# COMP547

## DEEP UNSUPERVISED LEARNING

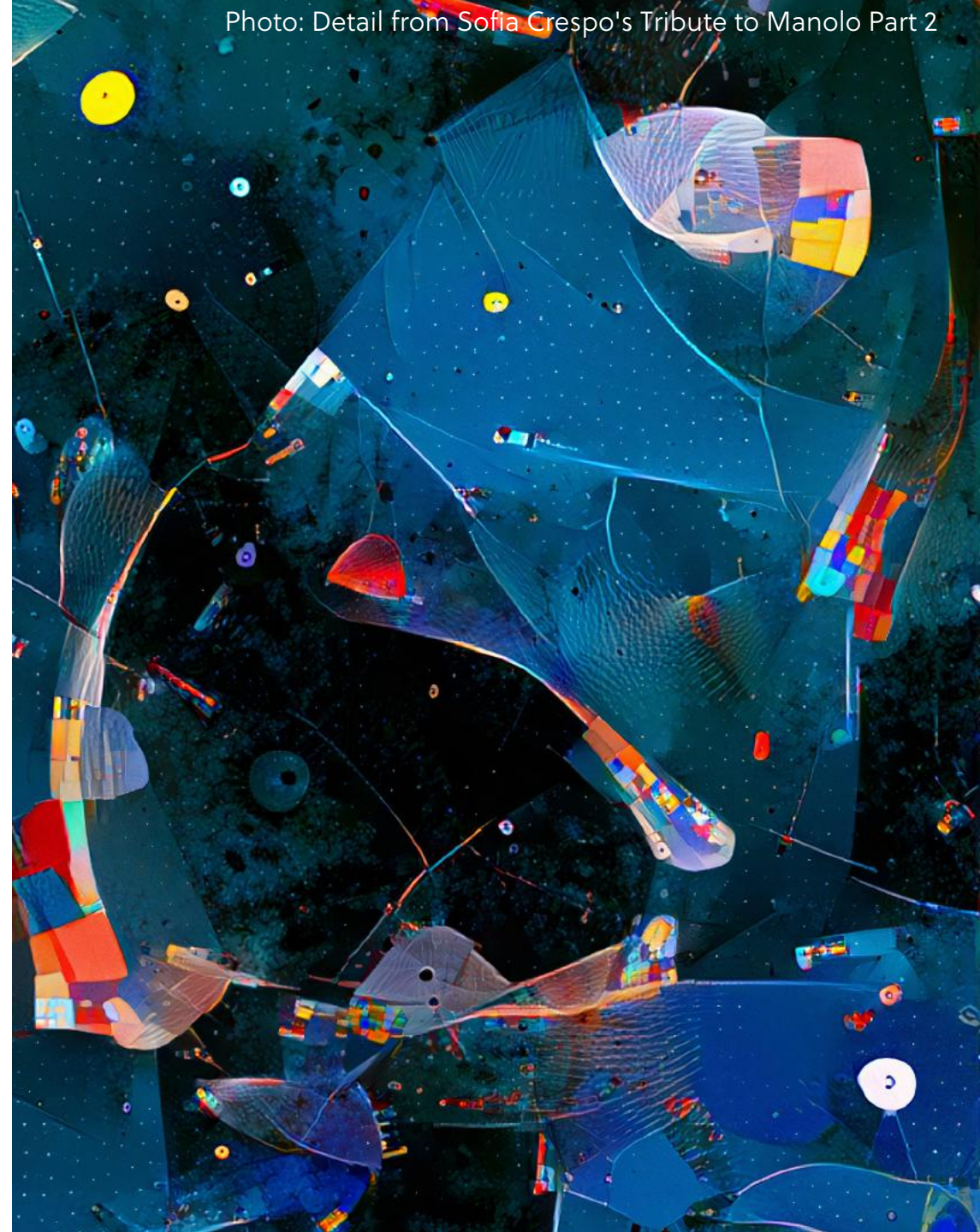## Lecture #2 – Neural Networks Basics and Spatial Processing with CNNs

**KOÇ UNIVERSITY**

Aykut Erdem // Koç University // Spring 2021

# Previously on COMP547

- course logistics

- course topics
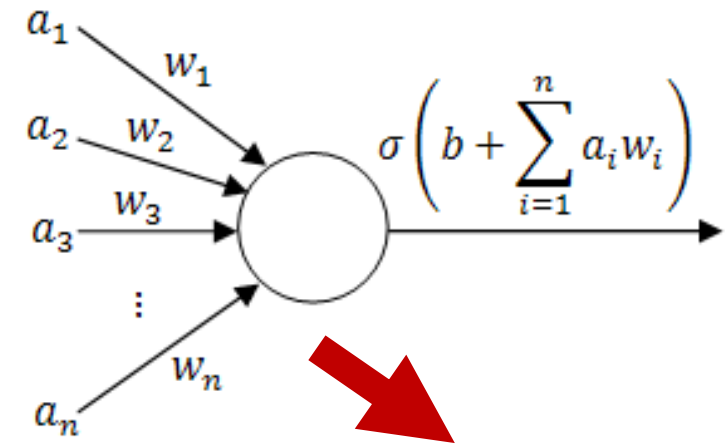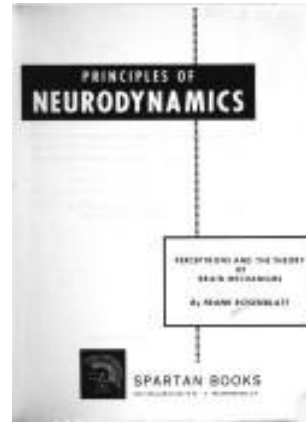
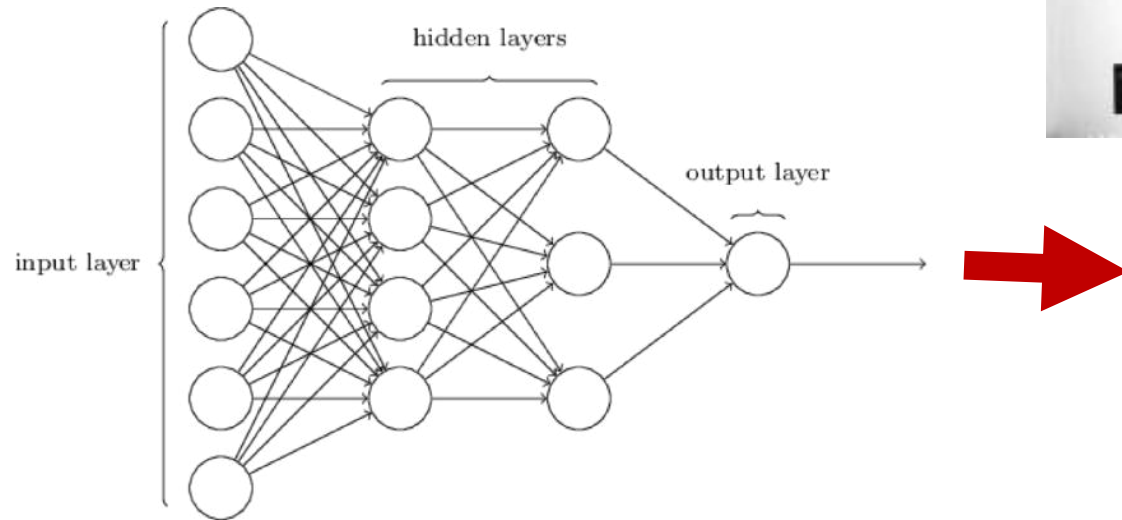- what is deep unsupervised learning

# Lecture overview

- deep learning

- computation in a neural net

- optimization

- backpropagation

- training tricks

- convolutional neural networks


- **Disclaimer:** Much of the material and slides for this lecture were borrowed from
  - —Costis Daskalakis and Aleksander Mądry's MIT 6.883 class
  - —Bill Freeman, Antonio Torralba and Phillip Isola's MIT 6.869 class

# Humble beginnings



- Perceptron [Rosenblatt '58]



- Criticism of Perceptrons (XOR affair) [Minsky Papert '69]
  - Effectively causes a "deep learning winter"

# (Early) Spring

...umelhart et al. '86, LeCun '85, Par...

• Convolutional layers [LeCun et al. '90]

...works/Long Short-Te...

...reiter Schmidhuber '97]

# Summer

- **2006:** First big suc

- **2012:** Breakthroug                                    et al. '12]



- **2015:** Deep learning-based vision models outperfo

# What enabled this success?

- Better architectures (e.g., ReLUs) and regularization techniques (e.g. Dropout)

- Sufficiently large datasets

- Enough computational power

# Deep learning

- Modeling the visual world is incredibly complicated. We need high capacity models.

- In the past, we didn't have enough data to fit these models. But now we do!

- We want a class of **high capacity models** that are **easy to optimize**.

**Deep neural networks!**

Classification units

PIT/AIT

V4/PIT

V2/V4

V1/V2

Serre, 2014

CAR   PERSON   ANIMAL   Output (object identity)

3rd hidden layer (object parts)

2nd hidden layer (corners and contours)

1st hidden layer (edges)

Visible layer (input pixels)

9

# Object recognition

Edges

Texture

Colors

Segments

Parts

"clown fish"

$\phi_k(x)$

$$f_\theta(x) = \sum_{k=1}^{K} \theta_k \phi_k(x)$$

Feature extractors

Classifier

# Object recognition

Edges

Texture

Colors

Segments

Parts

Learned

"clown fish"

$$\phi_k(x)$$

$$f_\theta(x) = \sum_{k=1}^{K} \theta_k \phi_k(x)$$

Feature extractors

Classifier

# Object recognition

Learned

"clown fish"

# Object recognition



Neural net

Learned

"clown fish"

# Object recognition



Learned

"clown fish"

Deep neural net

# Deep learning

$\mathbf{y}_i$

"clown fish"

$\mathbf{x}_i$



Loss

$\mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)$

$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6$

$$\theta^* = \arg\min_\theta \sum_{i=1}^{N} \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)$$

# Gradient descent

$$\theta^* = \arg\min_{\theta} \underbrace{\sum_{i=1}^{N} \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)}_{J(\theta)}$$

# Gradient descent



$$\theta^* = \arg \min_{\theta} J(\theta)$$

# Gradient descent

$$\theta^* = \arg\min_{\theta} \underbrace{\sum_{i=1}^{N} \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)}_{J(\theta)}$$

One iteration of gradient descent:

$$\theta^{t+1} = \theta^t - \eta_t \left.\frac{\partial J(\theta)}{\partial \theta}\right|_{\theta = \theta^t}$$

learning rate

# Computation in a neural net

Input
representation

Output
representation

# Computation in a neural net

### Linear layer

Input
representation

Output
representation

$x_i$

$w_{ij}$

$y_j$

$$y_j = \sum_i w_{ij} x_i$$

# Computation in a neural net

Linear layer

Input representation

Output representation

$x_i$

$w_{ij}$

$y_j$

$b_j$

1

weights

$$y_j = \sum_i w_{ij} x_i + b_j$$

bias

# Computation in a neural net

Linear layer

Input representation

Output representation

weights

$$y_j = \mathbf{x}^T \mathbf{w}_j + b_j$$

bias

$$\theta = \{\mathbf{W}, \mathbf{b}\}$$

parameters of the model

$\mathbf{x}$

$\mathbf{w}_j$

$y_j$

$b_j$

1

# Example: linear regression with a neural net

Linear layer

Input representation

Output representation



$$f_{\mathbf{w},b}(\mathbf{x}) = \mathbf{x}^T\mathbf{w} + b$$

# Computation in a neural net

"Perceptron"

$$g(y) = \begin{cases} 1, & \text{if} \quad y > 0 \\ 0, & \text{otherwise} \end{cases}$$

Input
representation

Output
representation

$\mathbf{x}$  $\mathbf{w}$

$y$   $g(y)$

$b$

$1$

Pointwise
Non-linearity



$g(y)$

$y$

# Example: linear classification with a perceptron



$$y = \mathbf{x}^T \mathbf{w} + b$$

# Example: linear classification with a perceptron



$$y = \mathbf{x}^T \mathbf{w} + b$$

$$g(y) = \begin{cases} 1, & \text{if} \quad y > 0 \\ 0, & \text{otherwise} \end{cases}$$

# Example: linear classification with a perceptron



$$y = \mathbf{x}^T \mathbf{w} + b$$

$$g(y) = \begin{cases} 1, & \text{if} \quad y > 0 \\ 0, & \text{otherwise} \end{cases}$$

# Example: linear classification with a perceptron

$$g(y)$$



$$y = \mathbf{x}^T \mathbf{w} + b$$

$$g(y) = \begin{cases} 1, & \text{if} \quad y > 0 \\ 0, & \text{otherwise} \end{cases}$$

# Example: linear classification with a perceptron



$$\hat{y} = \mathbf{x}^T \mathbf{w} + b$$

$$g(\hat{y}) = \begin{cases} 1, & \text{if} \quad \hat{y} > 0 \\ 0, & \text{otherwise} \end{cases}$$
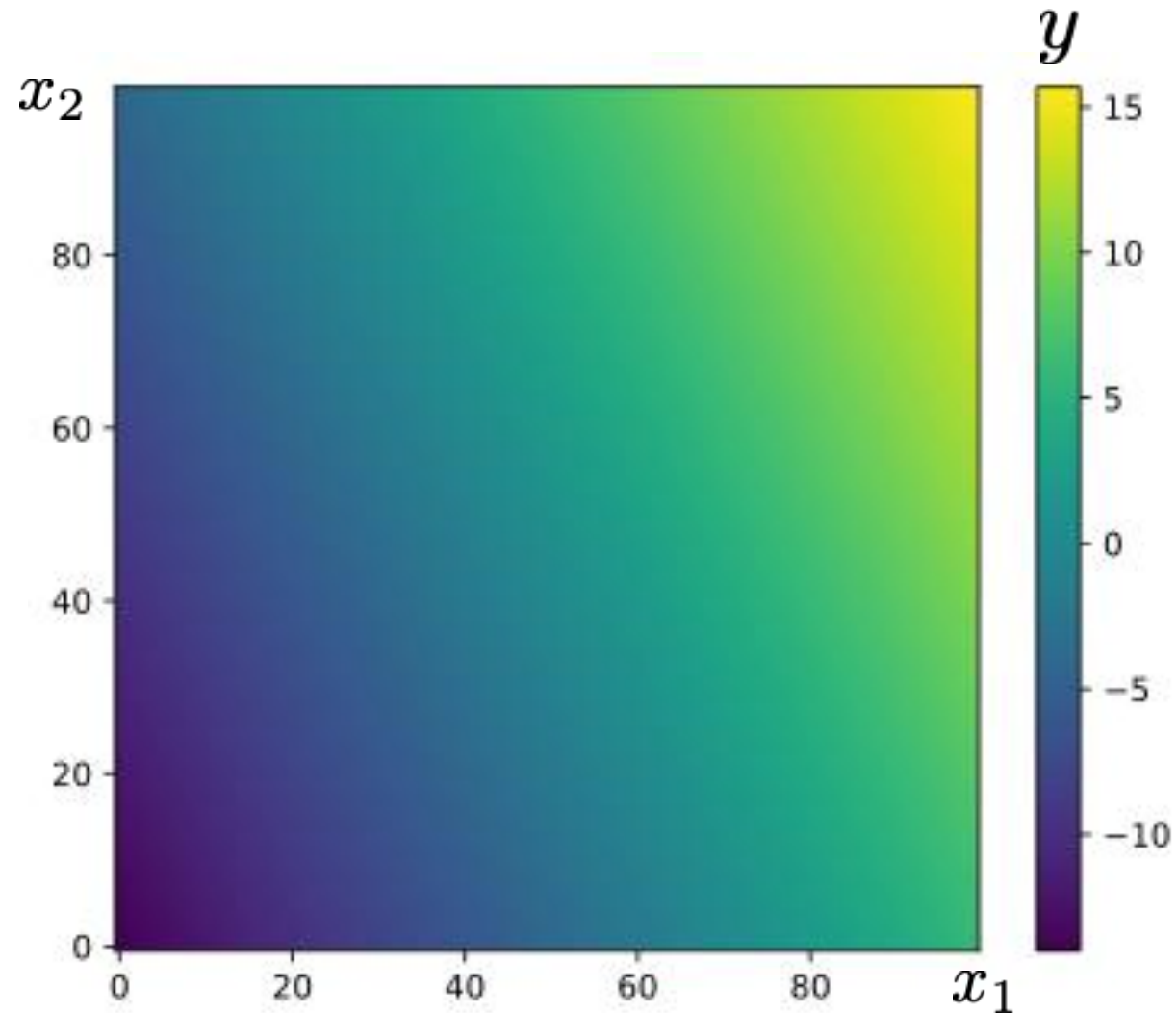
$$\mathbf{w}^*, b^* = \arg\min_{\mathbf{w},b} \mathcal{L}(g(\hat{y}), y_i)$$

# Example: linear classification with a perceptron



$$\hat{y} = \mathbf{x}^T \mathbf{w} + b$$

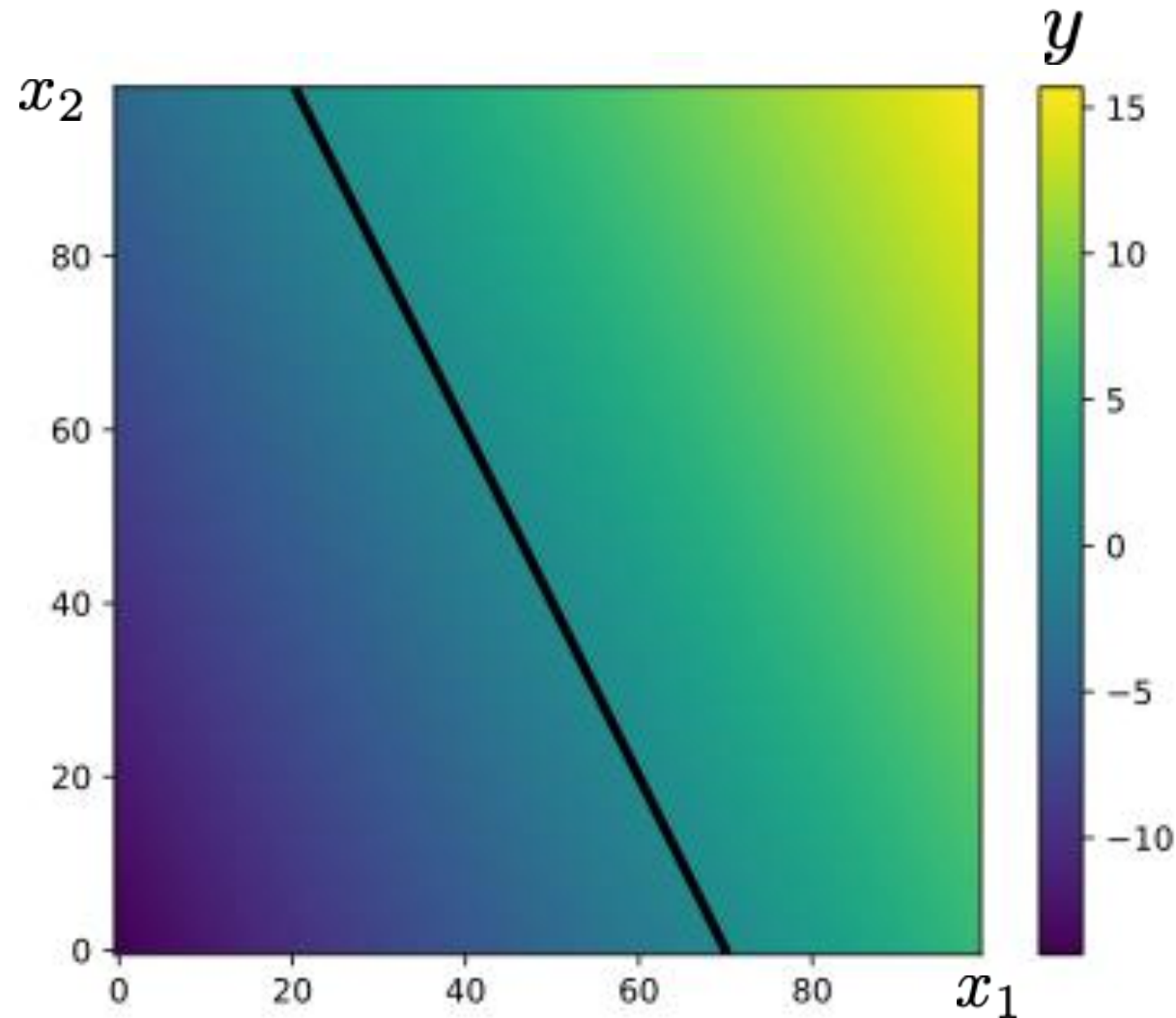$$g(\hat{y}) = \begin{cases} 1, & \text{if} \quad \hat{y} > 0 \\ 0, & \text{otherwise} \end{cases}$$

$$\mathbf{w}^*, b^* = \underset{\mathbf{w}, b}{\arg\min}\, \mathcal{L}(g(\hat{y}), y_i)$$

# Computation in a neural net



Input representation

Output representation

$\mathbf{x}$   $\mathbf{w}$   $y$   $g(y)$   $b$   $1$

$$g(y) = \begin{cases} 1, & \text{if} \quad y > 0 \\ 0, & \text{otherwise} \end{cases}$$

# Computation in a neural net – nonlinearity

Input
representation

Output
representation

Sigmoid

$$g(y) = \frac{1}{1 + e^{-y}}$$

# Computation in a neural net – nonlinearity

- Interpretation as firing rate of neuron

- Bounded between [0,1]

- Saturation for large +/- inputs

- Gradients go to zero

- Outputs centered at 0.5
  (poor conditioning)

- Not used in practice

Sigmoid

$$g(y) = \frac{1}{1 + e^{-y}}$$

$g(y)$

$y$

# Computation in a neural net – nonlinearity

- Bounded between [-1,+1]

- Saturation for large +/- inputs

- Gradients go to zero

- Outputs centered at 0

- Preferable to sigmoid

tanh(x) = 2 sigmoid(2x) −1

Tanh

$$g(y) = \frac{e^y - e^{-y}}{e^y + e^{-y}}$$

$g(y)$

# Computation in a neural net – nonlinearity

- Unbounded output (on positive side)

- Efficient to implement: $\dfrac{\partial g}{\partial y} = \begin{cases} 0, & \text{if} \quad y < 0 \\ 1, & \text{if} \quad y \geq 0 \end{cases}$

- Also seems to help convergence
(see 6x speedup vs tanh in [Krizhevsky et al.])

- Drawback: if strongly in negative region, unit is dead forever (no gradient).

- Default choice: widely used in current models.

Rectified linear unit (ReLU)

$$g(y) = \max(0, y)$$

$g(y)$

$y$

# Computation in a neural net – nonlinearity

- where α is small (e.g. 0.02)

- Efficient to implement: $\frac{\partial g}{\partial y} = \begin{cases} -a, & \text{if} \quad y < 0 \\ 1, & \text{if} \quad y \geq 0 \end{cases}$

- Also known as probabilistic ReLU (PReLU)

- Has non-zero gradients everywhere (unlike ReLU)

- α can also be learned (see Kaiming He et al. 2015).

Leaky ReLU

$$g(y) = \begin{cases} \max(0, y), & \text{if} \quad y \geq 0 \\ a \min(0, y), & \text{if} \quad y < 0 \end{cases}$$

# Stacking layers



Input representation

Intermediate representation

Output representation

$\mathbf{h}$

$\mathbf{x}$

$\mathbf{W}_j^{(1)}$

$b_j^{(1)}$

$1$

$\mathbf{h}$

$\mathbf{W}_j^{(2)}$

$b_j^{(2)}$

$1$

$\mathbf{y}$

$\mathbf{h}$ = "hidden units"

# Stacking layers



Input representation

Intermediate representation

Output representation

$$\theta = \{\mathbf{W}^{(1)}, \ldots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \ldots, \mathbf{b}^{(L)}\}$$

# Stacking layers

Input
representation

Intermediate
representation

Output
representation

$\mathbf{W}^{(1)}$

$\mathbf{h}$

$\mathbf{W}^{(2)}$

$\mathbf{x}$

positive

negative

$\mathbf{y}$

$1$

$\mathbf{b}^{(1)}$

$1$

$\mathbf{b}^{(2)}$

$$\theta = \{\mathbf{W}^{(1)}, \ldots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \ldots, \mathbf{b}^{(L)}\}$$

# Stacking layers



$$\theta = \{\mathbf{W}^{(1)}, \ldots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \ldots, \mathbf{b}^{(L)}\}$$

# Stacking layers

Input
representation

Intermediate
representation

Output
representation

$$\mathbf{W}^{(1)} \qquad \mathbf{h} \qquad \mathbf{W}^{(2)}$$

$$\mathbf{x} \qquad \qquad \mathbf{y}$$

$$1 \qquad \mathbf{b}^{(1)} \qquad 1 \qquad \mathbf{b}^{(2)}$$

positive

negative

$$\theta = \{\mathbf{W}^{(1)}, \ldots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \ldots, \mathbf{b}^{(L)}\}$$

# Stacking layers

Input
representation

Intermediate
representation

Output
representation

$\mathbf{W}^{(1)}$ $\mathbf{h}$ $\mathbf{W}^{(2)}$

$\mathbf{x}$

$1$ $\mathbf{b}^{(1)}$ $1$ $\mathbf{b}^{(2)}$
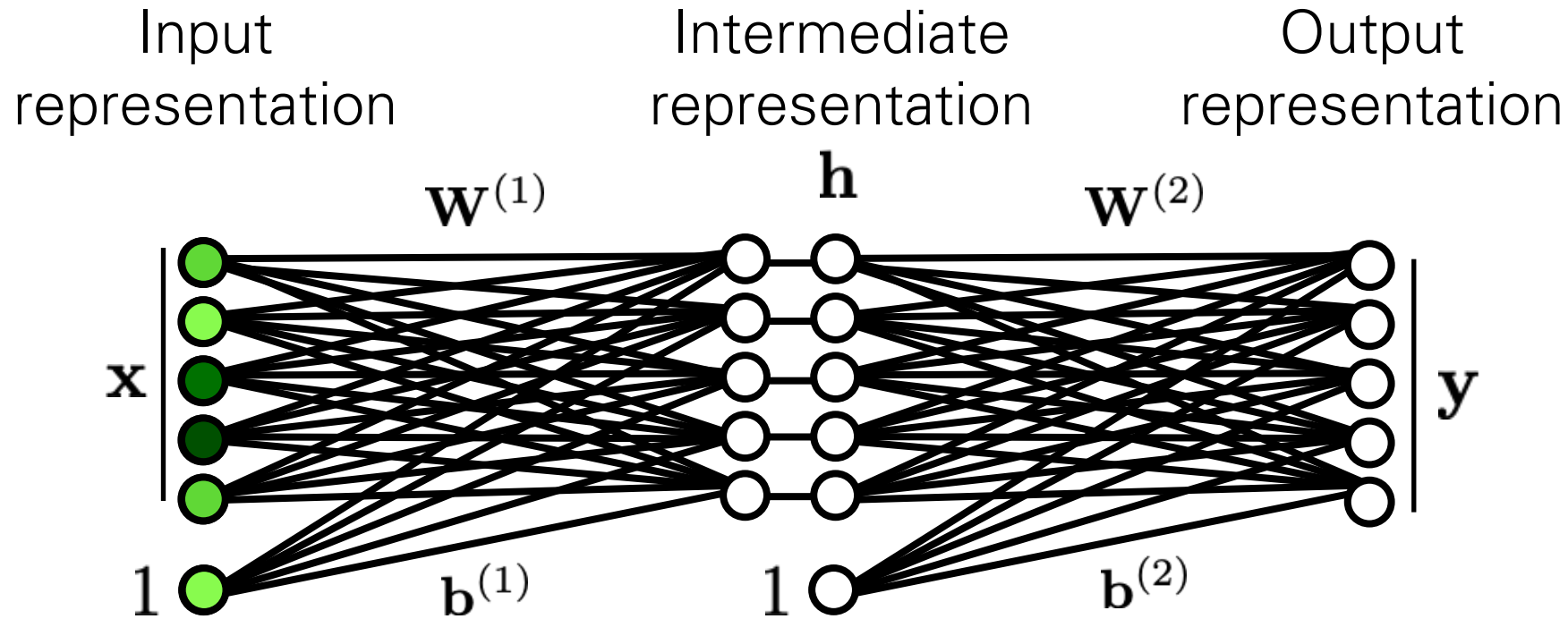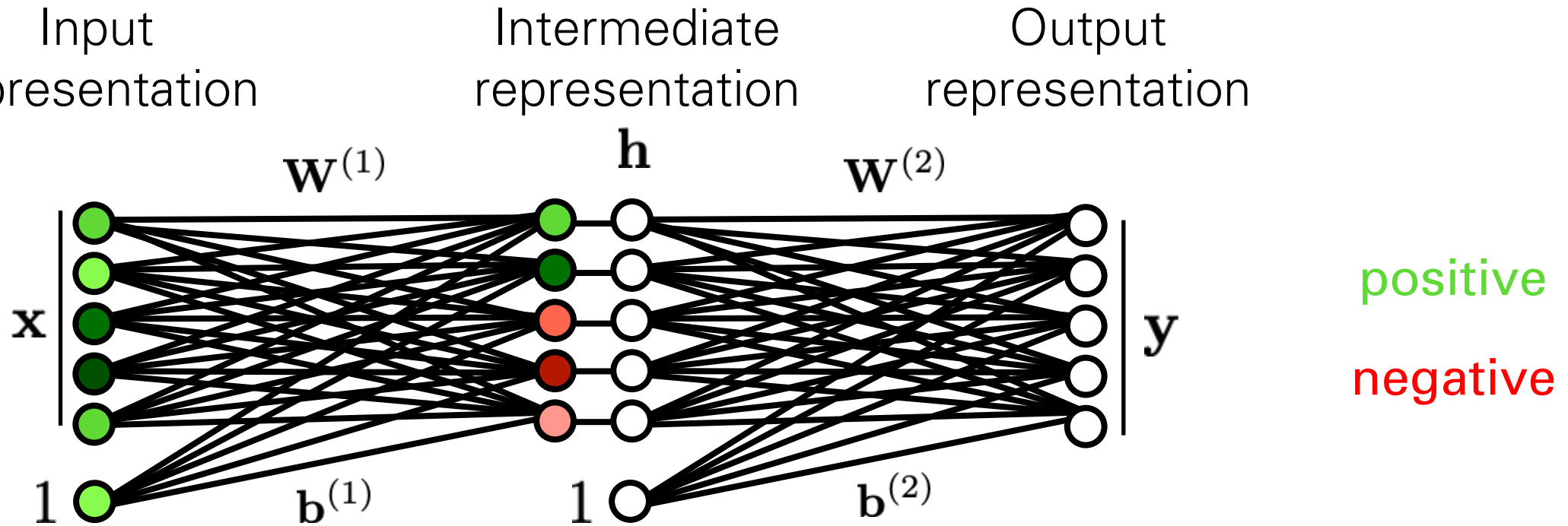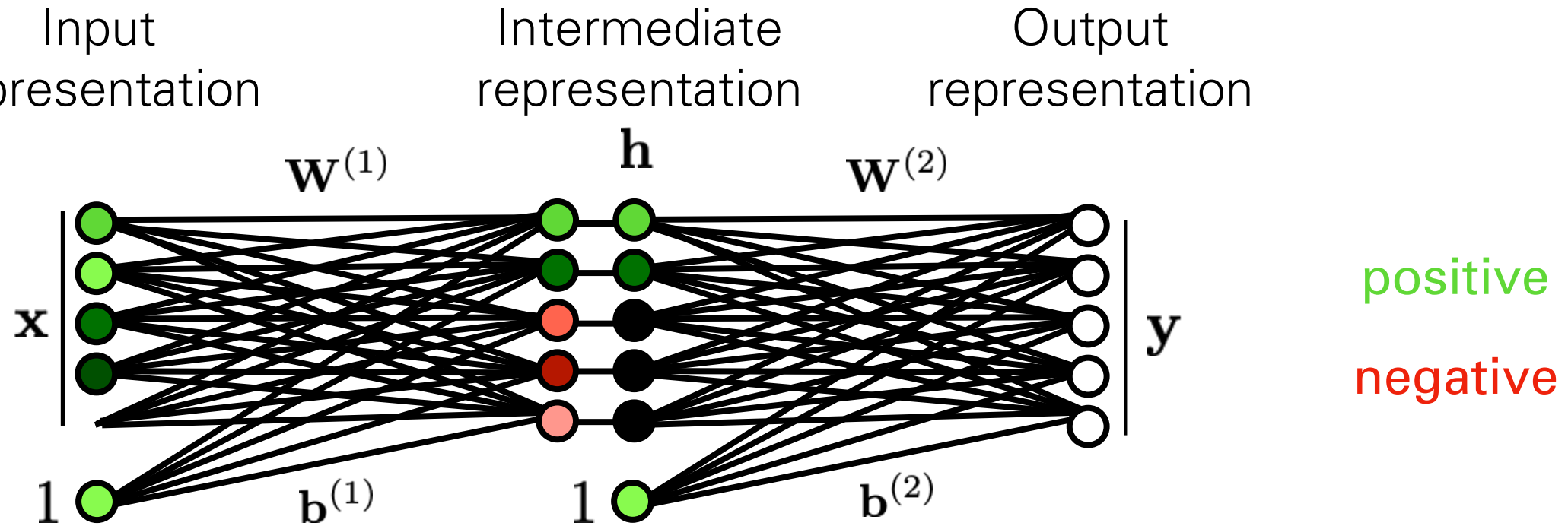
$\mathbf{y}$

positive

negative

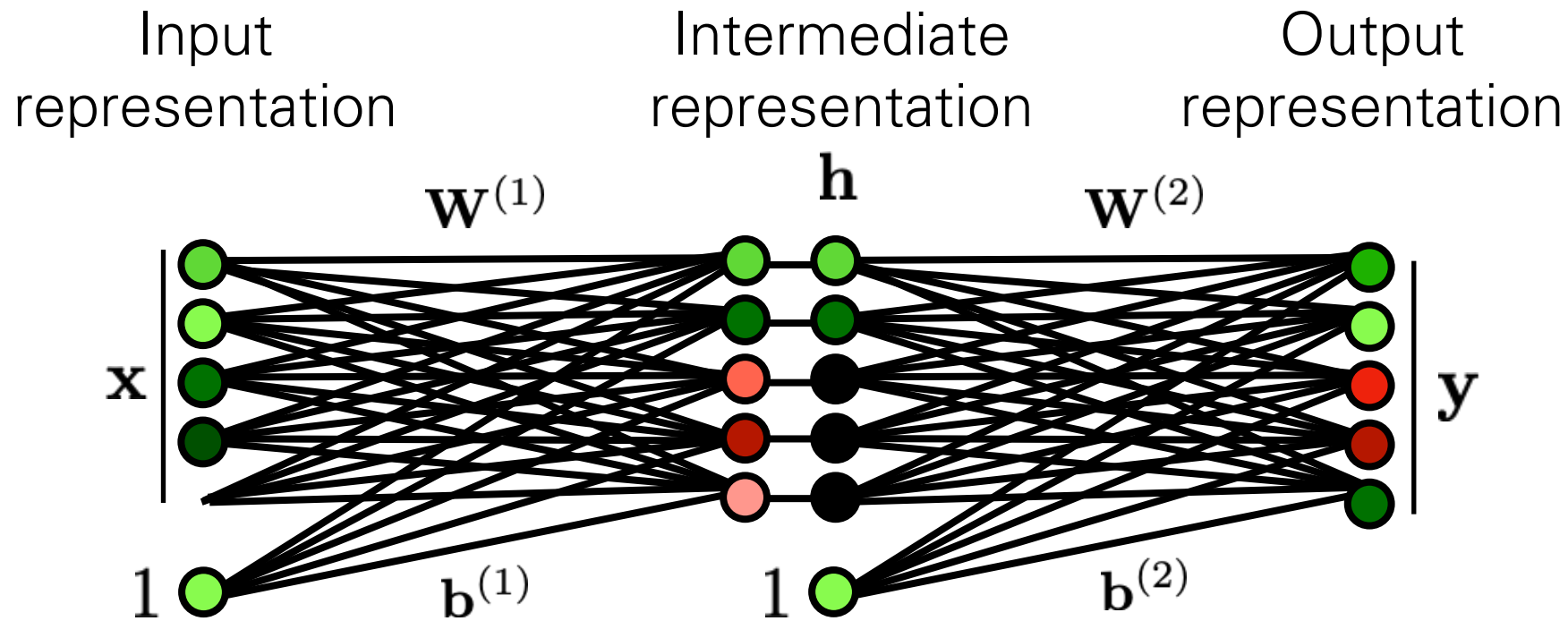$$\theta = \{\mathbf{W}^{(1)}, \ldots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \ldots, \mathbf{b}^{(L)}\}$$

# Representational power

- 1 layer? Linear decision surface.

- 2+ layers? In theory, can represent any function. Assuming non-trivial non-linearity.
  - Bengio 2009,
    http://www.iro.umontreal.ca/~bengioy/papers/ftml.pdf
  - Bengio, Courville, Goodfellow book
    http://www.deeplearningbook.org/contents/mlp.html
  - Simple proof by M. Neilsen
    http://neuralnetworksanddeeplearning.com/chap4.html
  - D. Mackay book
    http://www.inference.phy.cam.ac.uk/mackay/itprnn/ps/482.491.pdf

- But issue is efficiency: very wide two layers vs narrow deep model? In practice, more layers helps.

[http://playground.tensorflow.org]

# Deep nets



Linear

Non-linearity

Classify

"clown fish"

$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

# Classifier layer

Last layer



...

dolphin

cat

grizzly bear

angel fish

chameleon

**clown fish**

iguana

elephant

argmax → "clown fish"

# Loss function

Network output

Ground truth label



dolphin

cat

grizzly bear

angel fish

chameleon

**clown fish**

iguana

elephant

…

"clown fish"

Loss → error

# Loss function

Network output

Ground truth label

dolphin

cat

grizzly bear

angel fish

... 

chameleon

**clown fish**

iguana

elephant

"clown fish"

Loss ➡ **small**

# Loss function

Network output

Ground truth label



"grizzly bear"

dolphin

cat

grizzly bear

angel fish

... chameleon

Loss → **large**

**clown fish**

iguana

elephant

Prediction $\hat{\mathbf{y}}$

$f_\theta : X \to \mathbb{R}^K$

Ground truth label $\mathbf{y}$

## Network output    Ground truth label

$\hat{\mathbf{y}}$    $\mathbf{y}$

dolphin

cat

**grizzly bear** —

angel fish

`softmax`

chameleon

**clown fish**

iguana

elephant

Probability of the observed data under the model

$$H(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{k=1}^{K} y_k \log \hat{y}_k$$

# Deep learning

$\mathbf{y}_1$
"clown fish"

$\mathbf{x}_1$



Loss $\qquad \mathcal{L}(f_\theta(\mathbf{x}_1), \mathbf{y}_1)$

$\theta_1 \qquad \theta_2 \qquad \theta_3 \qquad \theta_4 \qquad \theta_5 \qquad \theta_6$

$$\theta^* = \arg\min_\theta \sum_{i=1}^{N} \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)$$

# Deep learning

$\mathbf{y}_2$

"grizzly bear"

$\mathbf{x}_2$



Loss

$\mathcal{L}(f_\theta(\mathbf{x}_2), \mathbf{y}_2)$

$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6$

$$\theta^* = \arg\min_\theta \sum_{i=1}^N \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)$$

# Deep learning

$\mathbf{y}_i$
"chameleon"

$\mathbf{x}_i$

Loss $\quad \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)$

$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6$

$$\theta^* = \arg\min_\theta \sum_{i=1}^{N} \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)$$

# Batch (parallel) processing

# Tensors

$$\mathbf{h}^{(1)} \in \mathbb{R}^{N_{\mathbf{batch}} \times C^{(1)}}$$

$N_{\mathbf{batch}}$

$C^{(1)}$

# "Tensor flow"

$$\mathbf{h}^{(1)} \in \mathbb{R}^{N_{\text{batch}} \times H^{(1)} \times W^{(1)} \times C^{(1)}}$$

$$\mathbf{h}^{(2)} \in \mathbb{R}^{N_{\text{batch}} \times H^{(2)} \times W^{(2)} \times C^{(2)}}$$



$H^{(1)} \times W^{(1)}$

$C^{(1)}$

$N_{\text{batch}}$

$H^{(2)} \times W^{(2)}$

$C^{(2)}$

$N_{\text{batch}}$

# Regularizing deep nets

Deep nets have millions of parameters!

On many datasets, it is easy to overfit — we may have more free parameters than data points to constrain them.

How can we regularize to prevent the network from overfitting?
1.  Fewer neurons, fewer layers
2.  Weight decay
3.  Dropout
4.  Normalization layers
5.  …

# Recall: regularized least squares

$$f_\theta(x) = \sum_{k=0}^{K} \theta_k x^k$$

$$R(\theta) = \lambda \|\theta\|_2^2 \longleftarrow$$

Only use polynomial terms if you really need them! Most terms should be zero

**ridge regression**, a.k.a., **Tikhonov regularization**

Probabilistic interpretation: R is a Gaussian **prior** over values of the parameters.

# Recall: regularized least squares

$$\theta^* = \arg\min_\theta \sum_{i=1}^{N} \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i) + R(\theta)$$

$$R(\mathbf{W}) = \lambda \|\mathbf{W}\|_2^2 \quad \longleftarrow \quad \text{weight decay}$$

"We prefer to keep weights small."

# Dropout

Input representation

Intermediate representation

Output representation



$$\theta = \{\mathbf{W}^{(1)}, \ldots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \ldots, \mathbf{b}^{(L)}\}$$

# Dropout



Input representation

Intermediate representation

Output representation

$$\theta = \{\mathbf{W}^{(1)}, \ldots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \ldots, \mathbf{b}^{(L)}\}$$

# Dropout

Input
representation

Intermediate
representation

Output
representation



$$\theta = \{\mathbf{W}^{(1)}, \ldots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \ldots, \mathbf{b}^{(L)}\}$$

# Dropout

Input representation

Intermediate representation
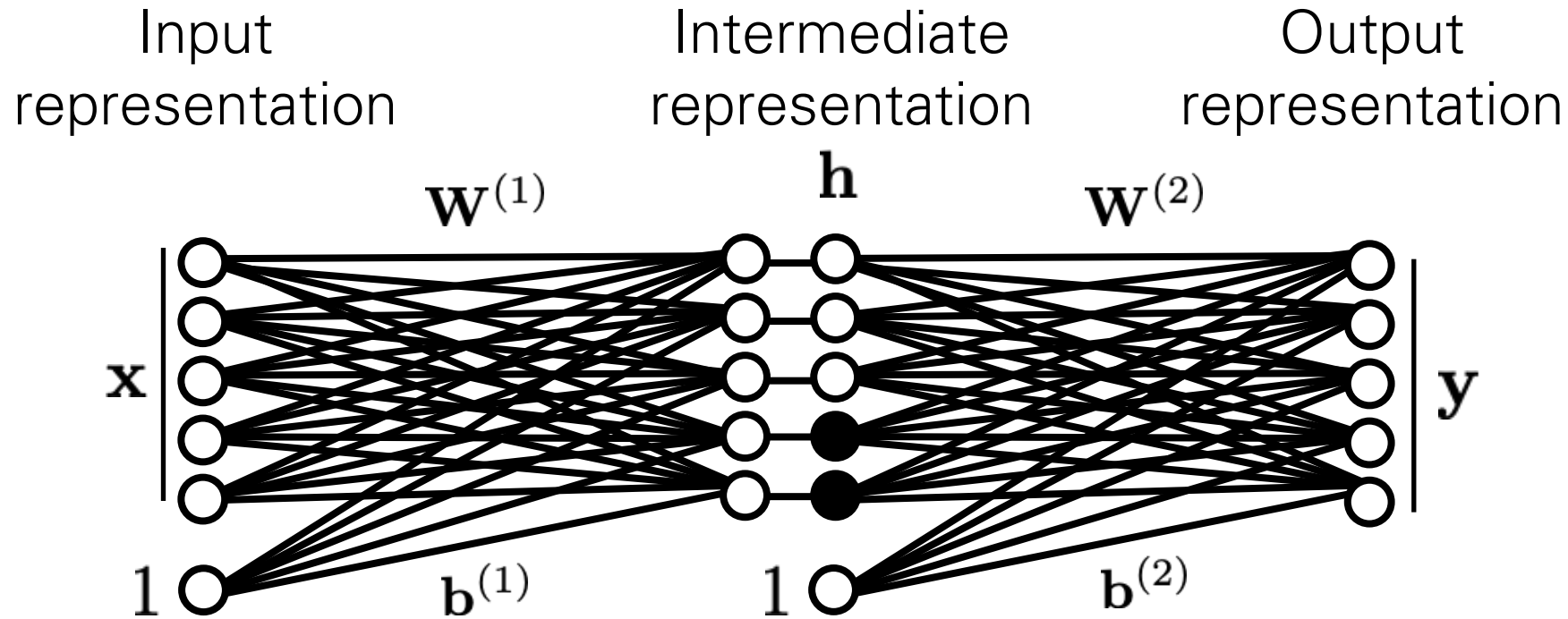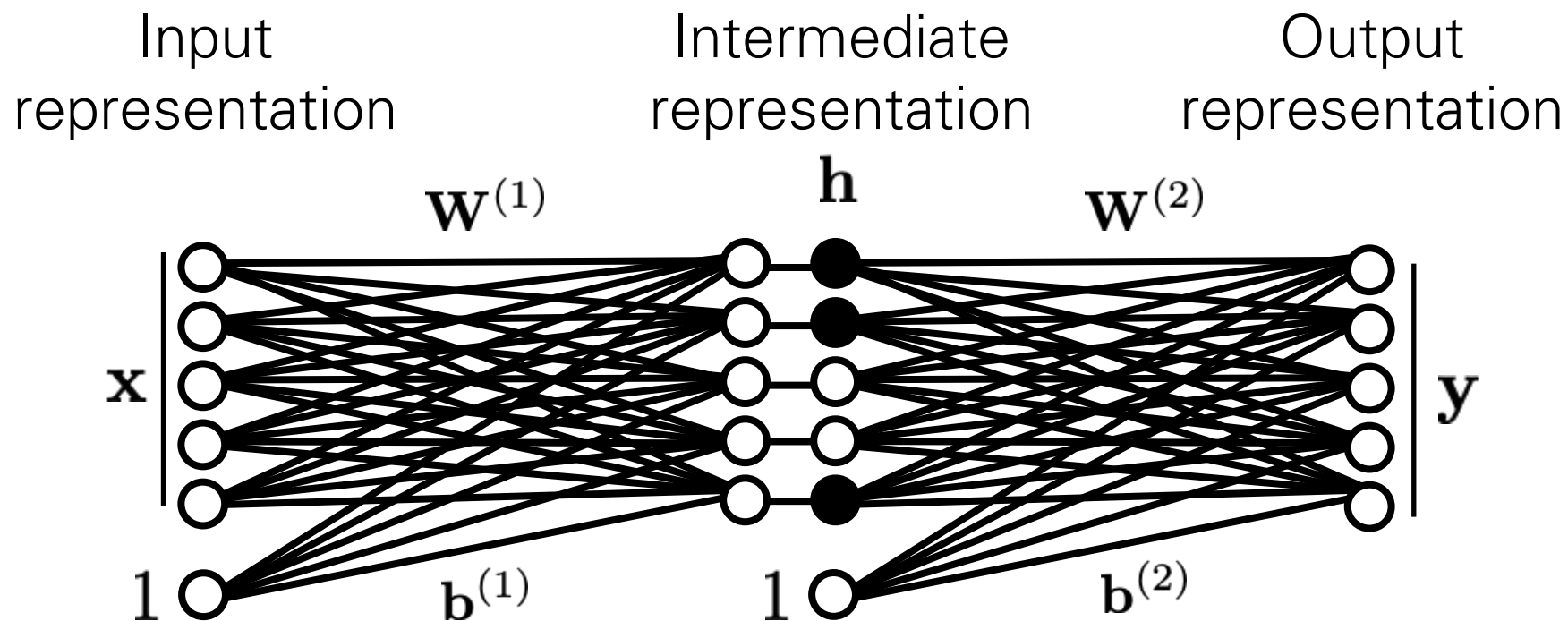
Output representation

$$\mathbf{W}^{(1)} \qquad \mathbf{h} \qquad \mathbf{W}^{(2)}$$



$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$
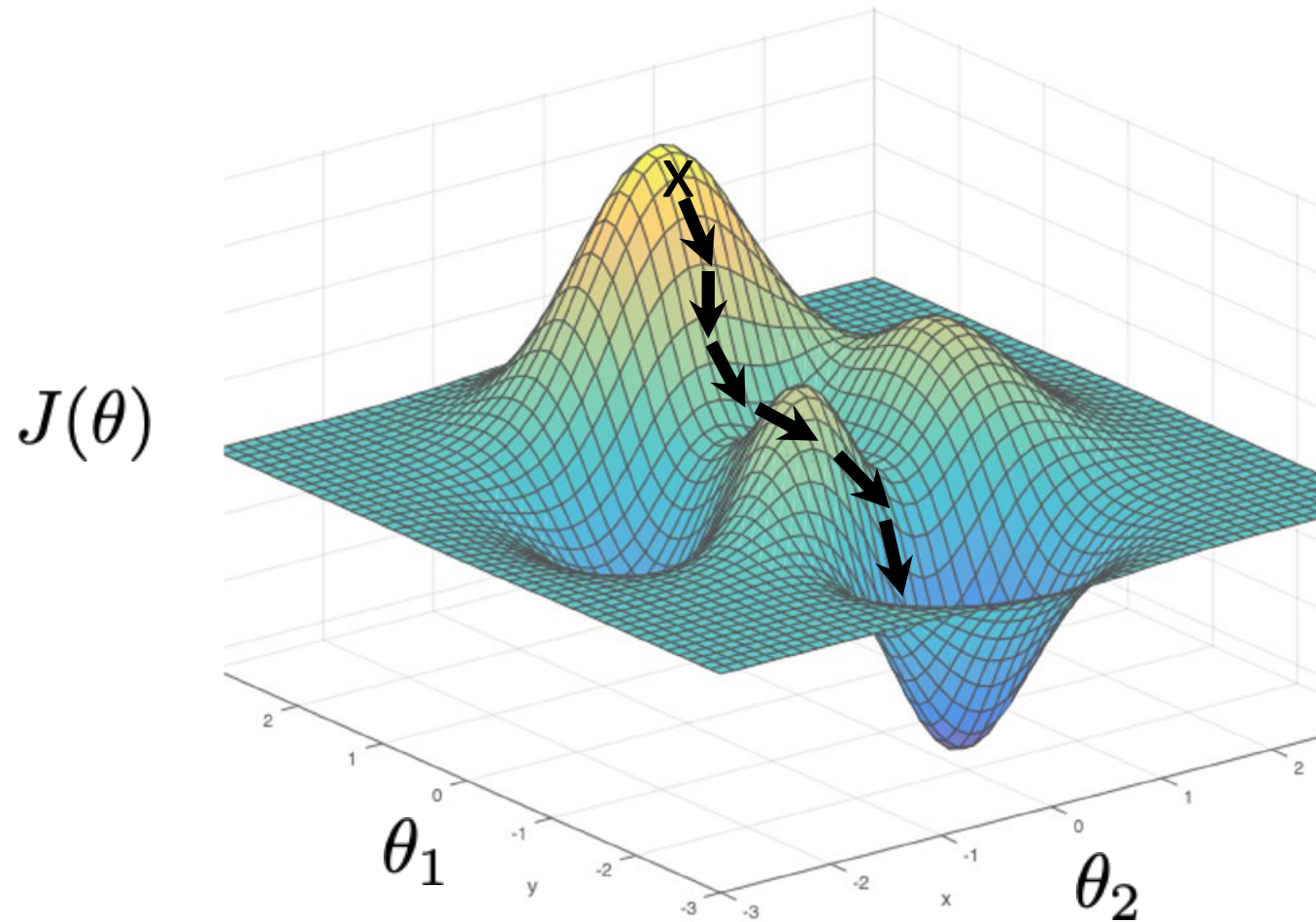
# Dropout

Randomly zero out hidden units.

Prevents network from relying too much on spurious correlations between different hidden units.

Can be understood as averaging over an exponential **ensemble** of subnetworks. This averaging smooths the function, thereby reducing the effective capacity of the network.

# Gradient descent



$J(\theta)$

$\theta_1$

$\theta_2$

$$\theta^* = \arg\min_\theta J(\theta)$$

# Gradient descent

$$\theta^* = \arg\min_{\theta} \sum_{i=1}^{N} \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{J(\theta)}$$

One iteration of gradient descent:

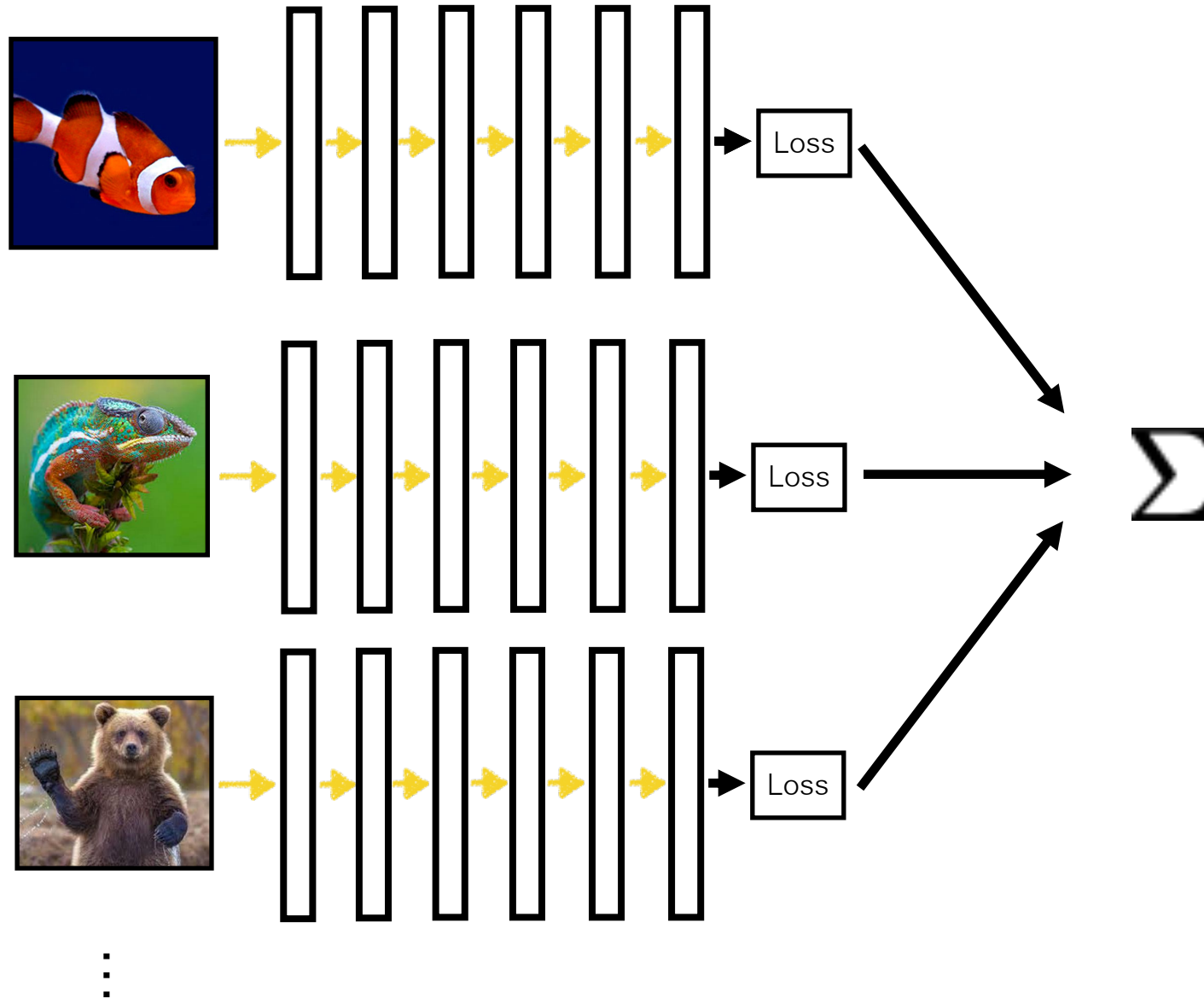$$\theta^{t+1} = \theta^t - \eta_t \left.\frac{\partial J(\theta)}{\partial \theta}\right|_{\theta=\theta^t}$$

learning rate

# Optimization

Params

$\theta$ ➡️ $\boxed{J}$ ➡️ $J(\theta)$
$\nabla_\theta J(\theta)$
$H_\theta(J(\theta))$

$$\theta^* = \arg\min_\theta J(\theta)$$

- What's the knowledge we have about J?
  - We can evaluate $J(\theta)$ ⟶ Gradient
  - We can evaluate $J(\theta)$ and $\nabla_\theta J(\theta)$ ⬅️ Black box optimization
  - We can evaluate $J(\theta)$, $\nabla_\theta J(\theta)$, and $H_\theta(J(\theta))$ ⬅️ First order optimization
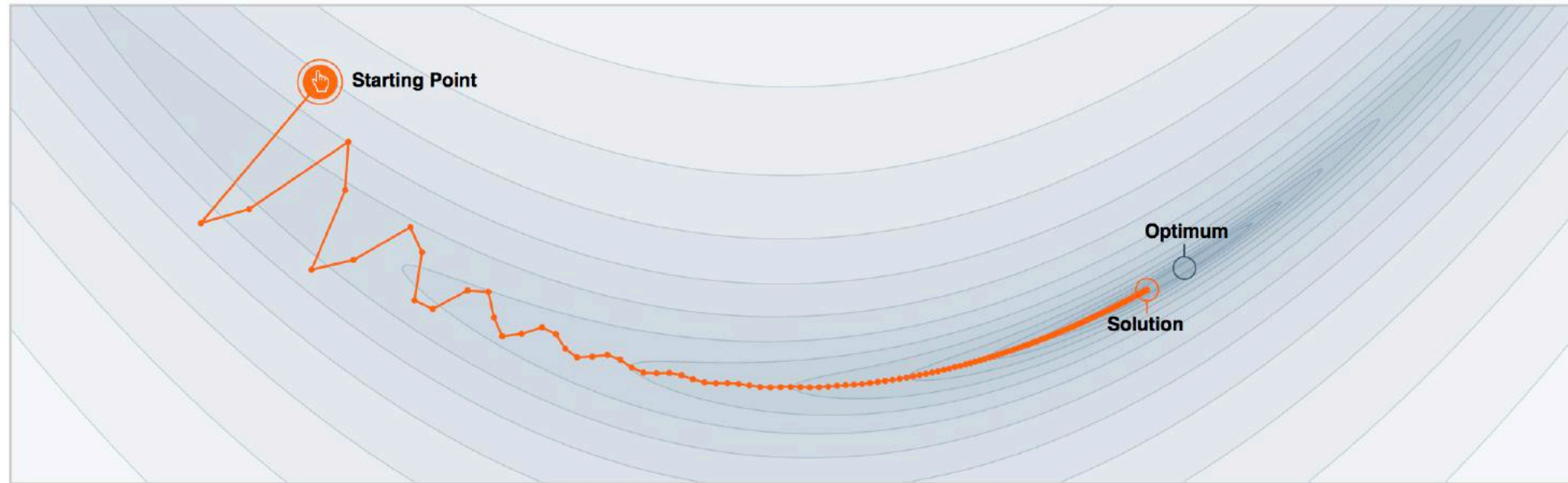  ⟶ Hessian ⬅️ Second order optimization

68

# Batch (parallel) processing

# Stochastic gradient descent (SGD)

- Want to minimize overall loss function $J$, which is sum of individual losses over each example.

- In Stochastic gradient descent, compute gradient on sub-set (batch) of data.

  If batchsize=1 then $\theta$ is updated after each example.

  If batchsize=N (full set) then this is standard gradient descent.

- Gradient direction is noisy, relative to average over all examples (standard gradient descent).

- Advantages

  - Faster: approximate total gradient with small sample

  - Implicit regularizer

- Disadvantages

  - High variance, unstable updates

# Momentum

- Basic idea: like a ball rolling down a hill, we should build up speed so as to make faster progress when "on a roll"

- Can dampen oscillations in SGD updates

- Common in popular variants of SGD

  - Nesterov's method

  - RMSProp

  - Adam

# Why Momentum Really Works



**Step-size α = 0.02**

0     0.003     0.006

**Momentum β = 0.99**

0.00     0.500     0.990

We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?
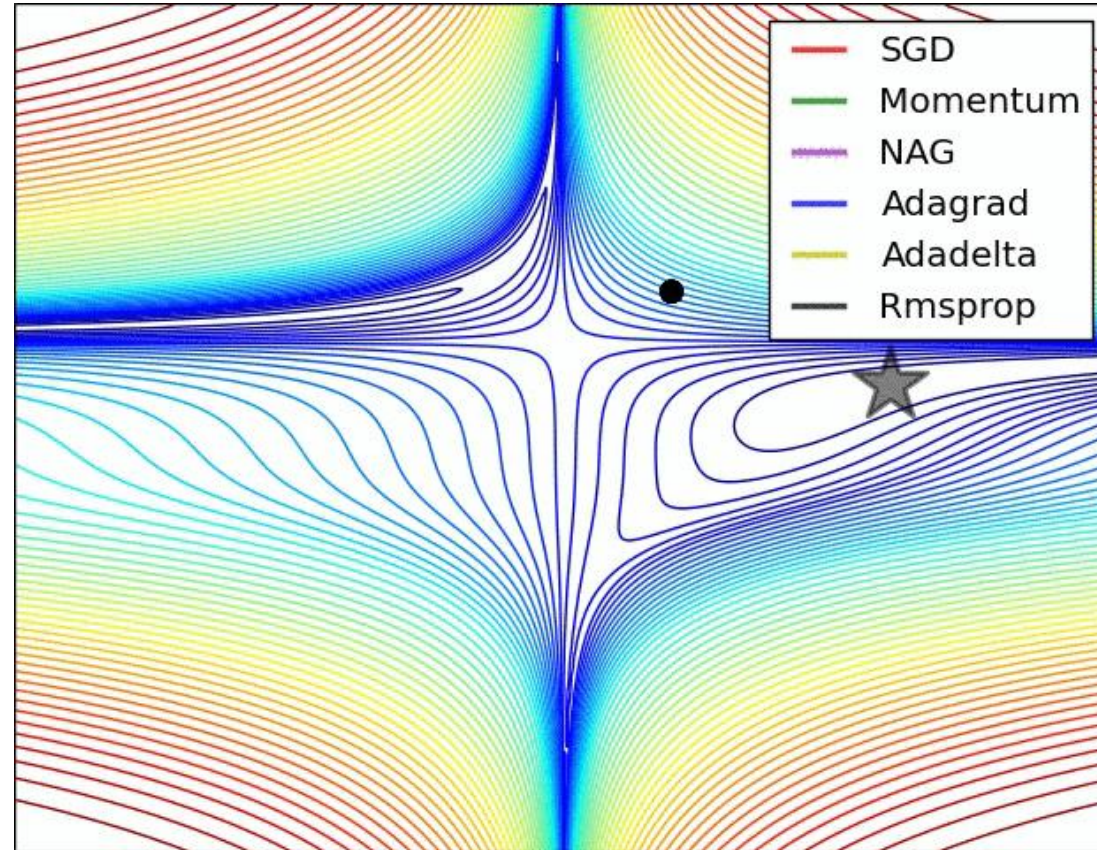
GABRIEL GOH    April. 4    Citation:
UC Davis       2017     Goh, 2017

[https://distill.pub/2017/momentum/]

# Comparison of gradient descent variants



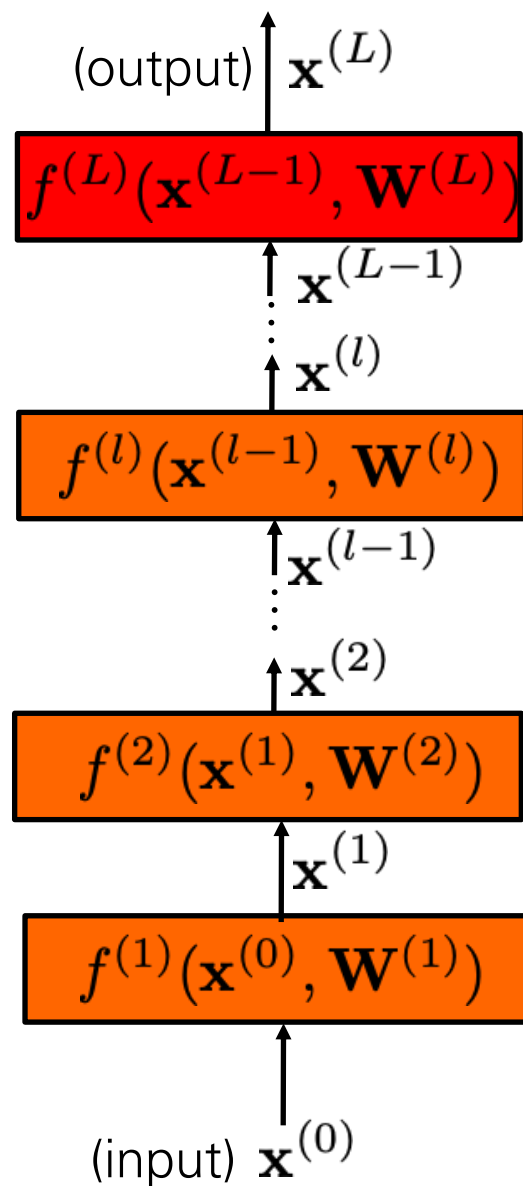[http://ruder.io/optimizing-gradient-descent/](http://ruder.io/optimizing-gradient-descent/)
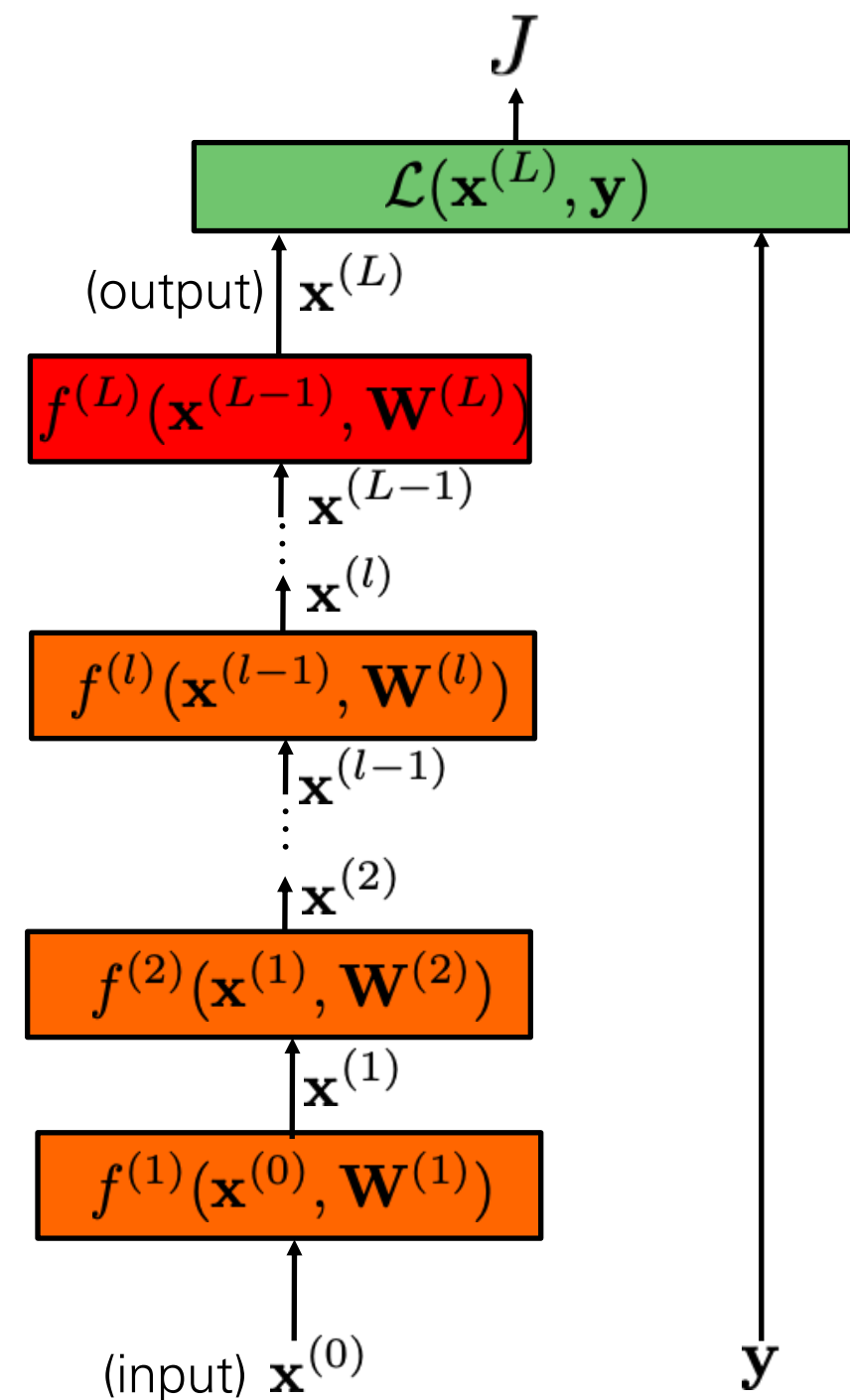
# Forward pass

- Consider model with $L$ layers. Layer $l$ has vector of weights $\mathbf{W}^{(l)}$

- **Forward pass:** takes input $\mathbf{x}^{(l-1)}$ and passes it through each layer $f^{(l)}$:

$$\mathbf{x}^{(l)} = f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})$$

- Output of layer $l$ is $\mathbf{x}^{(l)}$.

- Network output (top layer) is $\mathbf{x}^{(L)}$.

(output) $\mathbf{x}^{(L)}$

$f^{(L)}(\mathbf{x}^{(L-1)}, \mathbf{W}^{(L)})$

$\mathbf{x}^{(L-1)}$

$\mathbf{x}^{(l)}$

$f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})$

$\mathbf{x}^{(l-1)}$

$\mathbf{x}^{(2)}$

$f^{(2)}(\mathbf{x}^{(1)}, \mathbf{W}^{(2)})$

$\mathbf{x}^{(1)}$

$f^{(1)}(\mathbf{x}^{(0)}, \mathbf{W}^{(1)})$

(input) $\mathbf{x}^{(0)}$

# Forward pass

- Consider model with $L$ layers. Layer $l$ has vector of weights $\mathbf{W}^{(l)}$

- **Forward pass:** takes input $\mathbf{x}^{(l-1)}$ and passes it through each layer $f^{(l)}$:

$$\mathbf{x}^{(l)} = f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})$$

- Output of layer $l$ is $\mathbf{x}^{(l)}$.

- Network output (top layer) is $\mathbf{x}^{(L)}$.

- **Loss function** $\mathcal{L}$ compares $\mathbf{x}^{(L)}$ to $\mathbf{y}$.

- Overall energy is the sum of the cost over all training examples:

$$J = \sum_{i=1}^{N} \mathcal{L}(\mathbf{x}_i^{(L)}, \mathbf{y}_i)$$

# Gradient descent

- We need to compute gradients of the cost with respect to model parameters $\mathbf{W}^{(l)}$.

- By design, each layer is differentiable with respect to its parameters and input.

# Computing gradients

To compute the gradients, we could start by writing the full energy J as a function of the network parameters.

$$J(\mathbf{W}) = \sum_{i=1} \mathcal{L}(f^{(L)}(\dots f^{(2)}(f^{(1)}(\mathbf{x}_i^{(0)}, \mathbf{W}^{(1)}), \mathbf{W}^{(2)}), \dots \mathbf{W}^{(L)}), \mathbf{y}_i)$$

And then compute the partial derivatives… instead, we can use the chain rule to derive a compact algorithm:
**backpropagation**

# Computing gradients

The energy J is the sum of the losses associated to each training example $\{\mathbf{x}_i^{(0)}, \mathbf{y}_i\}$

$$J(\mathbf{W}) = \sum_{i=1}^{N} \mathcal{L}(\mathbf{x}_i^{(L)}, \mathbf{y}_i; \mathbf{W})$$

Its gradient with respect to each of the network's parameters w is:

$$\frac{\partial J(\mathbf{W})}{\partial w} = \sum_{i=1}^{N} \frac{\partial \mathcal{L}(\mathbf{x}_i^{(L)}, \mathbf{y}_i; \mathbf{W})}{\partial w}$$

is how much J varies when the parameter w is varied.

# Computing gradients

We could write the loss function to get the gradients as:

$$\mathcal{L}(\mathbf{x}^{(L)}, \mathbf{y}; \mathbf{W}) = \mathcal{L}(f^{(L)}(\mathbf{x}^{(L-1)}, \mathbf{W}^{(L)}), \mathbf{y})$$

If we compute the gradient with respect to the parameters of the last layer (output layer) W$^{(L)}$, using the **chain rule**:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial \mathbf{x}^{(L)}}{\partial \mathbf{W}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial f^{(L)}(\mathbf{x}^{(L-1)}, \mathbf{W}^{(L)})}{\partial \mathbf{W}^{(L)}}$$

(How much the cost changes when we change W$^{(L)}$ is the product between how much the loss changes when we change the output of the last layer and how much the output changes when we change the layer parameters.)

# Computing gradients: loss layer

If we compute the gradient with respect to the parameters of the last layer (output layer) W$^{(L)}$, using the chain rule:
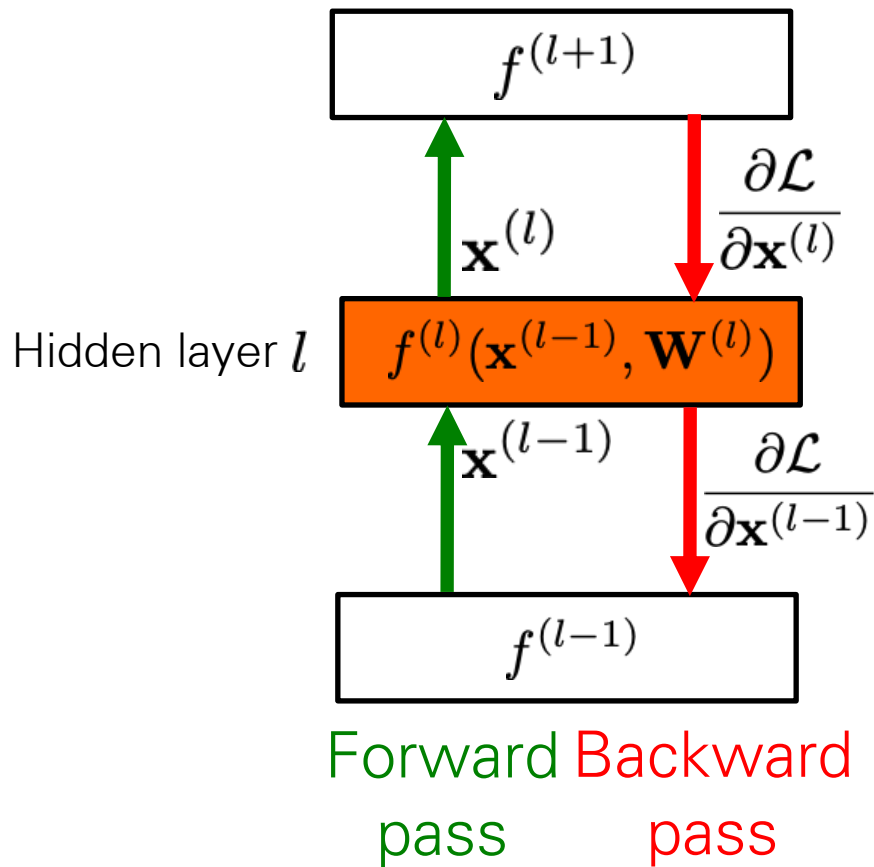
$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial \mathbf{x}^{(L)}}{\partial \mathbf{W}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial f^{(L)}(\mathbf{x}^{(L-1)}, \mathbf{W}^{(L)})}{\partial \mathbf{W}^{(L)}}$$
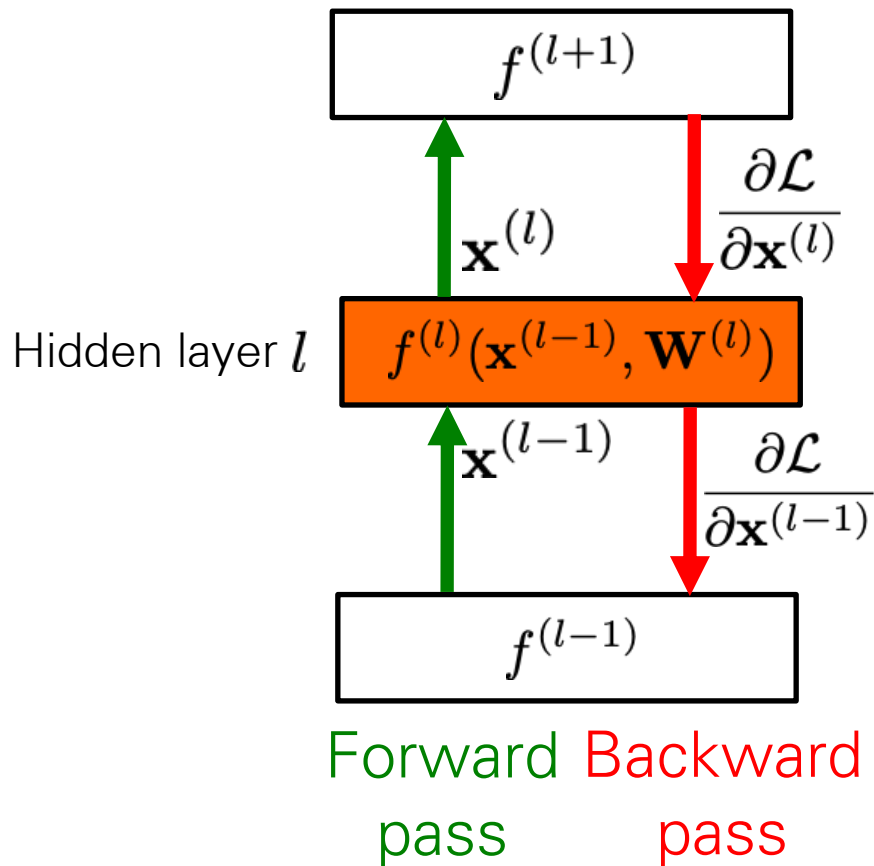
For example, for an Euclidean loss:

$$\mathcal{L}(\mathbf{x}^{(L)}, \mathbf{y}) = \frac{1}{2} \left\| \mathbf{x}^{(L)} - \mathbf{y} \right\|_2^2$$

Will depend on the layer structure and non-linearity.

The gradient is:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} = \mathbf{x}^{(L)} - \mathbf{y}$$

# Computing gradients: layer $l$

We could write the full loss function to get the gradients:

$$\mathcal{L}(\mathbf{x}^{(L)}, \mathbf{y}; \mathbf{W}) = \mathcal{L}(f^{(L)}(\ldots f^{(2)}(f^{(1)}(\mathbf{x}^{(0)}, \mathbf{W}^{(1)}), \mathbf{W}^{(2)}), \ldots \mathbf{W}^{(L)}), \mathbf{y})$$

If we compute the gradient with respect to $w_i$, using the chain rule:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial \mathbf{x}^{(L)}}{\partial \mathbf{x}^{(L-1)}} \cdot \frac{\partial \mathbf{x}^{(L-1)}}{\partial \mathbf{x}^{(L-2)}} \cdots \frac{\partial \mathbf{x}^{(l+1)}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{W}^{(l)}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}}$$

And this can be
computed iteratively!

$$\frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{W}^{(l)}}$$

This is easy.

# Backpropagation

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial \mathbf{x}^{(L)}}{\partial \mathbf{x}^{(L-1)}} \cdot \frac{\partial \mathbf{x}^{(L-1)}}{\partial \mathbf{x}^{(L-2)}} \cdots \frac{\partial \mathbf{x}^{(l+1)}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{W}^{(l)}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}}$$

$$\frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{W}^{(l)}}$$

If we have the value of $\dfrac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}}$ we can compute the gradient at the layer below as:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l-1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{x}^{(l-1)}}$$

Gradient layer l-1

Gradient layer l

$$\frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{x}^{(l-1)}}$$

# Backpropagation

— Goal: to update parameters of layer $l$

- Layer $l$ has two inputs (during training)

$$\mathbf{x}^{(l-1)} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}}$$

- We compute the outputs

$$\mathbf{x}^{(l)} = f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l-1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{x}^{(l-1)}}$$

- To compute the output, we need:

$$\frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{x}^{(l-1)}}$$

- To compute the weight update, we need:

$$\frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{W}^{(l)}}$$

$f^{(l+1)}$

$$\mathbf{x}^{(l)} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}}$$

Hidden layer $l$  $f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})$

$$\mathbf{x}^{(l-1)} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l-1)}}$$

$f^{(l-1)}$

Forward Backward
pass        pass

# Backpropagation — Goal: to update parameters of layer $l$

Hidden layer $l$

$f^{(l+1)}$

$\mathbf{x}^{(l)}$  $\dfrac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}}$

$f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})$

$\mathbf{x}^{(l-1)}$  $\dfrac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l-1)}}$

$f^{(l-1)}$

Forward pass   Backward pass

- Layer $l$ has two inputs (during training)

$\mathbf{x}^{(l-1)}$

$\dfrac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}}$

- We compute the outputs

$$\mathbf{x}^{(l)} = f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l-1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{x}^{(l-1)}}$$

- The weight update equation is:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(l)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{W}^{(l)}}$$

$$\mathbf{W}^{(l)} \leftarrow \mathbf{W}^{(l)} + \eta \left( \frac{\partial J}{\partial \mathbf{W}^{(l)}} \right)^T$$

(sum over all training examples to get J)

# Backpropagation Summary

- Forward pass: for each training example, compute the outputs for all layers:

$$\mathbf{x}^{(l)} = f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})$$

- Backwards pass: compute loss derivatives iteratively from top to bottom:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l-1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{x}^{(l-1)}}$$

- Compute gradients w.r.t. weights, and update weights:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(l)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{W}^{(l)}}$$

# Differentiable programming

Deep nets are popular for a few reasons:
1. High capacity
2. Easy to optimize (differentiable)
3. Compositional "block based programming"

An emerging term for general models with these properties is **differentiable programming.**

**Yann LeCun**
January 5 · 🌐

OK, Deep Learning has outlived its usefulness as a buzz-phrase. Deep Learning est mort. Vive Differentiable Programming!

**Thomas G. Dietterich**
@tdietterich          Following

DL is essentially a new style of programming--"differentiable programming"--and the field is trying to work out the reusable constructs in this style. We have some: convolution, pooling, LSTM, GAN, VAE, memory units, routing units, etc. 8/

8:02 AM - 4 Jan 2018

65 Retweets  194 Likes

6        65        194

# Differentiable programming

Deep learning

Differentiable programming

# Convolutional Neural Networks

# Convolutional Neural Networks

LeCun et al. 1989

Neural network with specialized connectivity

Tailored to processing natural signals with a grid topology (e.g., images).

# Image classification



image **x**

Classifier

"Fish"

label y

Photo credit: Fredo Durand

Classifier → "Bird"

Classifier → "Bird"

Bird

Classifier → "Sky"

| Sky | Sky | Sky | Sky | Sky | Sky | Sky | Bird |
|-----|-----|-----|-----|-----|-----|-----|------|
| Sky | Sky | Sky | Sky | Sky | Sky | Sky | Sky |
| Sky | Sky | Sky | Sky | Sky | Sky | Sky | Sky |
| Bird | Bird | Bird | Sky | Bird | Sky | Sky | Sky |
| Sky | Sky | Sky | Bird | Sky | Sky | Sky | Sky |

What's the object class of the center pixel?

"Bird"

"Bird"

"Sky"

"Sky"

Training data

**x**     *y*

{ "Bird" },

{ "Bird" },

{ "Sky" }
⋮

What's the object class of the center pixel?

f → "Bird"

f → "Bird"

f → "Sky"

f → "Sky"

(Colors represent one-hot codes)

This problem is called **semantic segmentation**

What's the object class of the center pixel?

$f$ → "Bird"

$f$ → "Bird"

$f$ → "Sky"

$f$ → "Sky"

Translation invariance: process each patch in the same way.

An equivariant mapping:

$f(\mathbf{translate}(x)) = \mathbf{translate}(f(x))$

**W** computes a weighted sum of all pixels in the patch





**W** $\rightarrow$

**W** $\rightarrow$

**W** $\rightarrow$

**W** is a convolutional kernel applied to the full image!

# Convolution



filter

# Fully-connected network

Fully-connected (fc) layer

# Locally connected network



Often, we assume output is a **local** function of input.

If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

# Convolutional neural network

Conv layer



Often, we assume output is a **local** function of input.

If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

# Weight sharing

Conv layer



Often, we assume output is a **local** function of input.

If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

**Toeplitz matrix**

$$\begin{pmatrix} a & b & c & d & e \\ f & a & b & c & d \\ g & f & a & b & c \\ h & g & f & a & b \\ i & h & g & f & a \end{pmatrix}$$

$$\mathbf{x}^{(l+1)} = \quad \ast \quad \mathbf{x}^{(l)}$$



e.g., pixel image

- Constrained linear layer (infinitely strong regularization)
- Fewer parameters —> easier to learn, less overfitting

$$\mathbf{x}^{(l+1)} = \quad * \quad \mathbf{x}^{(l)}$$

$$\mathbf{x}^{(l+1)} = \mathbf{M} * \mathbf{x}^{(l)}$$

Conv layers can be applied to arbitrarily-sized inputs
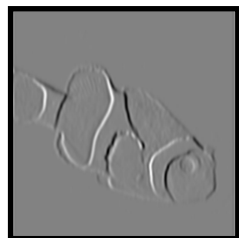
# Five views on convolutional layers

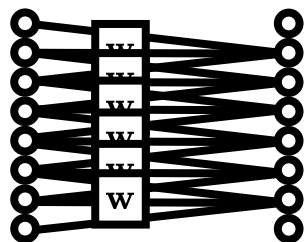1. Equivariant with translation (stationarity)   $f(\mathbf{translate}(x)) = \mathbf{translate}(f(x))$

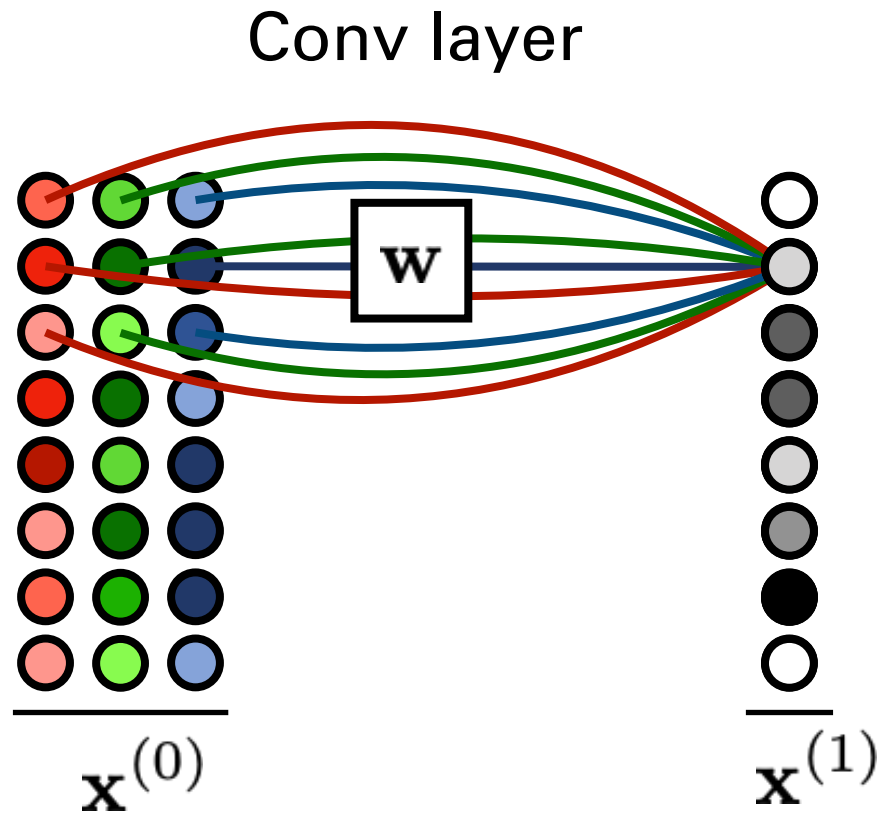2. Patch processing (Markov assumption)

3. Image filter

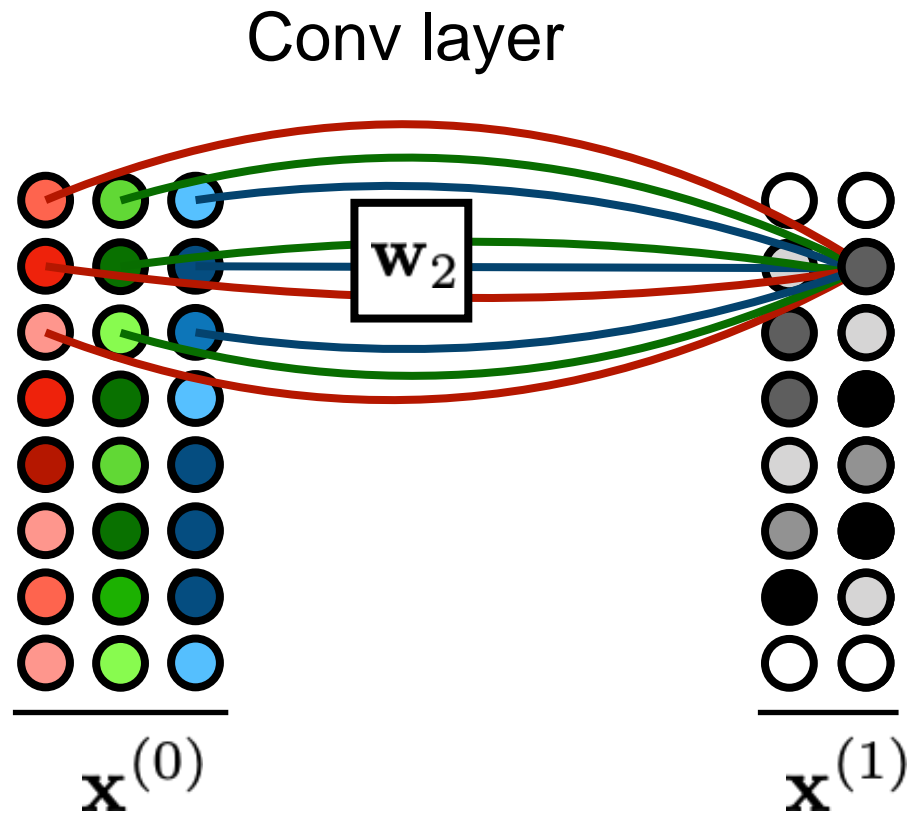4. Parameter sharing

5. A way to process variable-sized tensors

# Multiple channels



Conv layer

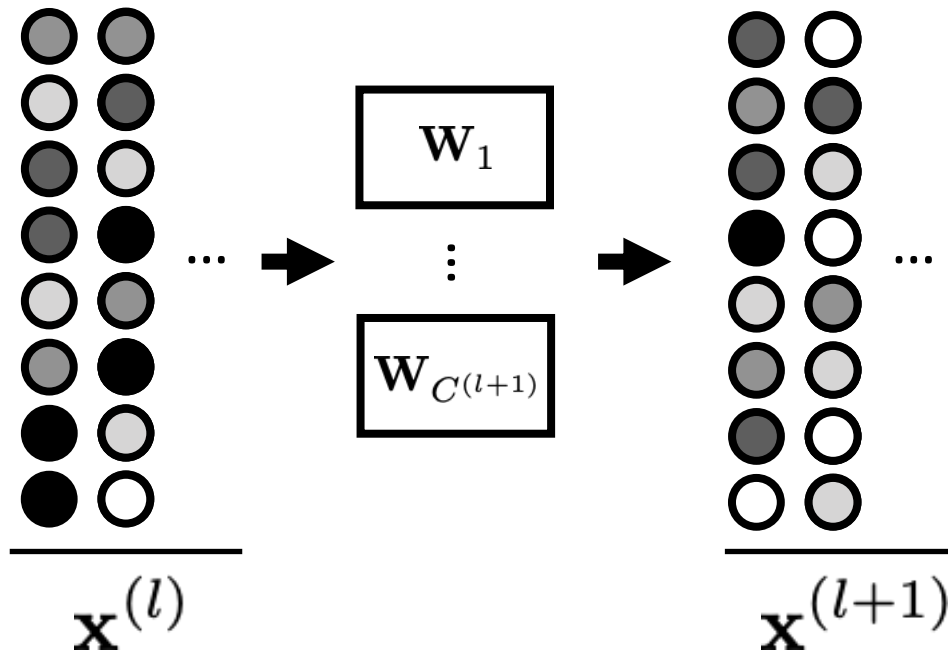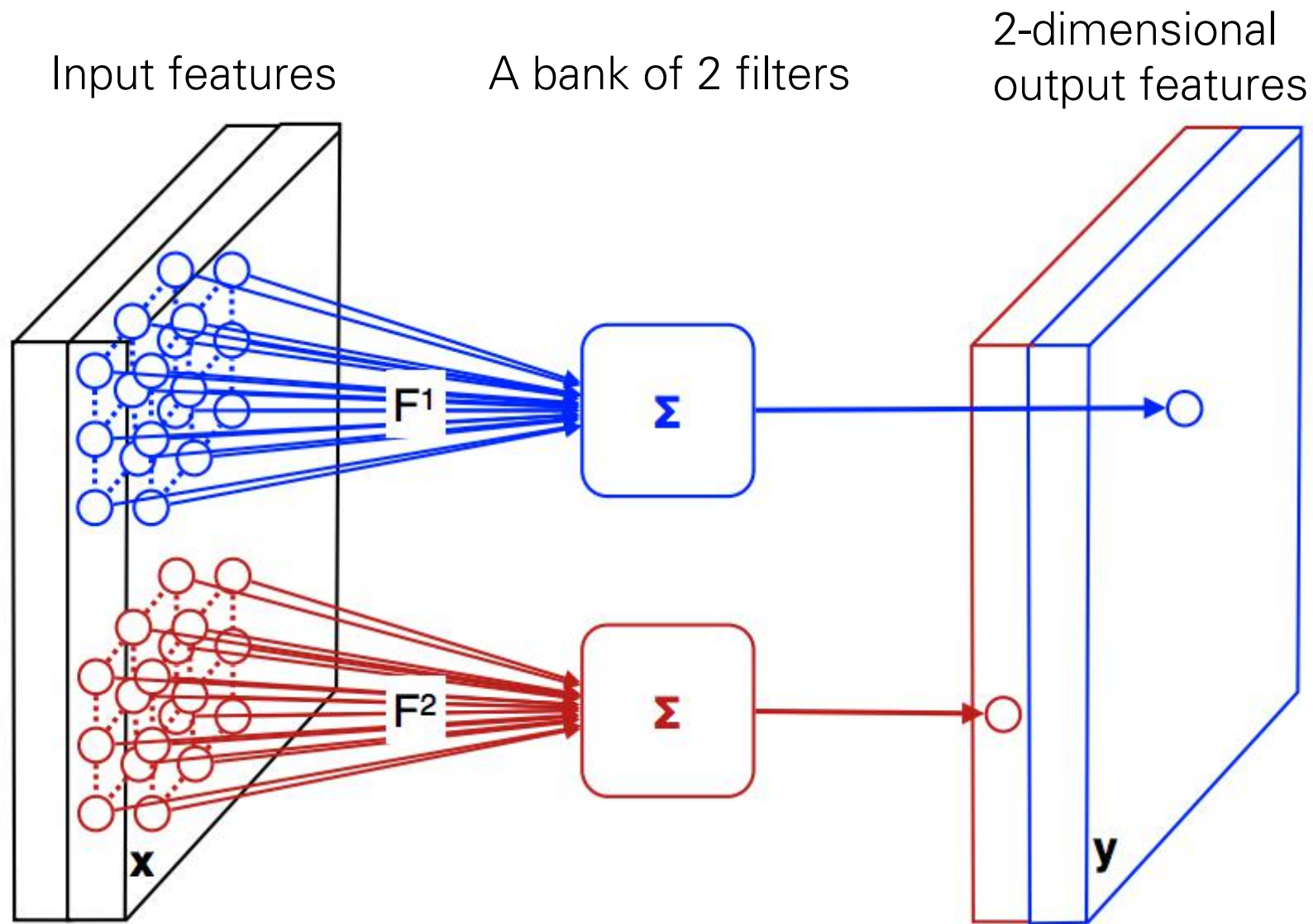$$\mathbb{R}^{N \times C} \to \mathbb{R}^{N \times 1}$$

# Multiple channels

Conv layer



$$\mathbb{R}^{N \times C^{(0)}} \to \mathbb{R}^{N \times C^{(1)}}$$

# Multiple channels

Conv layer



$$\mathbb{R}^{N \times C^{(l)}} \to \mathbb{R}^{N \times C^{(l+1)}}$$

Input features          A bank of 2 filters

2-dimensional
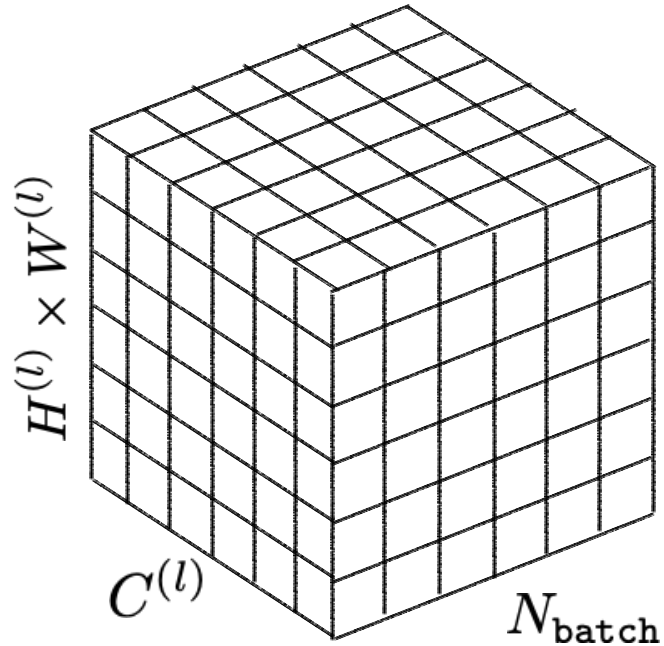output features



**F¹**

**Σ**

**F²**

**Σ**

**x**

**y**

$$\mathbb{R}^{H \times W \times C^{(l)}} \to \mathbb{R}^{H \times W \times C^{(l+1)}}$$
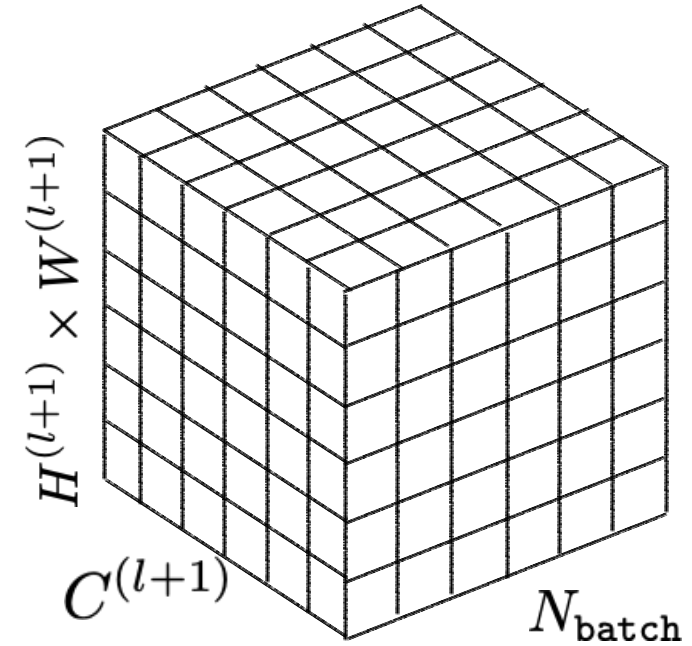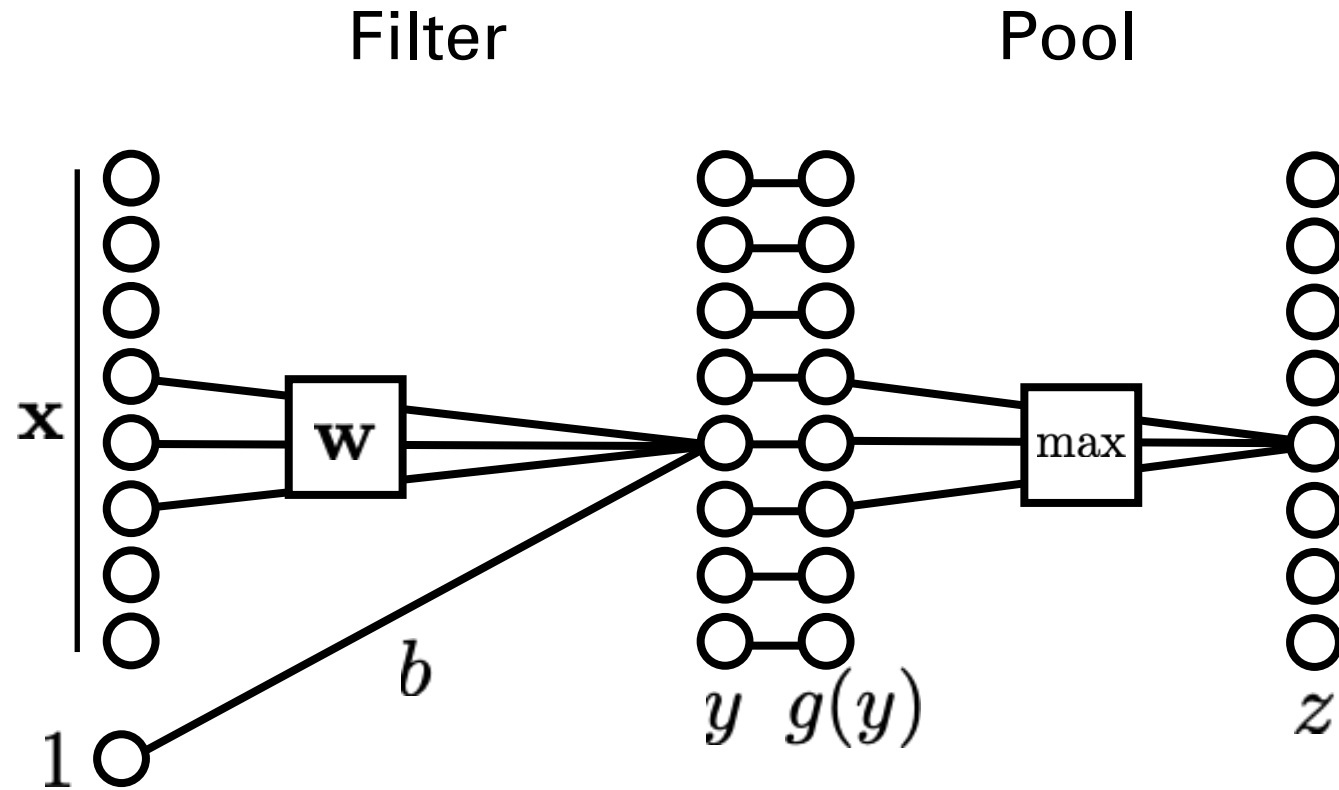
[Figure from Andrea Vedaldi]

117

# "Tensor flow"

$$\mathbf{x}^{(l)} \in \mathbb{R}^{N_{\mathrm{batch}} \times H^{(l)} \times W^{(l)} \times C^{(l)}}$$

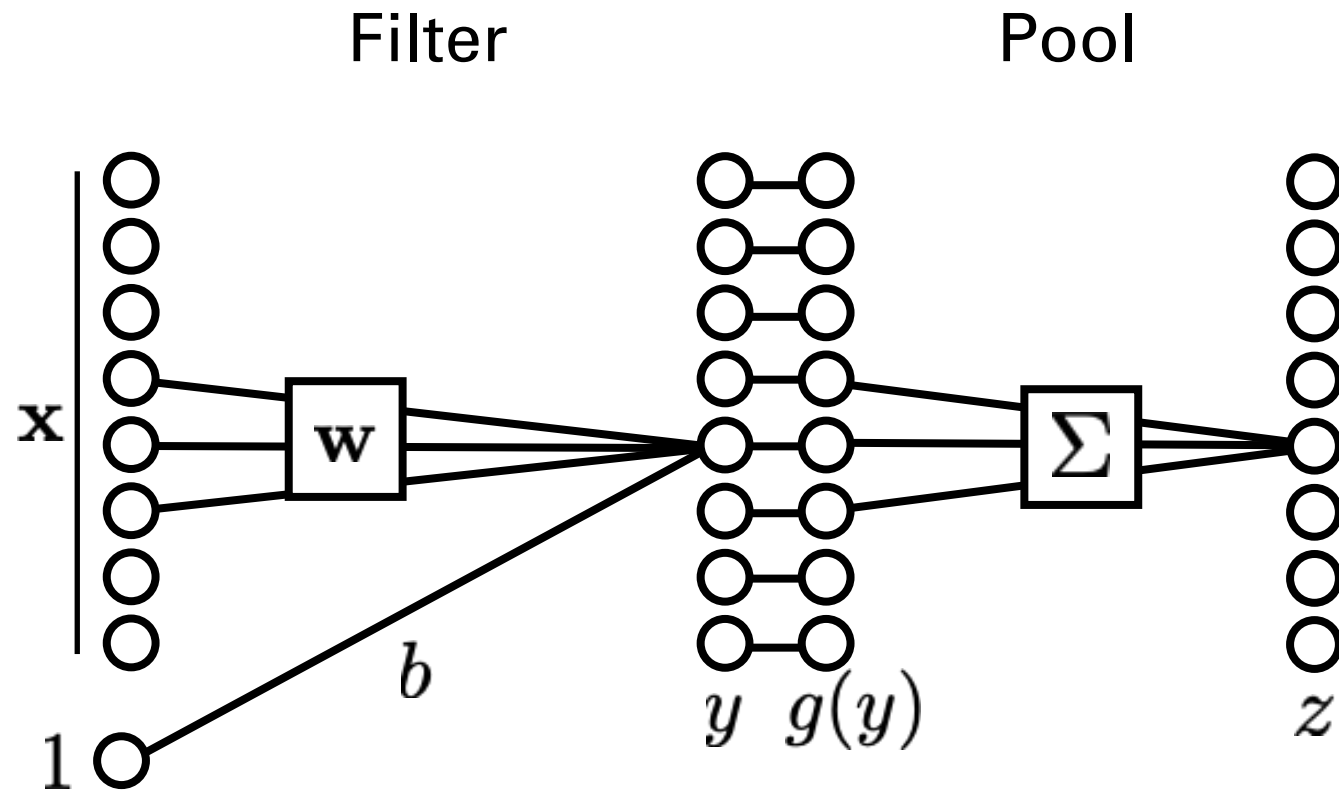$$\mathbf{x}^{(l+1)} \in \mathbb{R}^{N_{\mathrm{batch}} \times H^{(l+1)} \times W^{(l+1)} \times C^{(l+1)}}$$

# Pooling



Filter       Pool

Max pooling

$$z_k = \max_{j \in \mathcal{N}(j)} g(y_j)$$
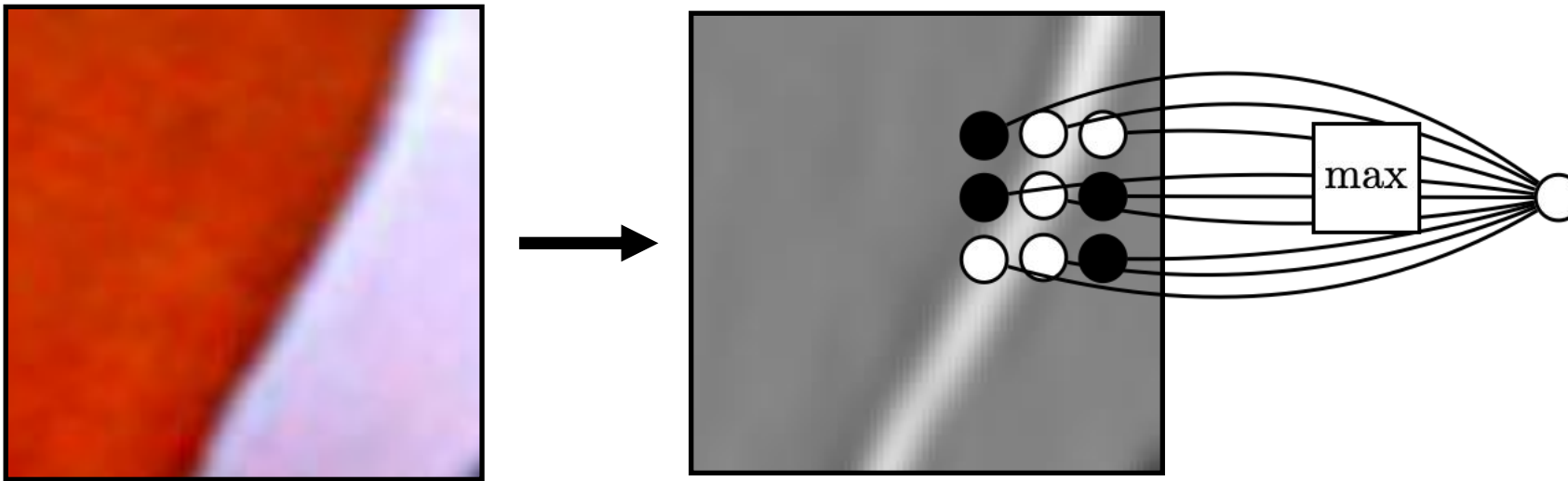
# Pooling



Filter

Pool

Max pooling

$$z_k = \max_{j \in \mathcal{N}(j)} g(y_j)$$

Mean pooling

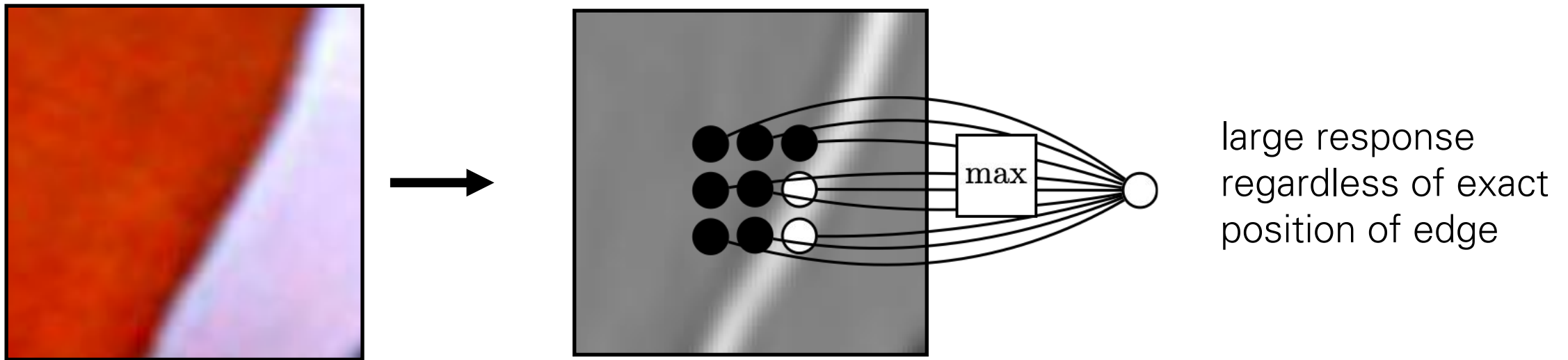$$z_k = \frac{1}{|\mathcal{N}|} \sum_{j \in \mathcal{N}(j)} g(y_j)$$
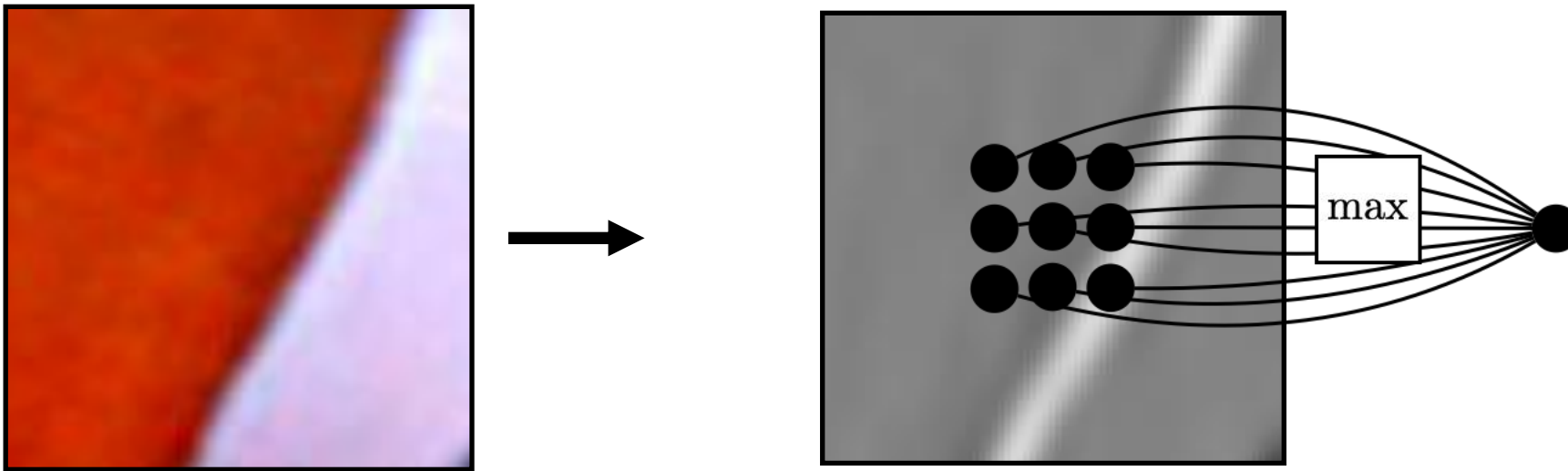
# Pooling – Why?

Pooling across spatial locations achieves
stability w.r.t. small translations:

# Pooling – Why?

Pooling across spatial locations achieves
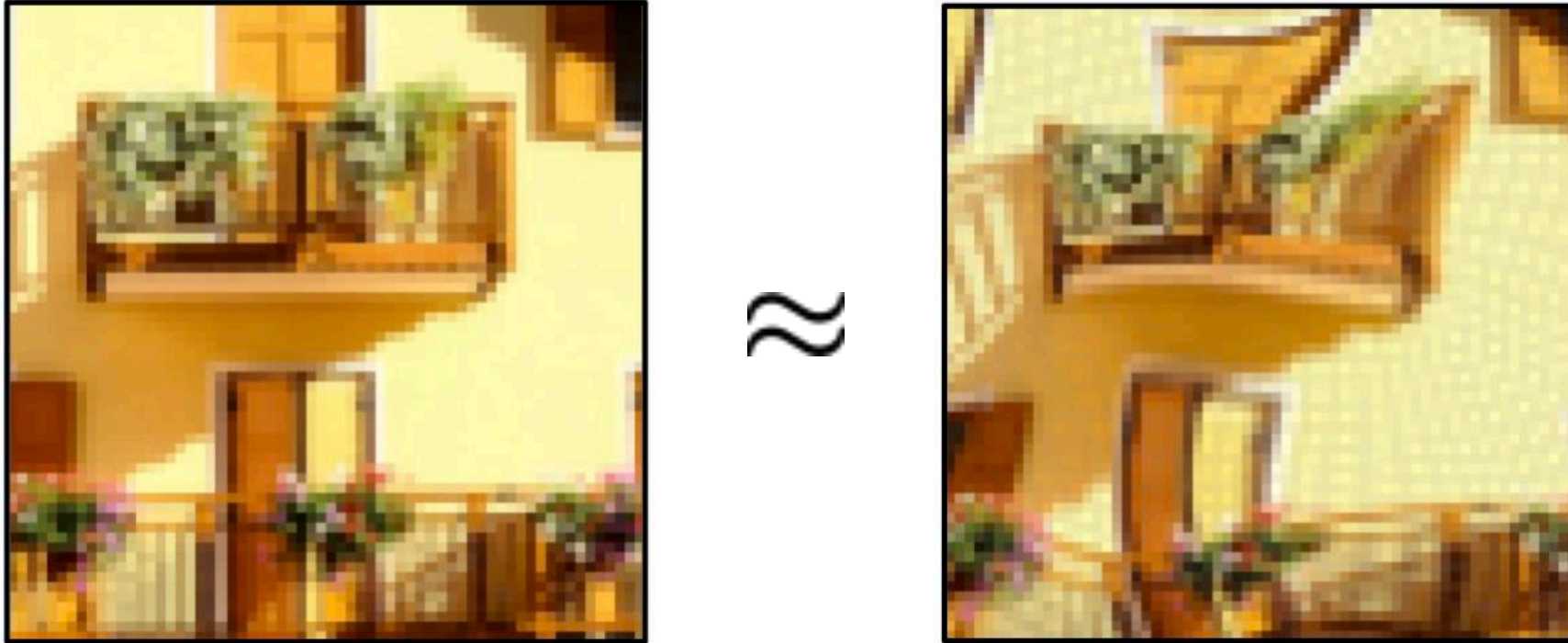stability w.r.t. small translations:



large response
regardless of exact
position of edge

# Pooling – Why?

Pooling across spatial locations achieves
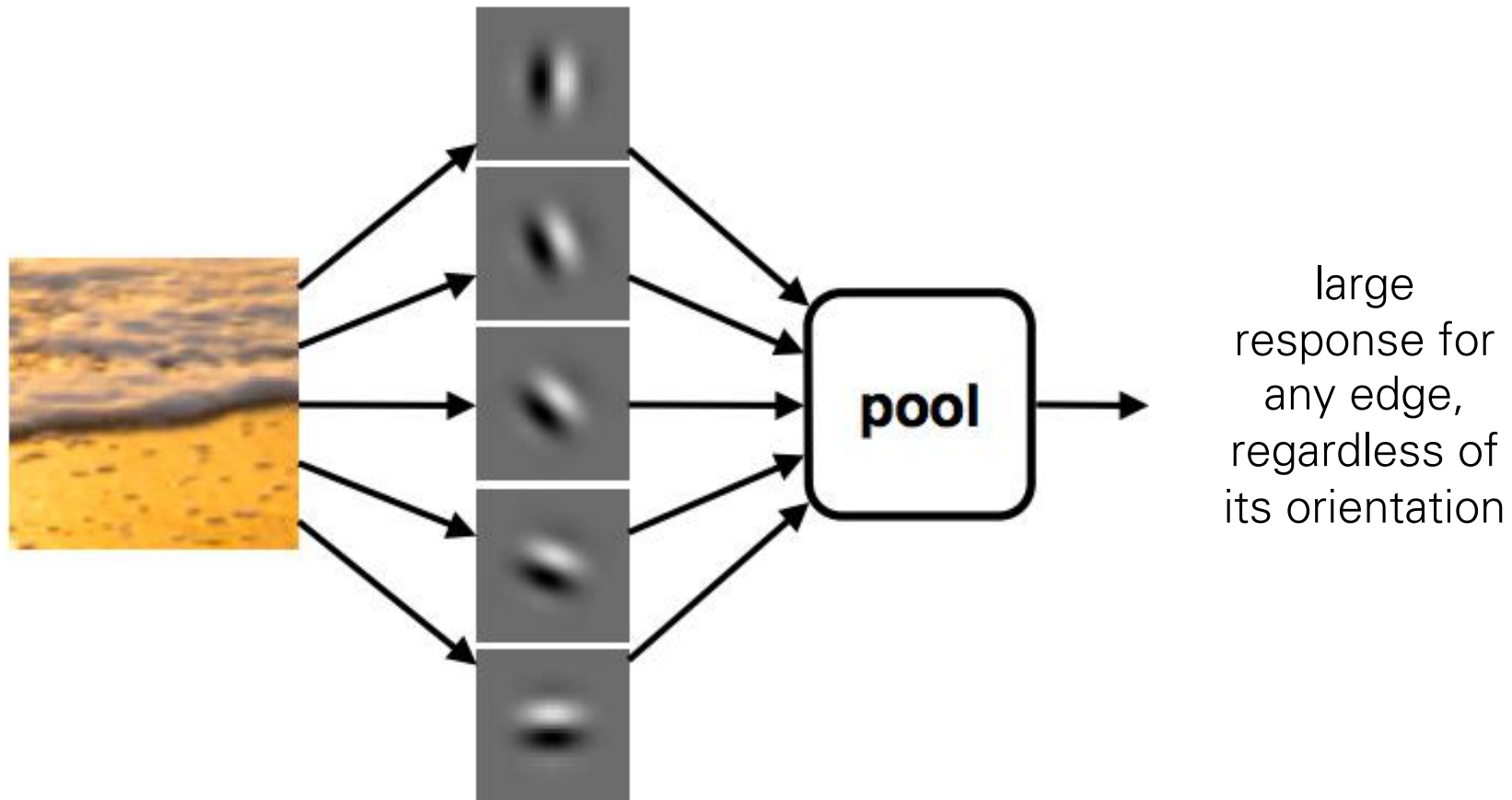stability w.r.t. small translations:

# CNNs are stable w.r.t. diffeomorphisms



$\approx$

["Unreasonable effectiveness of Deep Features as a Perceptual Metric", Zhang et al. 201
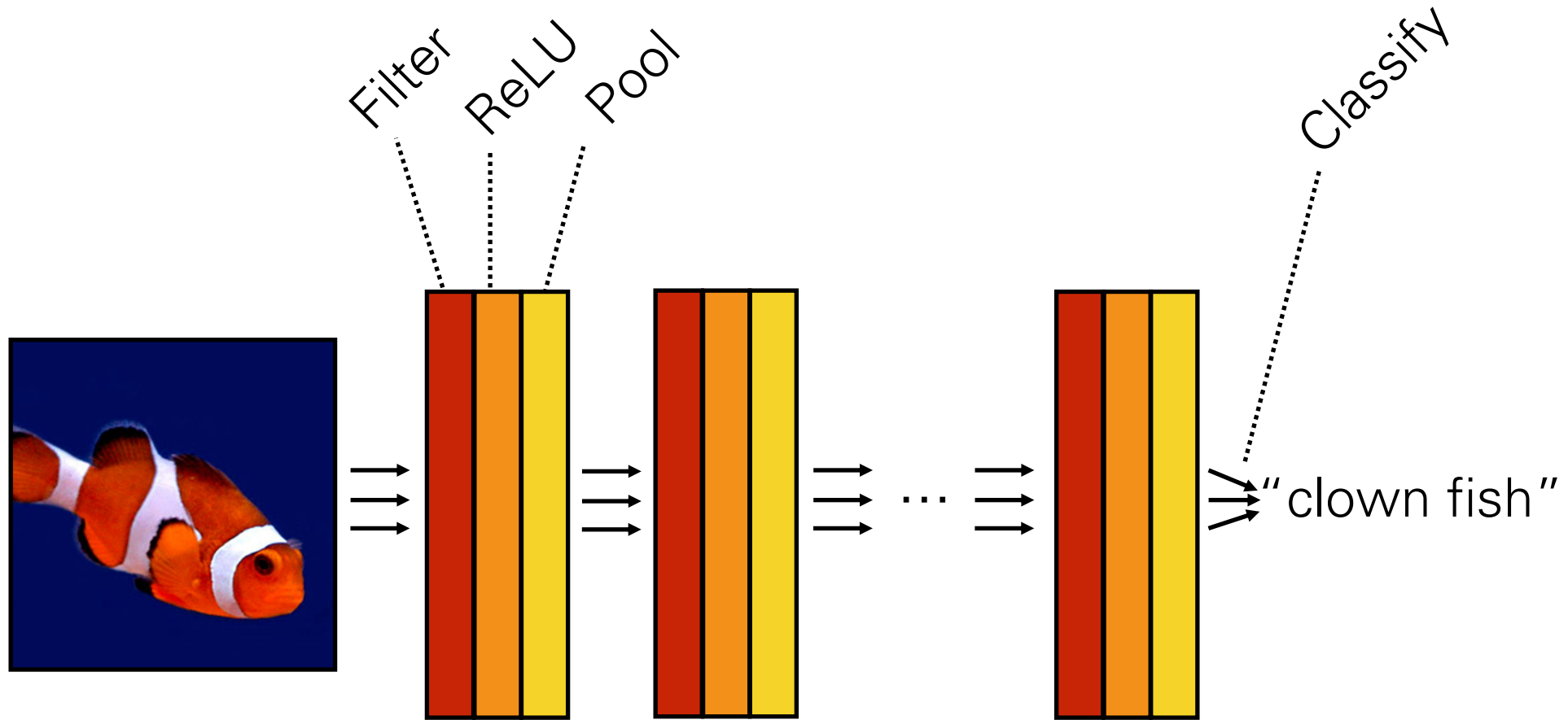
# Pooling – Why?

Pooling across feature channels (filter outputs)
can achieve other kinds of invariances:



large
response for
any edge,
regardless of
its orientation

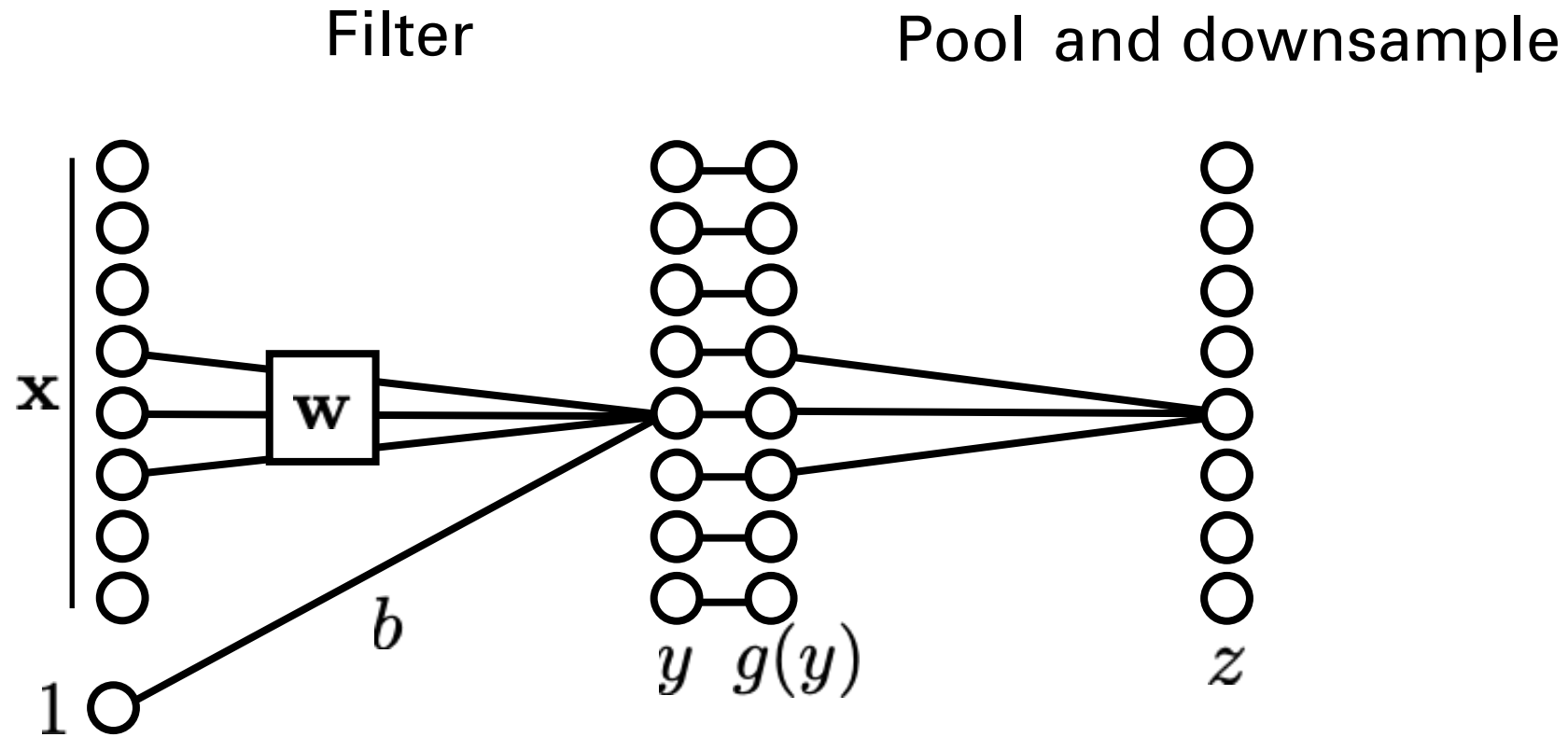[Derived from slide by Andrea Vedaldi]

# Computation in a neural net



$$f(\mathbf{x}) = f_L(\ldots f_2(f_1(\mathbf{x})))$$

# Downsampling



Filter          Pool and downsample

$\mathbf{x}$    $\mathbf{w}$    $b$    $1$    $y$    $g(y)$    $z$

# Downsampling

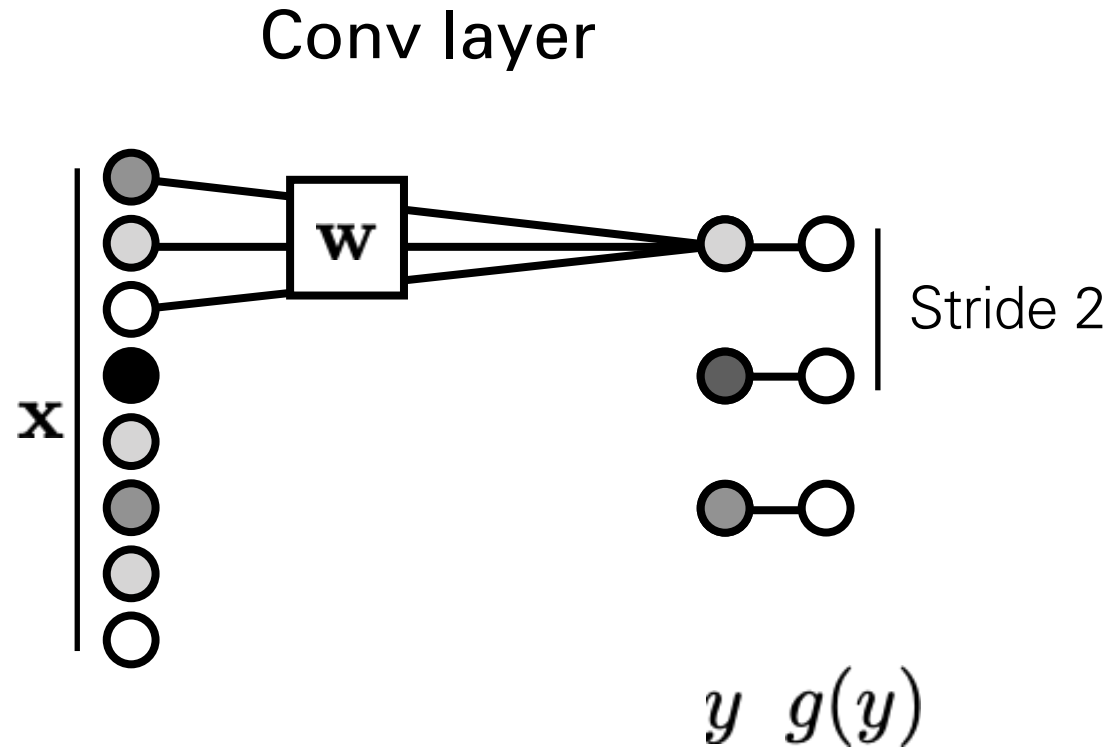Filter                              Downsample



$$\mathbb{R}^{H^{(l)} \times W^{(l)} \times C^{(l)}} \rightarrow \mathbb{R}^{H^{(l+1)} \times W^{(l+1)} \times C^{(l+1)}}$$
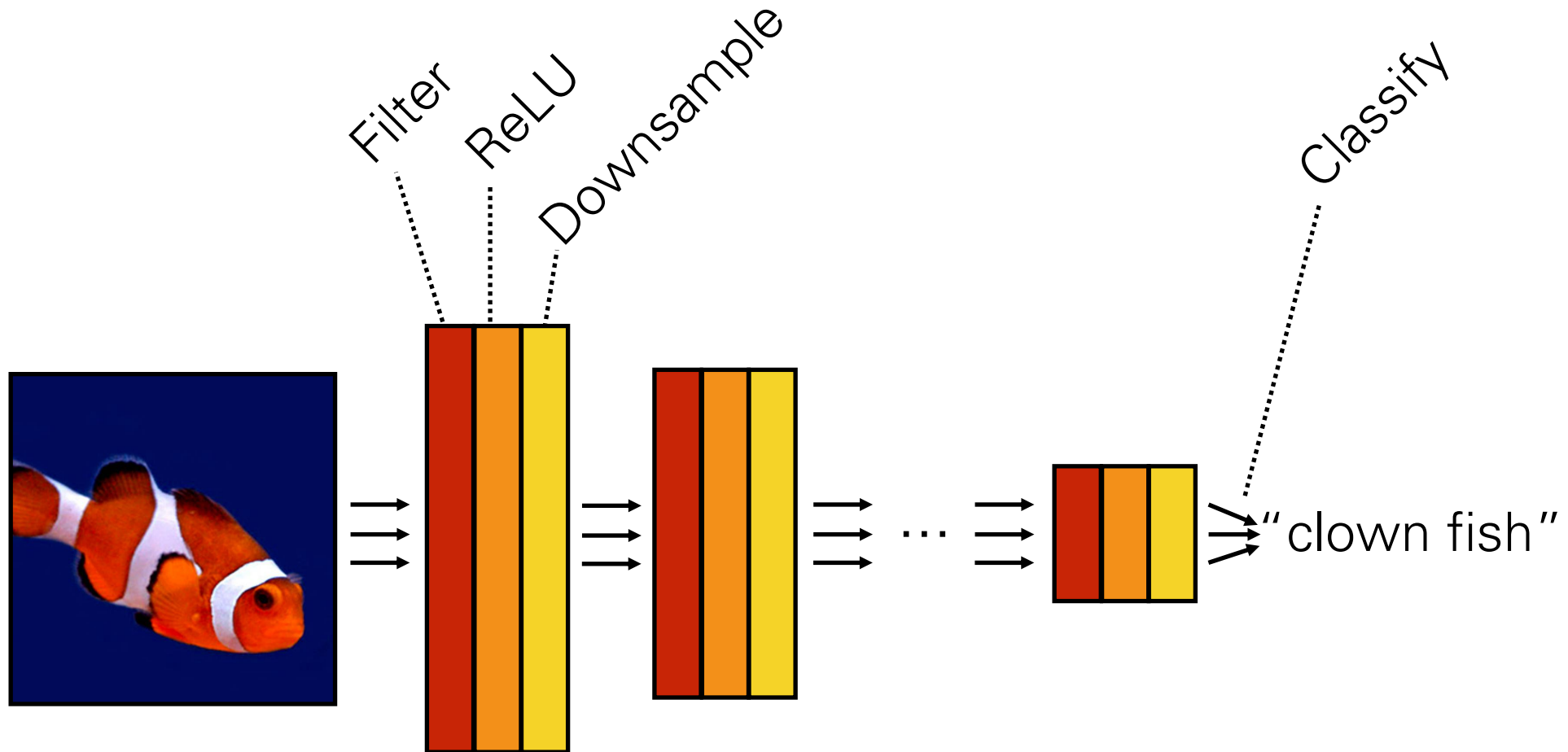
# Strided convolution



Conv layer

Stride 2

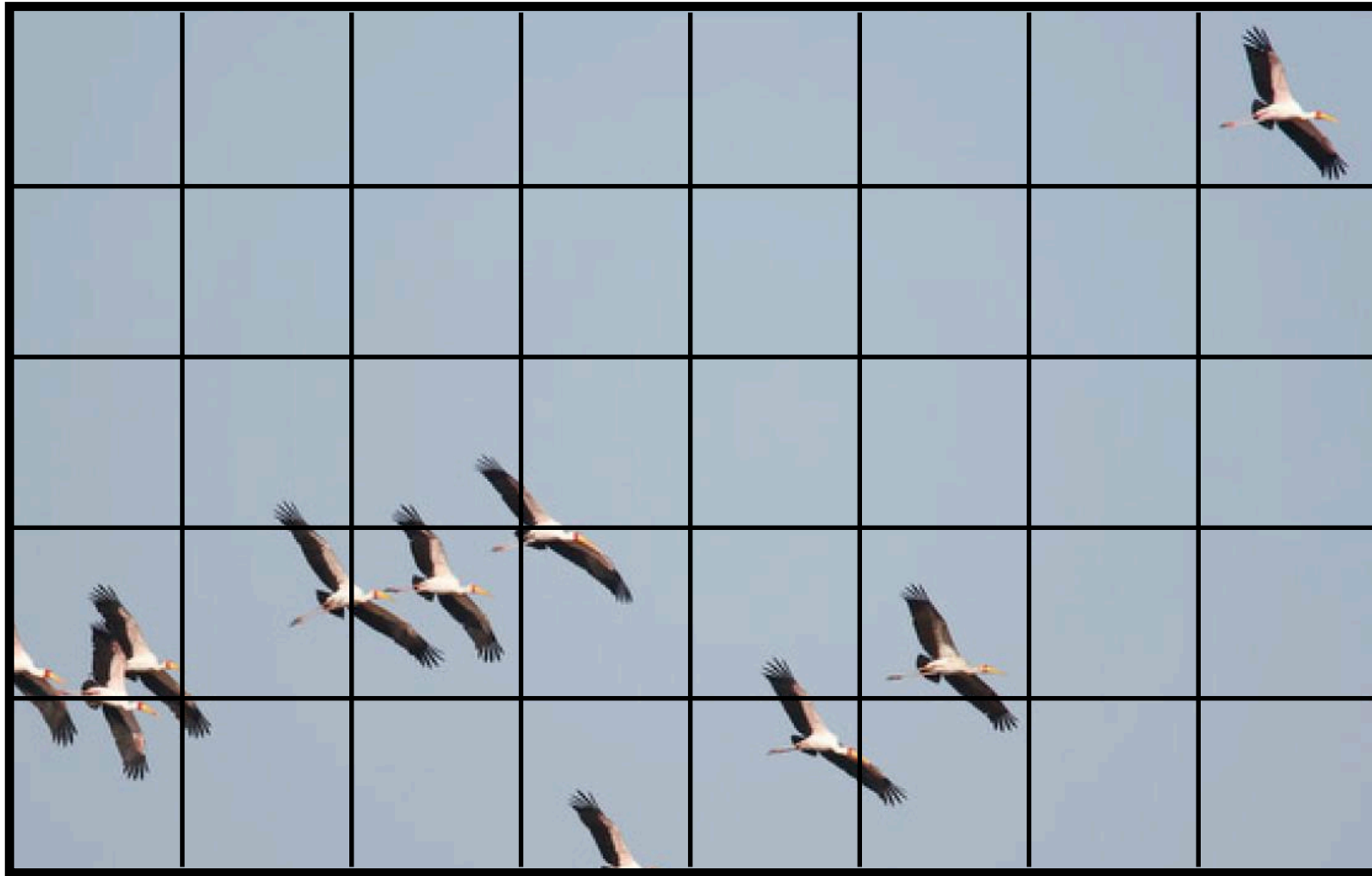**Strided convolutions** combine convolution and downsampling into a single operation.

$$y \quad g(y)$$

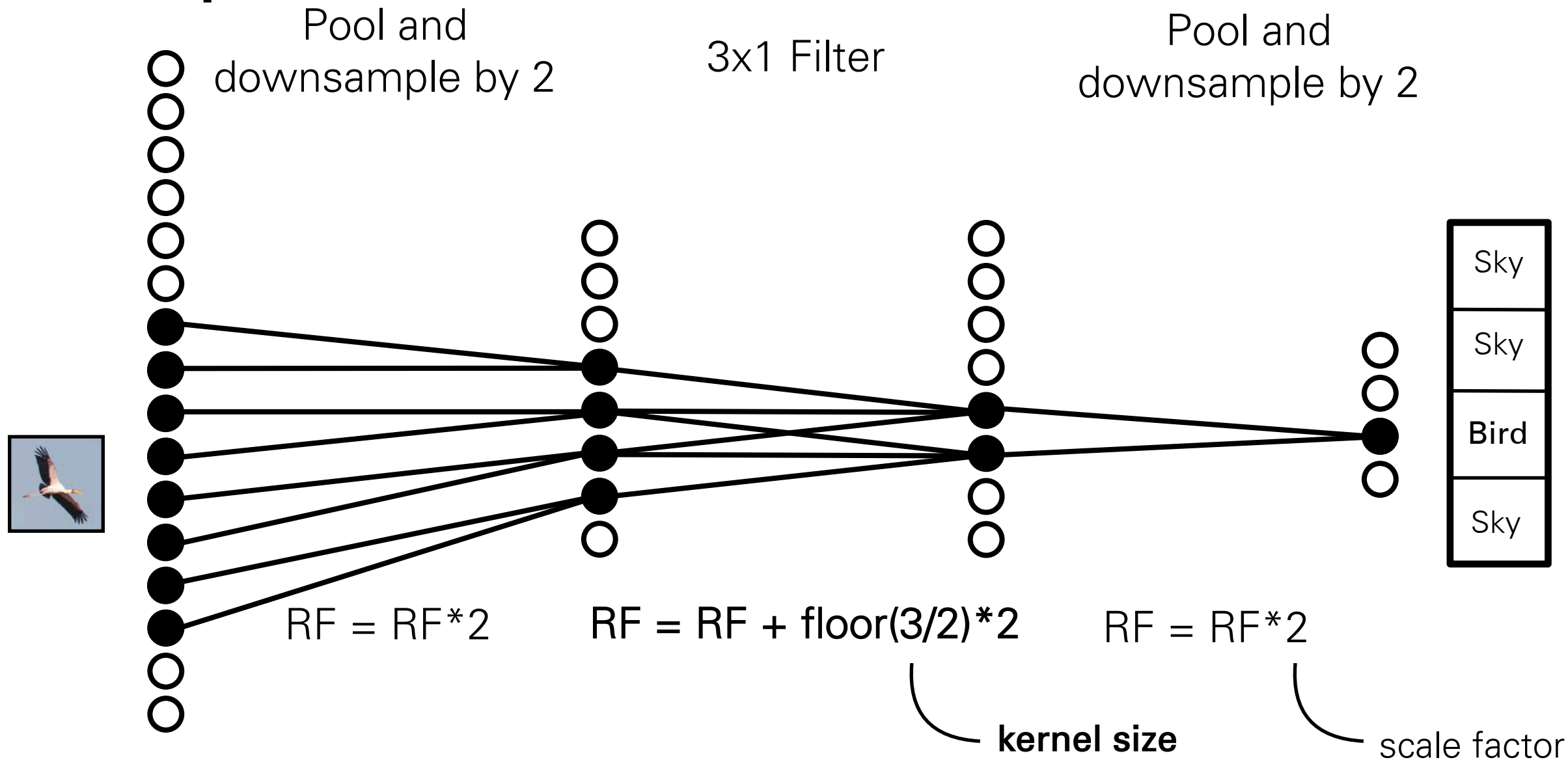# Computation in a neural net



Filter  ReLU  Downsample  Classify

"clown fish"

$$f(\mathbf{x}) = f_L(\ldots f_2(f_1(\mathbf{x})))$$

# Receptive fields

# Receptive fields

Pool and
downsample by 2

3x1 Filter

Pool and
downsample by 2

Sky

Sky

**Bird**

Sky

RF = RF*2

RF = RF + floor(3/2)*2

RF = RF*2

**kernel size**

scale factor
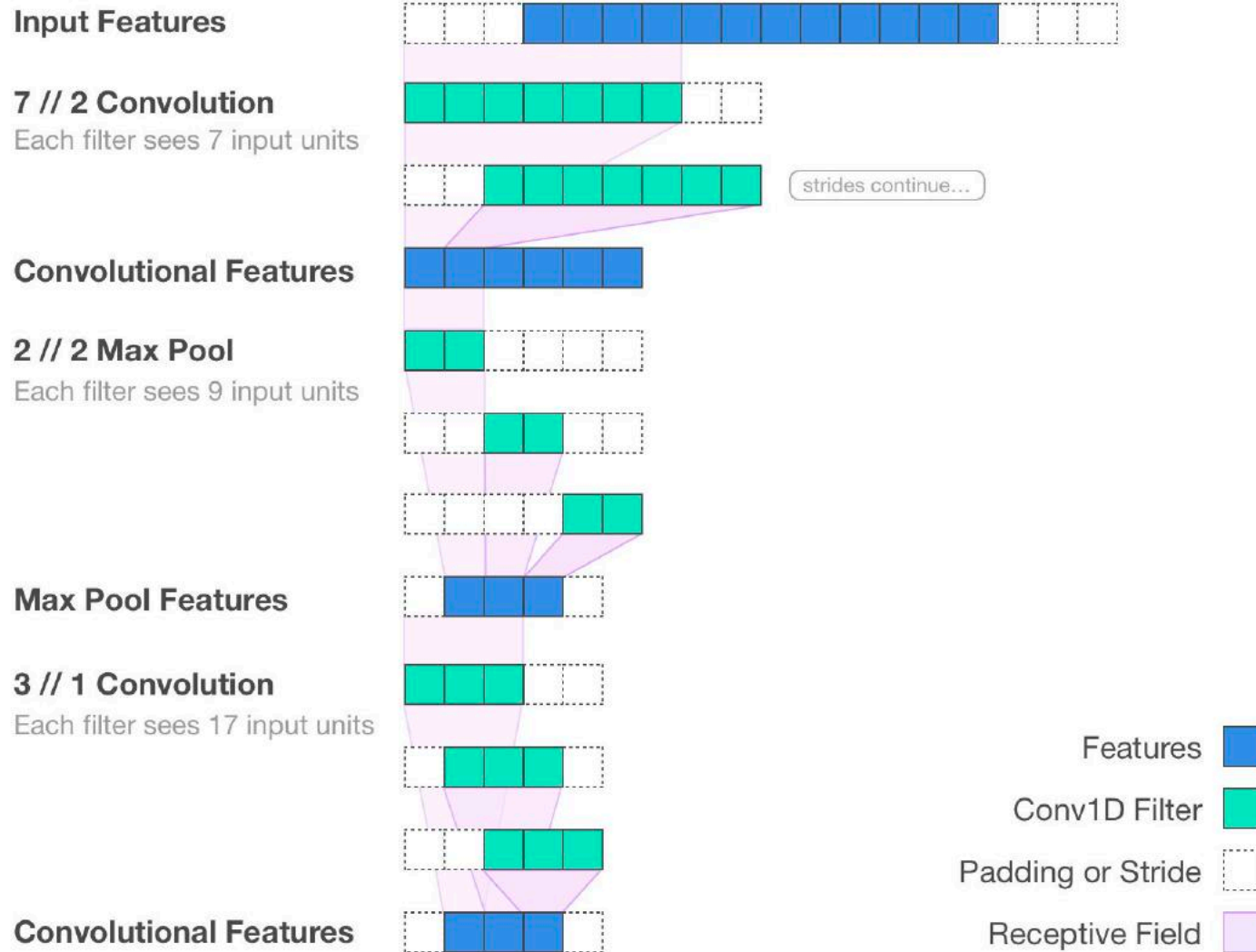
# Effective Receptive Field

Contributing input units to a convolutional filter.                    @jimmfleming // fomoro.com

**Input Features**

**7 // 2 Convolution**
Each filter sees 7 input units

strides continue…

**Convolutional Features**

**2 // 2 Max Pool**
Each filter sees 9 input units

**Max Pool Features**

**3 // 1 Convolution**
Each filter sees 17 input units

**Convolutional Features**

Features
Conv1D Filter
Padding or Stride
Receptive Field

[http://fomoro.com/tools/receptive-fields/index.html]

133

# Why CNNs?



**Fig. 1.** (a) Scatterplots of pairs of pixels at three different spatial displacements, averaged over five examples images. (b) Autocorrelation function. Photographs are of New York City street scenes, taken with a Canon 10D digital camera, and processed in RAW linear sensor mode (producing pixel intensities are in roughly proportional to light intensity). Correlations were computed on the logs of these sensor intensity values [41].

[http://6.869.csail.mit.edu/fa18/notes/simoncelli2005.pdf]

# Why CNNs?

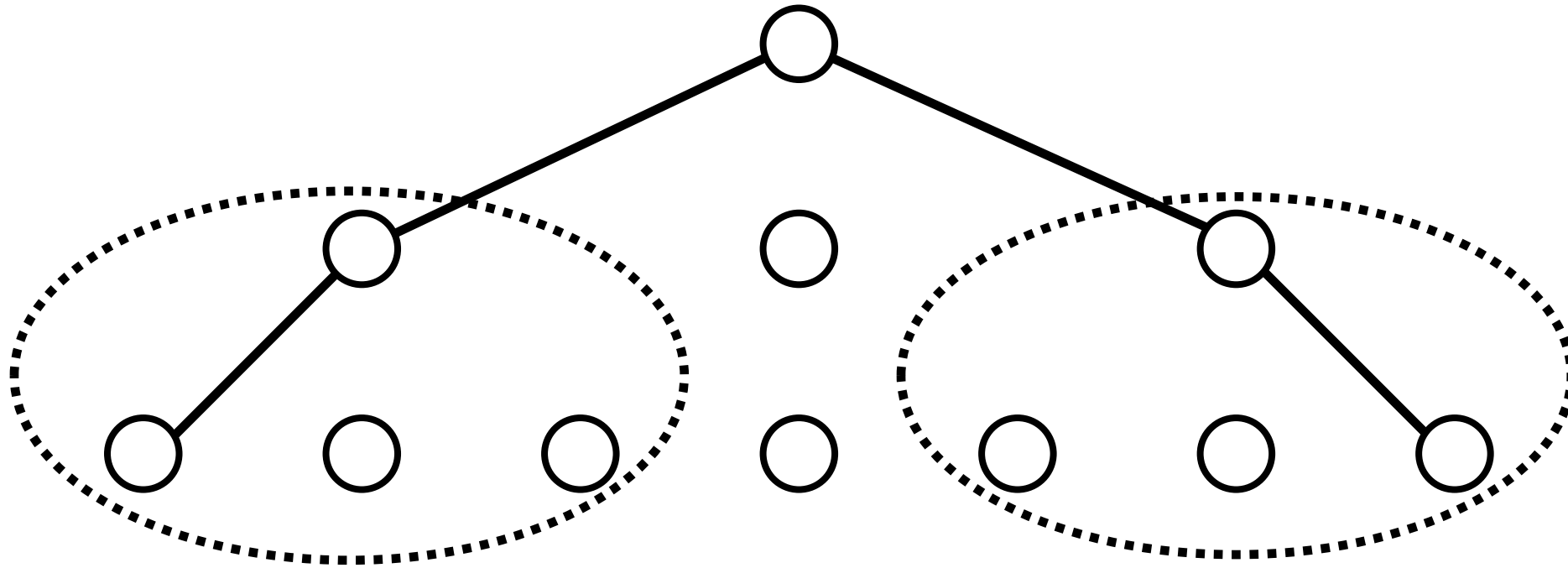Statistical dependences between pixels decay as a power law of distance between the pixels.

It is therefore often sufficient to model local dependences only. —> **Convolution**

More generally, we should allocate parameters that model dependences in proportion to the strength of those dependences. —> **Multiscale, hierarchical representations**

[For more discussion, see "Why does Deep and Cheap Learning Work So Well?", Lin et al. 2017]
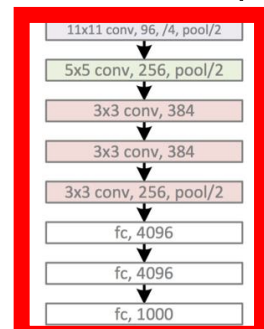
# Why CNNs?

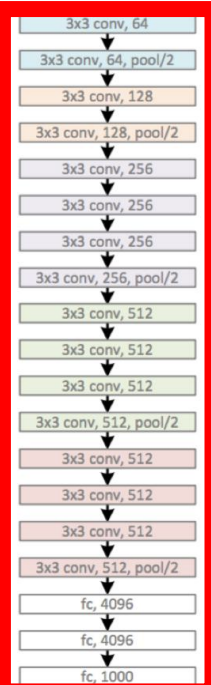Capturing long-range dependences:

# Deep Neural Networks for Visual Recognition

**2012: AlexNet**
5 conv. layers

| |
|---|
| 11x11 conv, 96, /4, pool/2 |
| 5x5 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384 |
| 3x3 conv, 256, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

Error: 15.3%

**2014: VGG**
16 conv. layers

| |
|---|
| 3x3 conv, 64 |
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 128 |
| 3x3 conv, 128, pool/2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

Error: 8.5%

**2015: GoogLeNet**
22 conv. layers

Error: 7.8%

**2016: ResNet**
>100 conv. layers

Error: 4.4%

**2012: AlexNet**
**5 conv. layers**



| 11x11 conv, 96, /4, pool/2 |
| 5x5 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384 |
| 3x3 conv, 256, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

Error: 15.3%

**2014: VGG
16 conv.
layers**

3x3 conv, 64

3x3 conv, 64, pool/2

3x3 conv, 128

3x3 conv, 128, pool/2

3x3 conv, 256

3x3 conv, 256

3x3 conv, 256

3x3 conv, 256, pool/2

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512, pool/2

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512, pool/2

fc, 4096

fc, 4096

fc, 1000

Softmax

Error: 8.5%

# VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE-SCALE IMAGE RECOGNITION

https://arxiv.org/pdf/1409.1556.pdf

Small convolutional kernels: 3x3
ReLu non-linearities
>100 million parameters.



ConvNet
D=64    Pool
D=128
D=256
D=512
D=512
D=4096  D=4096  D=1000

VGG

224x224  112x112  56x56  28x28  14x14  FC  FC  FC + Softmax

# Chaining convolutions



3x3

3x3

5x5

25 coefficients, but only
18 degrees of freedom

9 coefficients, but only
6 degrees of freedom.
Only separable filters… would this be enough?

# Dilated convolutions

3x3

5x5

| a | 0 | b | 0 | c |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| d | 0 | e | 0 | f |
| 0 | 0 | 0 | 0 | 0 |
| g | 0 | h | 0 | i |

○

=

7x7

25 coefficients
9 degrees of freedom

49 coefficients
18 degrees of freedom

(a) Input          (b) Dilation 2          (c) Output

What is lost?

[https://arxiv.org/pdf/1511.07122.pdf]

Figure 1: Systematic dilation supports exponential expansion of the receptive field without loss of resolution or coverage. (a) $F_1$ is produced from $F_0$ by a 1-dilated convolution; each element in $F_1$ has a receptive field of $3 \times 3$. (b) $F_2$ is produced from $F_1$ by a 2-dilated convolution; each element in $F_2$ has a receptive field of $7 \times 7$. (c) $F_3$ is produced from $F_2$ by a 4-dilated convolution; each element in $F_3$ has a receptive field of $15 \times 15$. The number of parameters associated with each layer is identical. The receptive field grows exponentially while the number of parameters grows linearly.

[https://arxiv.org/pdf/1511.07122.pdf]

**2016: ResNet
>100 conv. layers**

# Deep Residual Learning for Image Recognition

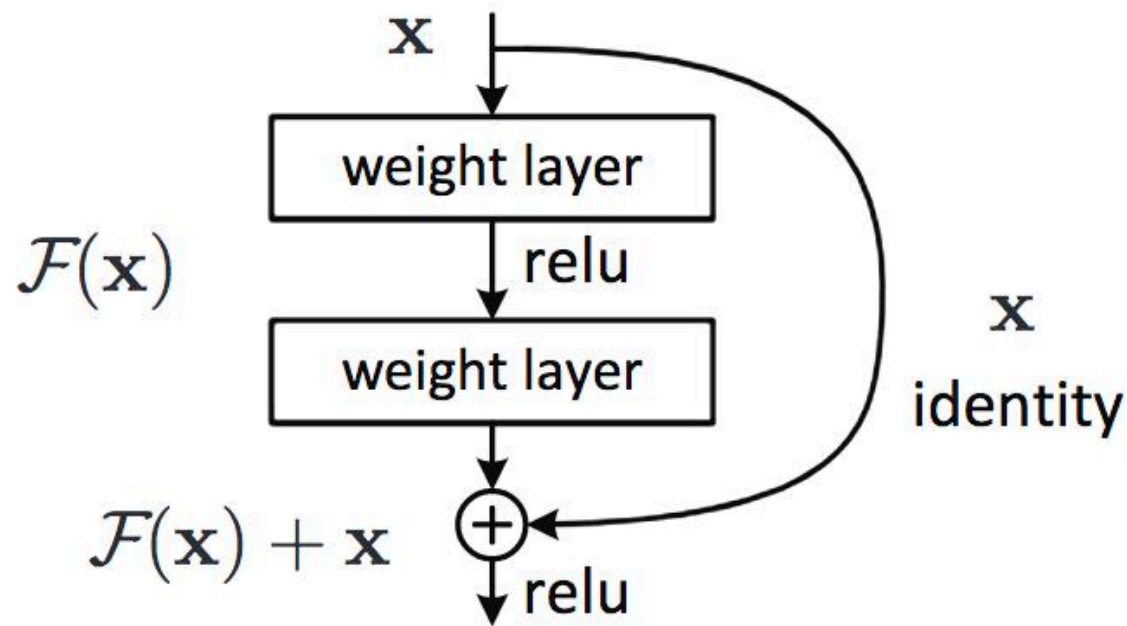https://arxiv.org/pdf/1512.03385.pdf



Figure 2. Residual learning: a building block.
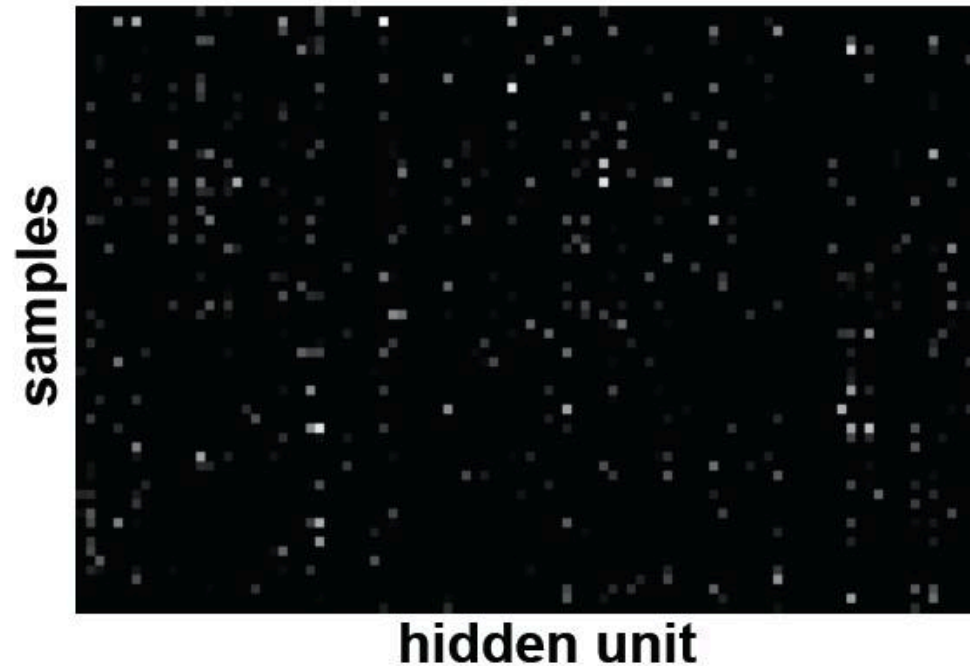
Error: 4.4%

If output has same size as input:



$$\mathcal{F}(\mathbf{x})$$

weight layer

relu

weight layer

$$\mathbf{x}$$ identity

$$\mathcal{F}(\mathbf{x}) + \mathbf{x}$$

relu

If output has a different size:



$$\mathcal{F}(\mathbf{x})$$

weight layer

relu

weight layer

weight layer

$$\mathcal{F}(\mathbf{x}) + W\,x$$

relu

# Other good things to know

- Check gradients numerically by finite differences
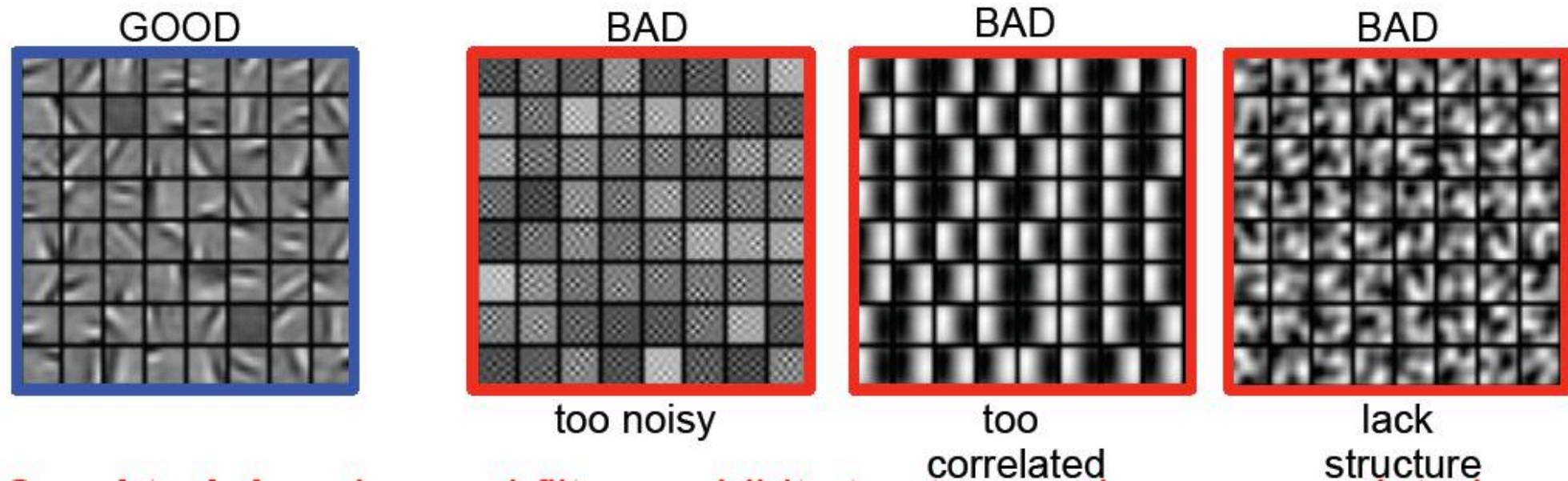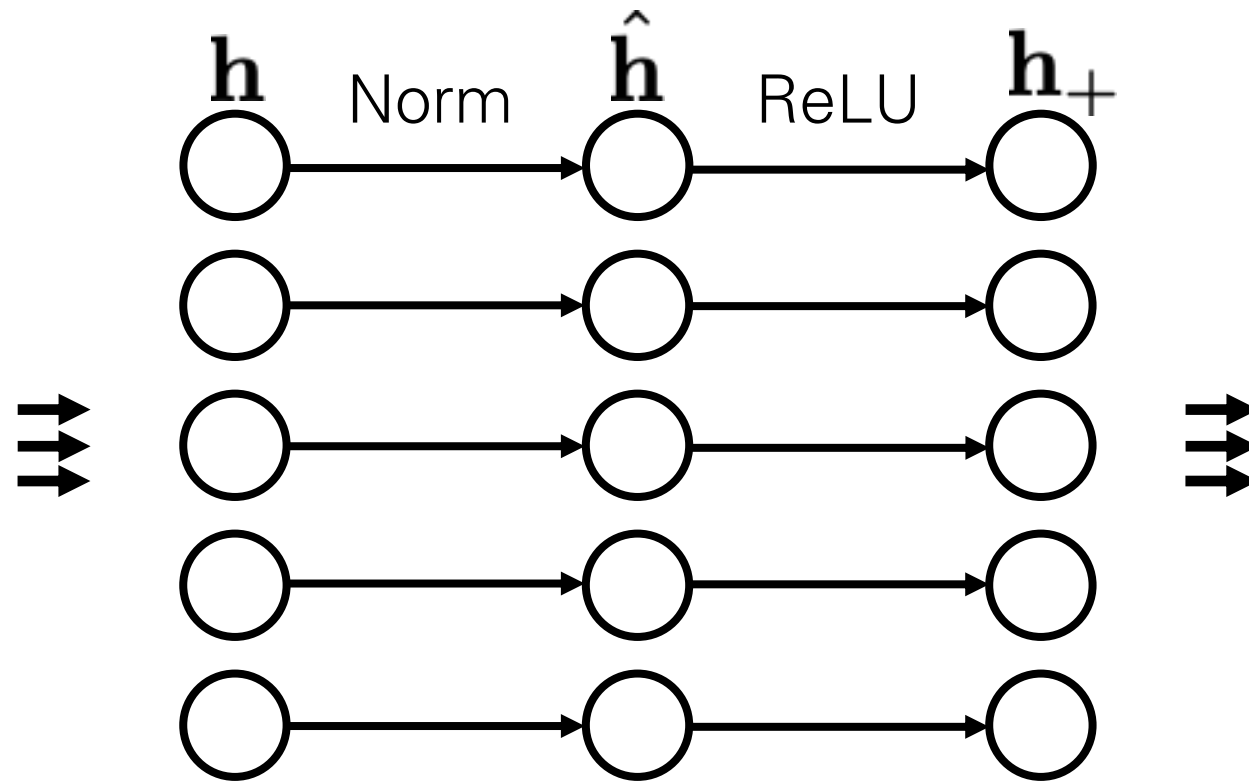- Visualize hidden activations — should be uncorrelated and high variance



**Good training:** hidden units are sparse across samples and across features.

[Derived from slide by Marc'Aurelio Ranzato]

# Other good things to know

- Check gradients numerically by finite differences
- Visualize hidden activations — should be uncorrelated and high variance



**Bad training:** many hidden units ignore the input and/or exhibit strong correlations.

[Derived from slide by Marc'Aurelio Ranzato]

# Other good things to know

- Check gradients numerically by finite differences
- Visualize hidden activations — should be uncorrelated and high variance
- Visualize filters



GOOD

BAD
too noisy

BAD
too correlated

BAD
lack structure

**Good training:** learned filters exhibit structure and are uncorrelated.

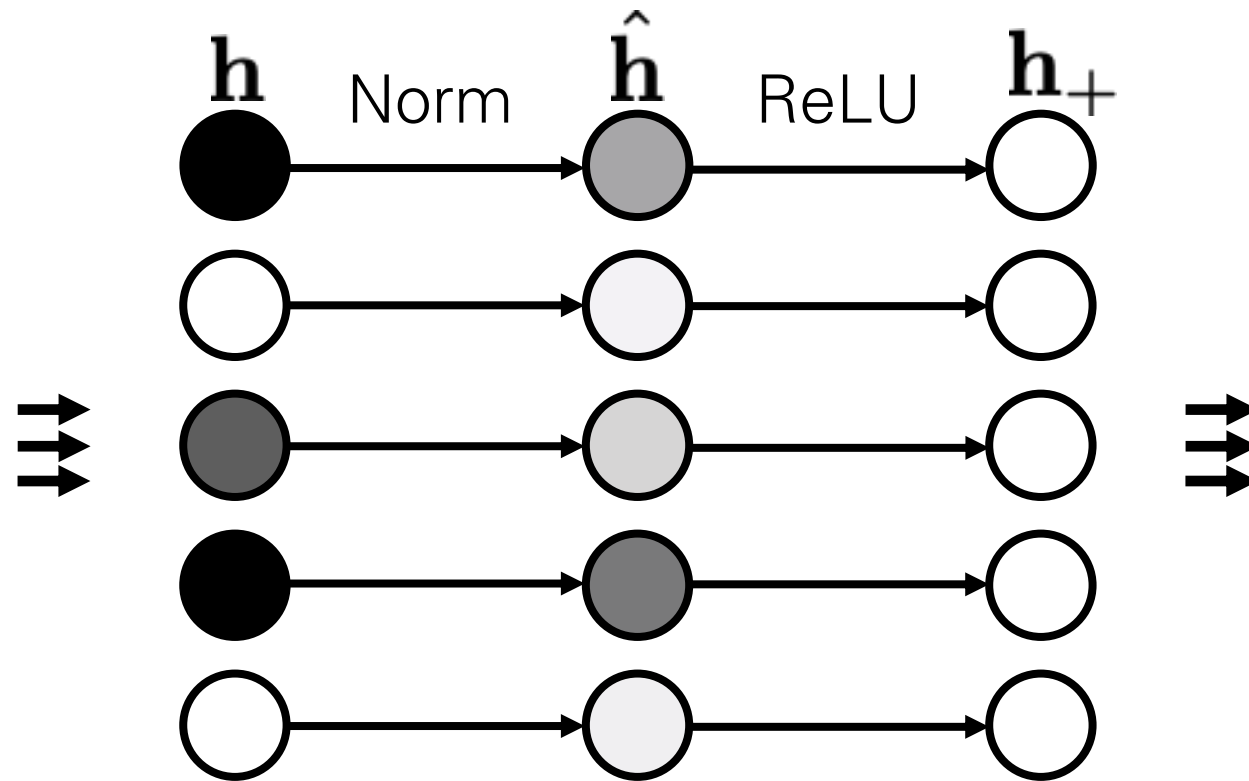[Derived from slide by Marc'Aurelio Ranzato]

# Normalization layers



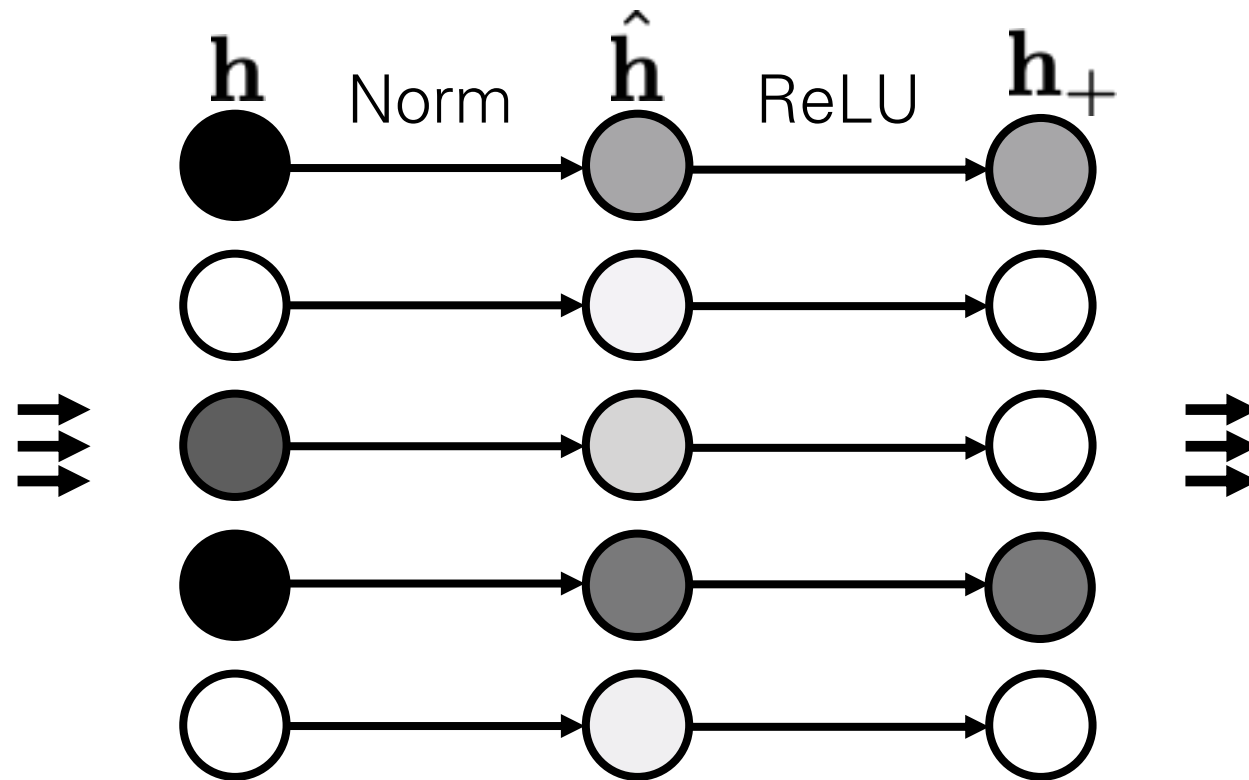$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\mathrm{Var}[h_k]}}$$

# Normalization layers



$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\mathrm{Var}[h_k]}}$$

# Normalization layers



$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\mathbf{Var}[h_k]}}$$

# Normalization layers



$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\mathbf{Var}[h_k]}}$$

# Normalization layers
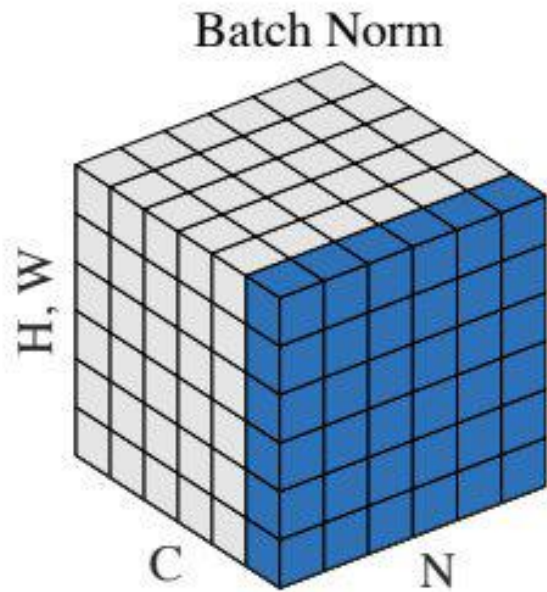
Keep track of mean and variance of a unit (or a population of units) over time.

Standardize unit activations by subtracting mean and dividing by variance.

Squashes units into a **standard range**, avoiding overflow.

Also achieves **invariance** to mean and variance of the training signal.

Both these properties reduce the effective capacity of the model, i.e. regularize the model.

# Normalization layers



Normalize w.r.t. a single hidden unit's pattern of activation over training examples (a batch of examples).

[Figure from Wu & He, arXiv 2018]

# Normalization layers



Normalize w.r.t. the mean and variance of the activations of all the hidden units (neurons) on this layer (c).
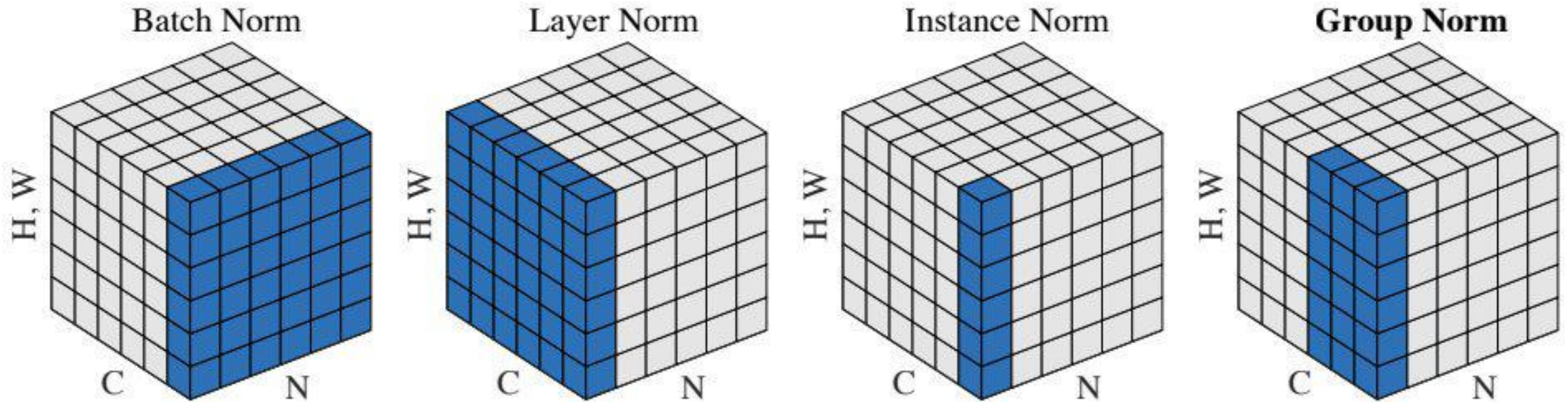
[Figure from Wu & He, arXiv 2018]

# Normalization layers



Normalize w.r.t. the mean and variance of the activations of all the hidden units (neurons) on this layer (c) that process this particular location (h,w) in the image.

[Figure from Wu & He, arXiv 2018]
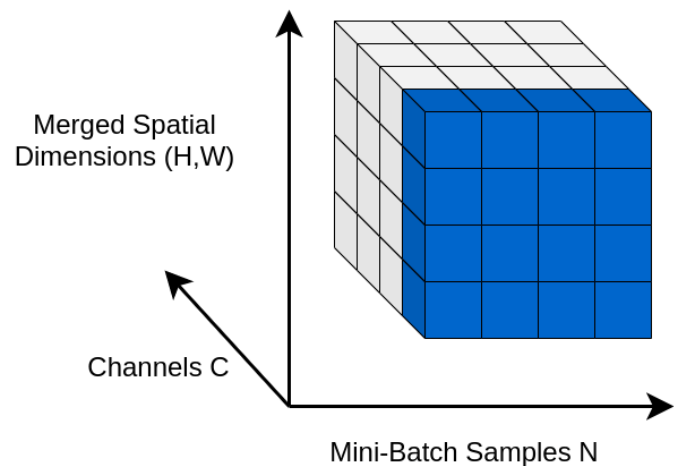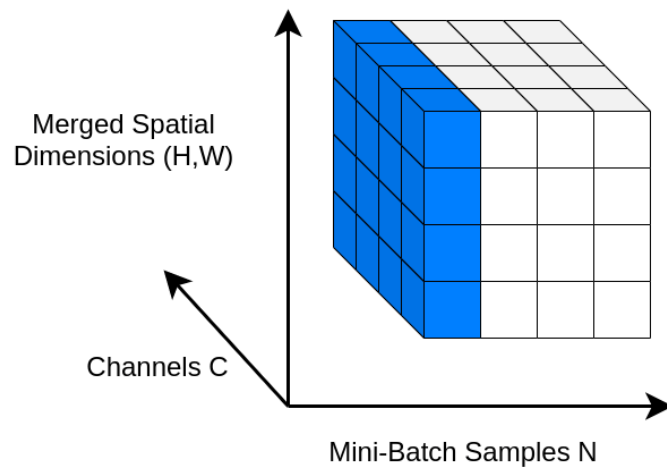
# Normalization layers



Batch Norm     Layer Norm     Instance Norm     Group Norm

Might as well…

[Figure from Wu & He, arXiv 2018]
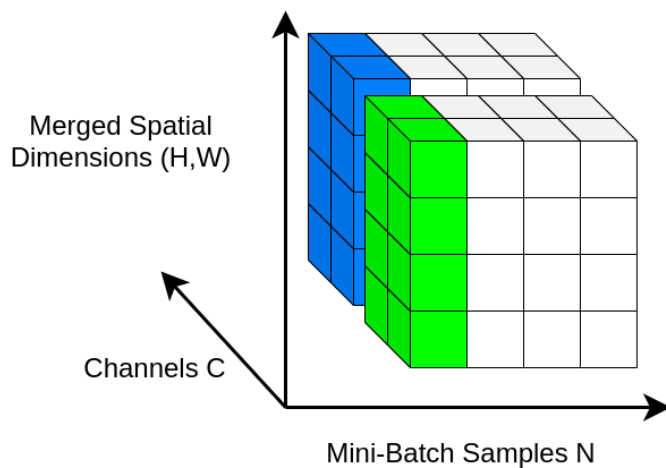
# Normalization layers



Batch Normalization (2015)

Layer Normalization (2016)
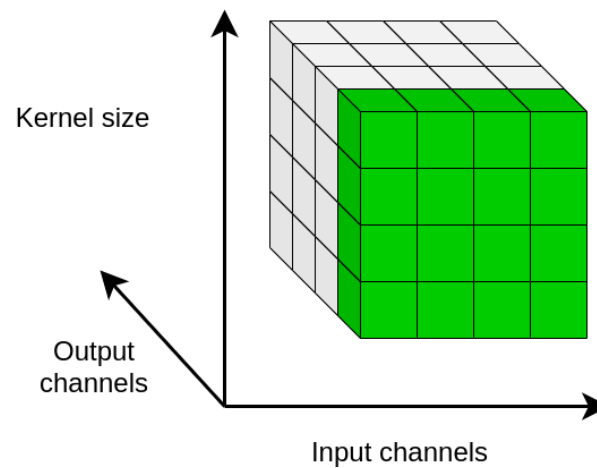
Instance Normalization (2016)

Group Normalization (2018)

Weight Standardization (2019)
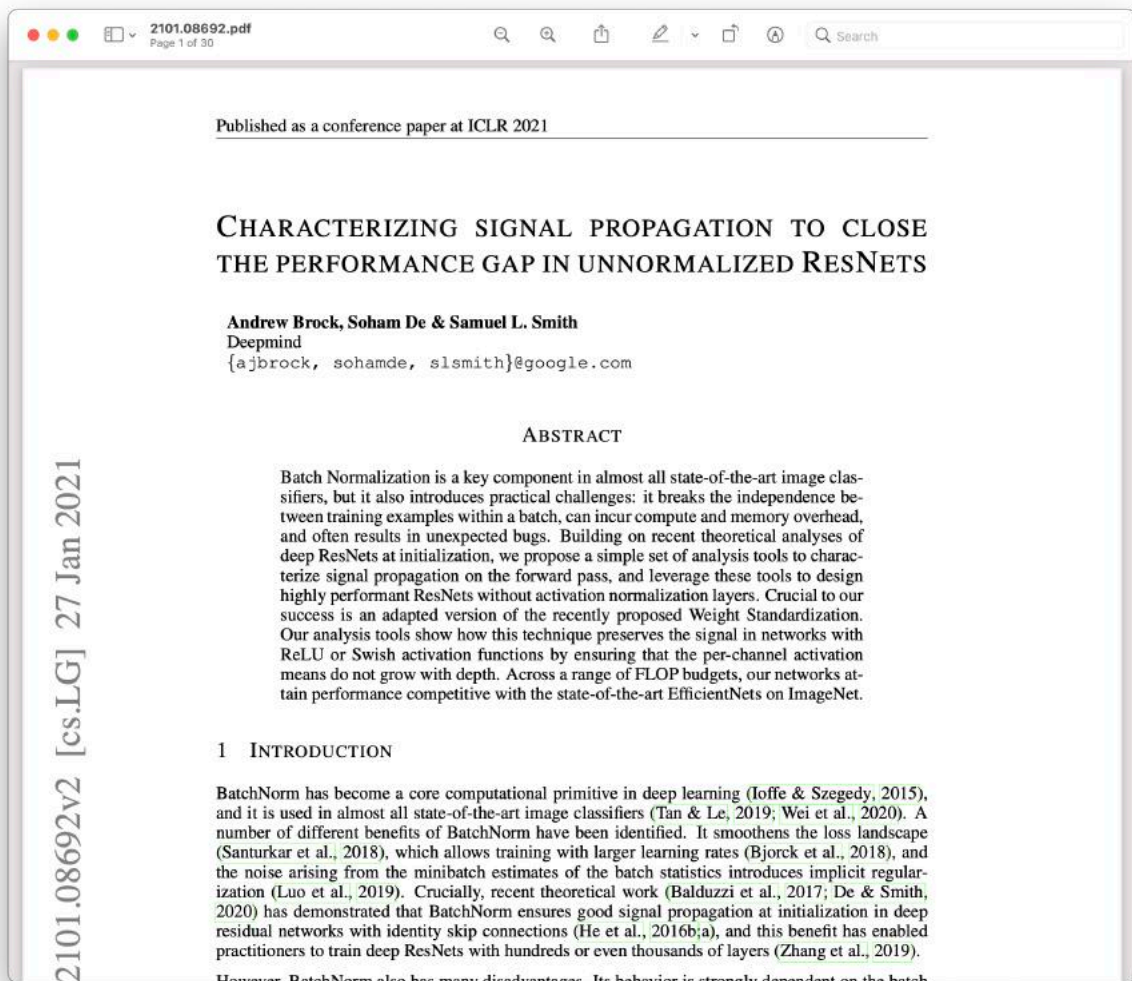
[https://theaisummer.com/normalization]

# No normalization layers

# Next Lecture:
# Sequential Processing with RNNs