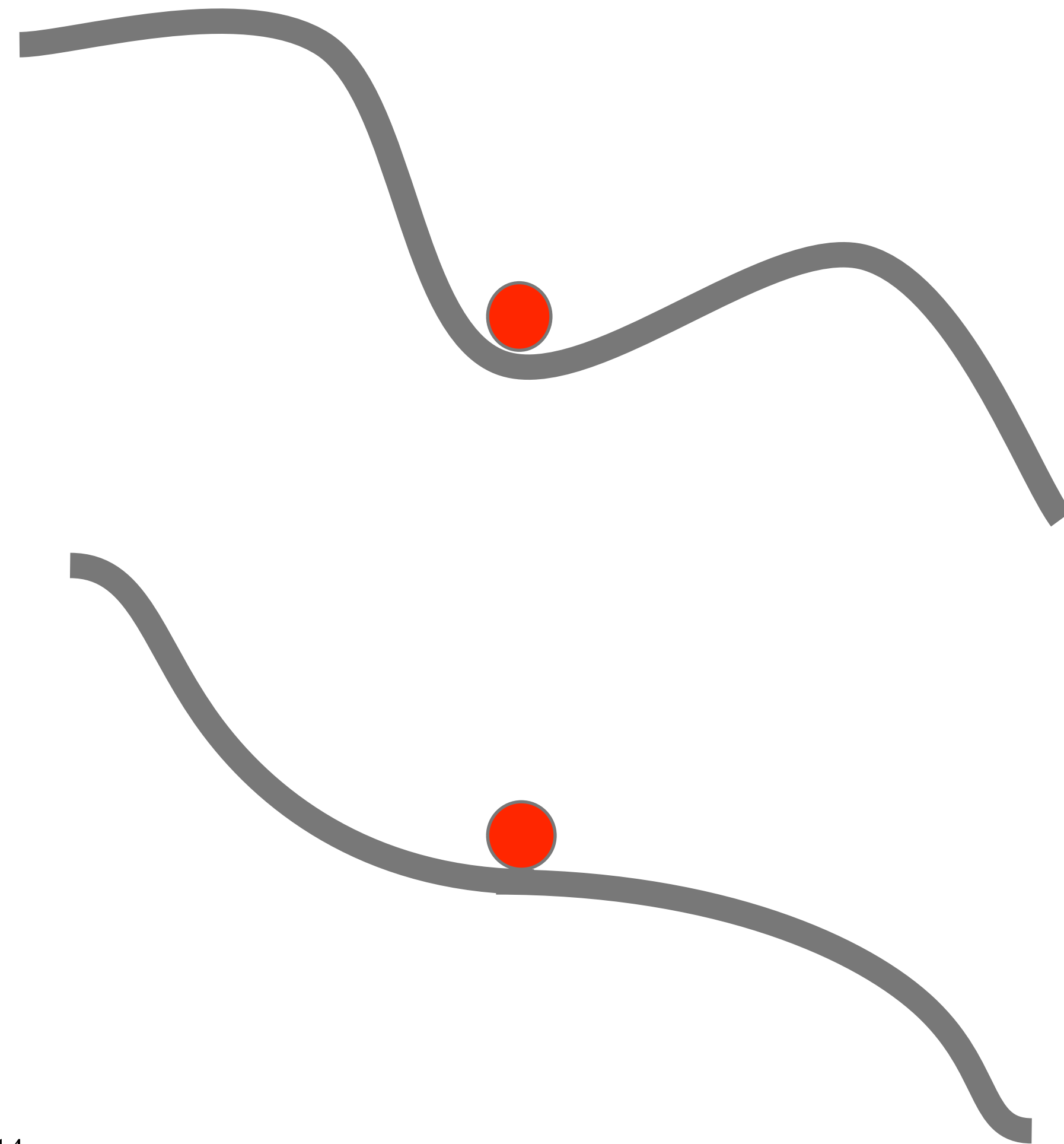
Zero gradient,  
gradient descent  
gets stuck



# Optimization: Problems with SGD

What if the loss function has a **local minima** or **saddle point** ?

Saddle points much more common in high dimension



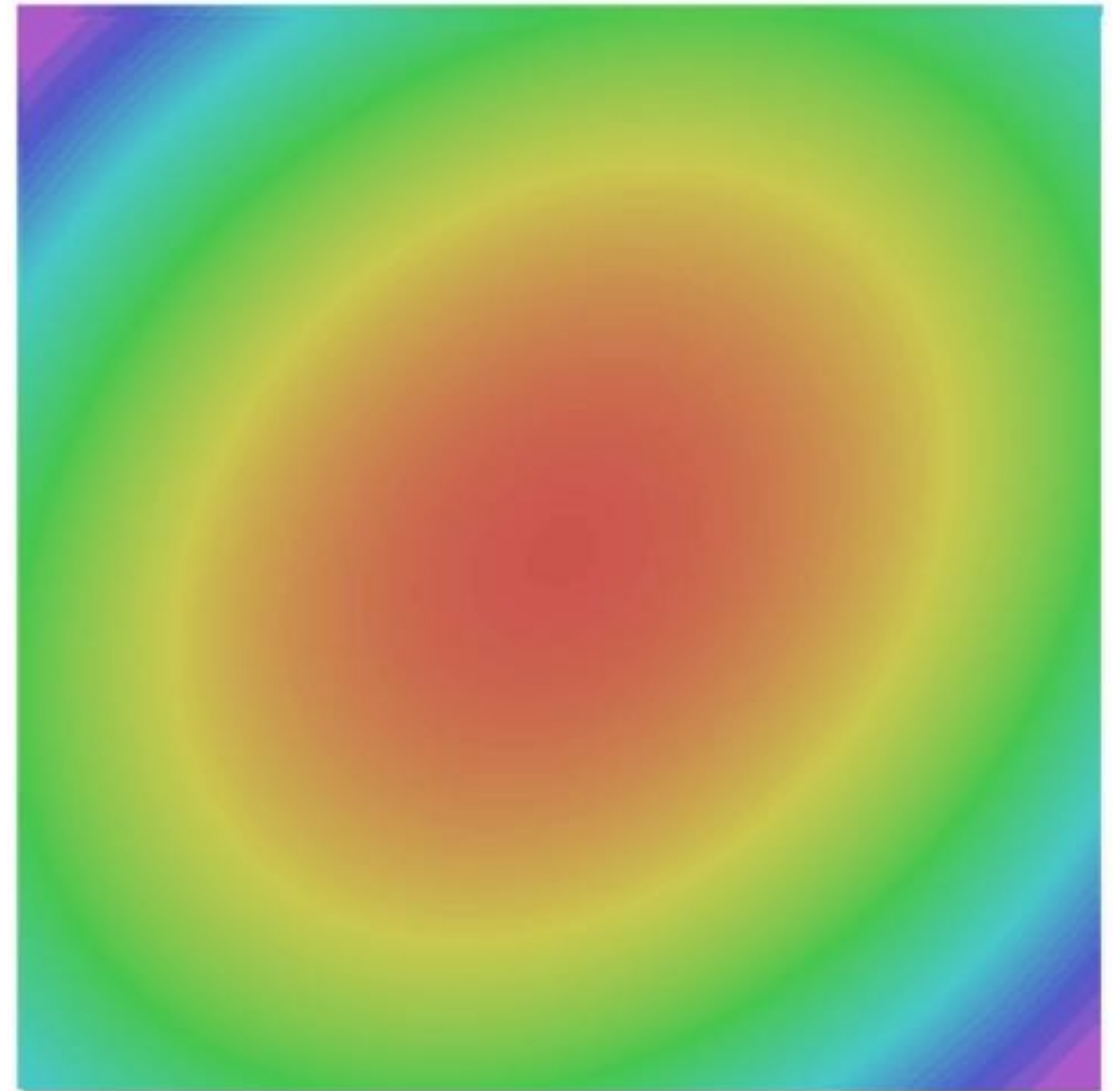
Dauphin et al, "Identifying and attacking the saddle point problem in high-dimensional non-convex optimization", NIPS 2014

# Optimization: Problems with SGD

Our gradients come from minibatches so they can be noisy!

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W)$$





# SGD + Momentum

## SGD

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

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

## SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1}$$

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

- Build up “velocity” as a running mean of gradients
- Rho gives “friction”; typically rho=0.9 or 0.99

Sutskever et al, “On the importance of initialization and momentum in deep learning”, ICML 2013

# SGD + Momentum

## SGD+Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t)$$

$$x_{t+1} = x_t + v_{t+1}$$

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

## SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1}$$

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

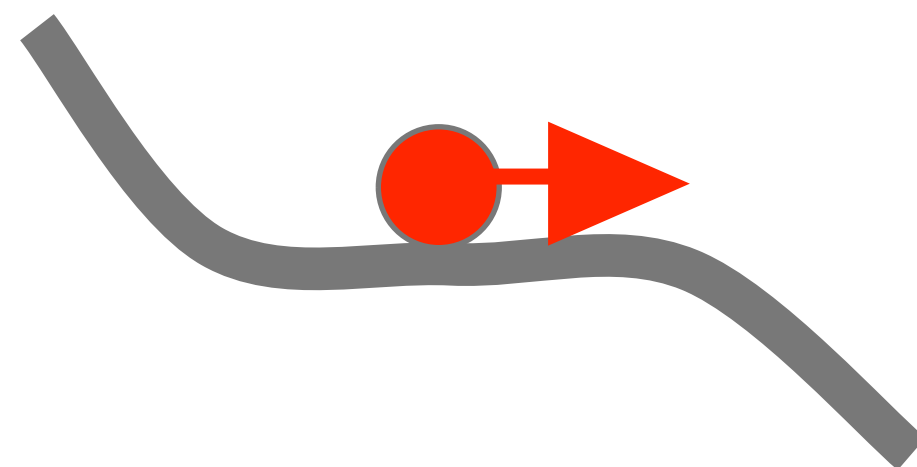
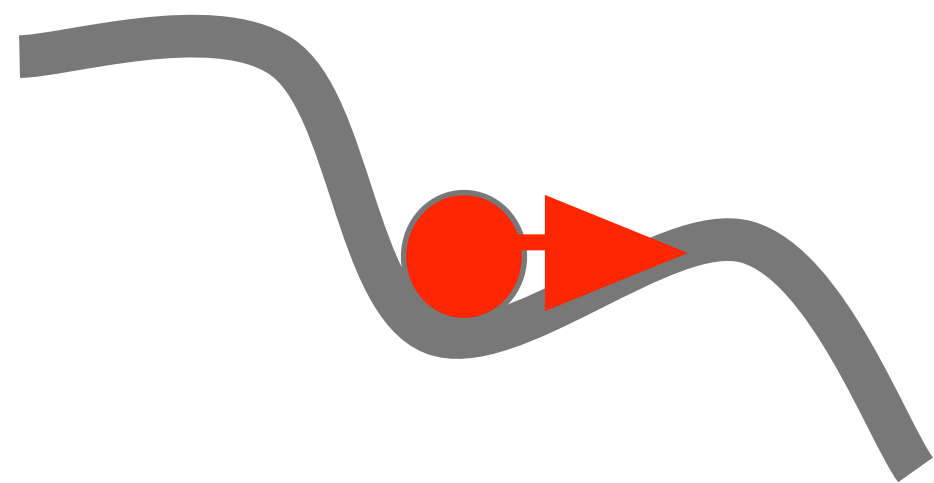
You may see SGD+Momentum formulated different ways,  
but they are equivalent - give same sequence of x

Sutskever et al, “On the importance of initialization and momentum in deep learning”, ICML 2013

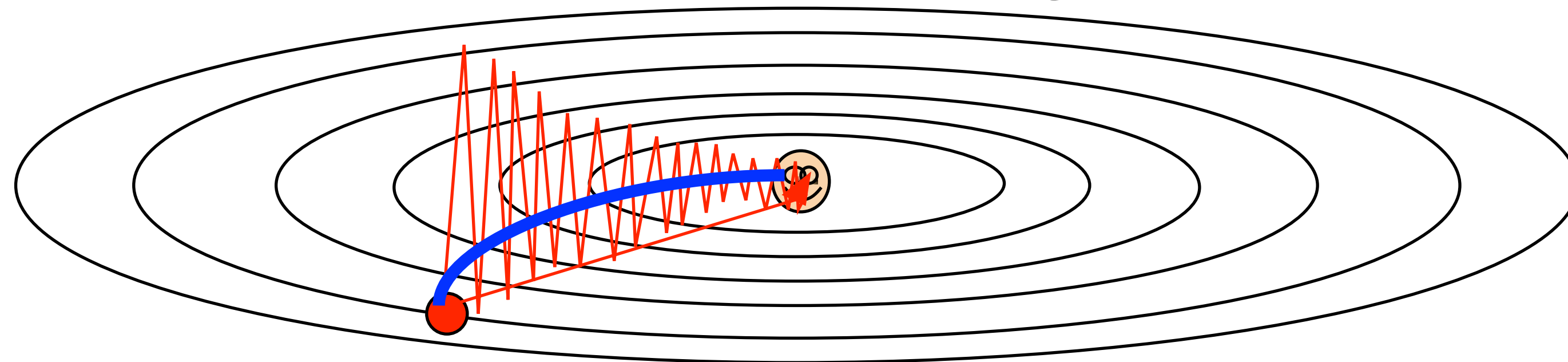
# SGD + Momentum

Local Minima

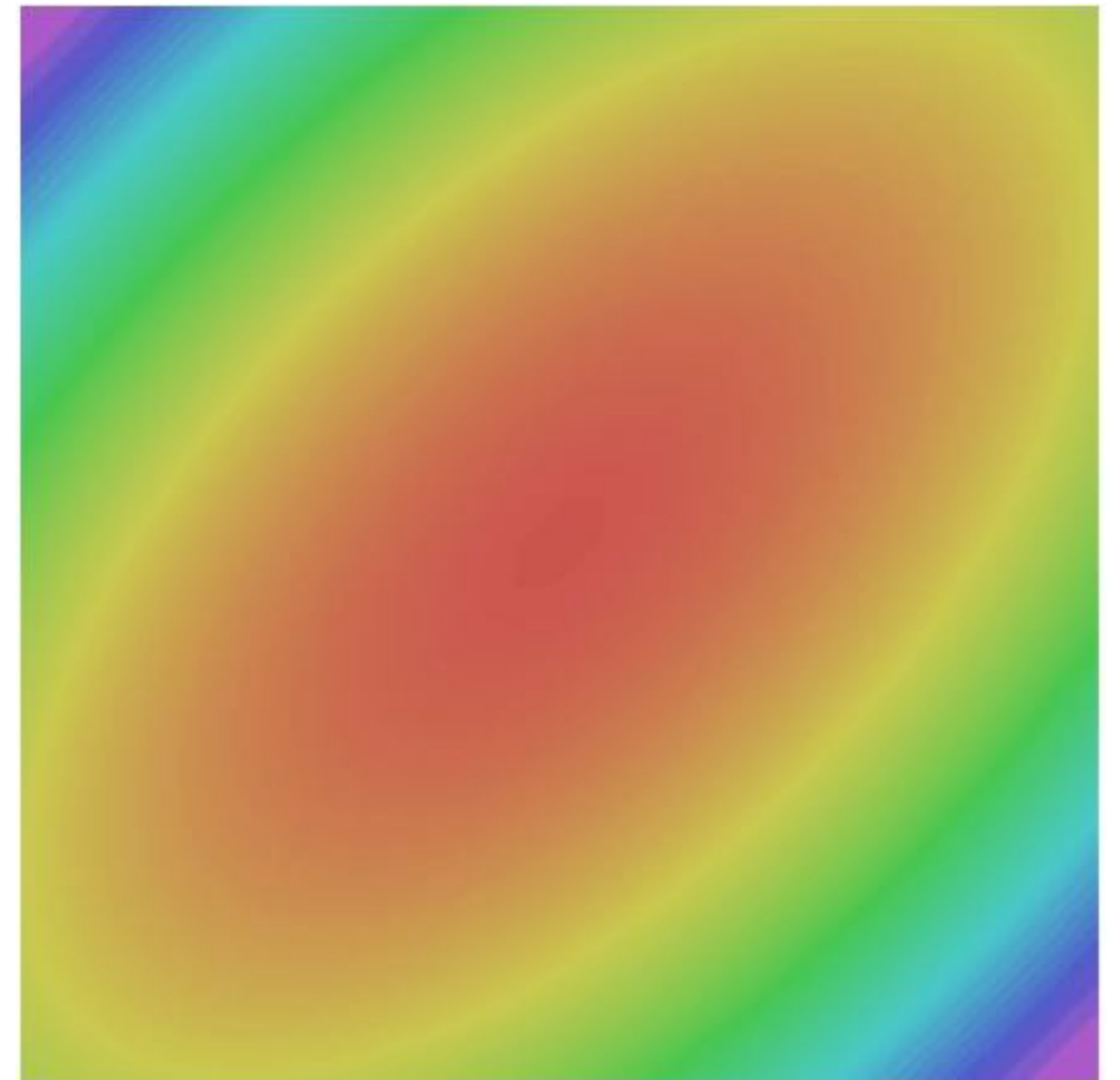
Saddle points



Poor Conditioning



Gradient Noise



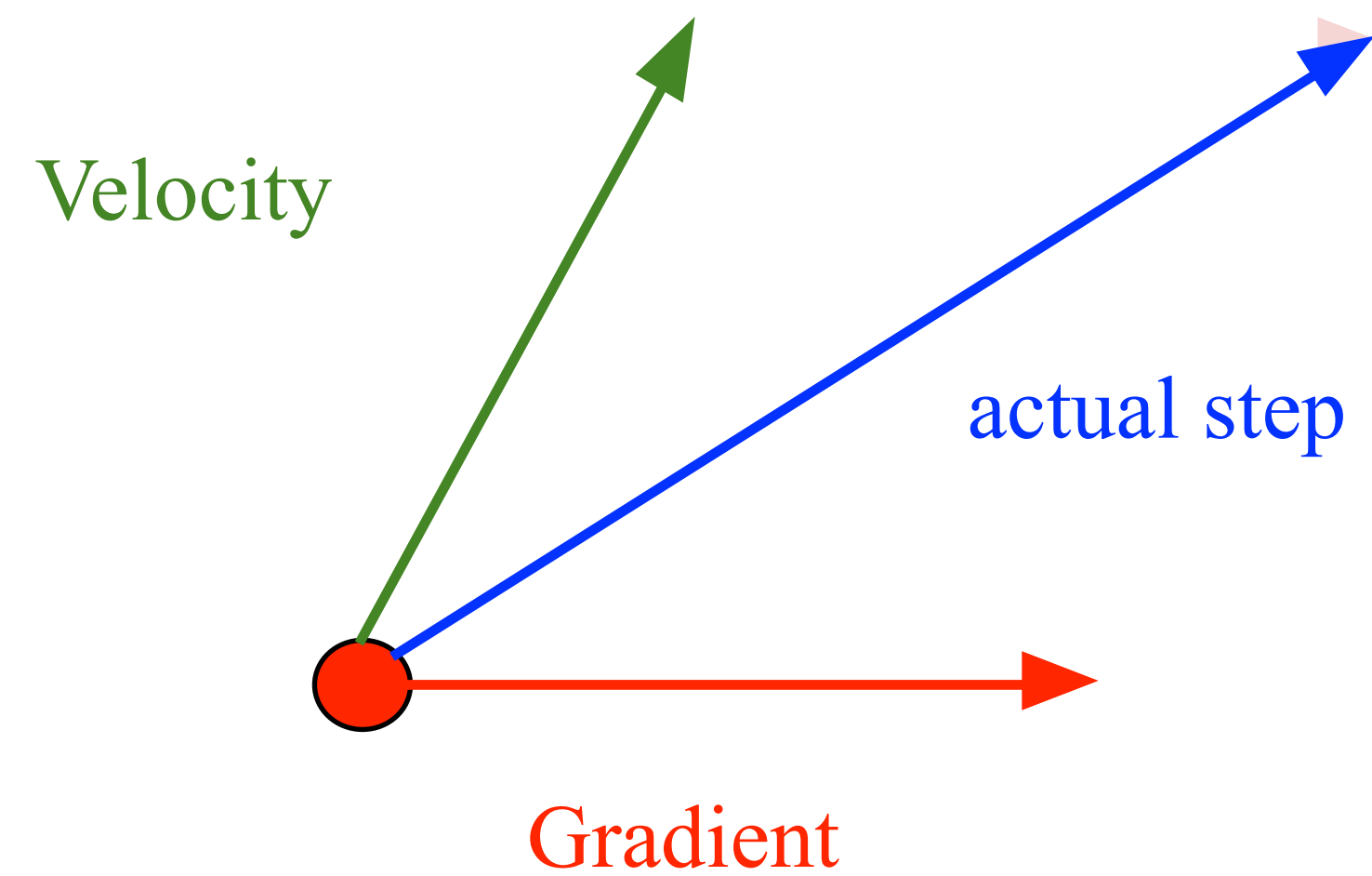
SGD

SGD+Momentum



# SGD+Momentum

Momentum update:



Combine gradient at current point with velocity to get step used to update weights

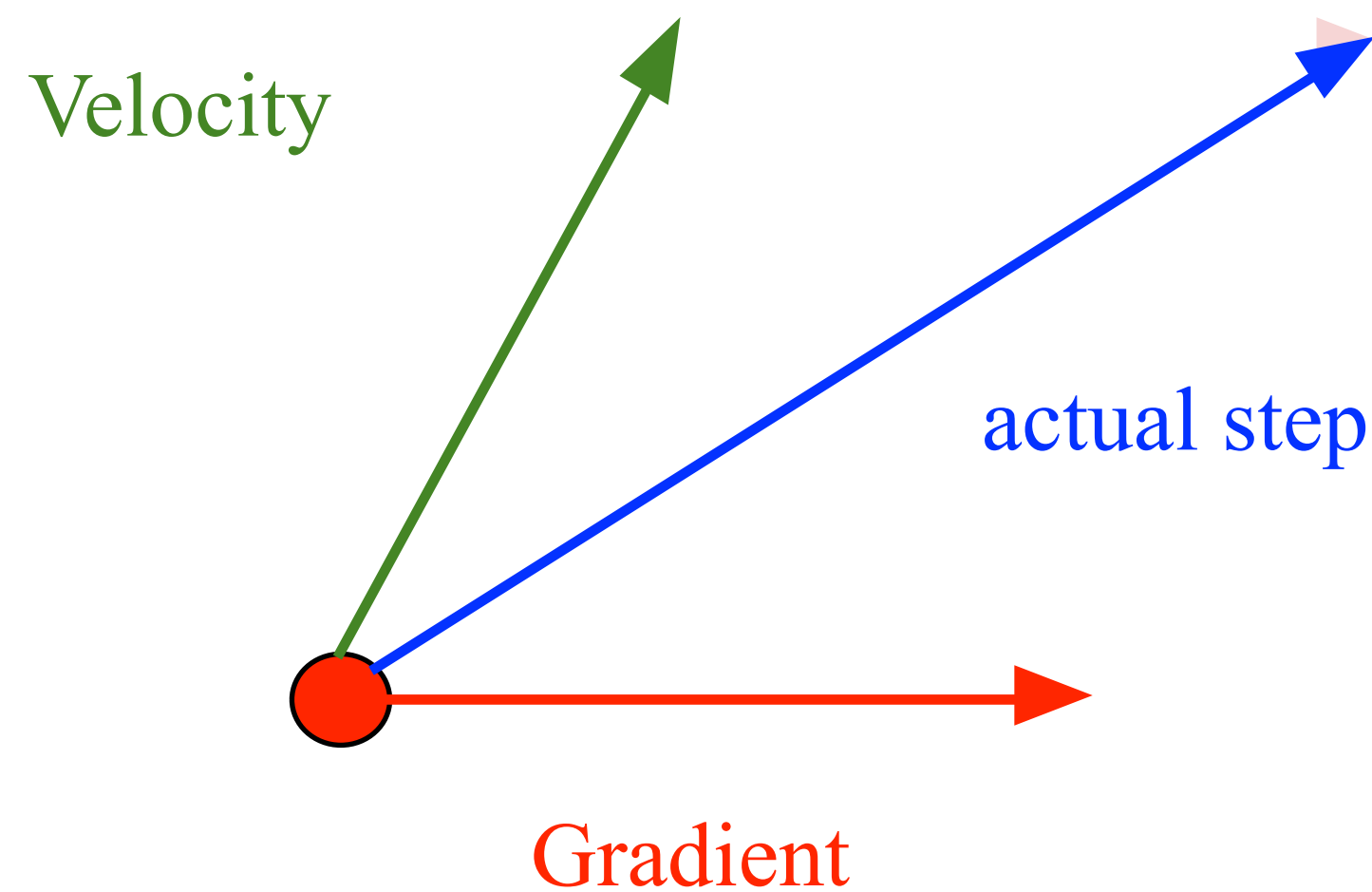
Nesterov, “A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ ”, 1983

Nesterov, “Introductory lectures on convex optimization: a basic course”, 2004

Sutskever et al, “On the importance of initialization and momentum in deep learning”, ICML 2013

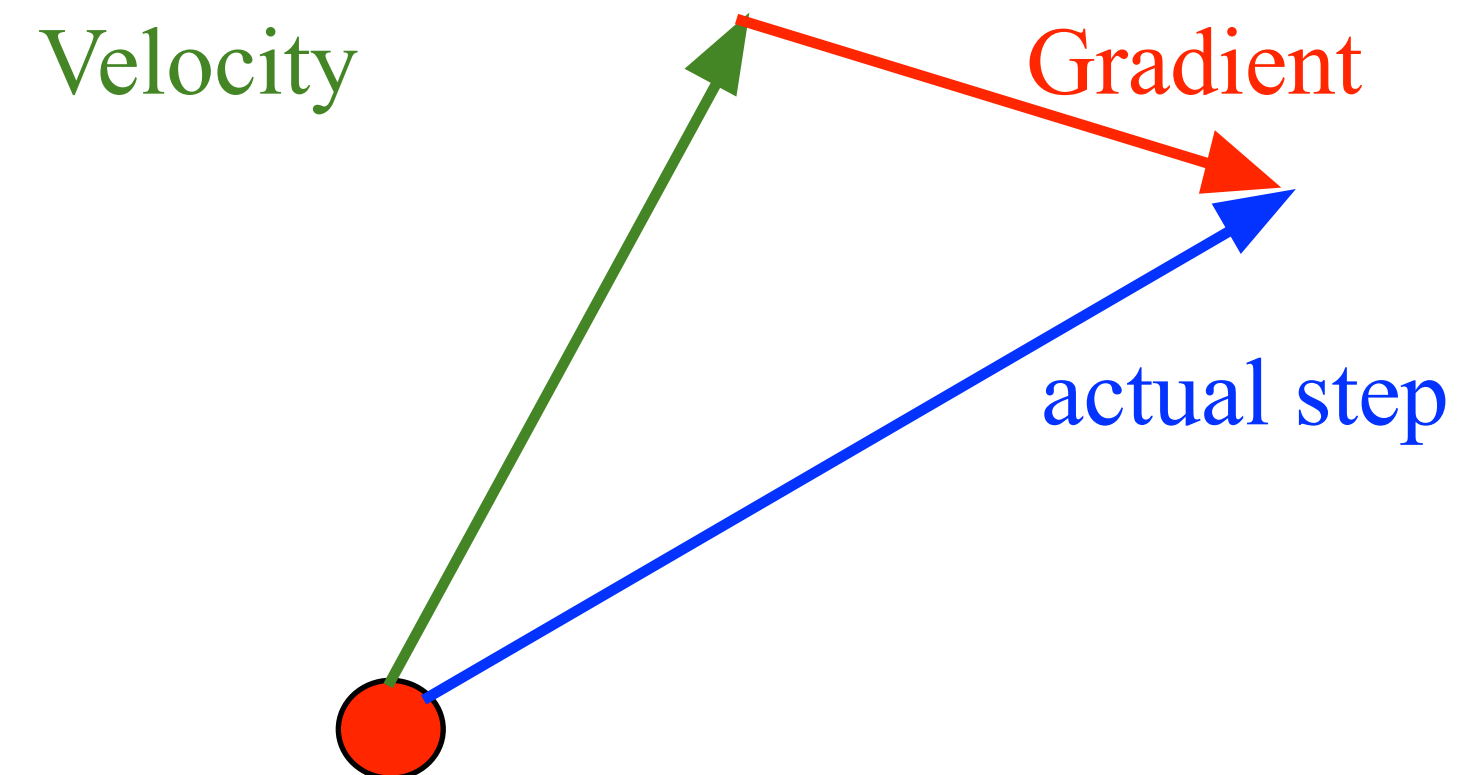
# Nesterov Momentum

Momentum update:



Combine gradient at current point with velocity to get step used to update weights

Nesterov Momentum



“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

Nesterov, “A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ ”, 1983

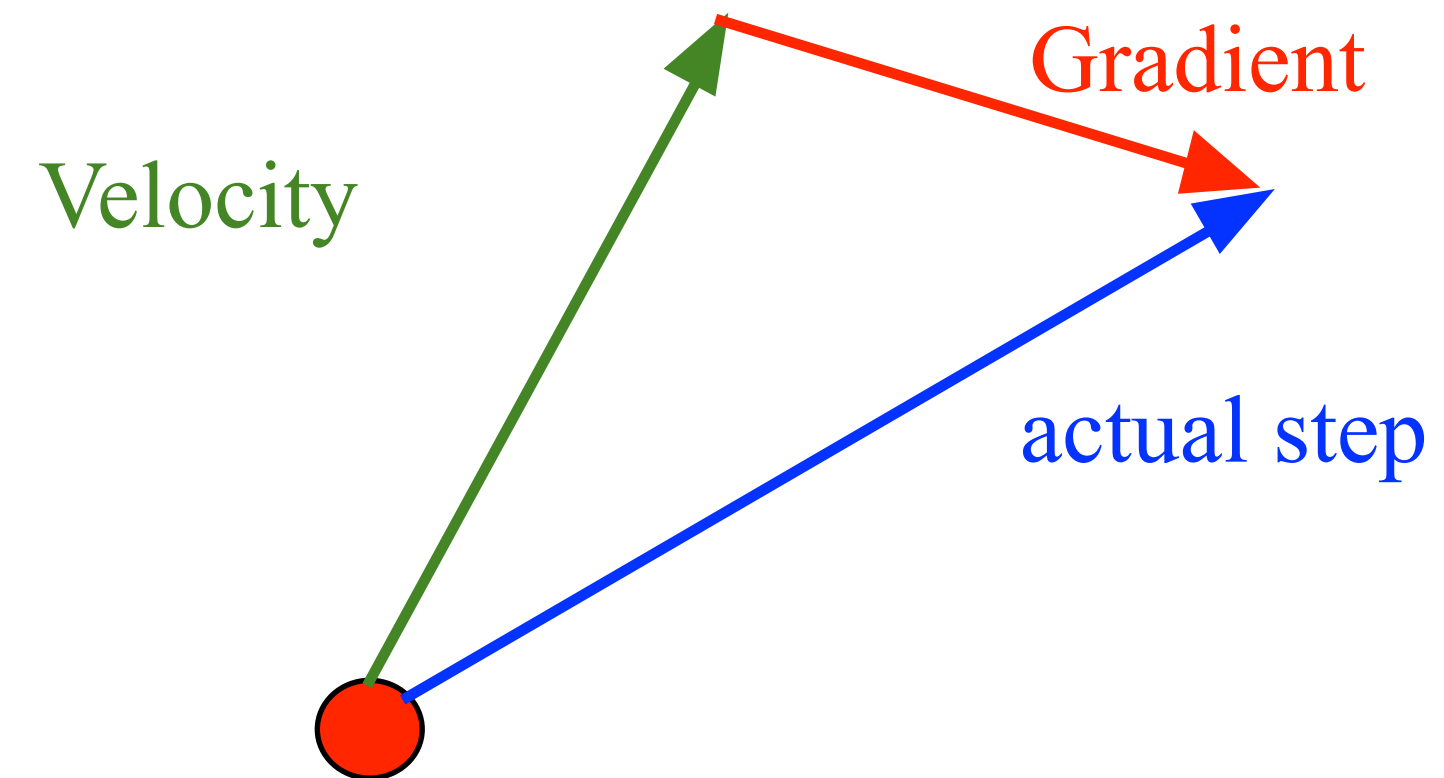
Nesterov, “Introductory lectures on convex optimization: a basic course”, 2004

Sutskever et al, “On the importance of initialization and momentum in deep learning”, ICML 2013

# Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$



“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

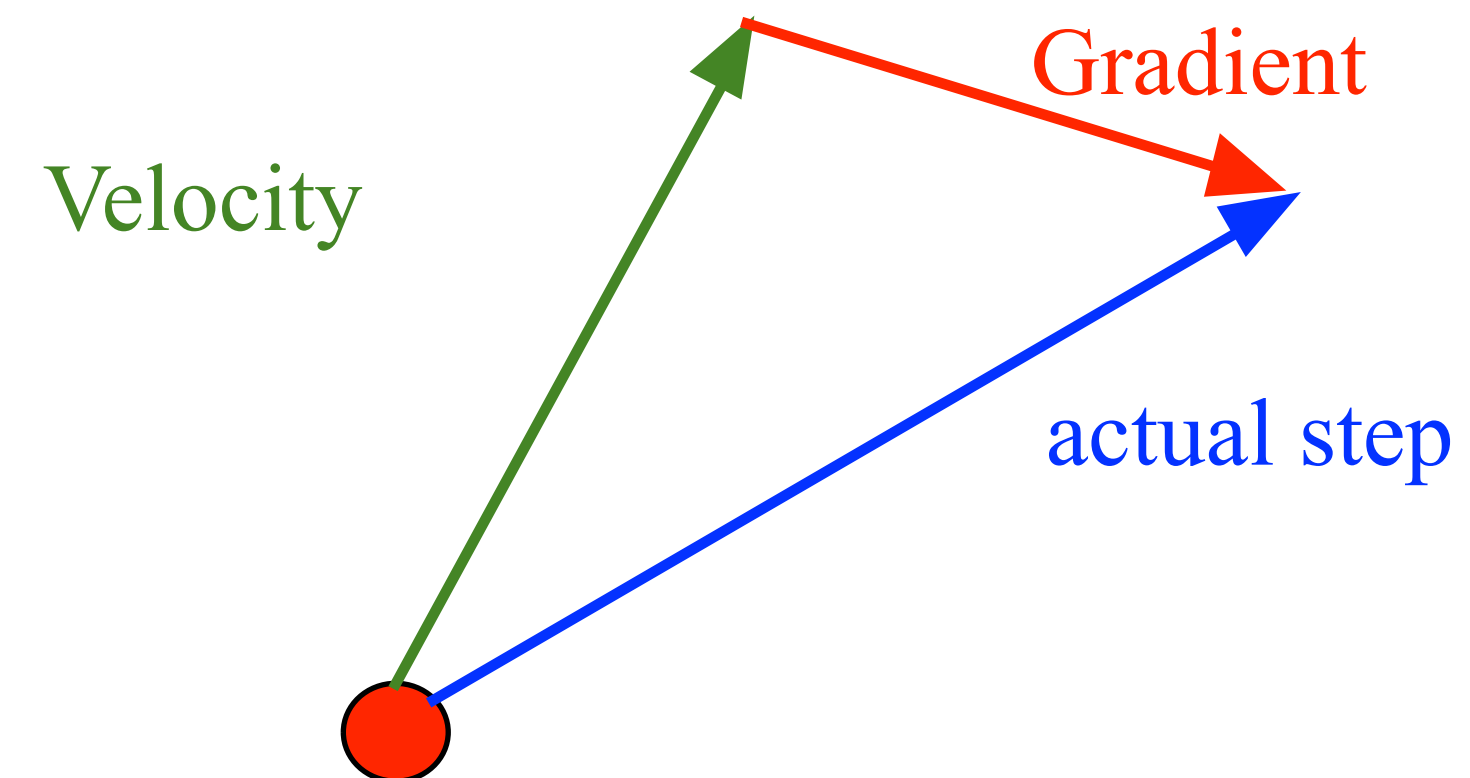


# Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$

Annoying, usually we want  
update in terms of  $x_t, \nabla f(x_t)$



“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

# Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

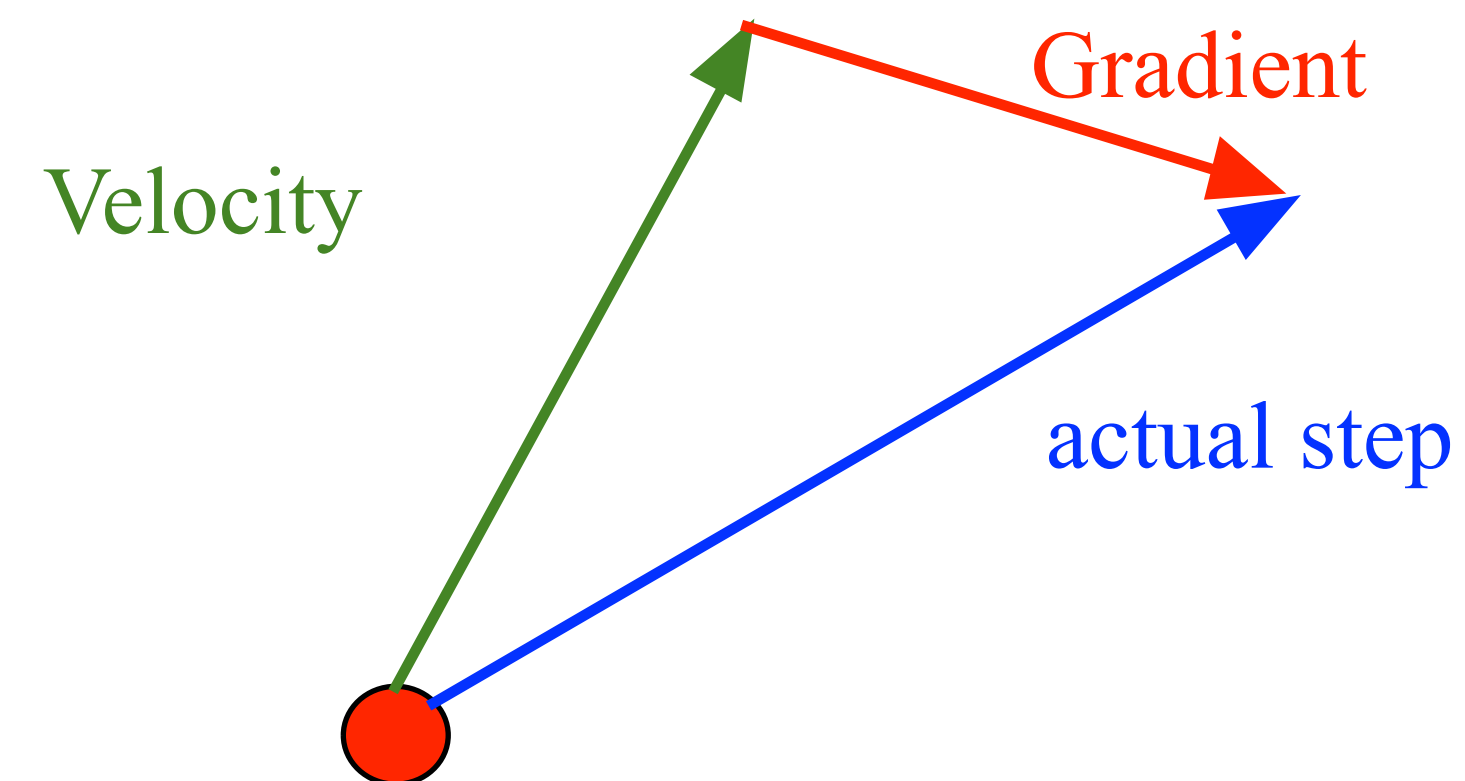
$$x_{t+1} = x_t + v_{t+1}$$

Change of variables  $\tilde{x}_t = x_t + \rho v_t$  and rearrange:

$$v_{t+1} = \rho v_t - \alpha \nabla f(\tilde{x}_t)$$

$$\begin{aligned}\tilde{x}_{t+1} &= \tilde{x}_t - \rho v_t + (1 + \rho)v_{t+1} \\ &= \tilde{x}_t + v_{t+1} + \rho(v_{t+1} - v_t)\end{aligned}$$

Annoying, usually we want update in terms of  $x_t, \nabla f(x_t)$



“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

# Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$

Change of variables  $\tilde{x}_t = x_t + \rho v_t$  and rearrange:

$$v_{t+1} = \rho v_t - \alpha \nabla f(\tilde{x}_t)$$

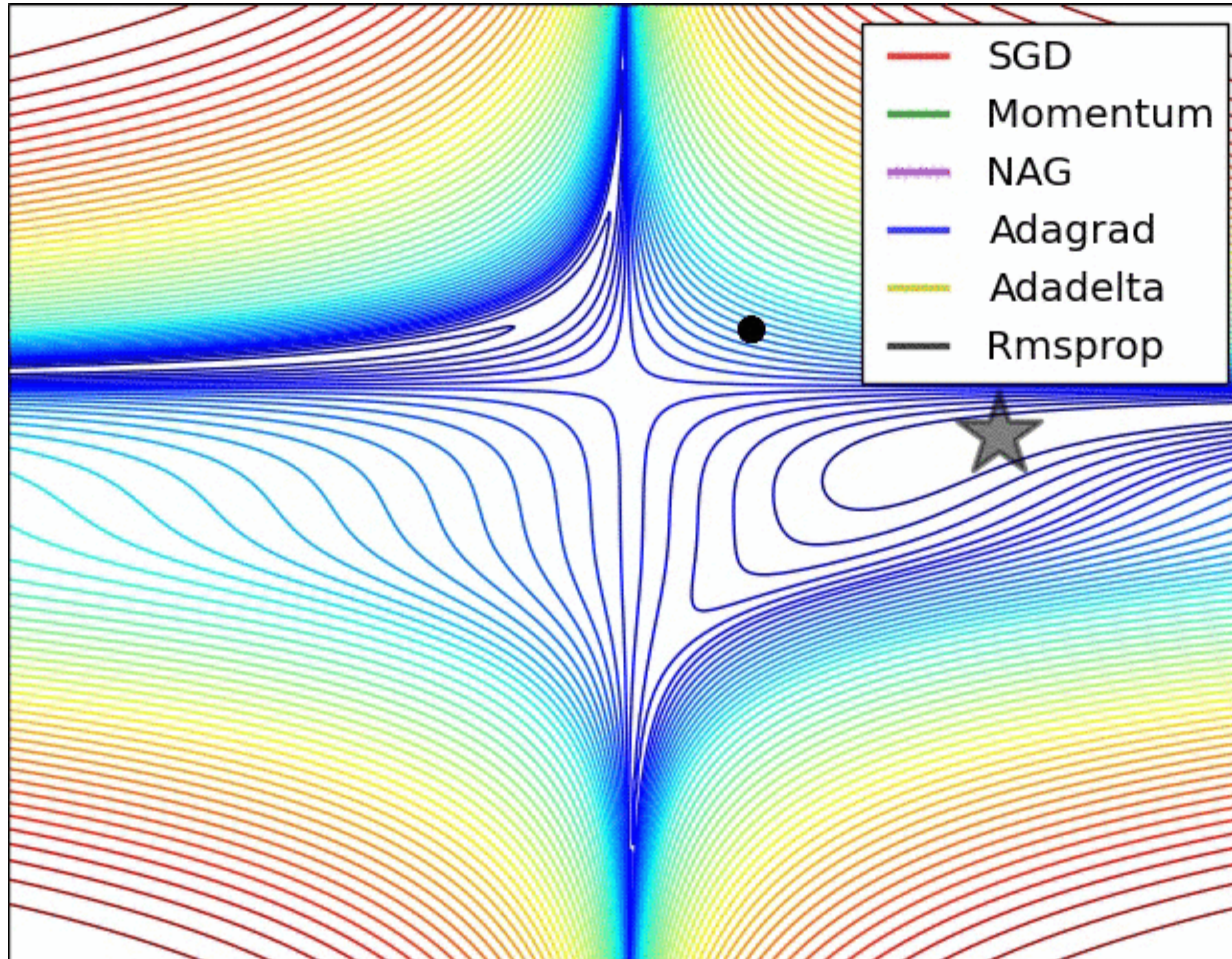
$$\begin{aligned}\tilde{x}_{t+1} &= \tilde{x}_t - \rho v_t + (1 + \rho)v_{t+1} \\ &= \tilde{x}_t + v_{t+1} + \rho(v_{t+1} - v_t)\end{aligned}$$

Annoying, usually we want update in terms of  $x_t, \nabla f(x_t)$

```
dx = compute_gradient(x)
old_v = v
v = rho * v - learning_rate * dx
x += -rho * old_v + (1 + rho) * v
```



# SGD vs Momentum vs Nesterov Momentum





# Case Study: AlexNet

*[Krizhevsky et al. 2012]*

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

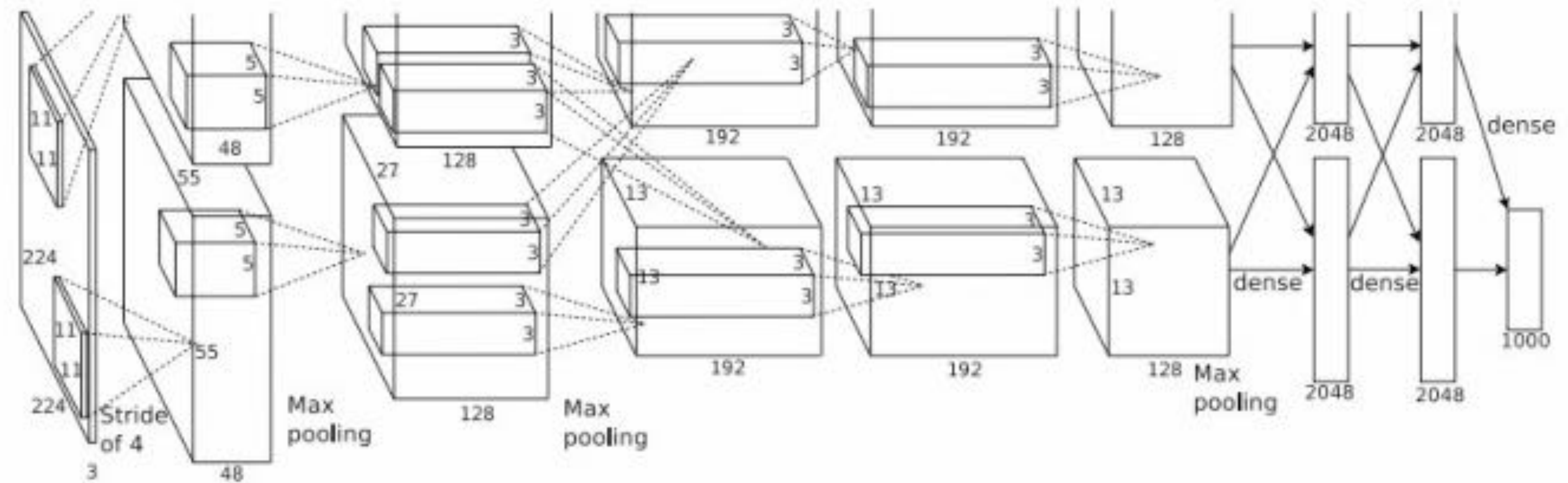
[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



## Details/Retrospectives:

- first use of ReLU
- used Local Response Normalization layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- **Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus**
- L2 weight decay 5e-4

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

# Case Study: AlexNet

*[Krizhevsky et al. 2012]*

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

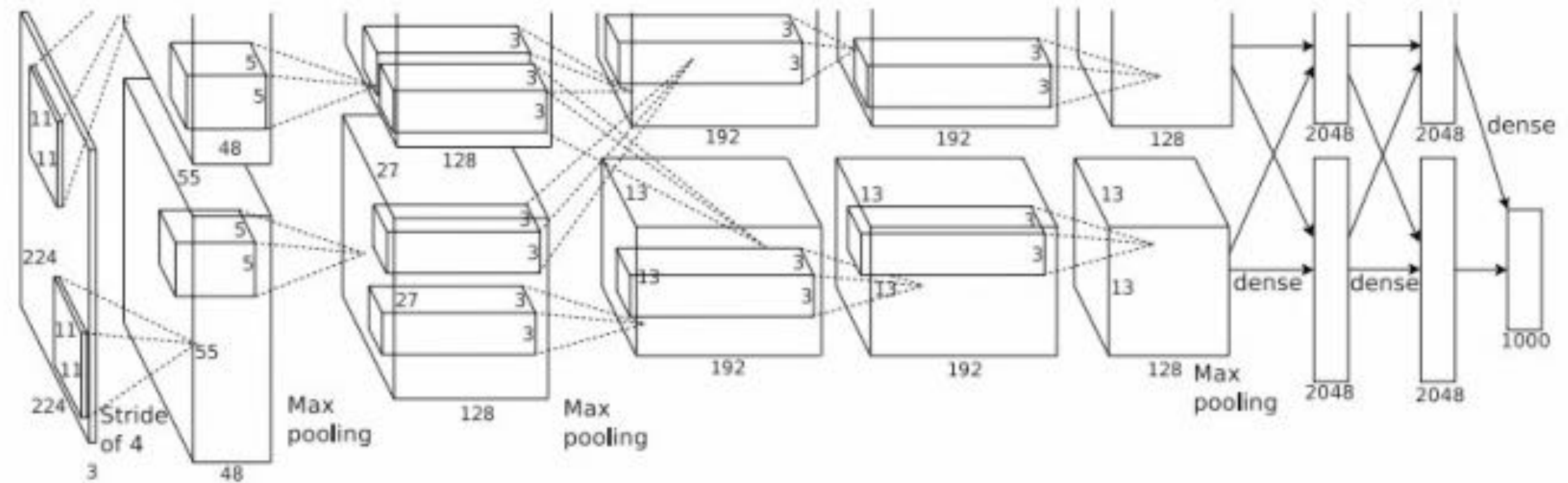
[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



## Details/Retrospectives:

- first use of ReLU
- used Local Response Normalization layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- **L2 weight decay 5e-4**

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



# Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

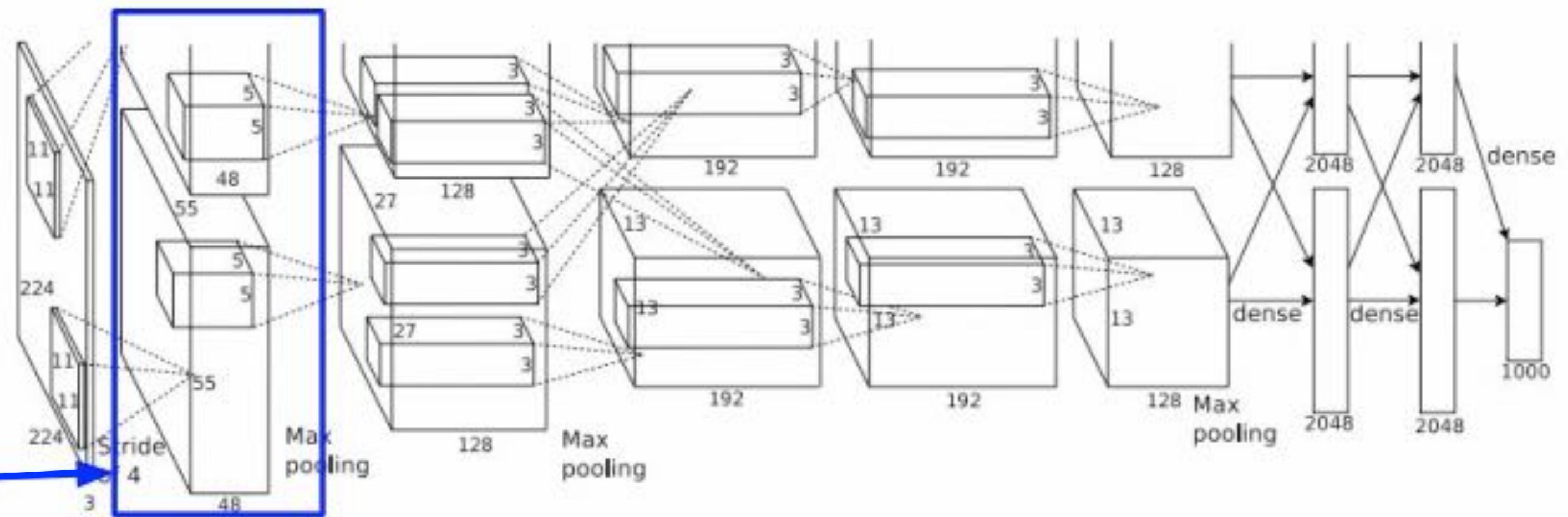
[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



[55x55x48] x 2

Historical note: Trained on GTX 580 GPU with only 3 GB of memory.  
Network spread across 2 GPUs, half the neurons (feature maps) on each GPU.

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

# Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

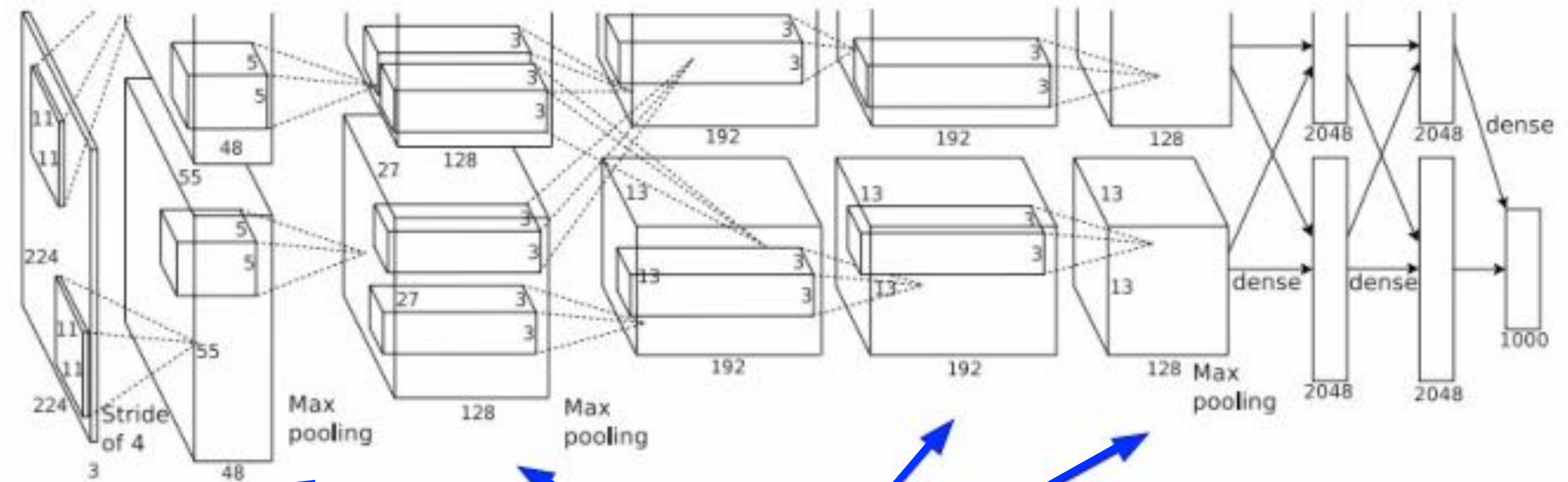
[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



**CONV1, CONV2, CONV4, CONV5:**  
Connections only with feature maps  
on same GPU

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



# Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

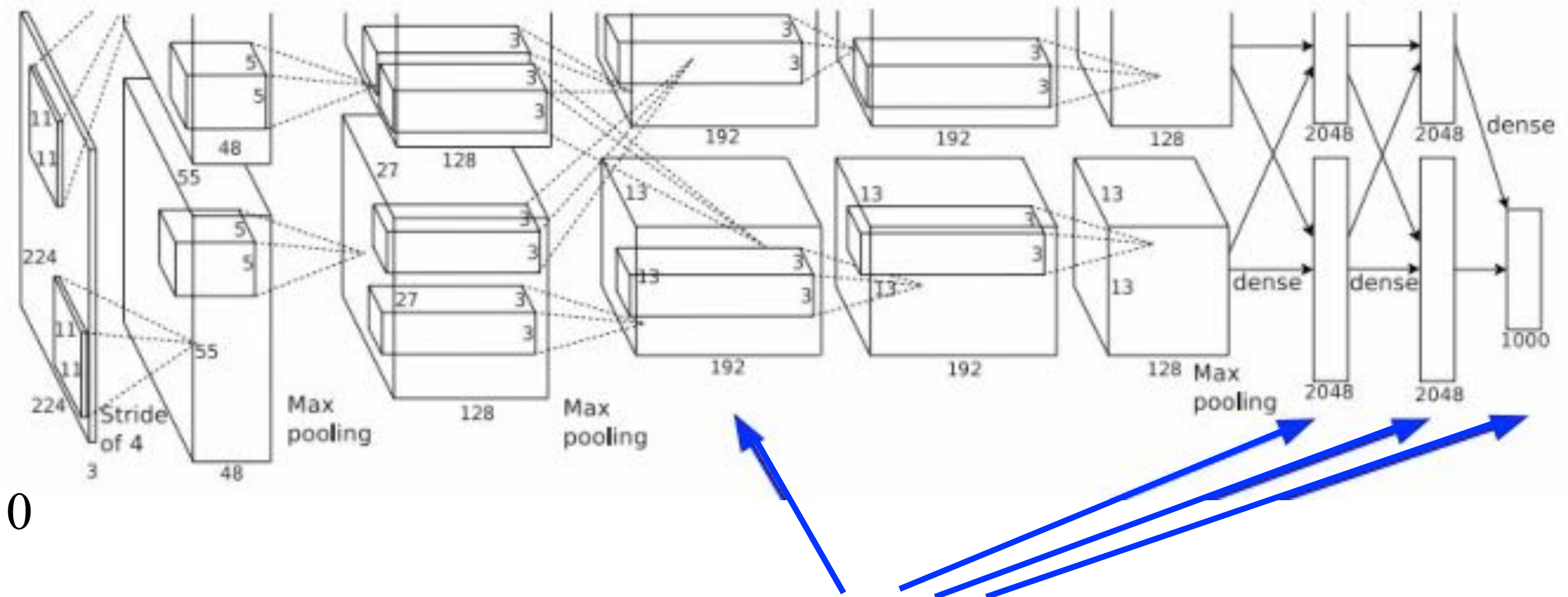
[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



**CONV3, FC6, FC7, FC8:**  
Connections with all feature maps in preceding layer, communication across GPUs

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



# Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] **CONV1** : 96 11x11 filters at stride 4, pad 0

[27x27x96] **MAX POOL1** : 3x3 filters at stride 2

[27x27x96] **NORM1**: Normalization layer

[27x27x256] **CONV2** : 256 5x5 filters at stride 1, pad 2

[13x13x256] **MAX POOL2**: 3x3 filters at stride 2

[13x13x256] **NORM2**: Normalization layer

[13x13x384] **CONV3** : 384 3x3 filters at stride 1, pad 1

[13x13x384] **CONV4** : 384 3x3 filters at stride 1, pad 1

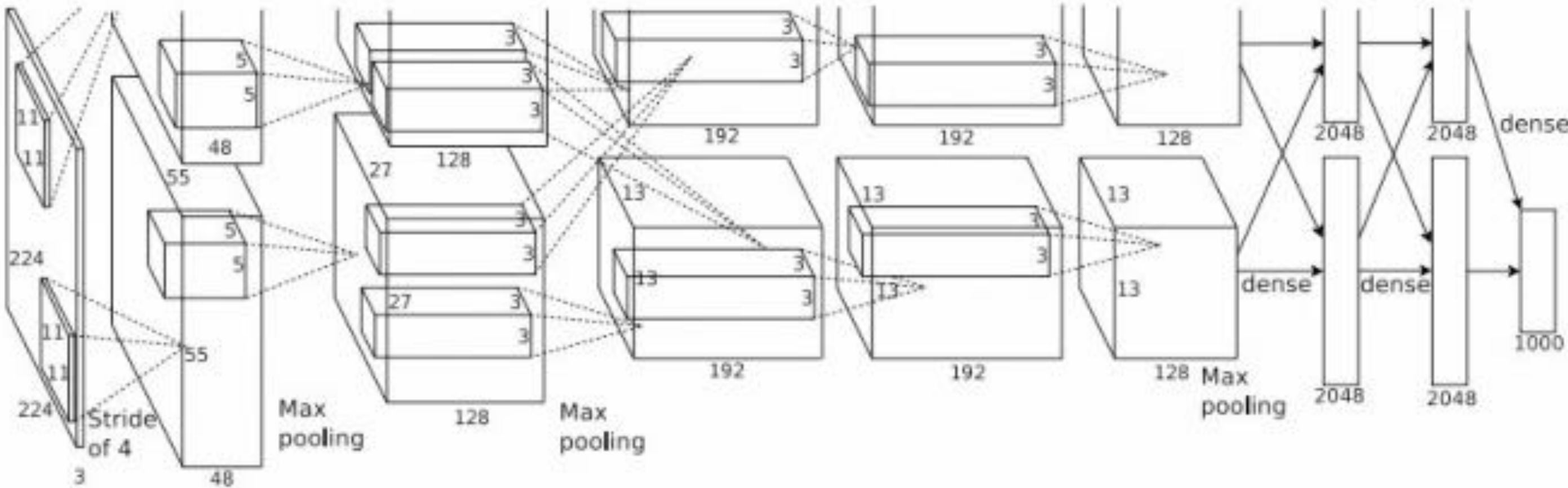
[13x13x256] **CONV5** : 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3** : 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



Model	Top-1 (val)	Top-5 (val)	Top-5 (test)
<i>SIFT + FVs</i> [7]	—	—	26.2%
1 CNN	40.7%	18.2%	—
5 CNNs	38.1%	16.4%	<b>16.4%</b>
1 CNN*	39.0%	16.6%	—
7 CNNs*	36.7%	15.4%	<b>15.3%</b>

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

# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

