# Today's Lecture

- unreasonable effectiveness of data

- deep learning

- computation in a neural net

- optimization

- backpropagation

- convolutional neural networks

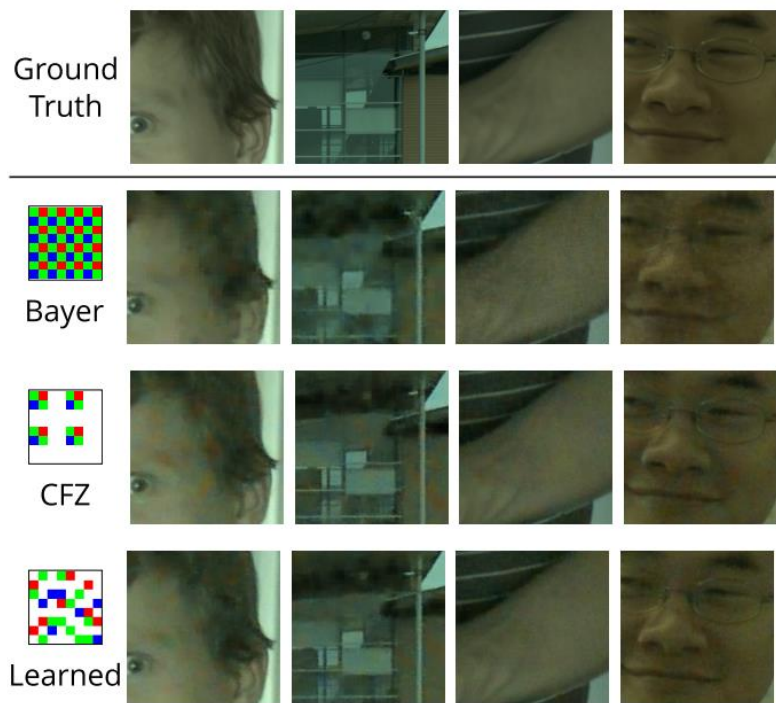- applications in computational photography

**Disclaimer:** 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

# Neural Networks in Computational Photography

- Now: learned pipelines for computational imaging
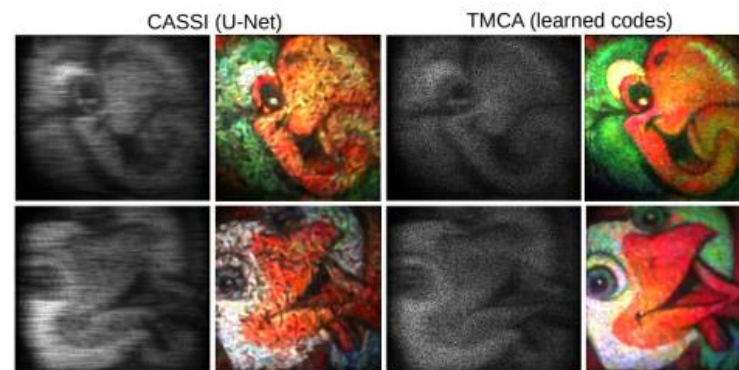


Learning CFAs



(b) Raw data via traditional pipeline     (c) Our result

Learning ISPs



Learning coded apertures

# Neural Networks in Computational Photography

- Now: learned pipelines for computational imaging
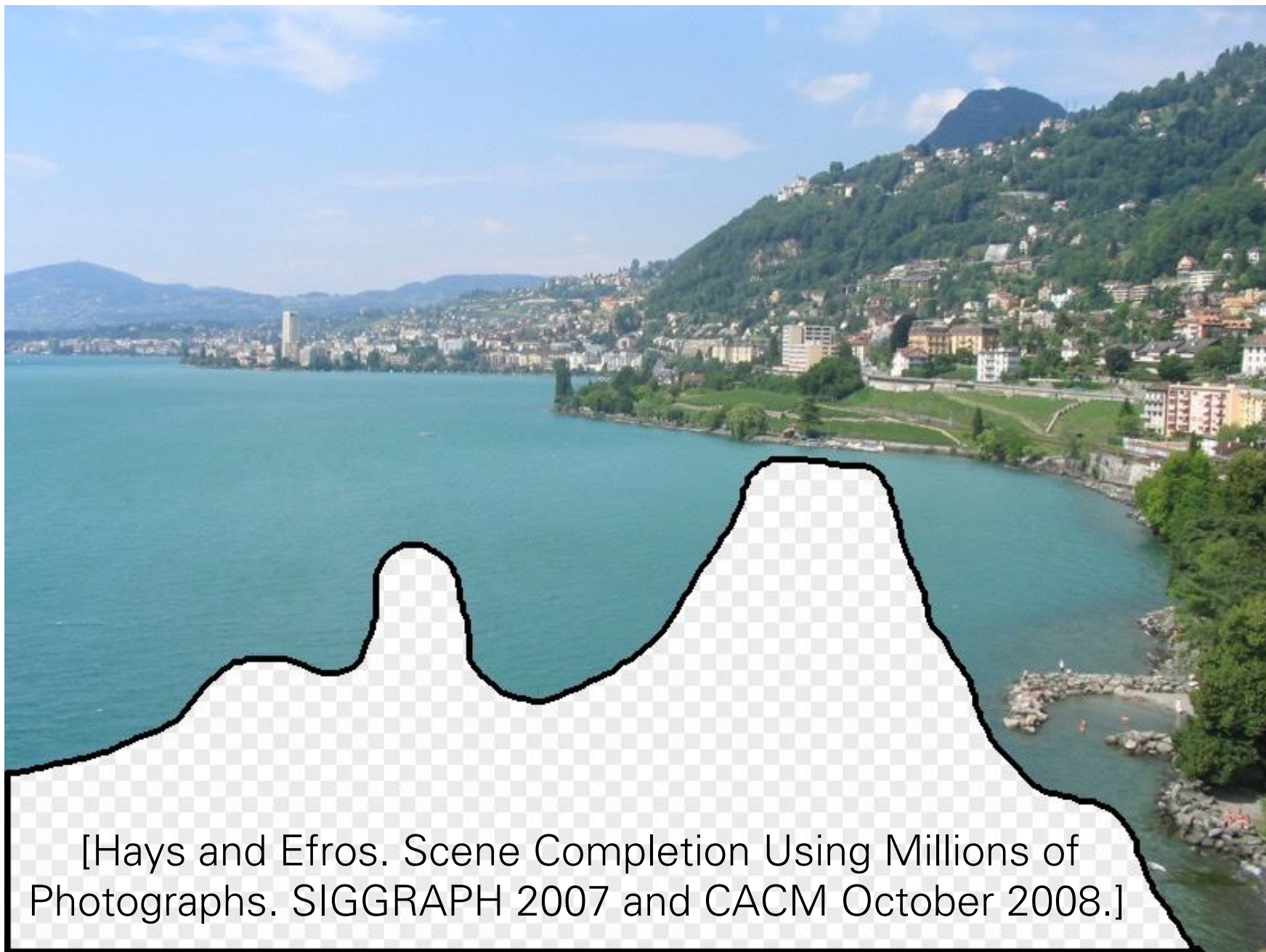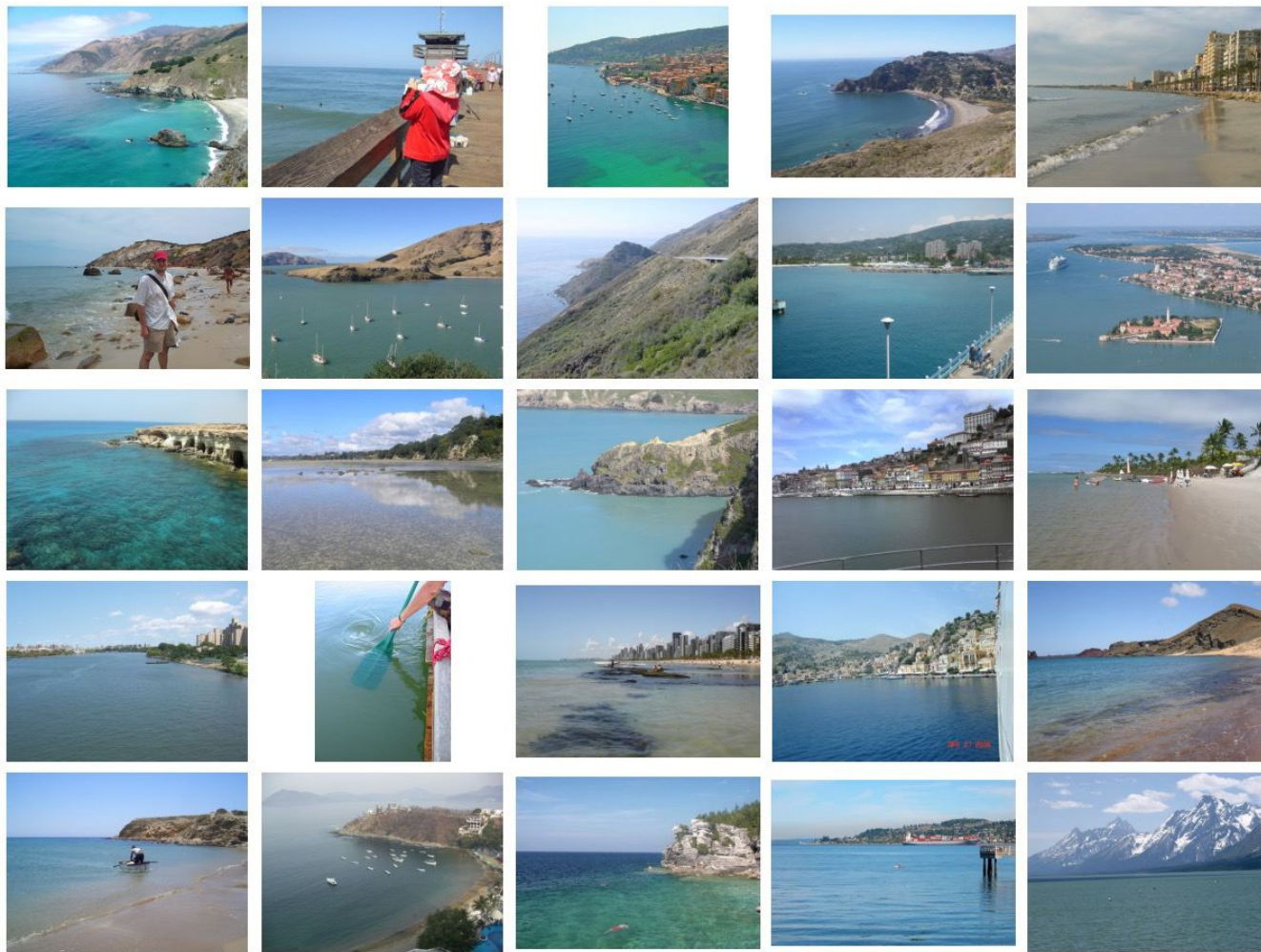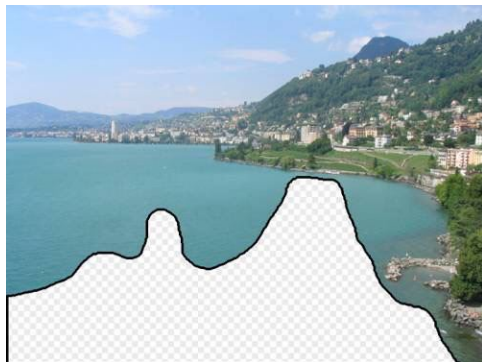


Learning denoising



Learning deblurring



HDR Imaging
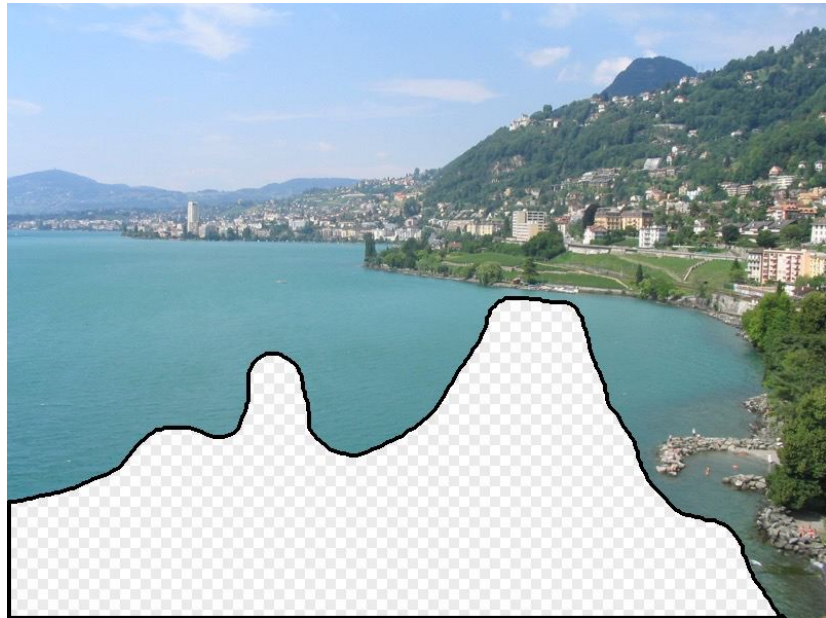
# Unreasonable Effectiveness of Data

[Hays and Efros. Scene Completion Using Millions of Photographs. SIGGRAPH 2007 and CACM October 2008.]
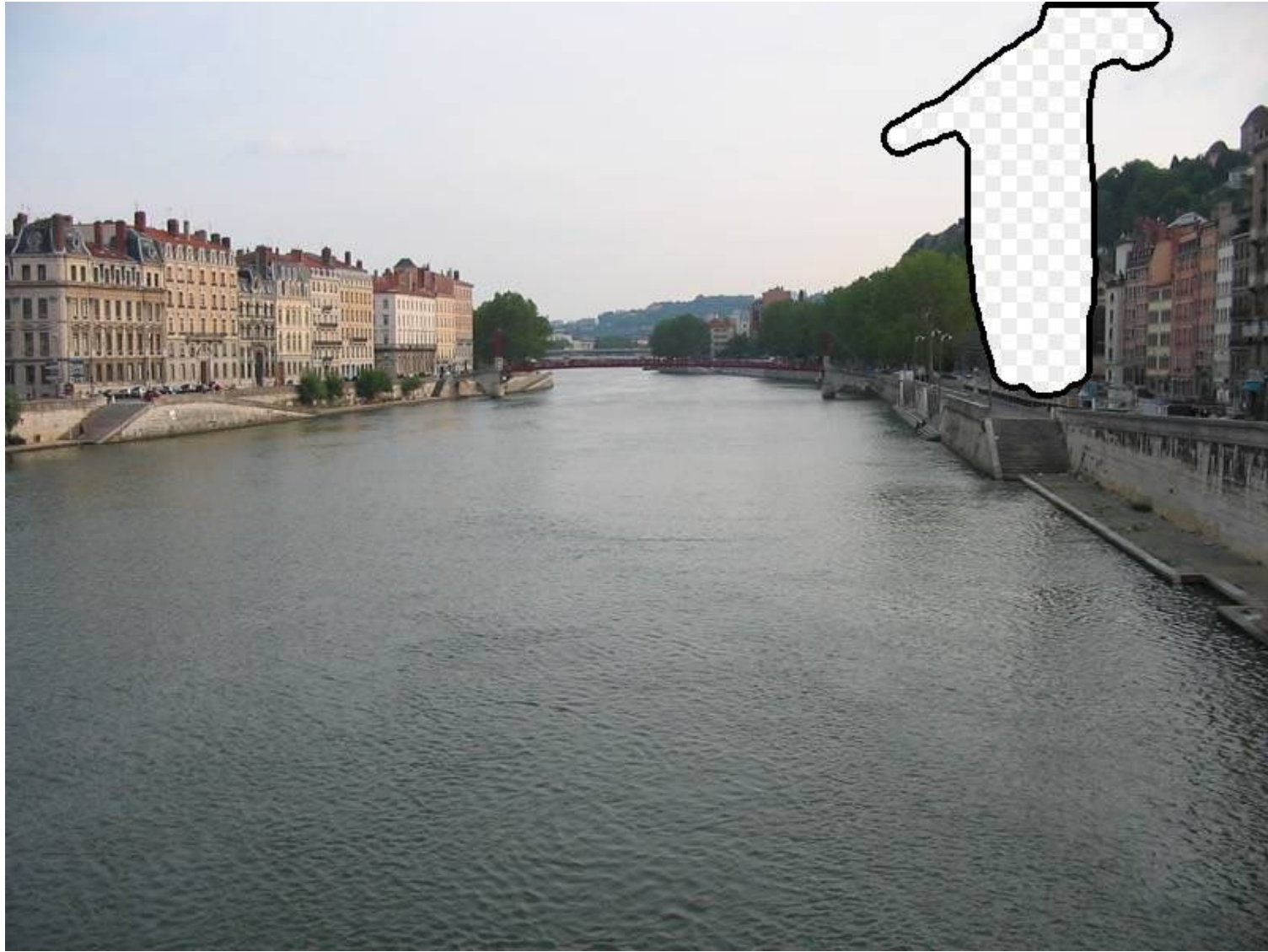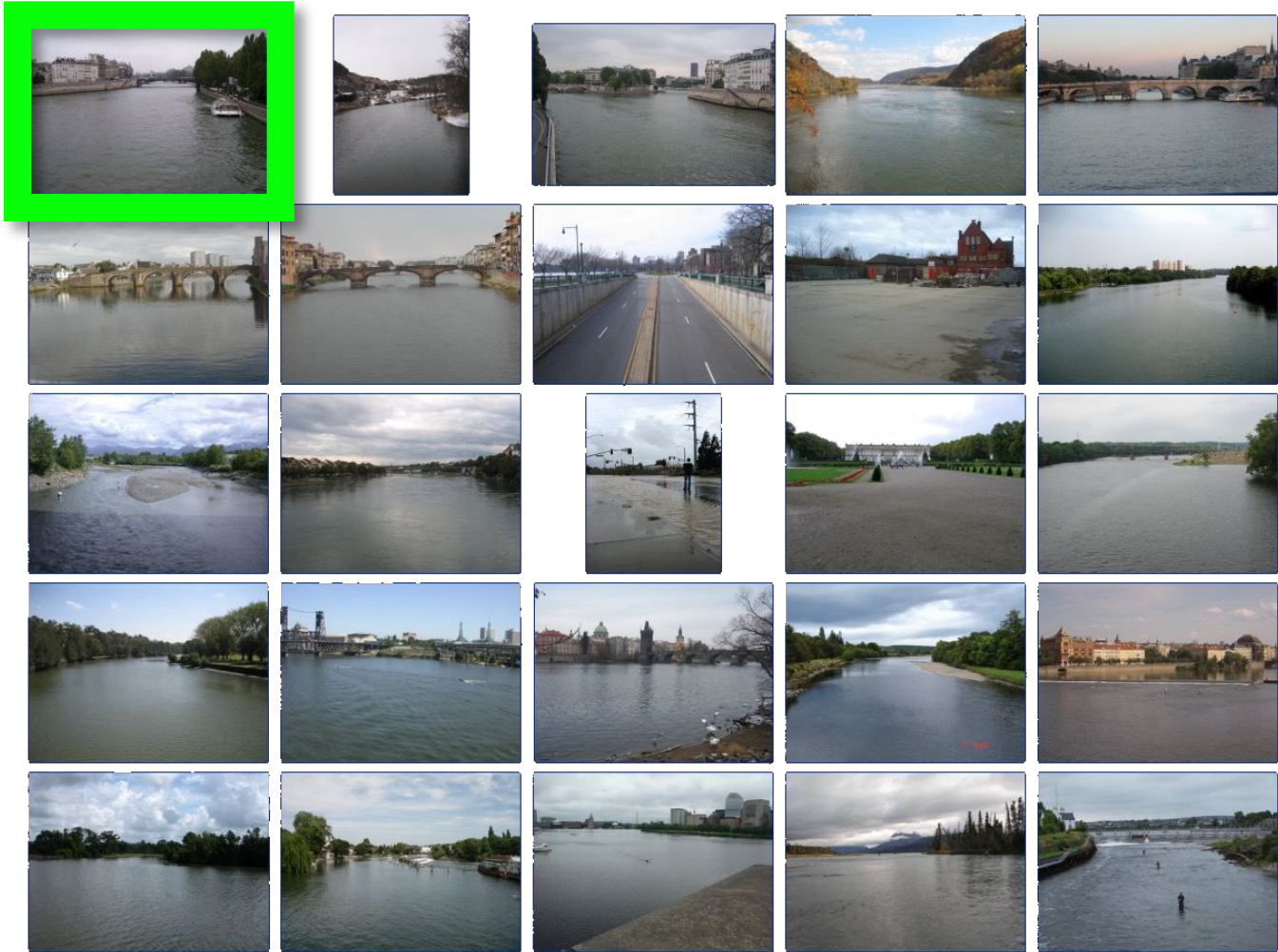
# 2 Million Flickr Images

… 200 total

14

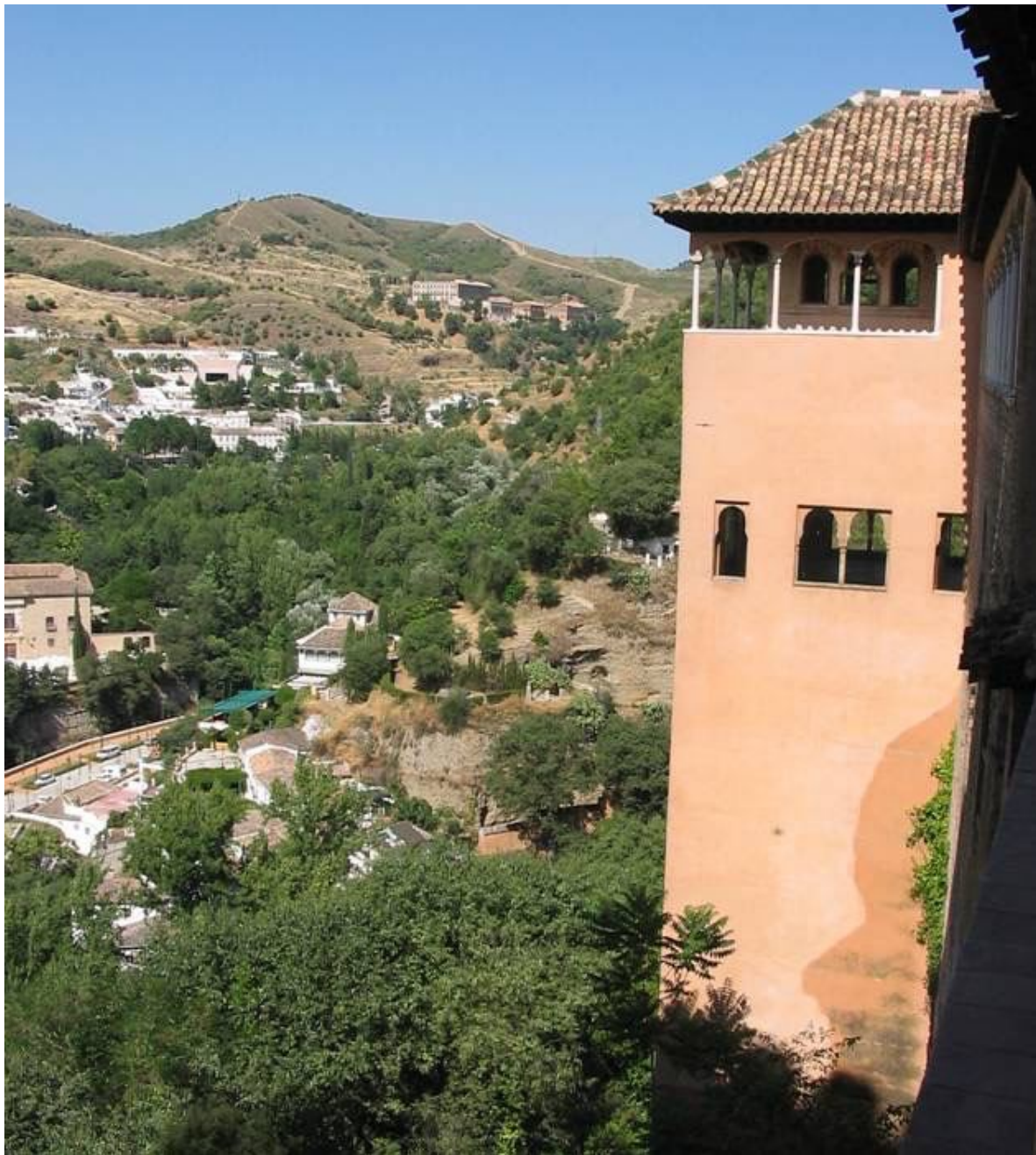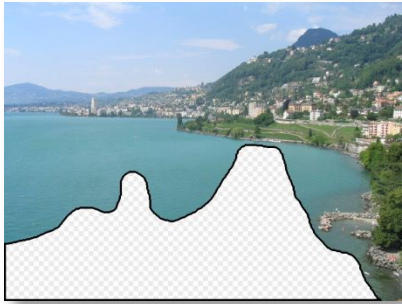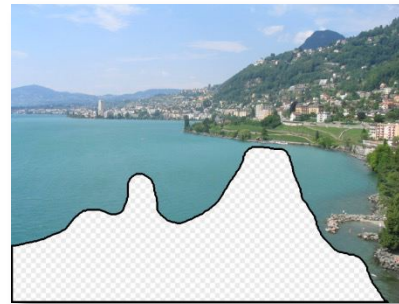… 200 scene matches

# Why does it work?

Nearest neighbors from a
collection of 20 thousand images

Nearest neighbors from a collection of 2 million images

# "Unreasonable Effectiveness of Data"

[Halevy, Norvig, Pereira 2009]

- Parts of our world can be explained by elegant mathematics

  physics, chemistry, astronomy, etc.

- But much cannot

  psychology, economics, genetics, etc.

- Enter <u>The Data</u>!

  Great advances in several fields:

  e.g., speech recognition, machine translation

  Case study: Google

"For many tasks, once we have a billion or so examples, we essentially have a closed set that represents (or at least approximates) what we need…"

Learning

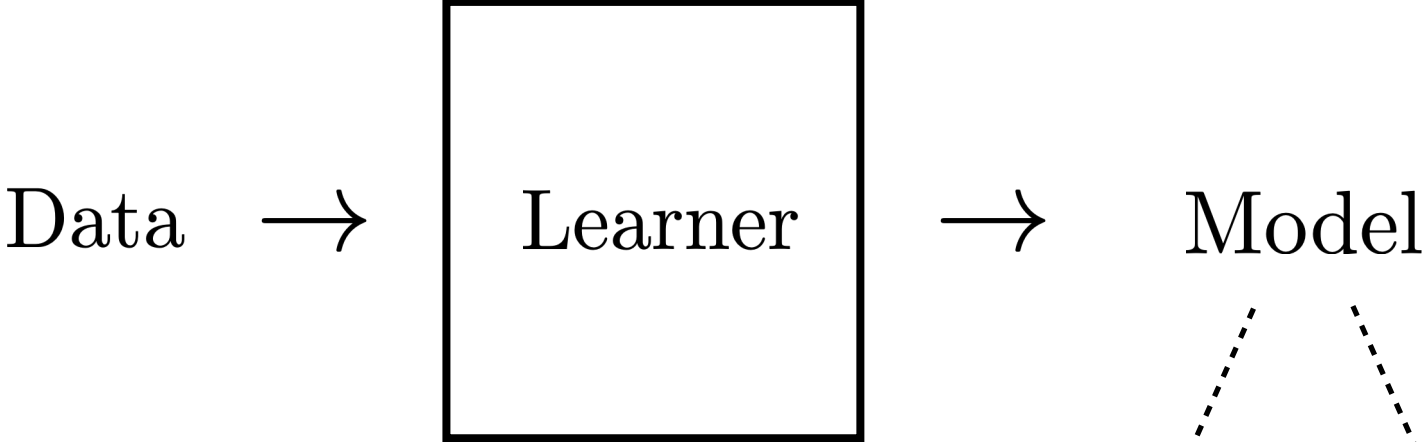Data → [ Learner ] → Model

Inference

Input → [ Model ] → Output

# A Brief History of Deep Learning

# Humble beginnings

$$\sigma\left(b + \sum_{i=1}^{n} a_i w_i\right)$$

- 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]

vorks/Long Short-Te
reiter Schmidhuber '97]

# Summer

- **2006:** First big suc

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

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



XRCE (2011)   AlexNet (2012)   ZFNet (2013)   GoogLeNet (2014)   ResNet (2015)   SENet (2017)

Human Performance Zone

["Mask RCNN", He et al. 2017]

| | | |
|---|---|---|
| *what color is the vase?* | *is the bus full of passengers?* | *is there a red shape above a circle?* |
| ```classify[color](     attend[vase])``` | ```measure[is](     combine[and](         attend[bus],         attend[full])``` | ```measure[is](     combine[and](         attend[red],         re-attend[above](             attend[circle])))``` |
| green (green) | yes (yes) | no (no) |

["Neural module networks", Andreas et al. 2017]

INPUT

OUTPUT

pix2pix

process

Ivy Tasi @ivymyt

Vitaly Vidmirov @vvid

["pix2pix", Isola et al. 2017]

# 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)

40

# Image transformations

X

Input image

Edge normals

Edge strength

3D orientation

Contact edges

Depth discontinuities

Y

Z

1. From pixels to edges

2. From edges to geometric primitives

# Object recognition



Edges

Texture

Colors

Segments

Parts

"clown fish"

$\phi_k(x)$

Feature extractors

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

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



Learned

"clown fish"

Neural net

# Object recognition

Learned

"clown fish"

Deep neural net

# Deep learning

$\mathbf{y}_i$

"clown fish"

$\mathbf{x}_i$

Loss

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

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

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

47

# 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



$J(\theta)$

$\theta_1$

$\theta_2$

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

# 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 Neural Nets

# 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

bias

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

# Computation in a neural net

## Linear layer

Input representation

Output representation

$\mathbf{x}$

$\mathbf{w}_j$

$b_j$

1

$y_j$

weights

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

bias

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

parameters of the model

# 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

$\mathbf{x}$

$\mathbf{w}$

$b$

$y \quad g(y)$

$g(y)$

Pointwise Non-linearity

1

# Example: linear classification with a perceptron



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

$$y = g(z)$$

One layer neural net (perceptron) can perform linear classification!

# Example: linear classification with a perceptron

Training data

Bad fit
7 misclassifications

Okay fit
4 misclassifications

Good fit
0 misclassifications

$$\mathbf{w}^*, b^* = \arg\min_{\mathbf{w},b} \sum_{i=1}^{N} \mathcal{L}(g(z^{(i)}), y^{(i)})$$

# Computation in a neural net

Input
representation

Output
representation

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

# Computation in a neural net – nonlinearity

Input
representation

Output
representation

Sigmoid

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

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

$b$

$1$

$g(y)$

# 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

$$\text{tanh}(x) = 2 \text{ sigmoid}(2x) - 1$$

Tanh

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

# Computation in a neural net – nonlinearity

- Unbounded output (on positive side)

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

- 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

$\mathbf{h}$

$\mathbf{x}$

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

$b_j^{(1)}$

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

$b_j^{(2)}$

$\mathbf{y}$

1

1

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

# Stacking layers

Input
representation

Intermediate
representation

Output
representation

$$\mathbf{z} \quad \mathbf{h} = g(\mathbf{z})$$



$\mathbf{x}$

$\mathbf{W}_{1_j}$

$b_{1_j}$

$\mathbf{W}_{2_j}$

$b_{2_j}$

$\mathbf{y}$

1

1

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

# Stacking layers

Input representation

Intermediate representation

Output representation

$$\mathbf{h} = g(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \qquad \mathbf{y} = g(\mathbf{W}_2 \mathbf{h} + \mathbf{b}_2)$$

$$\theta = \{\mathbf{W}_1, \ldots, \mathbf{W}_L, \mathbf{b}_1, \ldots, \mathbf{b}_L\}$$

# Stacking layers



Input representation

Intermediate representation

Output representation

$\mathbf{W}_1$

$\mathbf{h}$

$\mathbf{W}_2$

$\mathbf{x}$

$\mathbf{b}_1$

$\mathbf{b}_2$

$\mathbf{y}$

positive

negative

$$\mathbf{h} = g(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1) \qquad \mathbf{y} = g(\mathbf{W}_2\mathbf{h} + \mathbf{b}_2)$$

$$\theta = \{\mathbf{W}_1, \dots, \mathbf{W}_L, \mathbf{b}_1, \dots, \mathbf{b}_L\}$$

# Stacking layers

Input representation

Intermediate representation

Output representation

$$\mathbf{h}$$

$$\mathbf{W}_1 \qquad \mathbf{W}_2$$

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

positive

negative

$$1 \qquad \mathbf{b}_1 \qquad 1 \qquad \mathbf{b}_2$$

$$\mathbf{h} = g(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \qquad \mathbf{y} = g(\mathbf{W}_2 \mathbf{h} + \mathbf{b}_2)$$

$$\theta = \{\mathbf{W}_1, \dots, \mathbf{W}_L, \mathbf{b}_1, \dots, \mathbf{b}_L\}$$

# Stacking layers



Input representation

Intermediate representation

Output representation

$\mathbf{h}$

$\mathbf{W}_1$

$\mathbf{W}_2$

$\mathbf{x}$

$\mathbf{y}$

positive

negative

$1$

$\mathbf{b}_1$

$1$

$\mathbf{b}_2$

$$\mathbf{h} = g(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1) \qquad \mathbf{y} = g(\mathbf{W}_2\mathbf{h} + \mathbf{b}_2)$$

$$\theta = \{\mathbf{W}_1, \dots, \mathbf{W}_L, \mathbf{b}_1, \dots, \mathbf{b}_L\}$$

# Stacking layers

Input
representation

Intermediate
representation

Output
representation

$\mathbf{h}$

$\mathbf{W}_1$

$\mathbf{W}_2$

$\mathbf{x}$

$\mathbf{y}$

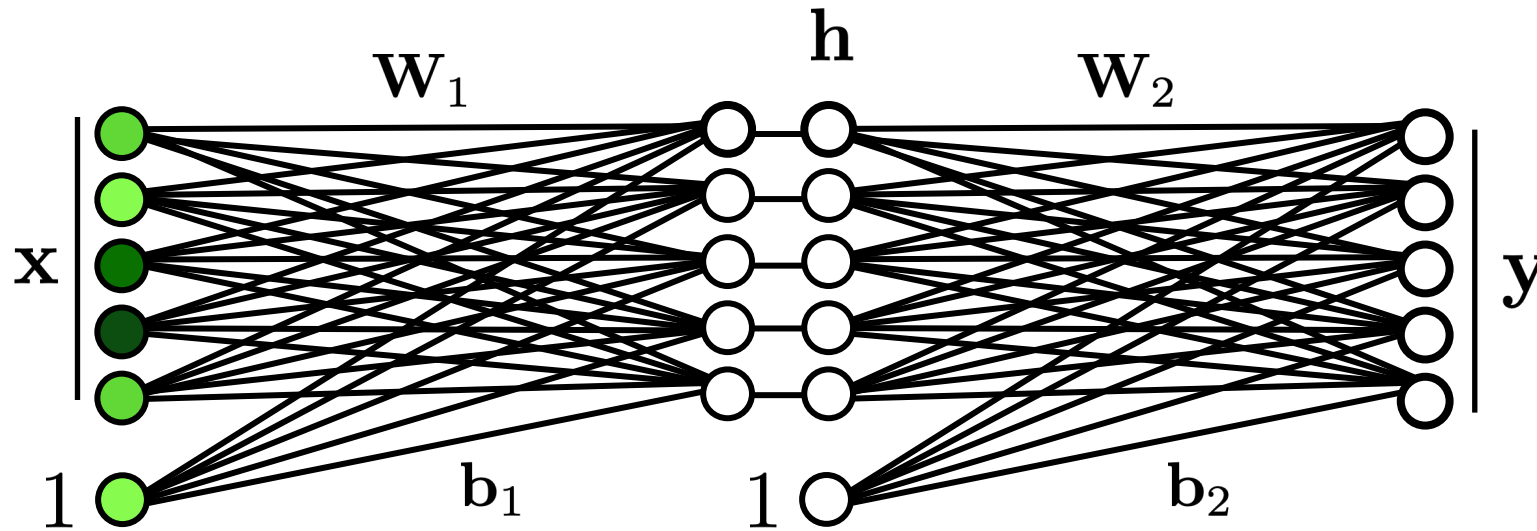positive

negative
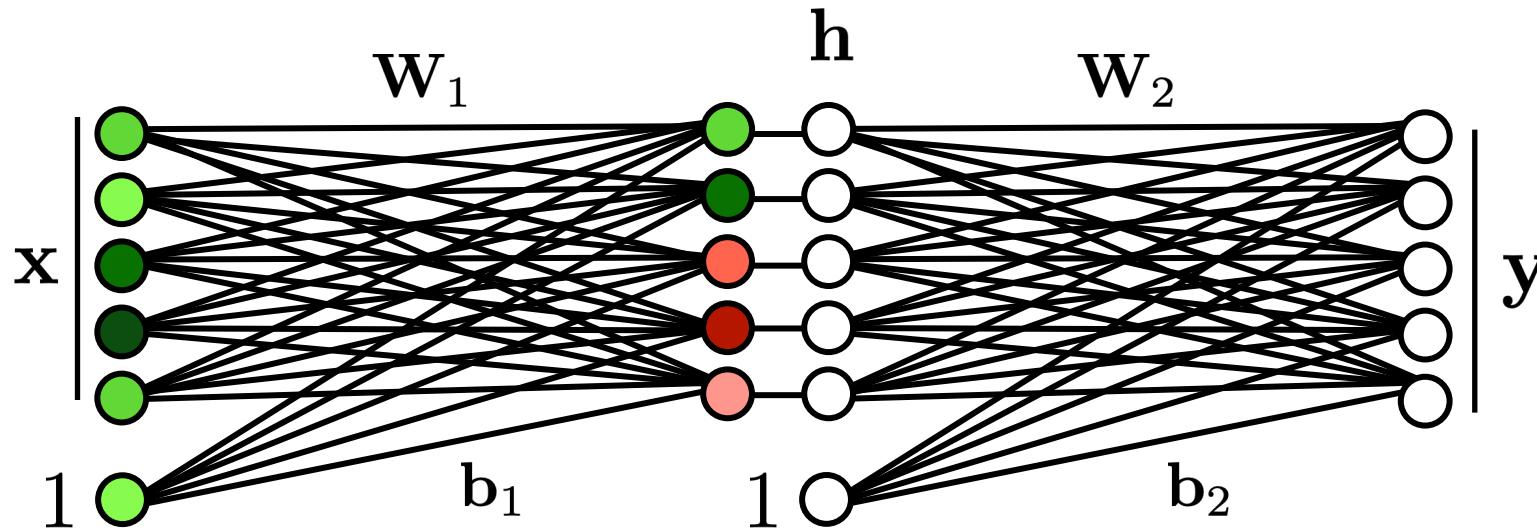
1

$\mathbf{b}_1$

1

$\mathbf{b}_2$

$$\mathbf{h} = g(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1) \qquad \mathbf{y} = g(\mathbf{W}_2\mathbf{h} + \mathbf{b}_2)$$

$$\theta = \{\mathbf{W}_1, \ldots, \mathbf{W}_L, \mathbf{b}_1, \ldots, \mathbf{b}_L\}$$

# Connectivity Patterns

Input
representation

Output
representation

$\mathbf{W}_1$

$\mathbf{x}$

$\mathbf{y}$

*Fully connected layer*

Input
representation

Output
representation

$\mathbf{W}_2$

$\mathbf{x}$

$\mathbf{y}$

*Locally connected layer
(Sparse W)*

[http://playground.tensorflow.org]

# Deep nets

Linear

Non-linearity

Classify

"clown fish"

$$f(\mathbf{x}) = f_L(f_{L-1}(\ldots f_2(f_1(\mathbf{x}))))$$

# Example: linear classification with a perceptron



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

$$y = g(z)$$

One layer neural net (perceptron) can perform linear classification!

# Example: nonlinear classification with a deep net



$$\mathbf{z} = \mathbf{W}_1\mathbf{x} + \mathbf{b}_1$$

$$\mathbf{h} = g(\mathbf{z})$$

$$z_3 = \mathbf{W}_2\mathbf{h} + b_2$$

$$y = 1(z_3 > 0)$$

# 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.

# Deep nets

Linear

Non-linearity

Classify

"clown fish"

$$f(\mathbf{x}) = f_L(f_{L-1}(\ldots f_2(f_1(\mathbf{x}))))$$

# Classifier layer

Last layer

dolphin

cat

grizzly bear

angel fish

··· ⇒
⇒
⇒

angel fish

chameleon                    argmax                    "clown fish"

**clown fish**

iguana

elephant

⋮

# 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

"clown fish"

dolphin

cat

grizzly bear

...

angel fish

Loss → small

chameleon

**clown fish**

iguana

elephant

# Loss function

Network output

Ground truth label



dolphin

cat

grizzly bear

"grizzly bear"

angel fish

...

chameleon

Loss → **large**

**clown fish**

iguana

elephant

:

# Loss function

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

Ground truth label  $\mathbf{y}$

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

$\mathbf{x}$

$f$

| Prediction | | Ground truth label |
|---|---|---|
| dolphin | | dolphin |
| cat | | cat |
| grizzly bear | | grizzly bear |
| angel fish | | angel fish |
| chameleon | ⊙ | chameleon |
| **clown fish** | | **clown fish** |
| iguana | | iguana |
| elephant | | elephant |
| ⋮ | | ⋮ |

0                                1        0                                1

# Loss function

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

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

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

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

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

86

# Deep learning

$\mathbf{y}_2$

"grizzly bear"

$\mathbf{x}_2$



Loss

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

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

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

87

# Deep learning



$\mathbf{y}_i$

"chameleon"

$\mathbf{x}_i$

Learned

Loss

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

$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6$
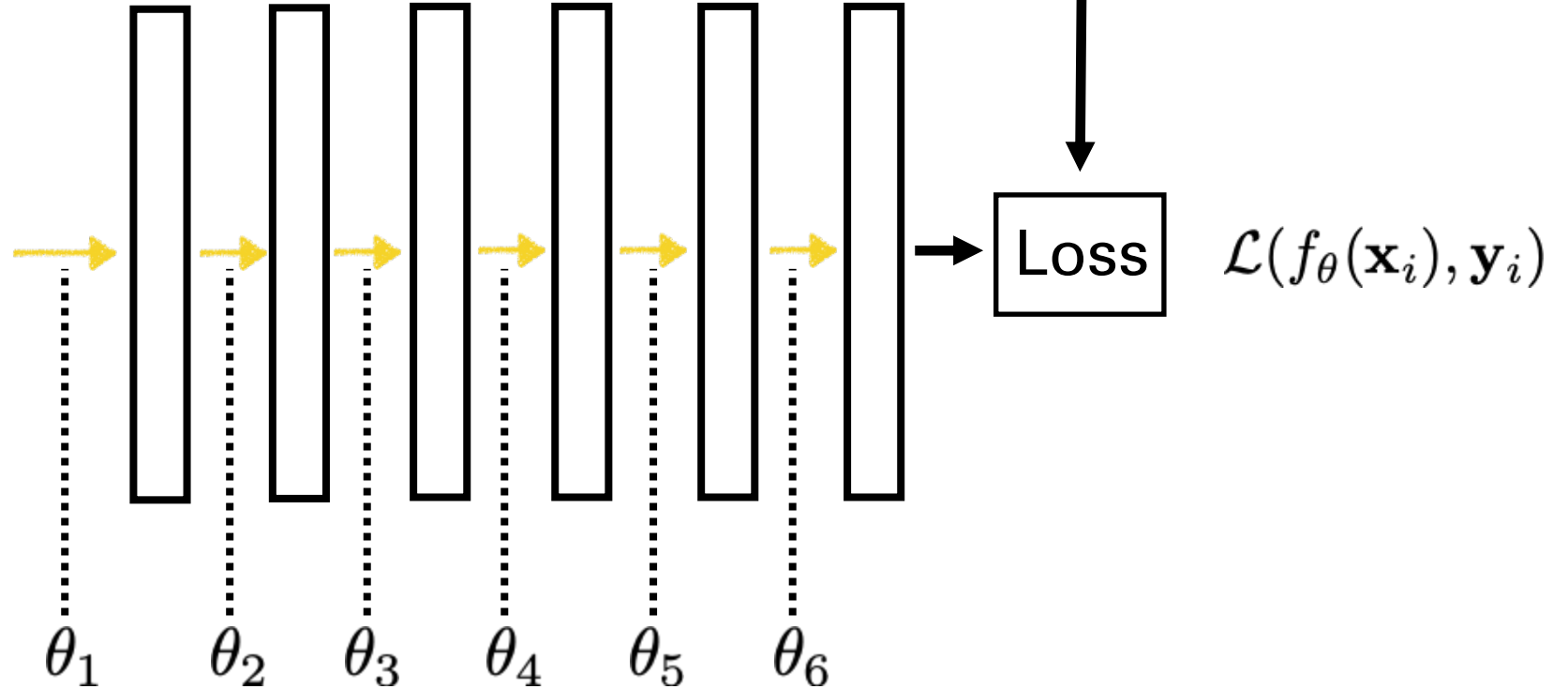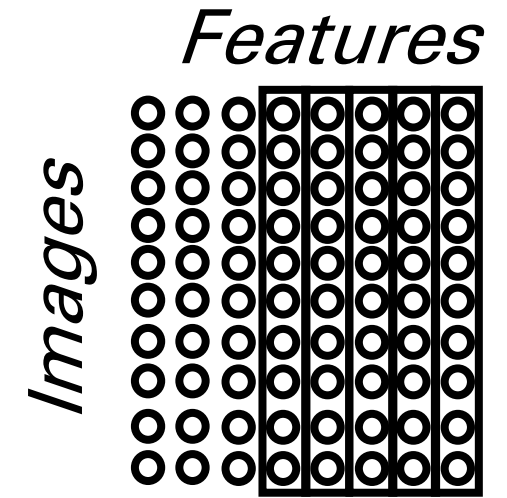
$$\theta^* = \arg\min_\theta \sum_{i=1}^{N} \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)$$

# Batch (parallel) processing

*Images*

Loss

Loss

Loss

$\Sigma$

# Tensors (multi-dimensional arrays)



*Each layer is a representation of the data*

# Everything is a tensor

$$\mathbf{z} = \mathbf{W}_1\mathbf{x} + \mathbf{b}_1$$

$$\mathbf{h} = g(\mathbf{z})$$

$$z_3 = \mathbf{W}_2\mathbf{h} + b_2$$

$$y = 1(z_3 > 0)$$

Tensor processing with batch size = 3:

# "Tensor flow"

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

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

# Regularizing deep nets

Deep nets have millions of parameters!

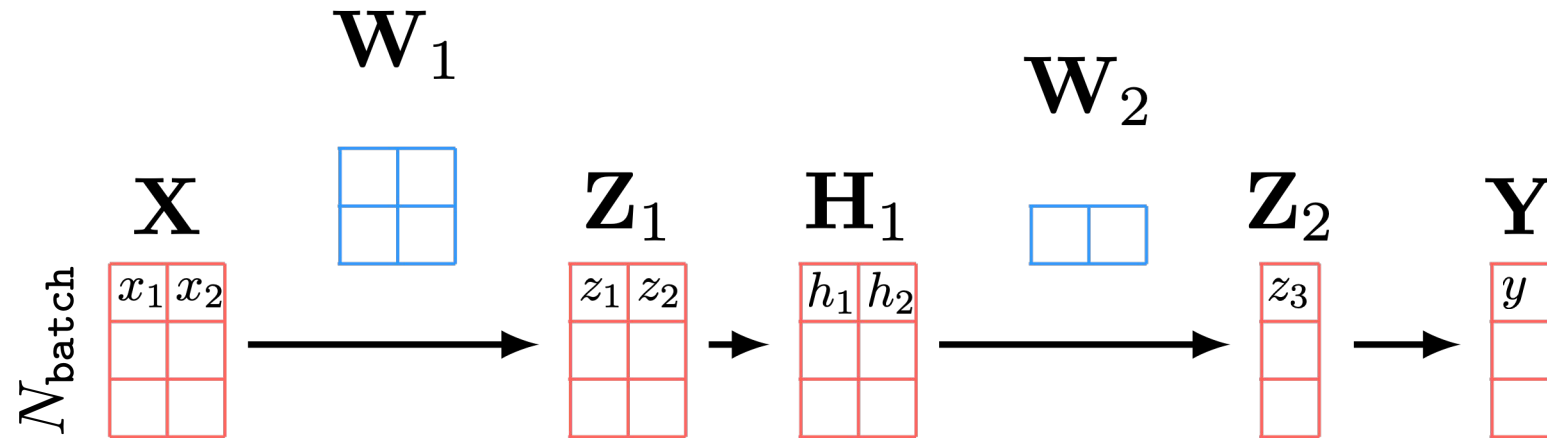On many datasets, it is easy to overfit — we may have more free parameters than data points to constrain them.

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

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

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

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

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

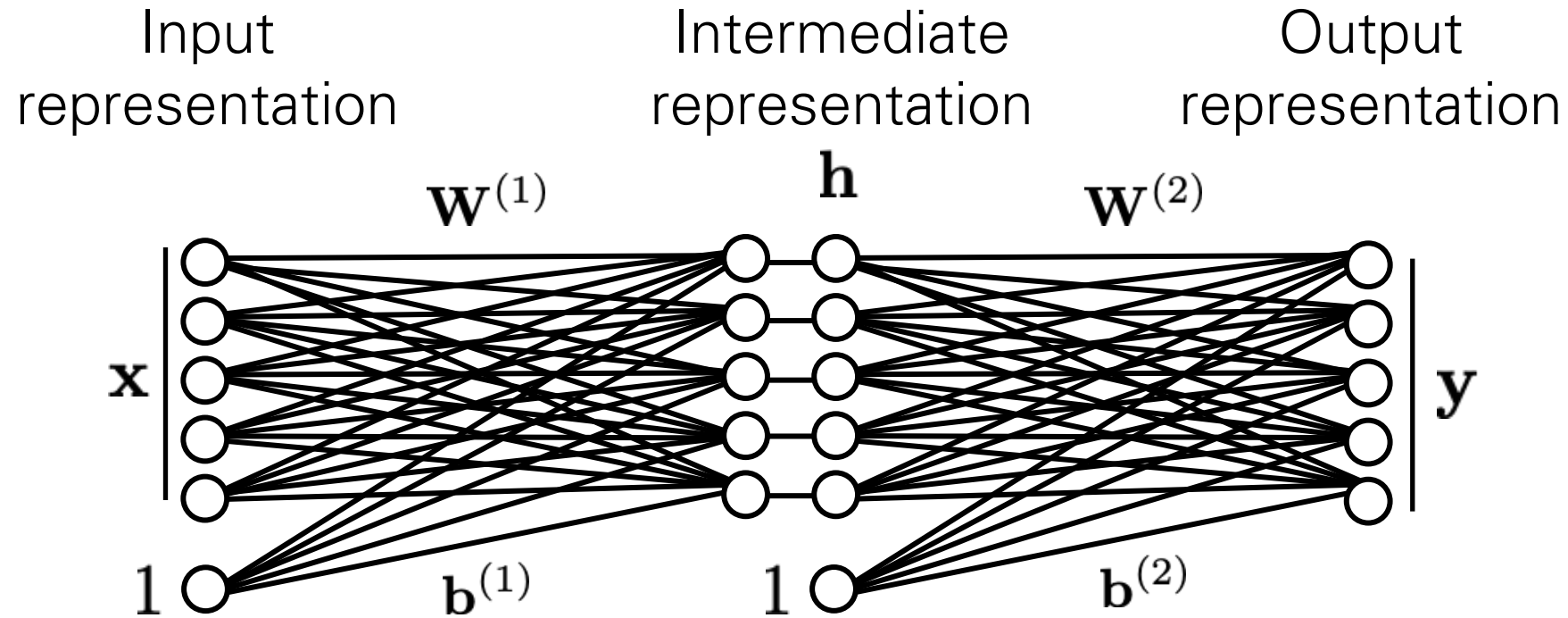Probabilistic interpretation: R is a Gaussian **prior** over values of the parameters.

# Recall: regularized least squares

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

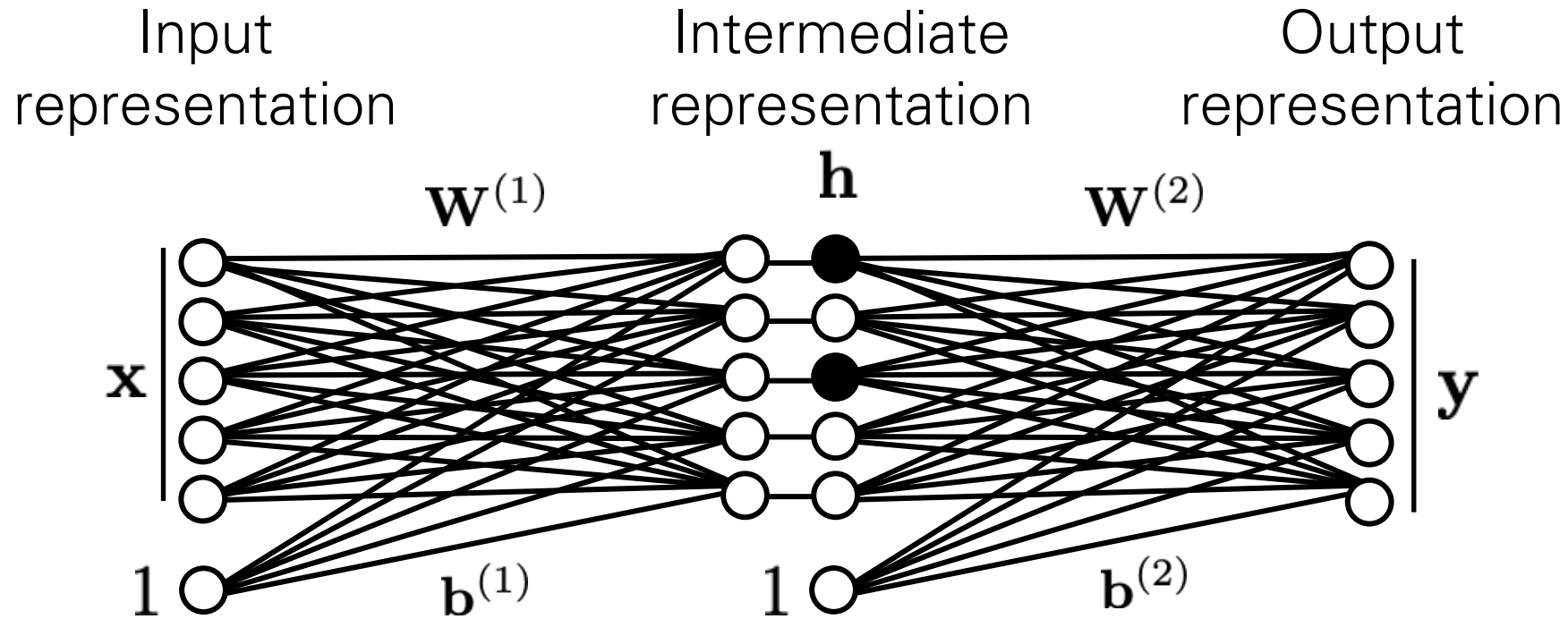$$R(\mathbf{W}) = \lambda \|\mathbf{W}\|_2^2 \quad \longleftarrow \quad \text{weight decay}$$

"We prefer to keep weights small."

# Dropout

Input
representation

Intermediate
representation

Output
representation

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



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

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

# 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

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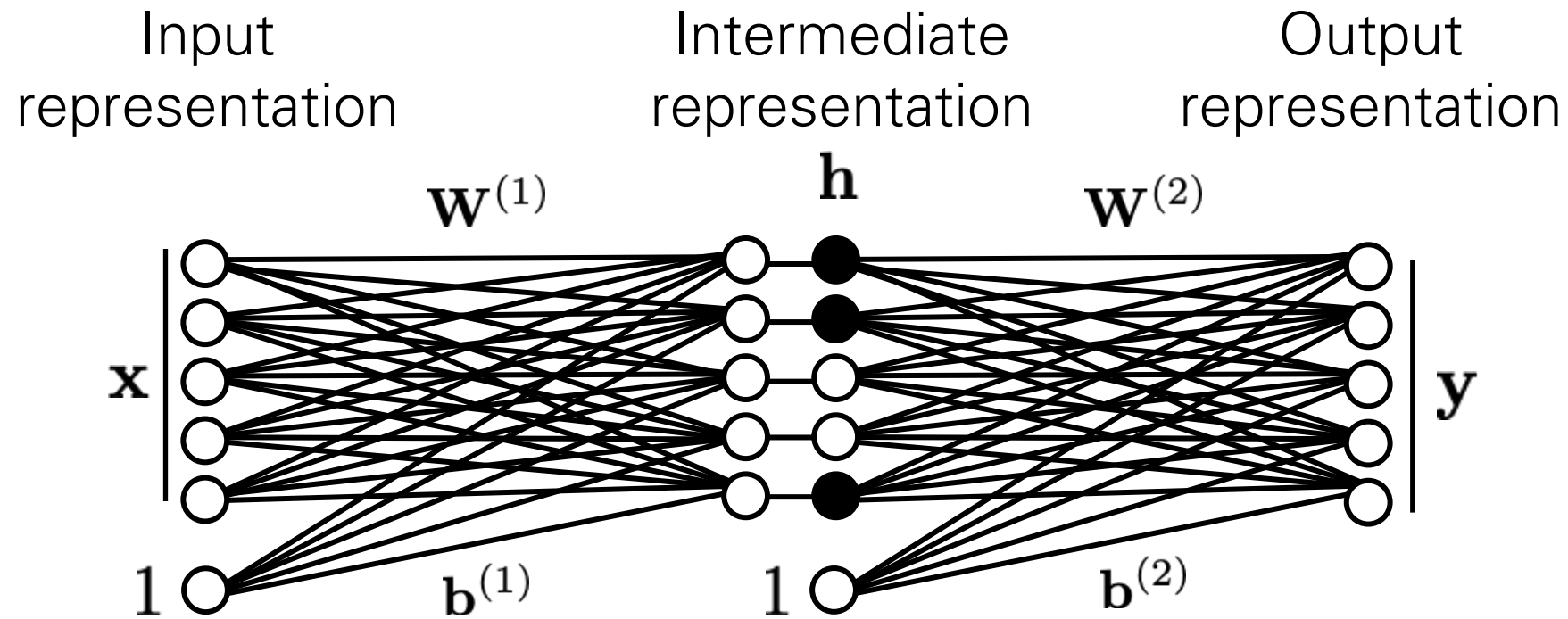
# Dropout

Input
representation

Intermediate
representation

Output
representation

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



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

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

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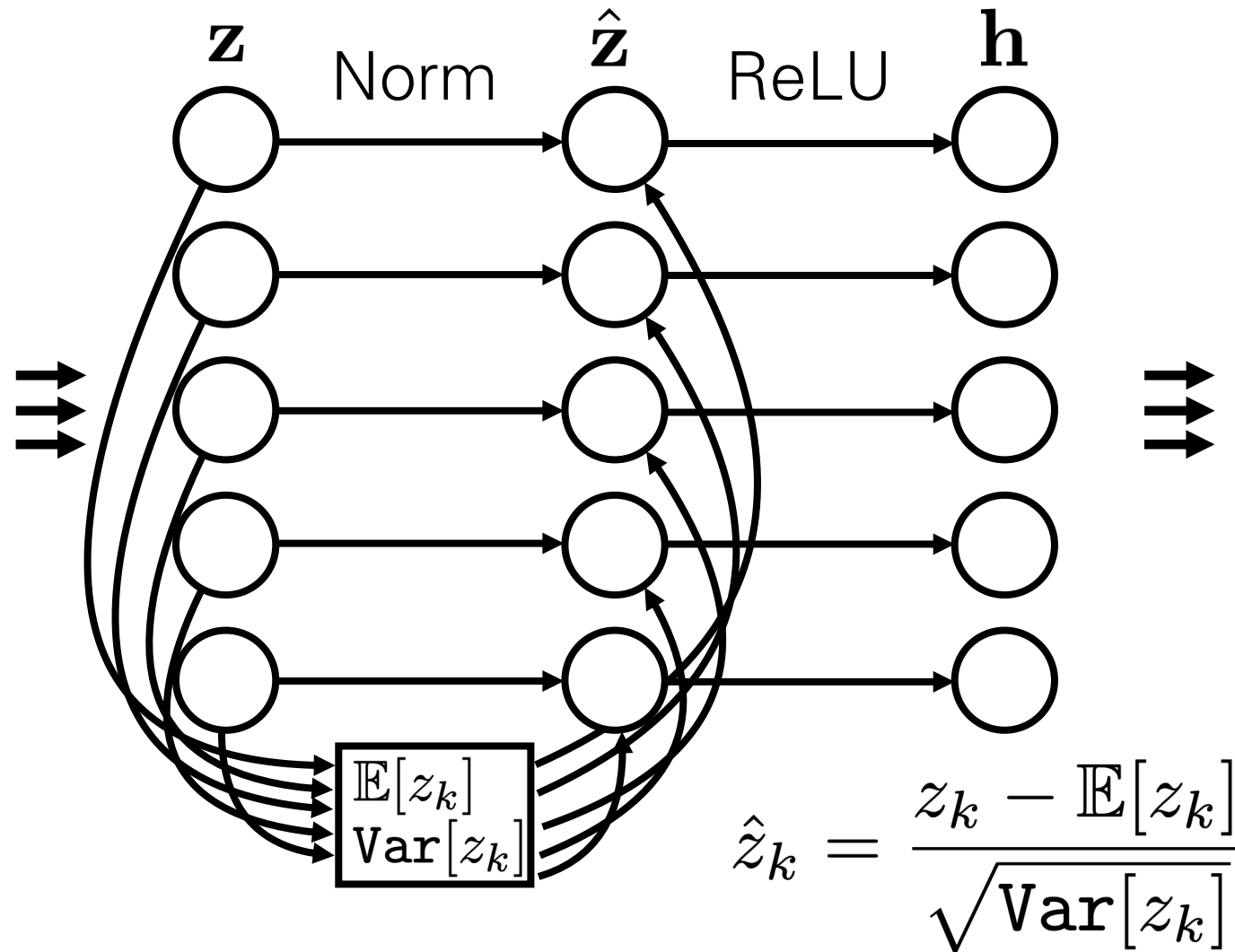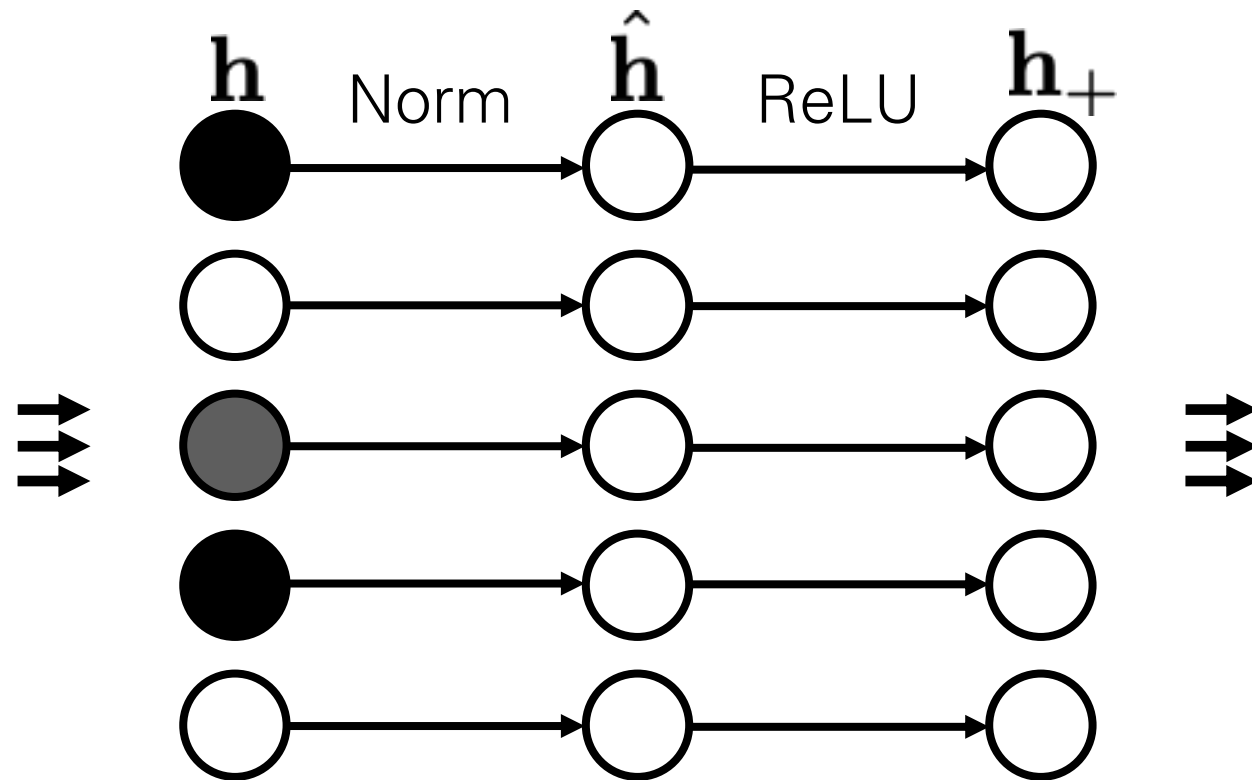
# Dropout

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.
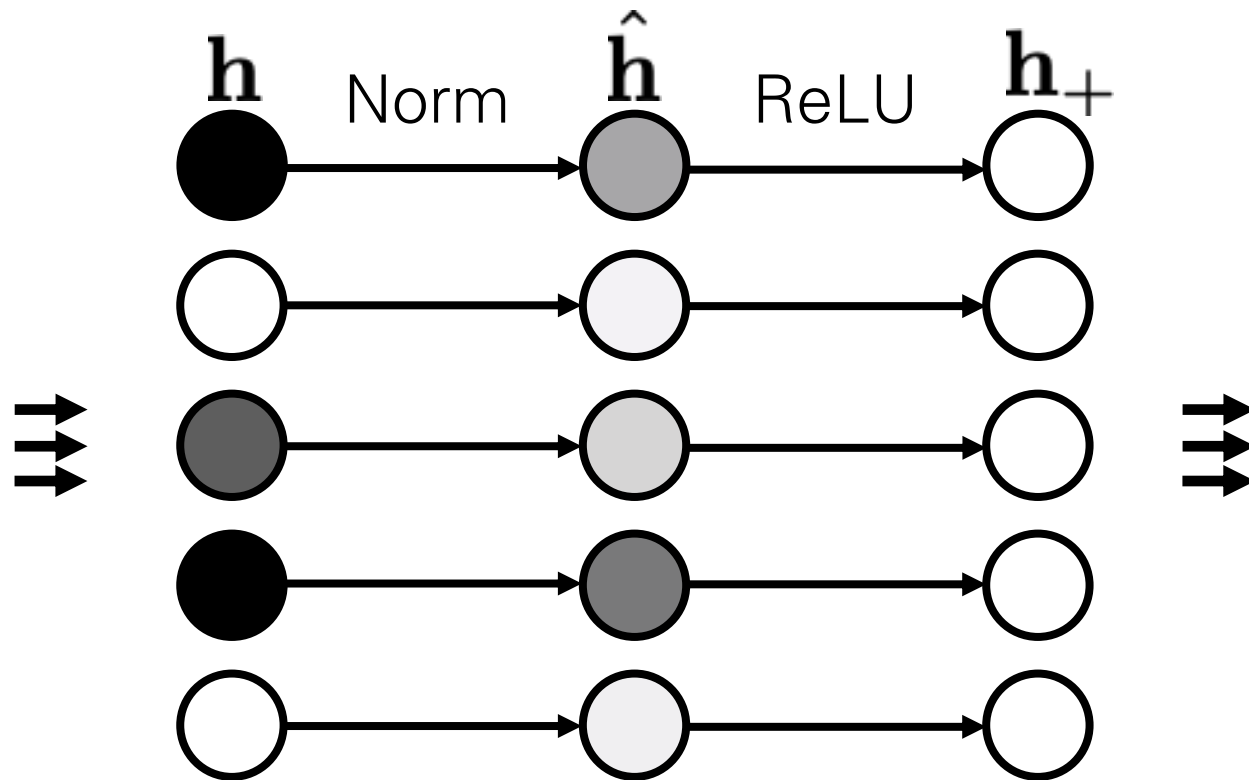
# Normalization layers



$$\hat{z}_k = \frac{z_k - \mathbb{E}[z_k]}{\sqrt{\text{Var}[z_k]}}$$

# Normalization layers

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

# Normalization layers

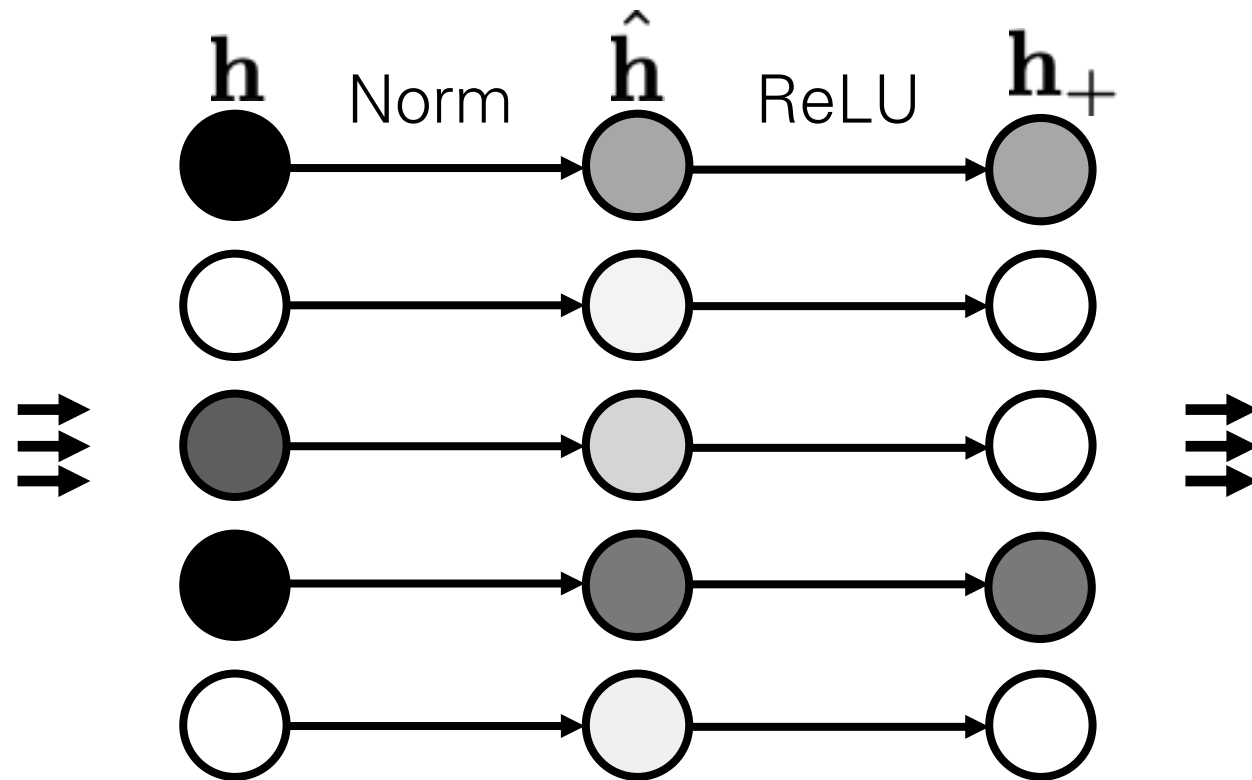

$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\mathrm{Var}[h_k]}}$$

# 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.
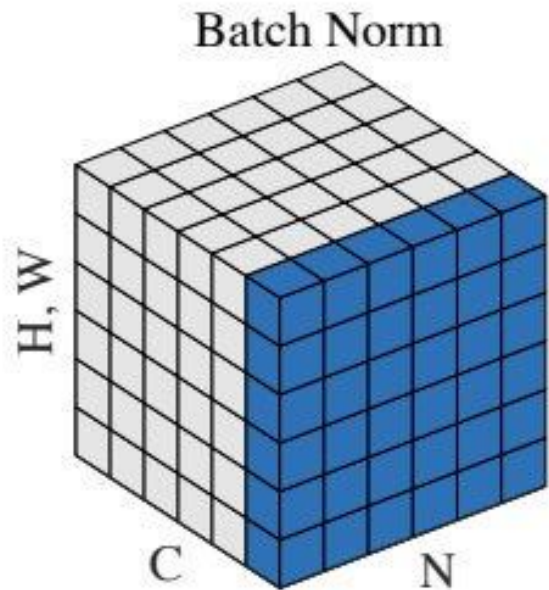
Standardize unit activations by subtracting mean and dividing by variance.

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

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

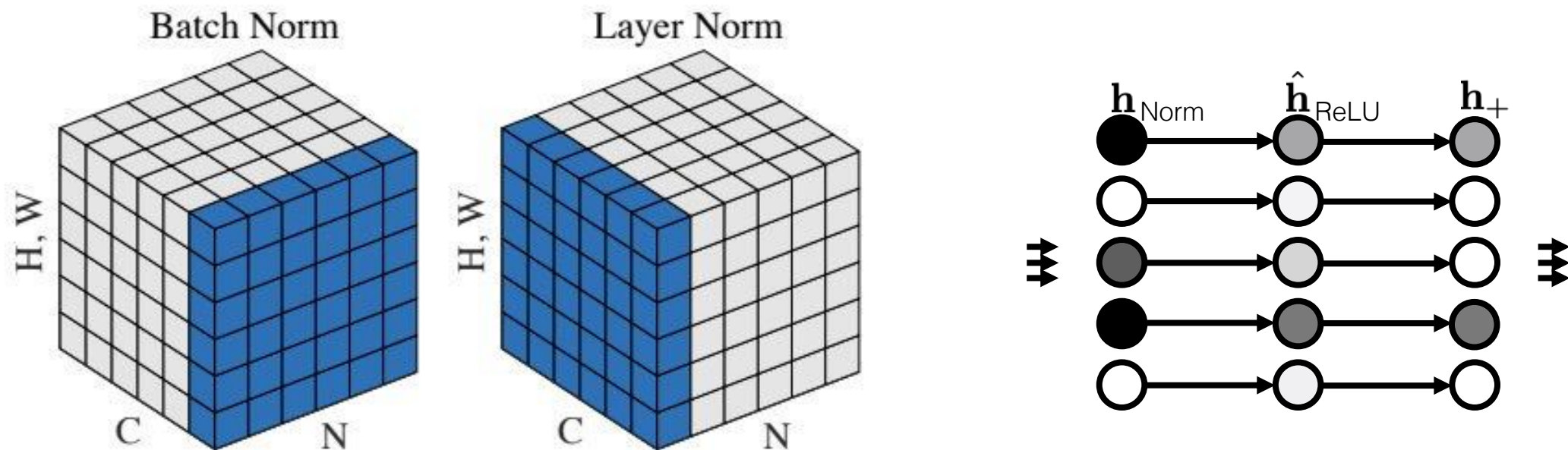Both these properties reduce the effective capacity of the model, i.e. regularize the model.

# Normalization layers

Batch Norm

H, W

C            N

Normalize w.r.t. a single hidden unit's pattern of activation over training examples (a batch of examples).
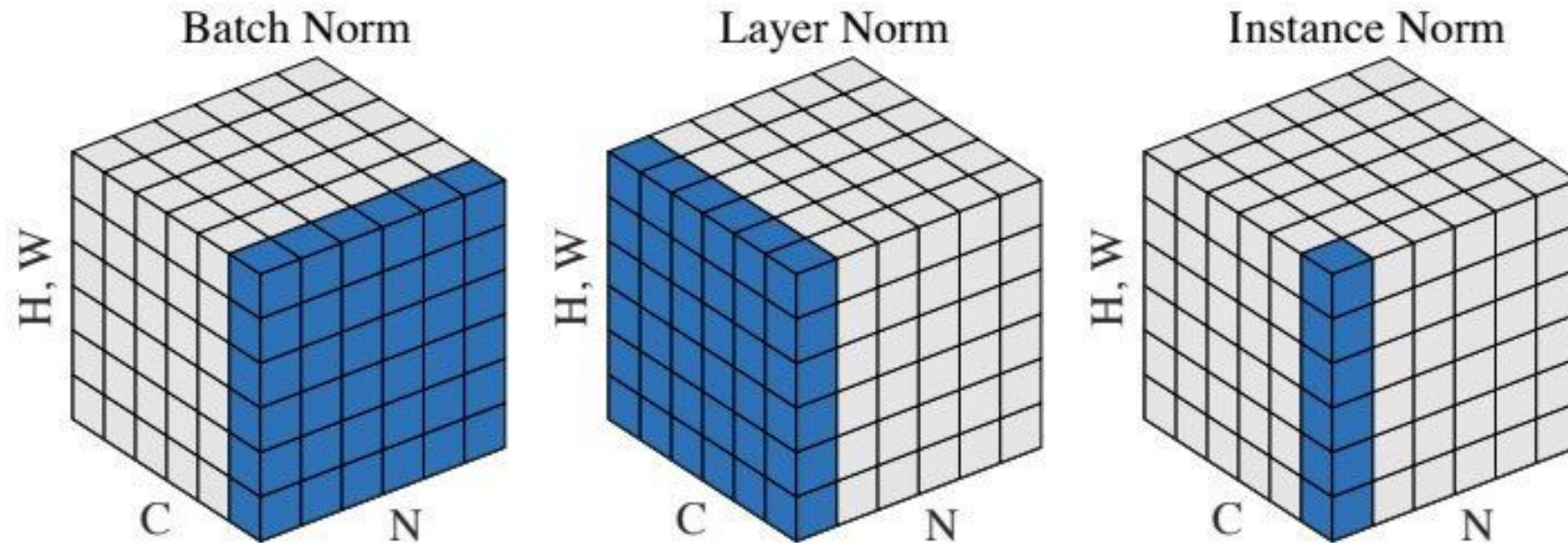
[Figure from Wu & He, arXiv 2018]

# Normalization layers



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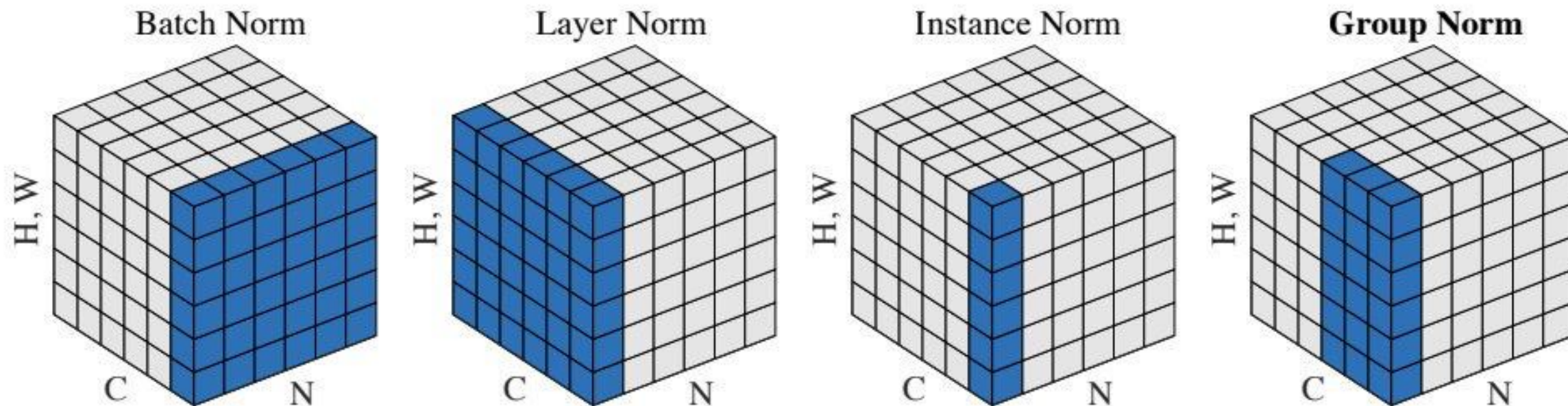
[Figure from Wu & He, arXiv 2018]

# 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.

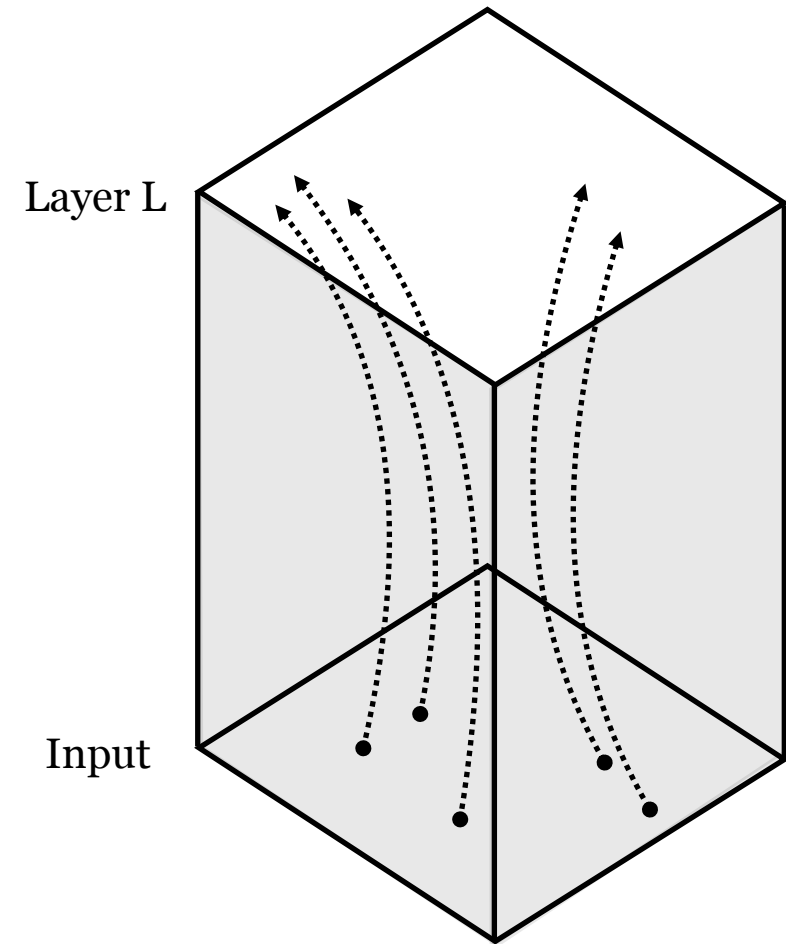[Figure from Wu & He, arXiv 2018]

# Normalization layers



Batch Norm     Layer Norm     Instance Norm     Group Norm
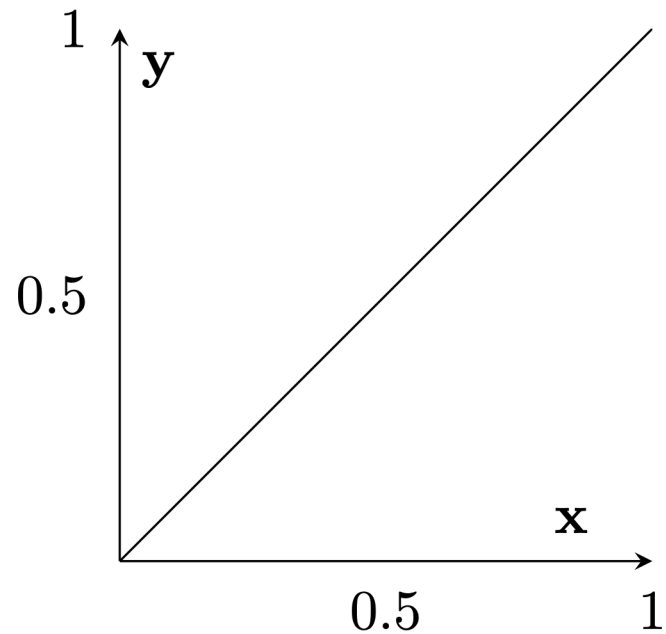
Might as well…

[Figure from Wu & He, arXiv 2018]
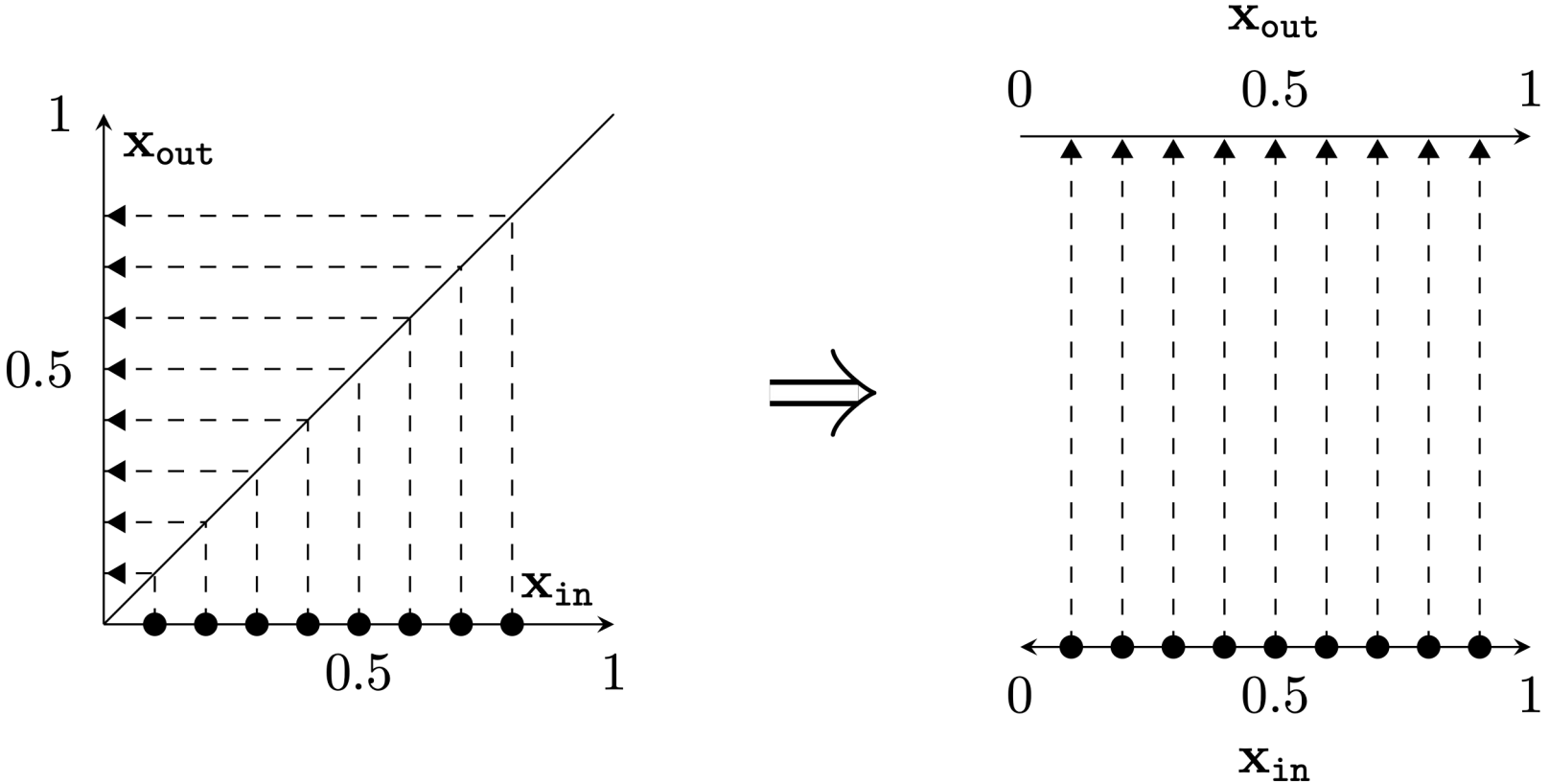
# Deep nets are data transformers

- Deep nets transform datapoints, layer by layer

- Each layer is a different representation of the data

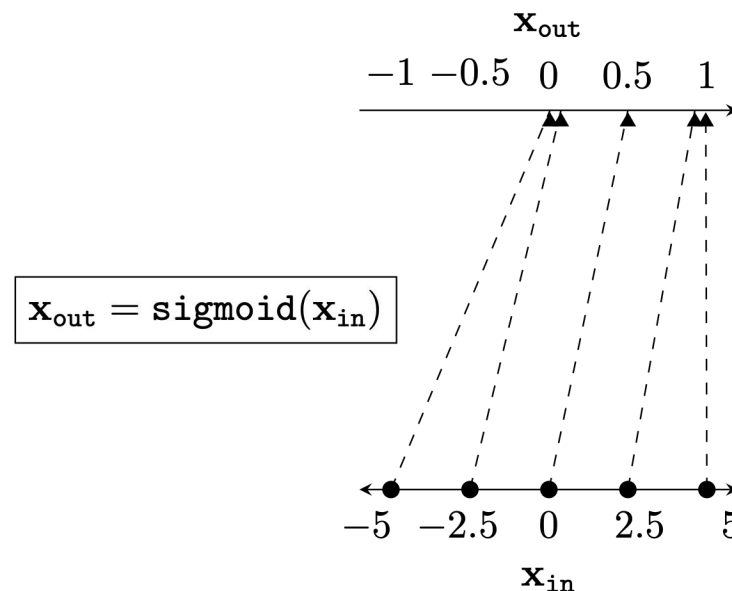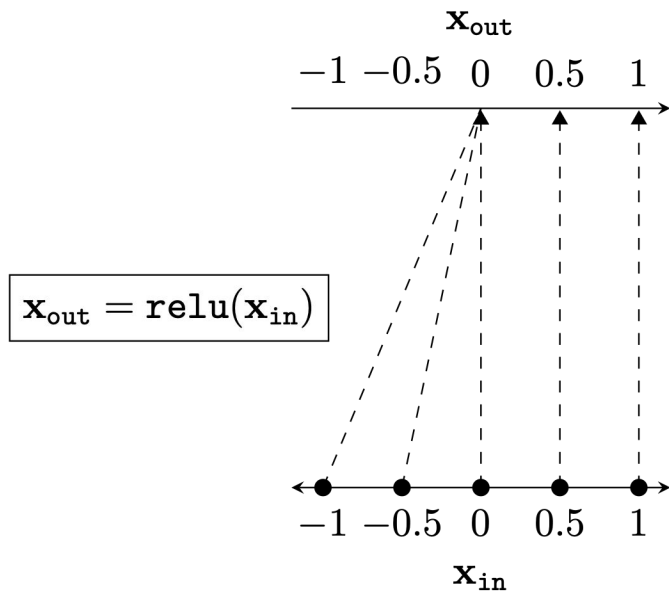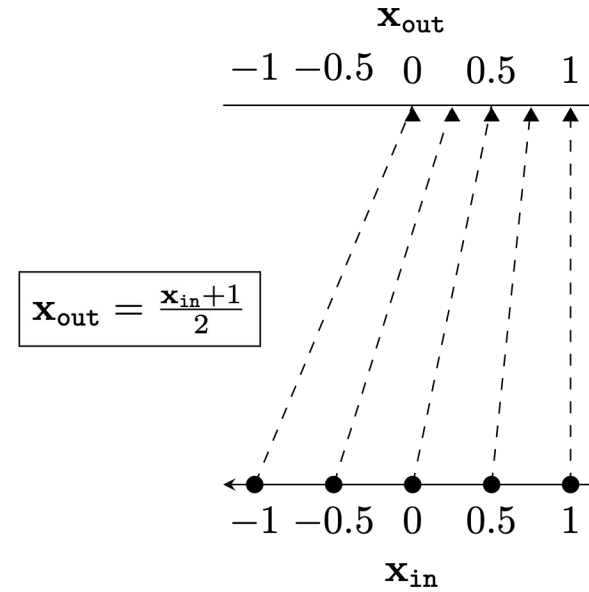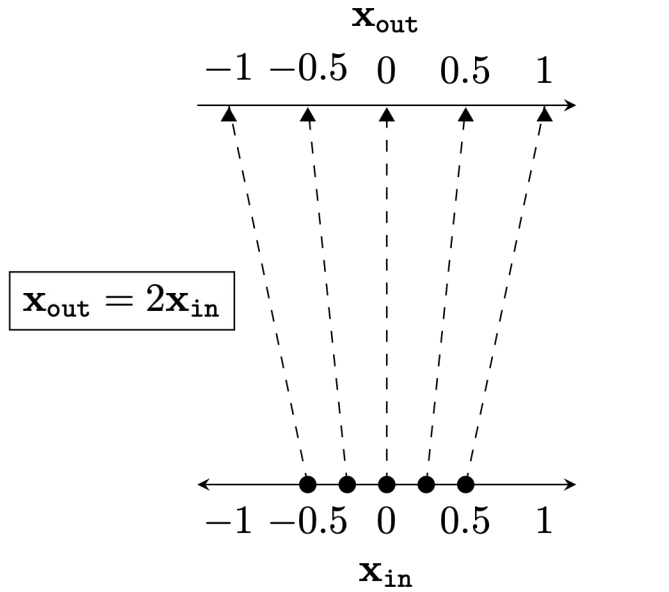- We call these representations **embeddings**

Layer L

Input

# Two different ways to represent a function

# Two different ways to represent a function

# Data transformations for a variety of neural net layers



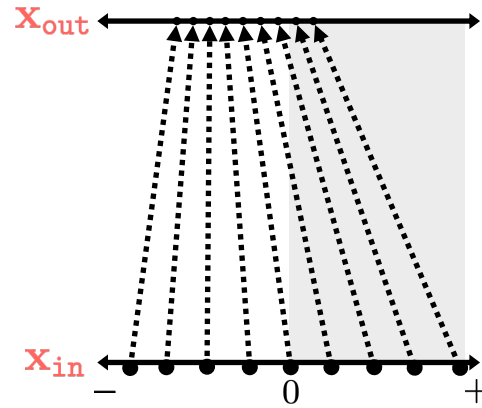$$\mathbf{x}_{\text{out}} = 2\mathbf{x}_{\text{in}}$$
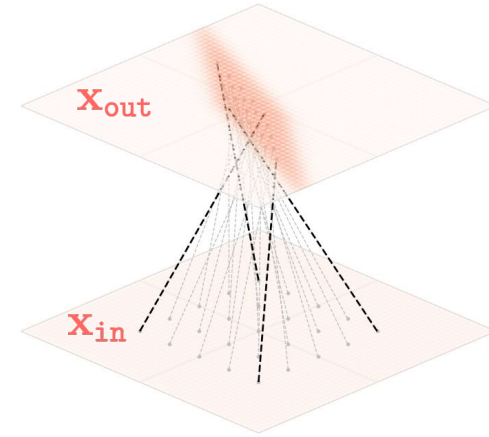
$$\mathbf{x}_{\text{out}} = \frac{\mathbf{x}_{\text{in}}+1}{2}$$

$$\mathbf{x}_{\text{out}} = \text{relu}(\mathbf{x}_{\text{in}})$$

$$\mathbf{x}_{\text{out}} = \text{sigmoid}(\mathbf{x}_{\text{in}})$$

Activations   Parameters

|               | Wiring graph | Equation | Mapping 1D | Mapping 2D |
|---------------|--------------|----------|------------|------------|

**linear**

$\mathbf{x}_{\text{in}}$  $\mathbf{W}$  $\mathbf{x}_{\text{out}}$  1  $\mathbf{b}$

$$\mathbf{x}_{\text{out}} = \mathbf{W}\mathbf{x}_{\text{in}} + \mathbf{b}$$

**relu**

$\mathbf{x}_{\text{in}}$  $\mathbf{x}_{\text{out}}$

$$x_{\text{out}_i} = \max(x_{\text{in}_i}, 0)$$

114

Activations    Parameters

| Wiring graph | Equation | Mapping | Matrix |

linear

$\mathbf{x_{in}}$      $\mathbf{x_{out}}$

$\mathbf{W}$

$\mathbf{b}$

1

$\mathbf{x_{out}} = \mathbf{W}\mathbf{x_{in}} + \mathbf{b}$

$\mathbf{x_{out}}$

$\mathbf{x_{in}}$

$-$    $0$    $+$

N+1

M    $\mathbf{W}$    $\mathbf{b}$

$\mathbf{x_{in}}$

1

$\rightarrow$    $\mathbf{x_{out}}$

relu

$\mathbf{x_{in}}$      $\mathbf{x_{out}}$

$x_{\mathbf{out}_i} = \max(x_{\mathbf{in}_i}, 0)$

$\mathbf{x_{out}}$

$\mathbf{x_{in}}$

$-$    $0$    $+$

logits

class probabilites

$x_1$ $z_{1_2}$ $h_{1_2}$ $z_{2_2}$ $y_2$
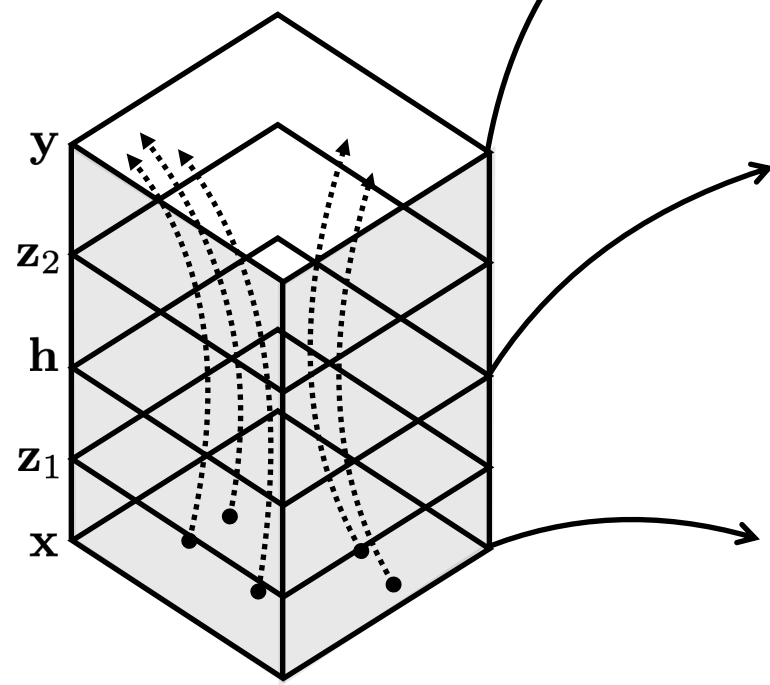
$x_2$ $z_{1_1}$ $h_{1_1}$ $z_{2_1}$ $y_1$

$\mathbf{W}_1$ relu $\mathbf{W}_2$ softmax

$\mathbf{y}$

$\mathbf{z}_2$

$\mathbf{h}$

$\mathbf{z}_1$

$\mathbf{x}$

$\mathbf{y}$

$\mathbf{h}$

$\mathbf{x}$

Training iteration

logits

class probabilites

$$W_1 \quad relu \quad W_2 \quad softmax$$

$$x_1 \rightarrow z_{1_2} \rightarrow h_{1_2} \rightarrow z_{2_2} \rightarrow y_2$$

$$x_2 \rightarrow z_{1_1} \rightarrow h_{1_1} \rightarrow z_{2_1} \rightarrow y_1$$

Training data

**x**

**y**

**h**

**z**$_1$

Training iteration

logits

class probabilites

$x_1$  $z_{1_2}$  $h_{1_2}$  $z_{2_2}$  $y_2$

$x_2$  $z_{1_1}$  $h_{1_1}$  $z_{2_1}$  $y_1$

$\mathbf{W}_1$  relu  $\mathbf{W}_2$  softmax

$\mathbf{y}$

$\mathbf{z}_2$

$\mathbf{h}$

$\mathbf{z}_1$

$\mathbf{x}$

Training iteration

Training iteration

softmax

linear

relu

linear

softmax

linear

relu

linear

softmax

linear

relu

linear

relu

linear

relu

linear

softmax

linear

relu

linear

relu

linear

relu

linear

softmax

linear

relu

linear

relu

linear

relu

linear

Training iteration

softmax

linear

relu

linear

linear

relu

linear

Layer 1 representation → Layer 6 representation

- structure, construction
- covering
- commodity, trade good, good
- conveyance, transport
- invertebrate
- bird
- hunting dog

[DeCAF, Donahue, Jia, et al. 2013]

[Visualization technique : t-sne, van der Maaten & Hinton, 2008]

# Optimization

# Gradient descent



$J(\theta)$

$\theta_1$

$\theta_2$

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

# 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

# Optimization

Params

$\theta$ → [ $J$ ] → $J(\theta)$
$\nabla_\theta J(\theta)$
$H_\theta(J(\theta))$

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

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

# 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?

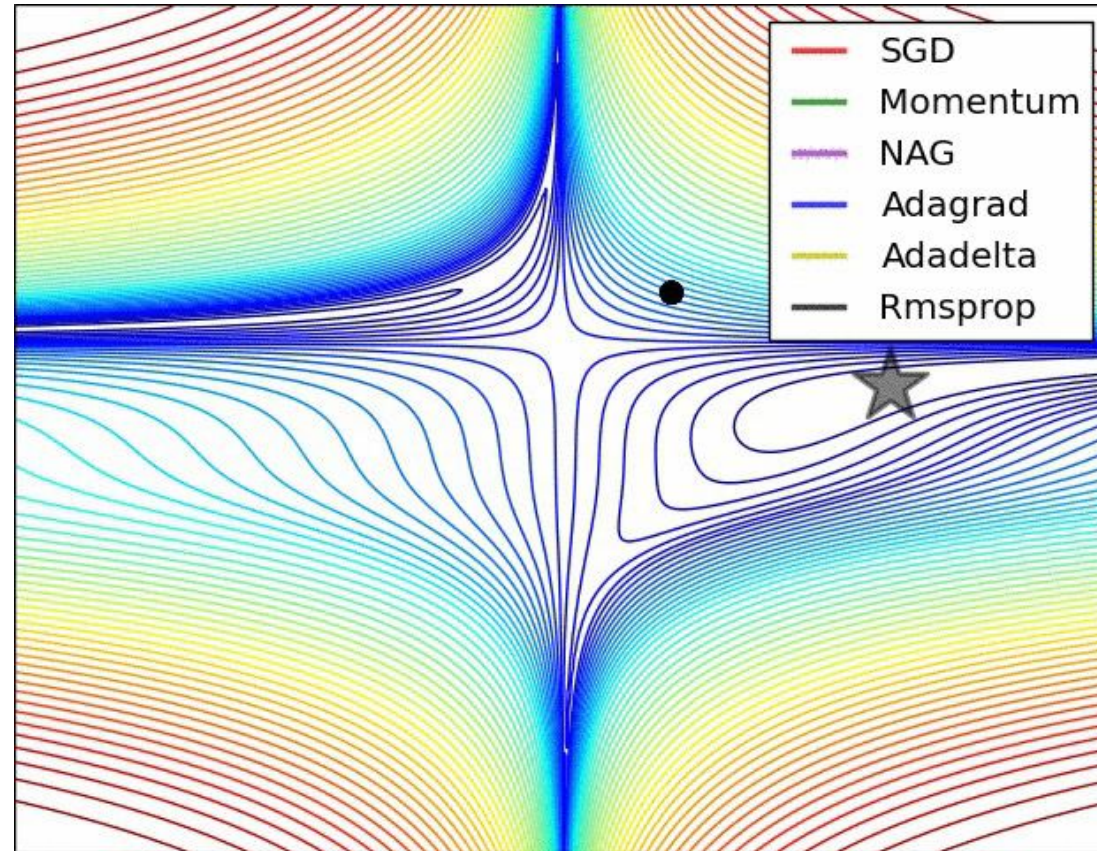GABRIEL GOH  |  April. 4  |  Citation:
UC Davis  |  2017  |  Goh, 2017

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

# Comparison of gradient descent variants



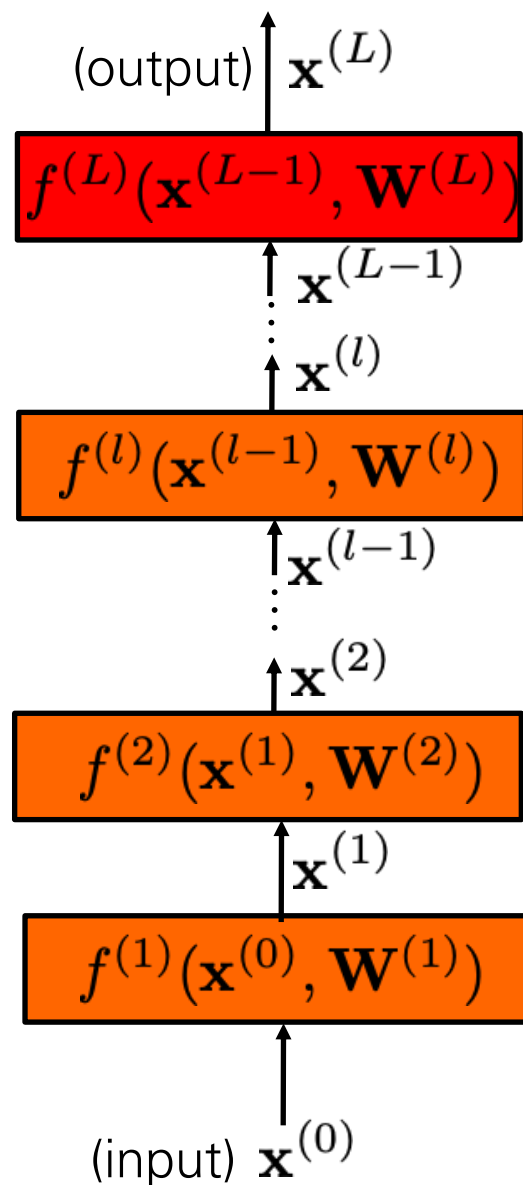[http://ruder.io/optimizing-gradient-descent/](http://ruder.io/optimizing-gradient-descent/)

# Backpropagation

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

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

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

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

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

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

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

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

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

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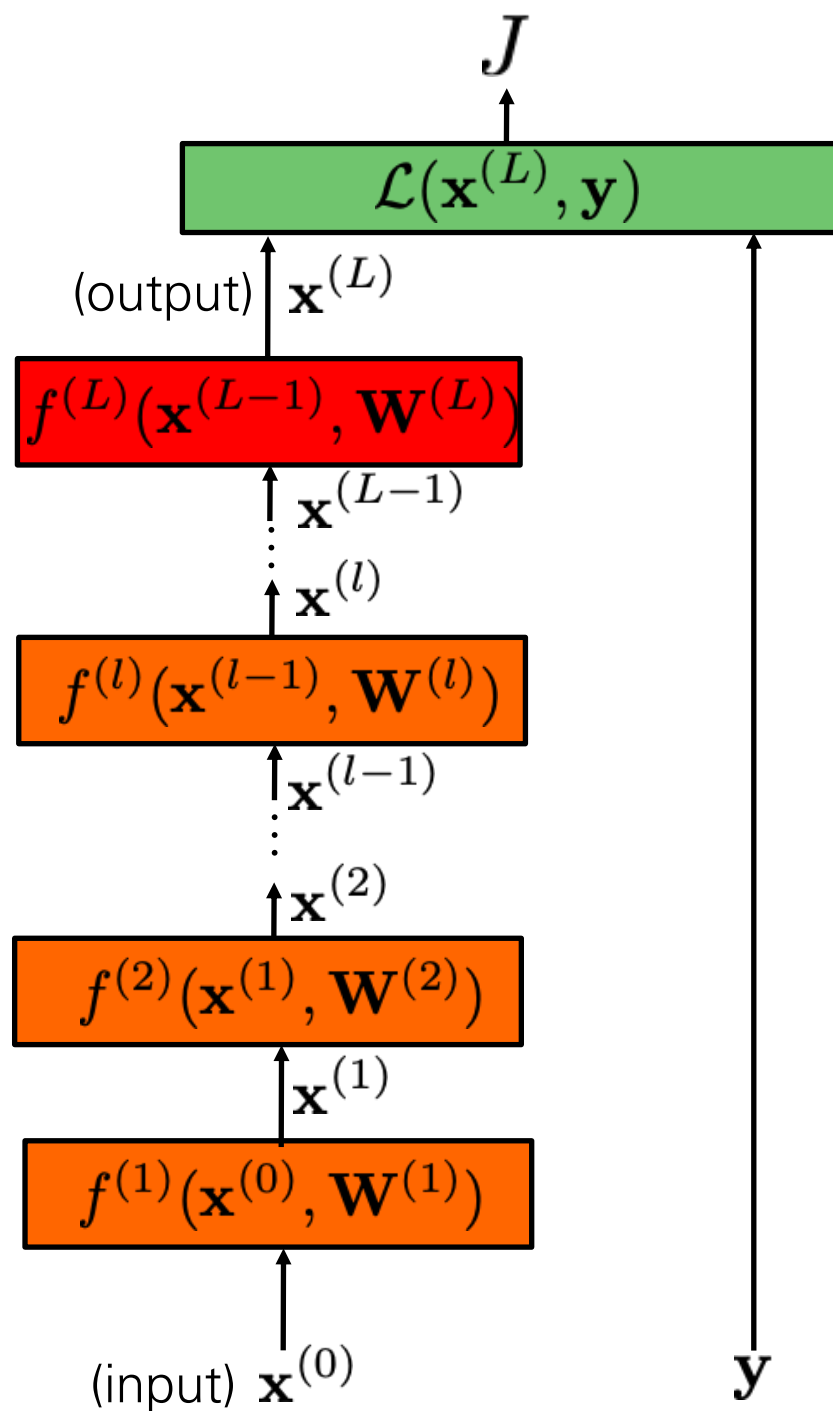
(input) $\mathbf{x}^{(0)}$

# 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:

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

$J$

$$\mathcal{L}(\mathbf{x}^{(L)}, \mathbf{y})$$

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

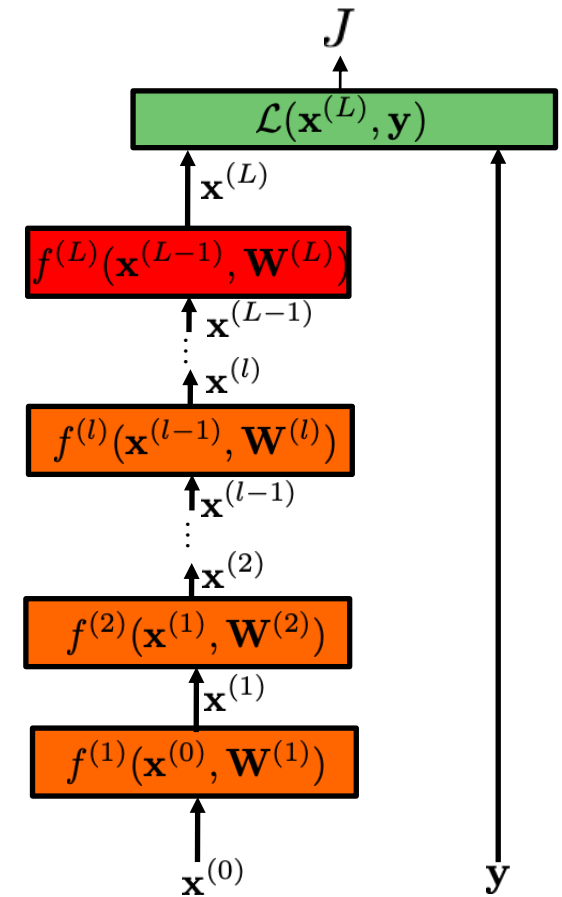$\mathbf{y}$

# 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.

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

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:

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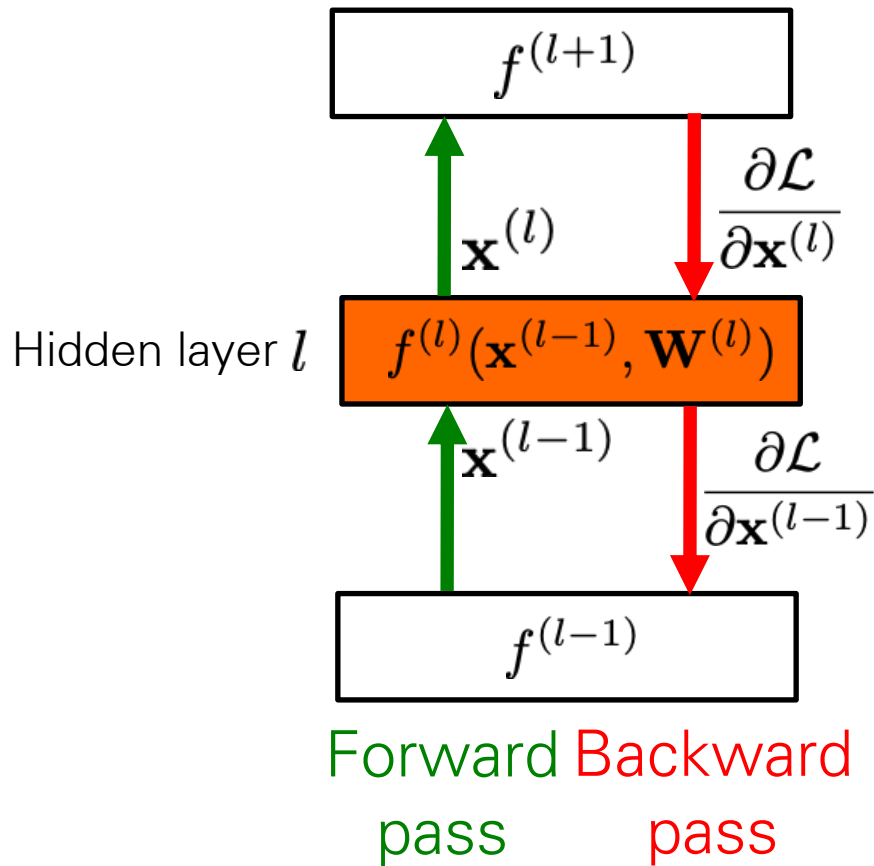)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial f^{(L)}(\mathbf{x}^{(L-1)}, \mathbf{W}^{(L)})}{\partial \mathbf{W}^{(L)}}$$

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:

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

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

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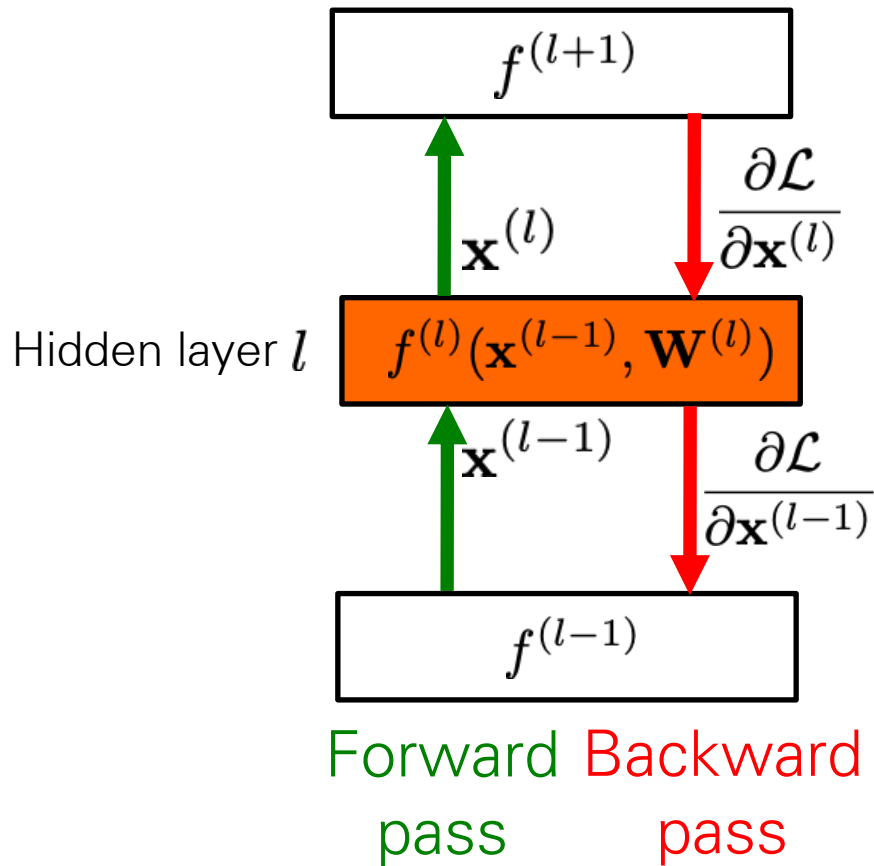frac{\partial \mathbf{x}^{(L)}}{\partial \mathbf{x}^{(L-1)}} \cdot \frac{\partial \mathbf{x}^{(L-1)}}{\partial \mathbf{x}^{(L-2)}} \cdots \frac{\partial \mathbf{x}^{(l+1)}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{W}^{(l)}}$$

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

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



Hidden layer $l$

Forward pass — Backward pass

— Goal: to update parameters of layer $l$

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

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

- 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:

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

# Backpropagation

Goal: to update parameters of layer $l$

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

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

- 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)

$$f^{(l+1)}$$

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

Hidden layer $l$

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

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

$$f^{(l-1)}$$

Forward pass    Backward pass

# 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



```
1  for i, data in enumerate(dataset):
2      iter_start_time = time.time()
3      if total_steps % opt.print_freq == 0:
4          t_data = iter_start_time - iter_data_time
5      visualizer.reset()
6      total_steps += opt.batch_size
7      epoch_iter += opt.batch_size
8      model.set_input(data)
9      model.optimize_parameters()
```

# Differentiable programming



[Figure from "Neural Module Networks", Andreas et al. 2017]

# 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"

What's the object class of the center pixel?

Training data

$\mathbf{x}$ $\quad y$

{ "Bird" },

{ "Bird" },

{ "Sky" }
⋮

$f$ → "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?

"Bird"

"Bird"

"Sky"

"Sky"

Translation invariance: process each patch in the same way.

An equivariant mapping:

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

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



**W**

**W**

**W**

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

# Convolution



filter

# Fully-connected network

Fully-connected (fc) layer



$\mathbf{W}$

$\mathbf{x}$

$1$

$\mathbf{b}$

$y \;\; g(y)$

# 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



$x$

$y \quad g(y)$

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.

# Linear system: $\mathbf{y} = f(\mathbf{x})$

A linear function f can be written as a matrix multiplication:



$$\mathbf{y} = \begin{bmatrix} & & \\ & h[n,k] & \\ & & \end{bmatrix} \mathbf{x}$$

n indexes rows,
k indexes columns

It can also be represented as a fully connected linear neural network



$\mathbf{x}$    $h[n,k]$    $\mathbf{y}$   ➡   $y[n] = \sum_{k=0}^{N-1} h[n,k]x[k]$

$h[n,k]$    Is the strength of the connection between x[k] and y[n]

# Convolution

A linear shift invariant (LSI) function f can be written as a matrix multiplication:



$h[n-k]$ n indexes rows, k indexes columns

It can also be represented as a convolutional layer of neural net:



$$y[n] = \sum_{k=-1}^{1} h[k]x[n-k]$$

$h[n-k]$ Is the strength of the connection between x[k] and y[n]

**Toeplitz matrix**

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

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

e.g., pixel image

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

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

$$\mathbf{x}^{(l+1)} = \qquad \text{[matrix]} \qquad * \quad \mathbf{x}^{(l)}$$

Conv layers can be applied to arbitrarily-sized inputs

# Five views on convolutional layers

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

2. Patch processing (Markov assumption)

3. Image filter

4. Parameter sharing

5. A way to process variable-sized tensors

# What if we have color?

## (aka multiple input channels?)

# Multiple channels

Conv layer

$$\mathbf{y} = \sum_c \mathbf{w}_c \circ \mathbf{x}_c$$

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

# Multiple channels

Conv layer



$$\mathbf{y}_k = \sum_c \mathbf{w}_{k_c} \circ \mathbf{x}_c$$

$$\mathbb{R}^{N \times C} \rightarrow \mathbb{R}^{N \times K}$$

# Multiple channels

Conv layer



$$\mathbf{y}_k = \sum_c \mathbf{w}_{k_c} \circ \mathbf{x}_c$$

$$\mathbb{R}^{N \times C} \longrightarrow \mathbb{R}^{N \times K}$$

Input features

A bank of 2 filters

2-dimensional output
**feature maps**

$F^1$

$F^2$

$\Sigma$

$\Sigma$

**x**

**y**

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

[Figure modified from Andrea Vedaldi]

# Feature maps



conv1

relu1

conv2

relu2

- Each layer can be thought of as a set of C **feature maps** aka **channels**
- Each feature map is an NxM image

# Multiple channels: Example

$\mathbf{x}_l$

128

128

3



Filter Bank with 3x3 filters

$\mathbf{x}_{(l+1)}$

128

128

96

How many parameters does each filter have?

(a) 9          (b) 27          (c) 96          (d) 864

# Multiple channels: Example

$\mathbf{x}_l$

128

128

3

Filter Bank with 3x3 filters

$\mathbf{x}_{(l+1)}$

128

128

96

How many filters are in the bank?

(a) 3    (b) 27    (c) 96    (d) can't say

# Filter sizes

When mapping from

$$\mathbf{x}_l \in \mathbb{R}^{H \times W \times C_l} \quad \longrightarrow \quad \mathbf{x}_{(l+1)} \in \mathbb{R}^{H \times W \times C_{(l+1)}}$$

using an filter of spatial extent $M \times N$

Number of parameters per filter: $M \times N \times C_l$

Number of filters: $C_{(l+1)}$

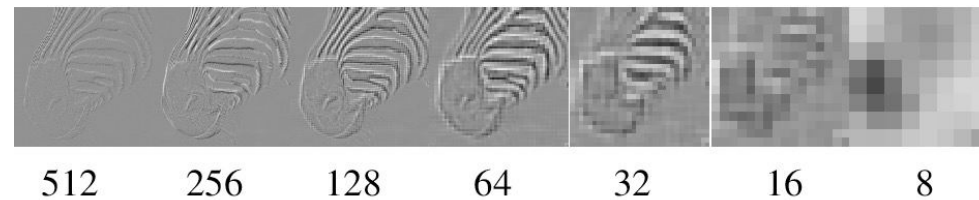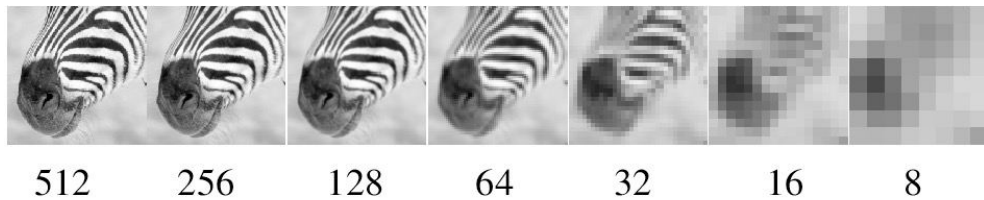# Pooling and downsampling

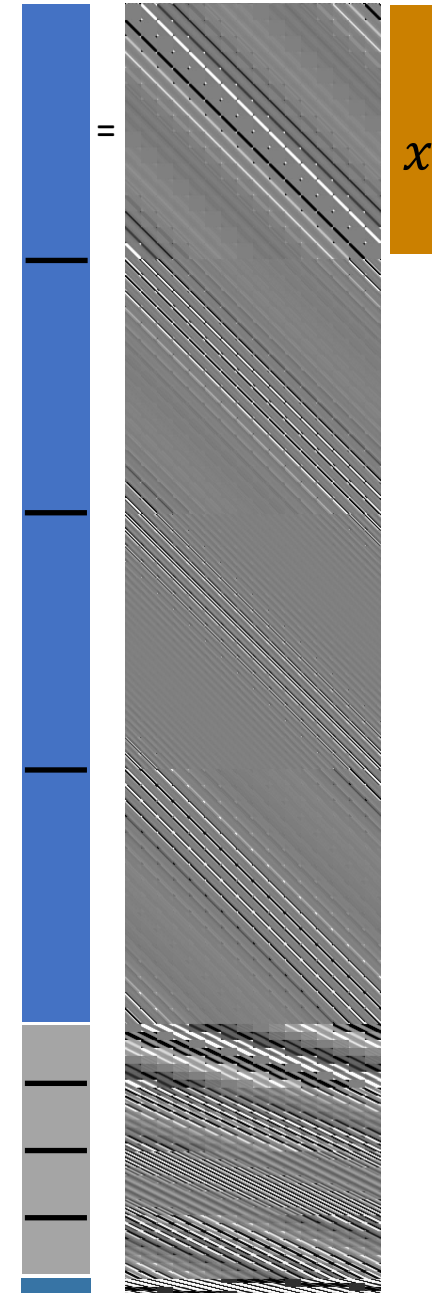We need translation and **scale** invariance

# Image pyramids

# Gaussian Pyramid

# Multiscale representations are great!



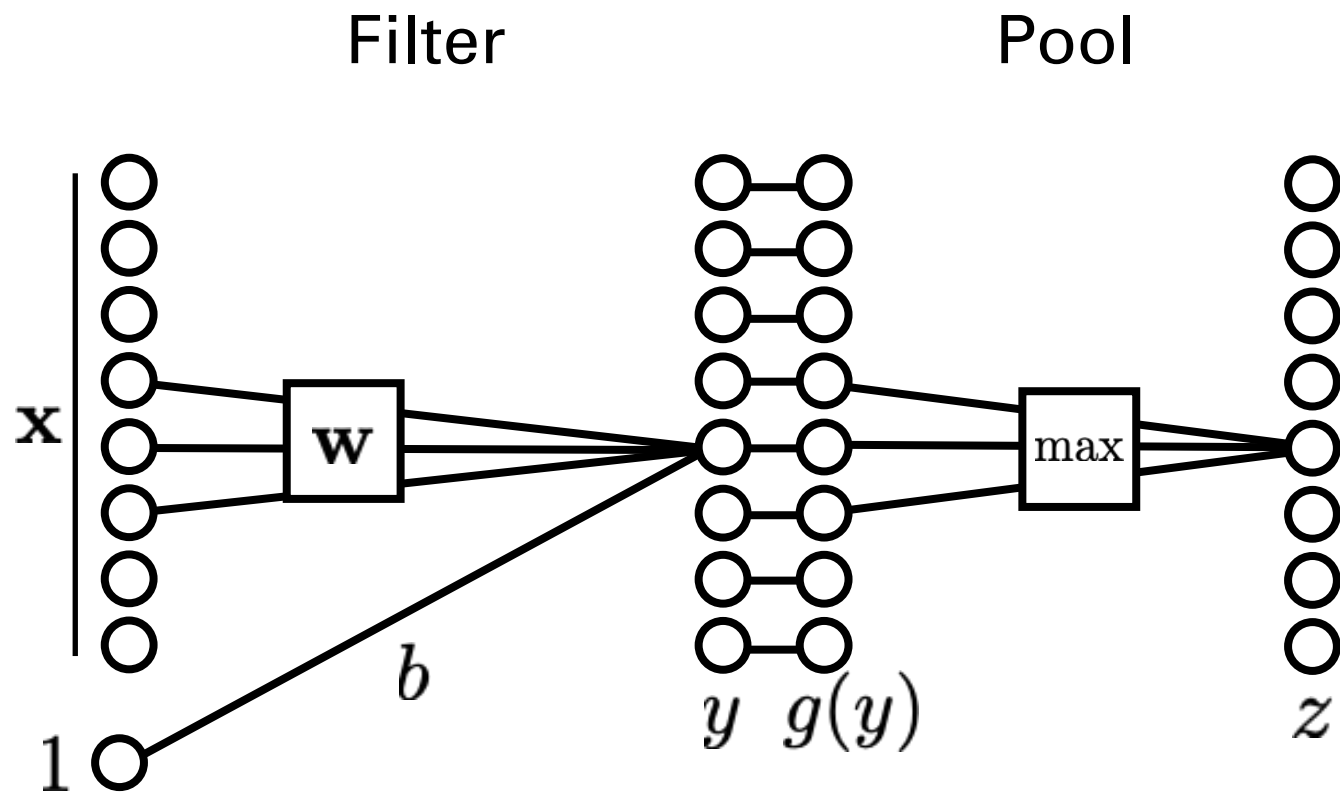512 256 128 64 32 16 8

Gaussian Pyr

512 256 128 64 32 16 8

Laplacian Pyr

How can we use multi-scale modeling in Convnets?

# Steerable Pyramid

# Pooling

Filter  Pool



Max pooling

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

# Pooling

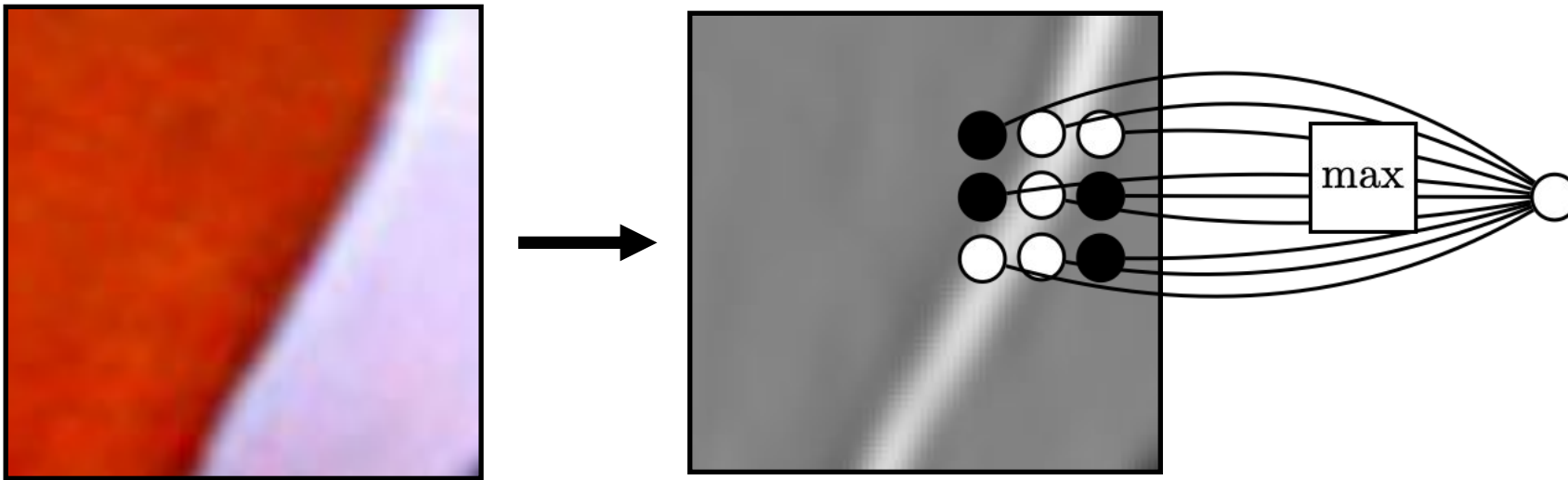Filter · · · · · · · · · · · · · · · Pool



Max pooling

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

Mean pooling

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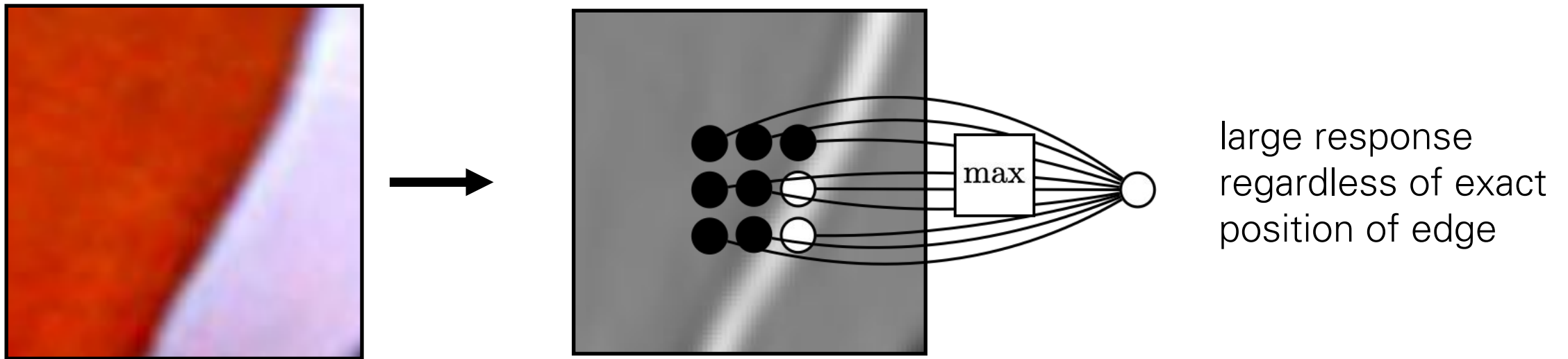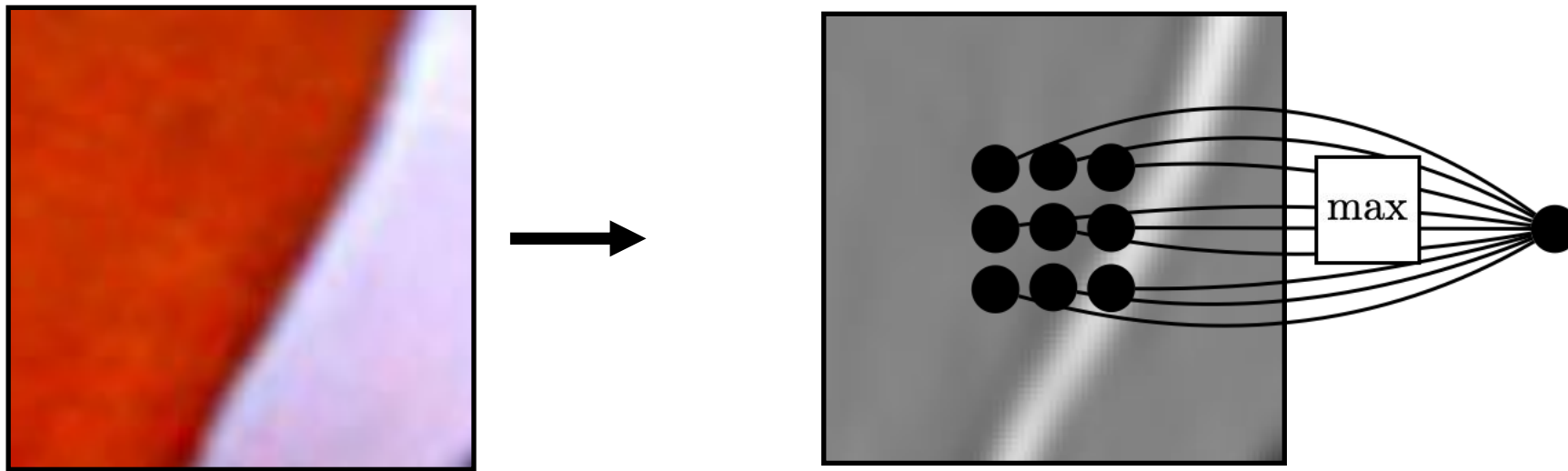
# Pooling – Why?

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

# Pooling – Why?

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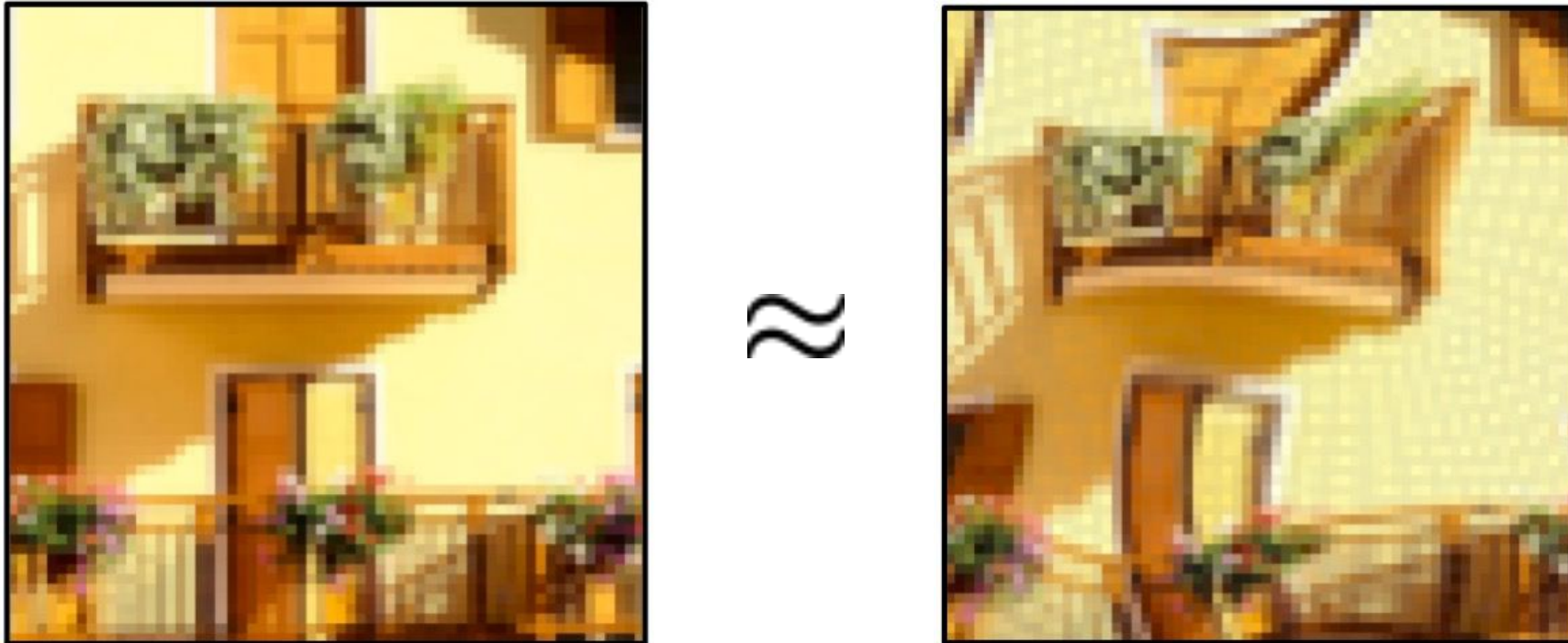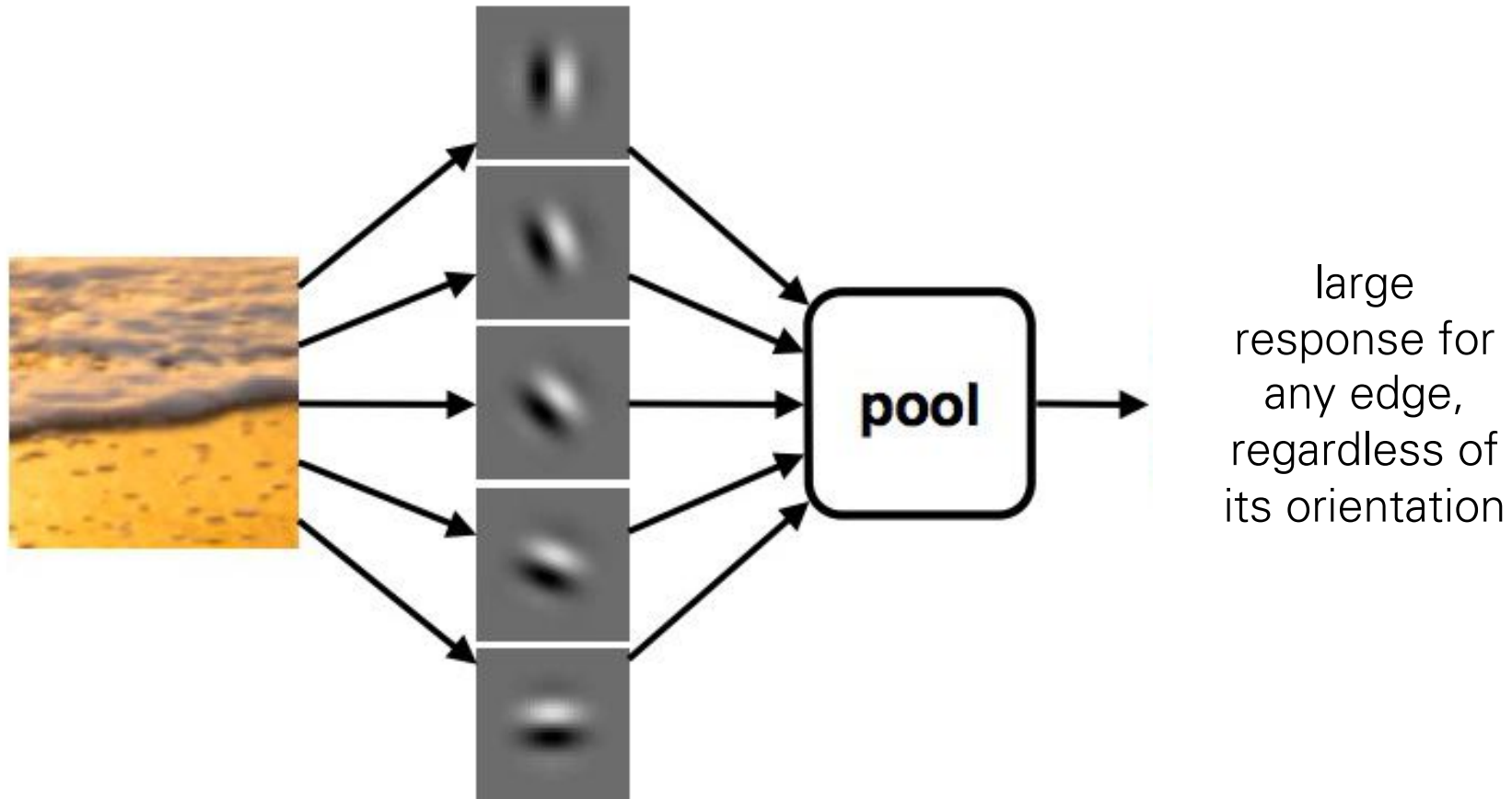stability w.r.t. small translations:



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



$\approx$

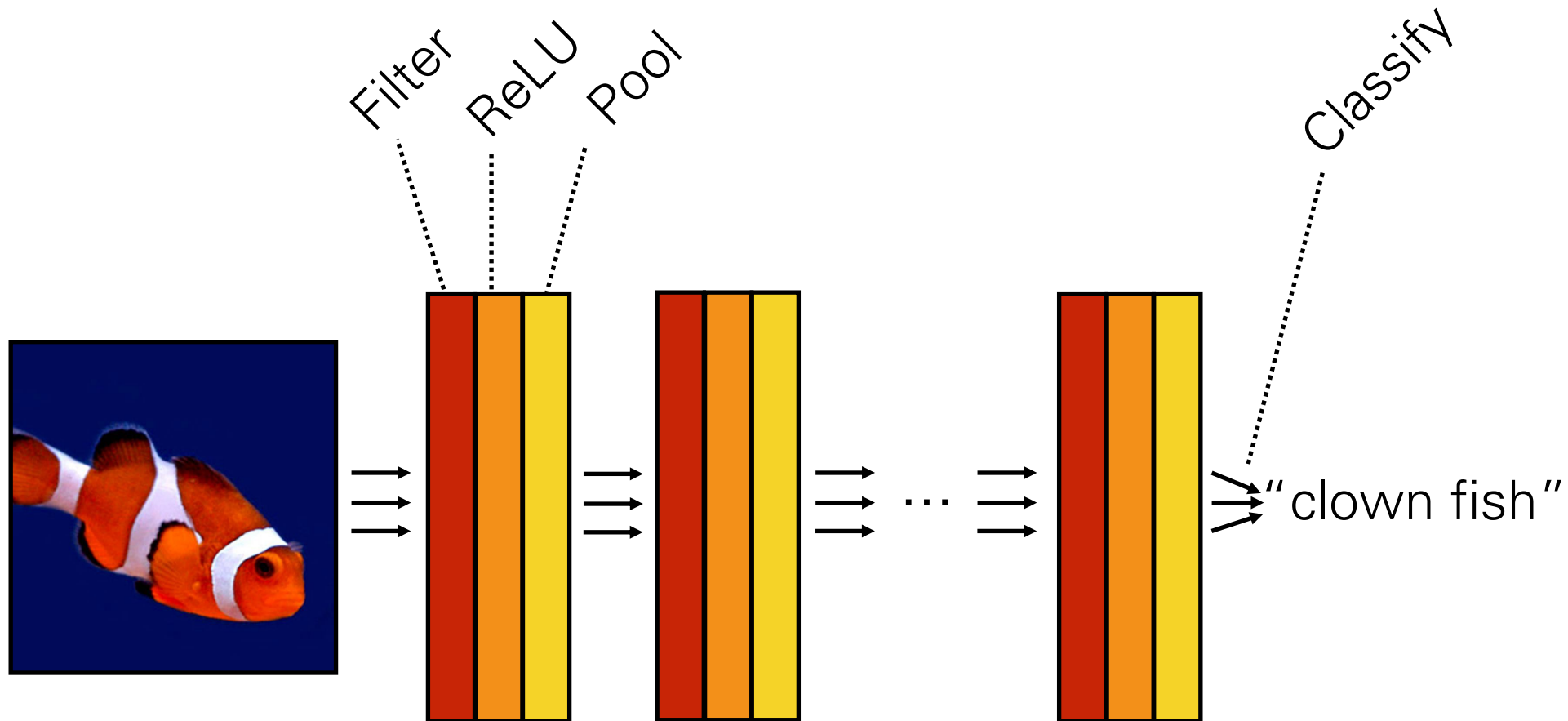["Unreasonable effectiveness of Deep Features as a Perceptual Metric", Zhang et al. 201

# Pooling – Why?

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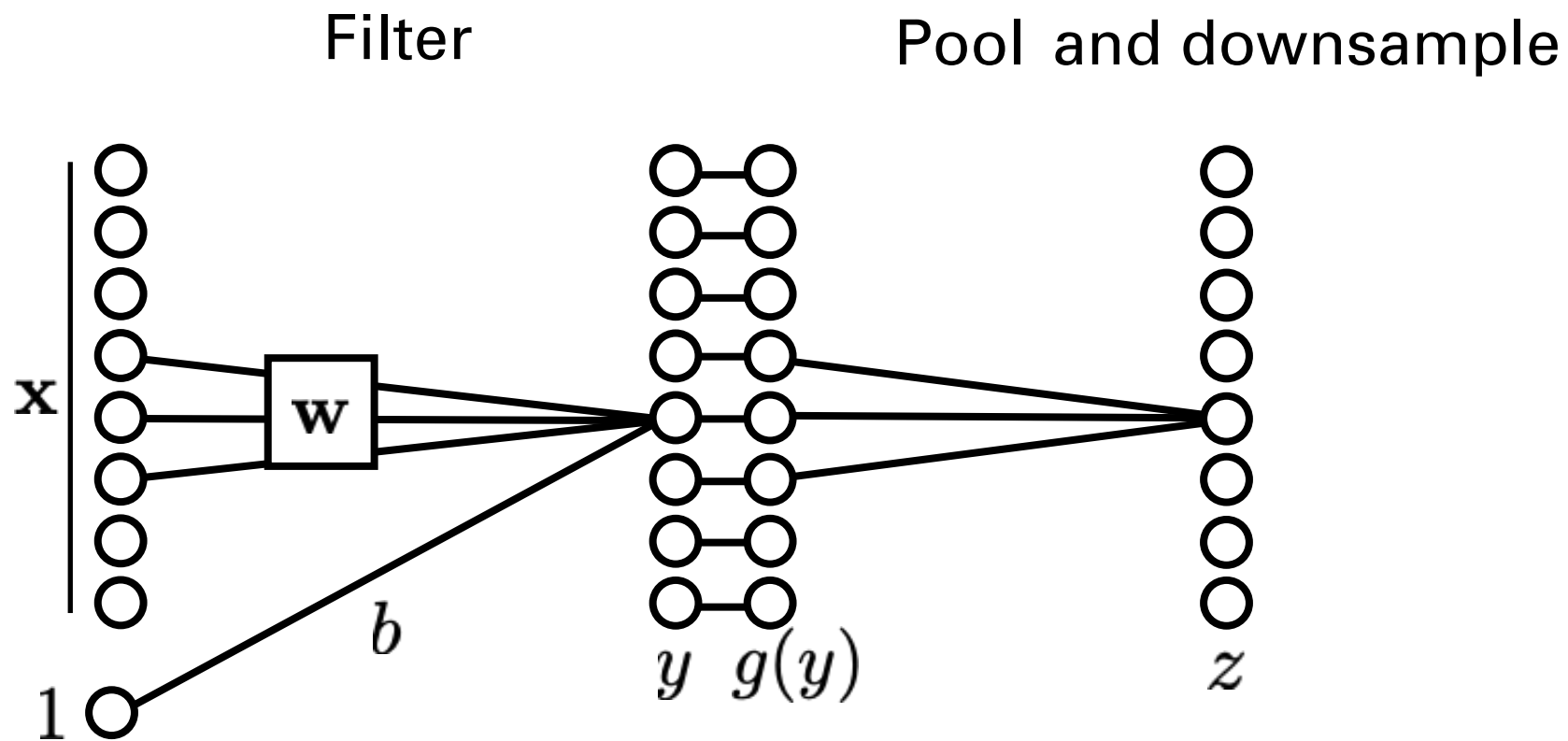can achieve other kinds of invariances:



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

$\mathbf{x}$

$\mathbf{w}$

$b$

$1$

$y \quad g(y)$

$z$

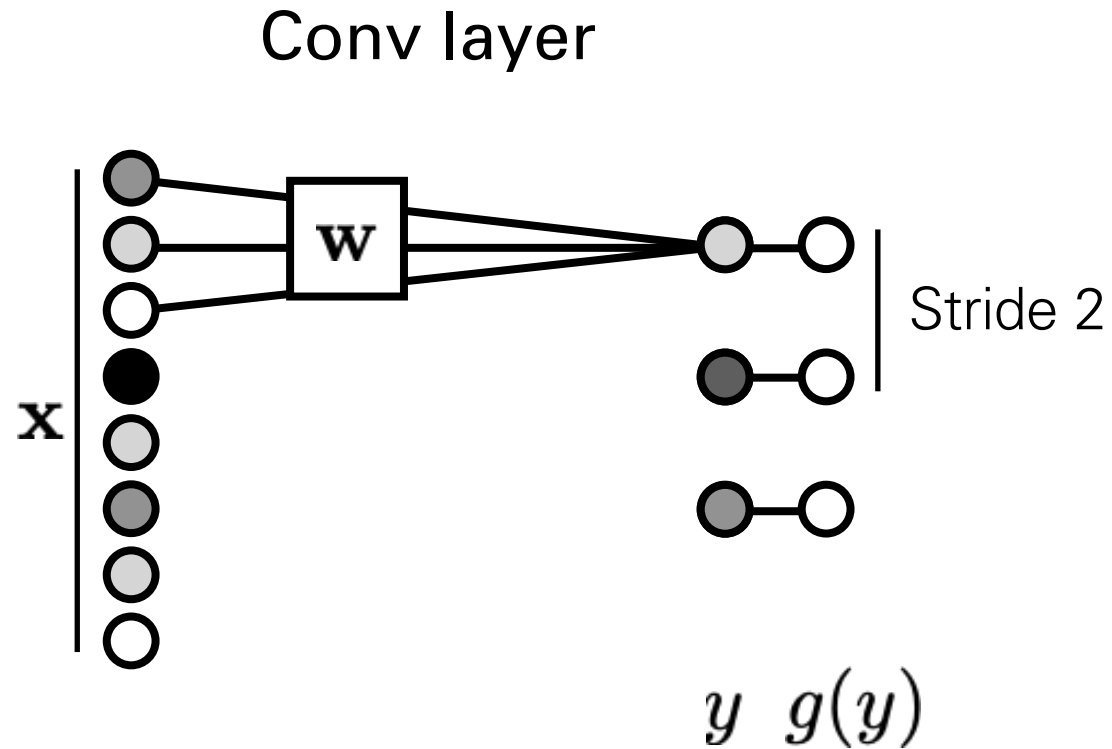# Downsampling
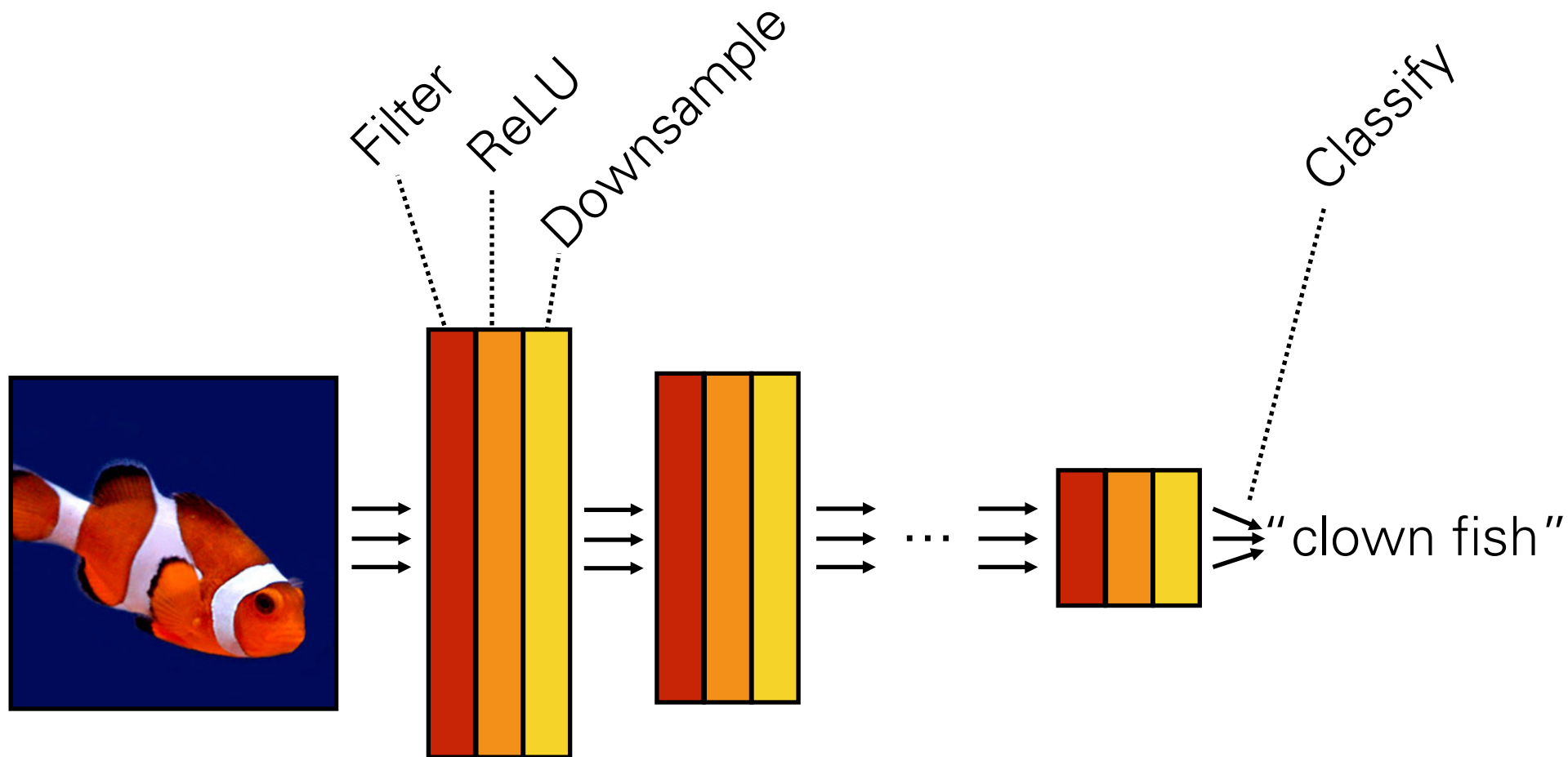
Filter                                    Downsample



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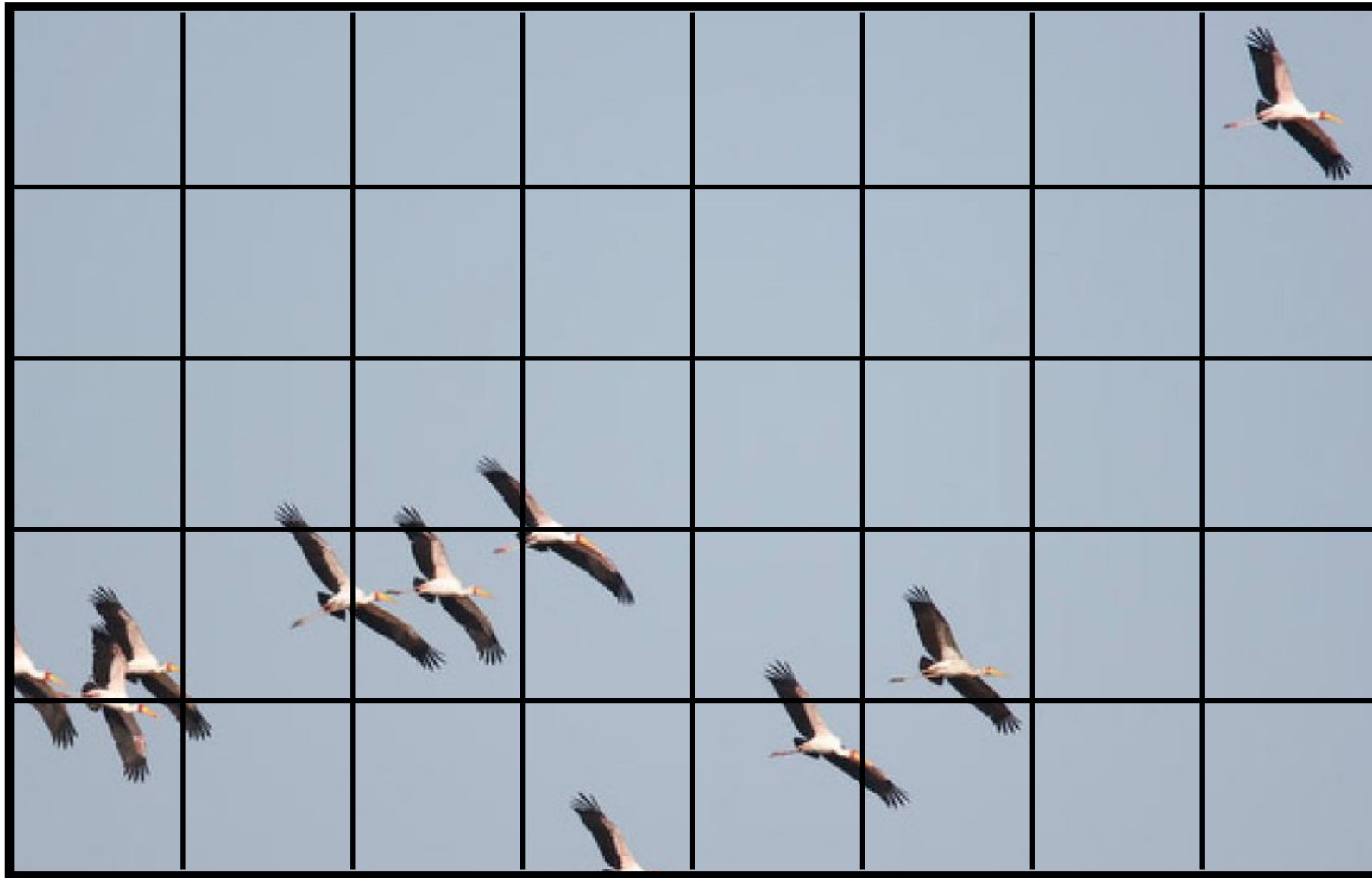
# Strided convolution

Conv layer



$x$

Stride 2

$y \;\; g(y)$

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

# Computation in a neural net

Filter

ReLU

Downsample

Classify

"clown fish"

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

# 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]
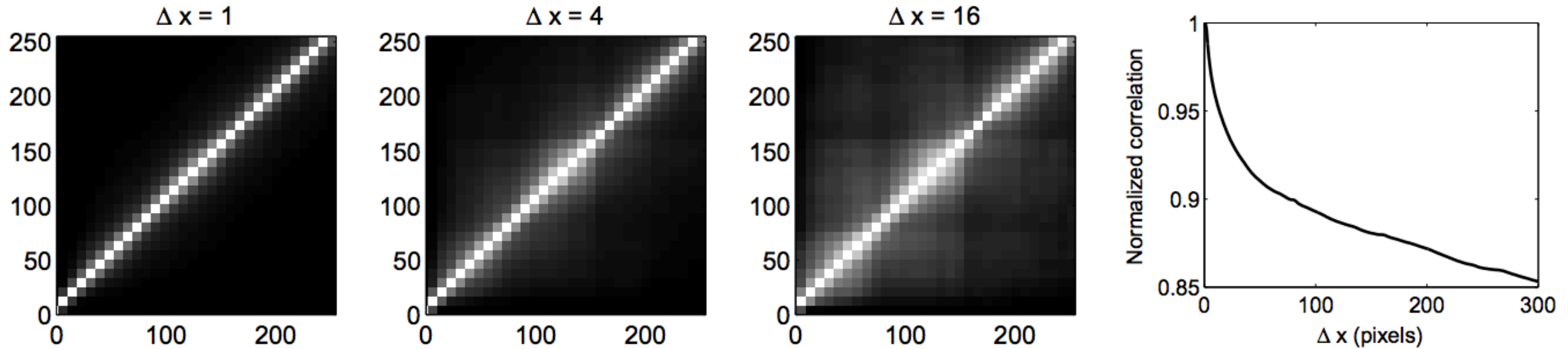
# CNNs – Why?



**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]

# CNNs – Why?

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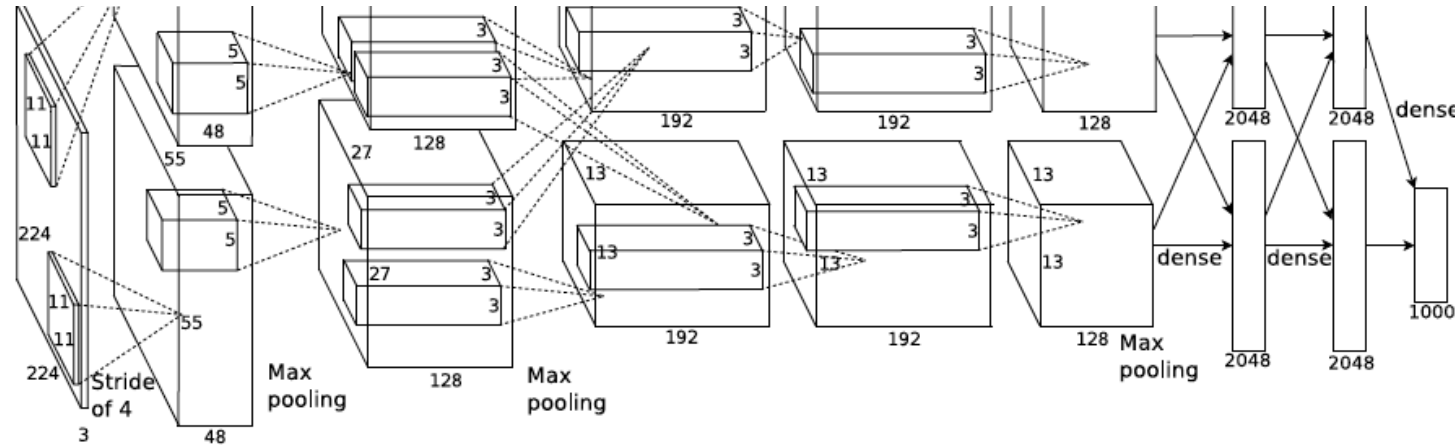
[For more discussion, see "Why does Deep and Cheap Learning Work So Well?", Lin et al. 2017]

# CNNs – Why?

Capturing long-range dependences:

# Alexnet — [Krizhevsky et al. NIPS 2012]



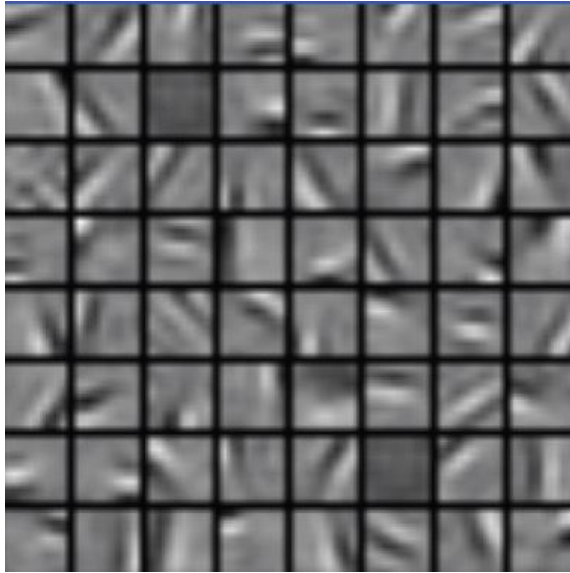| Layer |
|---|
| FULL CONNECT |
| FULL 4096/ReLU |
| FULL 4096/ReLU |
| MAX POOLING |
| CONV 3x3/ReLU 256fm |
| CONV 3x3ReLU 384fm |
| CONV 3x3/ReLU 384fm |
| MAX POOLING 2x2sub |
| LOCAL CONTRAST NORM |
| CONV 11x11/ReLU 256fm |
| MAX POOL 2x2sub |
| LOCAL CONTRAST NORM |
| CONV 11x11/ReLU 96fm |

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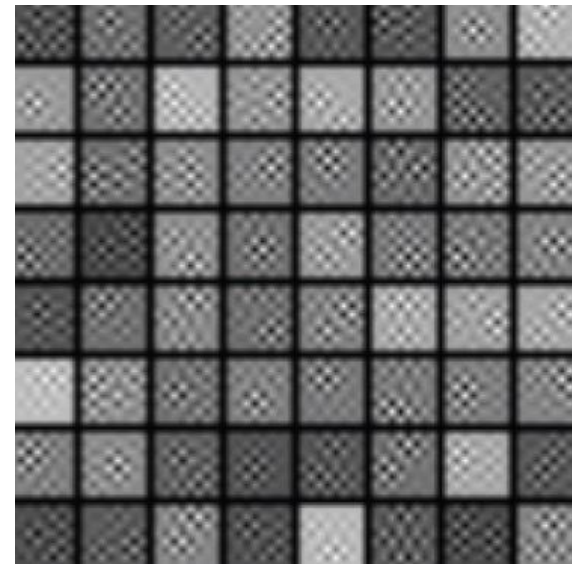[1000] FC8: 1000 neurons (class scores)

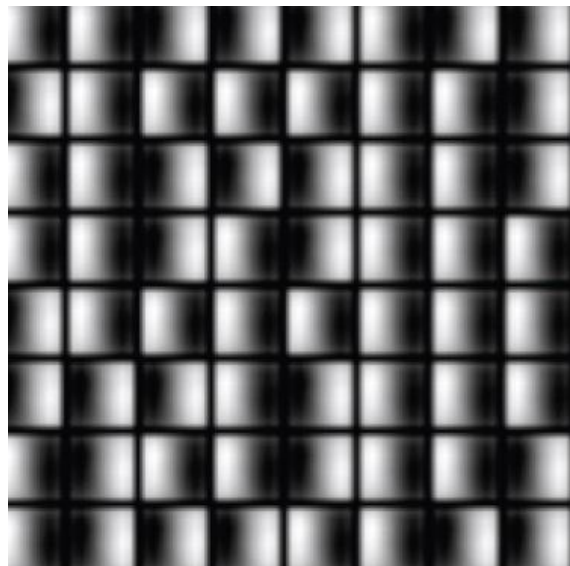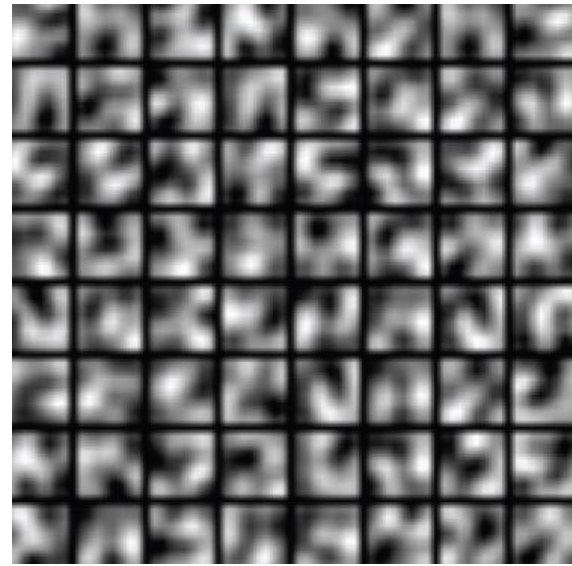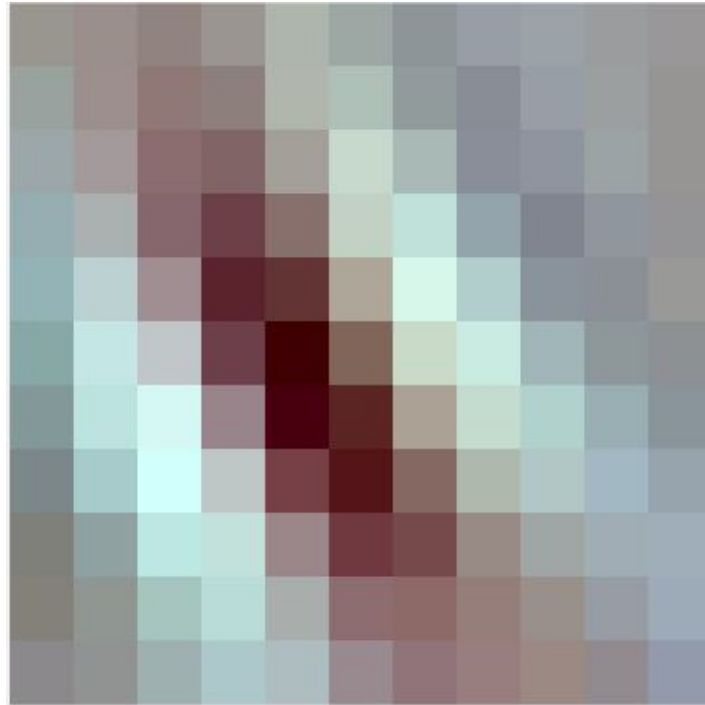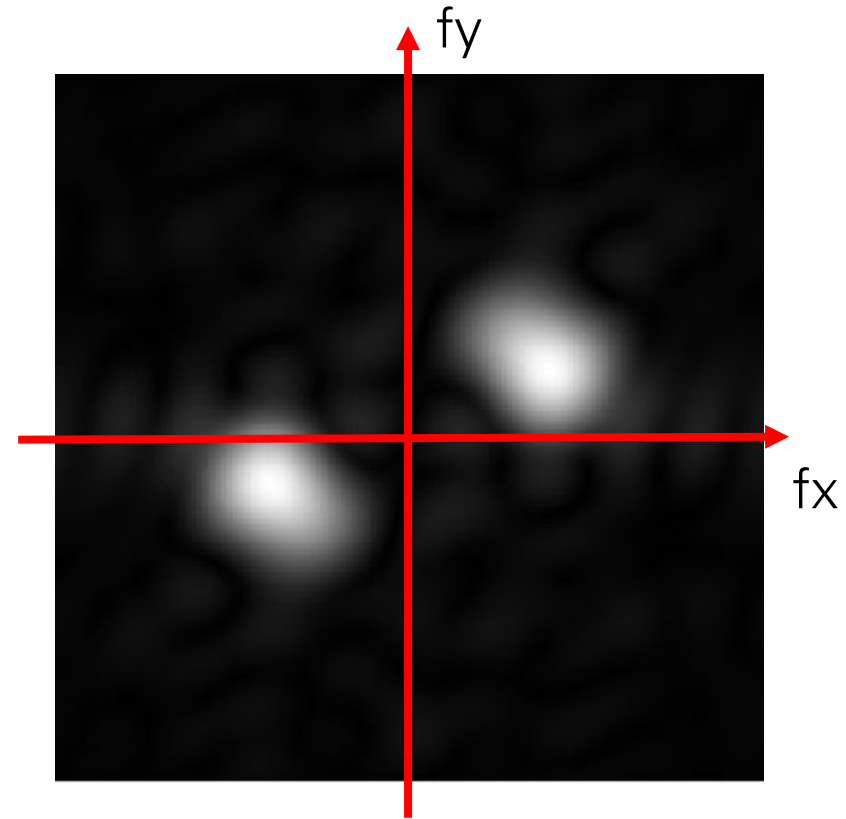# What filters are learned?
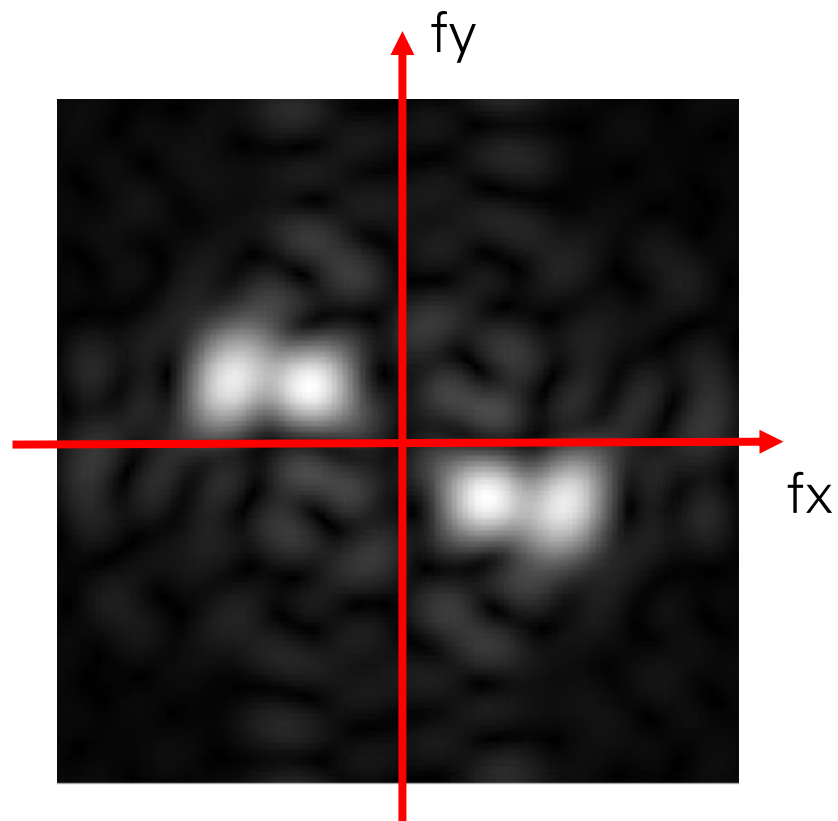
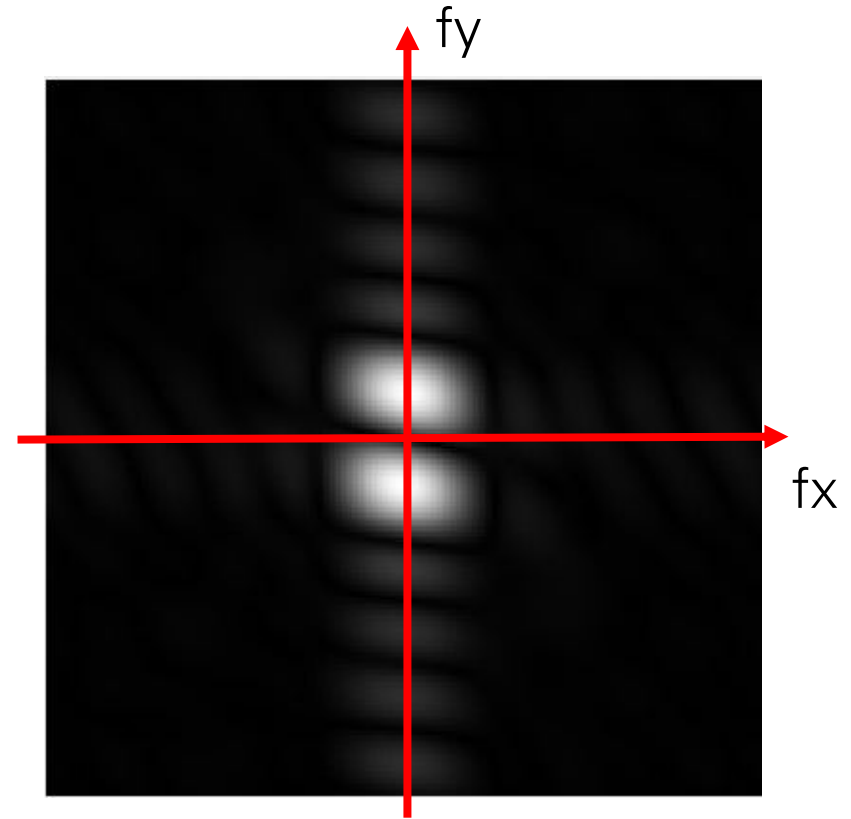# What filters are learned?

A



B



C



D

# Get to know your units
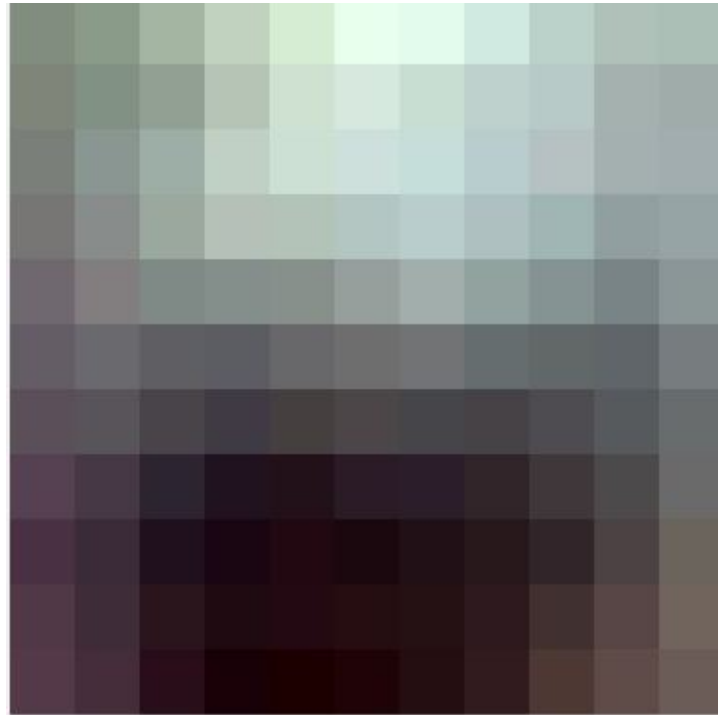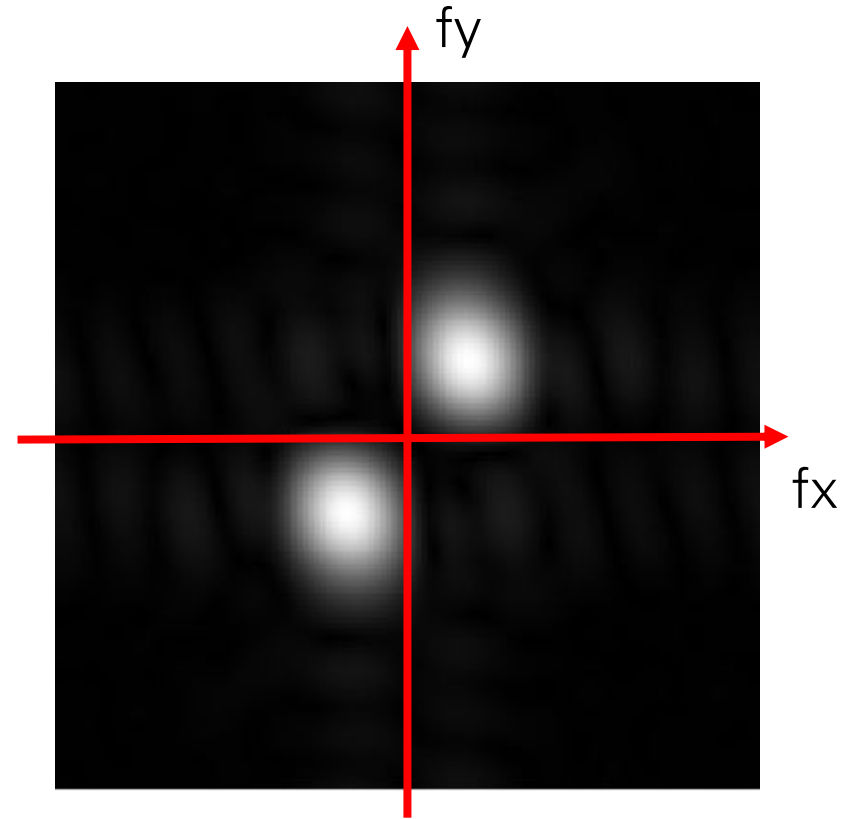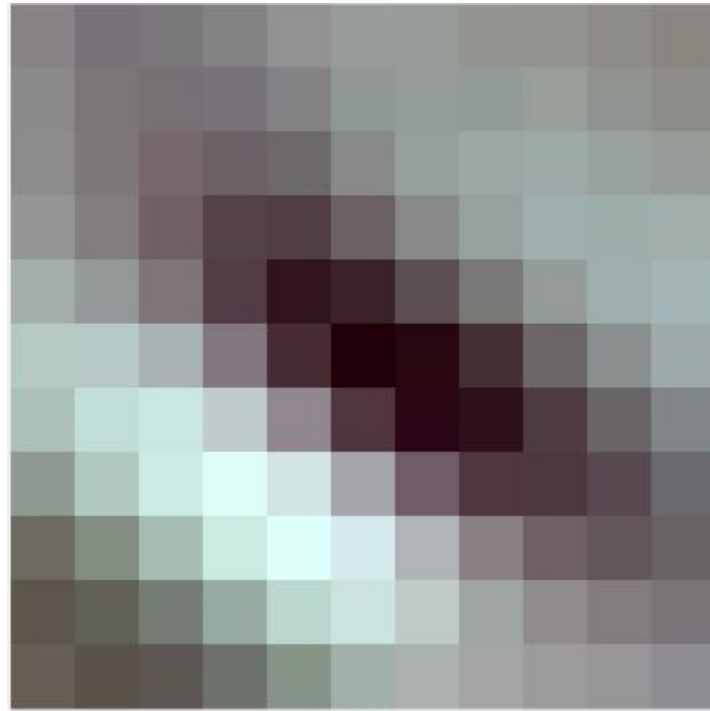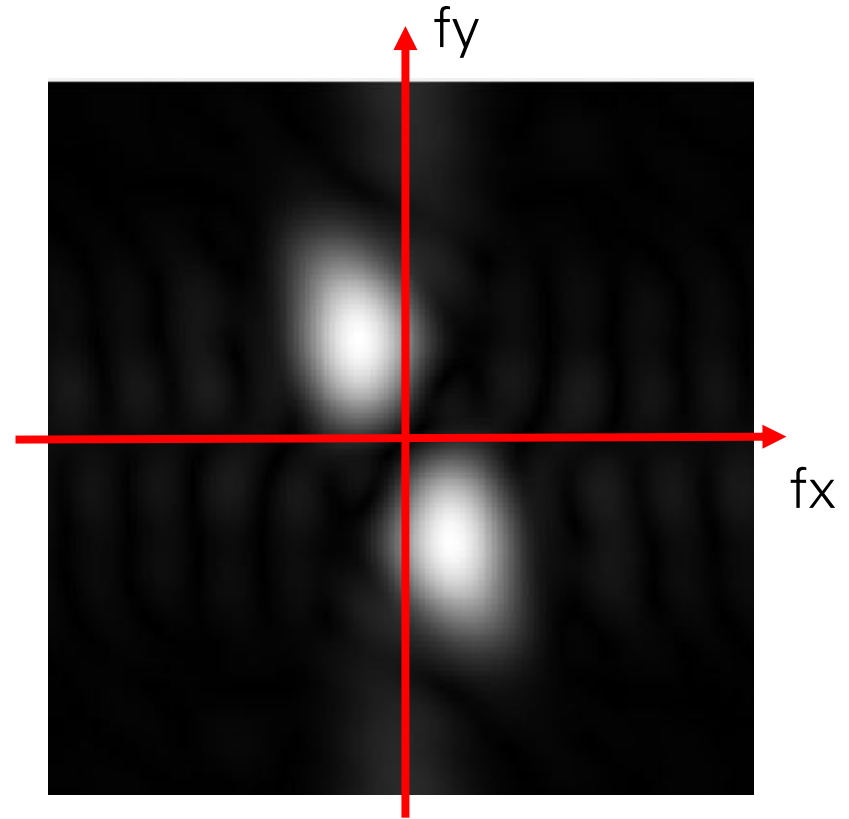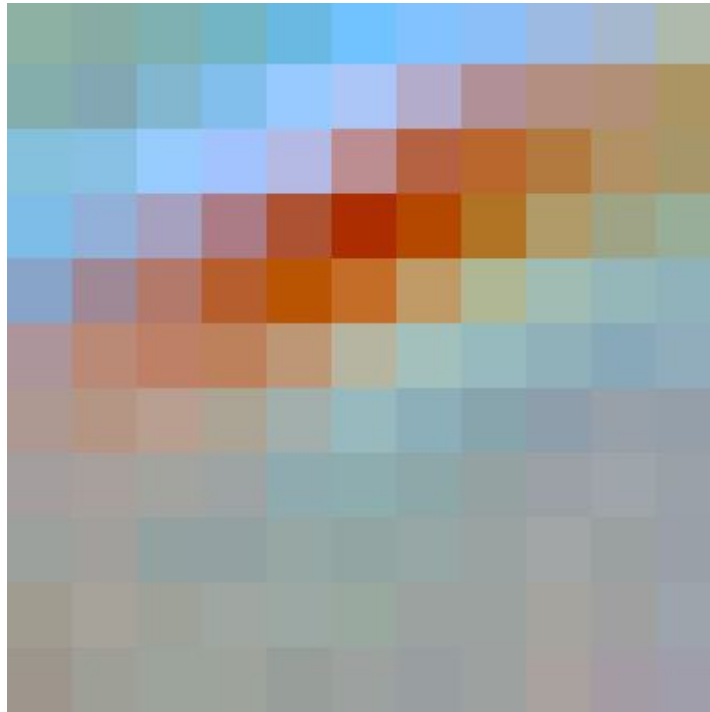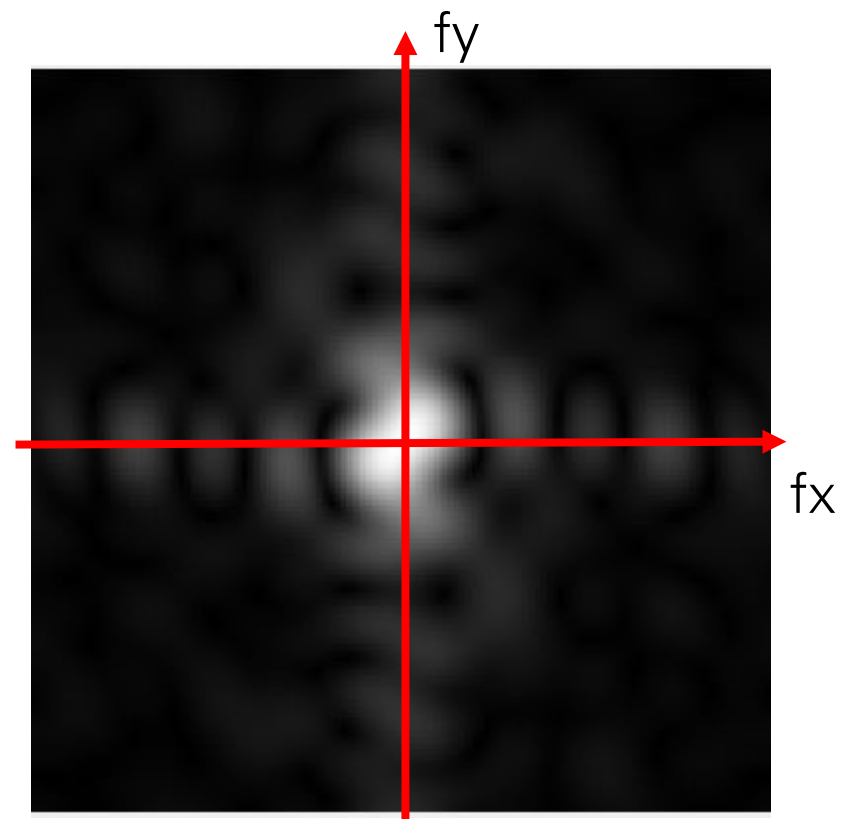


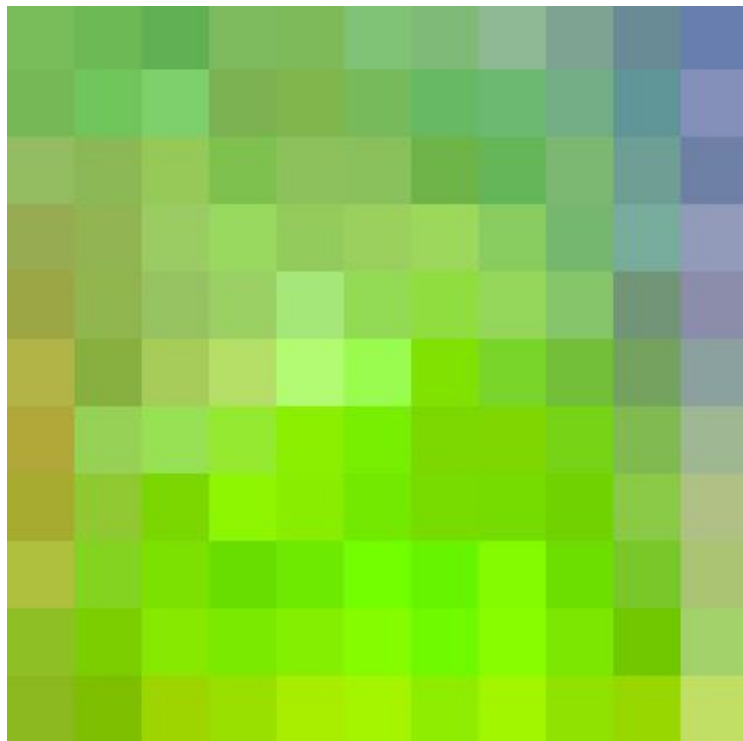11x11 convolution kernel
(3 color channels)

# Get to know your units

# Get to know your units
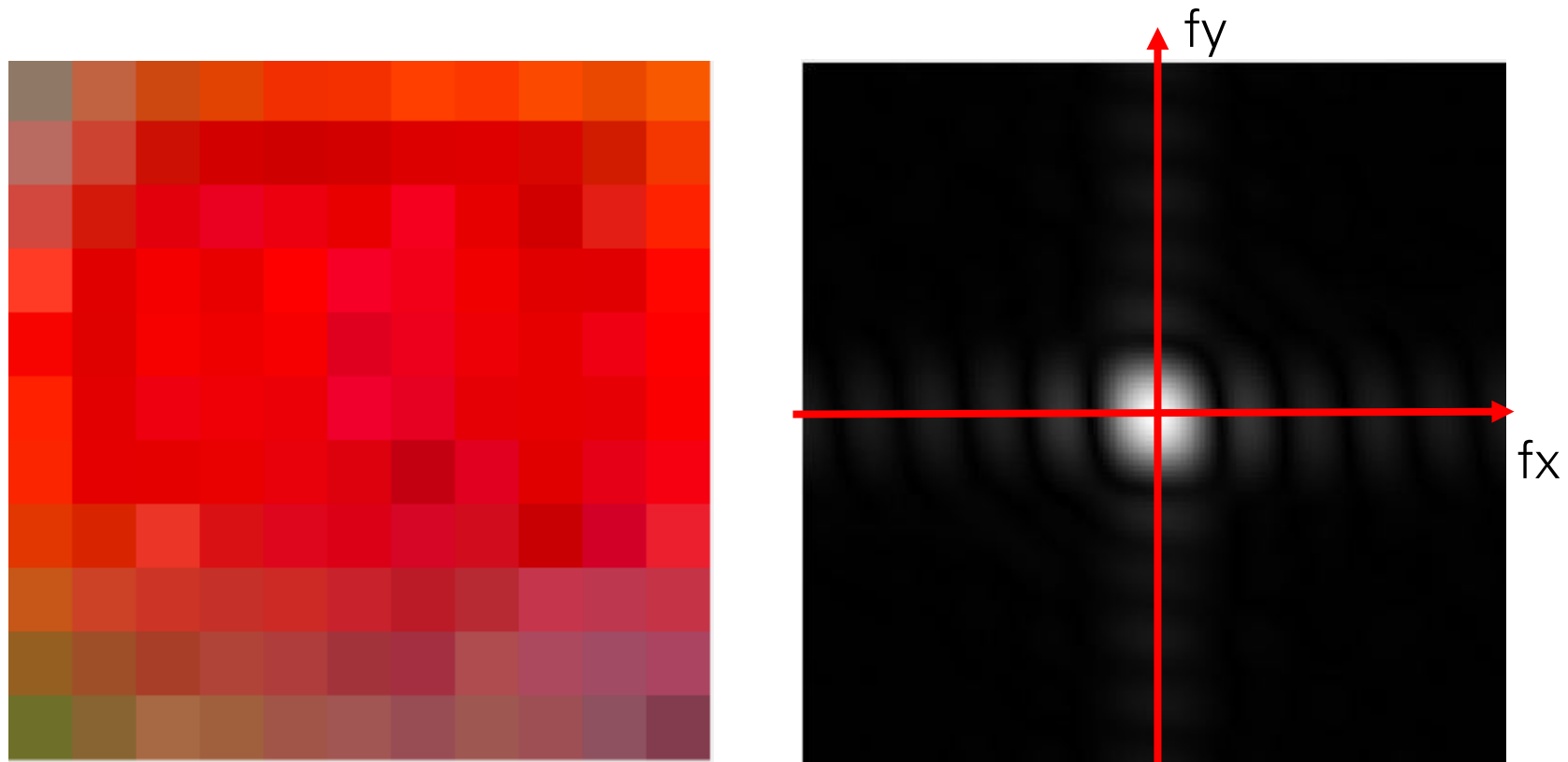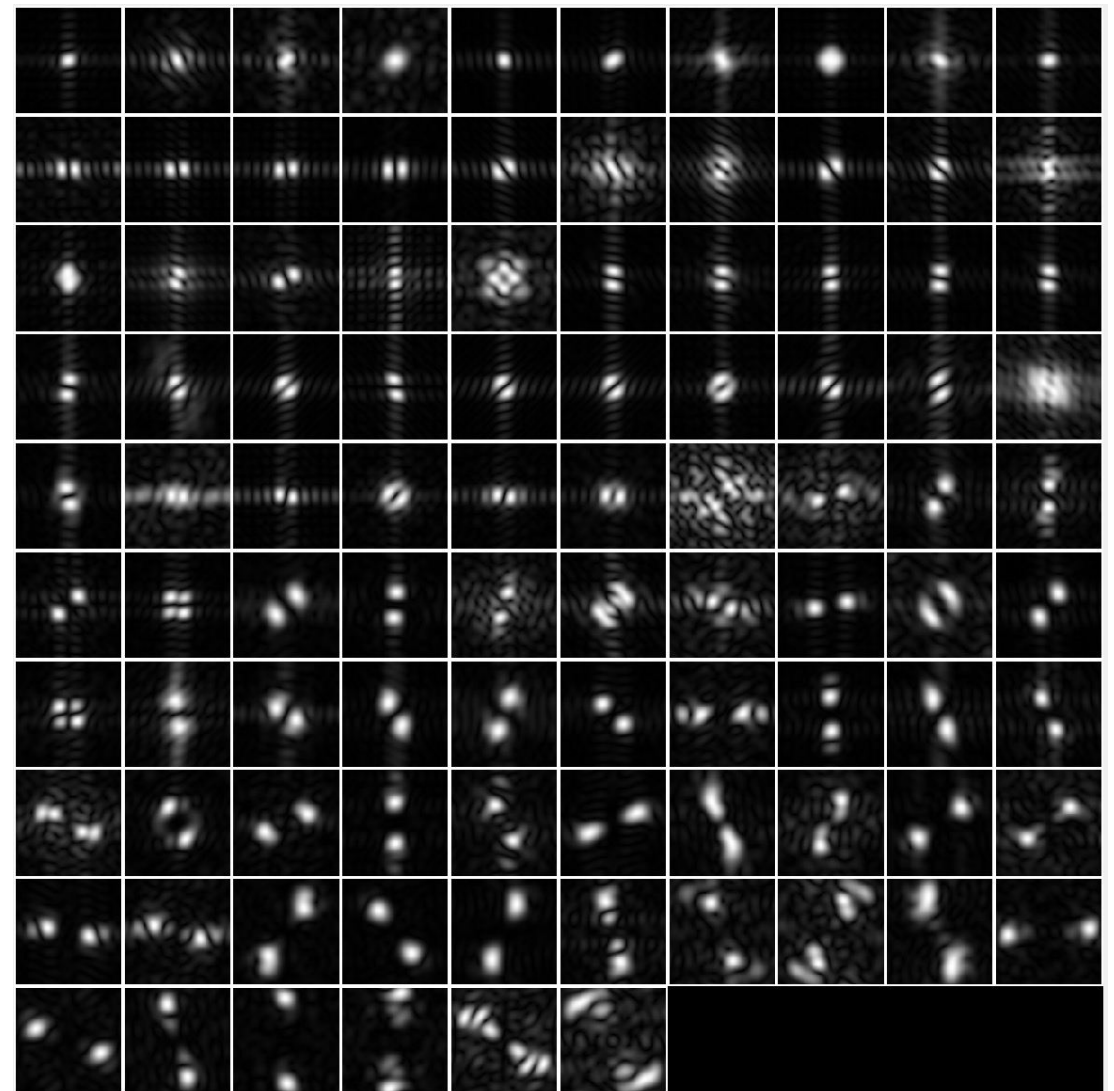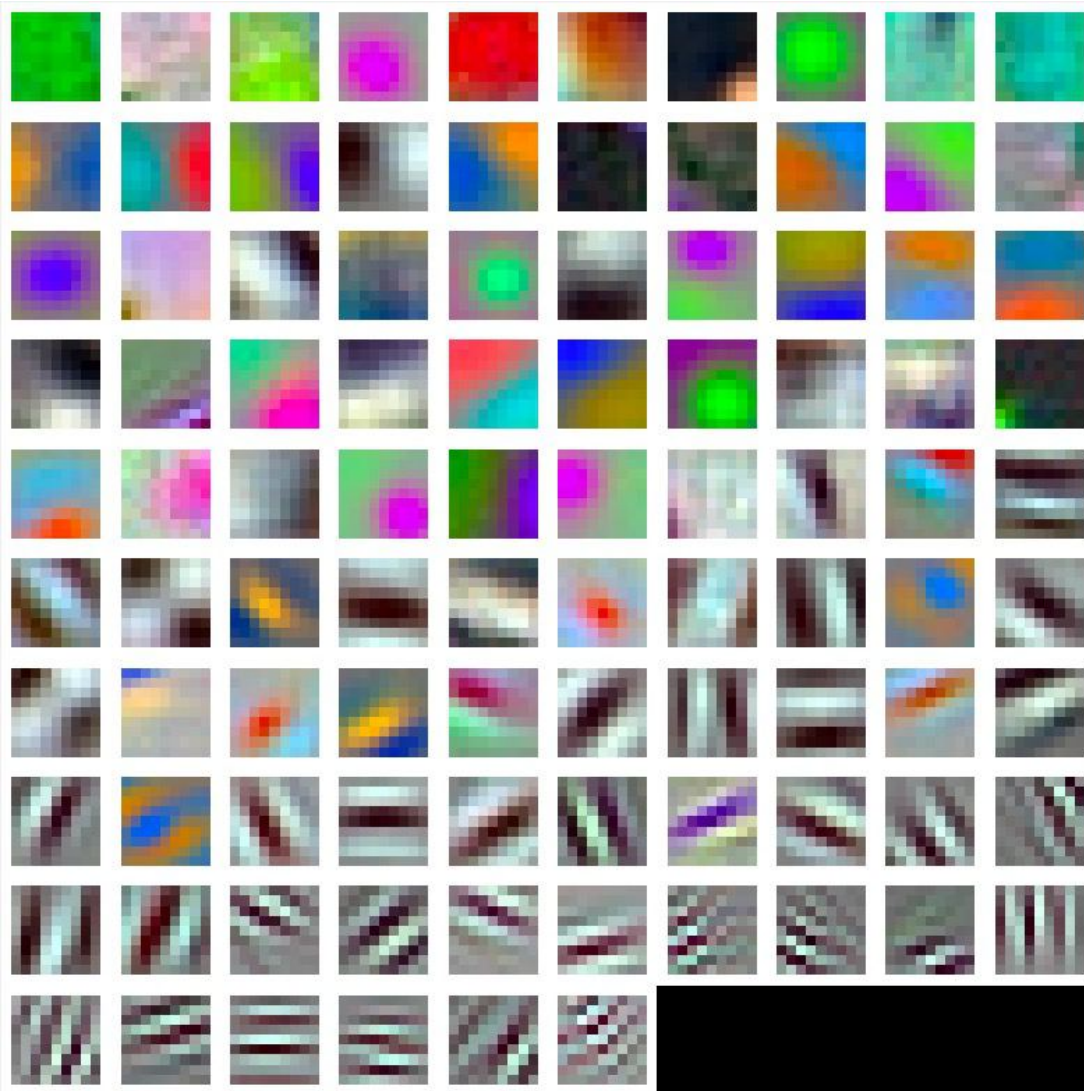
# Get to know your units
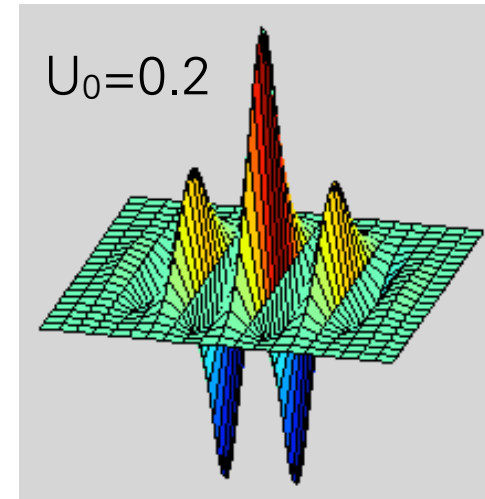
# Get to know your units

# Get to know your units

# Get to know your units

# Get to know your units



96 Units in conv1

224

# Gabor wavelets

$$\psi_c(x,y) = e^{-\frac{x^2+y^2}{2\sigma^2}} \cos(2\pi u_0 x)$$



u₀=0

U₀=0.1

U₀=0.2

$$\psi_s(x,y) = e^{-\frac{x^2+y^2}{2\sigma^2}} \sin(2\pi u_0 x)$$

b)

# Comparing Human and Machine Perception



FIGURE 1   Schematic overview of the processing done by the early visual system. On the left, are some of the major structures to be discussed; in the middle, are some of the major operations done at the associated structure; in the right, are the 2-D Fourier representations of the world, retinal image, and sensitivities typical of a ganglion and cortical cell.

John Daugman, 1988



**2D Receptive Field**

**2D Gabor Function**

**Difference**

Fig. 5. Top row: illustrations of empirical 2-D receptive field profiles measured by J. P. Jones and L. A. Palmer (personal communication) in simple cells of the cat visual cortex. Middle row: best-fitting 2-D Gabor elementary function for each neuron, described by (10). Bottom row: residual error of the fit, indistinguishable from random error in the Chi-squared sense for 97 percent of the cells studied.

228

# Deep Neural Networks for Visual Recognition



2012: AlexNet
5 conv. layers

Error: 15.3%

2014: VGG
16 conv. layers

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.

VGG

# Chaining convolutions

3x3　　　　　3x5　　　　　　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

7x7



$\circ$

| 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 |

=

25 coefficients
9 degrees of freedom

49 coefficients
18 degrees of freedom

What is lost?

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

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

233

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 $3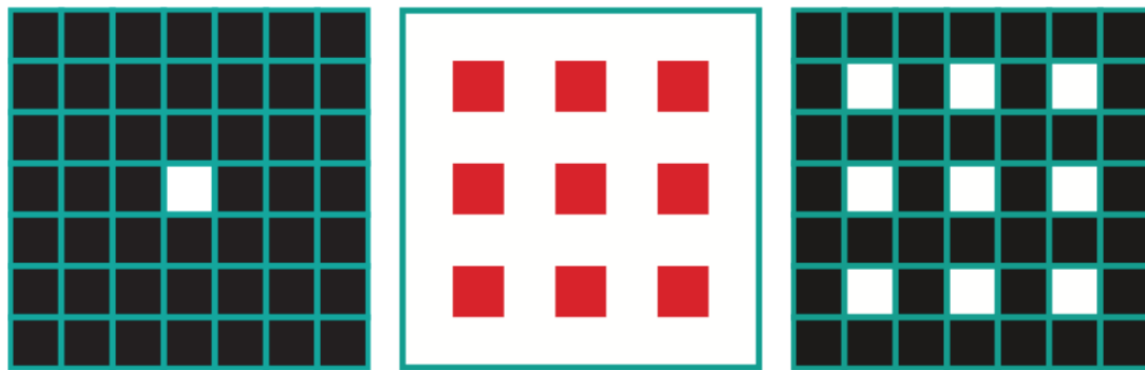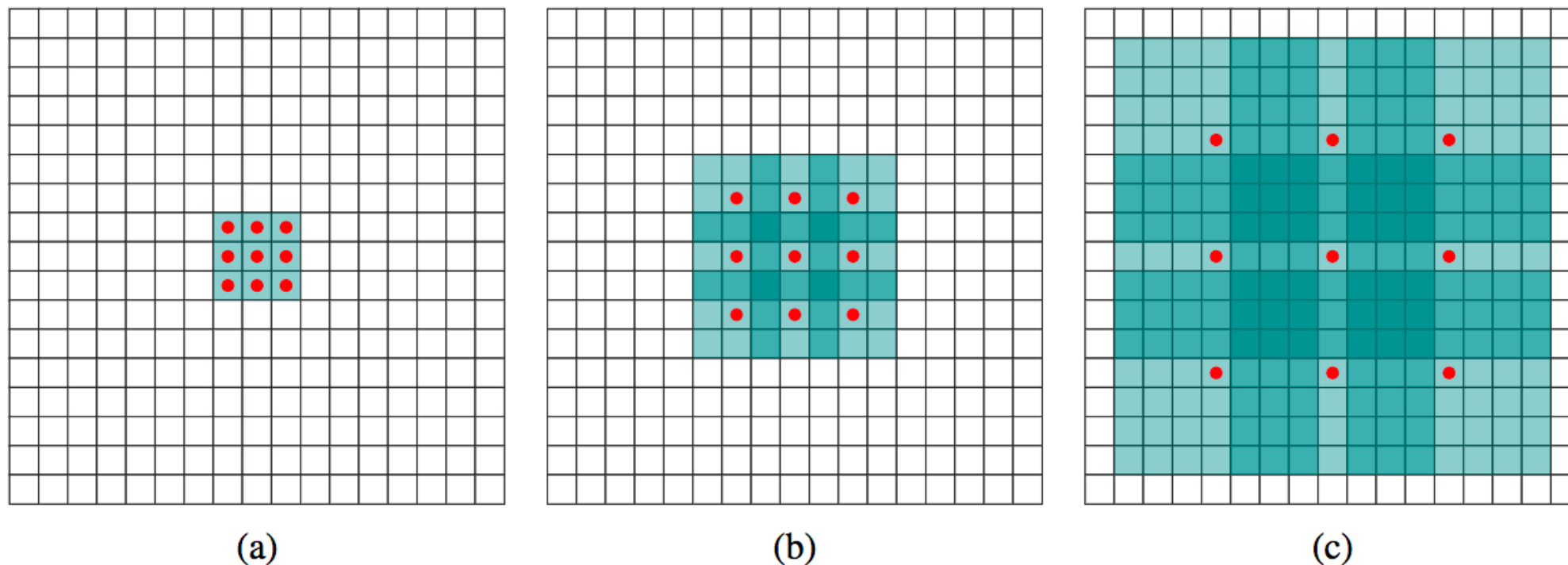 \times 3$. (b) $F_2$ is produced from $F_1$ by a 2-dilated convolution; each element in $F_2$ has a receptive field of $7 \times 7$. (c) $F_3$ is produced from $F_2$ by a 4-dilated convolution; each element in $F_3$ has a receptive field of $15 \times 15$. The number of parameters associated with each 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



34-layer residual

image

7x7 conv, 64, /2
pool, /2
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 128, /2
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 256, /2
3x3 conv, 256
3x3 conv, 256
3x3 conv, 256
3x3 conv, 256
3x3 conv, 256
3x3 conv, 256
3x3 conv, 256
3x3 conv, 256
3x3 conv, 256
3x3 conv, 256
3x3 conv, 256
3x3 conv, 512, /2
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
avg pool
fc 1000

Error: 4.4%

$$\mathbf{x}$$

weight layer

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

weight layer

$\mathbf{x}$
identity

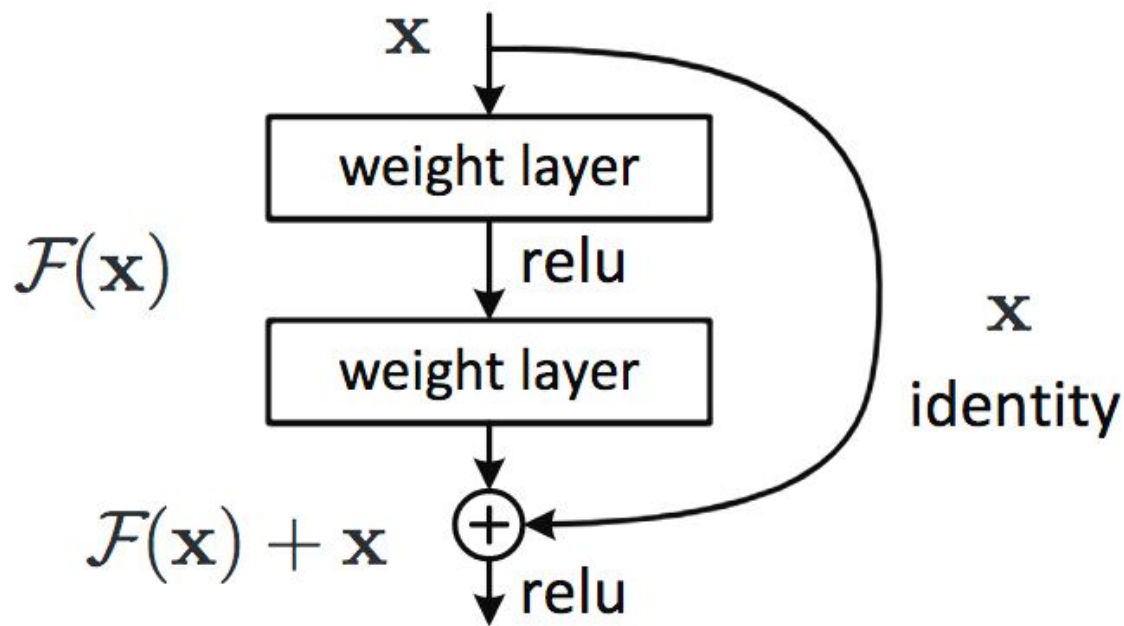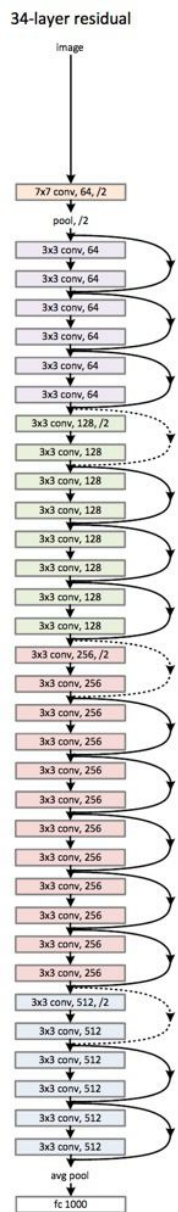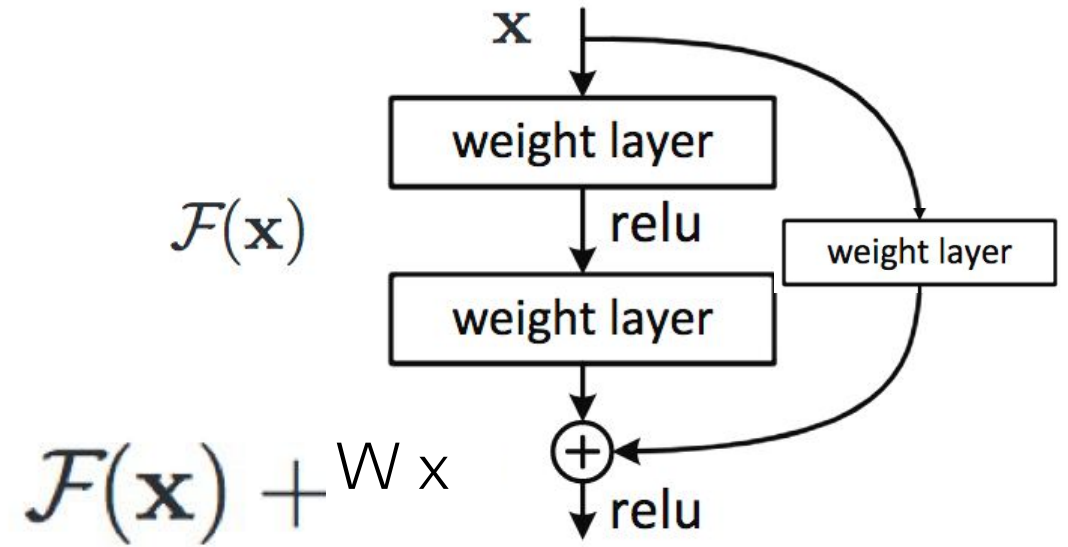$\mathcal{F}(\mathbf{x}) + \mathbf{x}$ $\oplus$

relu

Figure 2. Residual learning: a building block.
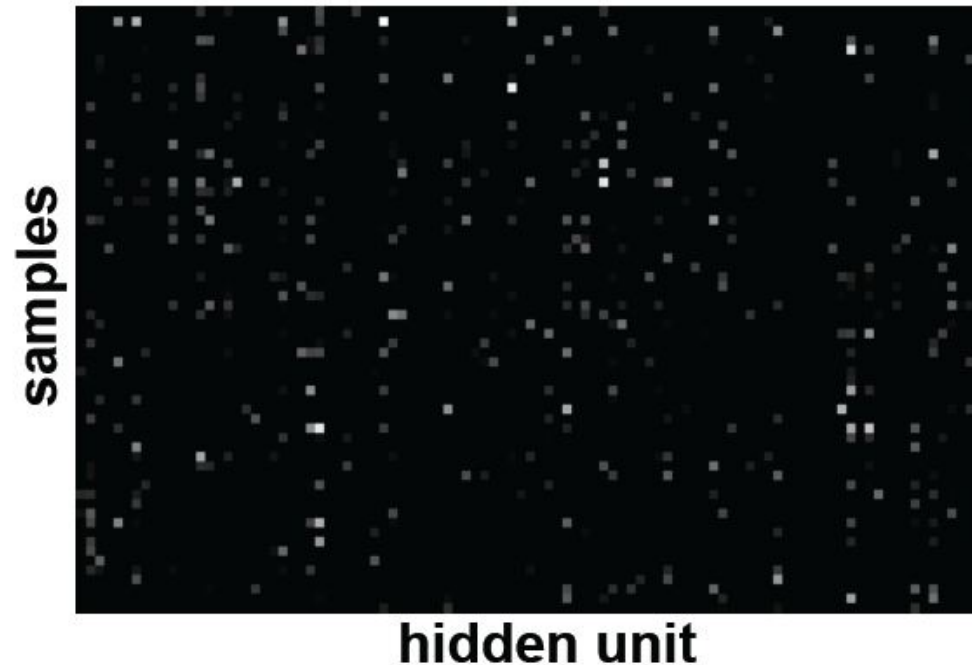
If output has same size as input:



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

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

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

relu

# Other good things to know

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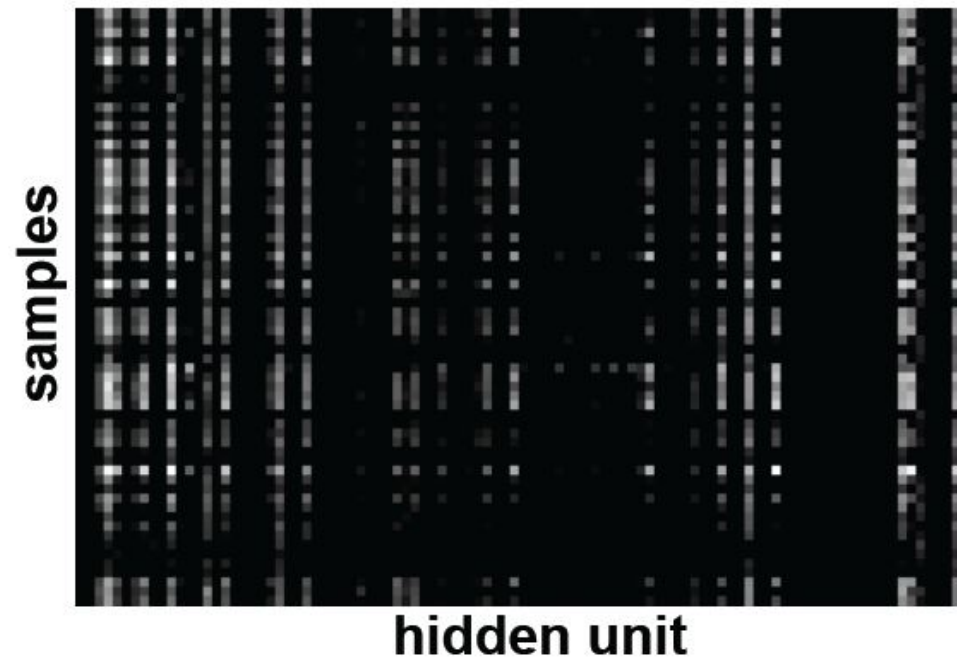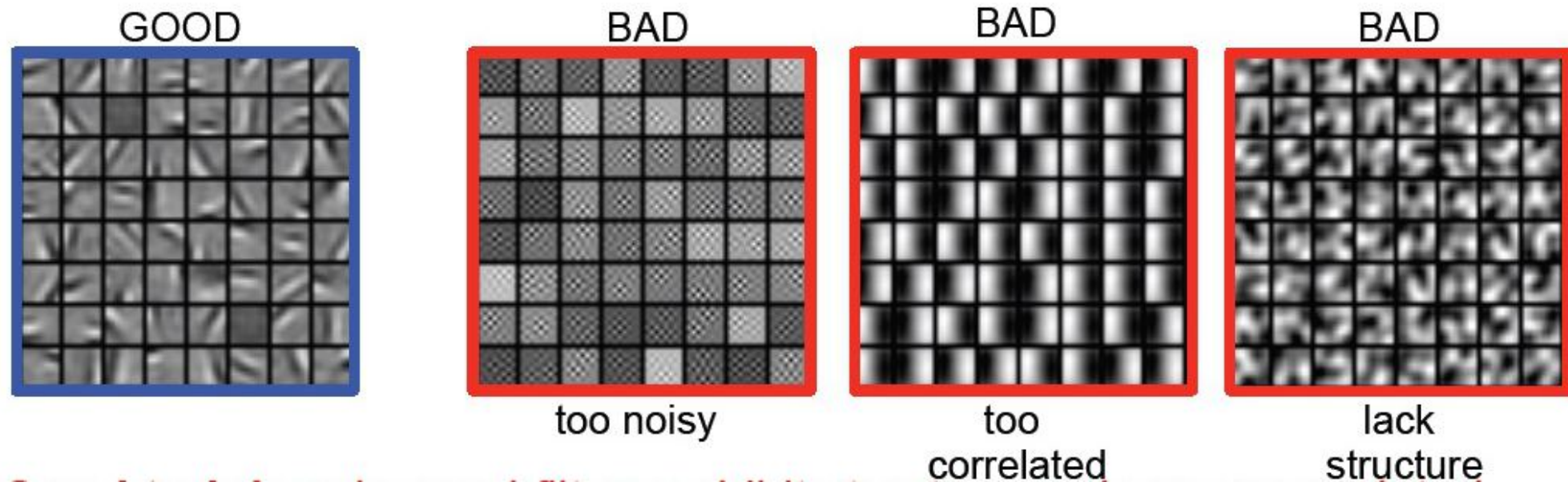- Visualize hidden activations — should be uncorrelated and high variance



**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]
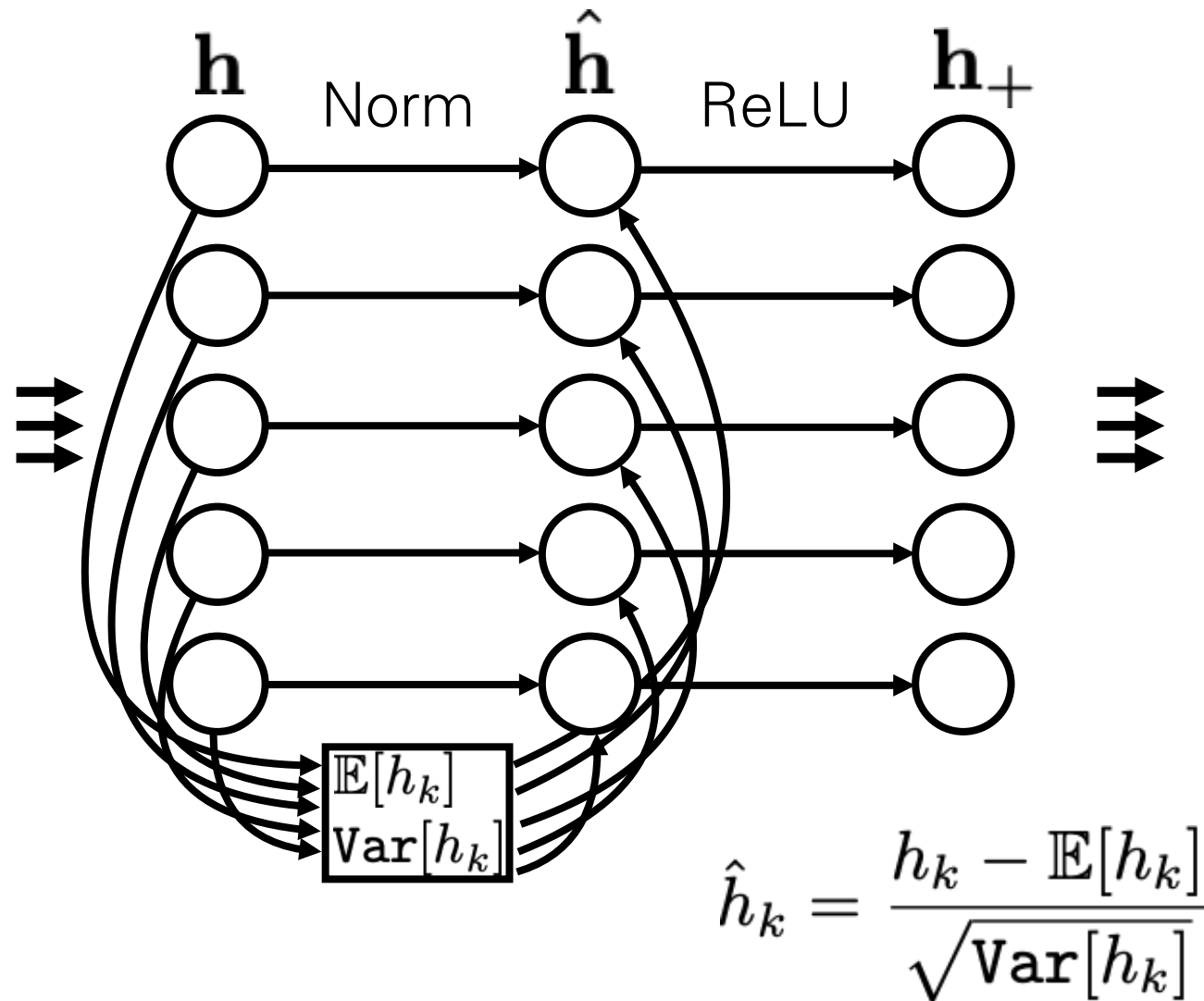
# Other good things to know

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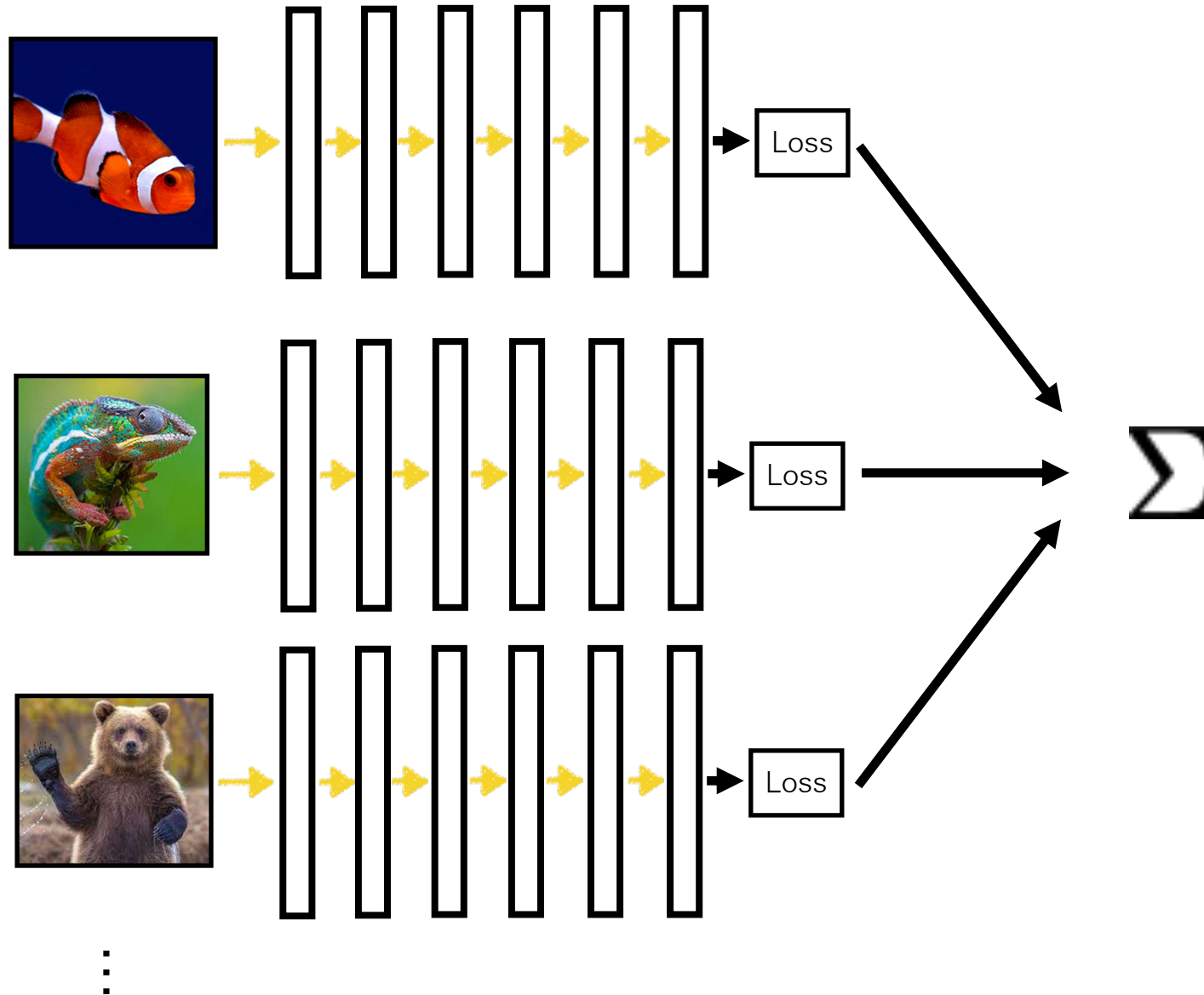- Visualize hidden activations — should be uncorrelated and high variance
- Visualize filters



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



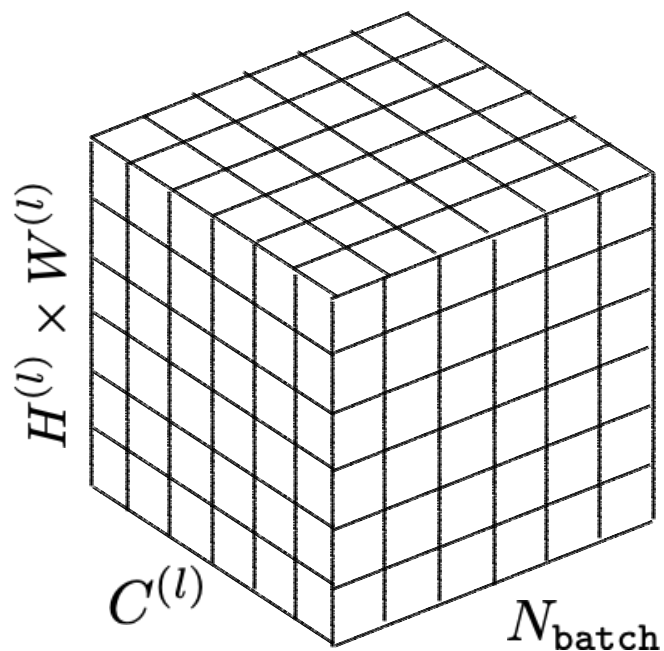$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\mathbf{Var}[h_k]}}$$
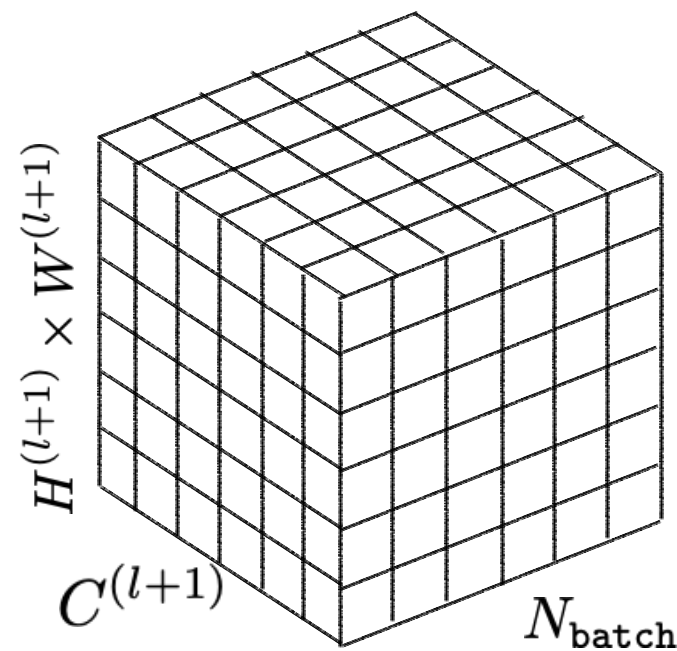
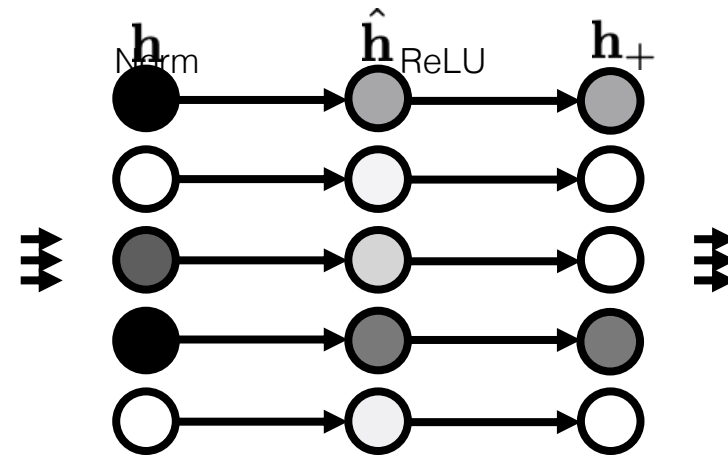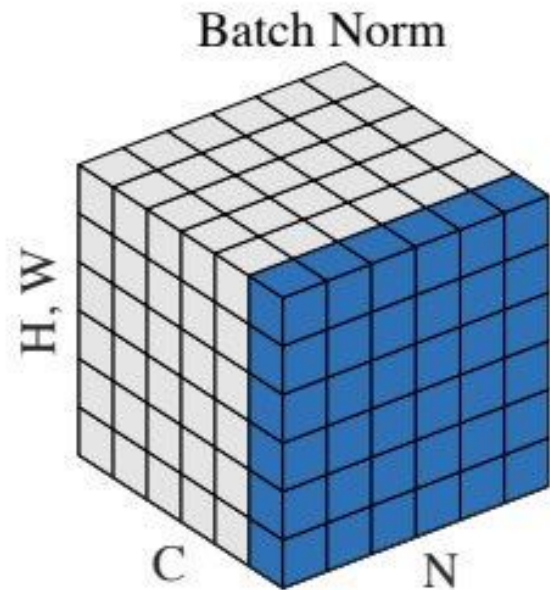# Batch processing

# "Tensor flow"

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

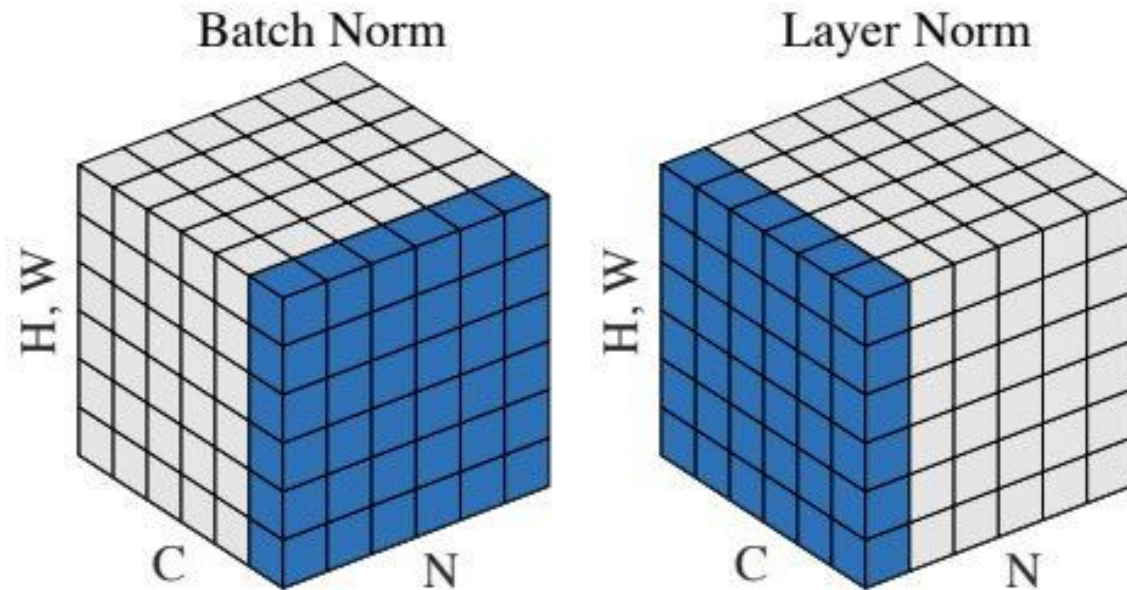$$\mathbf{x}^{(l+1)} \in \mathbb{R}^{N_{\text{batch}} \times H^{(l+1)} \times W^{(l+1)} \times C^{(l+1)}}$$

# Normalization layers



**Batch Norm**

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

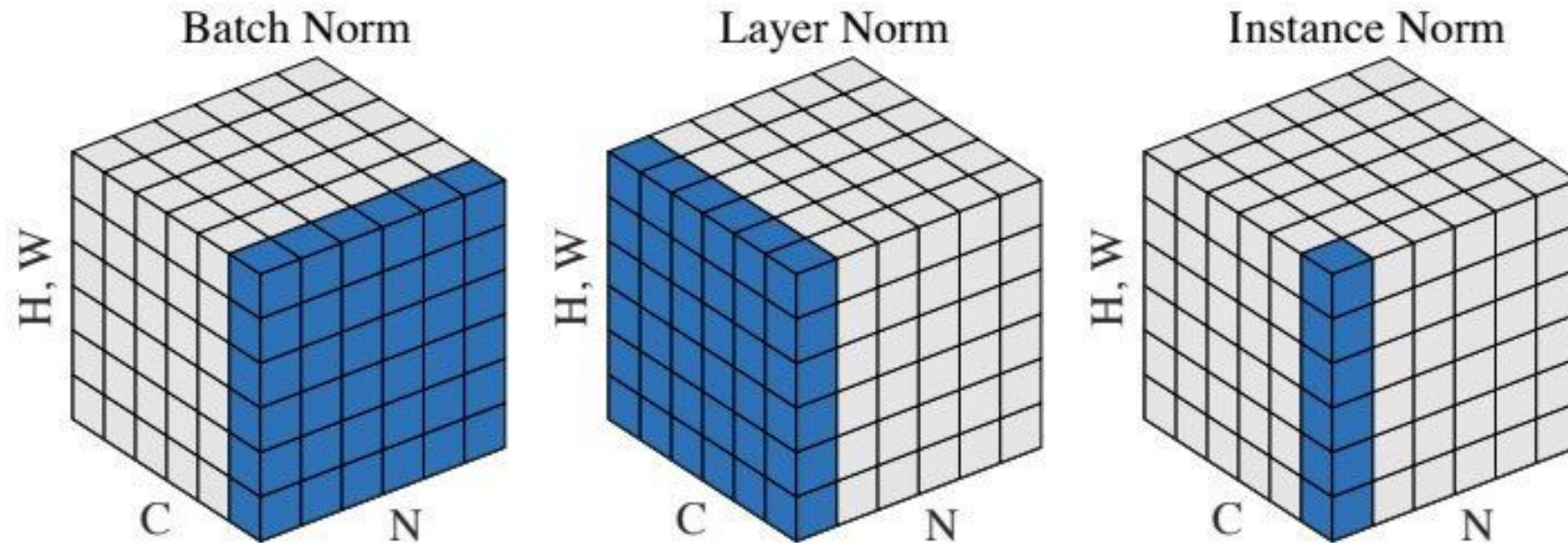[Figure from Wu & He, arXiv 2018]

# Normalization layers



Normalize w.r.t. the mean and variance of the activations of all the hidden units (neurons) on this layer (c).

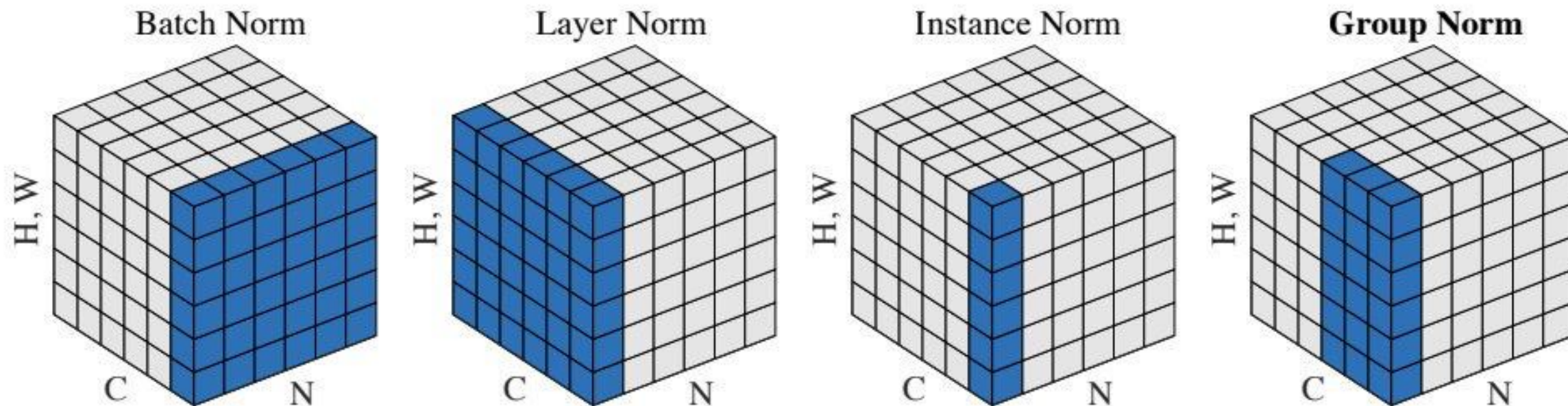[Figure from Wu & He, arXiv 2018]

# Normalization layers



Batch Norm      Layer Norm      Instance Norm

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]
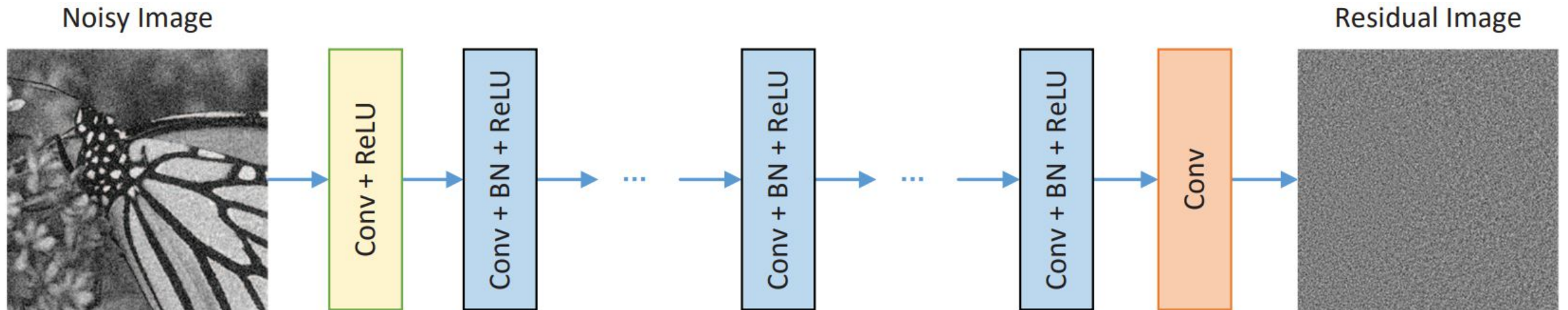
245

# Normalization layers



Batch Norm     Layer Norm     Instance Norm     Group Norm

Might as well…

[Figure from Wu & He, arXiv 2018]

# Applications in Computational Photography

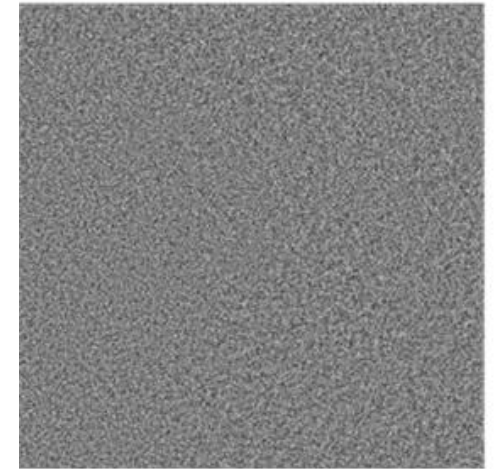# Image Denoising



Key idea: Residual learning

# Image Denoising



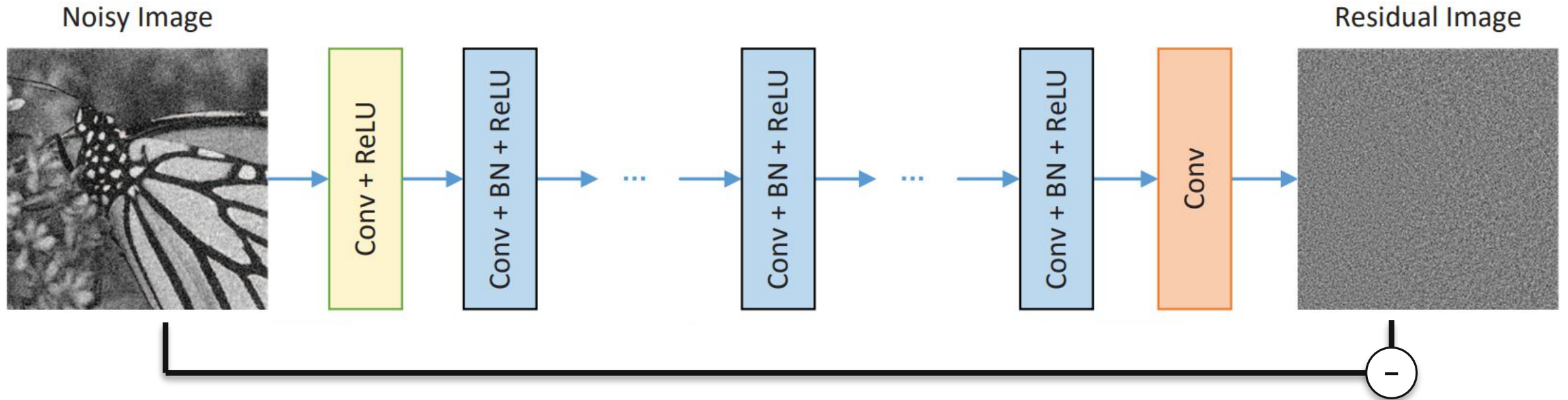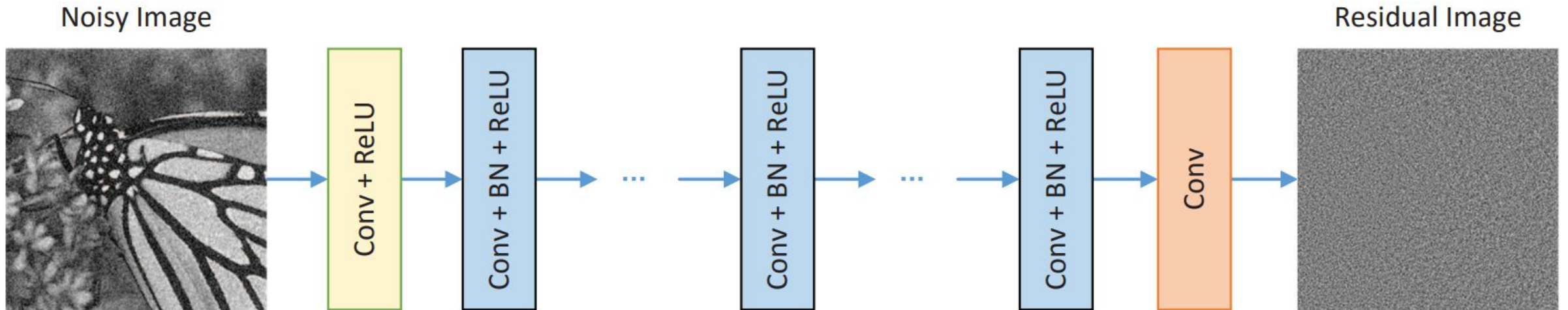Clean image     =     noisy image     −     estimated noise

(Zhang et al., IEEE TIP 2017)

# Image Denoising

# Image Denoising



No fully connected layers - can be applied to any input size

(Zhang et al., IEEE TIP 2017)

# Image Denoising



(a) Ground-truth  (b) Noisy / 17.25dB  (c) CBM3D / 25.93dB  (d) CDnCNN-B / 26.58dB

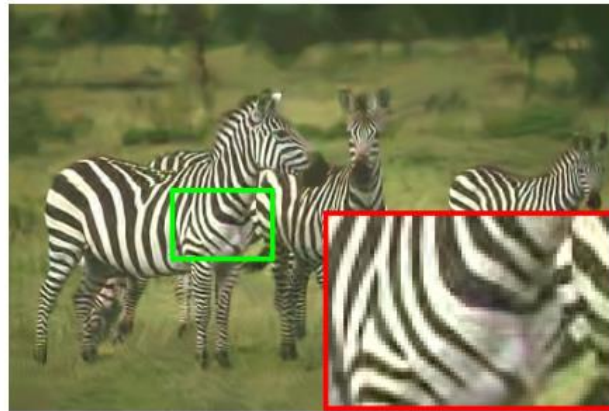(Zhang et al., IEEE TIP 2017)

# Image Denoising



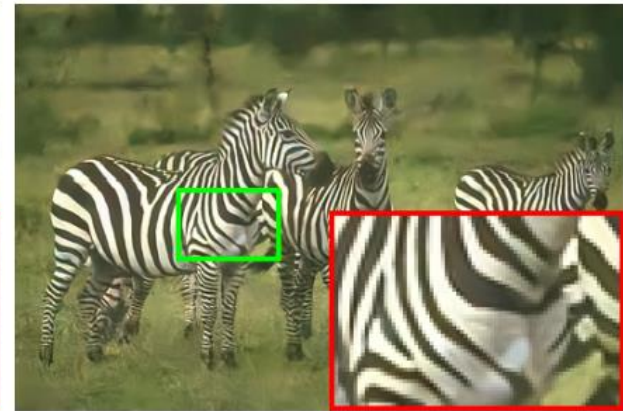(a) Ground-truth (b) Noisy / 15.07dB (c) CBM3D / 26.97dB (d) CDnCNN-B / 27.87dB

(Zhang et al., IEEE TIP 2017)

# Image Denoising



$X_{1...N}, \hat{\sigma}_p$

Input burst
and noise estimate

64    128    256    512    512    512    256    $K^2N$

□ convolution layer    □ average pooling layer
⋯ skip connection    □ bilinear upsampling layer

Per-pixel
Kernels

$X_{1...N}$

$\hat{Y}$
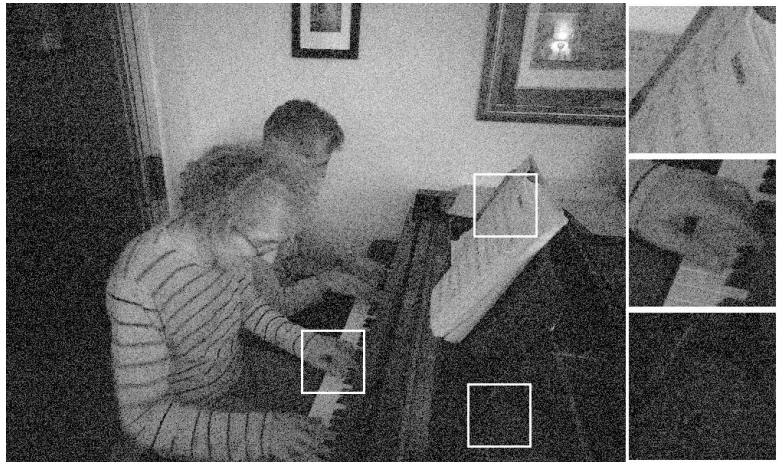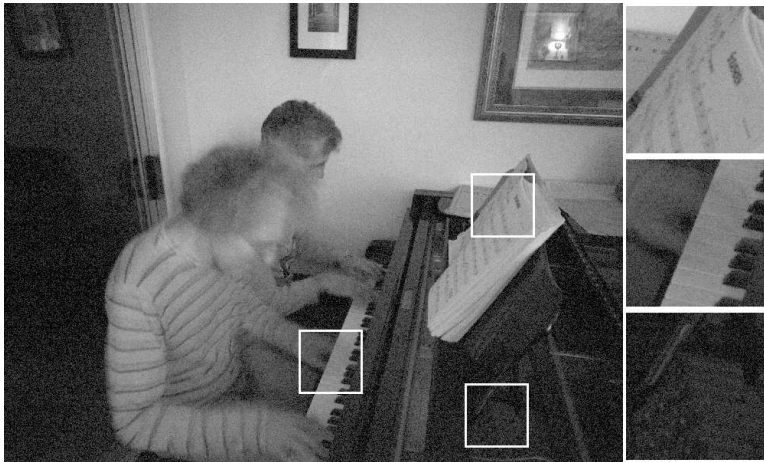
Key ideas:
- Perform denoising in RAW domain considering a burst of images
- Generates a stack of per-pixel filter kernels that jointly aligns, averages, and denoises a burst of images.
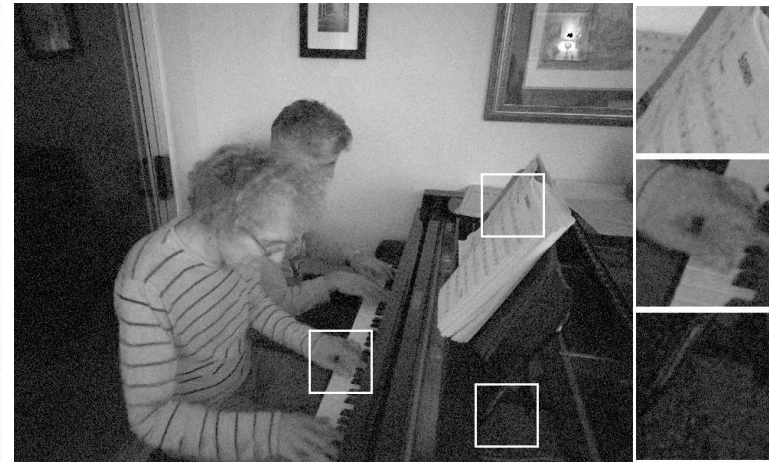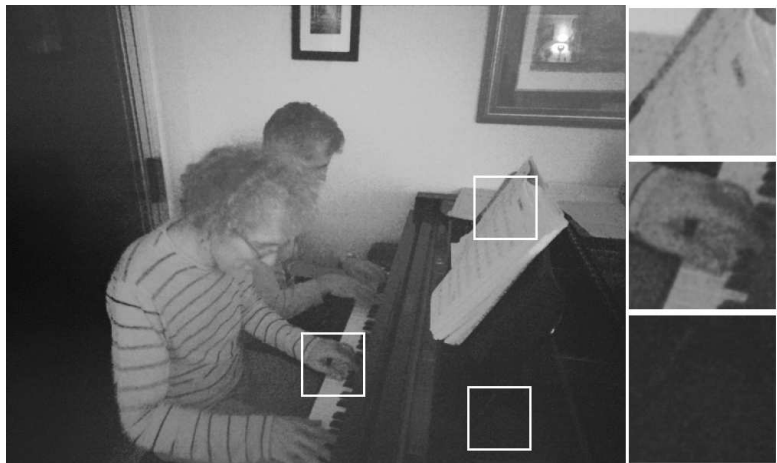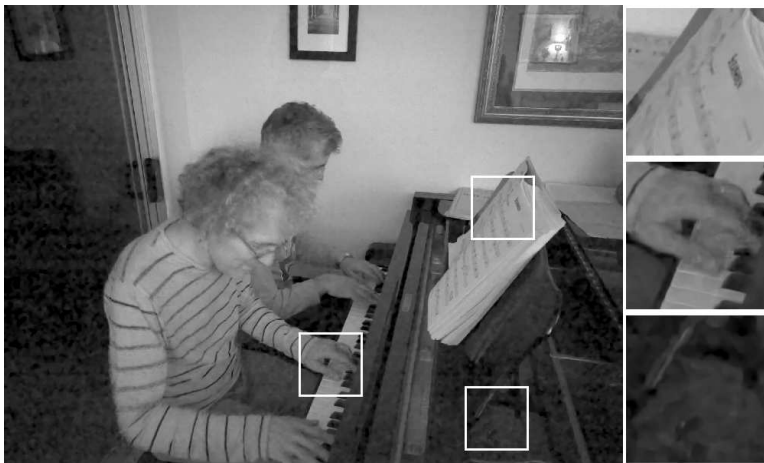
(Mildenhall et al., CVPR 2018

# Image Denoising



(a) Reference frame

(b) Burst average

(c) HDR+ [8]

(d) Non-local means [3]

(e) VBM4D [17]

(f) Our KPN model

(Mildenhall et al., CVPR 2018

# Image Denoising



Reference frame        (a) Reference     (b) Average     (c) HDR+     (d) NLM     (e) VBM4D     (f) Ours (KPN)

(Mildenhall et al., CVPR 2018
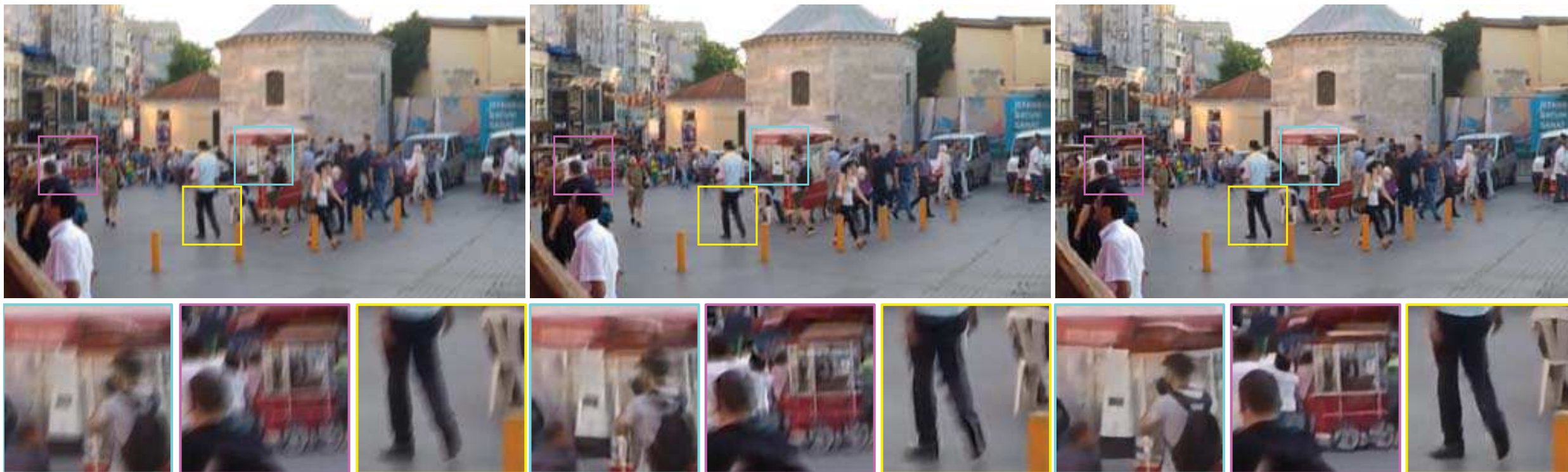
# Image Deblurring



Key idea: Multi-scale processing, i.e., use image pyramid to process and deblur

# Image Deblurring



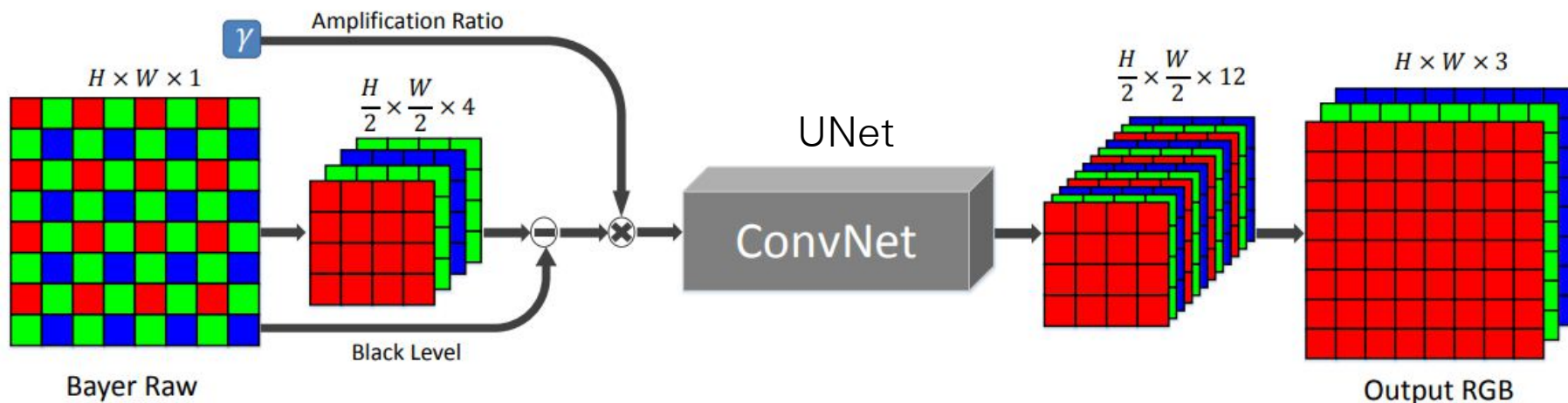Input                    Single-scale                    Multi-scale
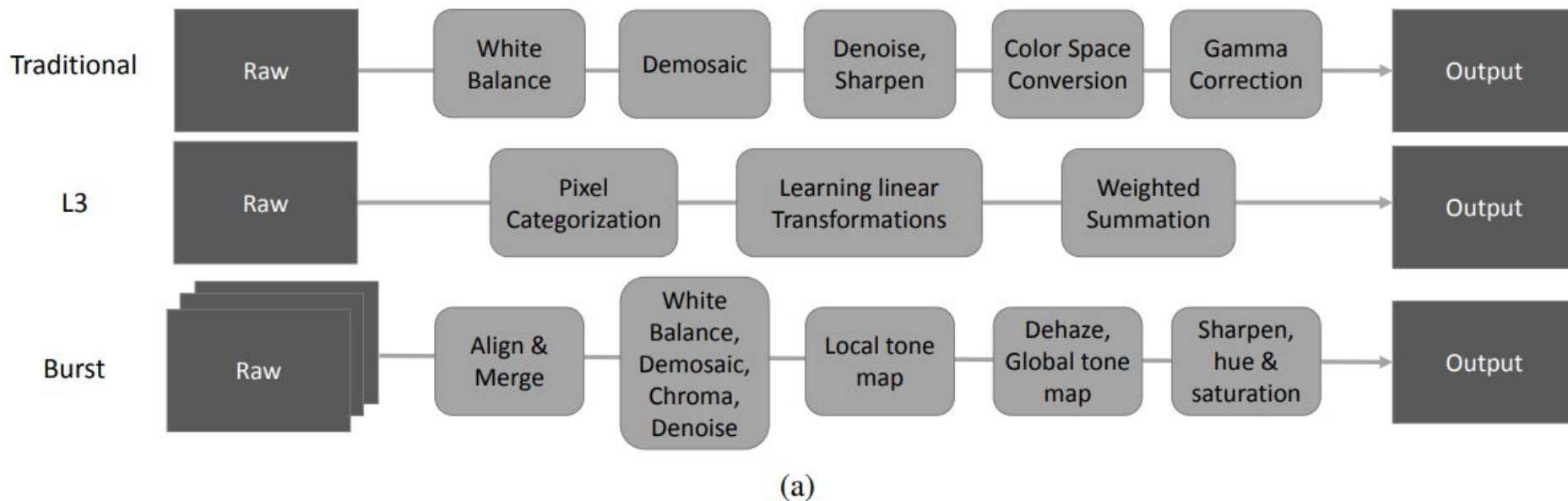
(Nah et al., CVPR 2017) 258

# Image Deblurring

Single-scale

Multi-scale



(Nah et al., CVPR 2017) 259

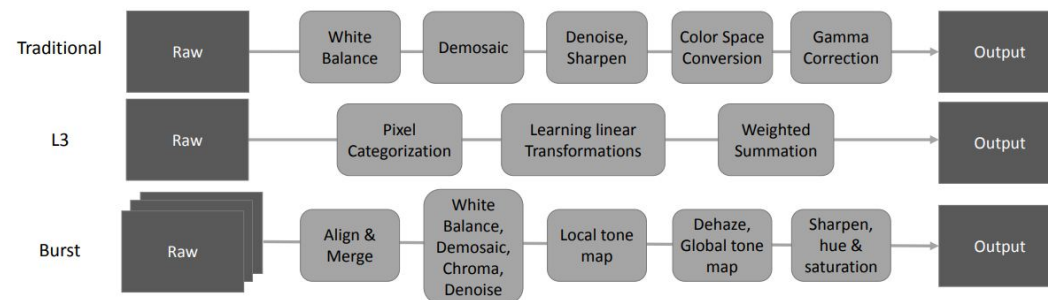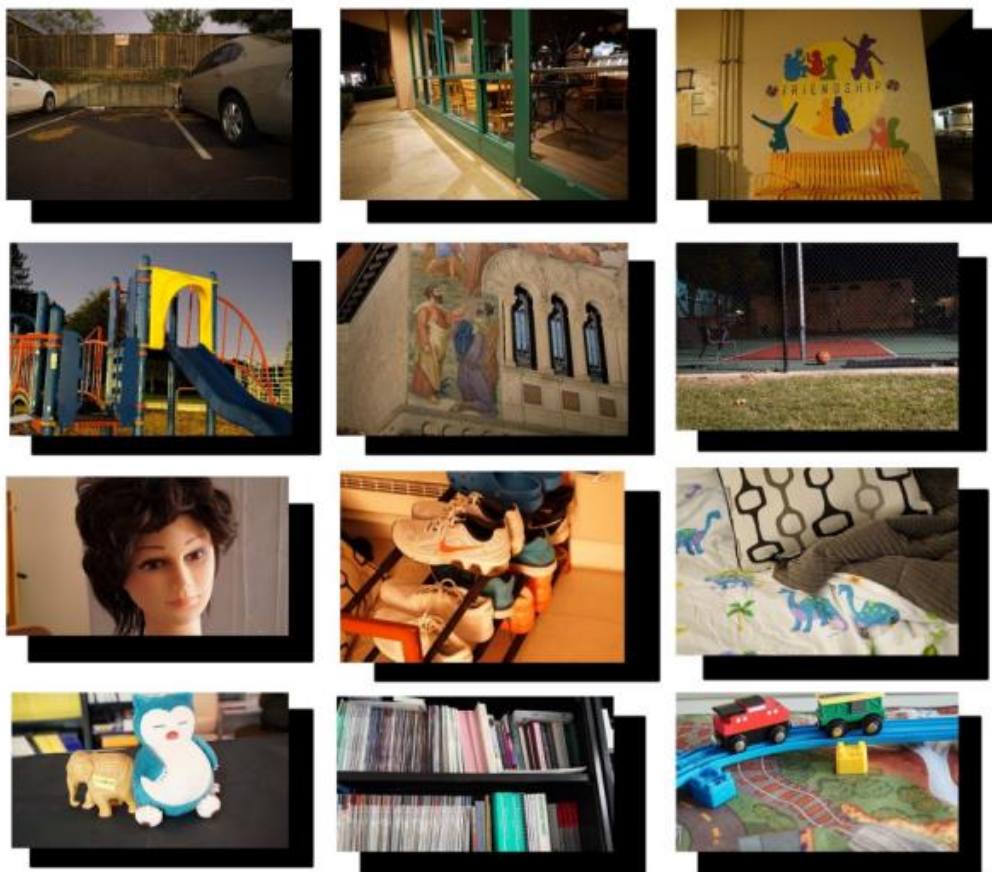# Learned ISP



(a)



Key idea: Learn a convnet to enhance extremely low-light images

(Chen et al., CVPR 2018)

# Learned ISP



Trained on short-exposure (noisy) / long-exposure image pairs

# Learned ISP

- Key ideas:

1. A frame-level enhancement network that works in a coarse-to-fine manner

2. Extension of this model to a burst of dark images

Coarse Network

$E_c$  $D_c$  $x^c$  $\hat{n}$

$x$

Fine Network

$E_f$  $D_f$  $\hat{y}$

$t$

$t_1$  $E_s$

$t_2$  $E_s$

max-pool

$D_s$  $\hat{y}$

$t_m$  $E_s$

Burst Network

(Karadeniz et al., IEEE TIP 2021)

Learned ISP

Noisy input

Learned ISP

Karadeniz et al. (single)

Chen et al, 2018

# Learned ISP

Traditional pipeline

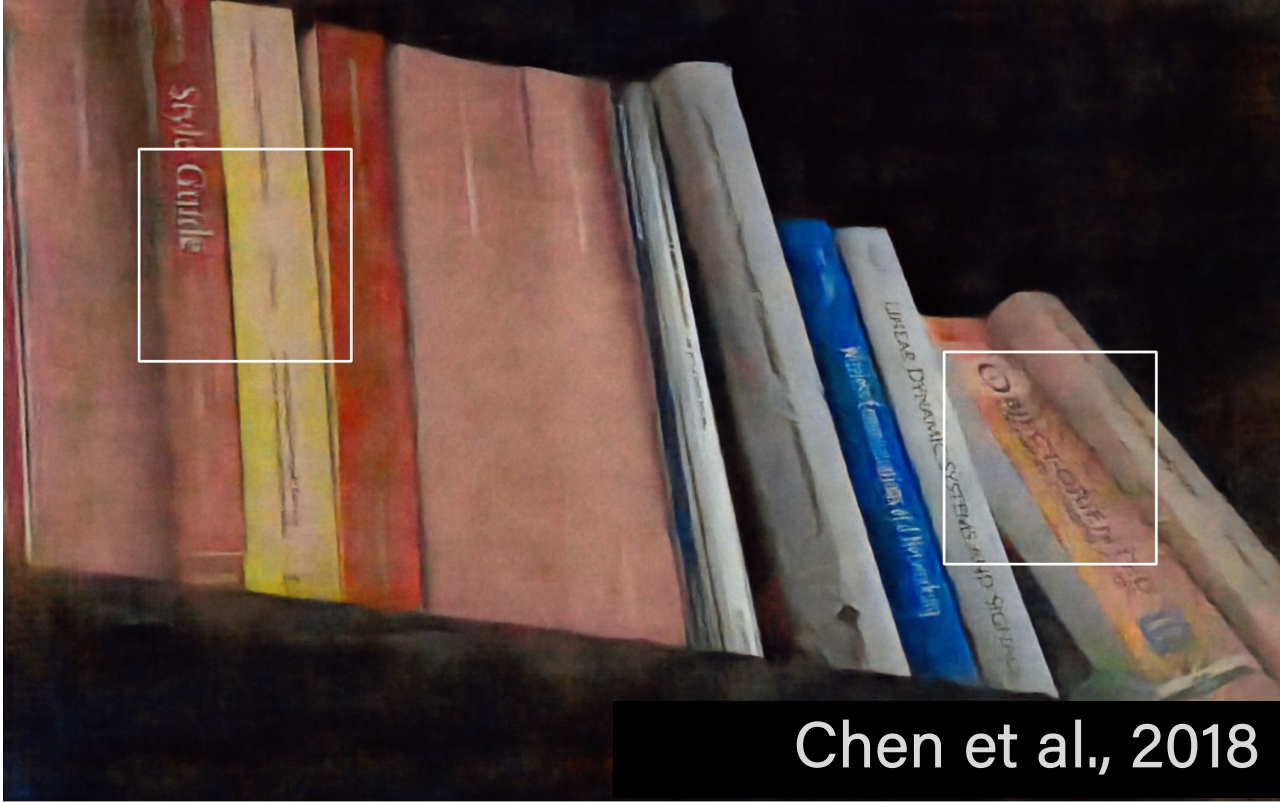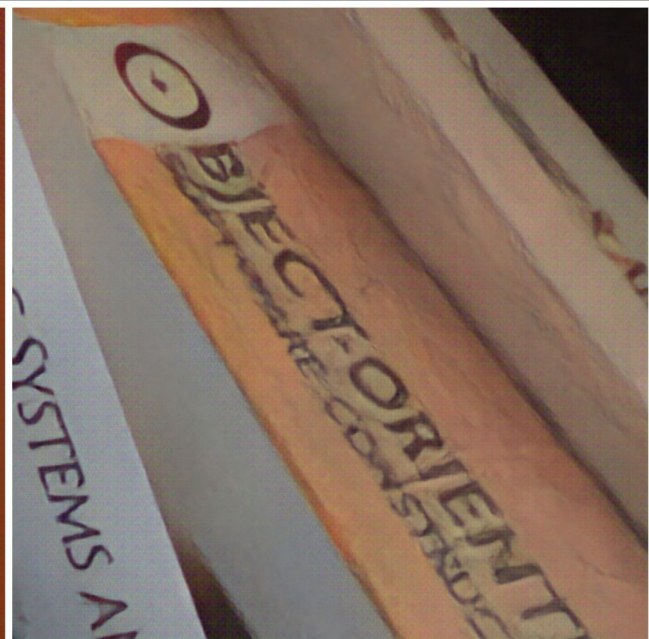Learned ISP

Traditional pipeline + Scaling

Learned ISP

Chen et al., 2018 (Ensemble)

Learned ISP

Karadeniz et al. (Burst)

Learned ISP

Ground truth

# Learned ISP



Key idea: deep neural networks as a controllable camera simulator to synthesize raw image data under different camera settings, including exposure time, ISO, and aperture.

(Ouyang et al., CVPR 2021)

# Learned ISP



(Ouyang et al., CVPR 2021)
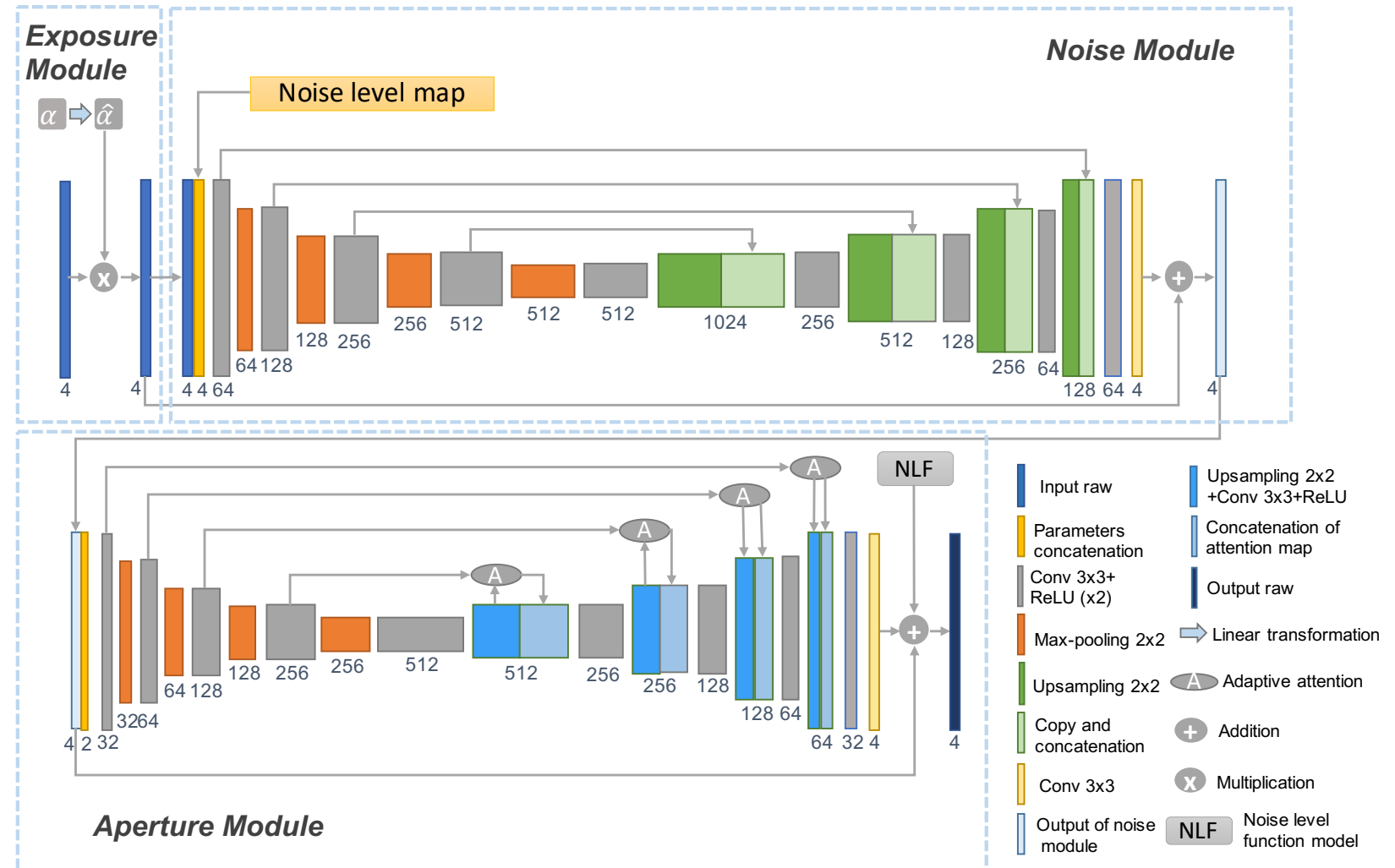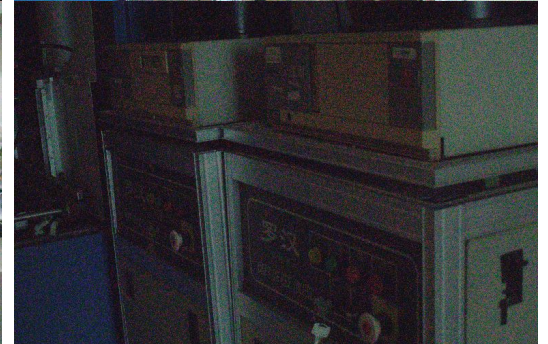
# Learned ISP



| Input | Exposure Module | Noise Module | Full Model | Ground Truth |

(Ouyang et al., CVPR 2021) 272
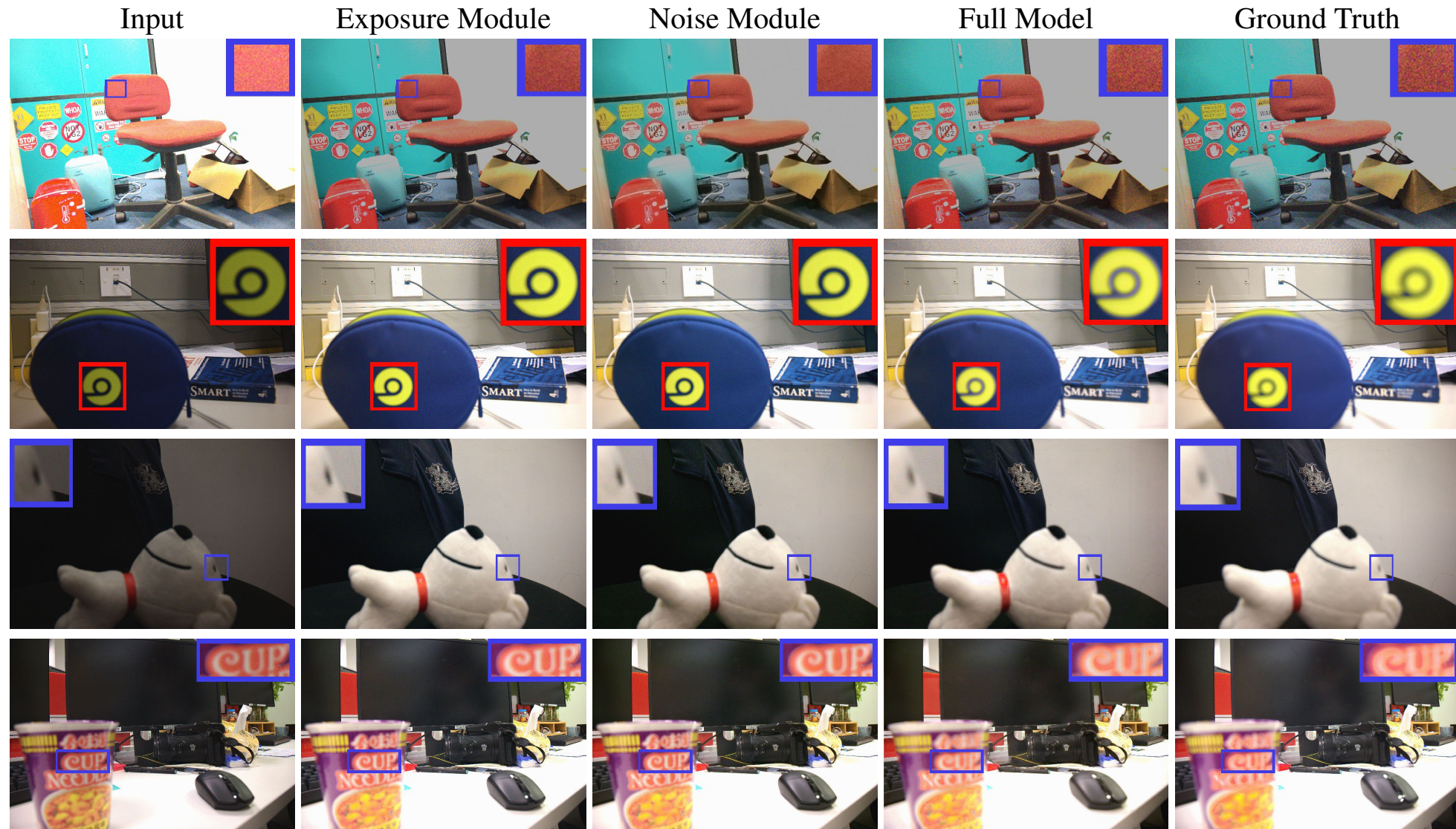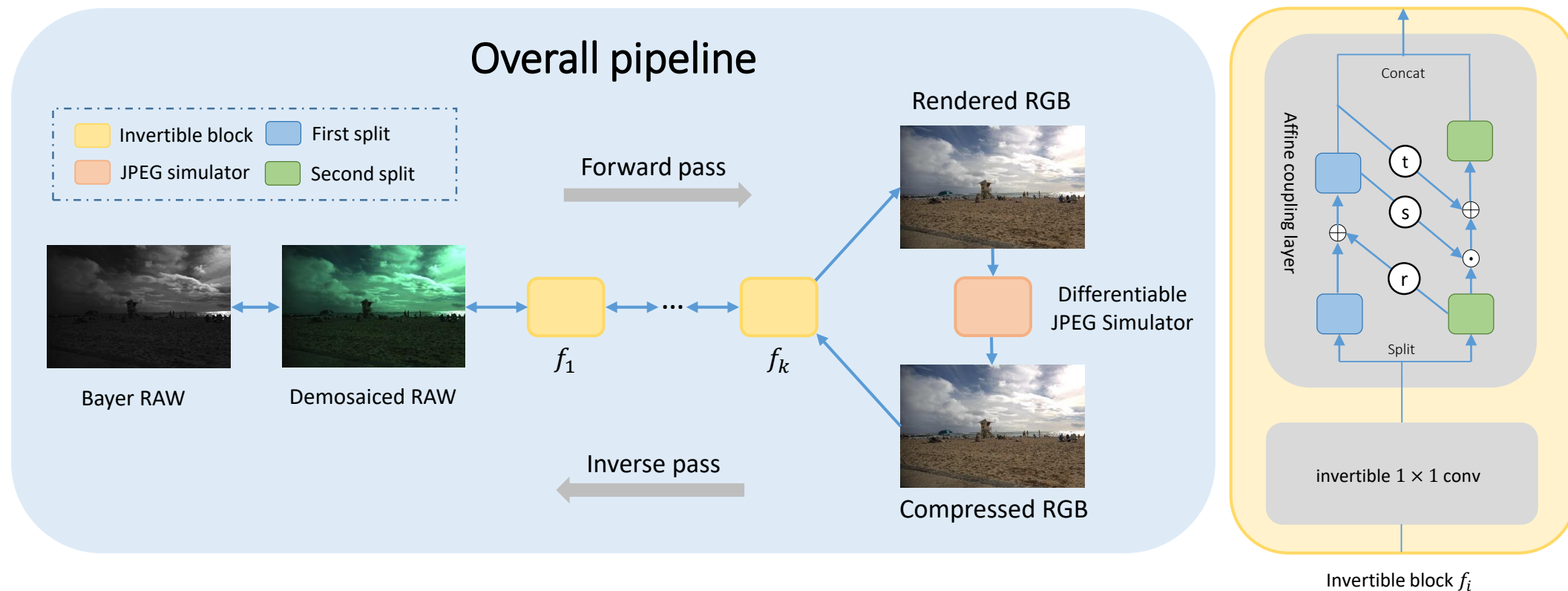
# Invertible ISP



$$L = ||f(\mathbf{x}) - \mathbf{y}||_1 + \lambda||f^{-1}(\mathbf{y}) - \mathbf{x}||_1,$$

Key idea: Use invertible neural networks to design the invertible structure and integrate a differentiable JPEG simulator to enhance the network stability to JPEG compression.

(Xing et al., CVPR 2021)

# Invertible ISP

Camera RAW

Our rendered RGB

$$\text{Our} \mathbin{|} \text{ISP } f$$

$$\text{Inverse} \mathbin{|} \text{ISP } f^{-1}$$

RAW error map

Our recovered RAW
(PSNR: 45.26)

## Applications

HDR reconstruction

Image retouching

Other application
(*i.e.*, RAW compression)

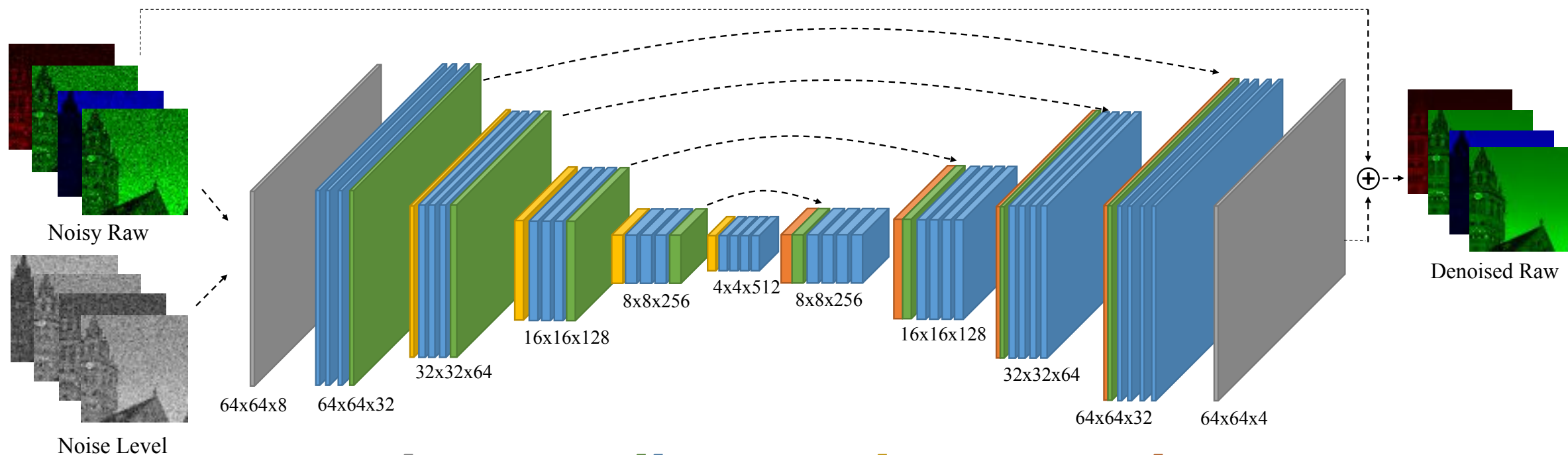(Xing et al., CVPR 2021)

# Image Denoising via Invertible ISP



Key idea: Transform sRGB images into RAW domain and perform denoising there

(Brooks et al., CVPR 2019)

# Image Denoising via Invertible ISP



Noisy Raw

Noise Level

64x64x8    64x64x32

32x32x64

16x16x128

8x8x256    4x4x512    8x8x256

16x16x128

32x32x64

64x64x32    64x64x4

Denoised Raw

Input/Output Layers    Convolutional Layers    2x Downsampling Layers    2x Upsampling Layers

(a) Noisy Input    (b) Our Model

(Brooks et al., CVPR 2019)
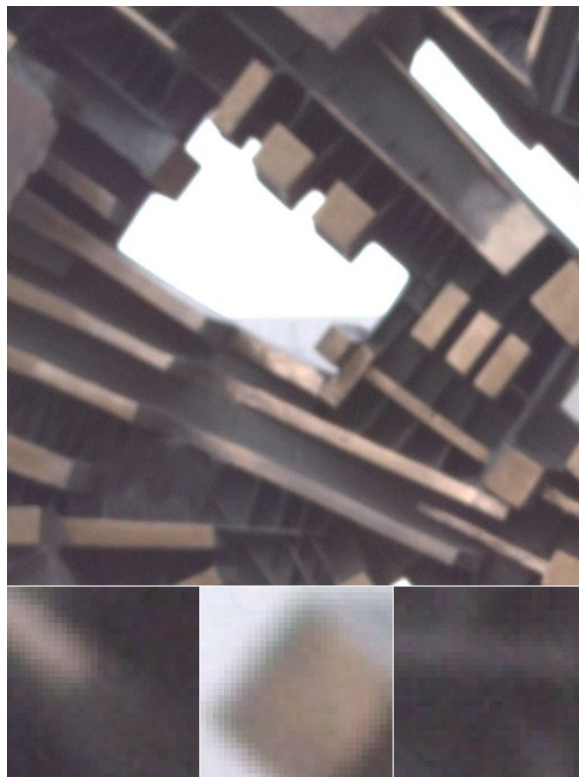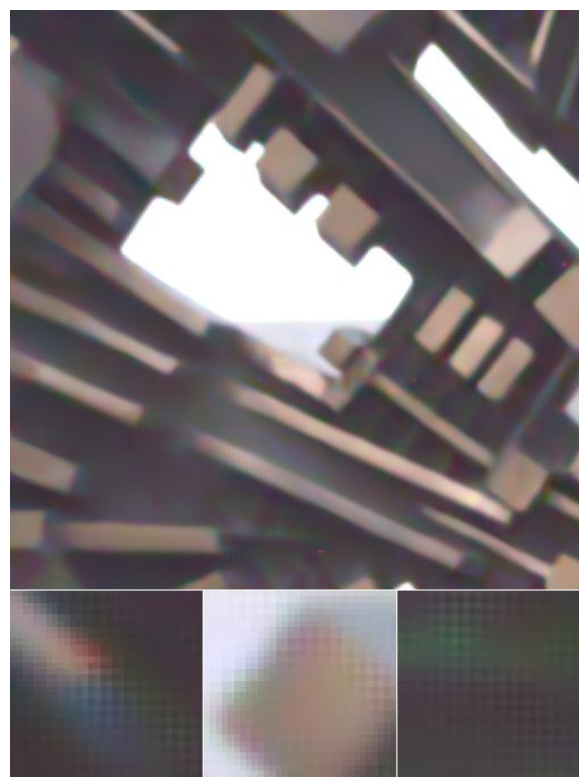
276

# Image Denoising via Invertible ISP



(a) Noisy Input, PSNR = 18.76    (b) Ground Truth    (c) N3Net [31], PSNR = 32.42    (d) Our Model, PSNR = 35.35

(Brooks et al., CVPR 2019)

# Image Enhancement



LOW-RES COEFFICIENT PREDICTION

§3.1.2 local features $L^i$

§3.2 bilateral grid of coefficients $A$

§3.1.4 fusion $F$

full-res input $I$

low-res input $\tilde{I}$

§3.1.1 low-level features $S^i$

§3.1.3 global features $G^i$

FULL-RES PROCESSING

pixel-wise network

slicing layer

apply coefficients

§3.4.1 guidance map $g$

§3.3 sliced coefficients $\bar{A}$

§3.4.2 full-res output $O$

Key idea: Learn to imitate a reference operator, work much faster

(Gharbi et al., SIGGRAPH 2017)

# Image Enhancement



HDR+ **32.7 dB**

Face brightening **38.9 dB**

Human retouch **33 dB**

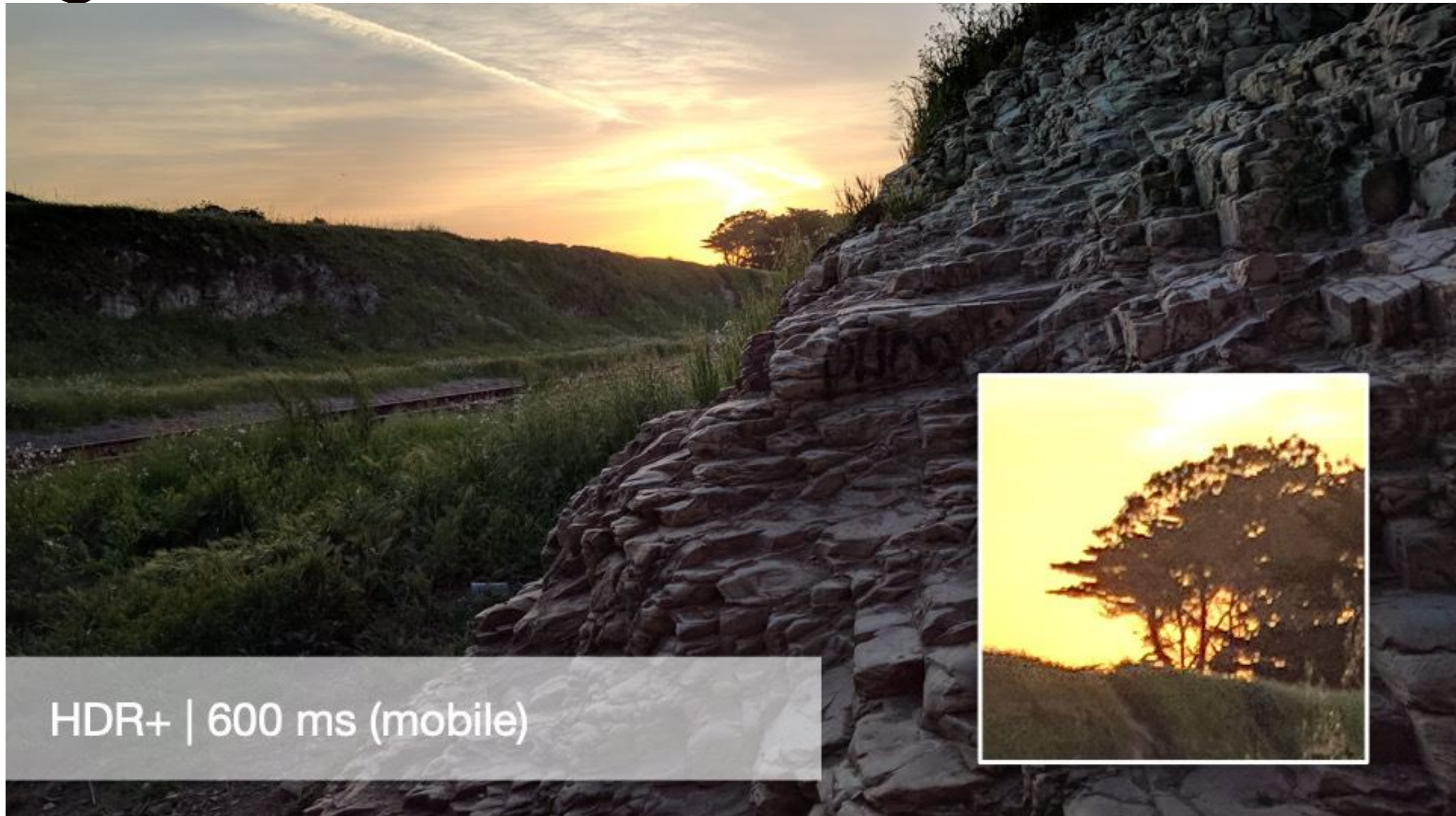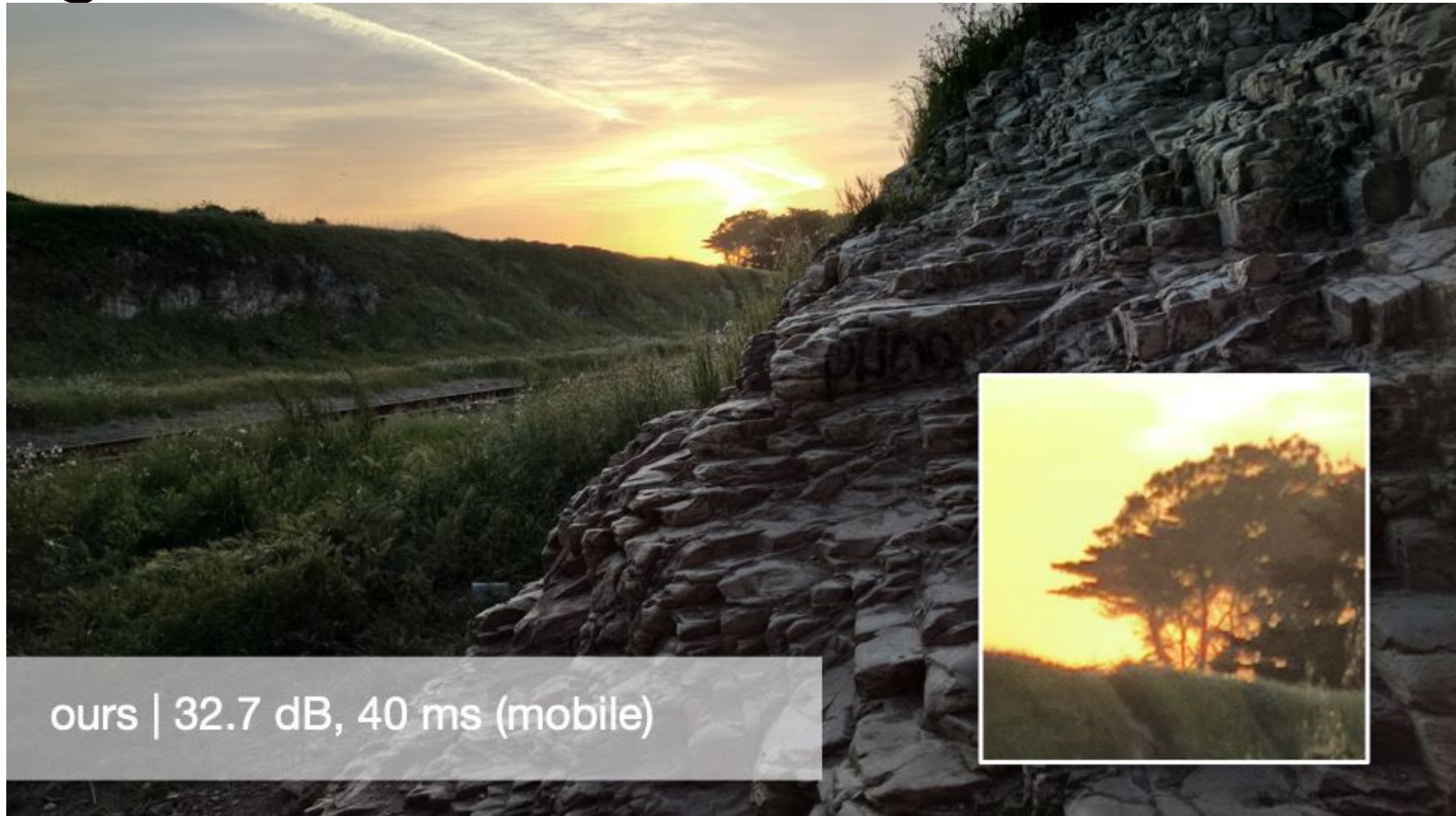input      reference      our output      reference (cropped)      our output (cropped)      difference

(Gharbi et al., SIGGRAPH 2017)

# Image Enhancement



full-res input | 12 Mpixels

(Gharbi et al., SIGGRAPH 2017)

# Image Enhancement



HDR+ | 600 ms (mobile)

(Gharbi et al., SIGGRAPH 2017)
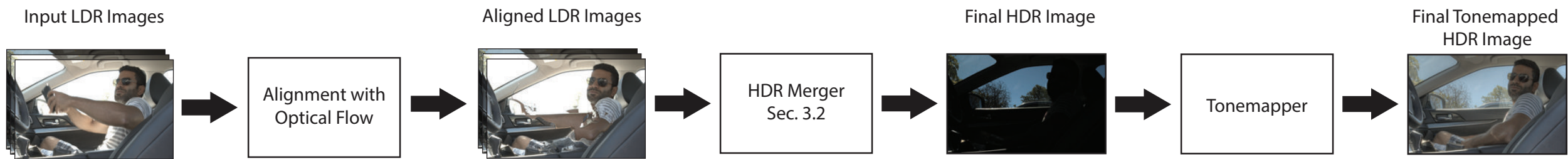
# Image Enhancement



ours | 32.7 dB, 40 ms (mobile)

(Gharbi et al., SIGGRAPH 2017)

# Deep HDR Reconstruction



Key idea: Cast HDR reconstruction problem as a learning problem

Kalantari and Ramamoorthi, SIGGRAPH

# Deep HDR Reconstruction



LDR Images | Our Tonemapped HDR Image | Kang (40.02 dB) | Sen (46.12 dB) | Ours (48.88 dB) | Ground Truth

(Kalantari and Ramamoorthi, SIGGRAPH 2017)

# Deep HDR Reconstruction



Kang et al.
37.24

Oh et al.
36.27

Sen et al.
41.85

Hu et al.
38.23

Ours
43.58

Ground
Truth

(Kalantari and Ramamoorthi, SIGGRAPH 2017)

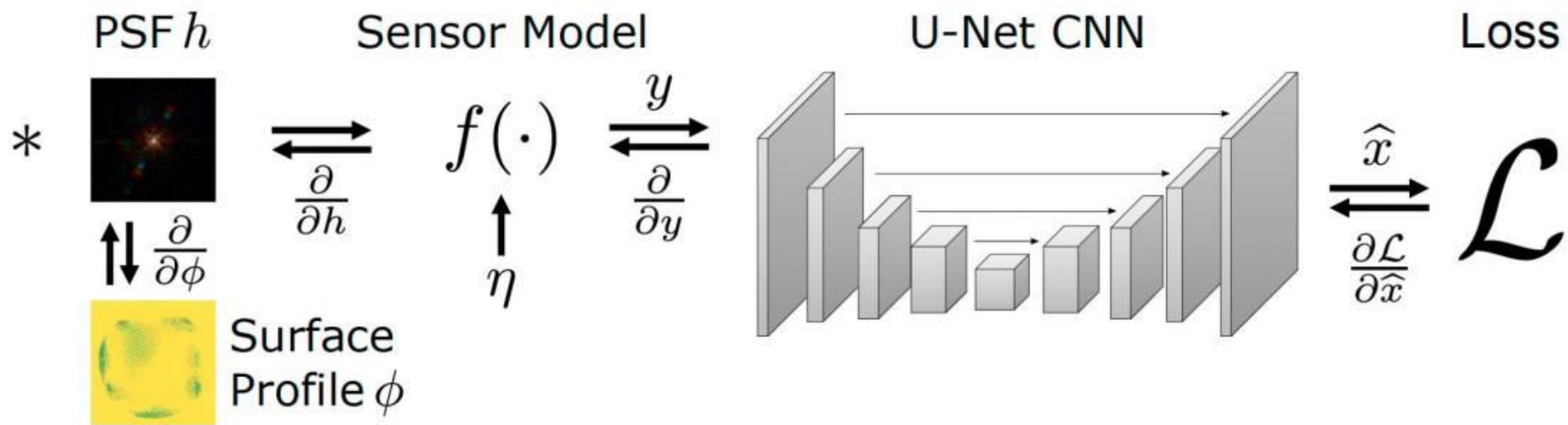# Deep Optics for HDR Imaging



(a) Star PSF

(b) E2E PSF

Key ideas:
- Minimize difference between reconstruction and tone-mapped GT images
- Jointly train an optical encoder and electronic decoder

(Metzler et al., CVPR 2020)

# Deep Optics for HDR Imaging



HDR Training Dataset $*$ PSF $h$ $\rightleftharpoons \frac{\partial}{\partial h}$ $\updownarrow \frac{\partial}{\partial \phi}$ Surface Profile $\phi$ Sensor Model $f(\cdot)$ $\eta$ $\rightleftharpoons^{y}_{\frac{\partial}{\partial y}}$ U-Net CNN $\rightleftharpoons^{\widehat{x}}_{\frac{\partial \mathcal{L}}{\partial \widehat{x}}}$ $\mathcal{L}$ Loss

# Deep Optics for HDR Imaging



LDR Image | E2E Measurement | E2E Reconstruction

(Metzler et al., CVPR 2020)

# Image Relighting

# Image Relighting



Residual connections                    bottleneck

# Image Relighting



Light-stage dataset capture (Google)

(Sun et al., ACM TOG 2019)

# Image Relighting

OLAT photos
(columns)

$$b = Ax$$
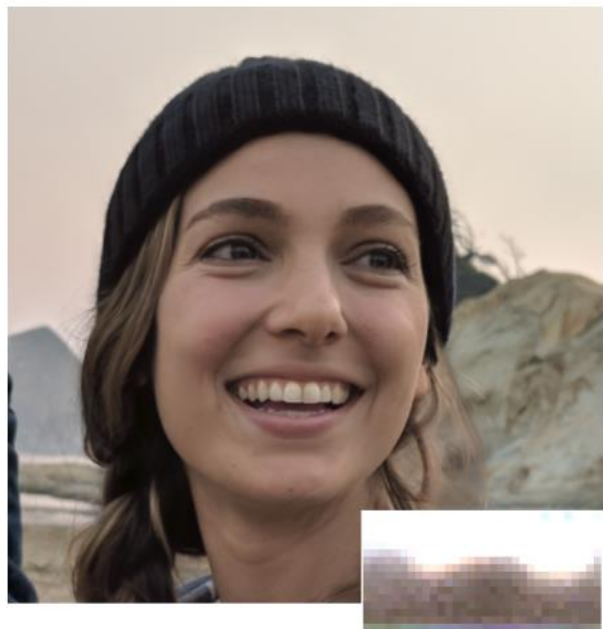
Re-rendered
image

Environment
map
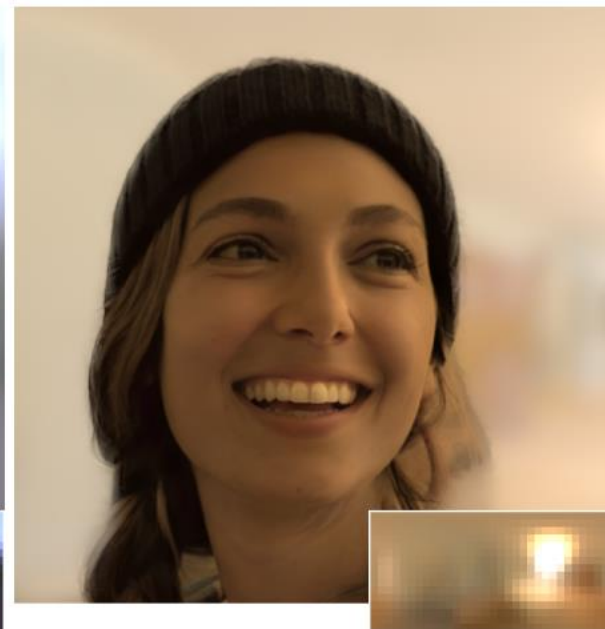


(a) OLAT images (7 cameras).

(b) Ground-truth renderings.

292

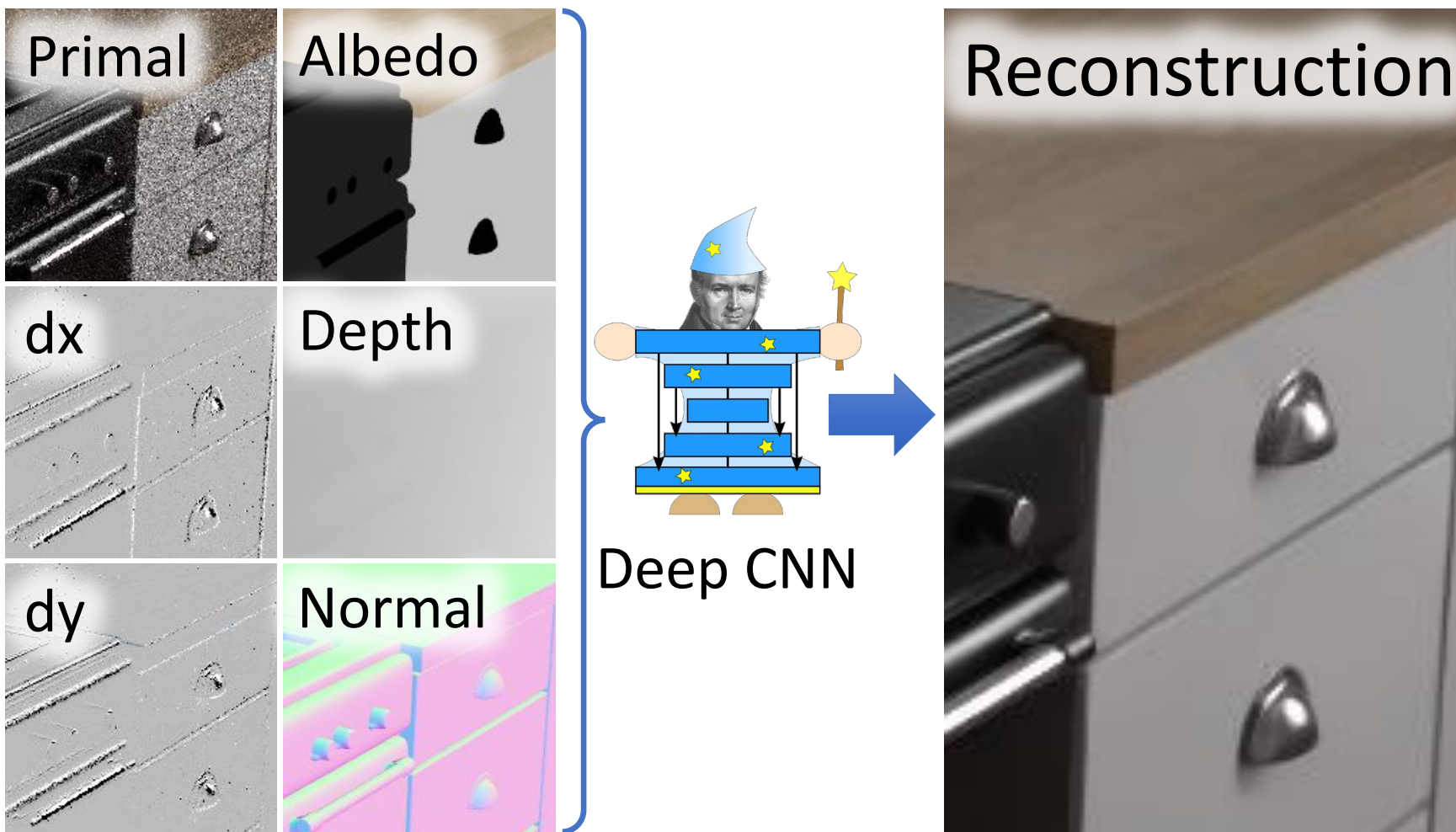# Image Relighting



(a) Input image and estimated lighting

(b) Rendered images from our method under three novel illuminations

(Sun et al., ACM TOG 2019)

# Image Relighting

Now a feature in Google Pixel phones



(Sun et al., ACM TOG 2019)

# Gradient-Domain Rendering



Key idea: Cast gradient-domain rendering as a learning problem

(Kettunen et al., SIGGRAPH 2019) 295

# Gradient-Domain Rendering



Based on DenseNet [Huang2016] and U-Net [Ronneberger2015]

(Kettunen et al., SIGGRAPH 2019)

# Gradient-Domain Rendering

Sponza (Easy)

(Kettunen et al., SIGGRAPH 2019)

# Gradient-Domain Rendering



Screened Poisson, 4 spp

Ours, 4 spp (Equal Time)

720p:   Rendering Time + 1 sec

Rendering Time + 0.3 sec

(Kettunen et al., SIGGRAPH 2019)

# Gradient-Domain Rendering



(Kettunen et al., SIGGRAPH 2019)

# **Next Lecture:**
# Deep Generative Models