

Shanghai Jiao Tong University

# Computer Vision

Instructor: Xu Zhao  
Class No.: C032703 F032528

Spring 2020



Xu Zhao @ Shanghai Jiao Tong university

---

## Lecture 12: Deep Neural Networks - Part 2

---

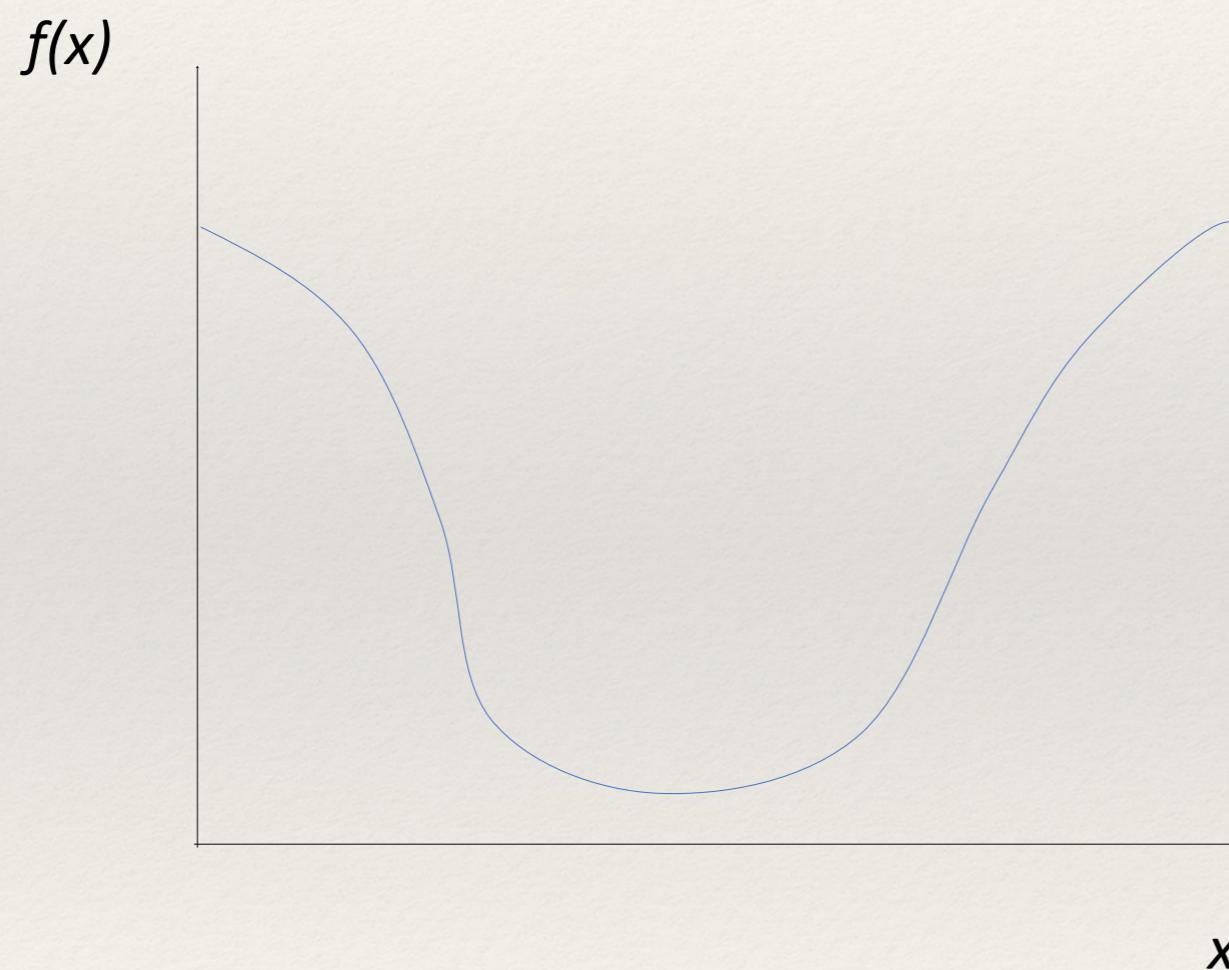
---

# Contents

---

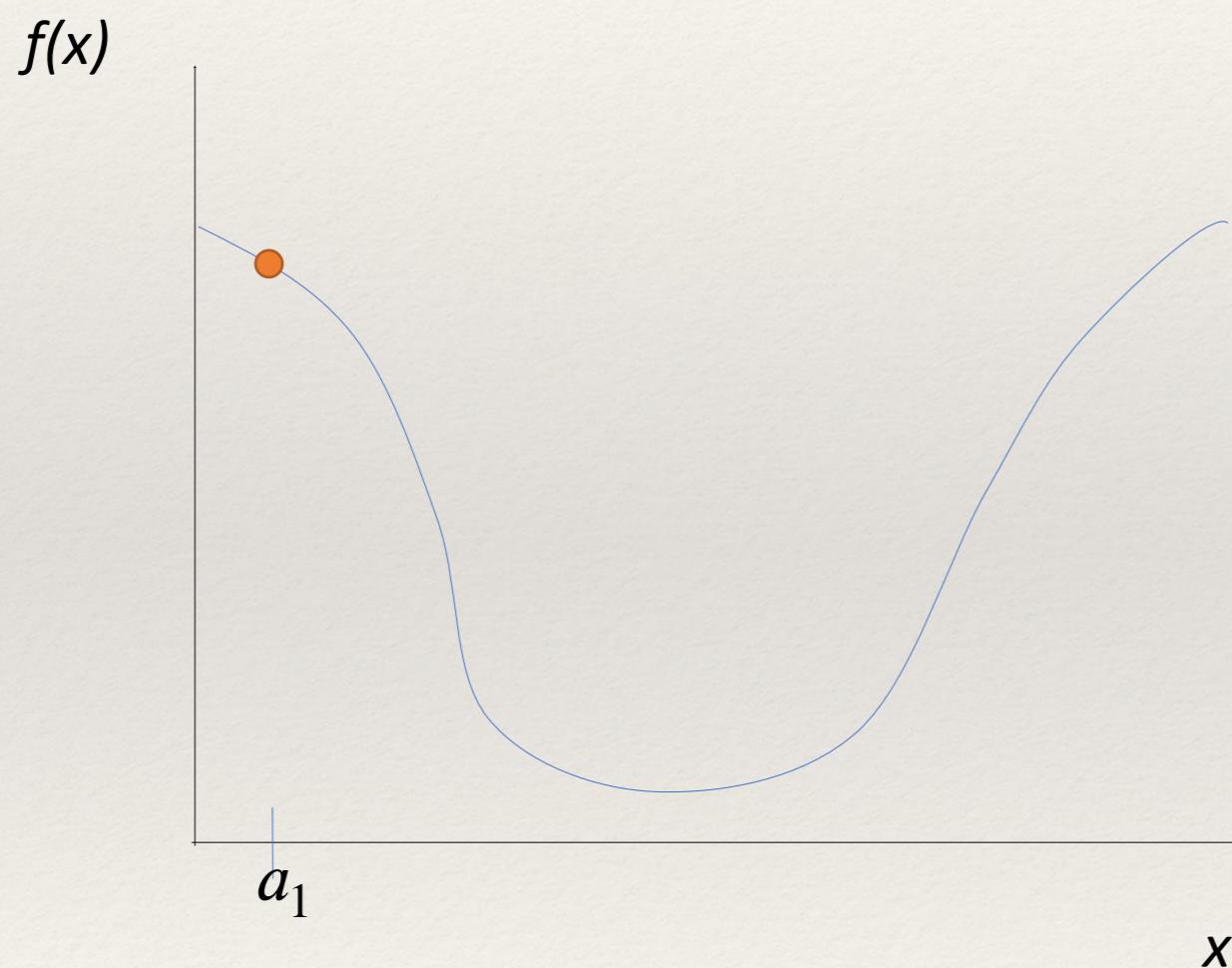
- ❖ **Training Neural Networks**
- ❖ **Specific Architectures for Diverse Applications**

# Gradient descent



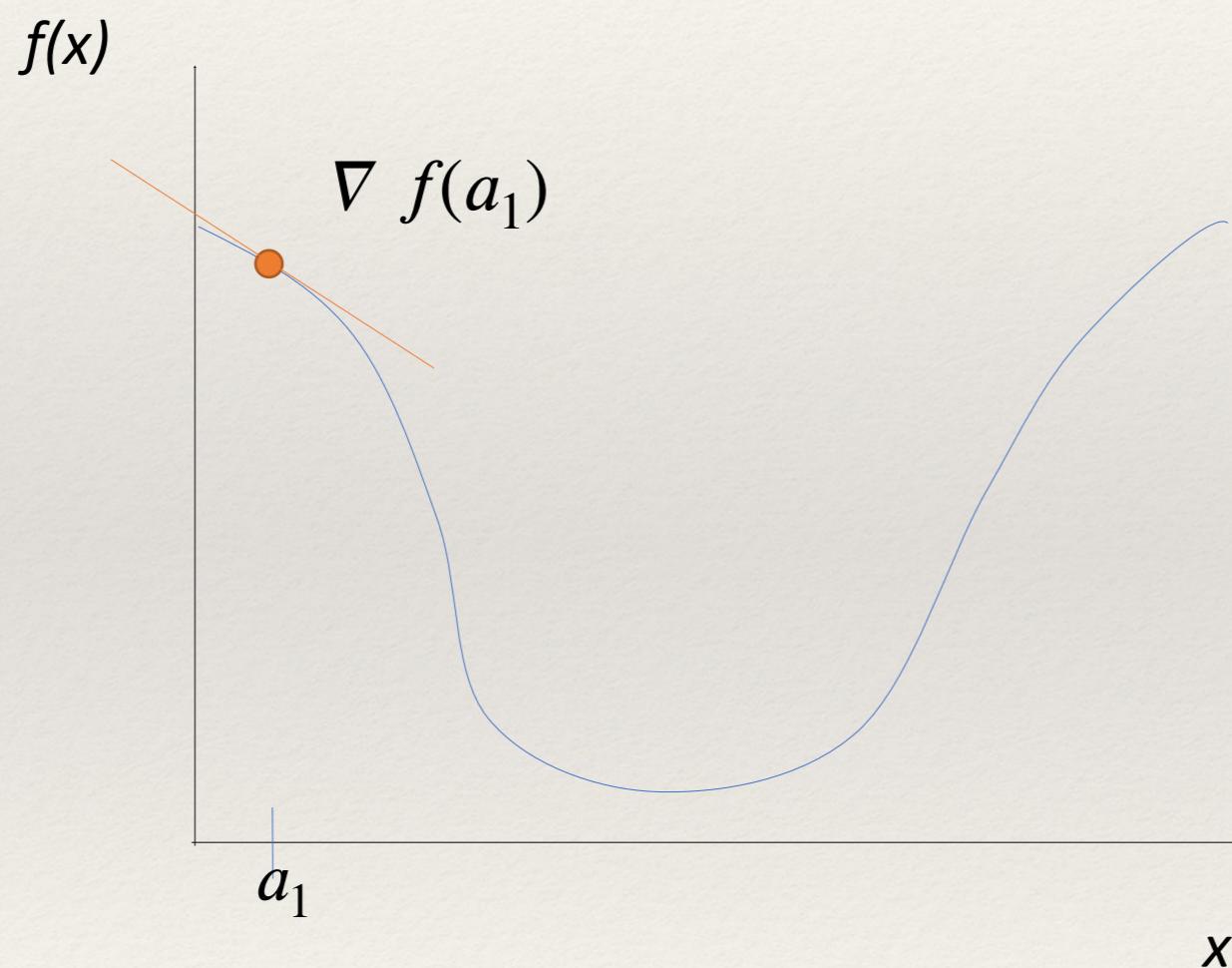
# General approach

Pick random starting point.



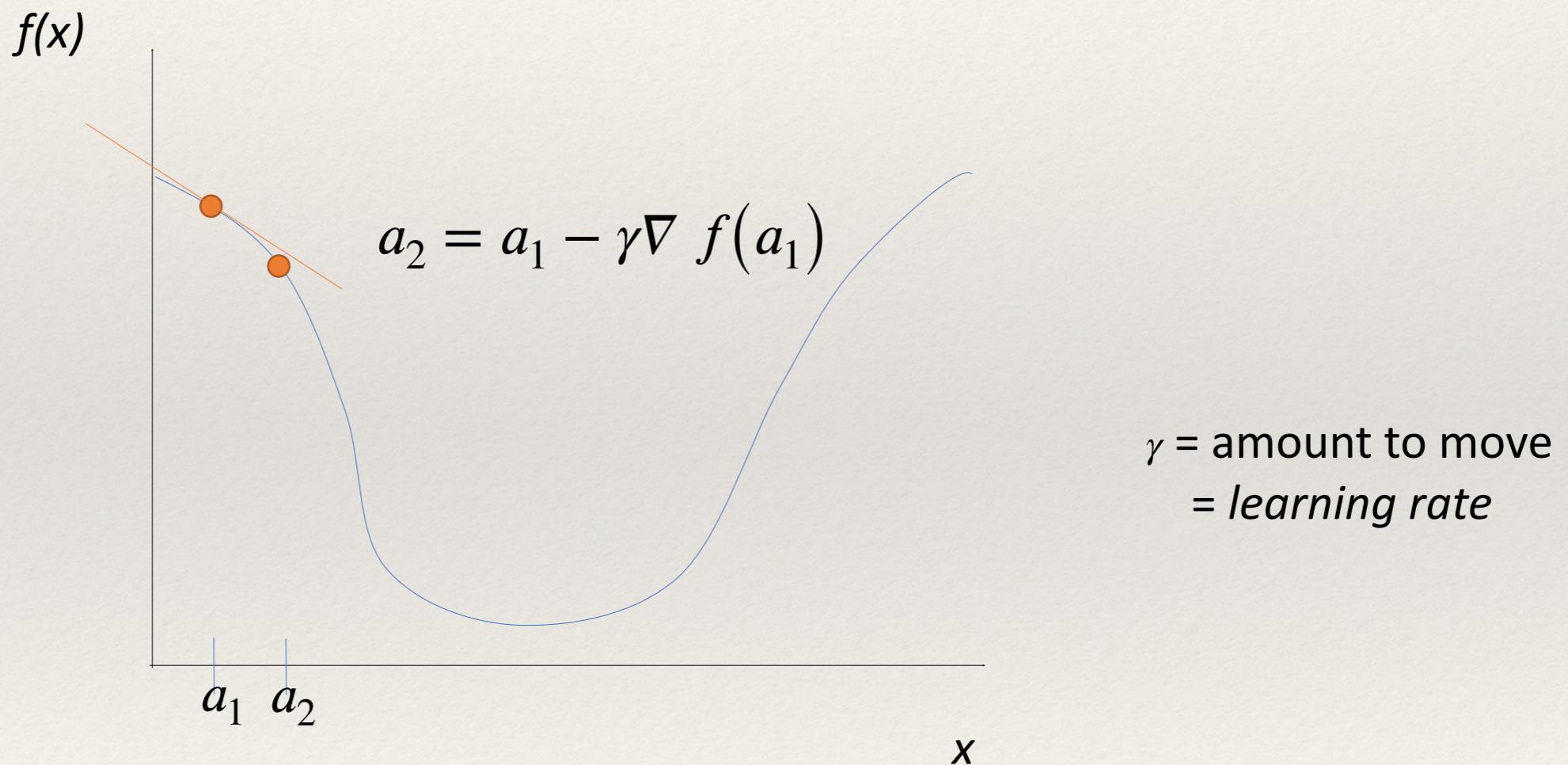
# General approach

Compute gradient at point (analytically or by finite differences)



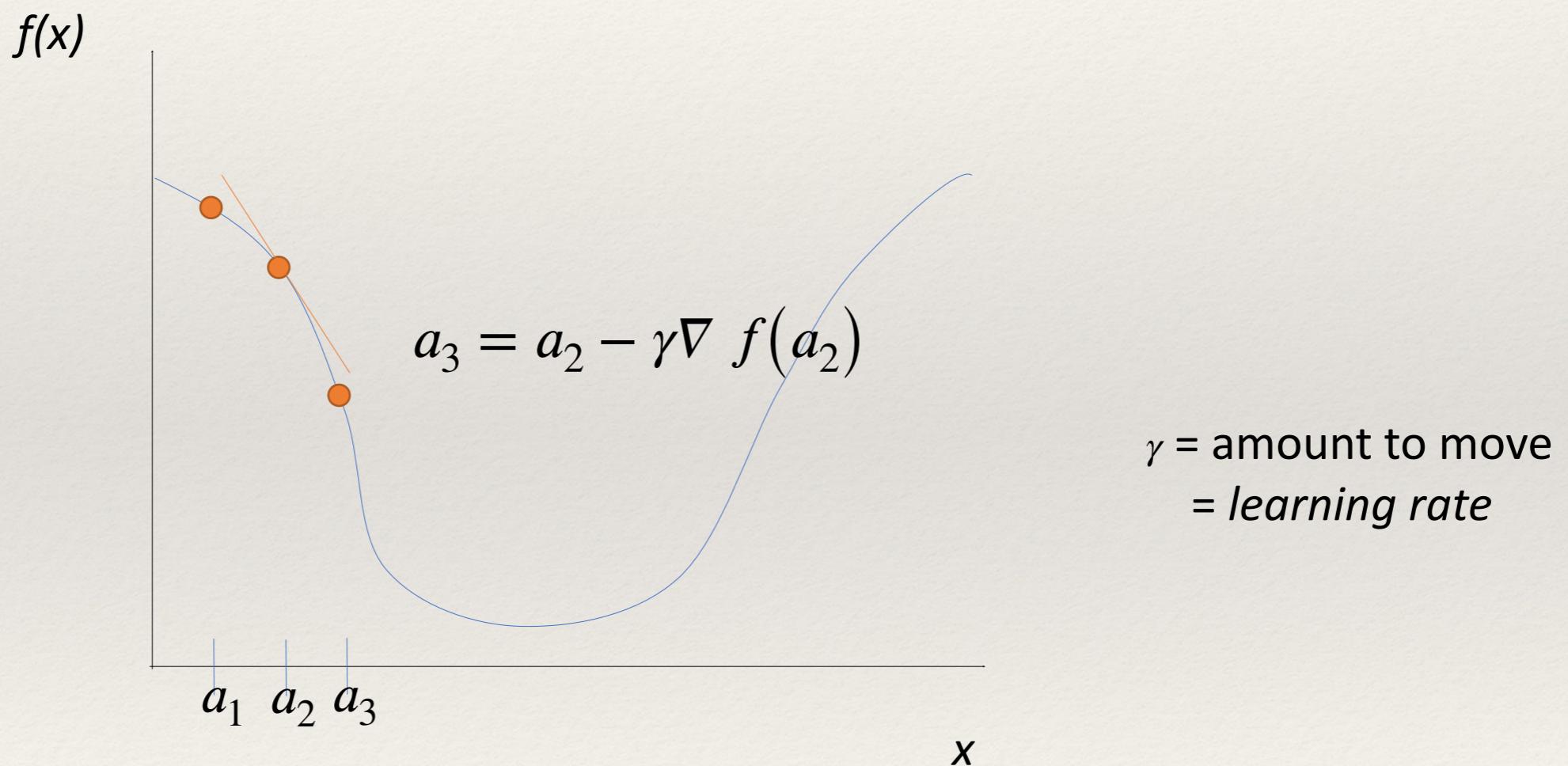
# General approach

Move along parameter space in direction of negative gradient



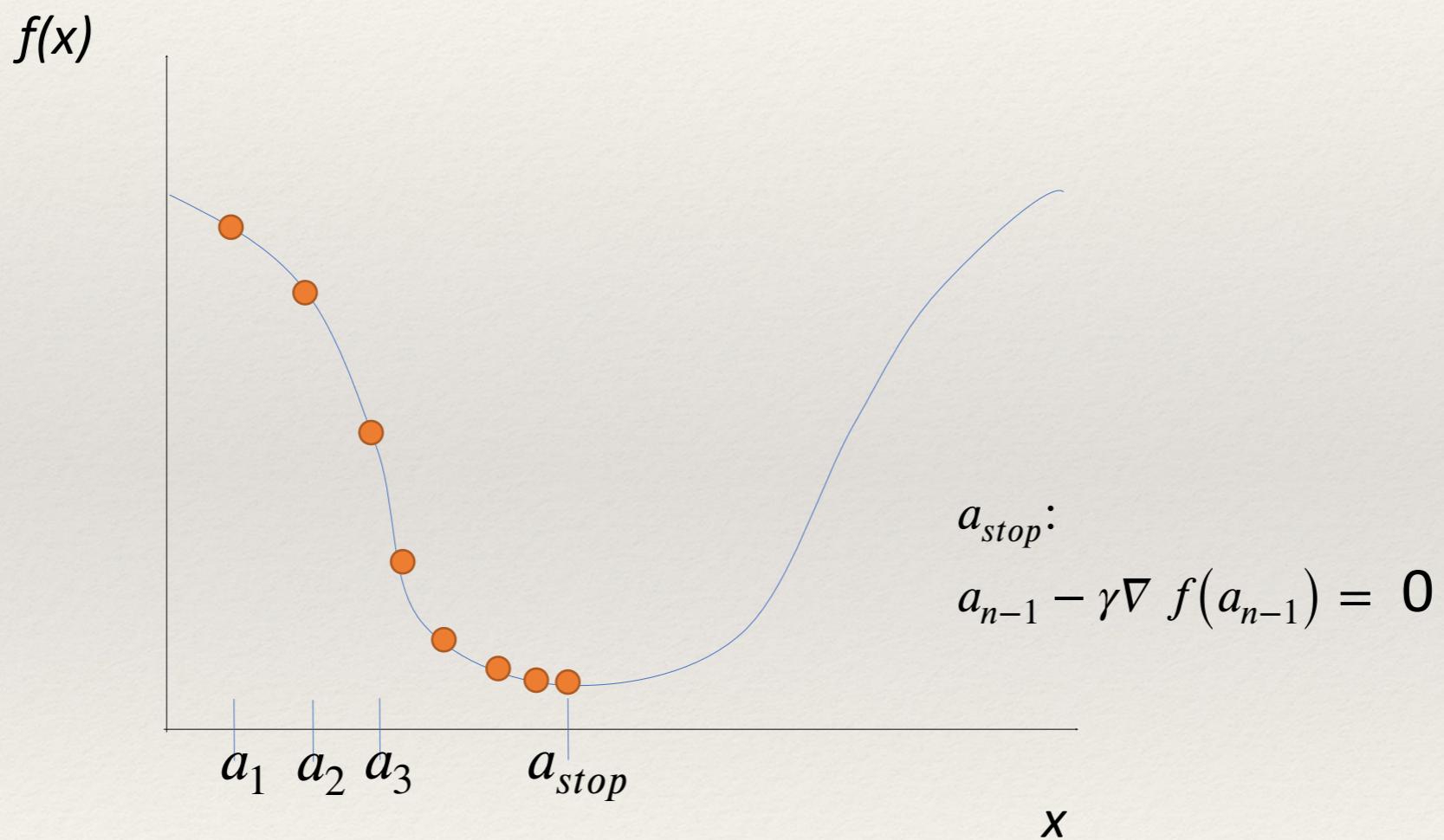
# General approach

Move along parameter space in direction of negative gradient.



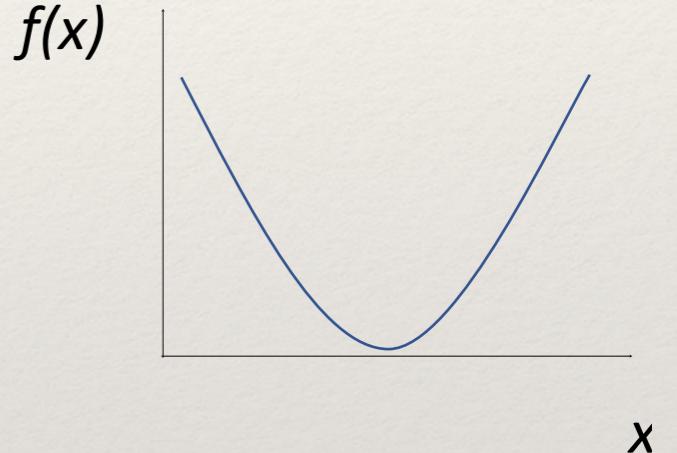
# General approach

Stop when we don't move any more.



# Gradient descent

- Optimizer for functions.
- Guaranteed to find optimum for convex functions.
  - Non-convex = find *local* optimum.
  - Most vision problems aren't convex.
- Works for multi-variate functions.
  - Need to compute matrix of *partial derivatives* ("Jacobian")



# Why would I use this over Least Squares?

---

If my function is convex,  
why can't I just use linear least squares?

$$A\mathbf{x} - \mathbf{b} = 0$$

$$F(\mathbf{x}) = \|A\mathbf{x} - \mathbf{b}\|^2.$$

$$\nabla F(\mathbf{x}) = 2A^T(A\mathbf{x} - \mathbf{b}).$$

Analytic solution = normal equations

$$\mathbf{x} = (A^T A)^{-1} A^T \mathbf{b}$$

You can, yes.

---

# Why would I use this over Least Squares?

---

But now imagine that I have 1,000,000 data points.

Matrices are huge.

Even for convex functions, gradient descent allows me to iteratively solve the solution without requiring very large matrices.

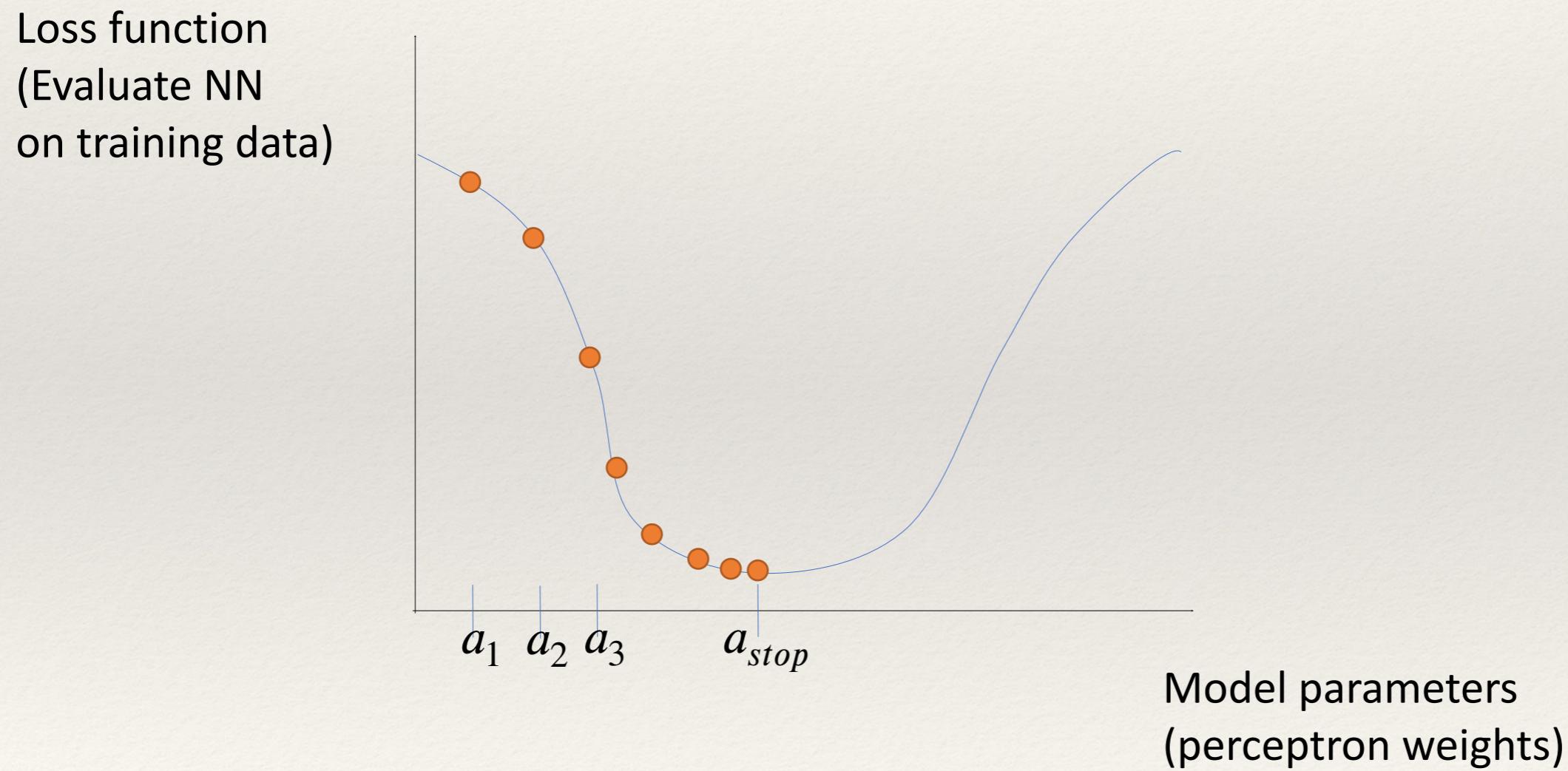
We'll see how.

# Train NN with Gradient Descent

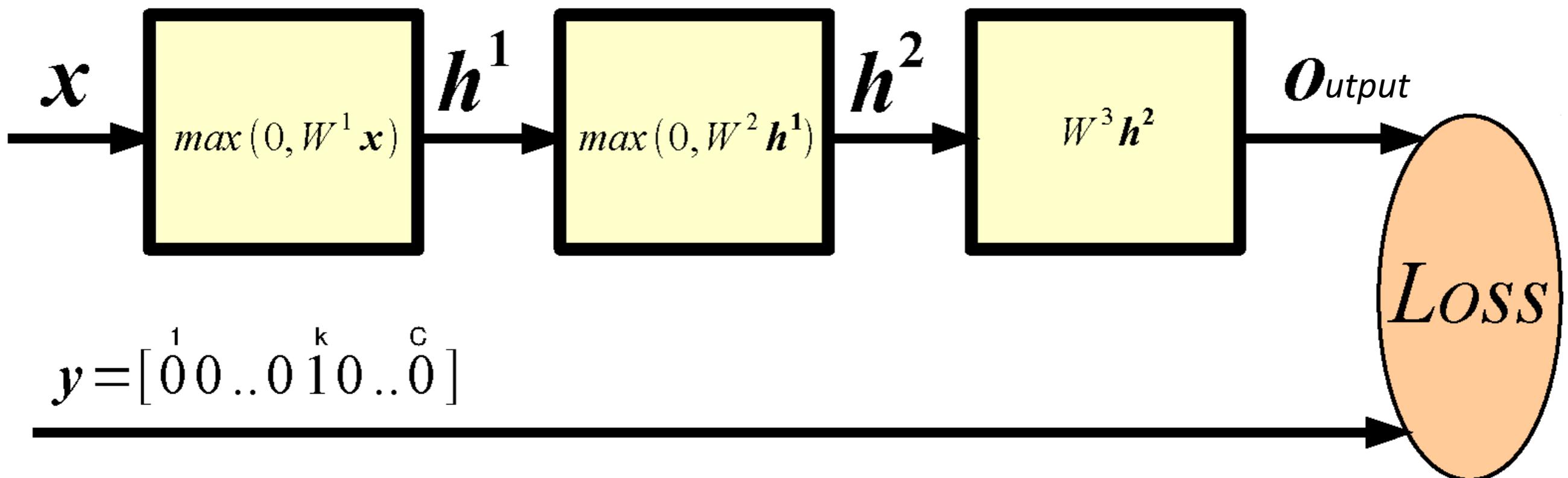
---

- $x^i, y^i = n$  training examples
- $f(\mathbf{x})$  = feed forward neural network
- $L(\mathbf{x}, y; \theta)$  = some *loss function*
- *Loss function* measures how ‘good’ our network is at classifying the training examples wrt. the parameters of the model (the perceptron weights).

# Train NN with Gradient Descent



# How Good is a Network?



# What is an appropriate loss?

- Define some output threshold on detection
- Classification: compare training class to output class
- Zero-one loss  $L$  (per class)

$y = \text{true label}$   $\hat{y} = \text{predicted label}$

$$L(\hat{y}, y) = I(\hat{y} \neq y),$$

- Is it good?
  - Nope – it's a step function.
  - I need to compute the *gradient* of the loss.
  - This loss is not differentiable, and 'flips' easily.

# Classification as probability

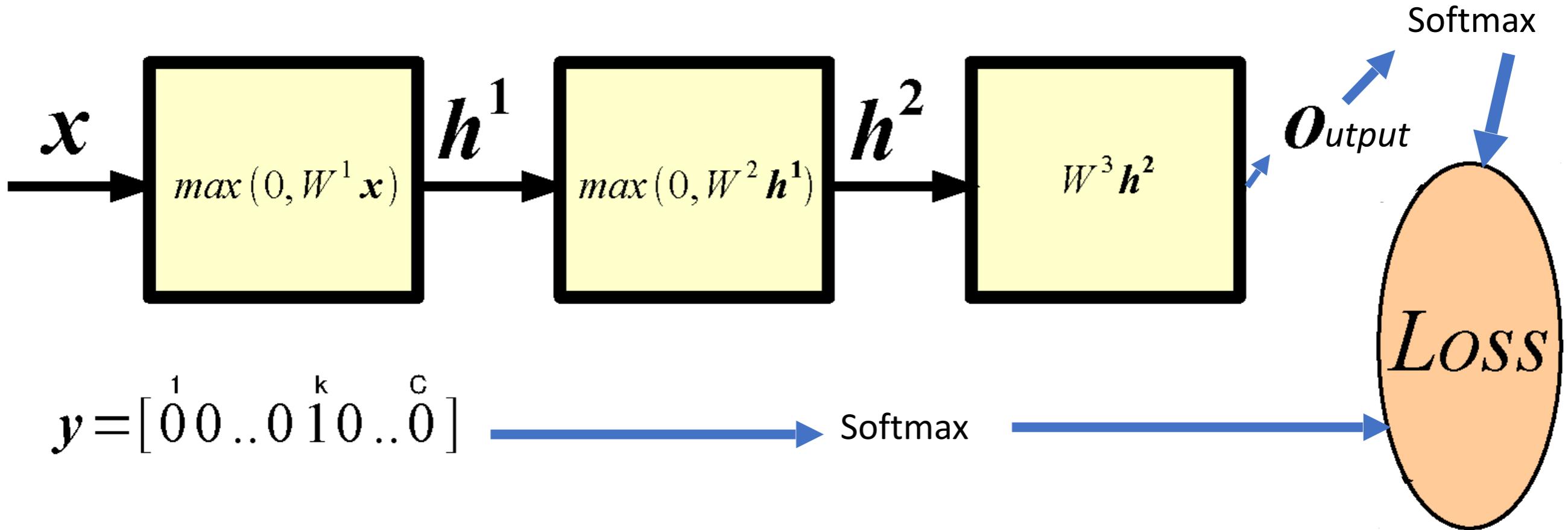
Special function on last layer - ‘Softmax’:

“Squashes” a  $C$ -dimensional vector  $\mathbf{O}$  of arbitrary real values to a  $C$ -dimensional vector  $\sigma(\mathbf{O})$  of real values in the range  $(0, 1)$  that add up to 1.

Turns the output into a probability distribution on classes.

$$p(c_k=1|x) = \frac{e^{o_k}}{\sum_{j=1}^C e^{o_j}}$$

# How Good is a Network?



Probability of class  $k$  given input (softmax):

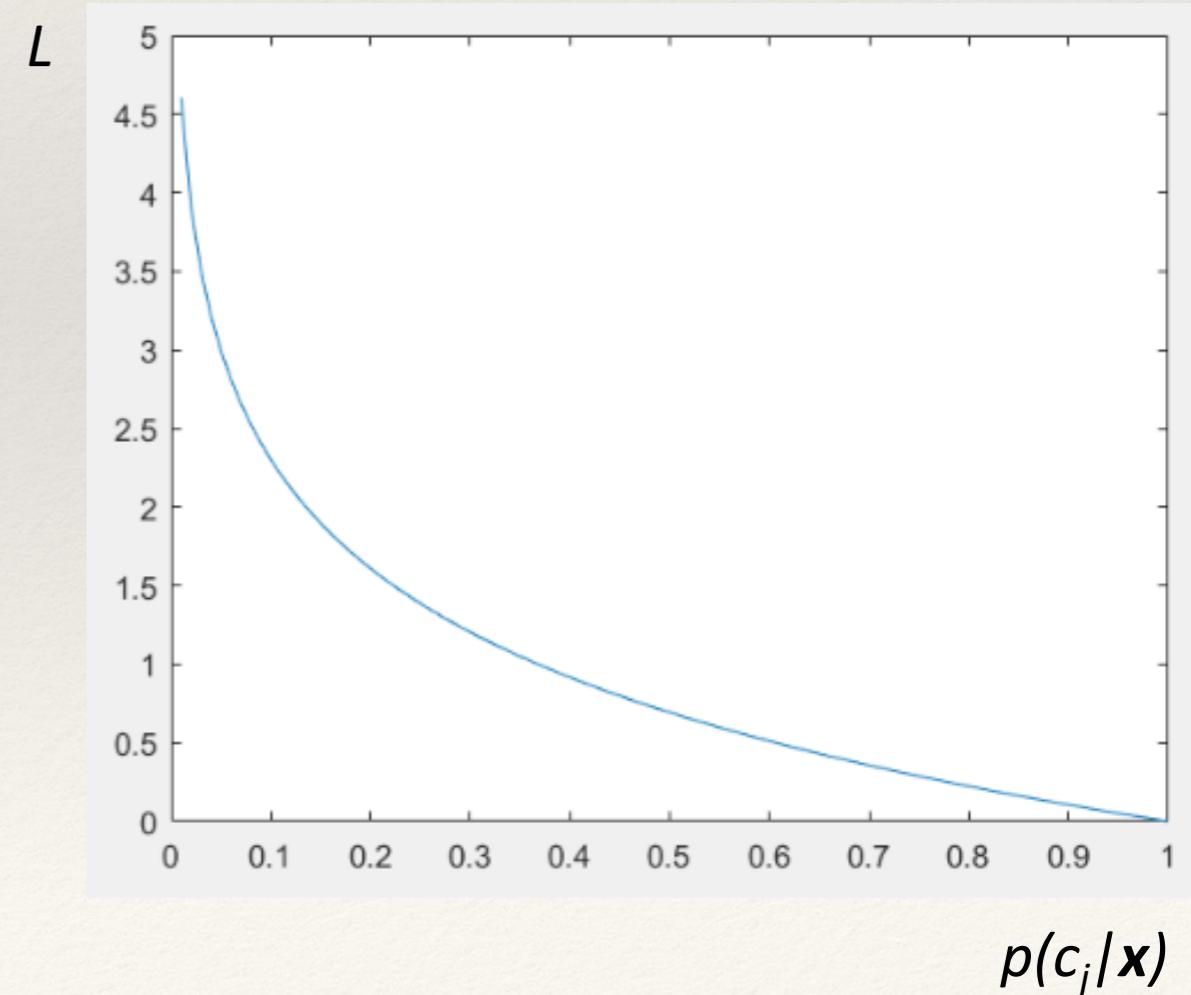
$$p(c_k=1|\mathbf{x}) = \frac{e^{o_k}}{\sum_{j=1}^C e^{o_j}}$$

# Cross-entropy loss function

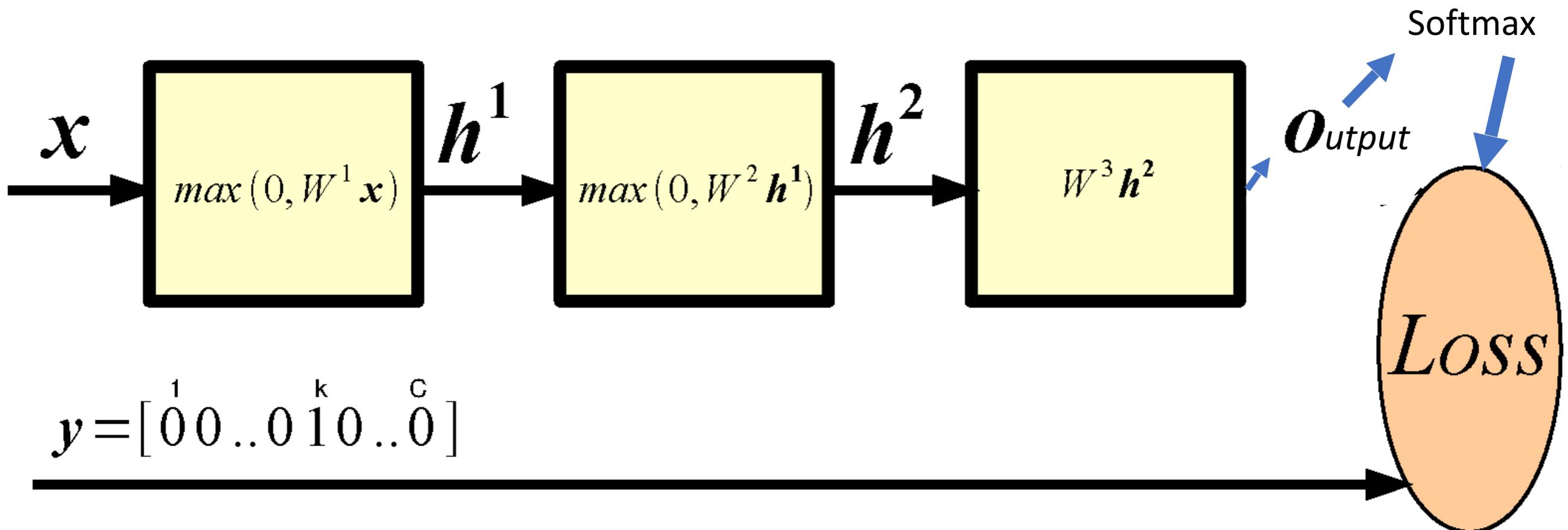
- Negative log-likelihood

$$L(\mathbf{x}, y; \theta) = - \sum_j y_j \log p(c_j | \mathbf{x})$$

- Minimizing this is equivalent to minimizing KL-divergence on predicted and target probability distributions
- Is it a good loss?
  - Differentiable
  - Cost decreases as probability increases



# How Good is a Network?



Probability of class  $k$  given input (softmax):

$$p(c_k=1|x) = \frac{e^{o_k}}{\sum_{j=1}^C e^{o_j}}$$

(Per-sample) **Loss**; e.g., negative log-likelihood (good for classification of small number of classes):

$$L(x, y; \theta) = -\sum_j y_j \log p(c_j|x)$$

# Training

**Learning** consists of minimizing the loss (plus some regularization term) w.r.t. parameters over the whole training set.

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{n=1}^P L(\mathbf{x}^n, y^n; \boldsymbol{\theta})$$

# Training

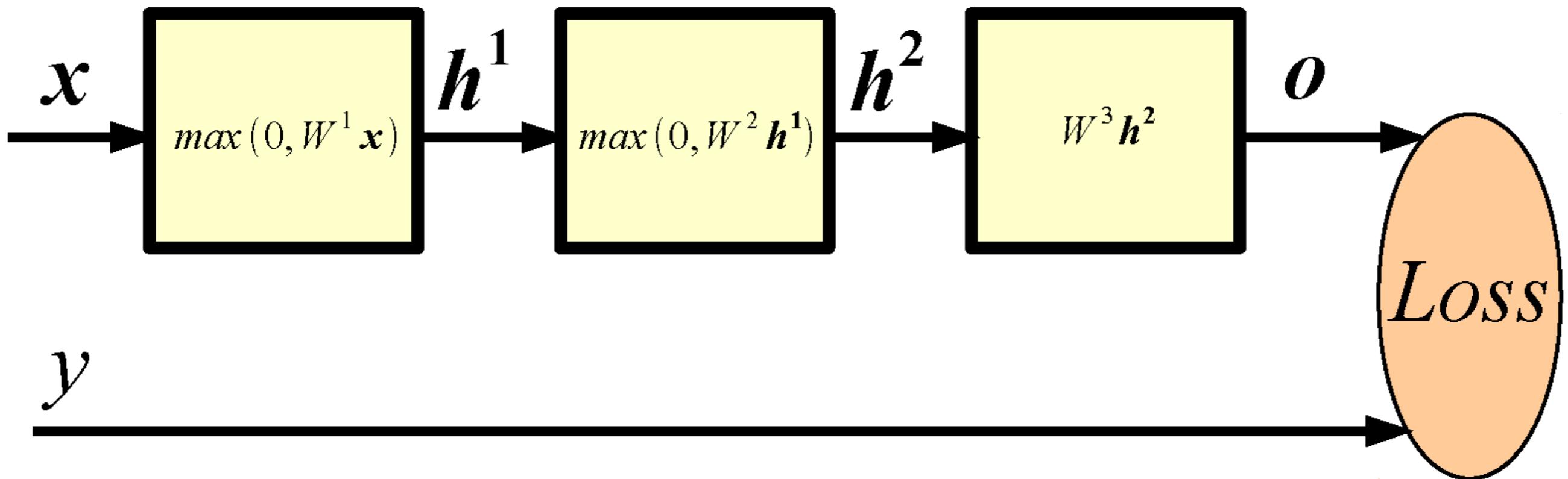
**Learning** consists of minimizing the loss (plus some regularization term) w.r.t. parameters over the whole training set.

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{n=1}^P L(\mathbf{x}^n, y^n; \boldsymbol{\theta})$$

**Question:** How to minimize a complicated function of the parameters?

**Answer:** Chain rule, a.k.a. **Backpropagation**! That is the procedure to compute gradients of the loss w.r.t. parameters in a multi-layer neural network.

# Key Idea: Wiggle To Decrease Loss

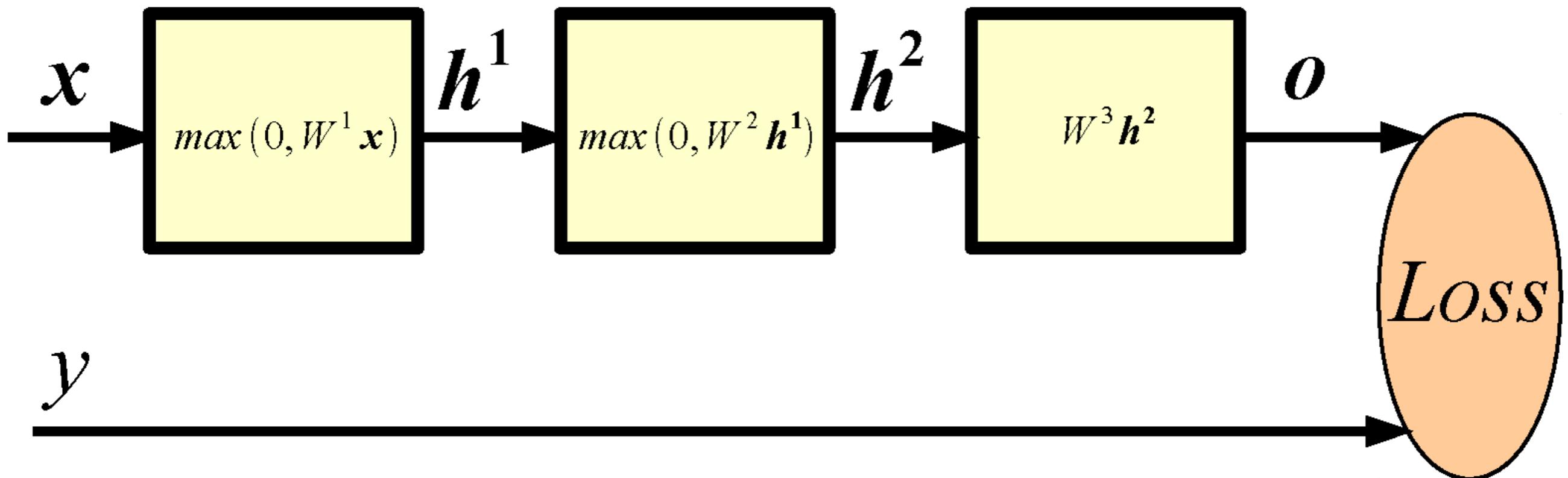


Let's say we want to decrease the loss by adjusting  $W_{i,j}^1$ .  
We could consider a very small  $\epsilon = 1e-6$  and compute:

$$L(\mathbf{x}, y; \boldsymbol{\theta})$$

$$L(\mathbf{x}, y; \boldsymbol{\theta} \setminus W_{i,j}^1, W_{i,j}^1 + \epsilon)$$

# Key Idea: Wiggle To Decrease Loss



Let's say we want to decrease the loss by adjusting  $W_{i,j}^1$ .  
We could consider a very small  $\epsilon = 1e-6$  and compute:

$$L(\mathbf{x}, y; \boldsymbol{\theta})$$

$$L(\mathbf{x}, y; \boldsymbol{\theta} \setminus W_{i,j}^1, W_{i,j}^1 + \epsilon)$$

Then, update:

$$W_{i,j}^1 \leftarrow W_{i,j}^1 + \epsilon \operatorname{sgn}(L(\mathbf{x}, y; \boldsymbol{\theta}) - L(\mathbf{x}, y; \boldsymbol{\theta} \setminus W_{i,j}^1, W_{i,j}^1 + \epsilon))$$

# Derivative w.r.t. Input of Softmax

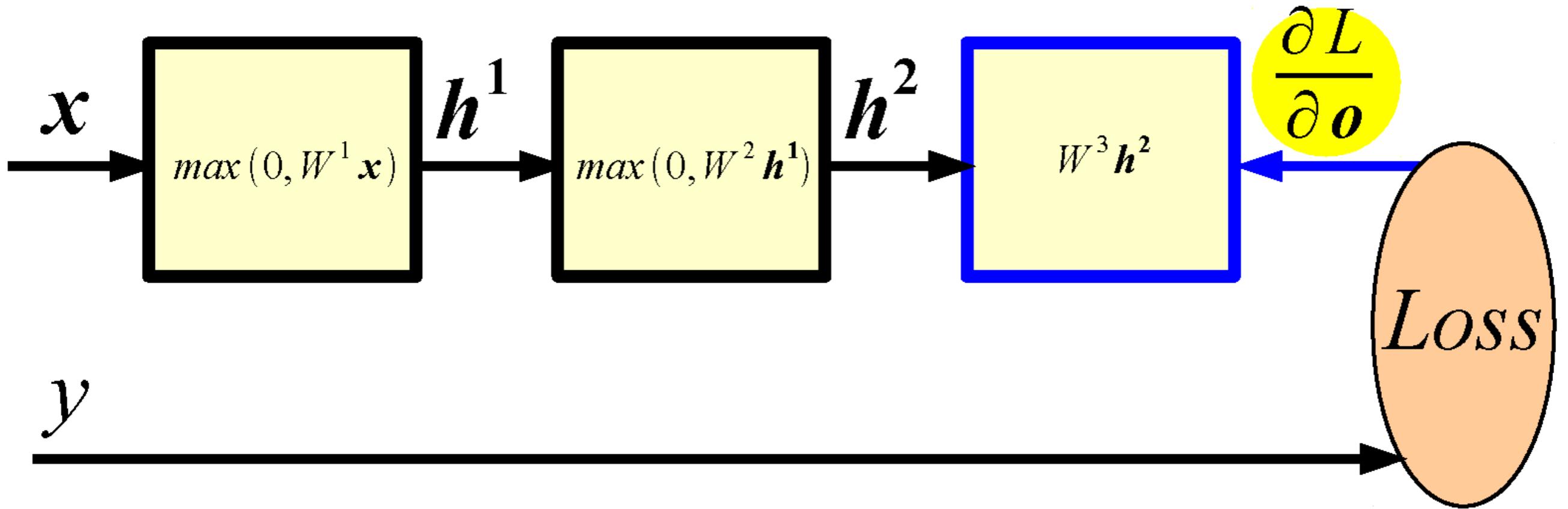
$$p(c_k=1|x) = \frac{e^{o_k}}{\sum_j e^{o_j}}$$

$$L(x, y; \theta) = -\sum_j y_j \log p(c_j|x) \quad y = [0^1 0^0 .. 0^k 1^0 .. 0^c]$$

By substituting the first formula in the second, and taking the derivative w.r.t.  $\theta$  we get:

$$\frac{\partial L}{\partial \theta} = p(c|x) - y$$

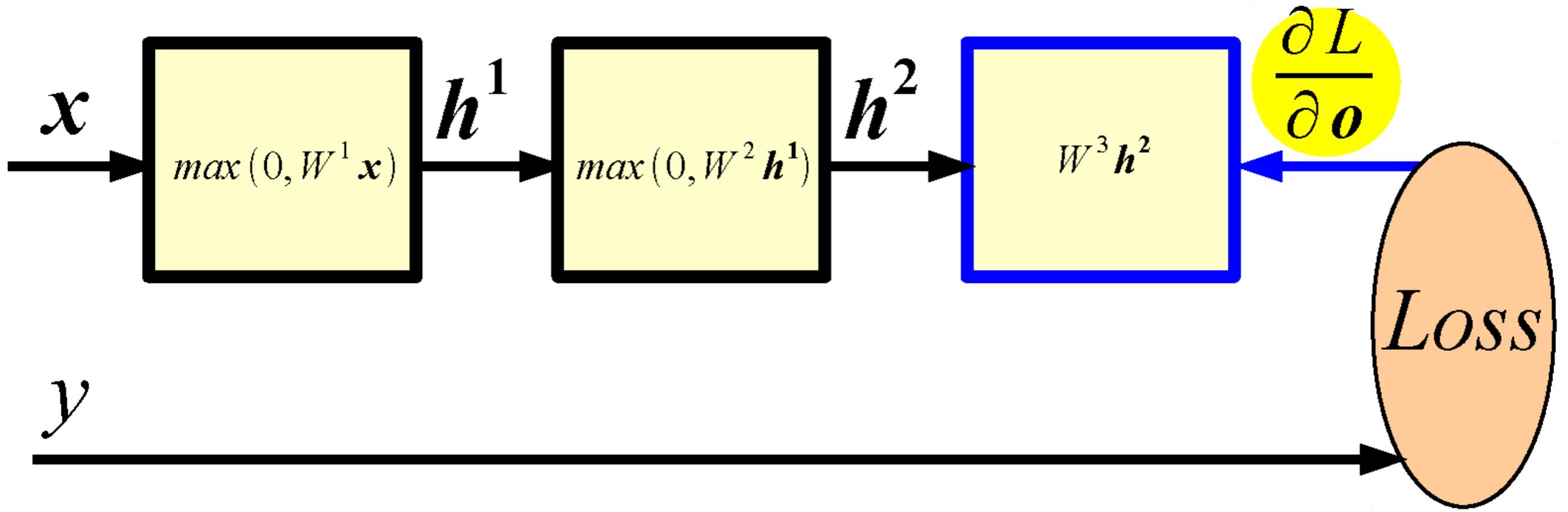
# Backward Propagation



Given  $\frac{\partial L}{\partial o}$  and assuming we can easily compute the Jacobian of each module, we have:

$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial o} \frac{\partial o}{\partial W^3}$$

# Backward Propagation

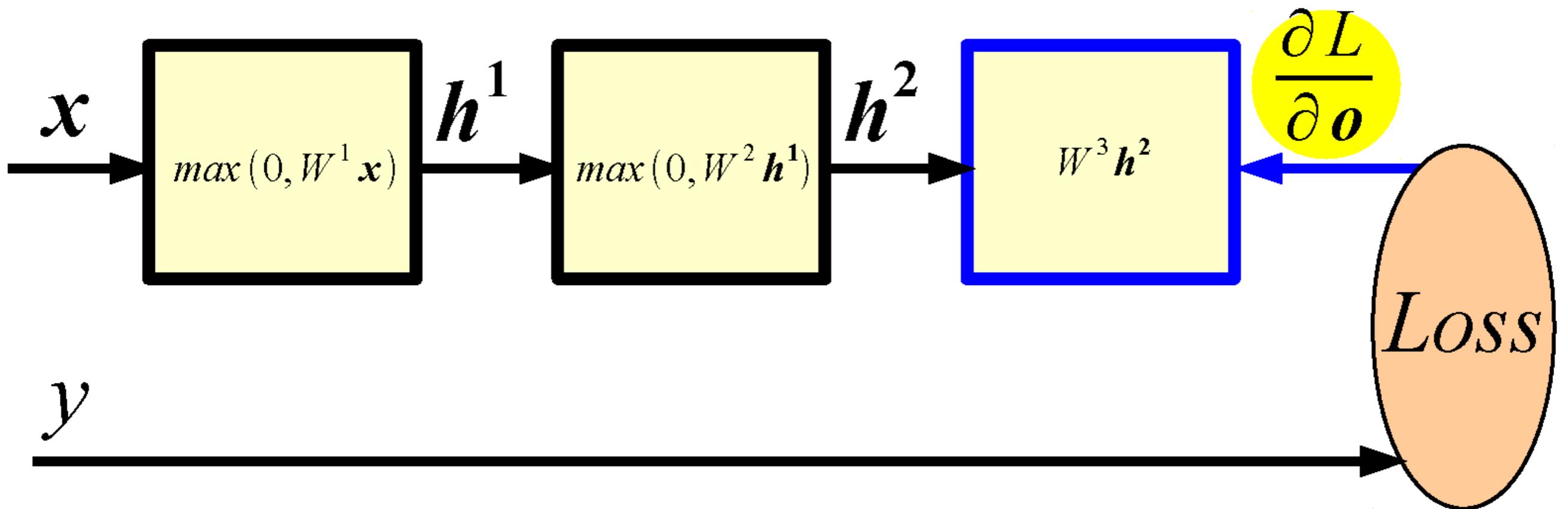


Given  $\frac{\partial L}{\partial \mathbf{o}}$  and assuming we can easily compute the Jacobian of each module, we have:

$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial W^3}$$

$$\frac{\partial L}{\partial h^2} = \frac{\partial L}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial h^2}$$

# Backward Propagation



Given  $\frac{\partial L}{\partial o}$  and assuming we can easily compute the Jacobian of each module, we have:

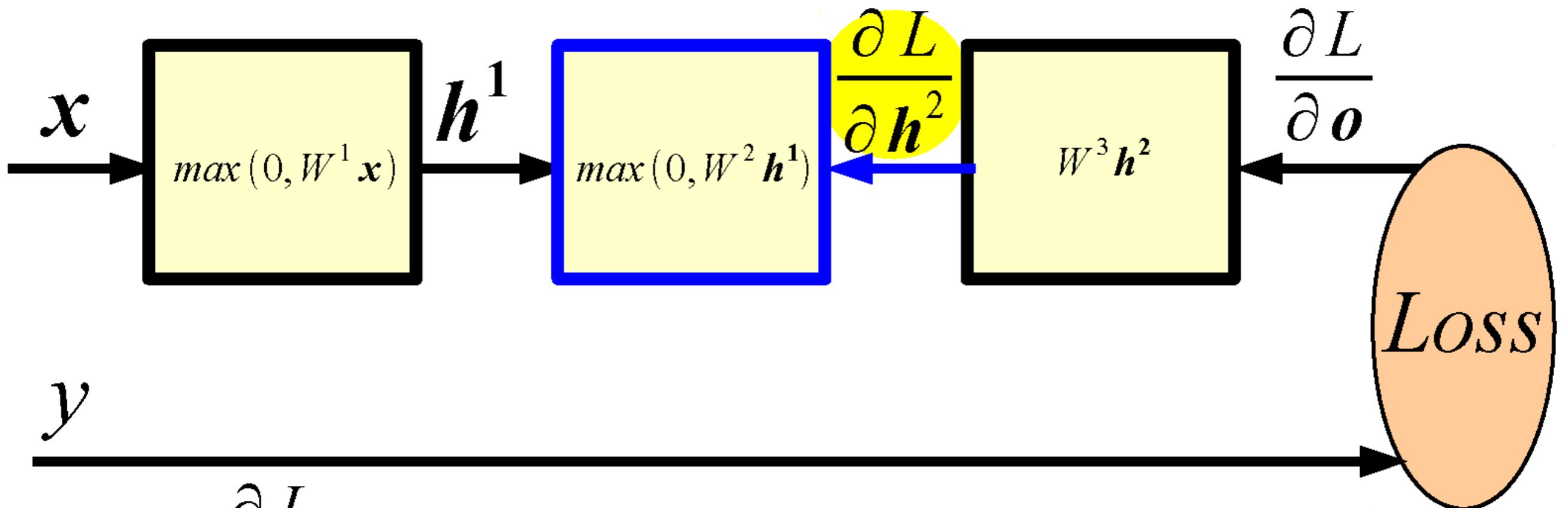
$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial o} \frac{\partial o}{\partial W^3}$$

$$\frac{\partial L}{\partial h^2} = \frac{\partial L}{\partial o} \frac{\partial o}{\partial h^2}$$

$$\frac{\partial L}{\partial W^3} = (p(c|x) - y) h^{2T}$$

$$\frac{\partial L}{\partial h^2} = W^{3T} (p(c|x) - y)$$

# Backward Propagation

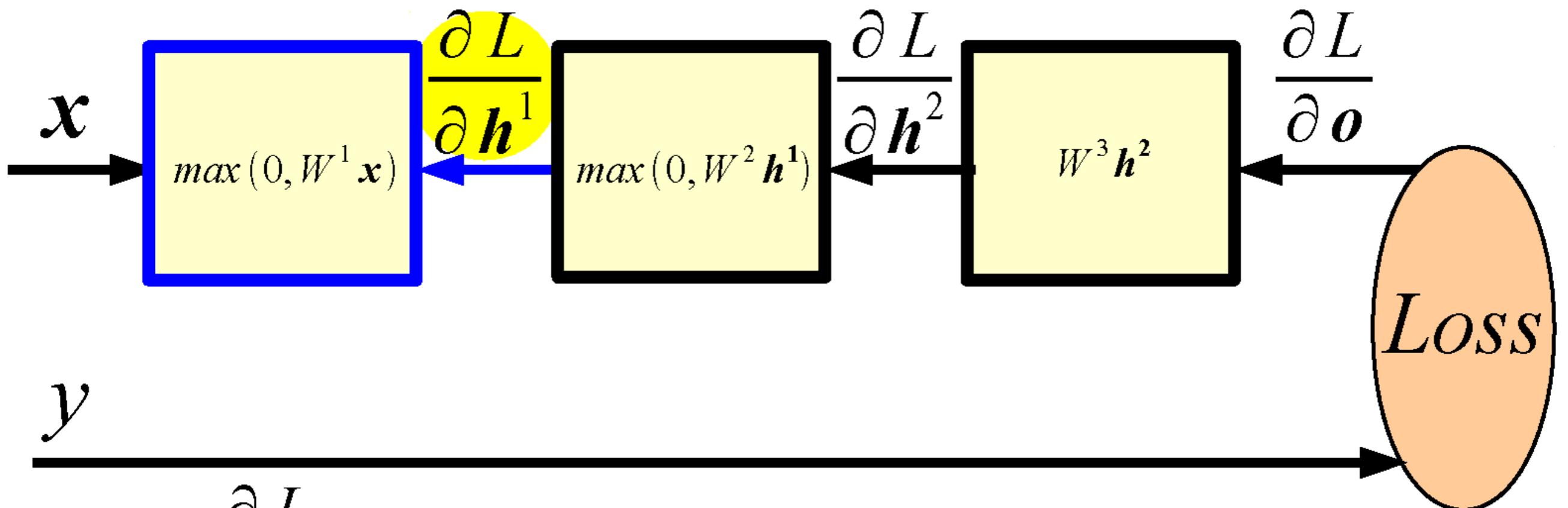


Given  $\frac{\partial L}{\partial h^2}$  we can compute now:

$$\frac{\partial L}{\partial W^2} = \frac{\partial L}{\partial h^2} \frac{\partial h^2}{\partial W^2}$$

$$\frac{\partial L}{\partial h^1} = \frac{\partial L}{\partial h^2} \frac{\partial h^2}{\partial h^1}$$

# Backward Propagation



Given  $\frac{\partial L}{\partial h^1}$  we can compute now:

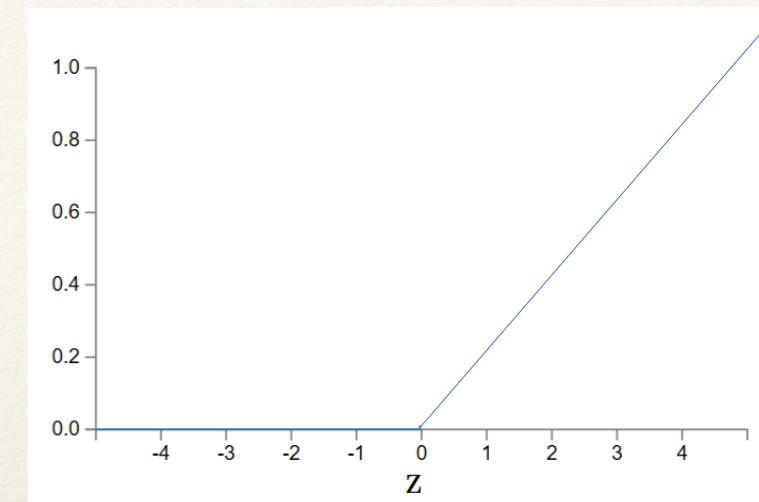
$$\frac{\partial L}{\partial W^1} = \frac{\partial L}{\partial h^1} \frac{\partial h^1}{\partial W^1}$$

# Backward Propagation

**Question:** Does BPROP work with ReLU layers only?

**Answer:** Nope, any a.e. differentiable transformation works.

*But the ReLU is not differentiable at 0!*



- '0' is the best place for this to occur, because we don't care about the result (it is no activation).
- 'Dead' perceptrons
- ReLU has unbounded positive response:
  - Potential faster convergence / overstep

# Backward Propagation

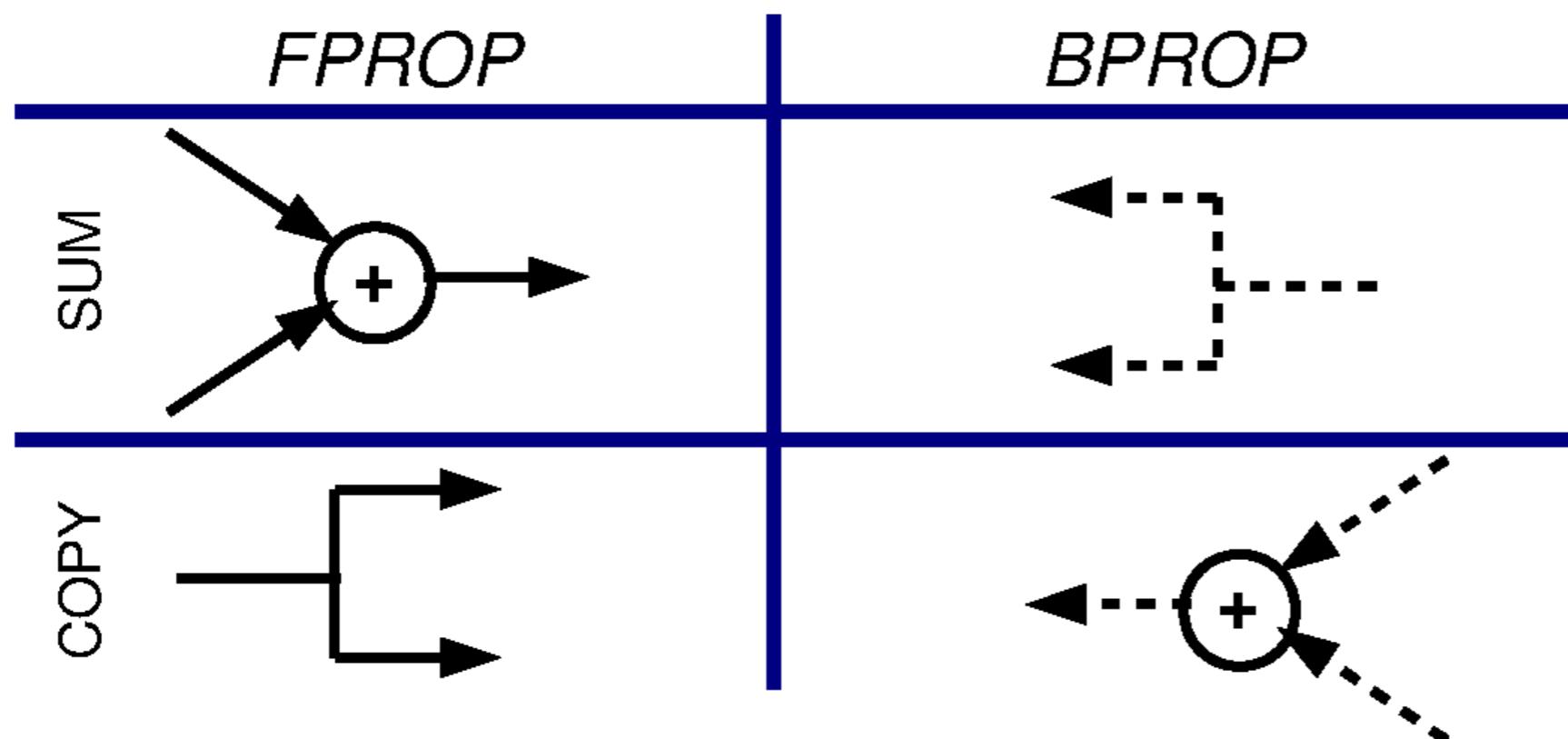
**Question:** Does BPROP work with ReLU layers only?

**Answer:** Nope, any a.e. differentiable transformation works.

**Question:** What's the computational cost of BPROP?

**Answer:** About twice FPROP (need to compute gradients w.r.t. input and parameters at every layer).

**Note:** FPROP and BPROP are dual of each other. E.g.,:



# Toy Code (Matlab): Neural Net Trainer

```
% F-PROP
for i = 1 : nr_layers - 1
    [h{i} jac{i}] = nonlinearity(W{i} * h{i-1} + b{i});
end
h{nr_layers-1} = W{nr_layers-1} * h{nr_layers-2} + b{nr_layers-1};
prediction = softmax(h{l-1});

% CROSS ENTROPY LOSS
loss = - sum(sum(log(prediction) .* target)) / batch_size;

% B-PROP
dh{l-1} = prediction - target;
for i = nr_layers - 1 : -1 : 1
    Wgrad{i} = dh{i} * h{i-1}';
    bgrad{i} = sum(dh{i}, 2);
    dh{i-1} = (W{i}' * dh{i}) .* jac{i-1};
end

% UPDATE
for i = 1 : nr_layers - 1
    W{i} = W{i} - (lr / batch_size) * Wgrad{i};
    b{i} = b{i} - (lr / batch_size) * bgrad{i};
end
```

# Stochastic Gradient Descent

- Dataset can be too large to strictly apply gradient descent.
- Instead, randomly sample a data point, perform gradient descent per point, and iterate.
  - True gradient is approximated only
  - Picking a subset of points: “*mini-batch*”

Pick starting  $W$  and learning rate  $\gamma$

While not at minimum:

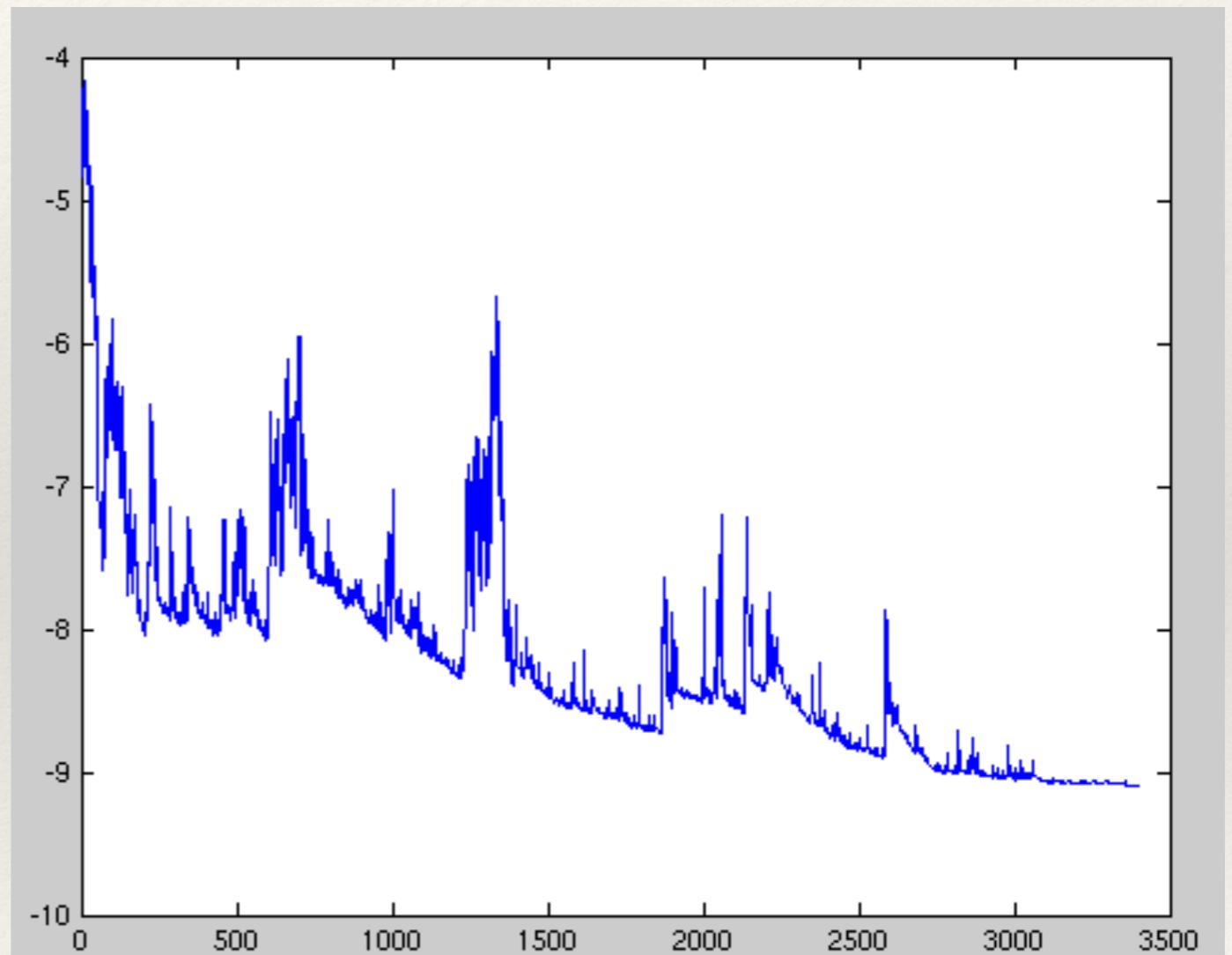
- Shuffle training set
- For each data point  $i=1\dots n$  (*maybe as mini-batch*)
  - *Gradient descent*

“Epoch”

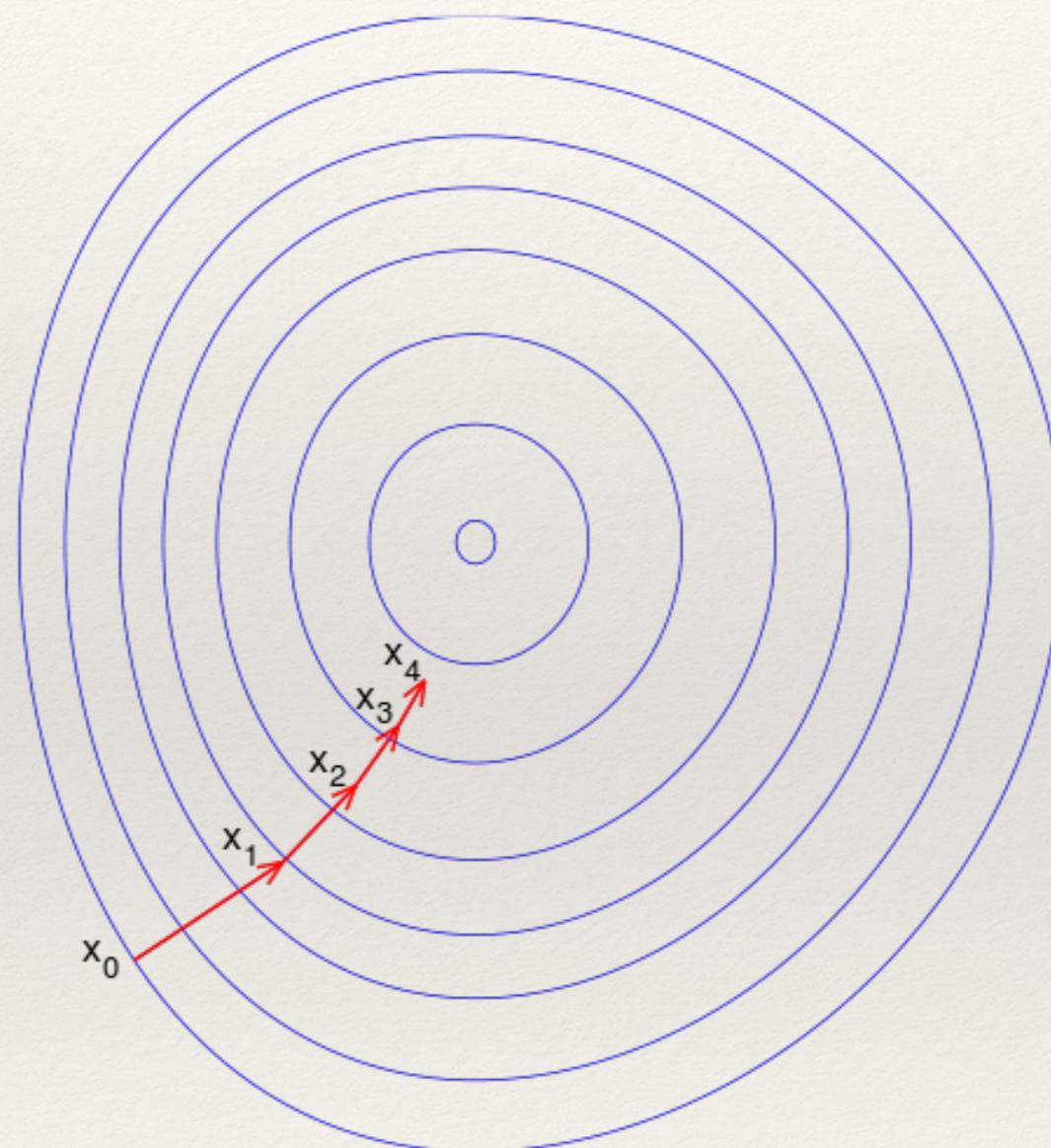
# Stochastic Gradient Descent

Loss will not always decrease (locally) as training data point is random.

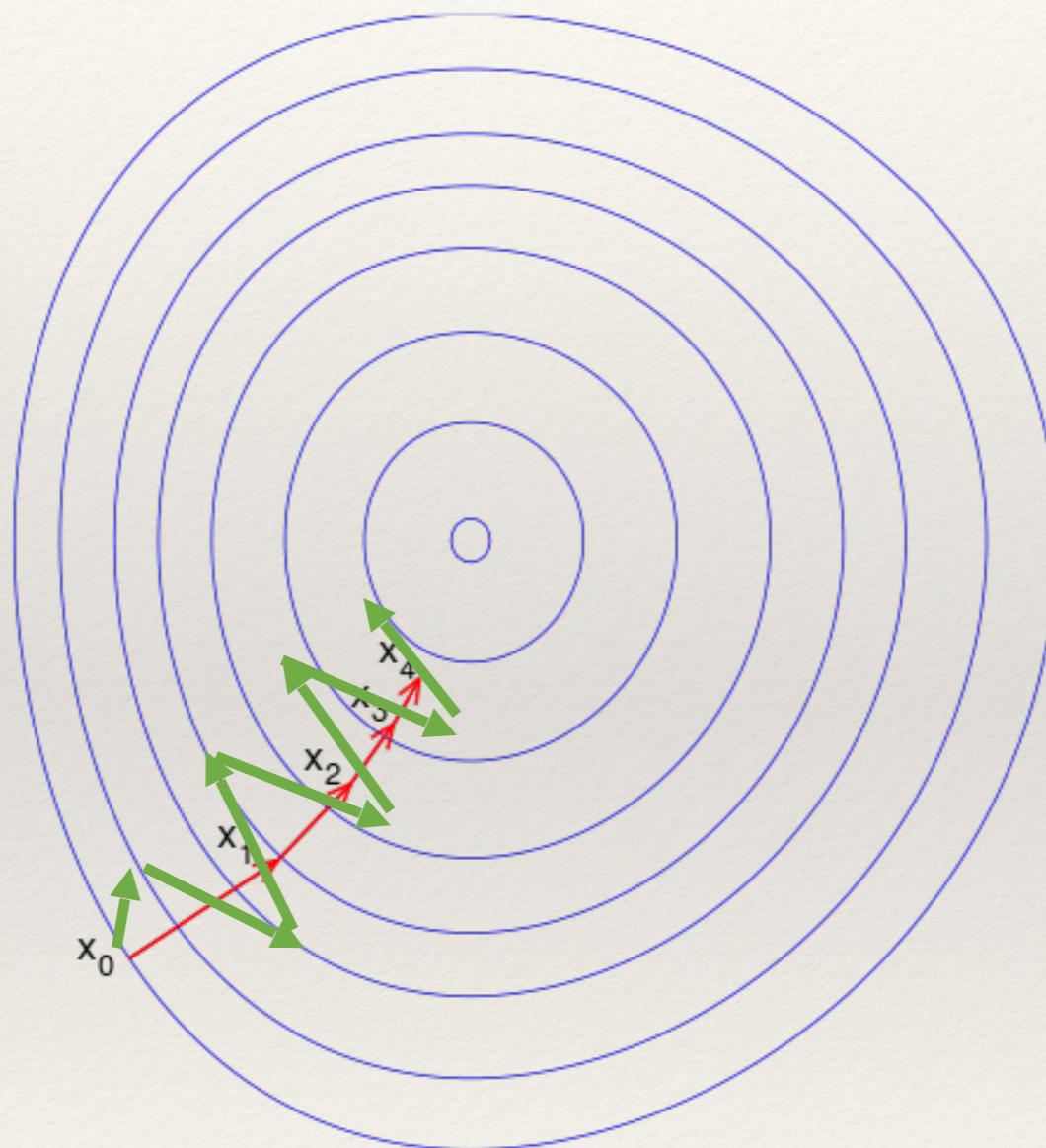
Still converges over time.



# Gradient descent oscillations



# Gradient descent oscillations



Slow to  
converge to the  
(local)  
optimum

# Momentum

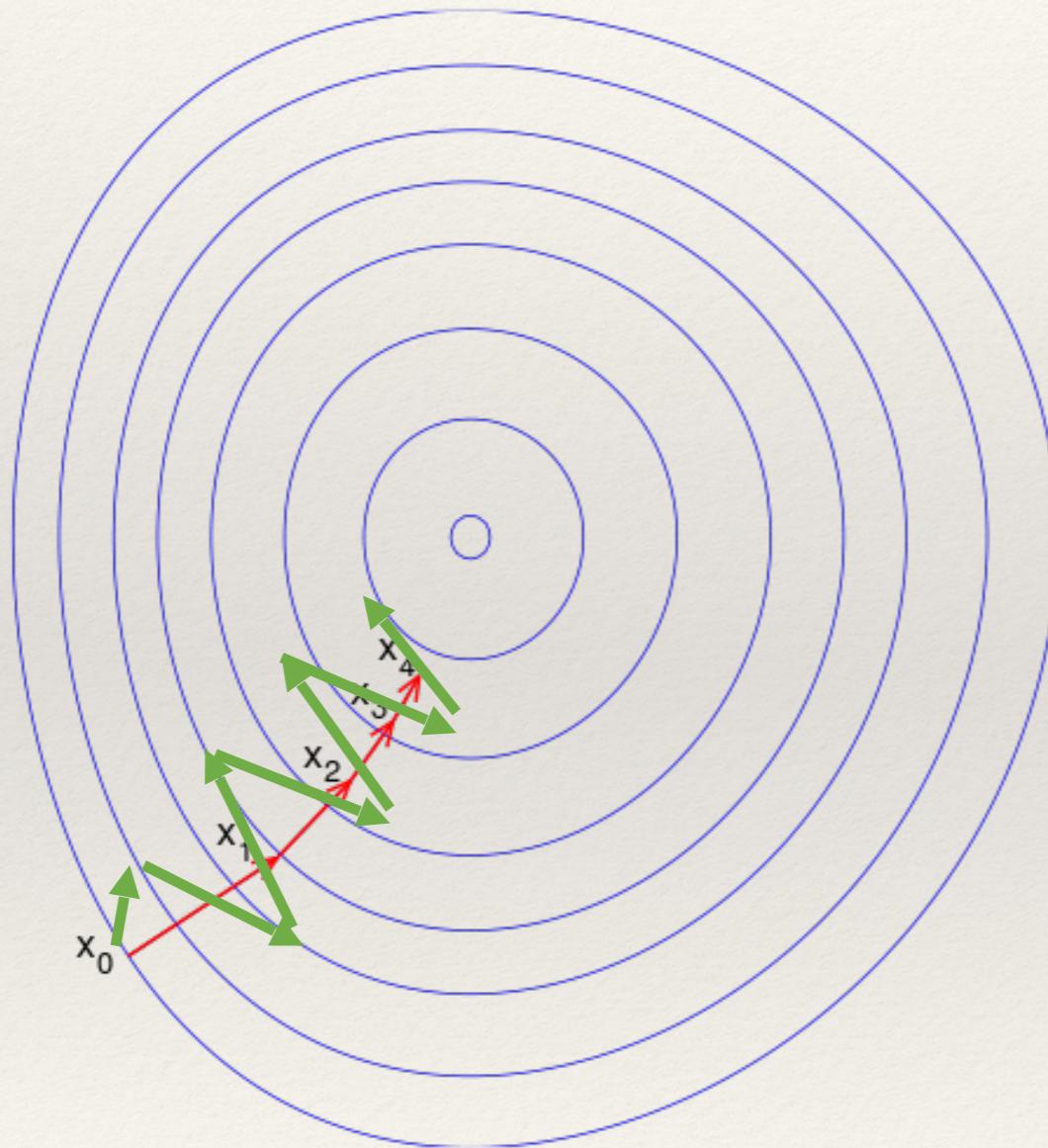
- Adjust the gradient by a weighted sum of the previous amount plus the current amount.

- Without momentum:  $\theta_{t+1} = \theta_t - \gamma \frac{\partial L}{\partial \theta}$

- With momentum (new  $\alpha$  parameter):

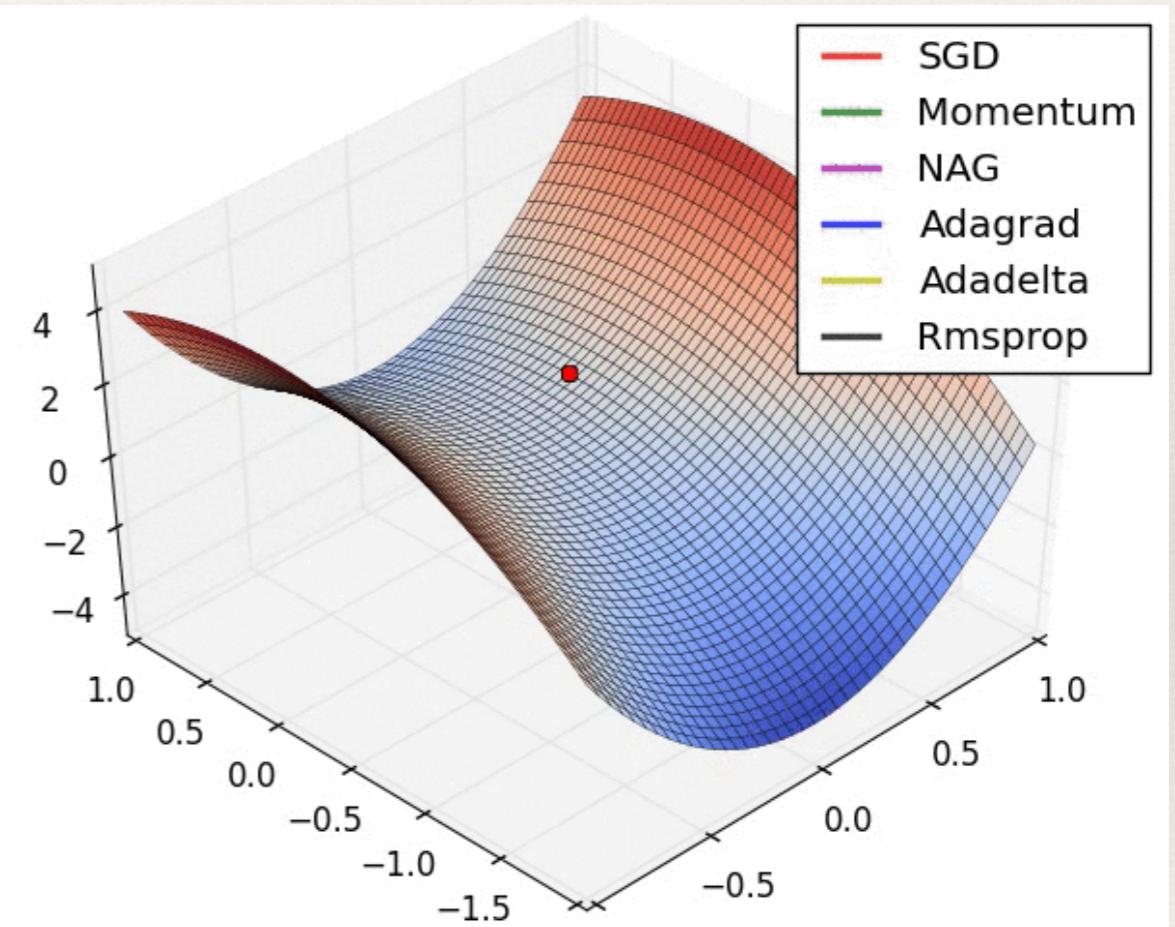
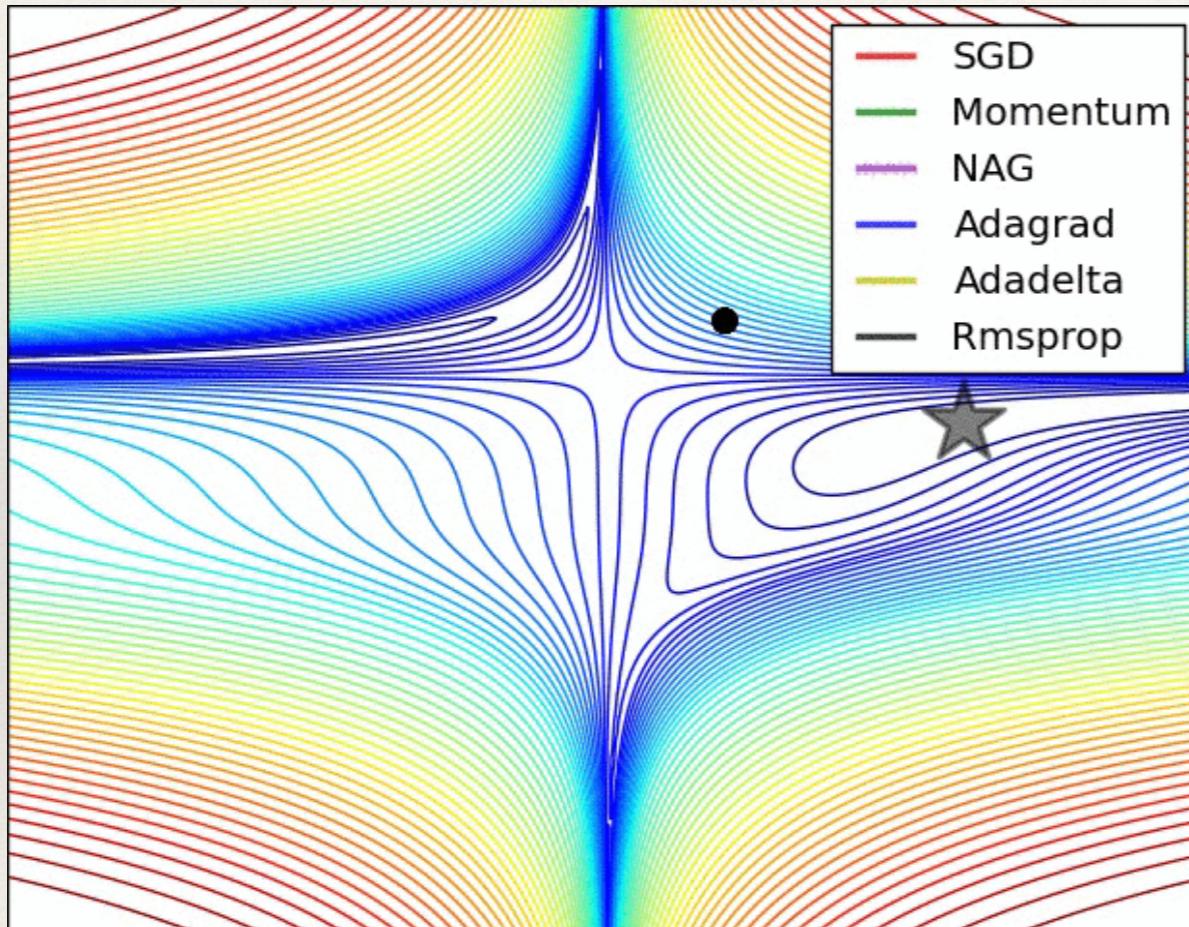
$$\theta_{t+1} = \theta_t - \gamma \left( \alpha \left[ \frac{\partial L}{\partial \theta} \right]_{t-1} + \left[ \frac{\partial L}{\partial \theta} \right]_t \right)$$

# Lowering the learning rate = smaller steps in SGD



- Less ‘ping pong’
- Takes longer to get to the optimum

# Flat regions in energy landscape



---

# Problem of fitting

---

- Too many parameters = overfitting
- Not enough parameters = underfitting
- More data = less chance to overfit
- How do we know what is required?

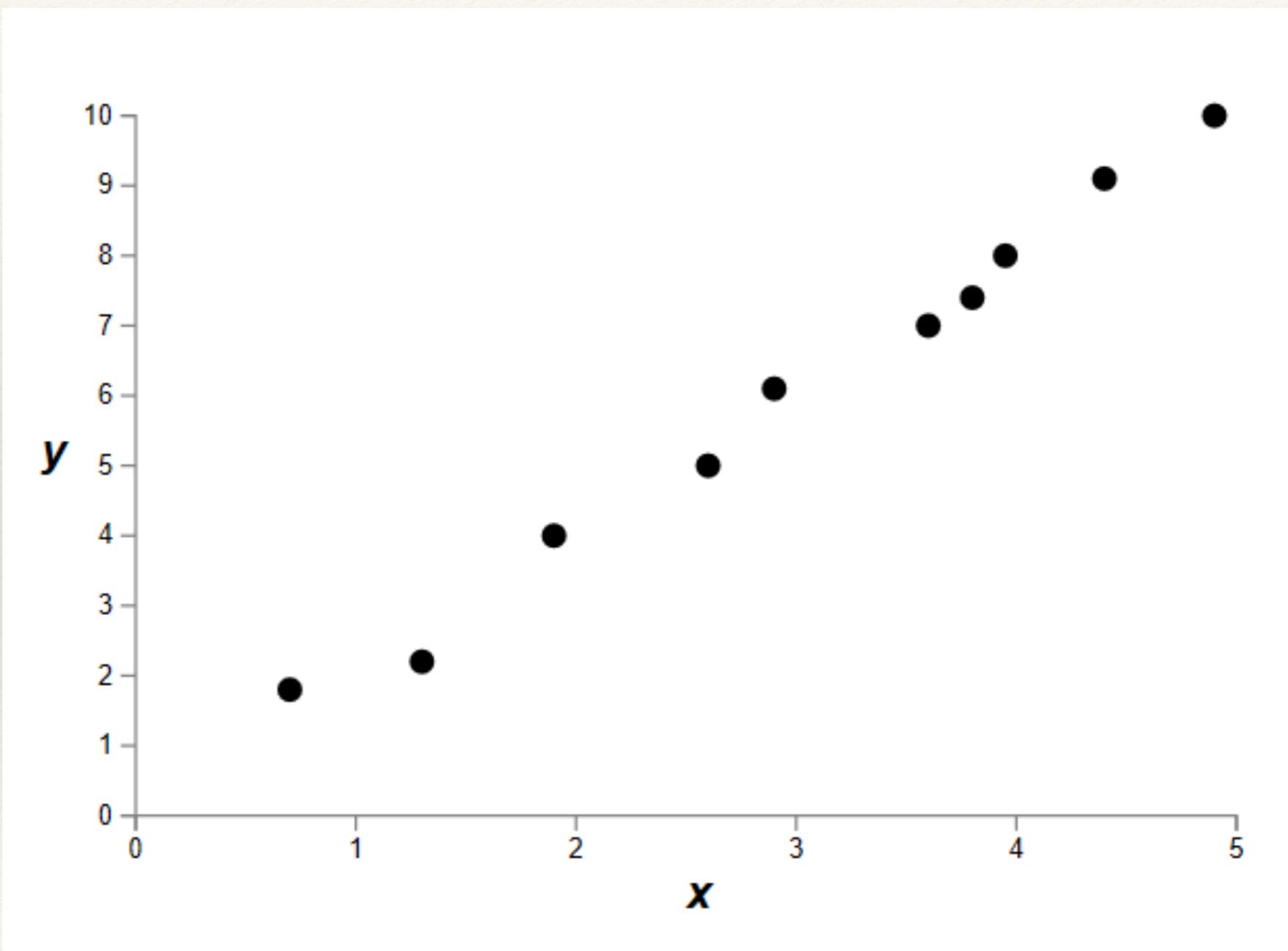
---

# Regularization

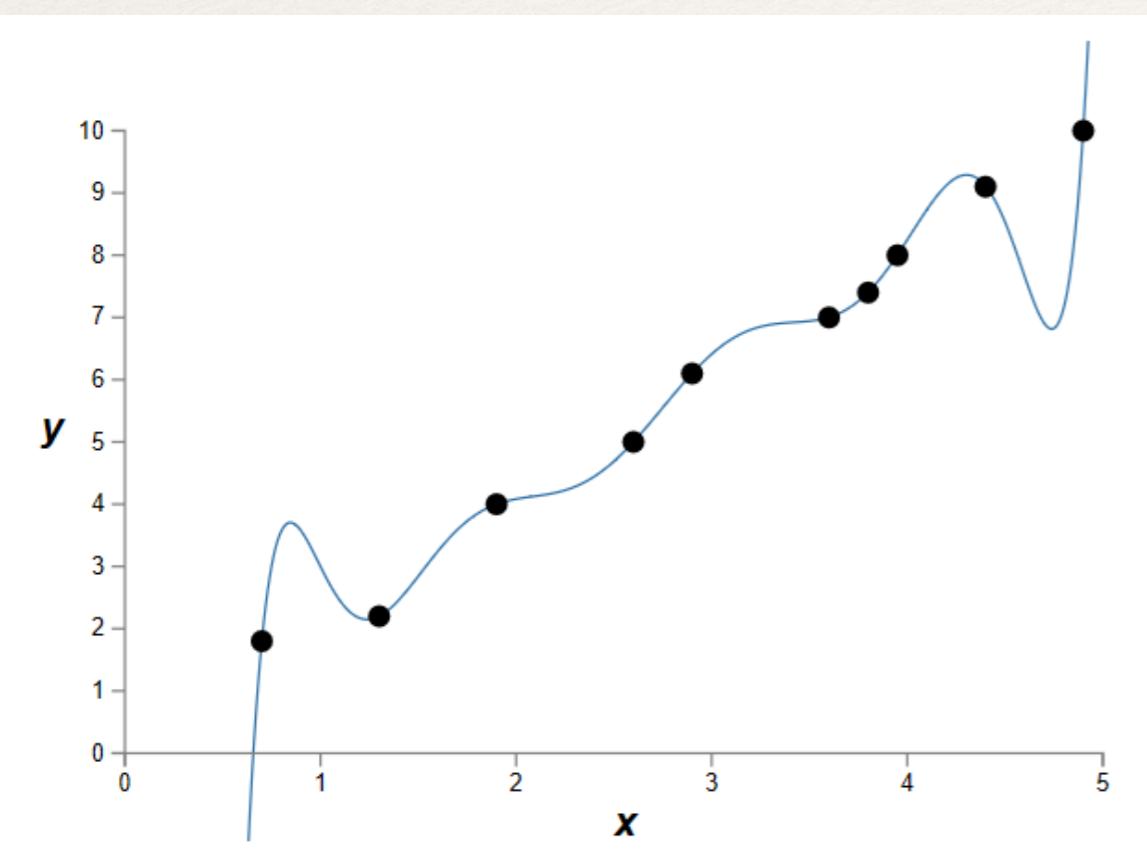
---

- Attempt to guide solution to *not overfit*
- But still give freedom with many parameters

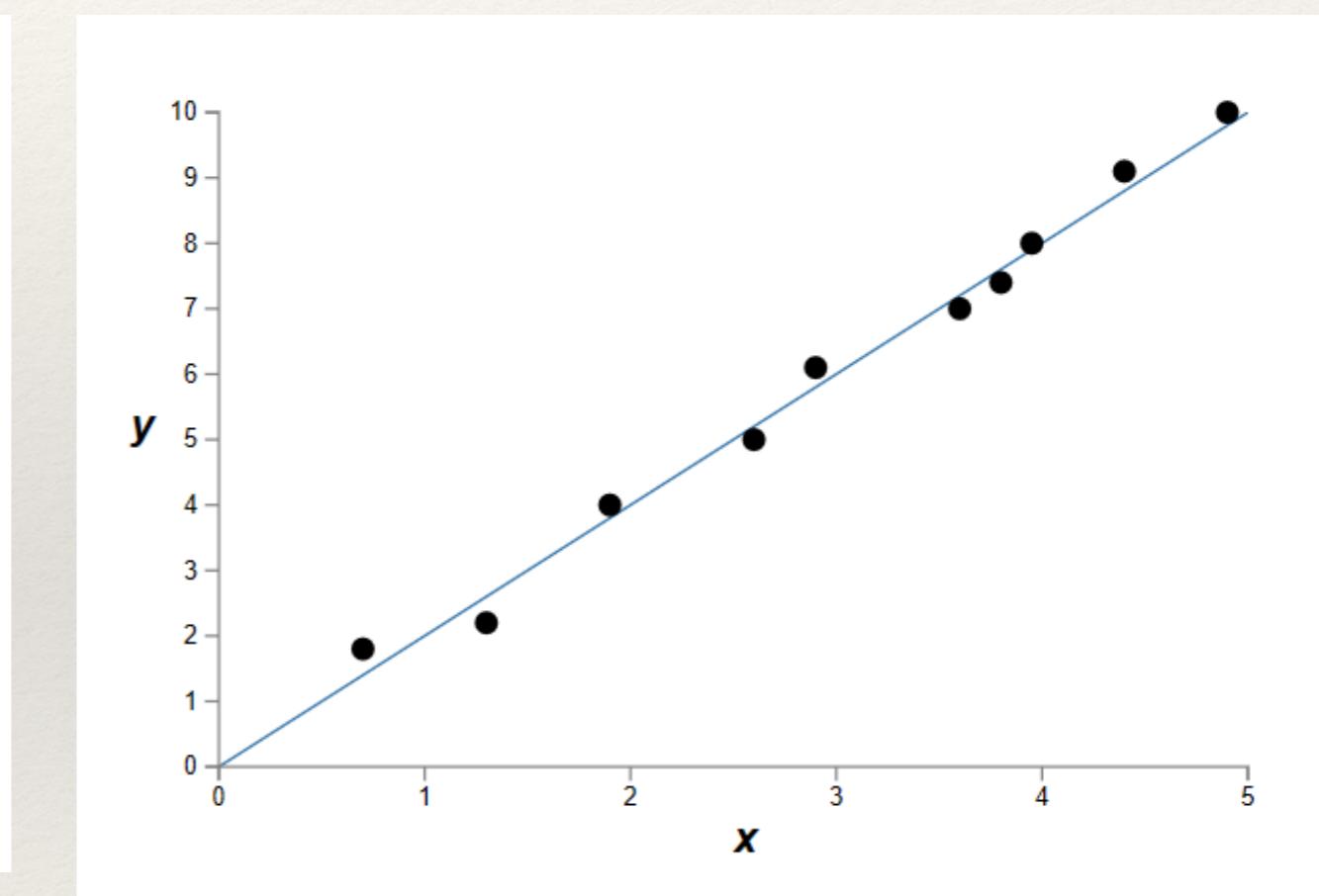
# Data fitting problem



Which is better?  
Which is better a priori?



9<sup>th</sup> order polynomial



1<sup>st</sup> order polynomial

---

# Regularization

---

- Attempt to guide solution to *not overfit*
- But still give freedom with many parameters
- Idea:  
*Penalize the use of parameters to prefer small weights.*

# Regularization:

---

- Idea: add a cost to having high weights
- $\lambda$  = regularization parameter

$$C = C_0 + \lambda \sum_w w^2,$$

# Both can describe the data...

---

- ...but one is simpler.
- Occam's razor:  
*"Among competing hypotheses, the one with the fewest assumptions should be selected"*

For us:

Large weights cause large changes in behaviour in response to small changes in the input.

Simpler models (or smaller changes) are more robust to noise.

# Regularization

- Idea: add a cost to having high weights
- $\lambda$  = regularization parameter

$$C = C_0 + \lambda \sum_w w^2,$$

$$C = -\frac{1}{n} \sum_{xj} \left[ y_j \ln a_j^L + (1 - y_j) \ln(1 - a_j^L) \right] + \lambda \sum_w w^2.$$



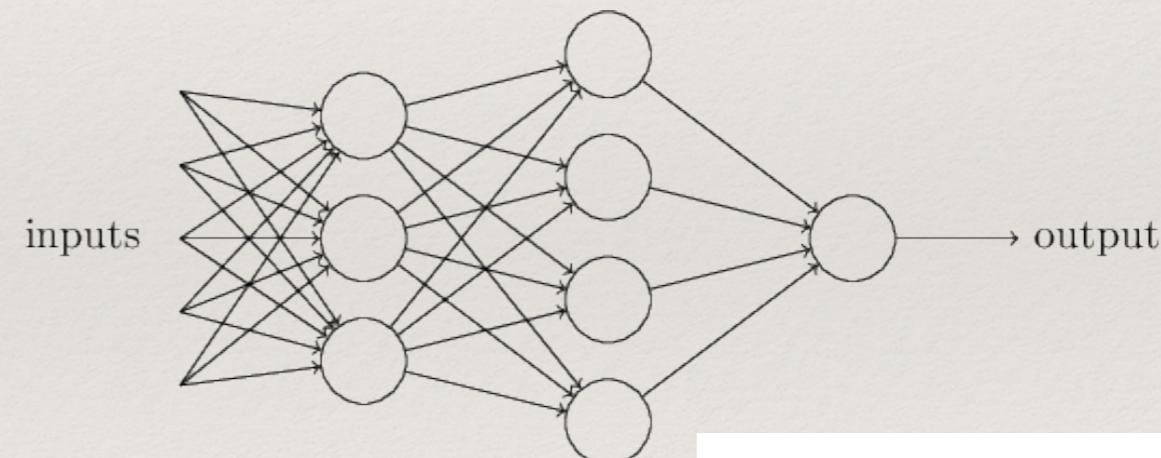
Normal cross-entropy  
loss (binary classes)

Regularization term

[Nielson]

# Regularization: Dropout

- Our networks typically start with random weights.
- Every time we train = slightly different outcome.
- Why random weights?
- If weights are all equal, response across filters will be equivalent.
  - Network doesn't train.



$$\mathbf{w} \cdot \mathbf{x} \equiv \sum_j w_j x_j;$$

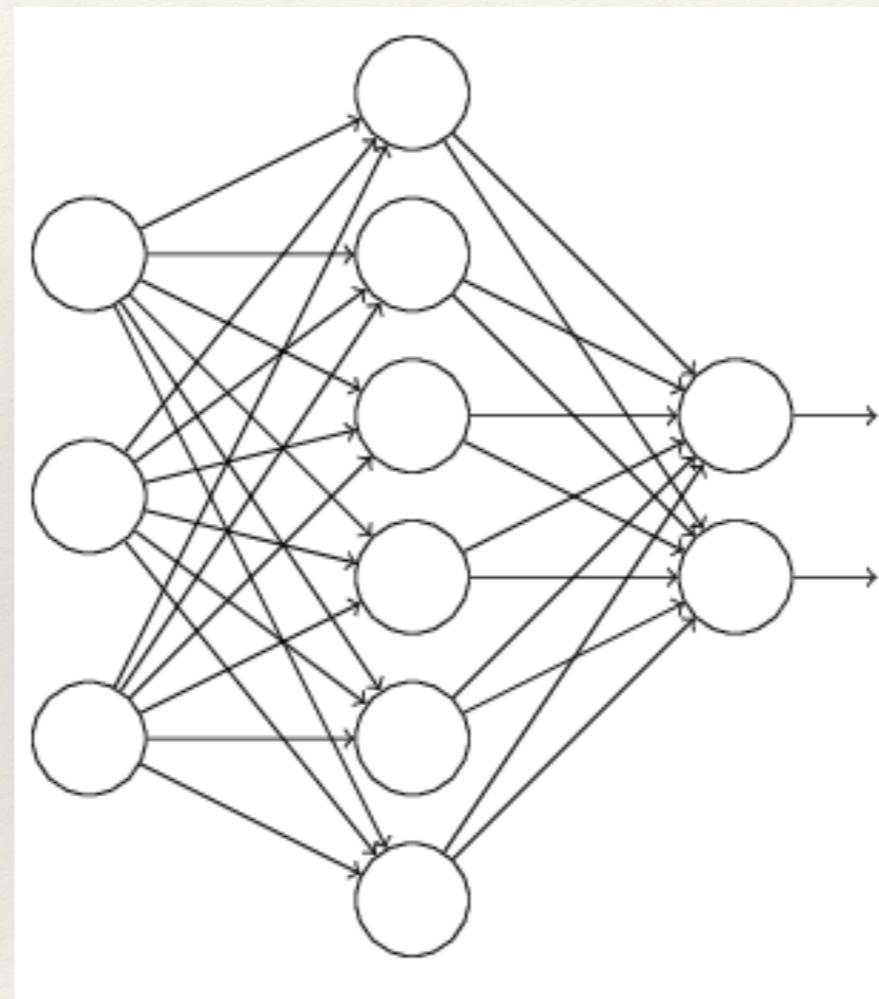
---

# Regularization

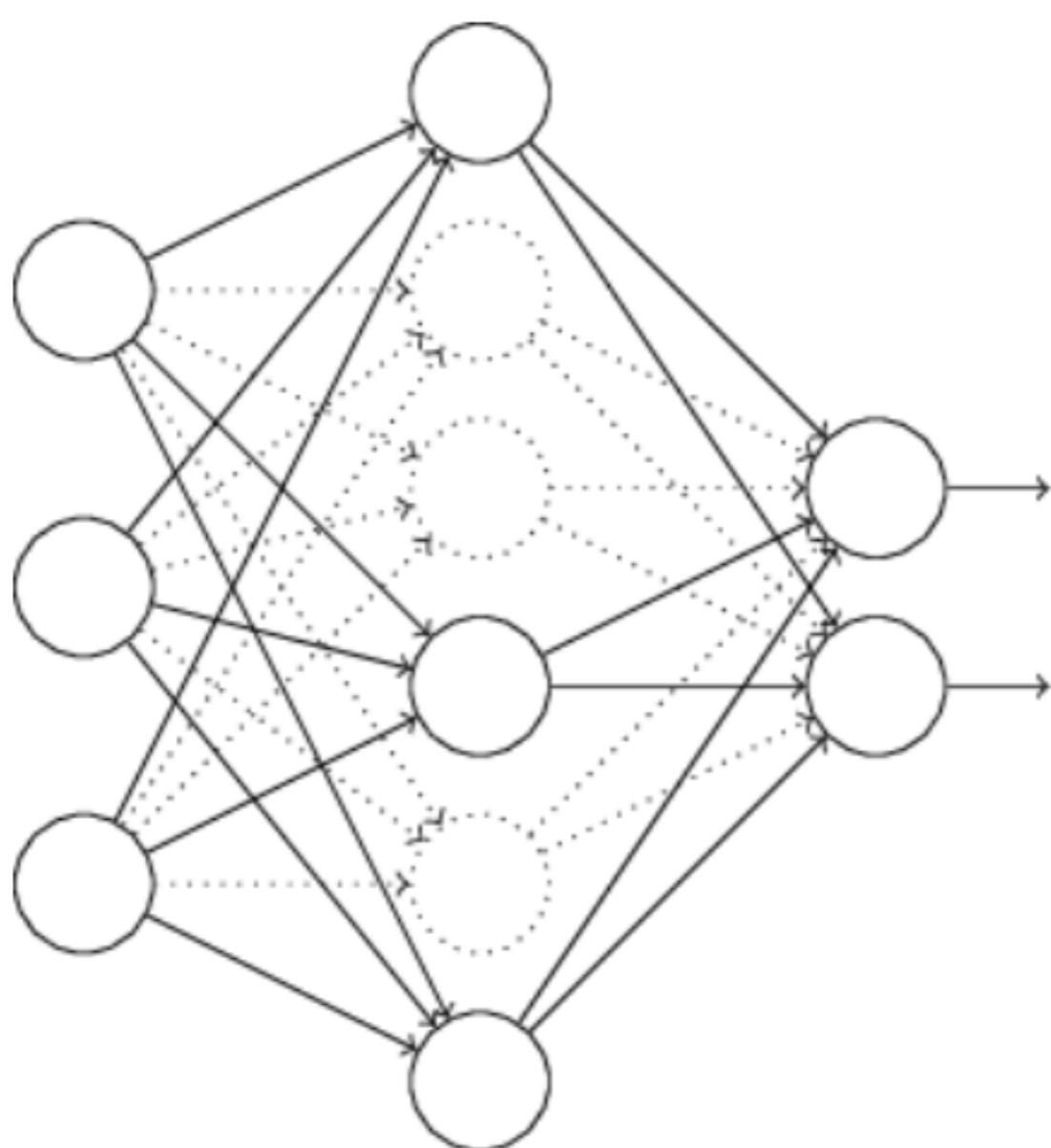
---

- Our networks typically start with random weights.
- Every time we train = slightly different outcome.
- Why not train 5 different networks with random starts and vote on their outcome?
  - Works fine!
  - Helps generalization because error is averaged.

# Regularization: Dropout



# Regularization: Dropout



At each mini-batch:

- Randomly select a subset of neurons.
- Ignore them.

On test: half weights outgoing to compensate for training on half neurons.

Effect:

- Neurons become less dependent on output of connected neurons.
- Forces network to learn more robust features that are useful to more subsets of neurons.
- Like averaging over many different trained networks with different random initializations.
- Except cheaper to train.

---

# Many forms of 'regularization'

---

- Adding more data is a kind of regularization
- Pooling is a kind of regularization
- Data augmentation is a kind of regularization

---

# Network Architecture: big space of designs

---

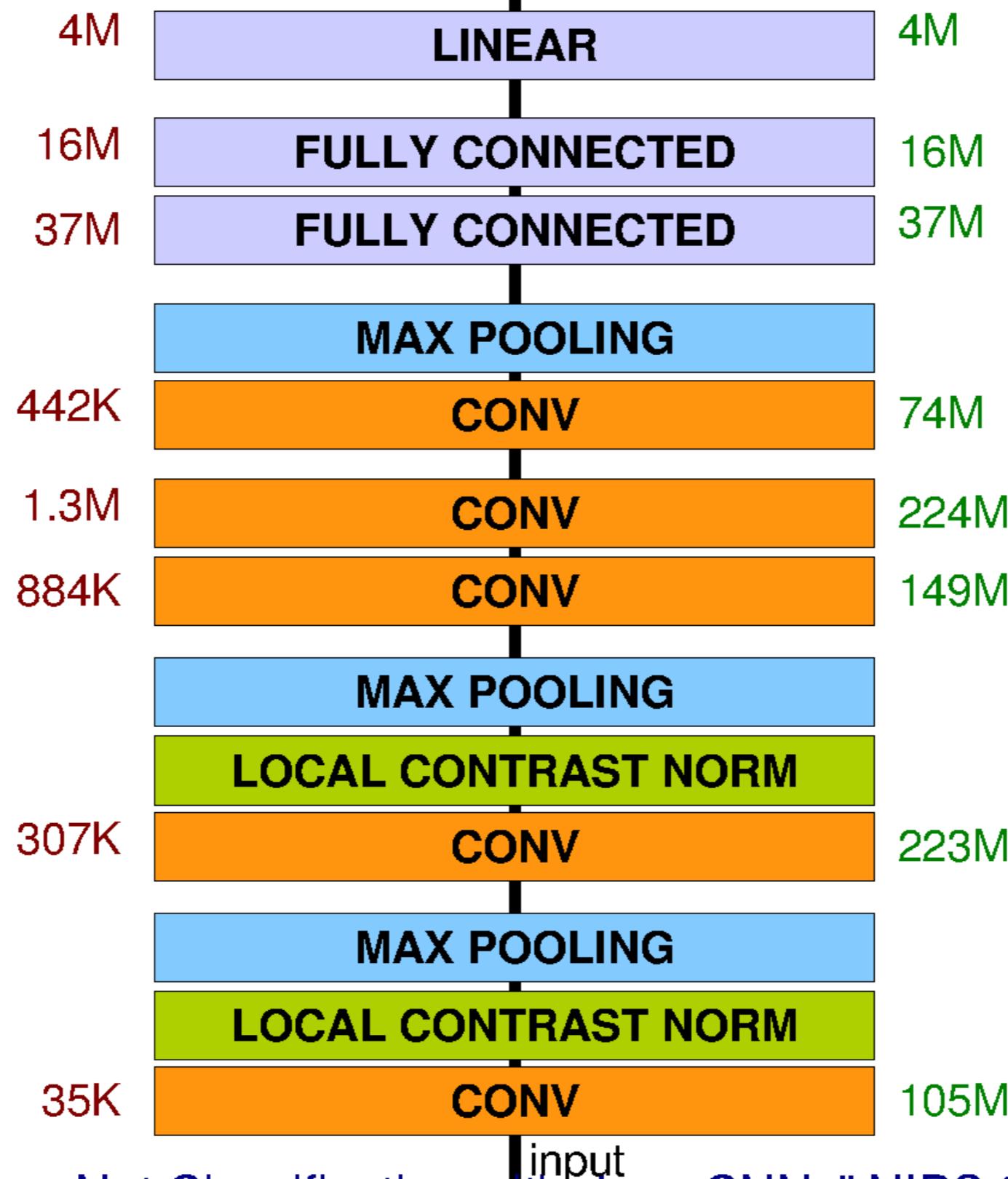
- ❖ But we still don't even know how many layers we need.

# Architecture for Classification

Total nr. params: 60M

category  
prediction

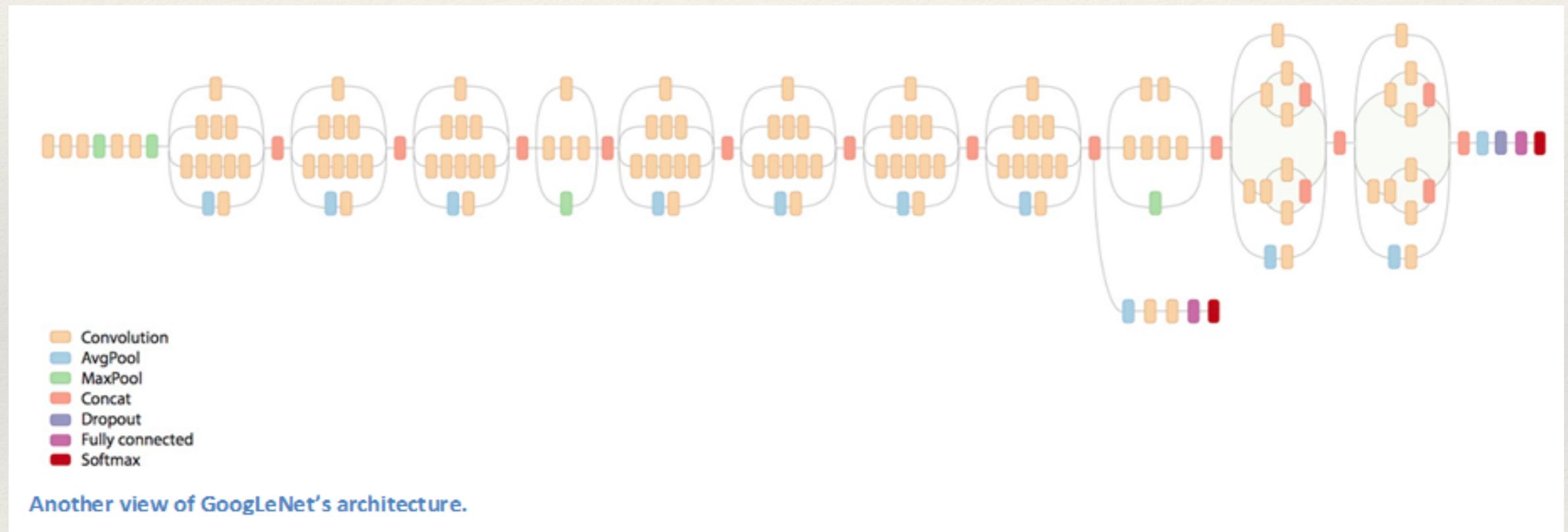
Total nr. flops: 832M



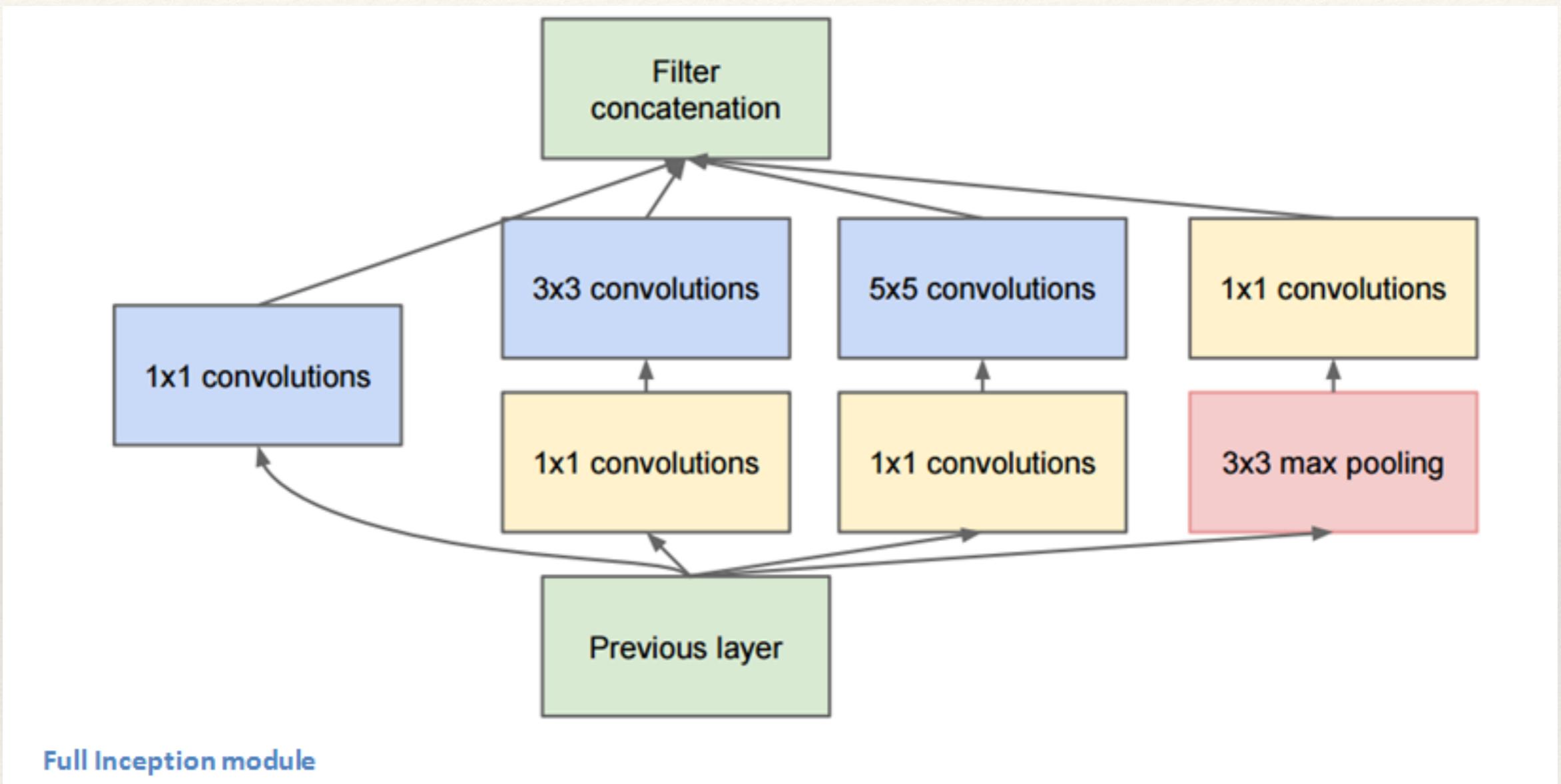
# Google LeNet (2014)



# Inception!

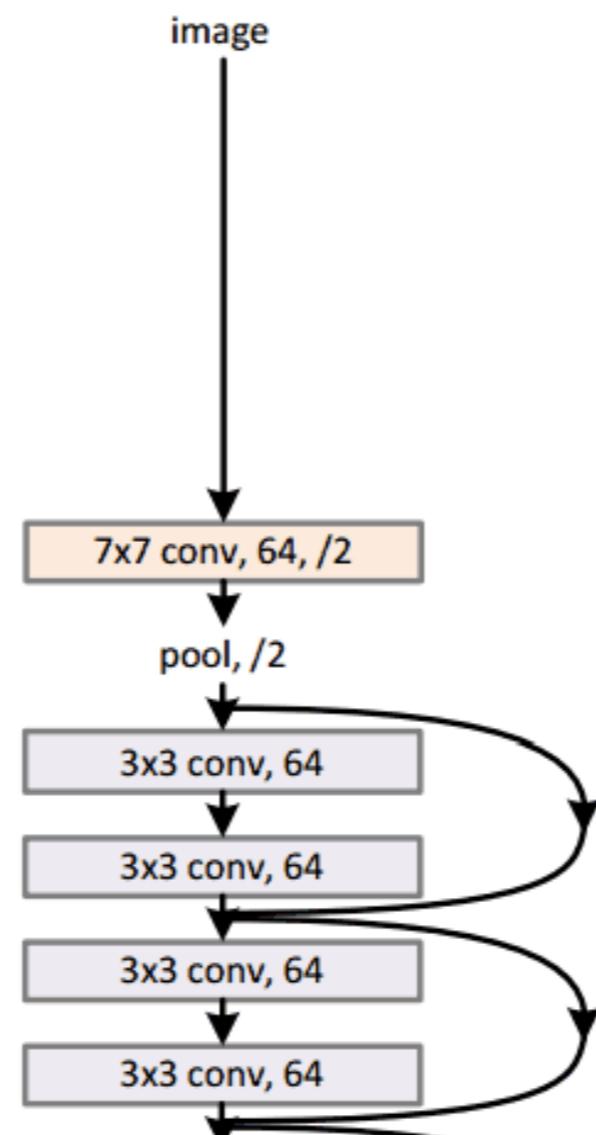


# Parallel layers



# ResNet (He et al., 2015)

34-layer residual

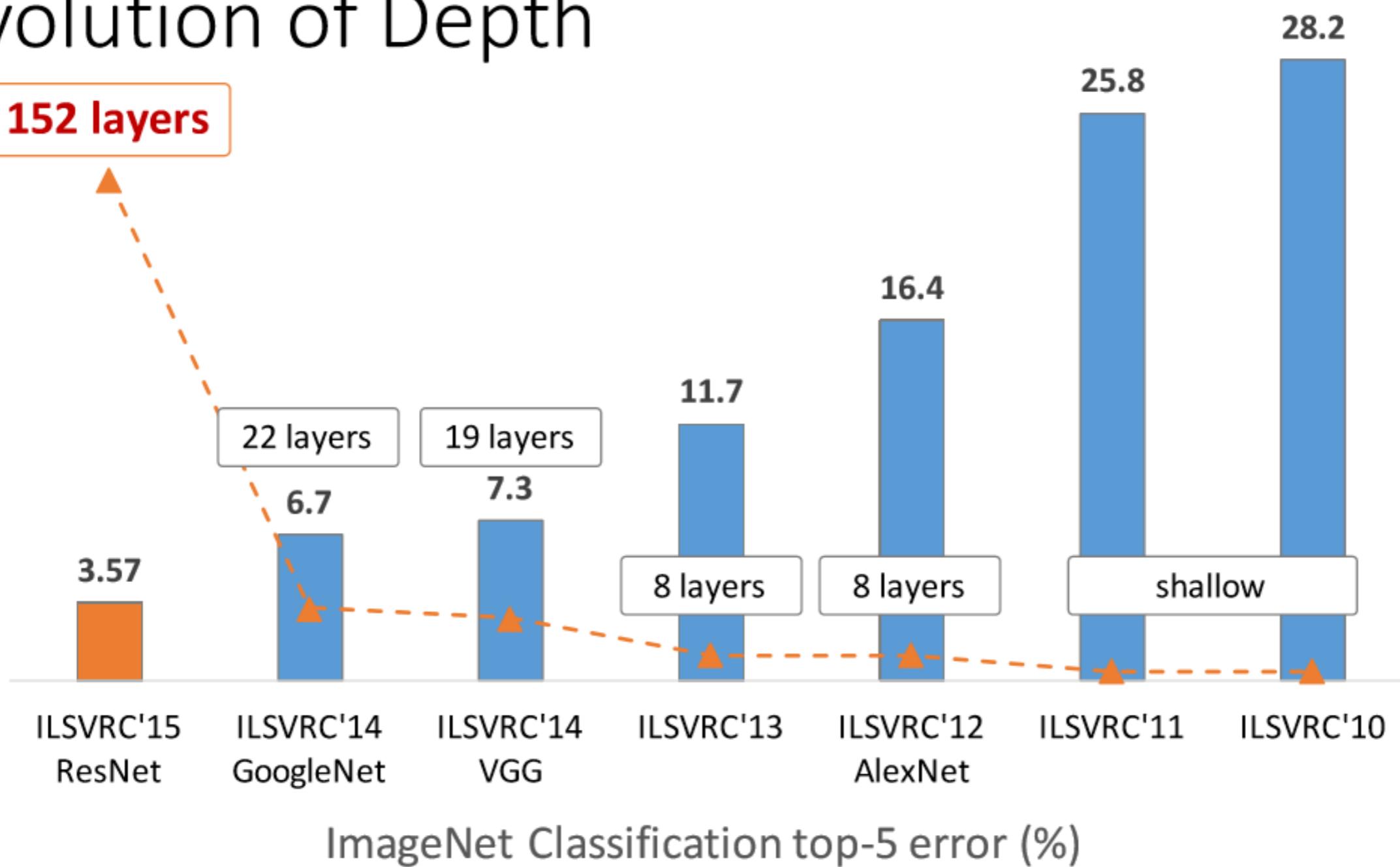


ResNet won ILSVRC 2015 with a top-5 error rate of 3.6%

Depending on their skill and expertise, humans generally hover around a 5-10% error.

~~Superhuman performance!~~  
*But the task is arguably not well defined.*

# Revolution of Depth



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Revolution of Depth

AlexNet, 8 layers  
(ILSVRC 2012)



VGG, 19 layers  
(ILSVRC 2014)



ResNet, **152 layers**  
(ILSVRC 2015)



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Vanishing/exploding gradient problem

---

Backpropagation:

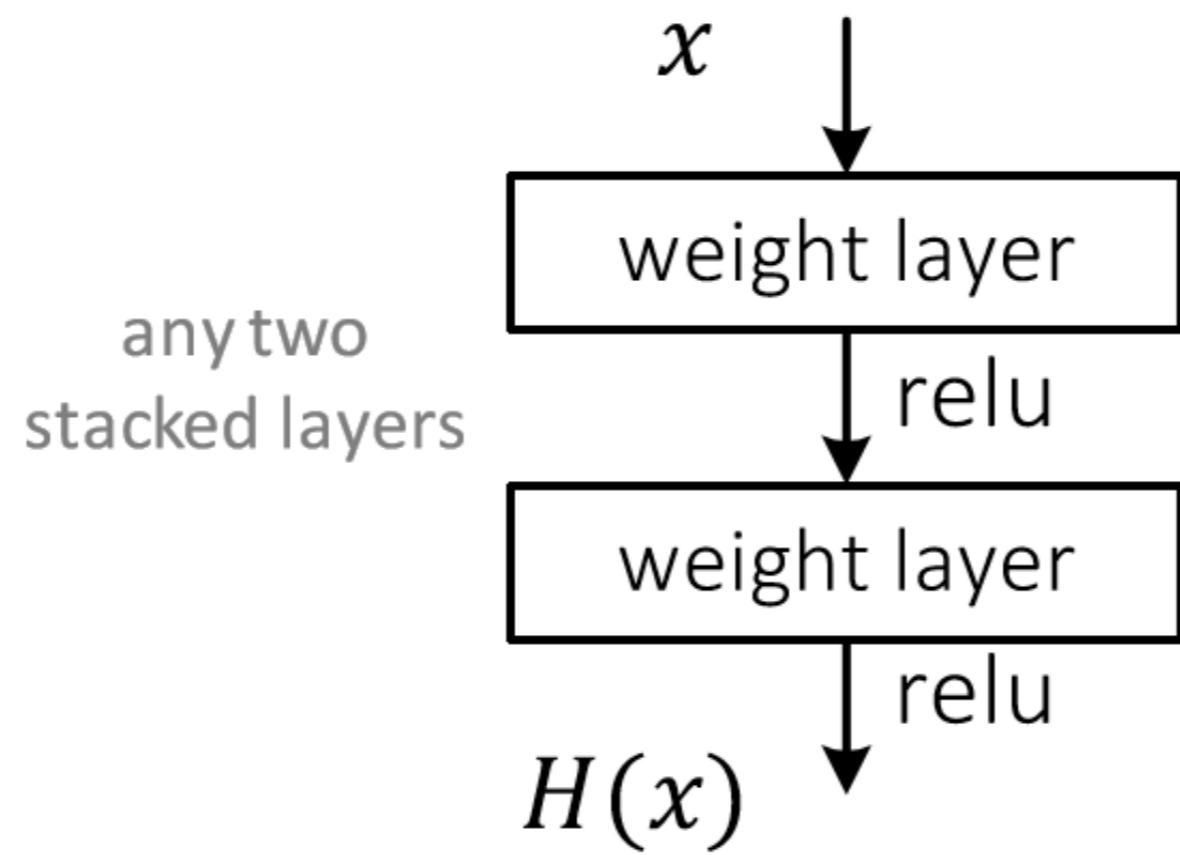
- Compute gradient update for every neuron which was involved in the output across layers

Involves chaining partial derivates over many layers!

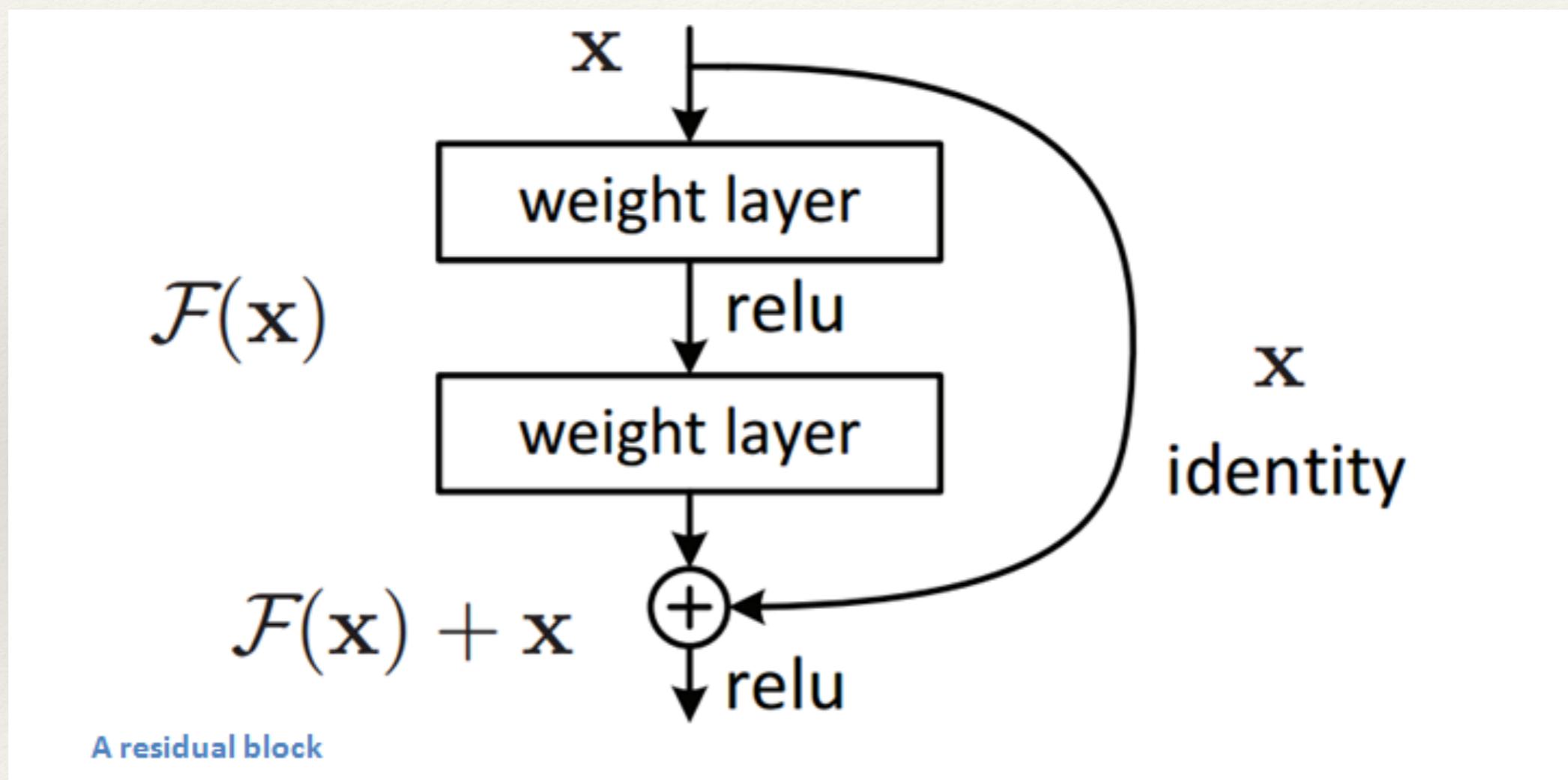
- If derivative  $< 1$ , gradient gets smaller and smaller as we go deeper and deeper -> *vanishing gradients!*
- If derivative  $> 1$ , gradient gets larger and larger as we go deeper and deeper -> *exploding gradients!*

# Regular net

$H(x)$  is any desired mapping,  
hope the 2 weight layers fit  $H(x)$

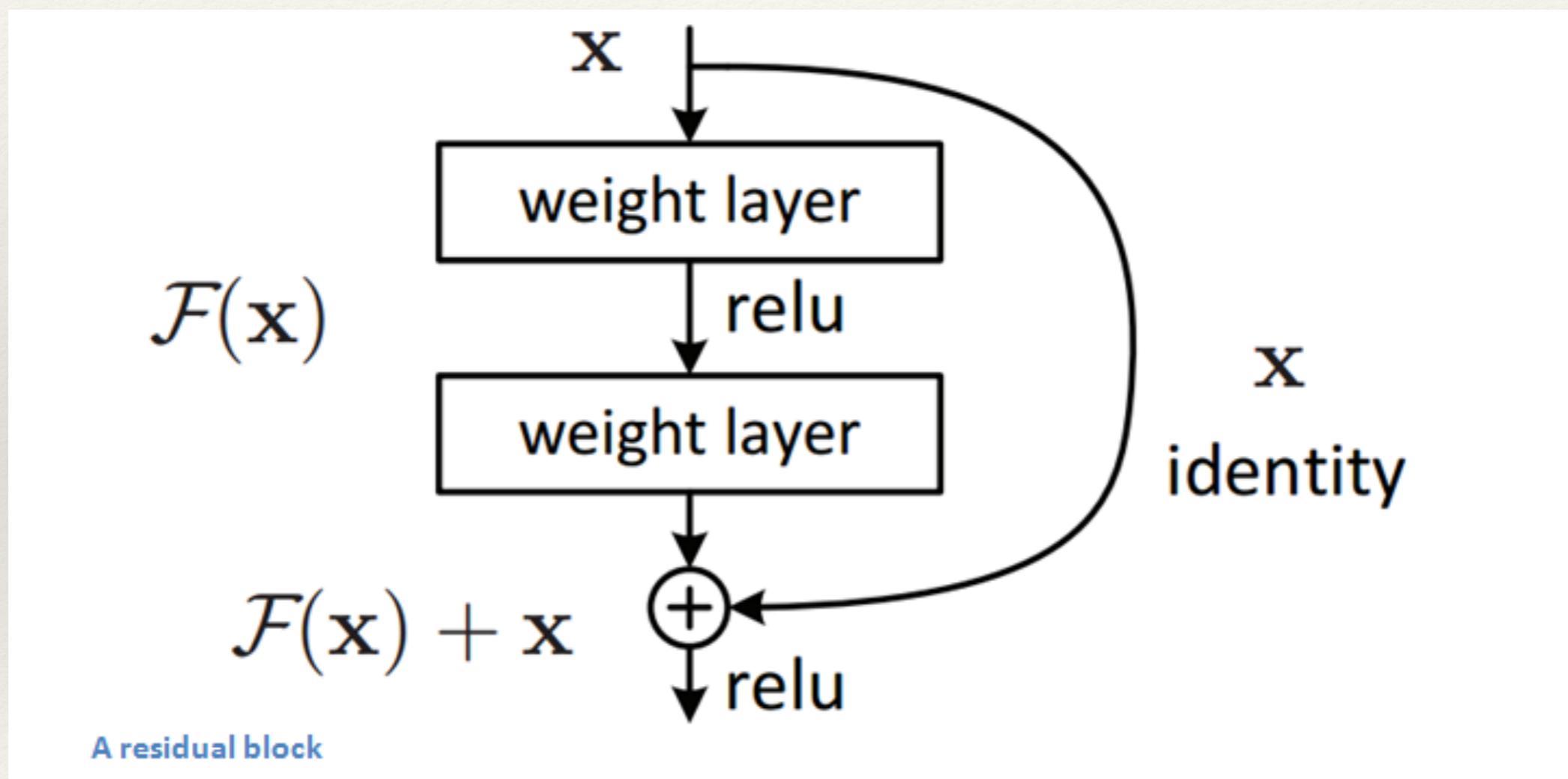


# Residual Unit



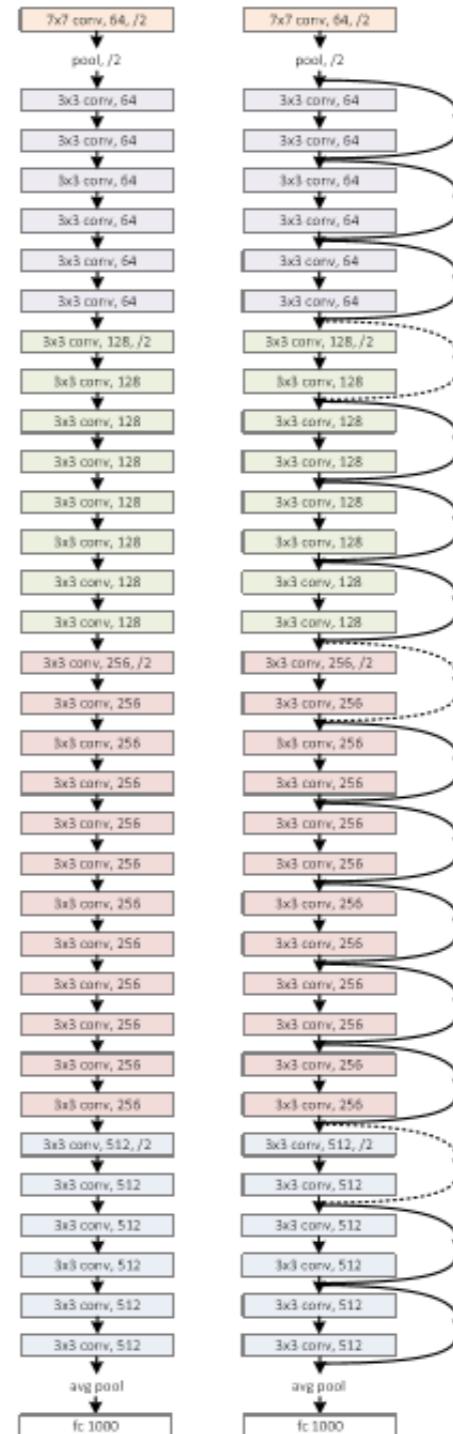
# Residual Unit

The inputs of a lower layer is made available to a node in a higher layer.

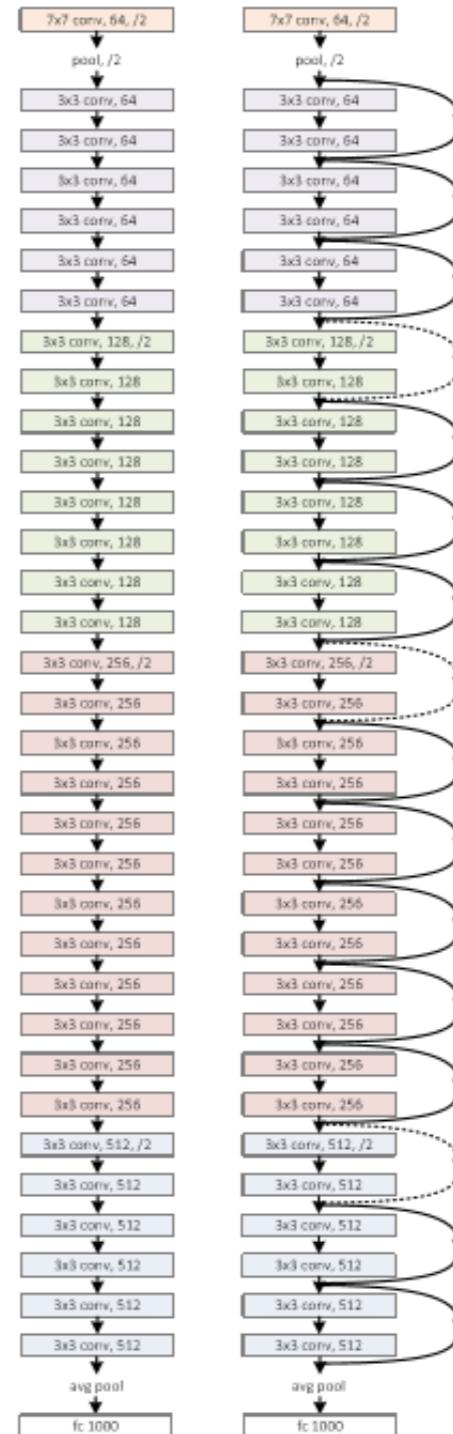


# Network “Design”

plain net



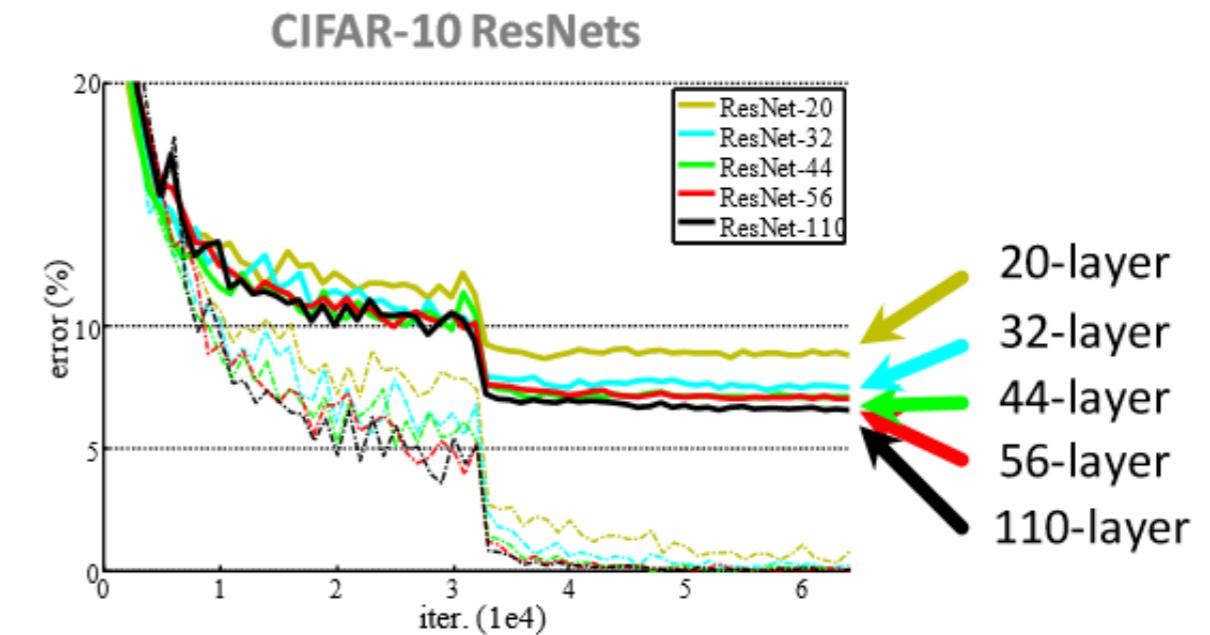
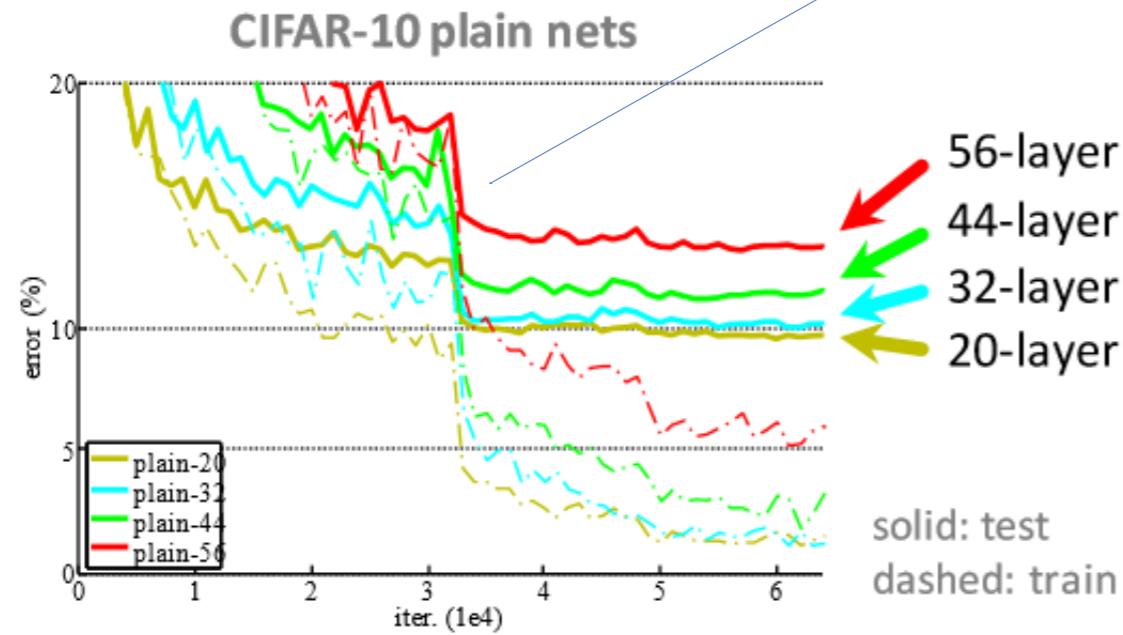
ResNet



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. “Deep Residual Learning for Image Recognition”. CVPR 2016.

# CIFAR-10 experiments

Why so steep?



- Deep ResNets can be trained without difficulties
- Deeper ResNets have **lower training error**, and also lower test error

# CIFAR-10

- ❖ 60,000 32x32 color images, 10 classes

Here are the classes in the dataset, as well as 10 random images from each:

**airplane**



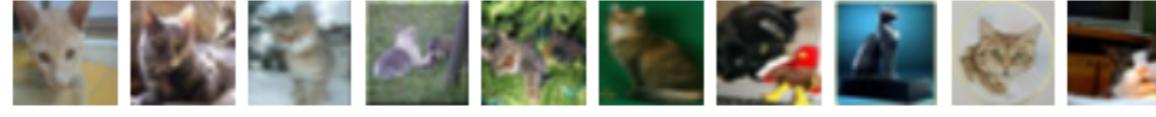
**automobile**



**bird**



**cat**



**deer**



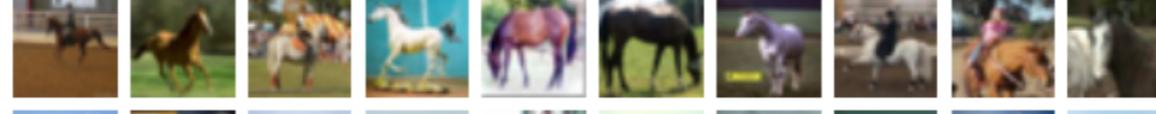
**dog**



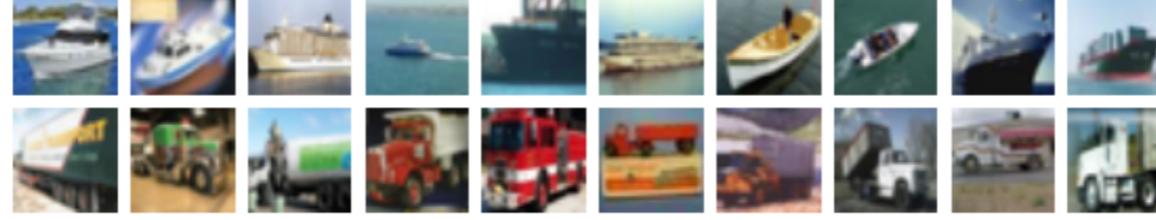
**frog**



**horse**

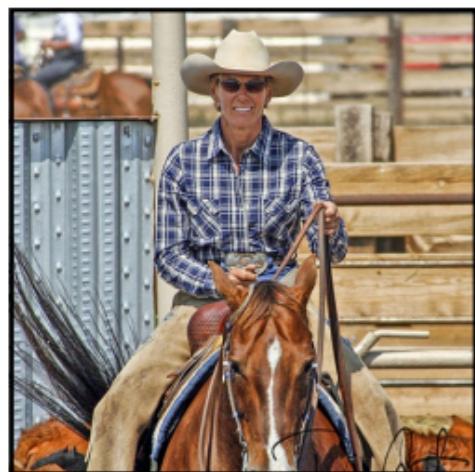


**ship**



**truck**

# R-CNN: Region-based CNN



1. Input image

2. Extract region proposals (~2k)

3. Compute CNN features

4. Classify regions

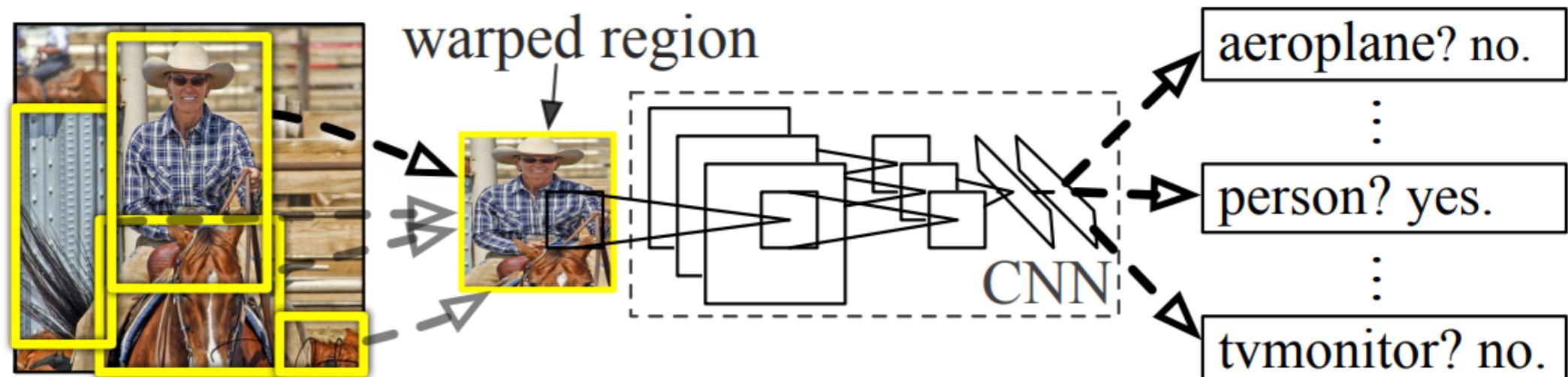
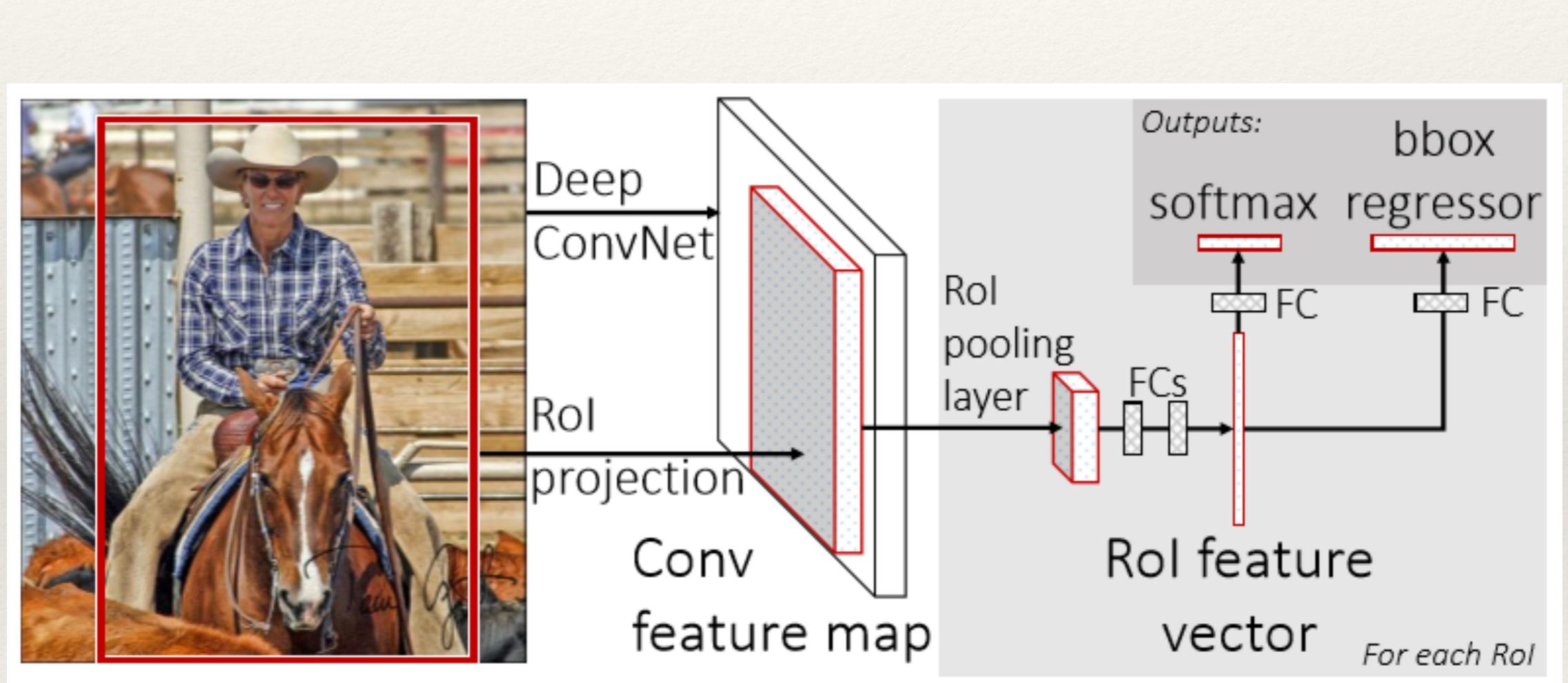


Figure: Girshick et al.

10,000 proposals with recall 0.991 is better....  
but still takes 17 seconds per image to generate them.  
Then I have to test each one!

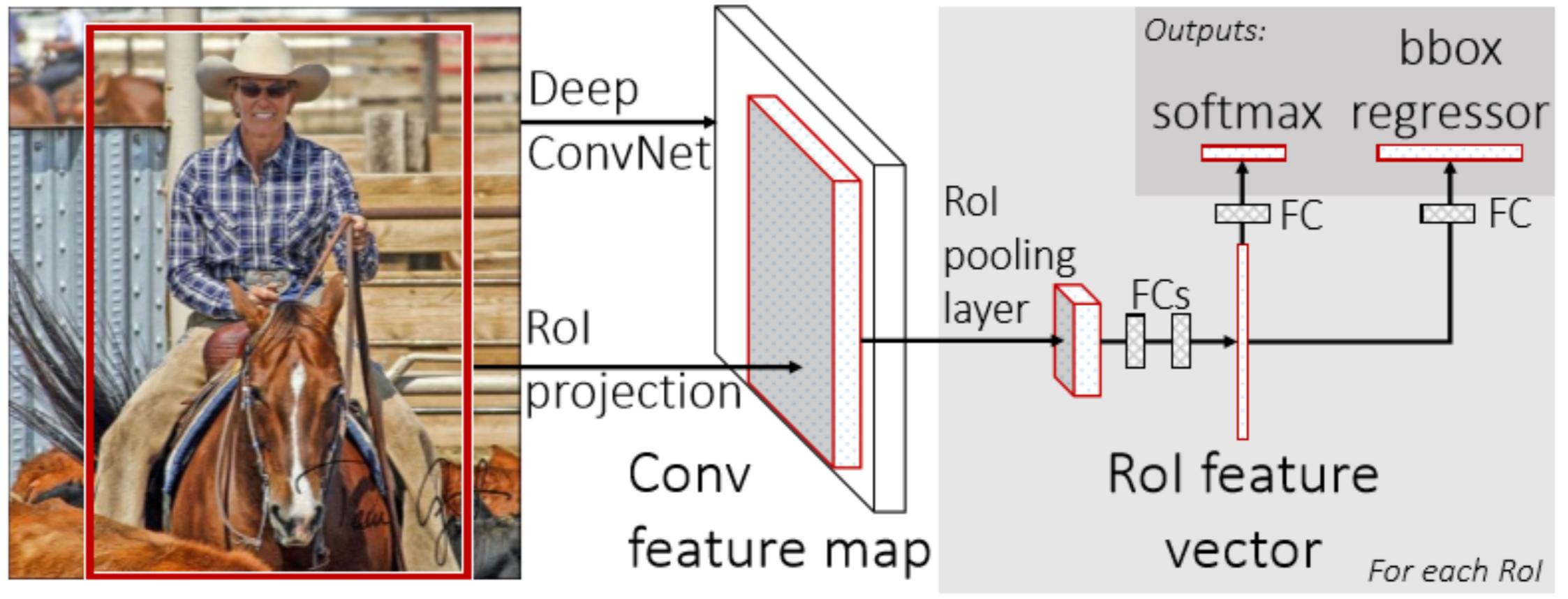
# Fast R-CNN



RoI = Region of Interest

Figure: Girshick et al.

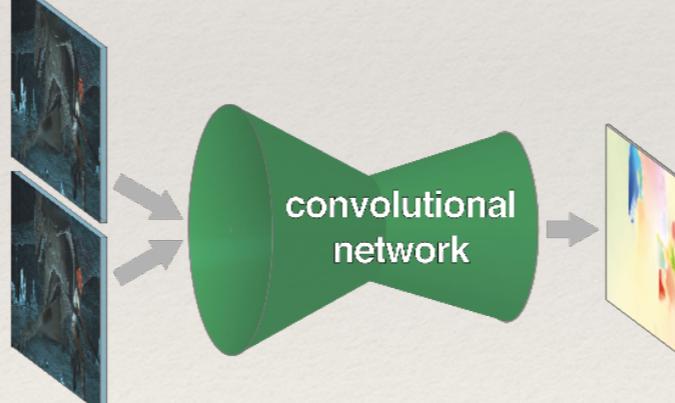
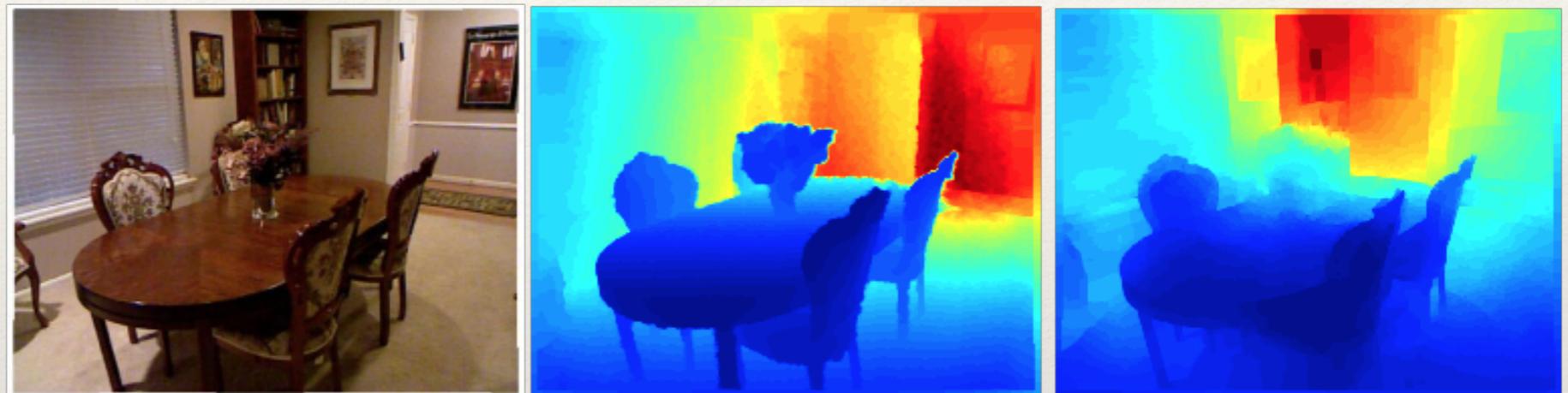
# Fast R-CNN



- Convolve whole image into feature map (many layers; abstracted)
- For each candidate ROI:
  - Squash feature map weights into fixed-size ‘ROI pool’ – adaptive subsampling!
    - Divide ROI into  $H \times W$  subwindows, e.g.,  $7 \times 7$ , and max pool
  - Learn classification on ROI pool with own fully connected layers (FCs)
  - Output classification (softmax) + bounds (regressor)

# What if we want pixels out?

semantic segmentation

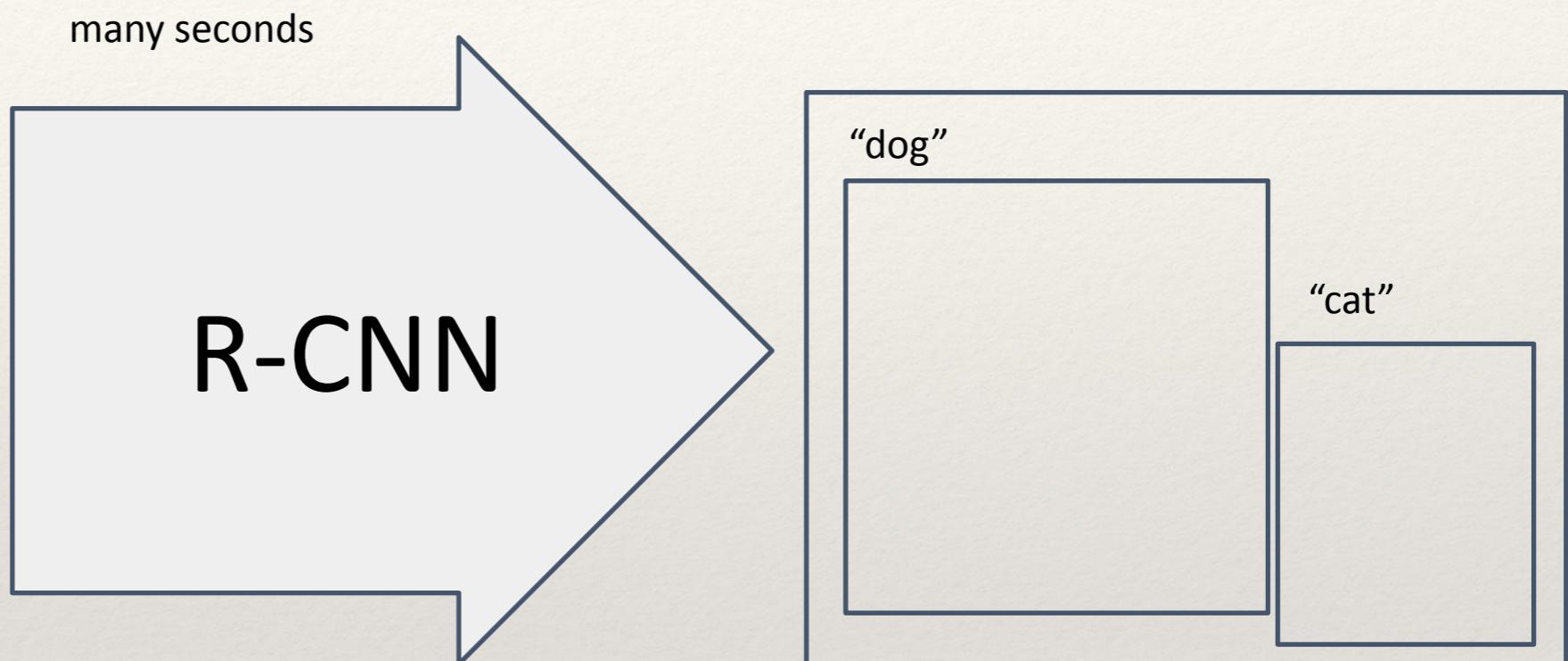


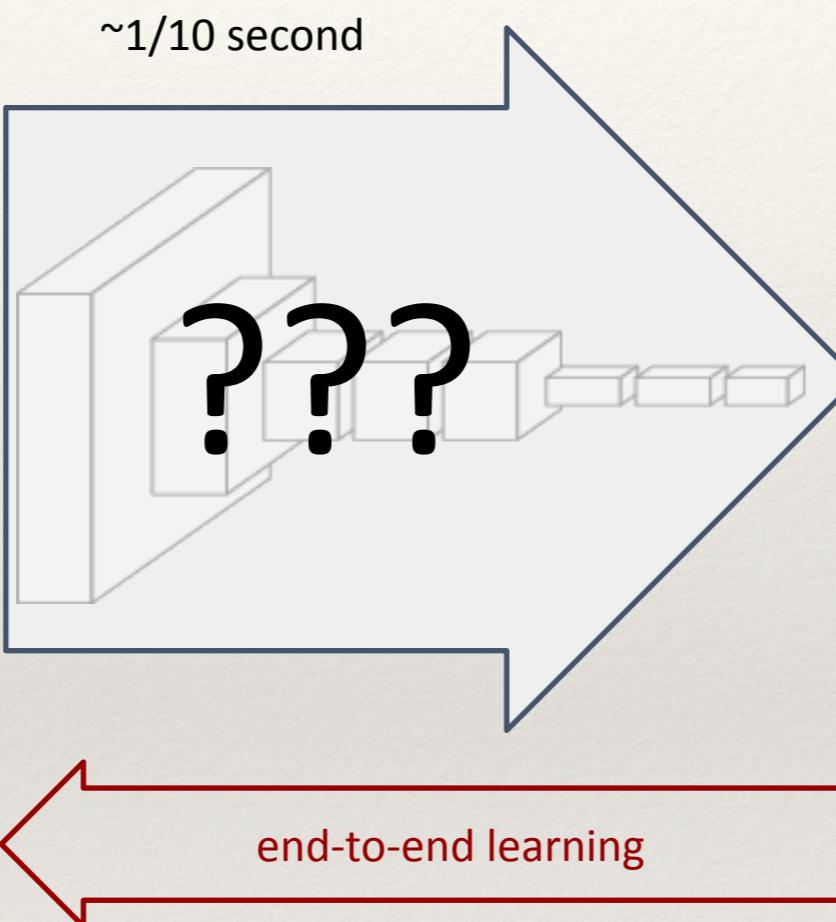
optical flow Fischer et al. 2015



boundary prediction Xie & Tu 2015

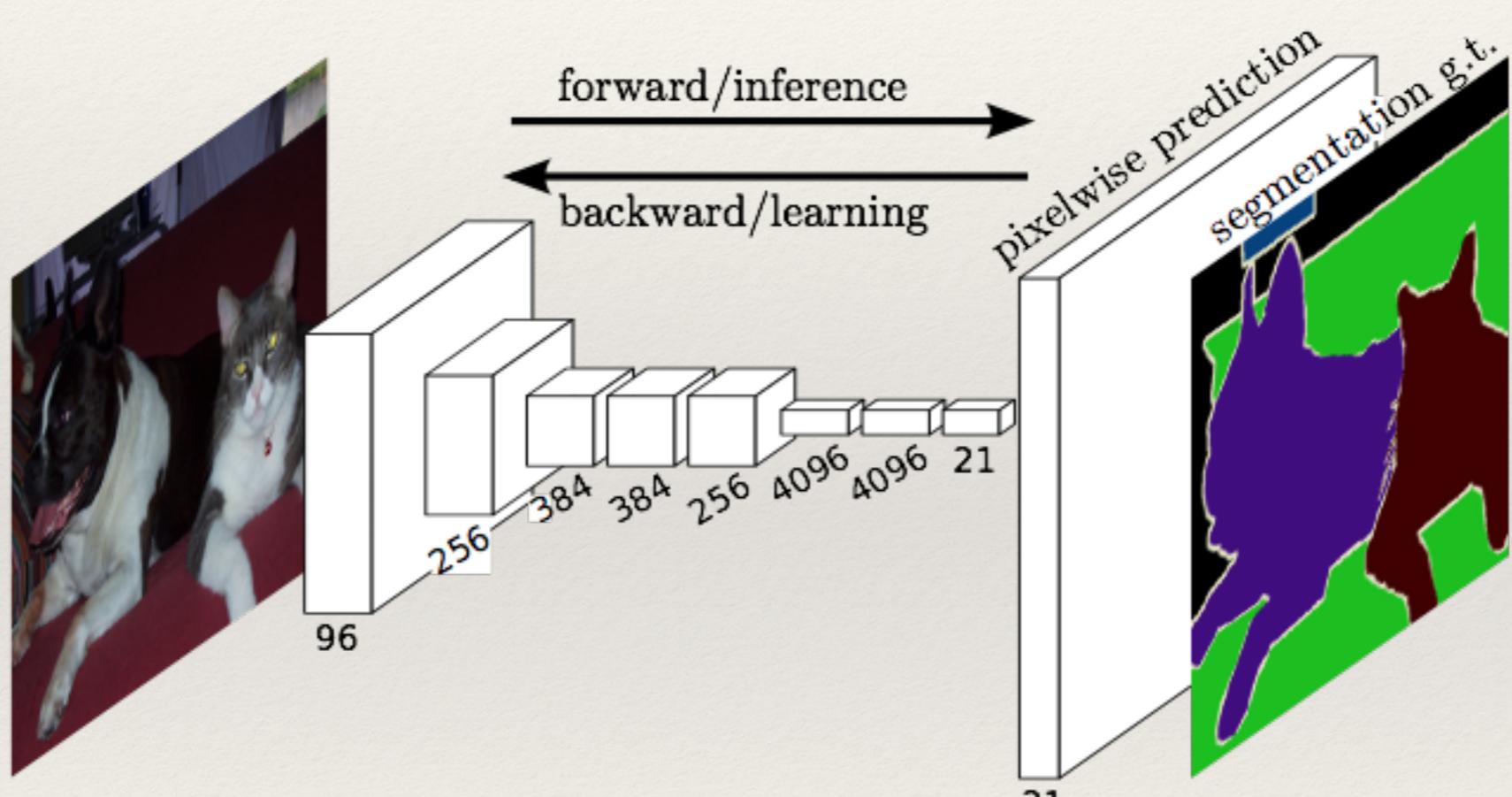
# R-CNN does detection





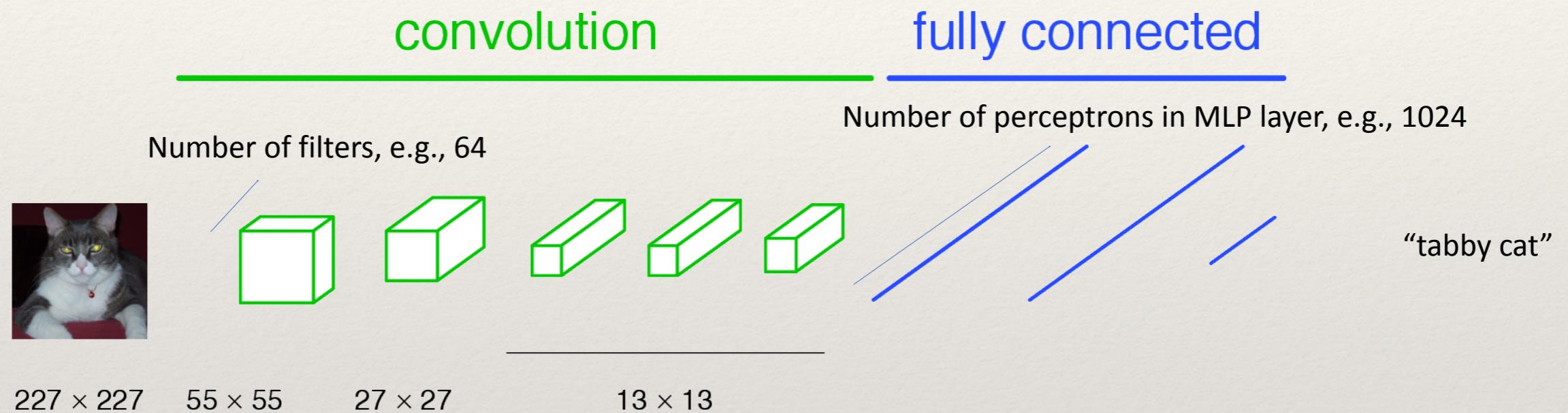
# Fully Convolutional Networks for Semantic Segmentation

- ❖ UC Berkeley

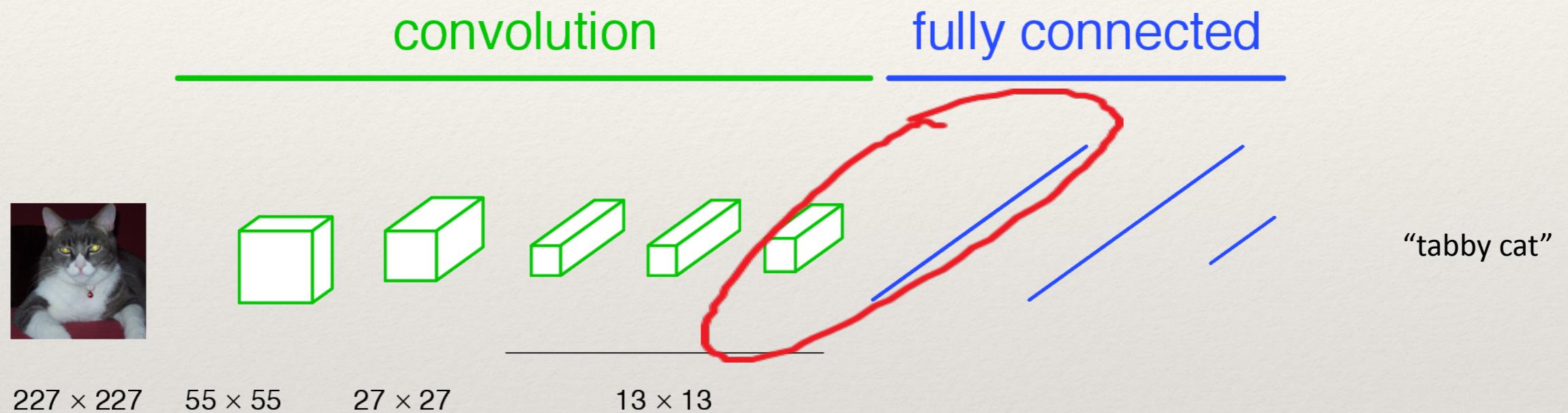


Jonathan Long\* Evan Shelhamer\* Trevor Darrell

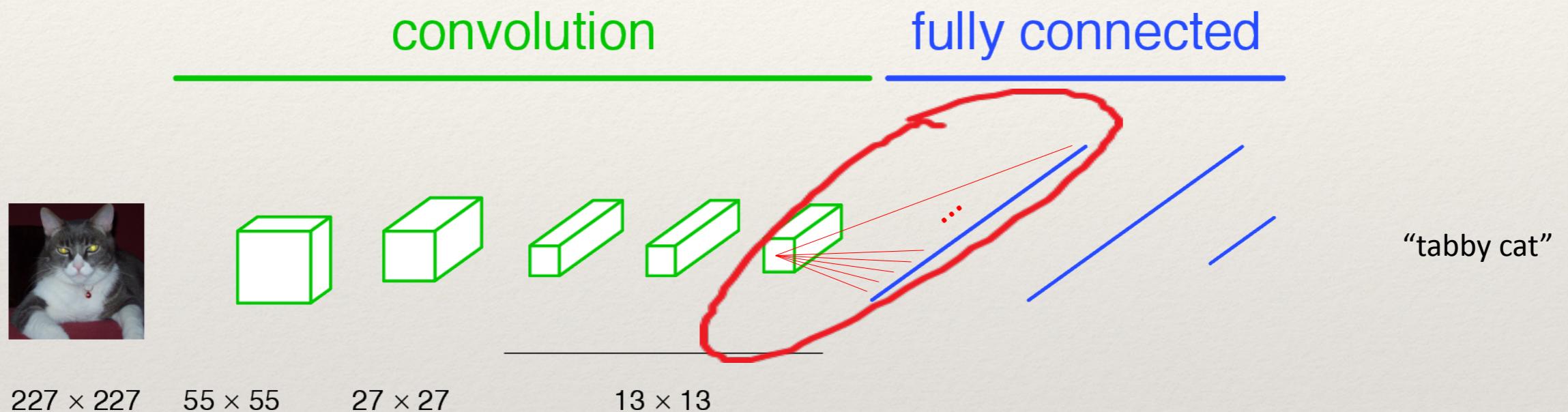
# A classification network...



# A classification network...

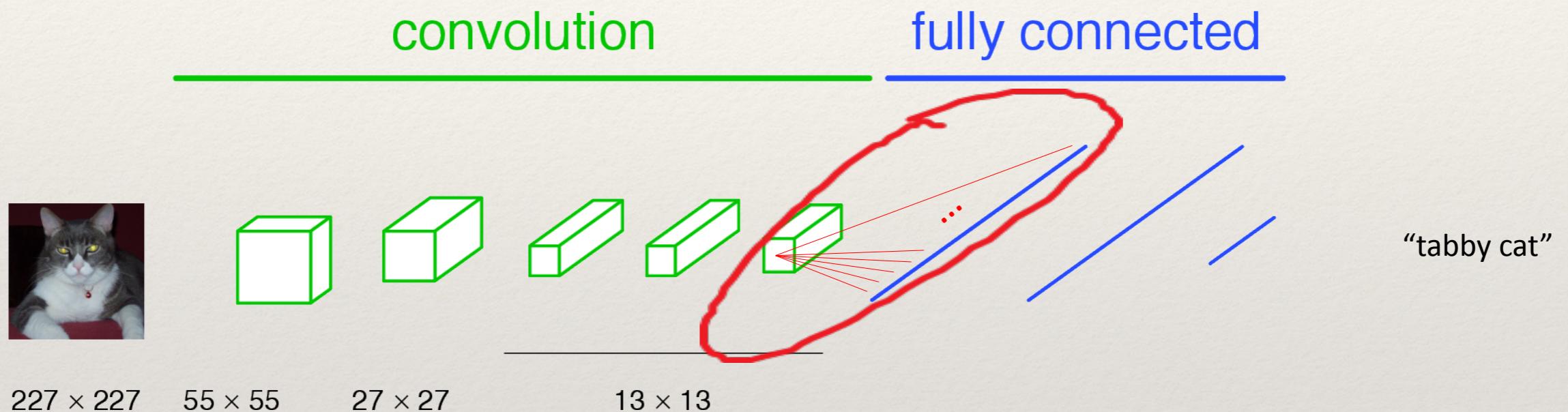


# A classification network...



The response of every kernel across all positions are attached densely to the array of perceptrons in the fully-connected layer.

# A classification network...



The response of every kernel across all positions are attached densely to the array of perceptrons in the fully-connected layer.

AlexNet: 256 filters over 6x6 response map

Each 2,359,296 response is attached to one of 4096 perceptrons, leading to 37 mil params.

# Problem

---

We want a label at every pixel

Current network gives us a label for the whole image.

Approach:

- ❖ Make CNN for every sub-image size ?
- ❖ ‘Convolutionalize’ *all layers* of network, so that we can treat it as one (complex) filter and slide around our full image.

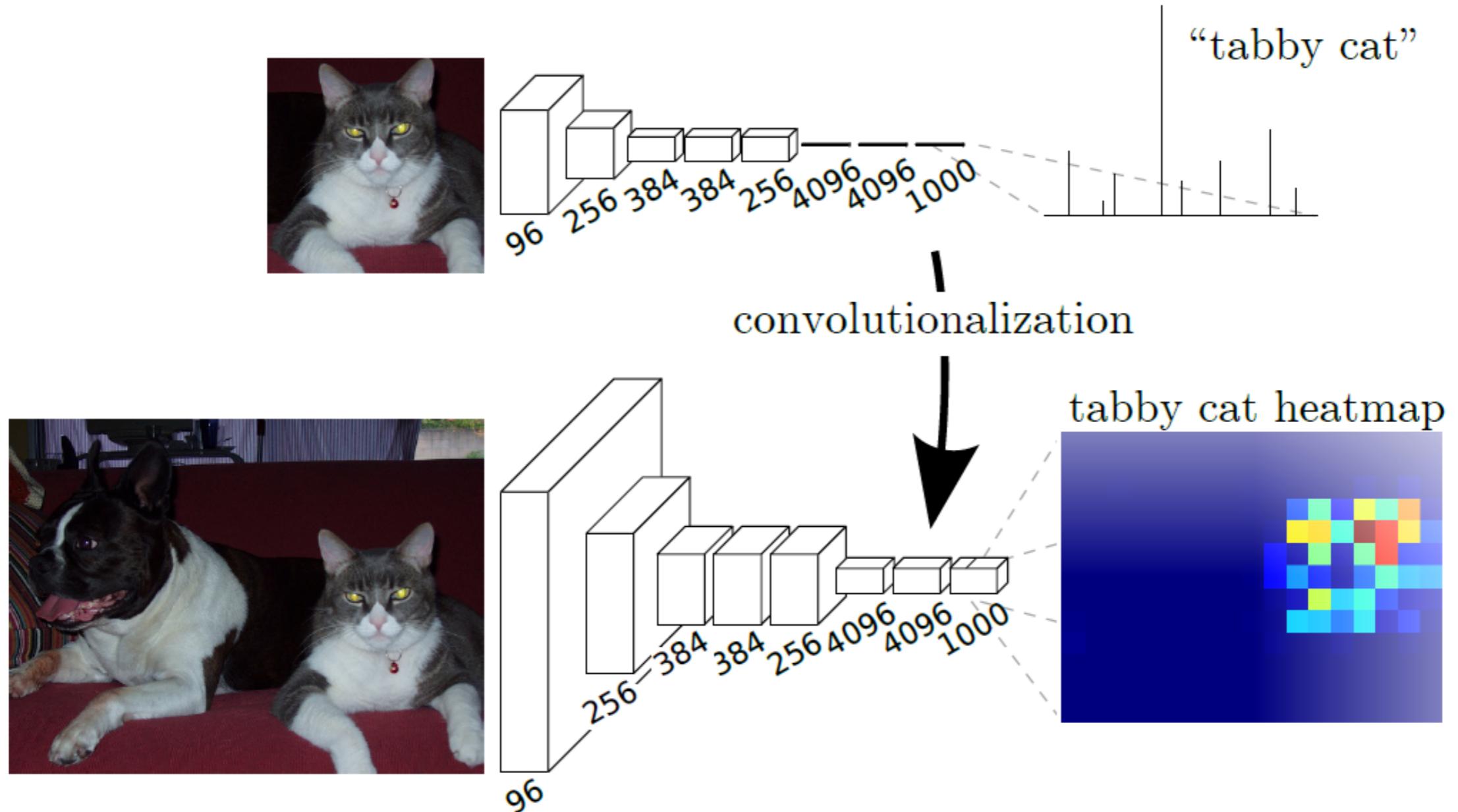


Figure 2. Transforming fully connected layers into convolution layers enables a classification net to output a heatmap. Adding layers and a spatial loss (as in Figure 1) produces an efficient machine for end-to-end dense learning.



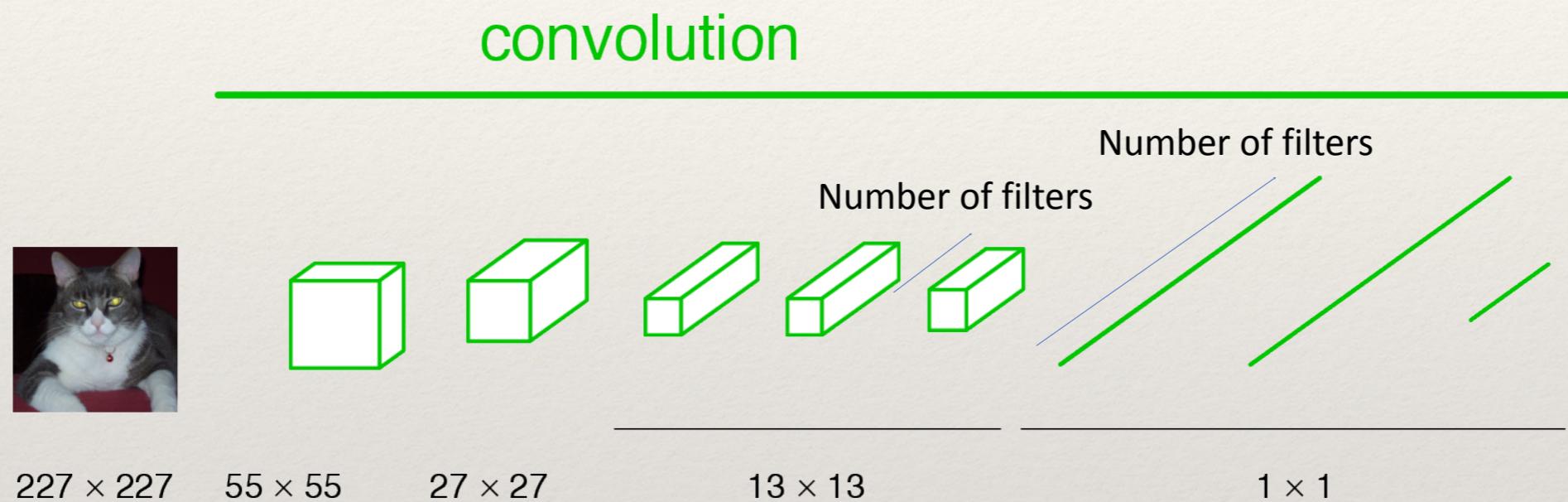
**Yann LeCun**

6 April 2015 ·

Follow

In Convolutional Nets, there is no such thing as "fully-connected layers". There are only convolution layers with 1x1 convolution kernels and a full connection table.

# Convolutionalization

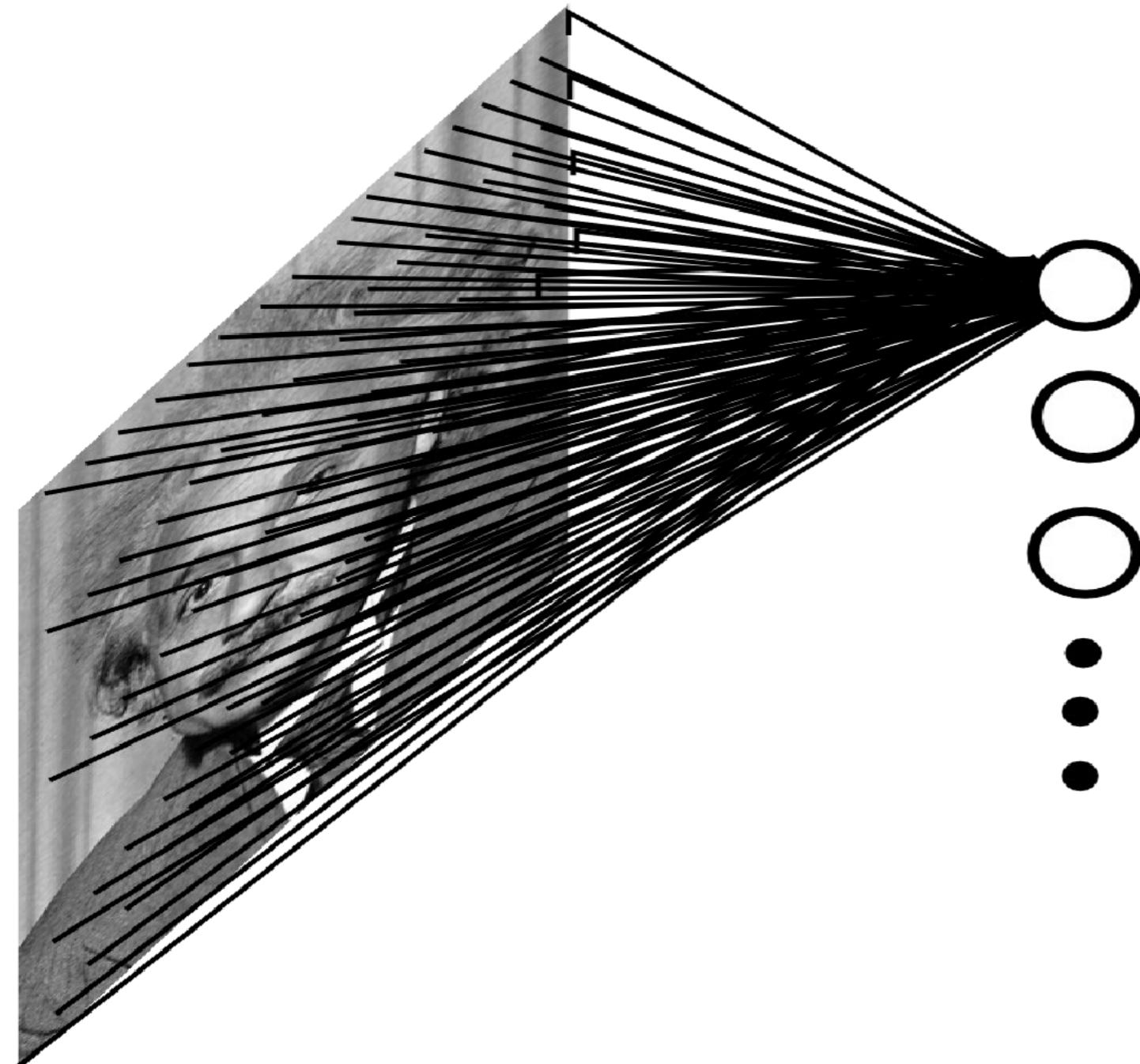


1x1 convolution operates across all filters in the previous layer, and is slid across all positions.

# Back to the fully-connected perceptron...

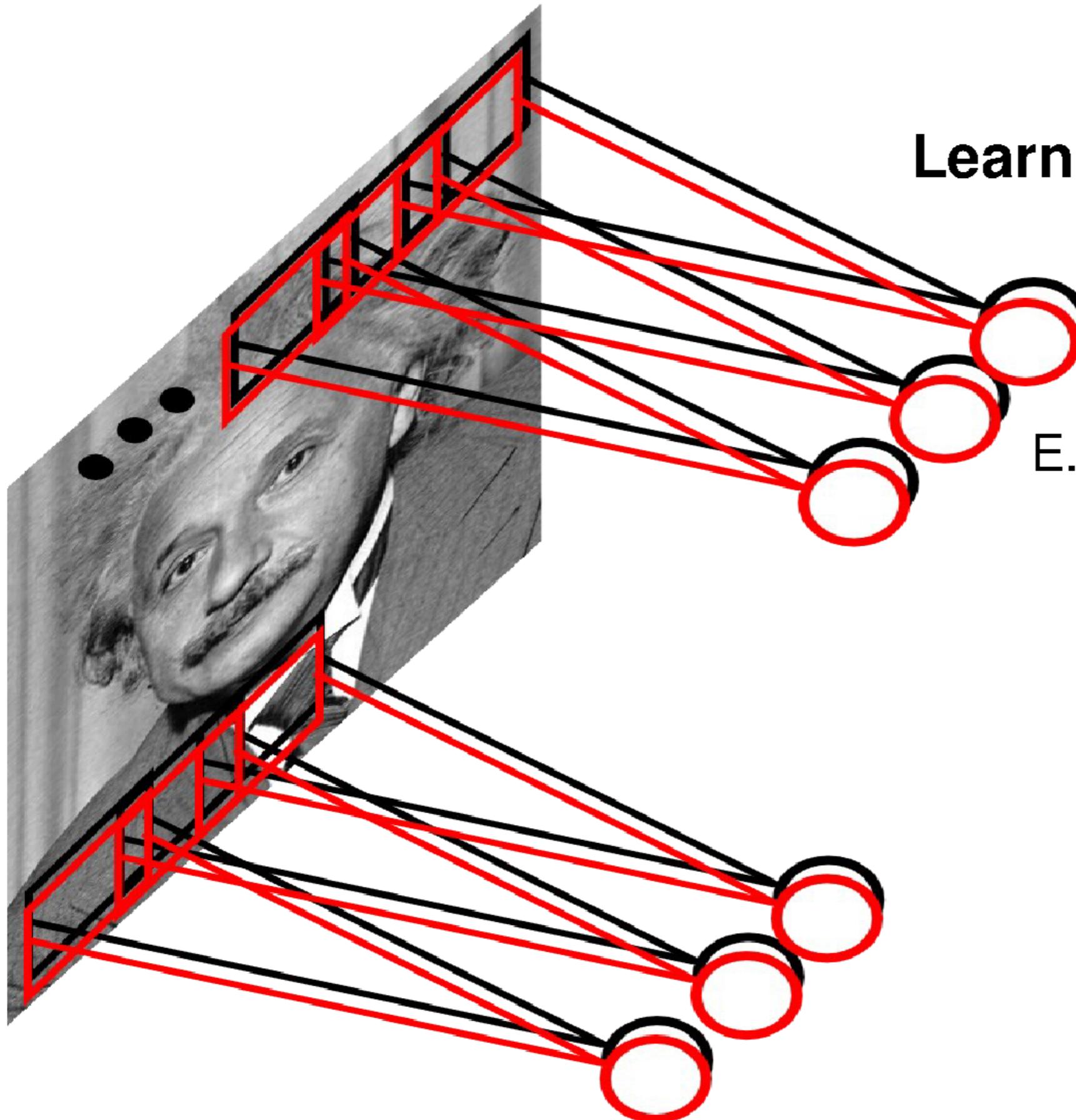
$$\text{output} = \begin{cases} 0 & \text{if } w \cdot x \leq 0 \\ 1 & \text{if } w \cdot x > 0 \end{cases}$$

$$w \cdot x \equiv \sum_j w_j x_j$$



Perceptron is connected to every value in the previous layer (across all channels; 1 visible).

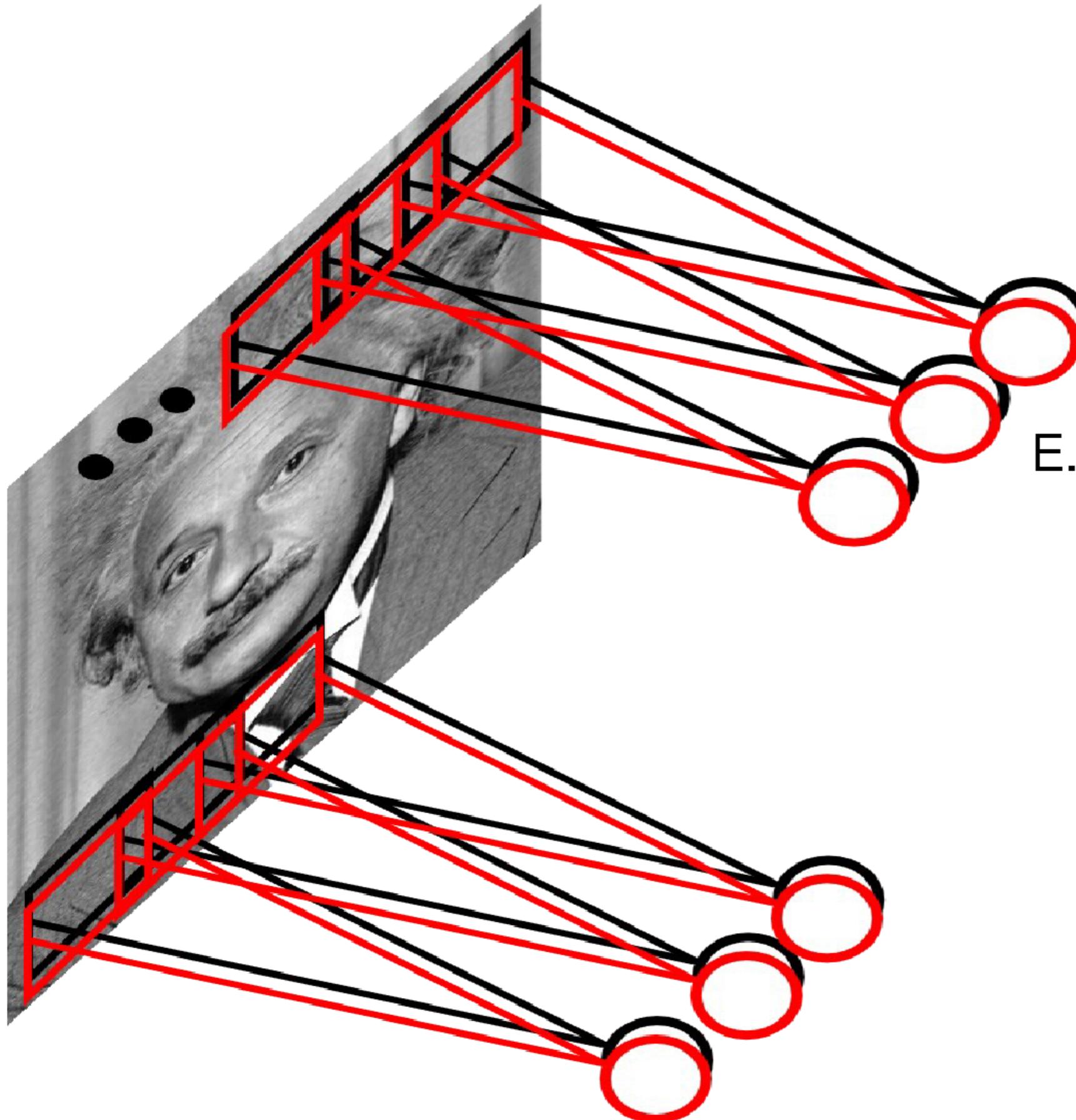
# Convolutional Layer



Learn multiple filters.

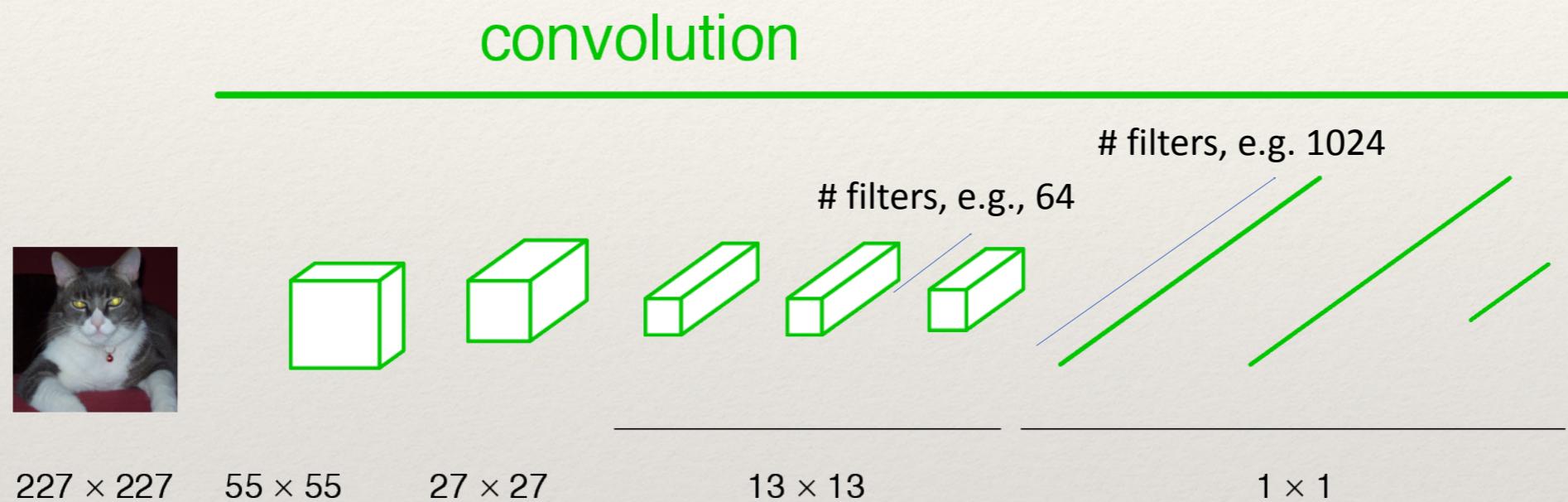
E.g.: 200x200 image  
100 Filters  
Filter size: 10x10  
10K parameters

# Convolutional Layer



E.g.: 200x200 image  
100 Filters  
Filter size: 1x1  
100 parameters

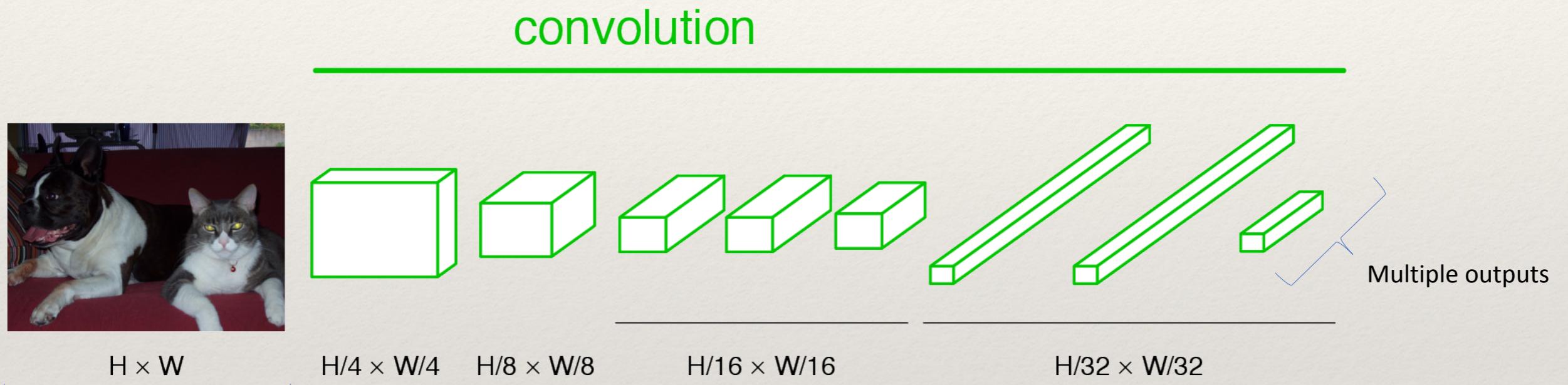
# Convolutionalization



1x1 convolution operates across all filters in the previous layer, and is slid across all positions.

e.g.,  $64 \times 1 \times 1$  kernel, with shared weights over  $13 \times 13$  output,  $\times 1024$  filters = 11mil params.

# Becoming fully convolutional



Arbitrary-  
sized image

When we turn these operations into a convolution, the  $13 \times 13$  just becomes another parameter and our output size adjust dynamically.

Now we have a *vector/matrix* output, and our network acts itself like a complex filter.

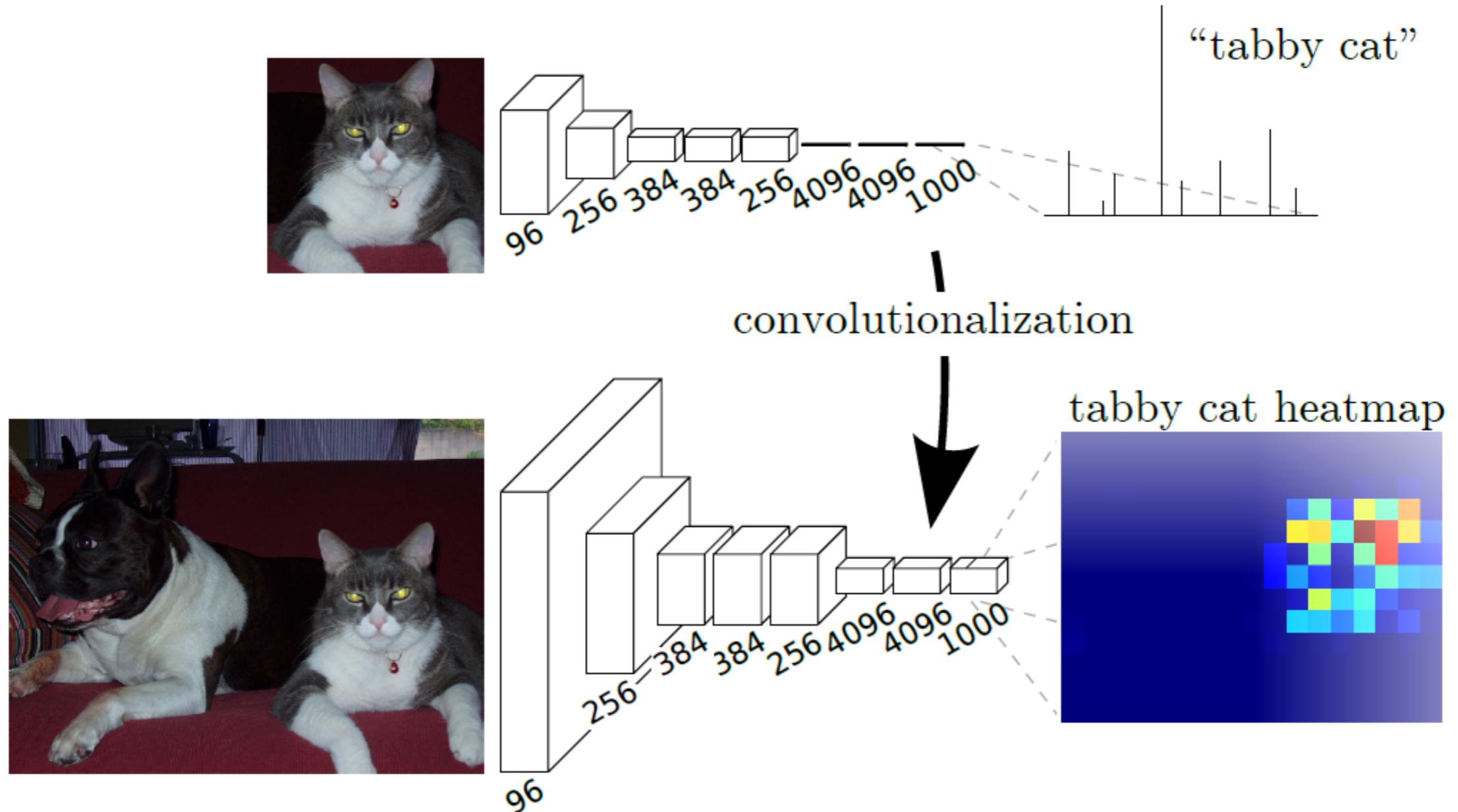
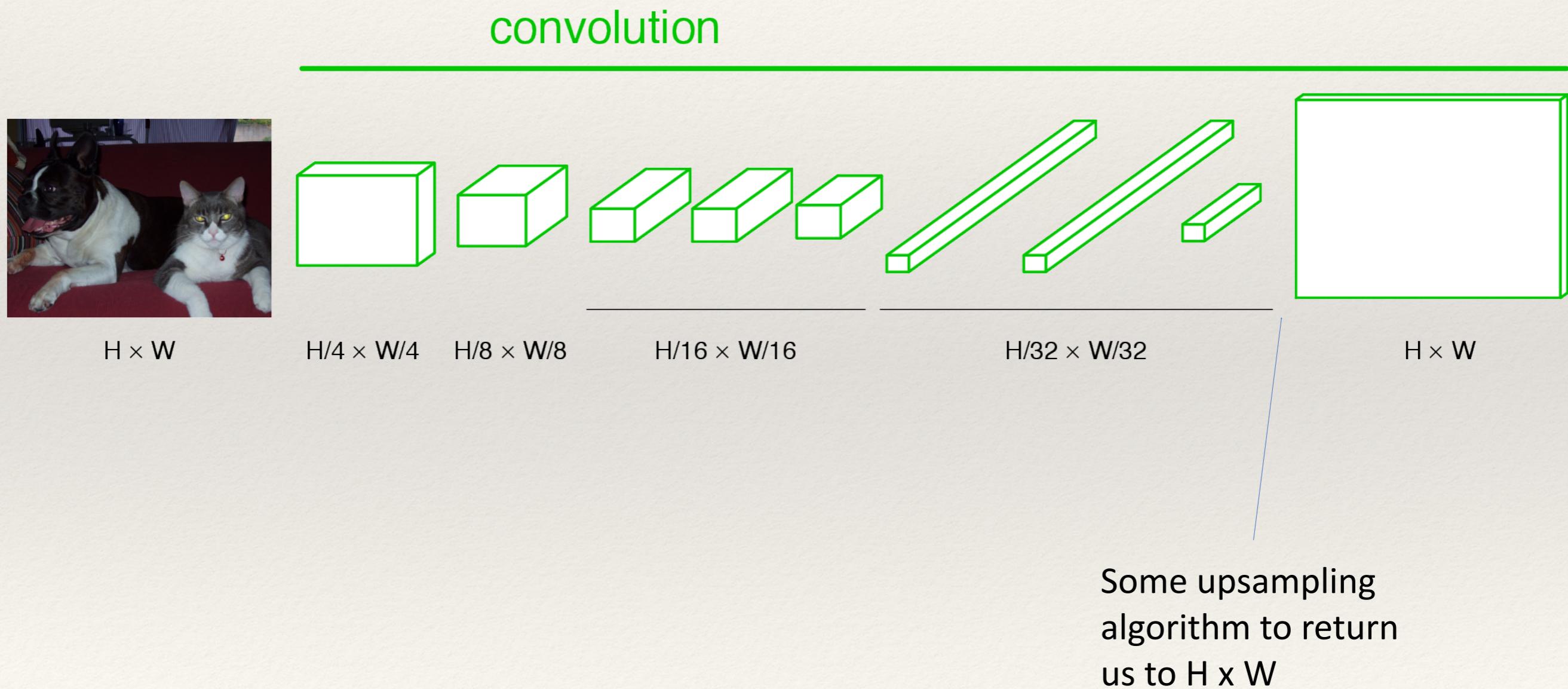
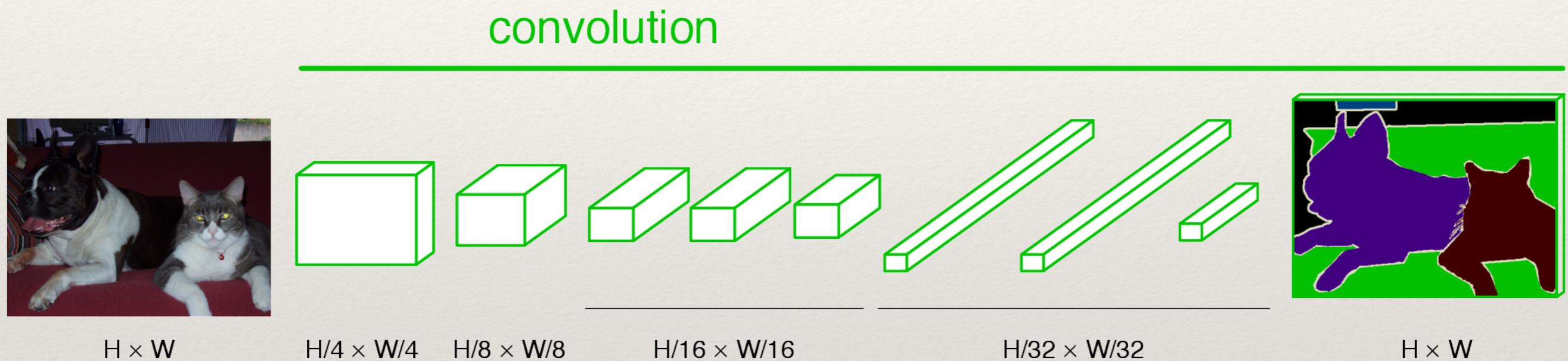


Figure 2. Transforming fully connected layers into convolution layers enables a classification net to output a heatmap. Adding layers and a spatial loss (as in Figure 1) produces an efficient machine for end-to-end dense learning.

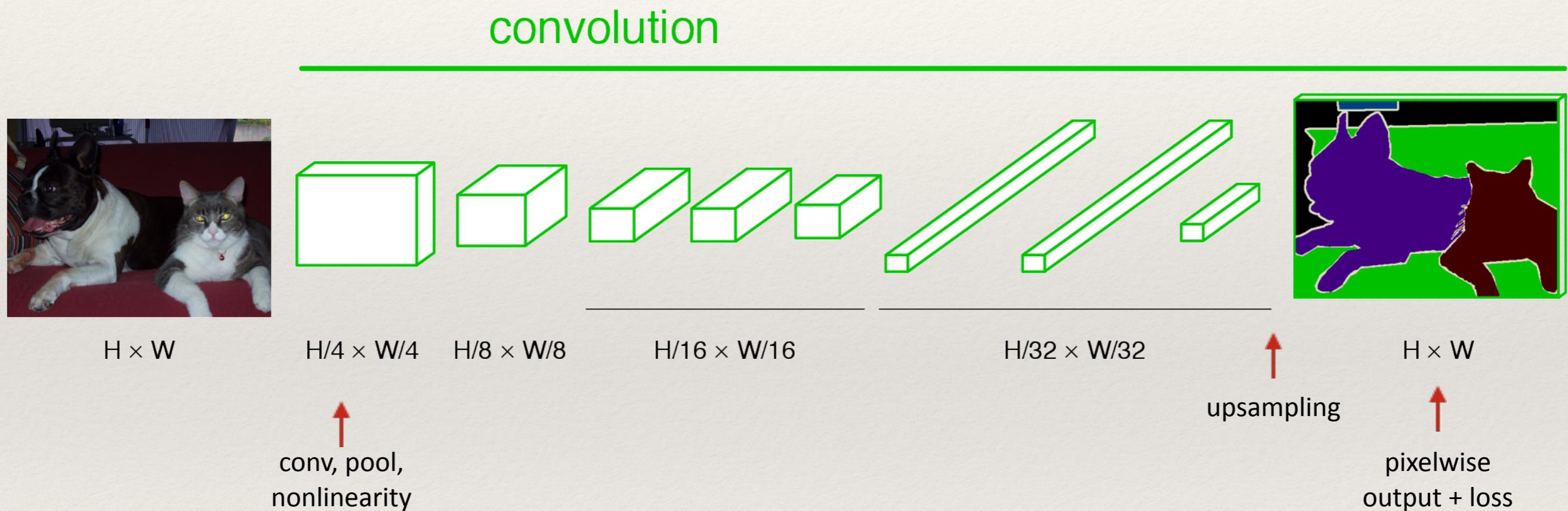
# Upsampling the output



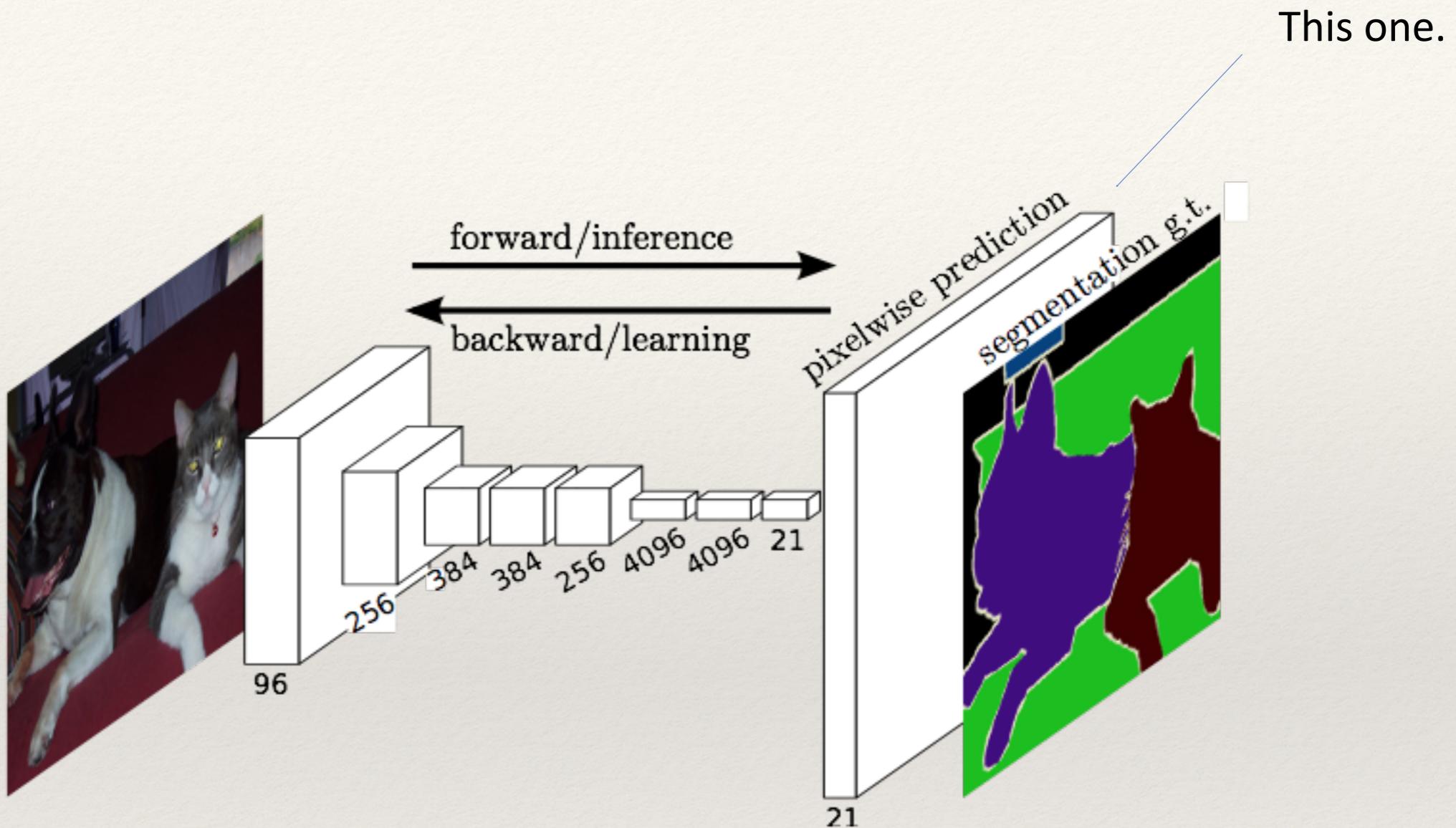
# End-to-end, pixels-to-pixels network



# End-to-end, pixels-to-pixels network

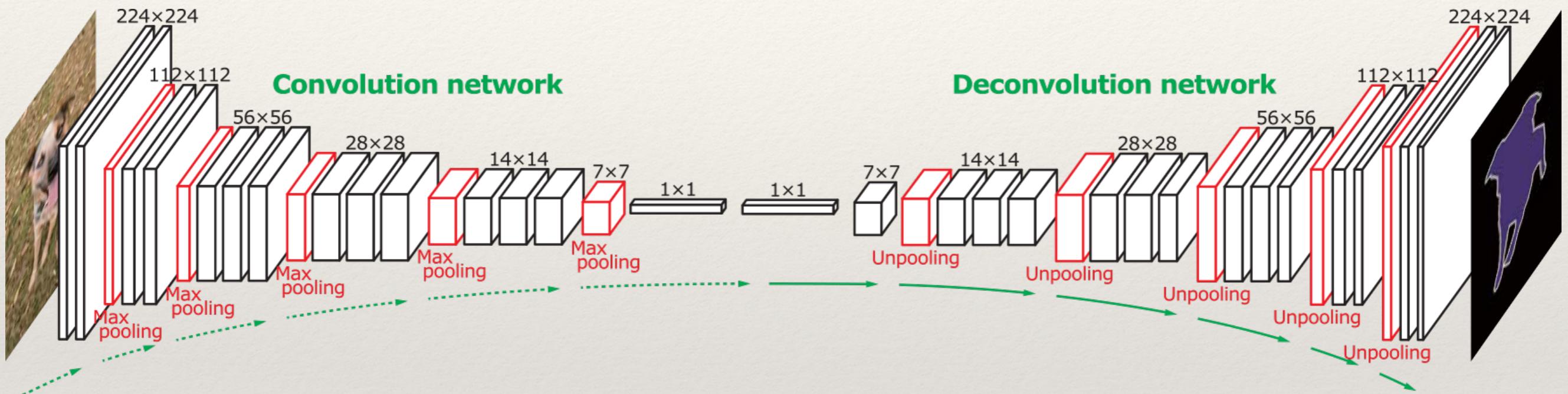


# What is the upsampling layer?



Hint: it's actually an upsampling *network*

# 'Deconvolution' networks learn to upsample



Often called “deconvolution”, but misnomer.

Not the deconvolution that we saw in deblurring -> that is division in the Fourier domain.

‘Transposed convolution’ is better.