

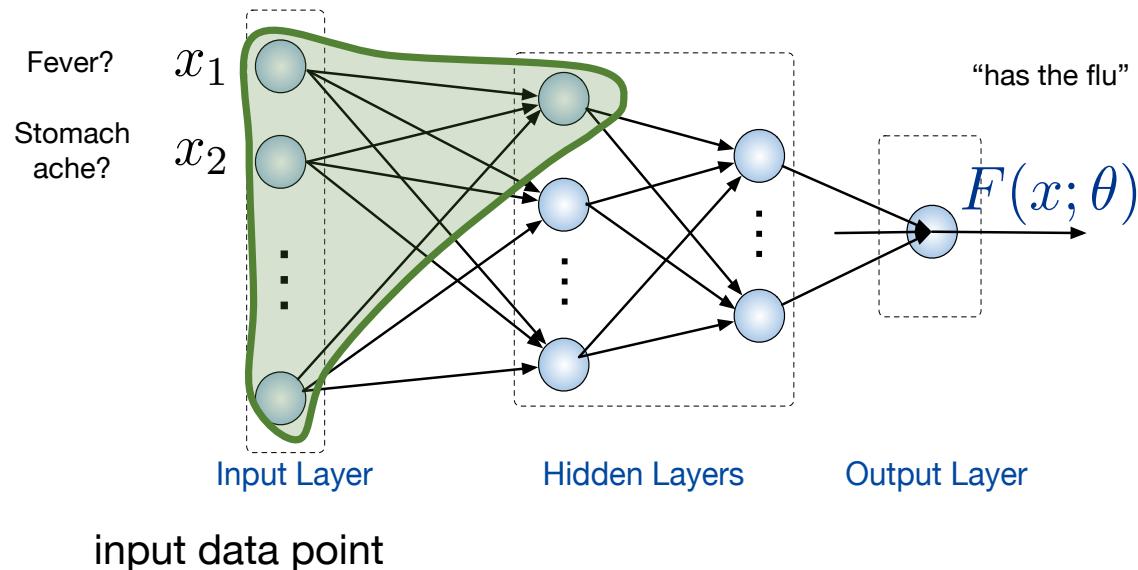


Massachusetts
Institute of
Technology

Deep Learning for Images: Convolutional Neural Networks

Stefanie Jegelka
MIT

Monday: Fully connected NN



- Hierarchy of simple “detectors”: increasingly complex decisions
- One unit = linear classifier:
- Trained by stochastic gradient descent

$$\text{output} = \text{nonlinearity} \left(\sum_{j=1}^k x_j w_j + b \right)$$

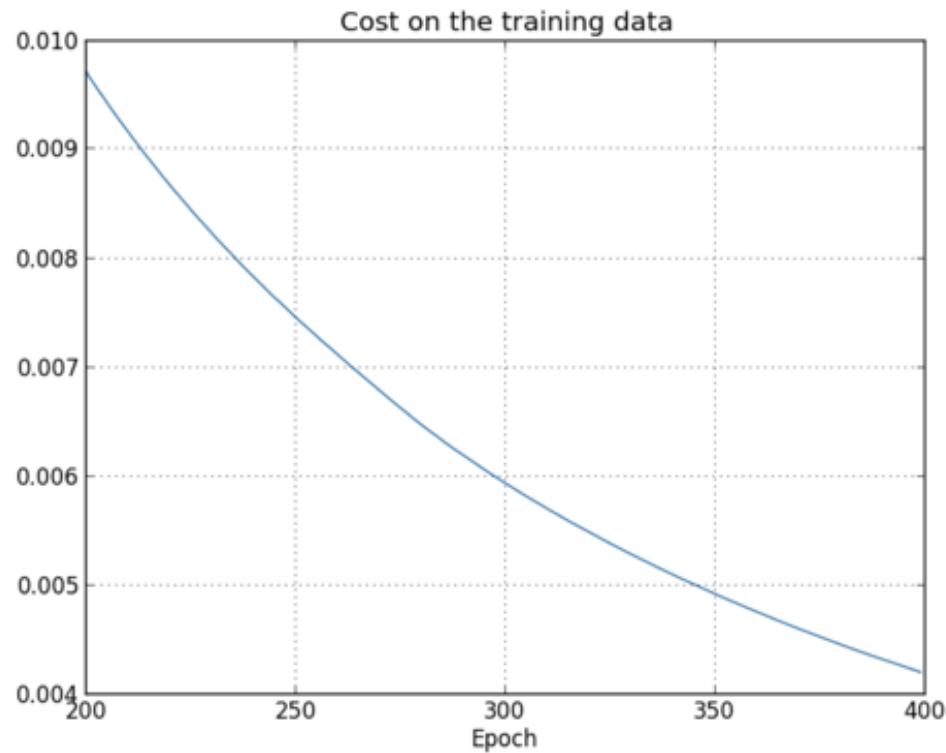
↑
e.g. ReLu, sigmoid

Learned weights and bias

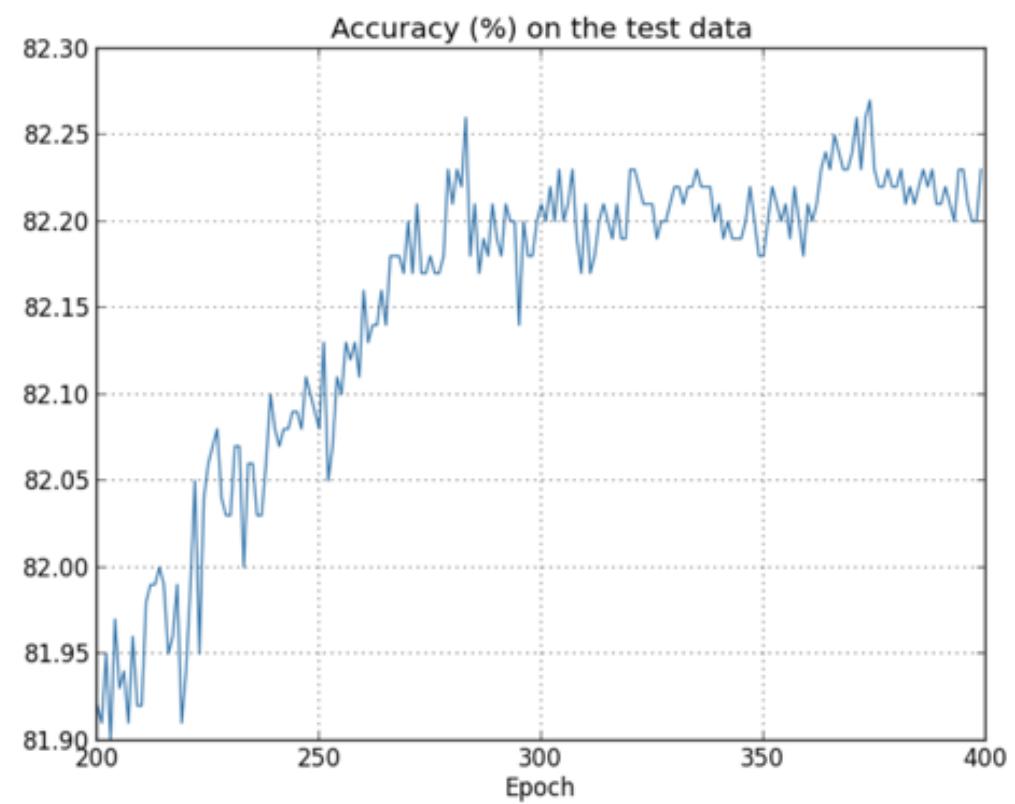
Input: data x or output of previous layer

Regularization

Loss on the training data



Accuracy on the test data



Regularization

- **Squared norm / weight decay**

add $+ \frac{\lambda}{2} \|\theta\|^2$ to training loss

$$\sum_{j=1}^d \theta_j^2$$


Regularization

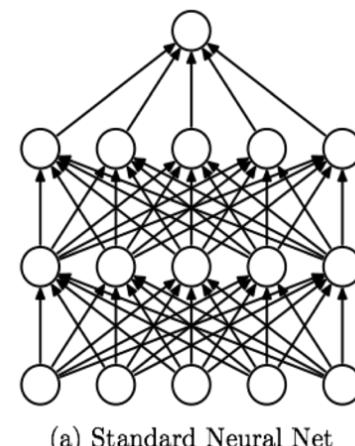
- **Squared norm / weight decay**
add $+\frac{\lambda}{2} \|\theta\|^2$ to training loss
- **Early stopping**
Check error on validation set after each epoch

Regularization

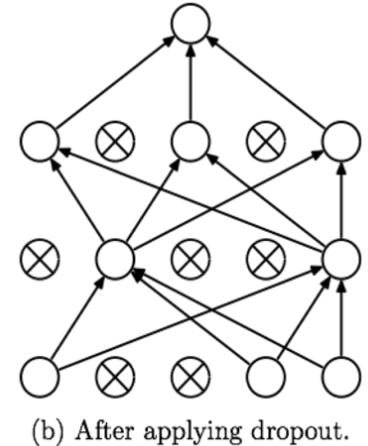
- **Squared norm / weight decay**
add $+\frac{\lambda}{2} \|\theta\|^2$ to training loss
- **Early stopping**
Check error on validation set after each epoch
- **Data augmentation**
Perturbations (rotation, noise,...), averaging

Regularization

- **Squared norm / weight decay**
add $+\frac{\lambda}{2} \|\theta\|^2$ to training loss
- **Early stopping**
Check error on validation set after each epoch
- **Data augmentation**
Perturbations (rotation, noise,...), averaging
- **Dropout**
for each training data point, randomly switch off $\frac{1}{2}$ of the units, don't update them either;
at test time, scale weights by $\frac{1}{2}$



(a) Standard Neural Net



(b) After applying dropout.

Figure: (Srivastava et al. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. JMLR, 2014)

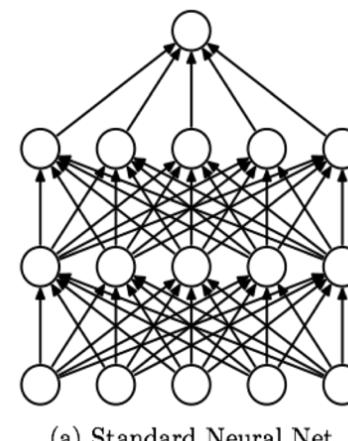
Regularization

$$\nabla(\text{loss} + \frac{\lambda}{2}\|\theta\|^2) =$$

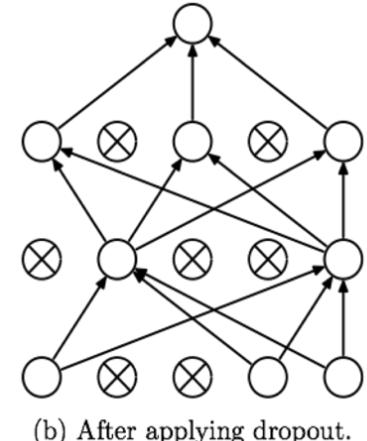
$$\nabla \text{loss} - \lambda \theta$$



- **Squared norm / weight decay**
add $+\frac{\lambda}{2}\|\theta\|^2$ to training loss
- **Early stopping**
Check error on validation set after each epoch
- **Data augmentation**
Perturbations (rotation, noise,...), averaging
- **Dropout**
for each training data point, randomly switch off $\frac{1}{2}$ of the units, don't update them either;
at test time, scale weights by $\frac{1}{2}$
- **Batch normalization**
normalize inputs of each layer



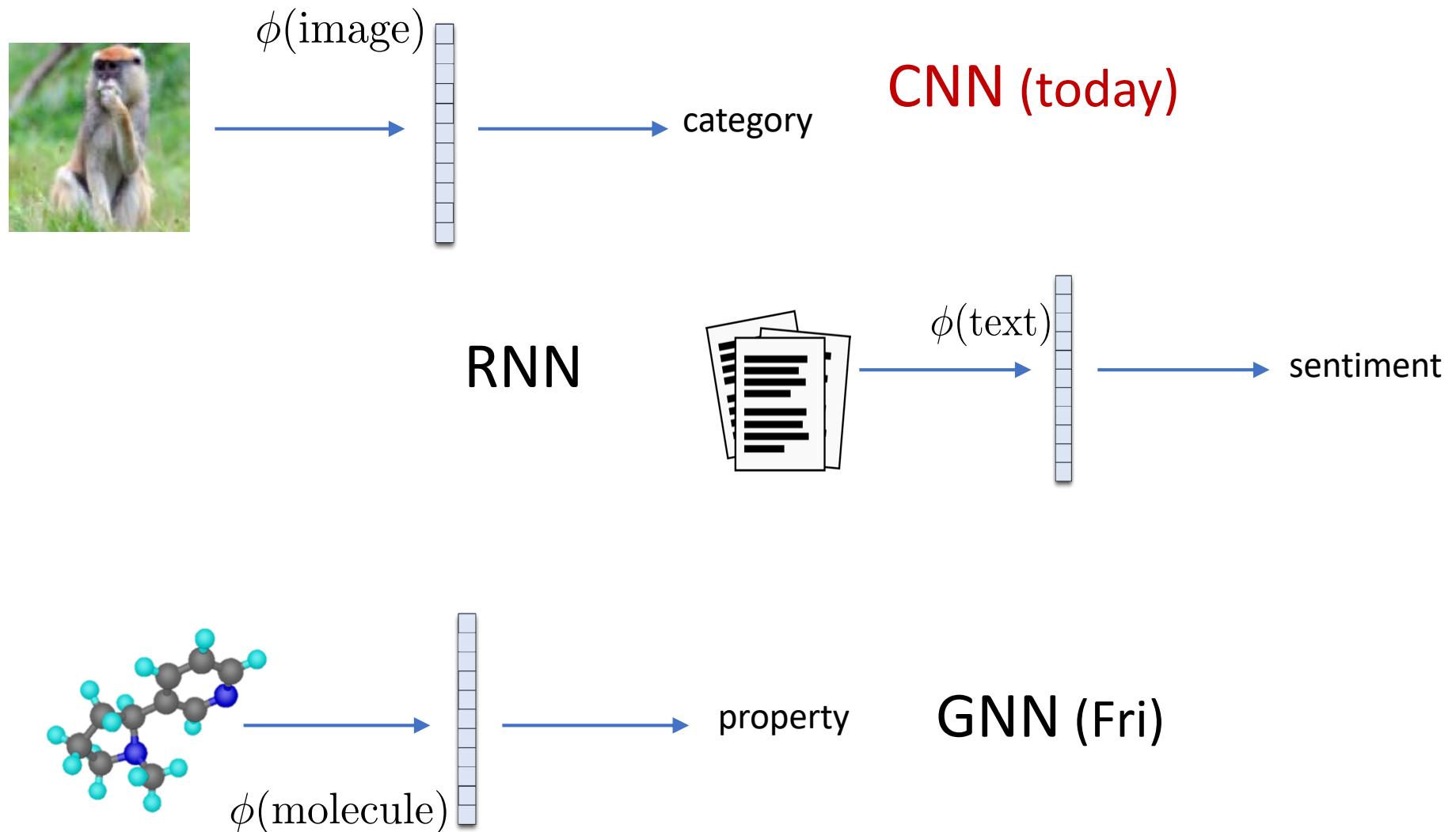
(a) Standard Neural Net



(b) After applying dropout.

Figure: (Srivastava et al. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. JMLR, 2014)

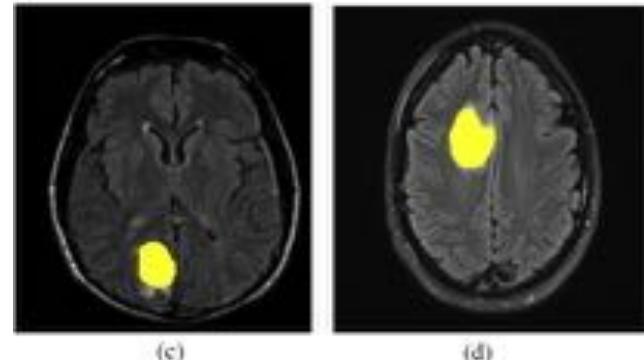
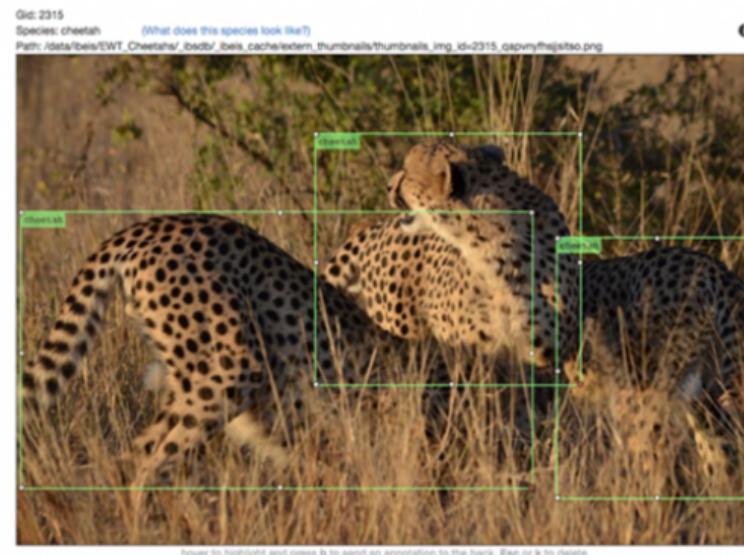
Specialized methods to encode



Outline

- **Mon:** Feedforward neural networks
- **Today:** Neural networks for images: convolutional neural networks (CNNs)
 - Convolutions & Pooling
 - What do CNNs learn?
 - Improving training
- **Thu:** Neural networks for graphs

Image Inputs

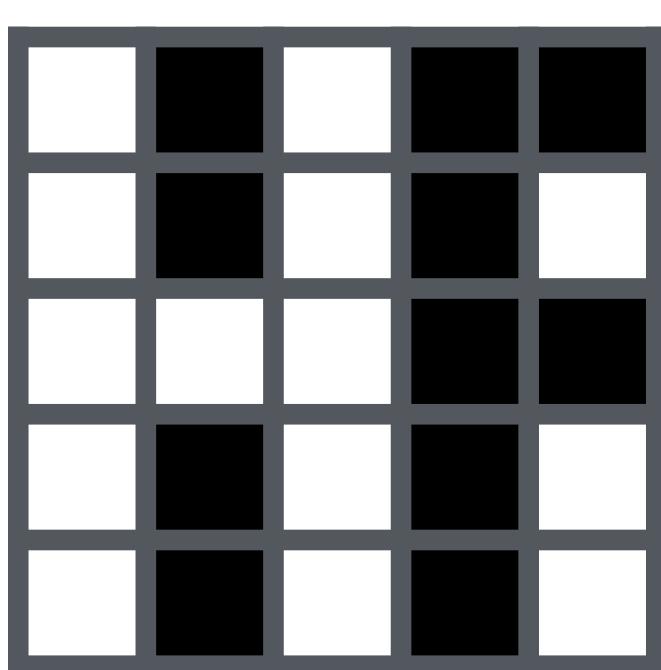


A screenshot of a Google search results page for the query "tiger". The search bar shows "tiger". Below the search bar, there are tabs for All, News, Images (which is selected), Videos, Books, More, Settings, and Tools. The main area displays several images of tigers. At the top of the image grid, there are six circular thumbnails with labels: "wallpaper", "cute", "roaring", "baby", "snow", and "ar". Below these are three rows of three images each. The first row shows a tiger walking, a tiger sitting, and a tiger standing. The second row shows a close-up of a tiger's face, a tiger being held by a person, and a tiger's face again. The third row shows a close-up of a tiger's face, a tiger's face, and a tiger's face. There are also some small text links next to the images, such as "Tiger - Wikipedia en.wikipedia.org", "Tiger | Species | WWF worldwildlife.org", and "Tiger | Smithsonian's National Zoo nationalzoo.si.edu".

Outline: CNNs

- Spatial locality & invariance
- Convolutional Neural Networks
 - Convolution and filters
 - Max Pooling
- CNN training
 - Illustrations: what do CNNs (not) learn?

Images are matrices (tensors)



1	0	1	0	0
1	0	1	0	1
1	1	1	0	0
1	0	1	0	1
1	0	1	0	1



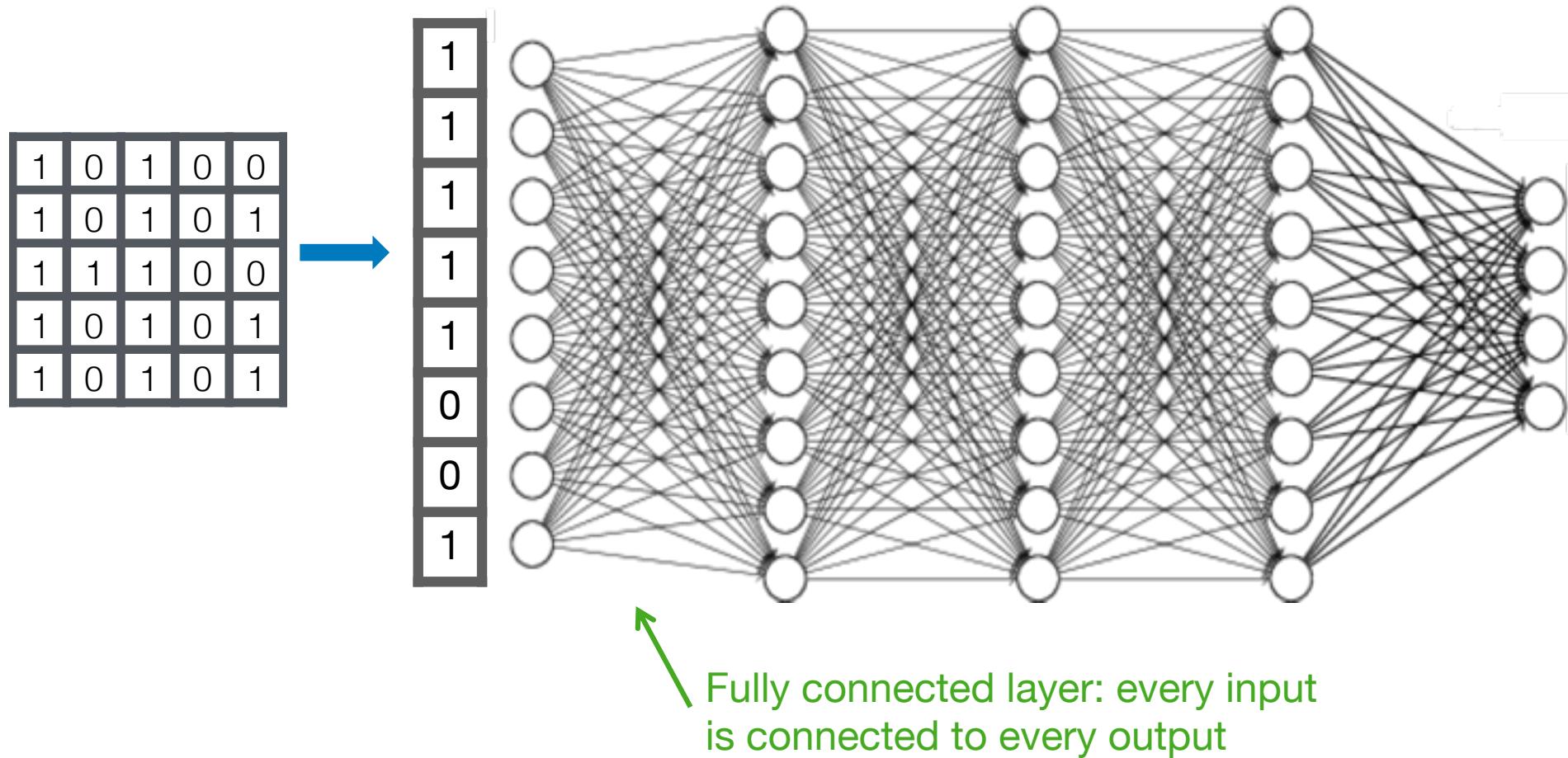
1
1
1
1
1
0
0
1
0
0

How do we input an image into a neural network?

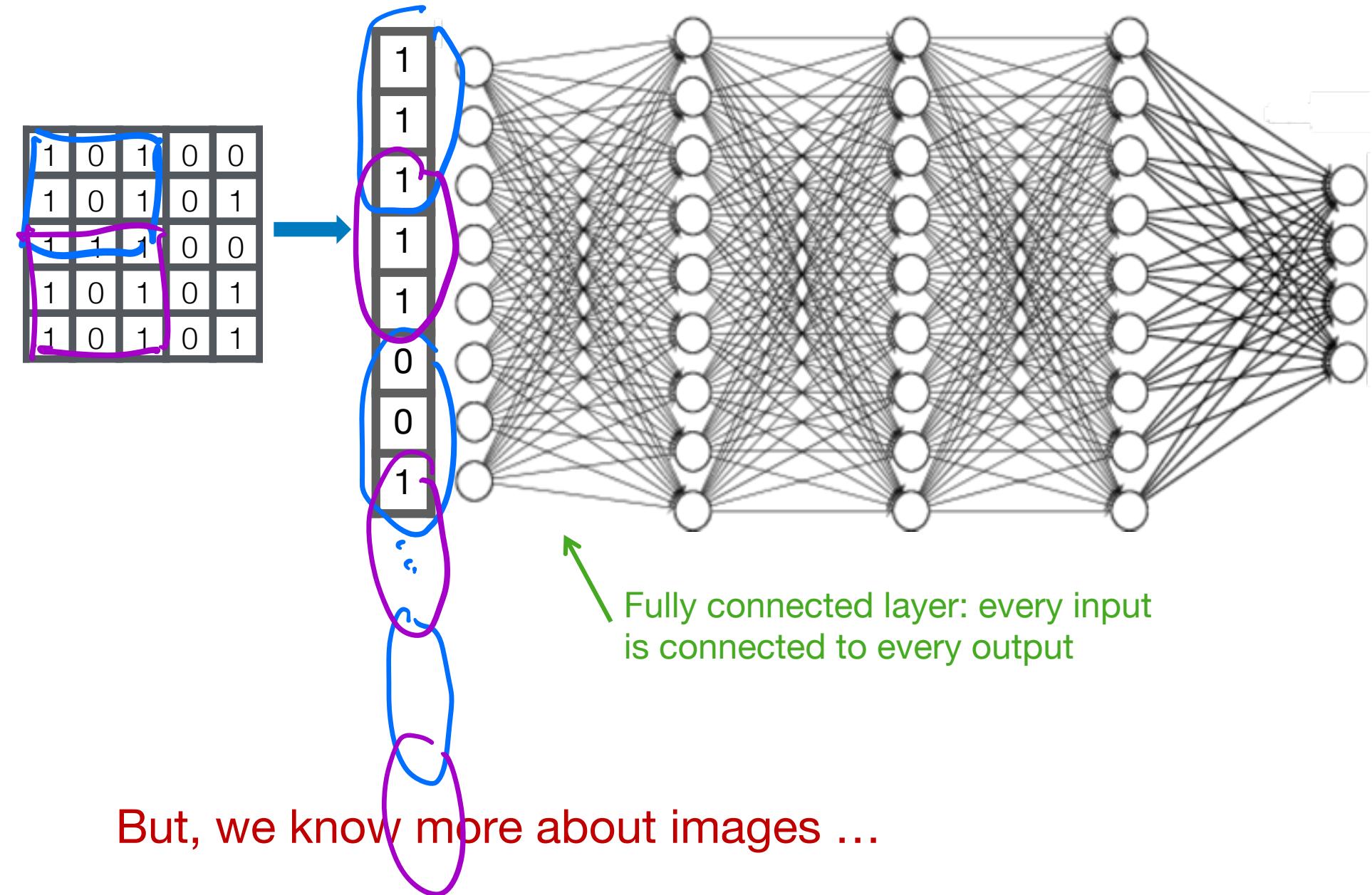
So far: make it a vector and use fully connected network

...

Fully connected network?



Fully connected network?



Locality and translation invariance



Locality and translation invariance



- Spatial locality: objects occur in local patches
- Translation invariance: a shifted cup is still a cup

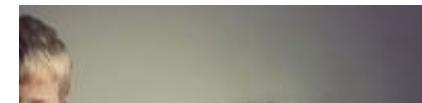
Locality and translation invariance



Locality and translation invariance



Locality and translation invariance



Idea:

1. learn a patch-size cup detector and *slide it* across the image
2. say “yes” if it fires *anywhere*
3. Hierarchy: edges -> shapes / textures -> parts -> objects -> scenes

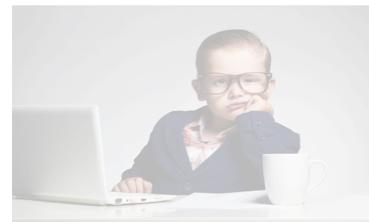
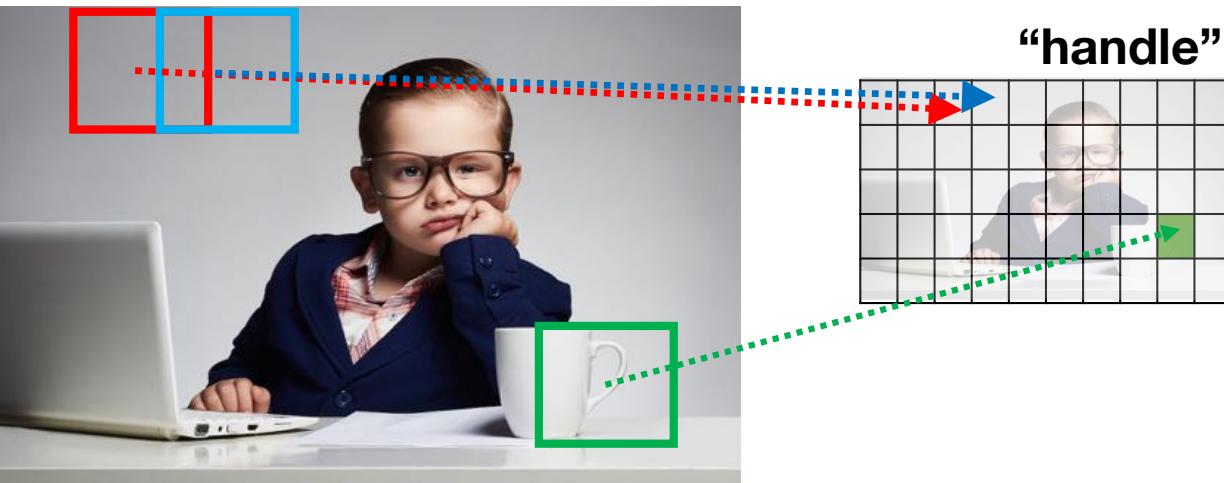
How do we implement this?



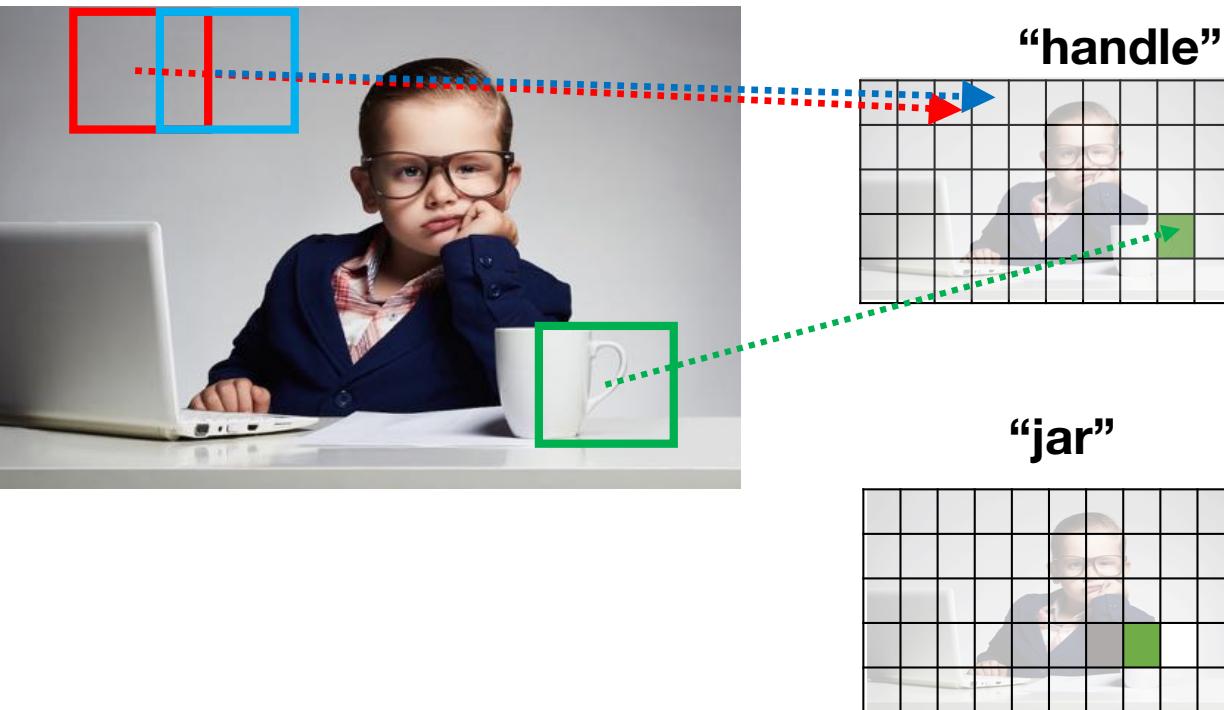
How do we implement this?

filter

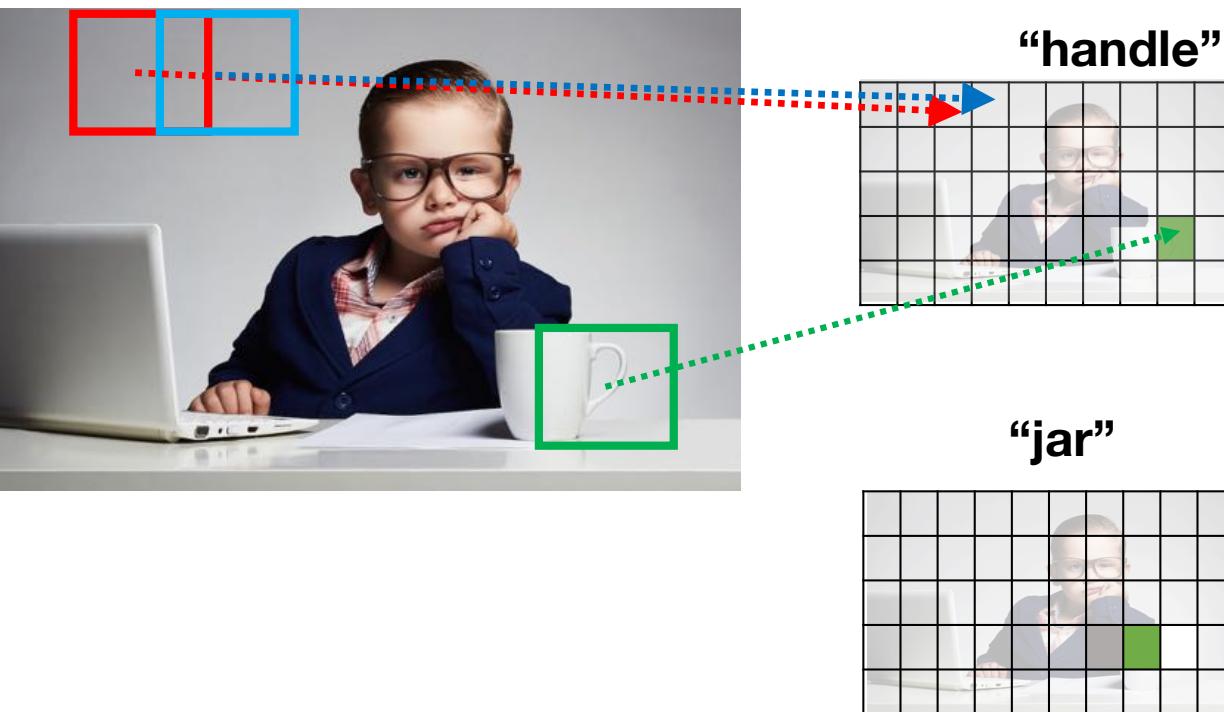
“handle”



How do we implement this?

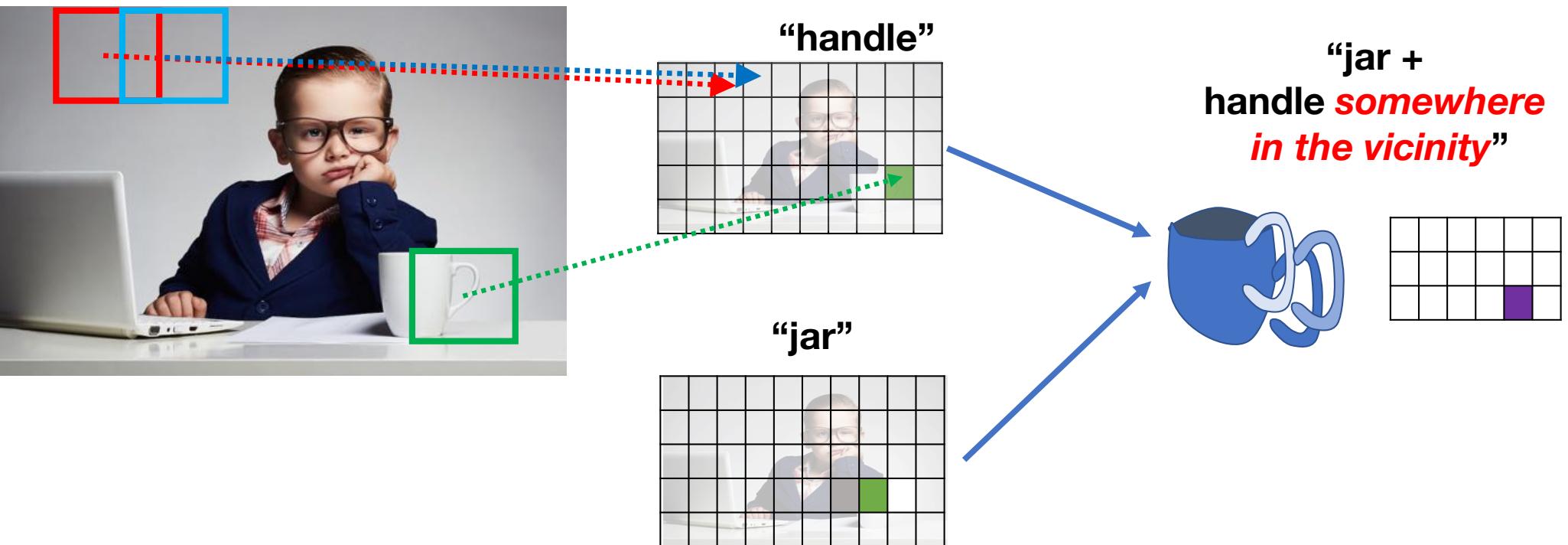


How do we implement this?



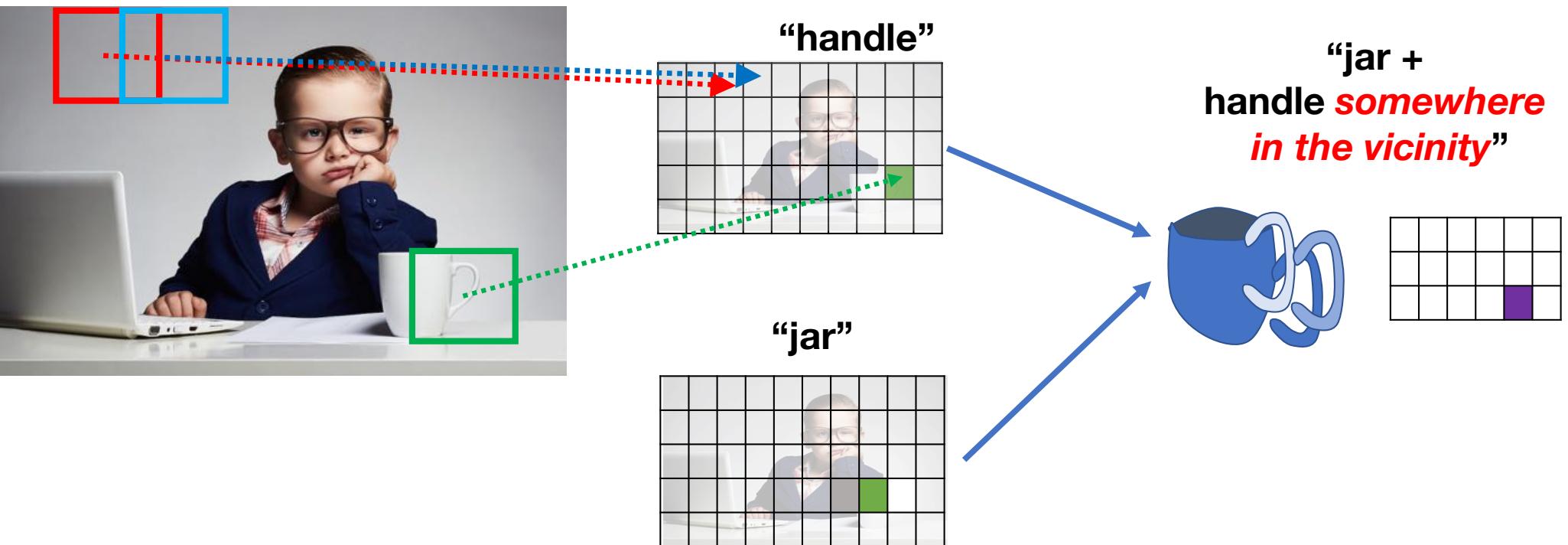
Convolution:
Slide a pattern
detector

How do we implement this?



Convolution:
Slide a pattern
detector

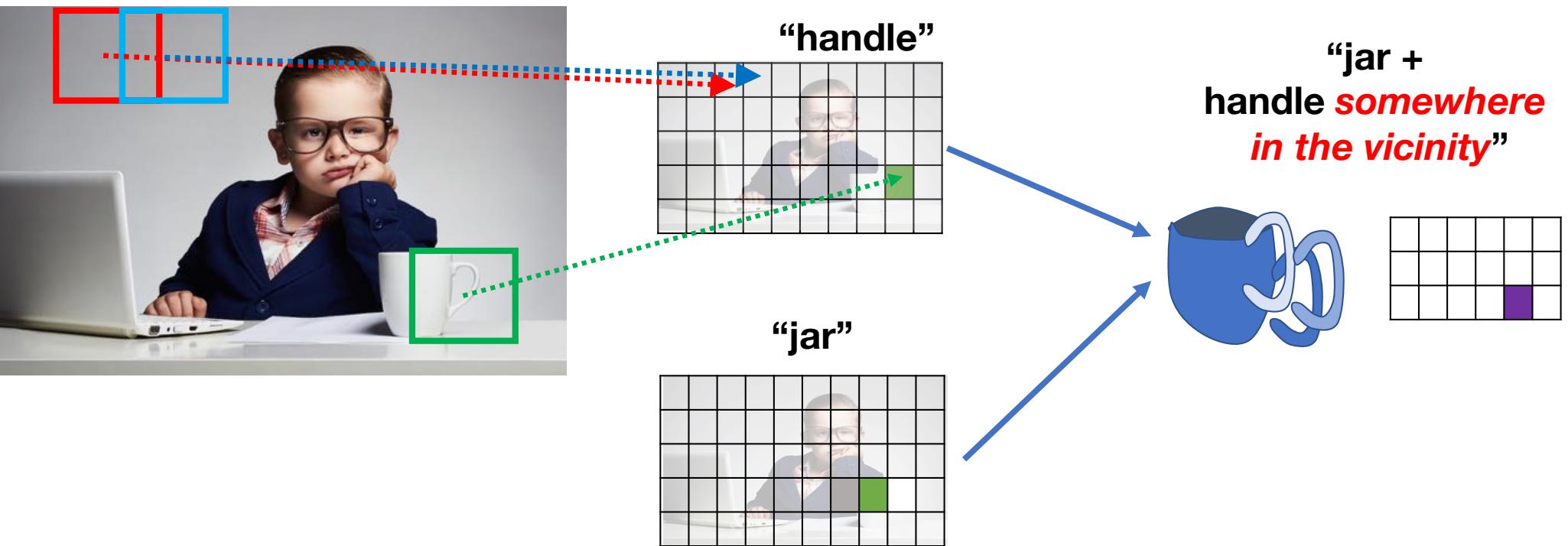
How do we implement this?



Convolution:
Slide a pattern
detector

Pooling:
“forget” exact
locations

How do we implement this?

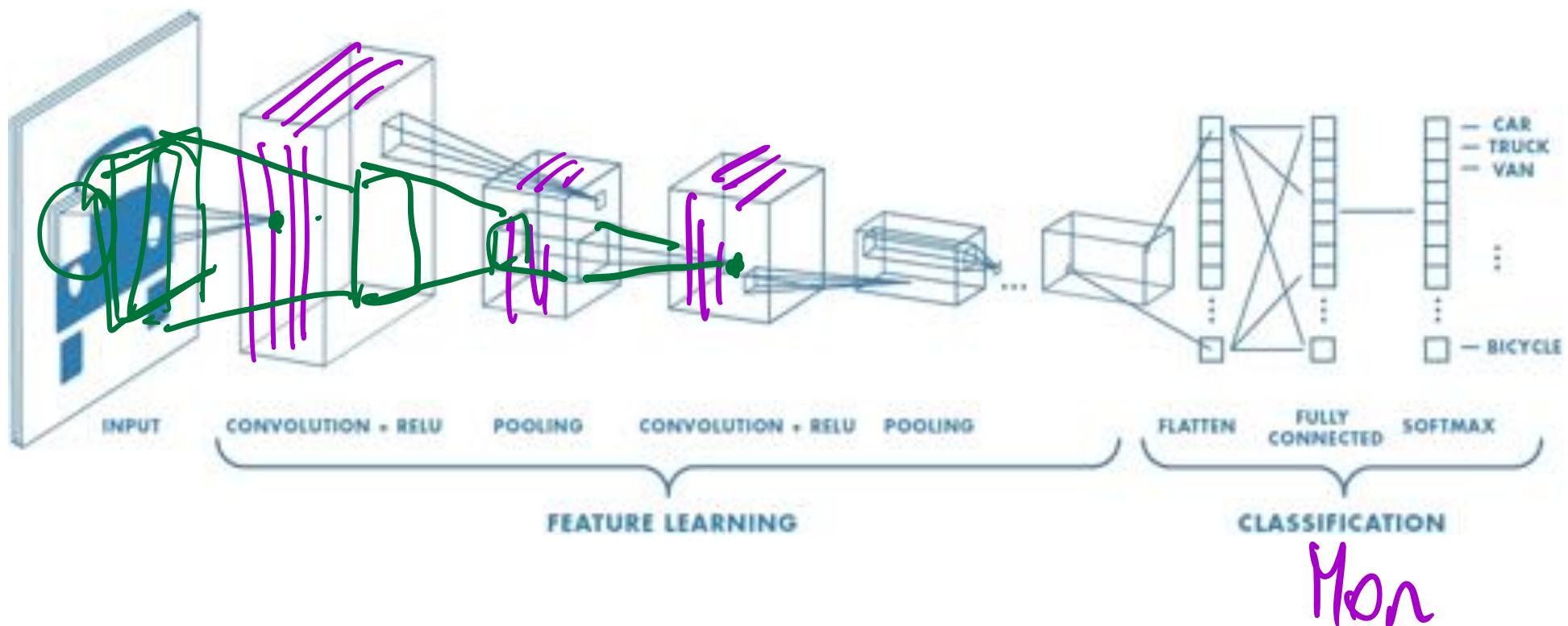


Convolution:
Slide a pattern
detector

Pooling:
“forget” exact
locations

Convolution:
Slide a pattern
detector

Hierarchy from Depth



CNNs: main ideas

- **Convolution: “local detectors”**
spatial locality

1. *learn a patch-size cup detector and slide it across the image*

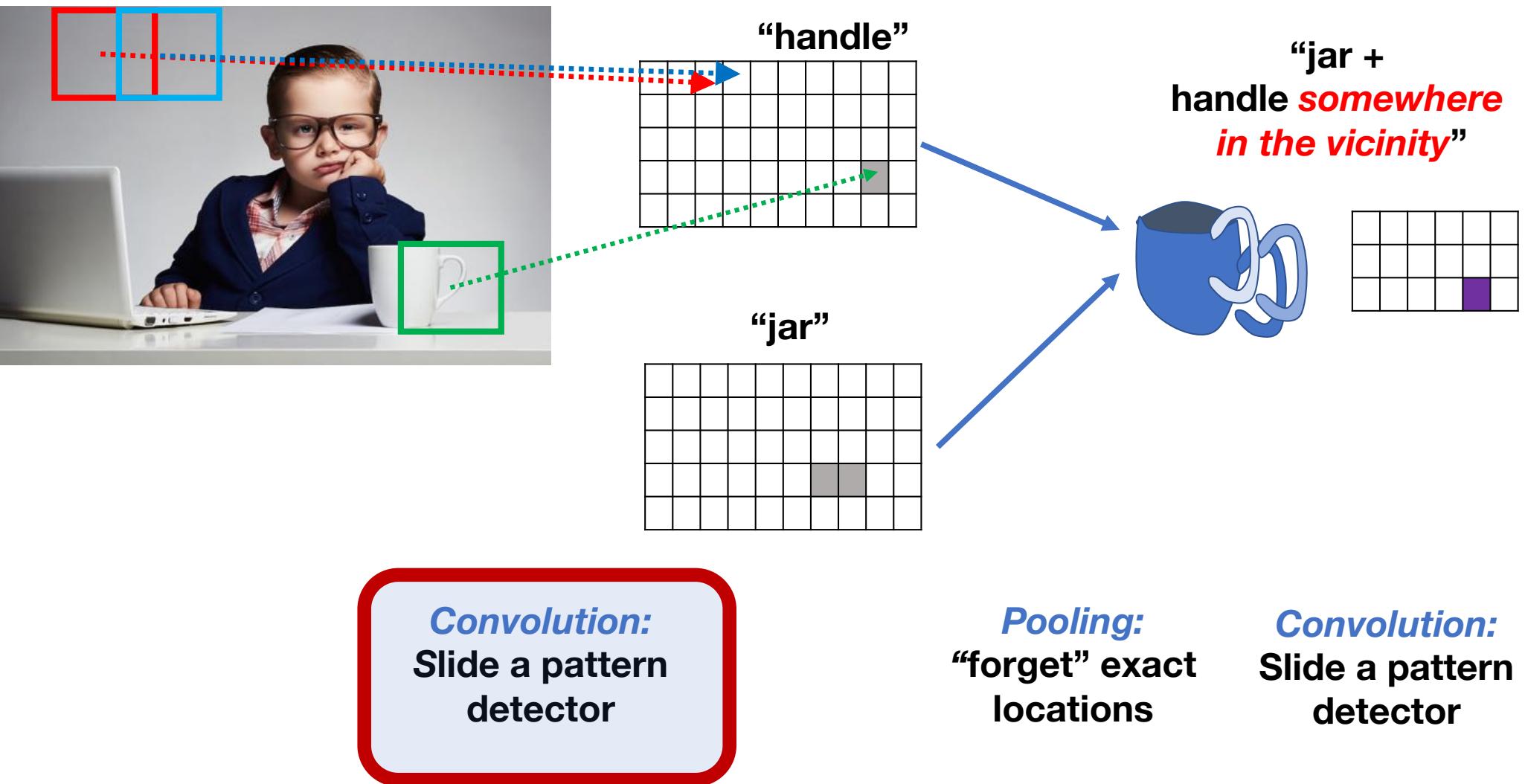
CNNs: main ideas

- **Convolution:** “local detectors”
spatial locality
 - 1. *learn a patch-size cup detector and slide it across the image*
- **Weight sharing:**
apply same detector to all image patches
efficiency, translation invariance

CNNs: main ideas

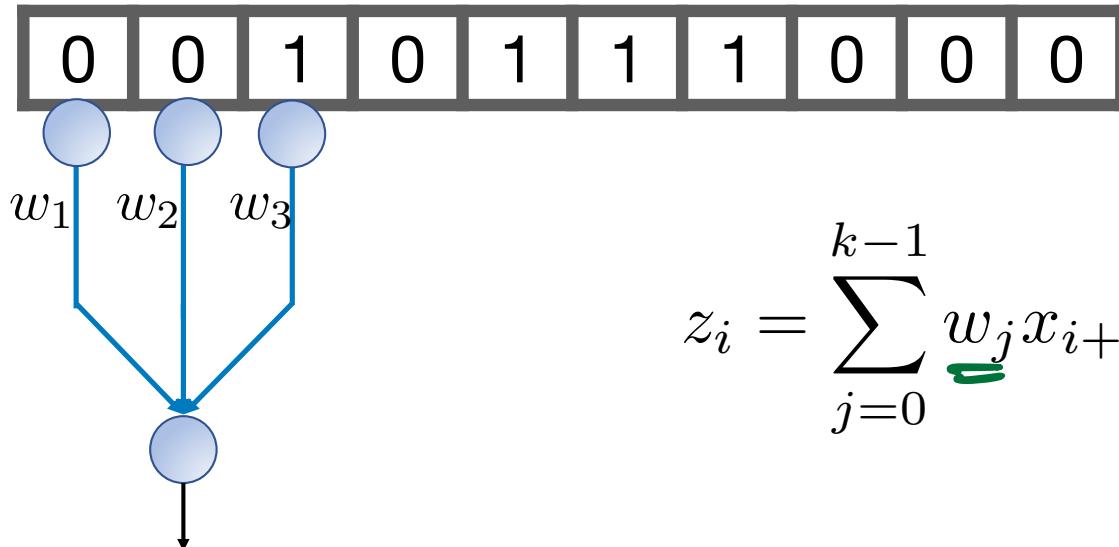
- **Convolution:** “local detectors”
spatial locality
 - 1. *learn a patch-size cup detector and slide it across the image*
- **Weight sharing:**
apply same detector to all image patches
efficiency, translation invariance
- **Pooling**
abstract away locality
 - 2. *say “yes” if it fires anywhere in a region*

How do we implement this?



Convolutional Layer: 1D example

- 1D image
- Filter
(weights, bias)
- After
convolution



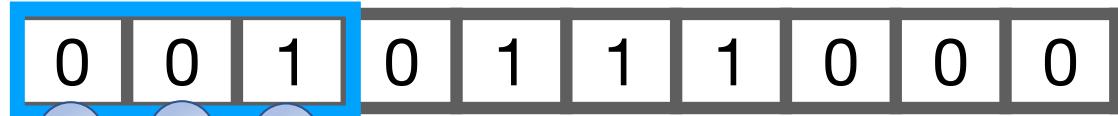
$$z_i = \sum_{j=0}^{k-1} w_j x_{i+j} + b_i$$

A **filter** is like a
locally connected neuron

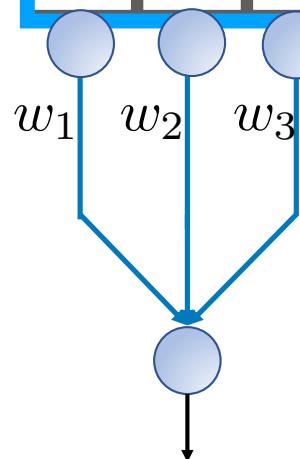
(here, will use bias=0)

Convolutional Layer: 1D example

- 1D image



- Filter
(weights, bias)
- After
convolution



$$z_i = \sum_{j=0}^{k-1} w_j x_{i+j}$$

A **filter** is like a
locally connected neuron

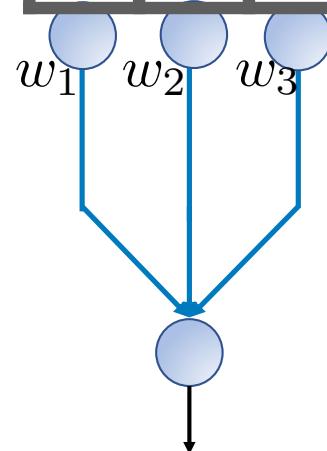
(here, will use bias=0)

Convolutional Layer: 1D example

- 1D image



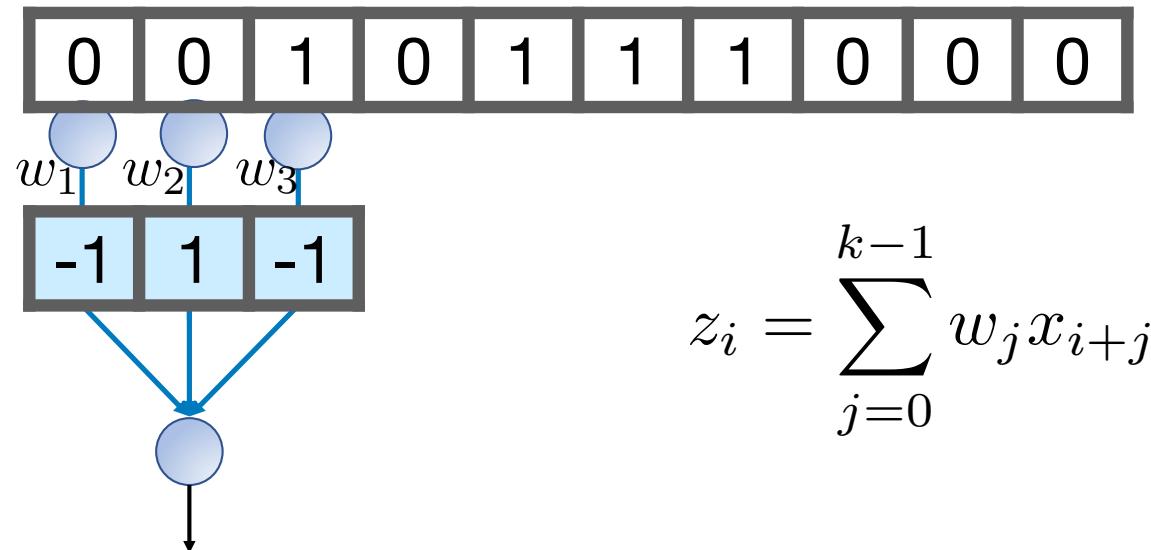
- Filter
(weights, bias)
- After
convolution



$$z_i = \sum_{j=0}^{k-1} w_j x_{i+j}$$

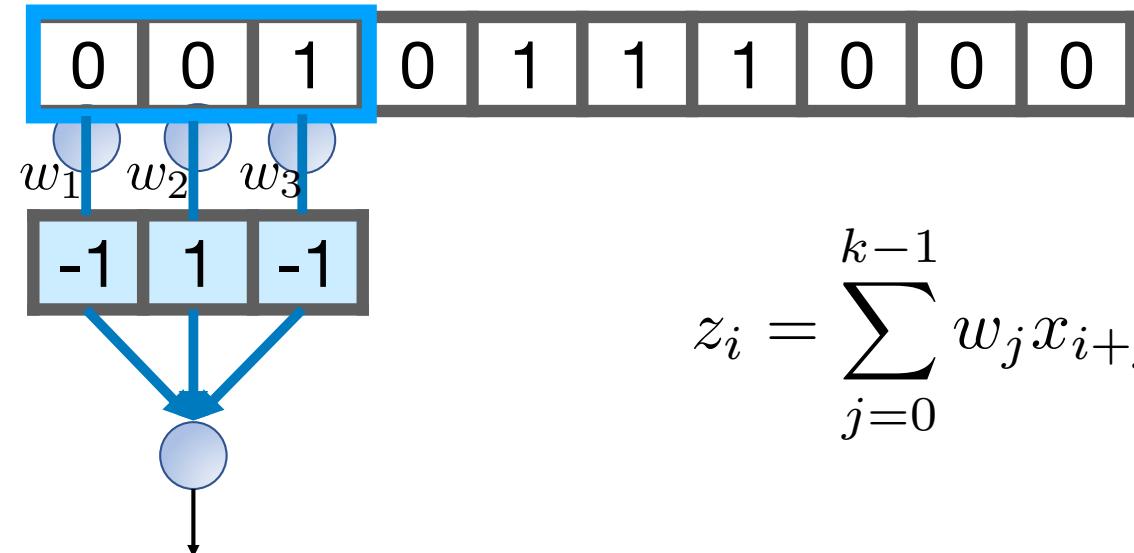
Convolutional Layer: 1D example

- 1D image
- Filter
(weights, bias)
- After
convolution



Convolutional Layer: 1D example

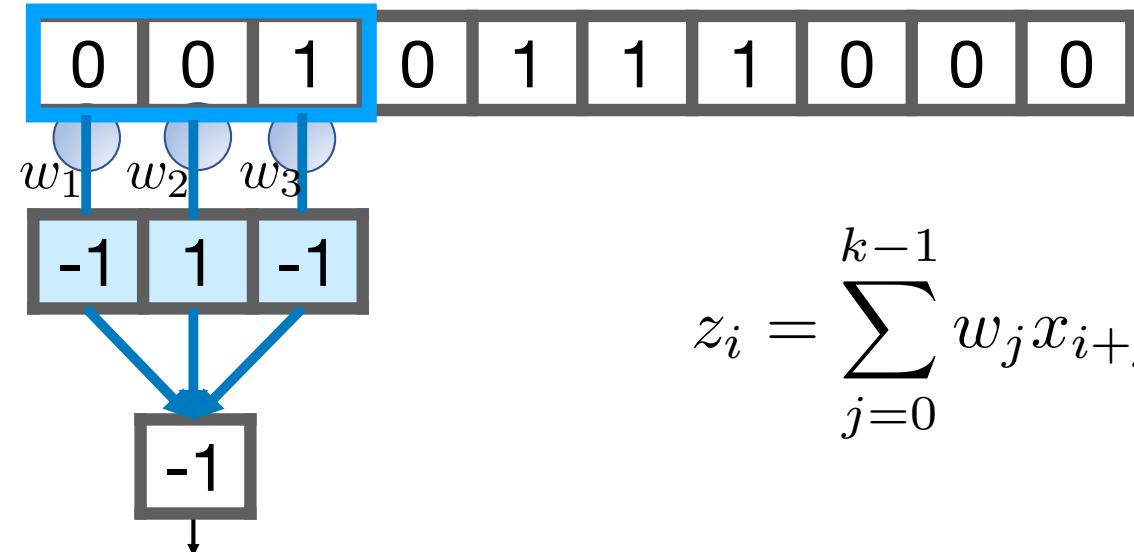
- 1D image
- Filter (*weights, bias*)
- After convolution



$$z_i = \sum_{j=0}^{k-1} w_j x_{i+j}$$

Convolutional Layer: 1D example

- 1D image
- Filter (*weights, bias*)
- After convolution

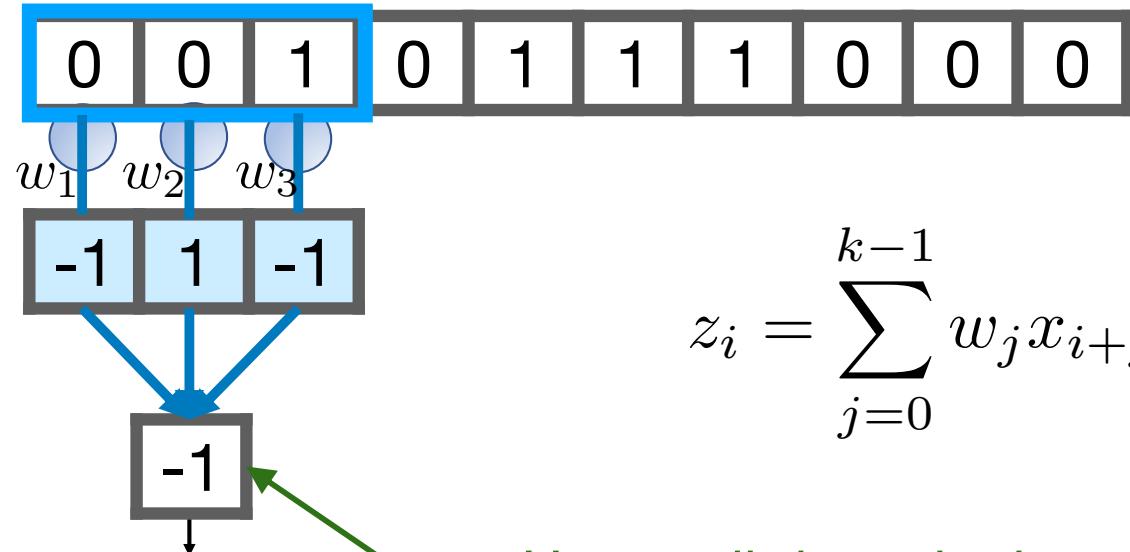


$$z_i = \sum_{j=0}^{k-1} w_j x_{i+j}$$

Example: $(-1) \cdot 0 + 1 \cdot 0 + (-1) \cdot 1 = -1$

Convolutional Layer: 1D example

- 1D image
- Filter (*weights, bias*)
- After convolution



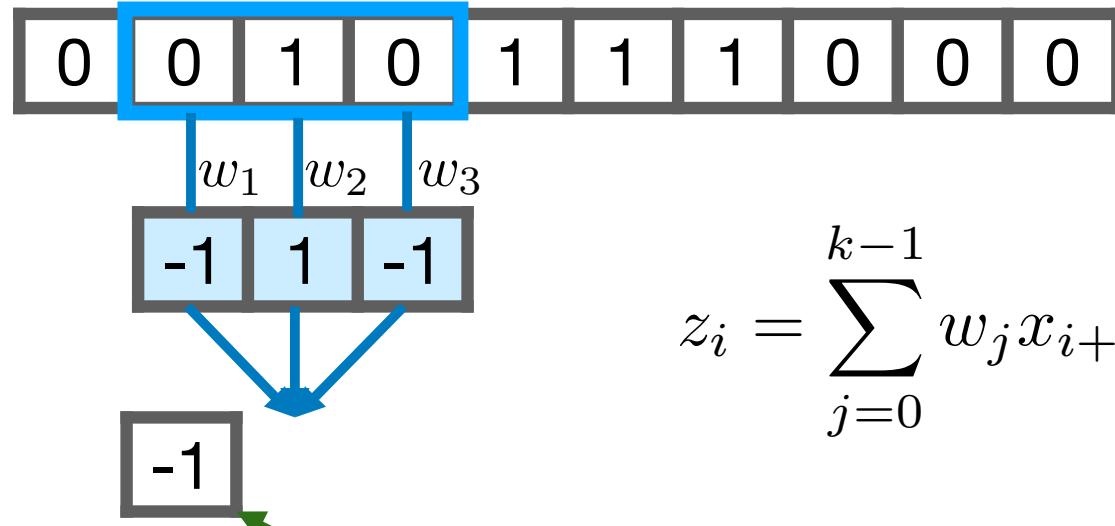
$$z_i = \sum_{j=0}^{k-1} w_j x_{i+j}$$

*How well does the input match
the filter pattern (low-high-low)?
larger == better*

Example: $(-1) \cdot 0 + 1 \cdot 0 + (-1) \cdot 1 = -1$

Convolutional Layer: 1D example

- 1D image
- Filter
- After convolution

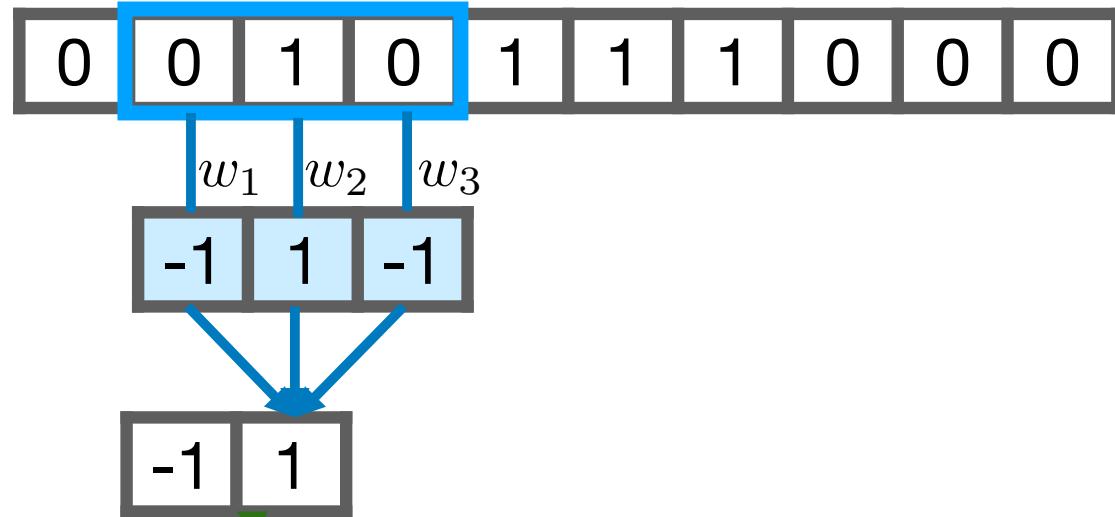


$$z_i = \sum_{j=0}^{k-1} w_j x_{i+j}$$

*How well does the input match
the filter pattern (low-high-low)?
larger == better*

Convolutional Layer: 1D example

- 1D image



- Filter

- After convolution

*How well does the input match the filter pattern (low-high-low)?
larger == better*

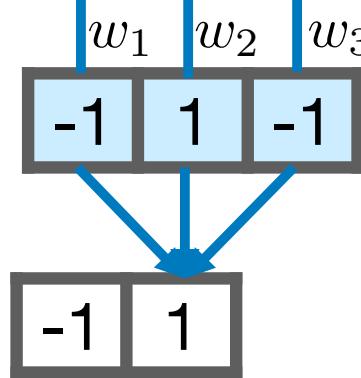
Example: $(-1) \cdot 0 + 1 \cdot 1 + (-1) \cdot 0 = 1$

Convolutional Layer: 1D example

- 1D image



- Filter



- After convolution

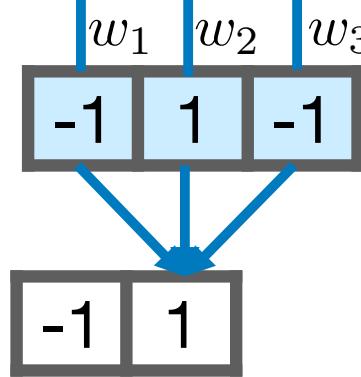
Stride: by how much do we shift the filter

Convolutional Layer: 1D example

- 1D image



- Filter



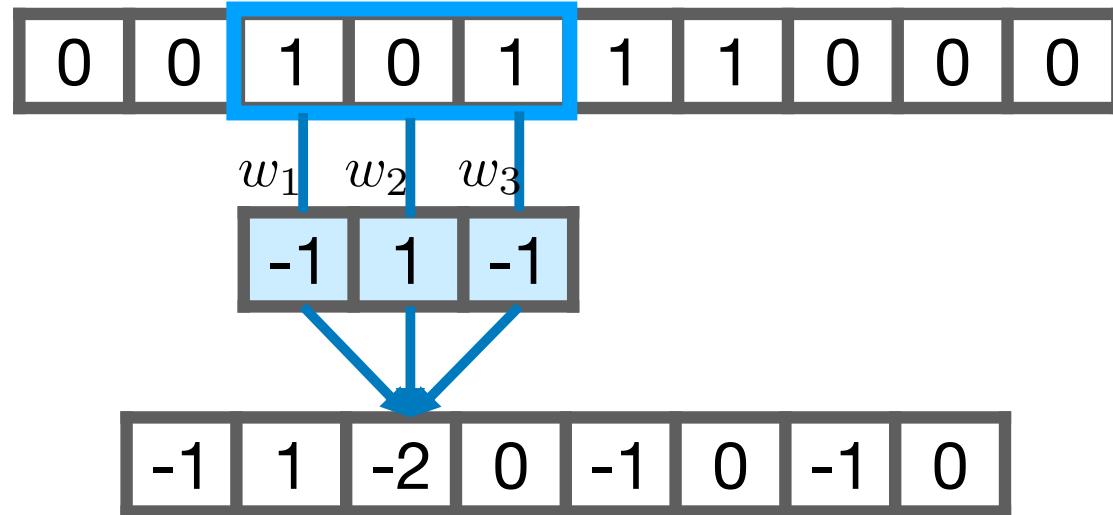
- After convolution

Stride: by how much do we shift the filter

Weight sharing: same filter applied to both (all) patches

Convolutional Layer: 1D example

- 1D image



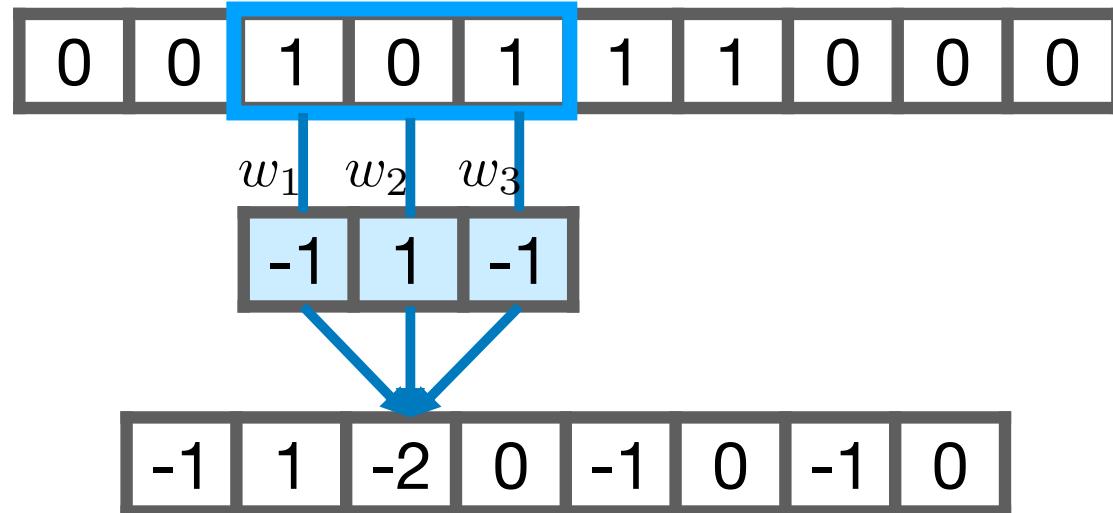
- Filter

- After convolution

Big advantage: due to **weight sharing**, needs much **fewer weights than a fully connected network**

Convolutional Layer: 1D example

- 1D image



- Filter

- After convolution

Big advantage: due to **weight sharing**, needs much **fewer weights than a fully connected network**

Example:

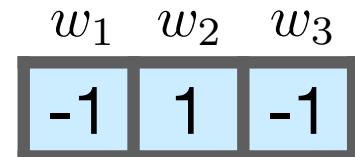
1. CNN: 5x5filters -> 26 parameters per filter (with bias).
For 20 filters: **520** parameters per layer – independent of input size!
2. Fully connected, 28x28 inputs, 30 hidden units:
 $28 \times 28 \times 30 + 30 = \mathbf{23,550}$ parameters for one layer

Convolutional Layer: 1D example

- 1D image



- Filter



- After convolution



- After ReLu



$$\text{ReLU}(x) = \max\{0, x\}$$

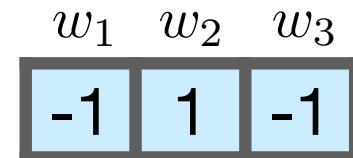
$$= \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{otherwise} \end{cases}$$

Convolutional Layer: 1D example

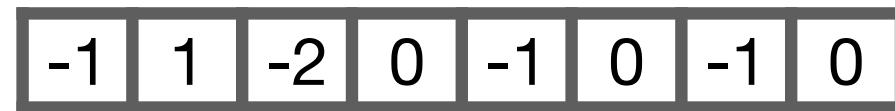
- 1D image



- Filter



- After convolution



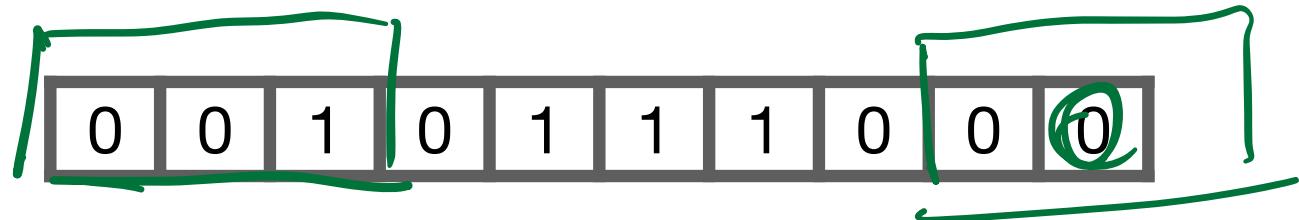
- After ReLu



Detection for “low-high-low” at position 2

Padding

- 1D image



- Filter



- After convolution



- After ReLu



Detection for “low-high-low” at position 2

Output is smaller! (why?)

Padding

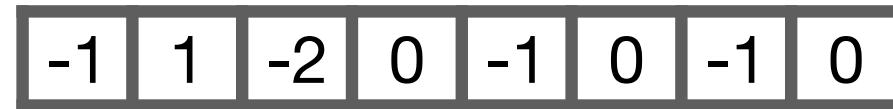
- 1D image



- Filter



- After convolution



- After ReLu



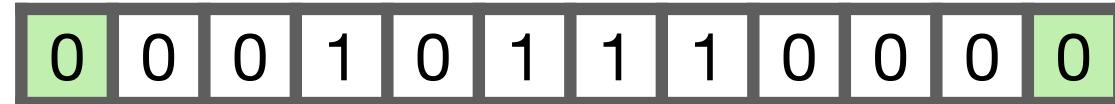
Detection for “low-high-low” at position 2

Output is smaller! (why?)

Remedy: pad input with zeros

Padding

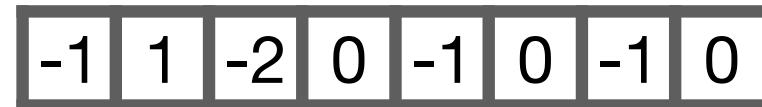
- 1D image



- Filter



- After convolution



- After ReLu



Output is smaller! (why?)

Remedy: pad input with zeros

Padding

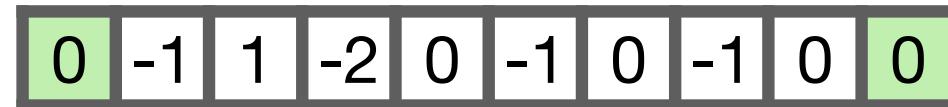
- 1D image



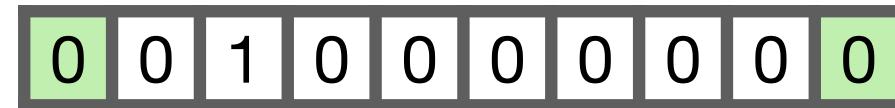
- Filter



- After convolution



- After ReLu

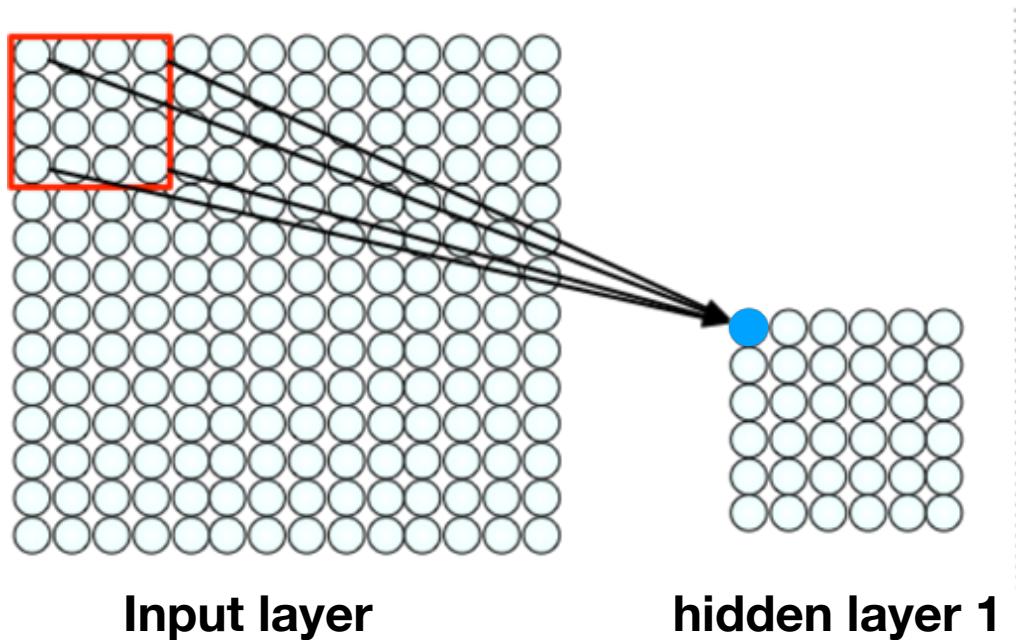


Output is smaller! (why?)

Remedy: pad input with zeros

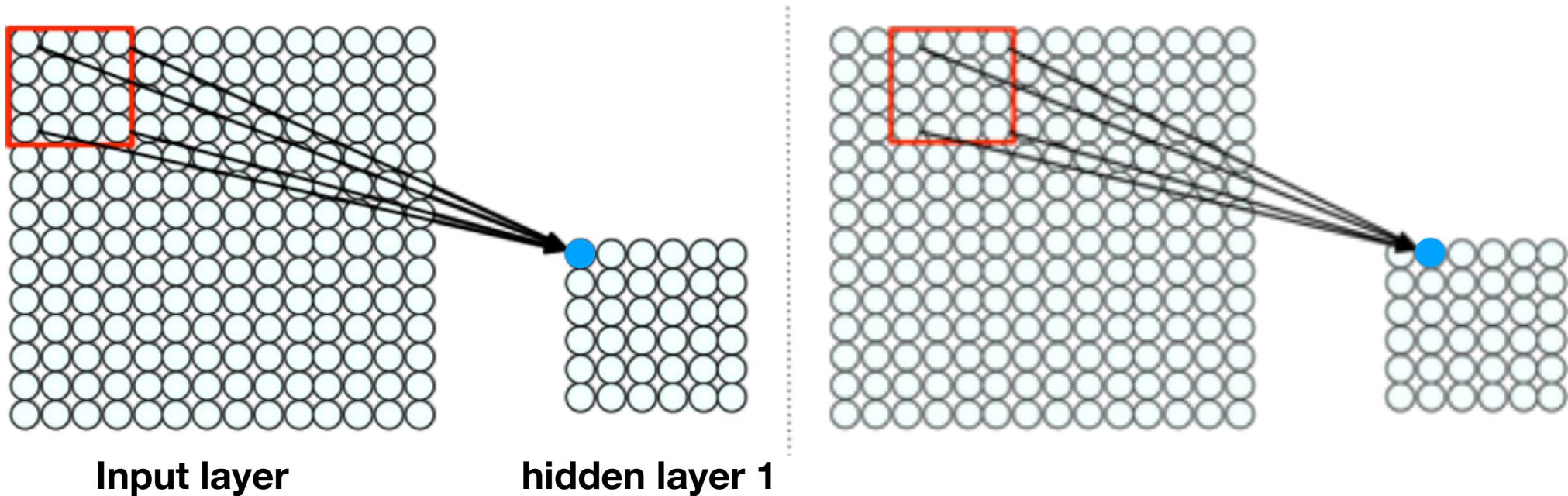
2D Convolution

- In 2D: shift filter in 2 directions (right, down)
- Output is an “image” too
- Example of a 2D convolution, 4x4 filter, stride 2



2D Convolution

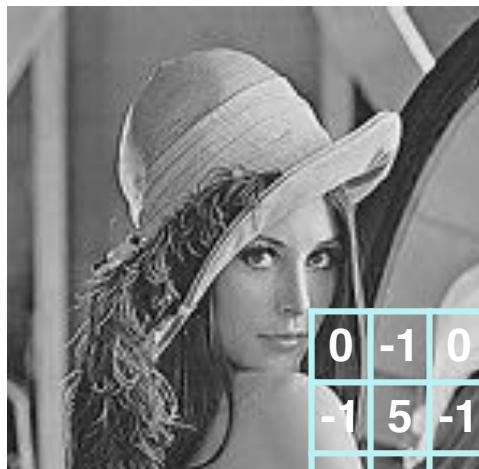
- In 2D: shift filter in 2 directions (right, down)
- Output is an “image” too
- Example of a 2D convolution, 4x4 filter, stride 2



Examples of Convolutions



orig



sharpen

$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & 1 & 0 \end{bmatrix}$$



edge detect

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$



strong edges

$$\begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix}$$

Learned filters (1st layer)

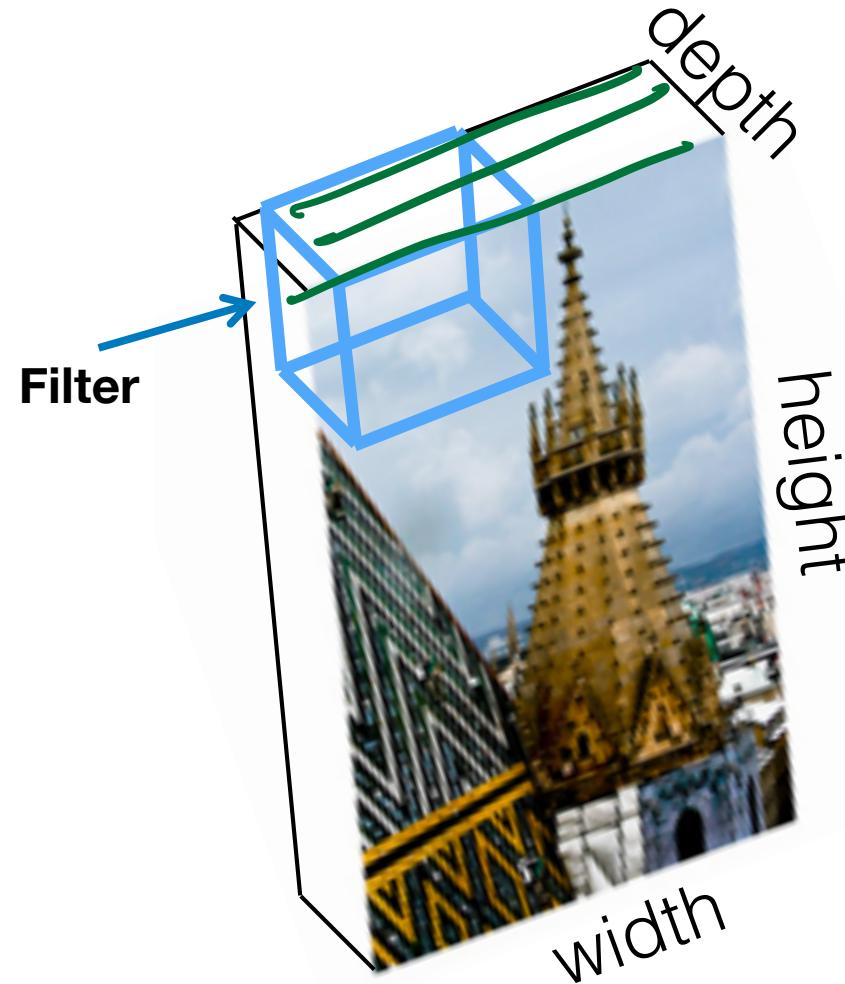


Size: $11 \times 11 \times 3$

First layer learns edge detectors!

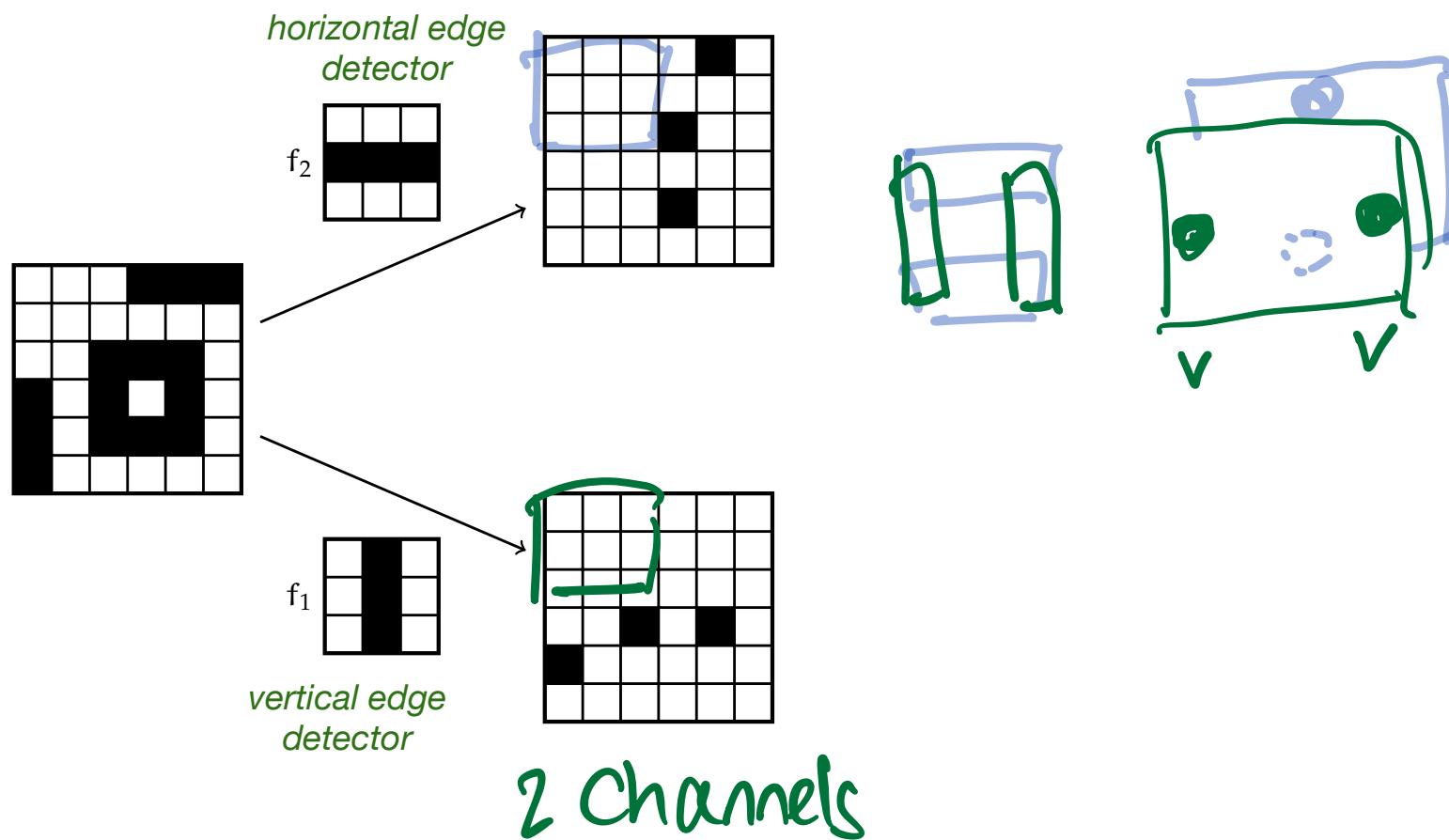
3D convolution

- Color image is actually a tensor (3 matrices): RGB
- Filter is a tensor (generalization of matrix) too



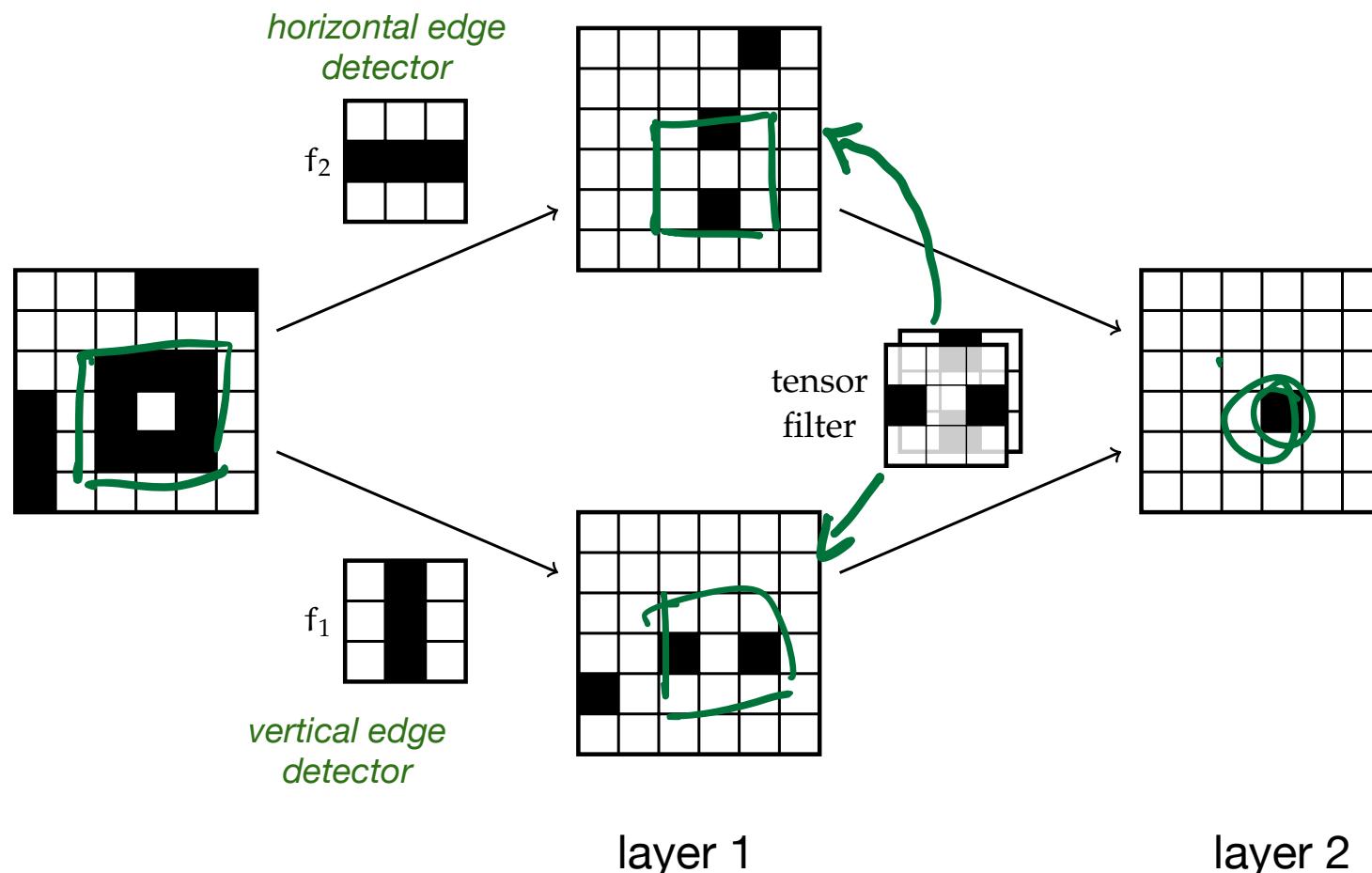
Multiple channels/filters

- Typically, we do not only want to learn one edge detector:
shape = specific combination of multiple types of edges
- Use **multiple filters** (*filter bank*) → multiple output “images”

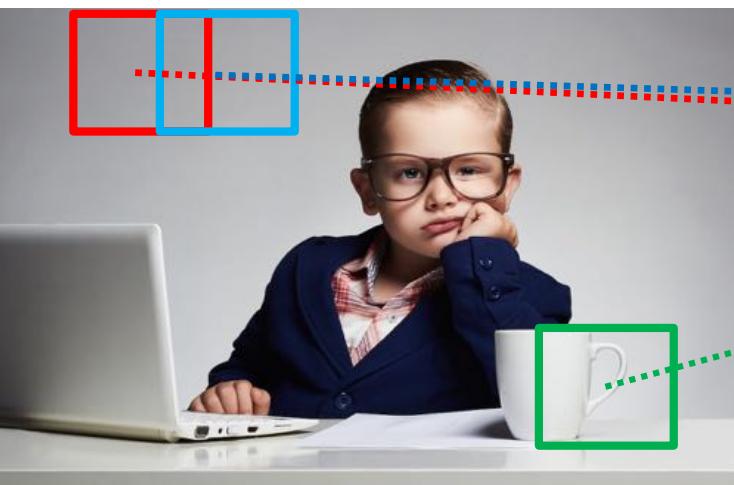


Multi-channel convolution

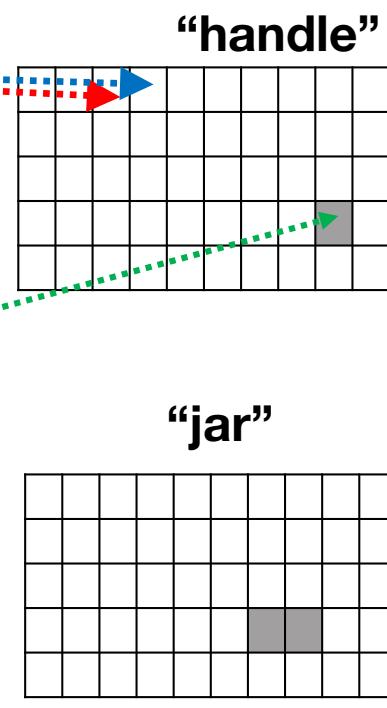
- Next layer can combine inputs from multiple channels:
3D (tensor) convolution
e.g. edges \rightarrow shapes



How do we implement this?



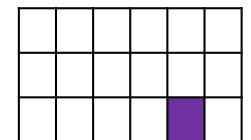
Convolution:
Slide a pattern
detector



channels

Pooling:
“forget” exact
locations

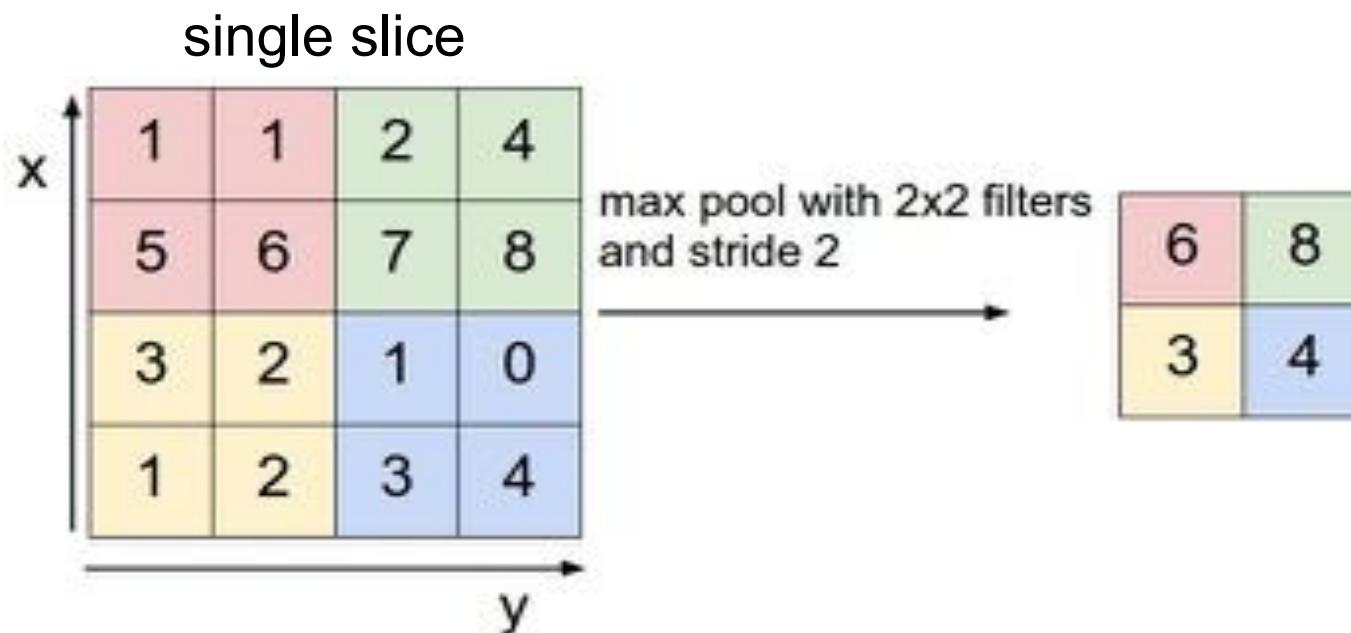
“jar +
handle **somewhere**
in the vicinity”



Convolution:
Slide a pattern
detector

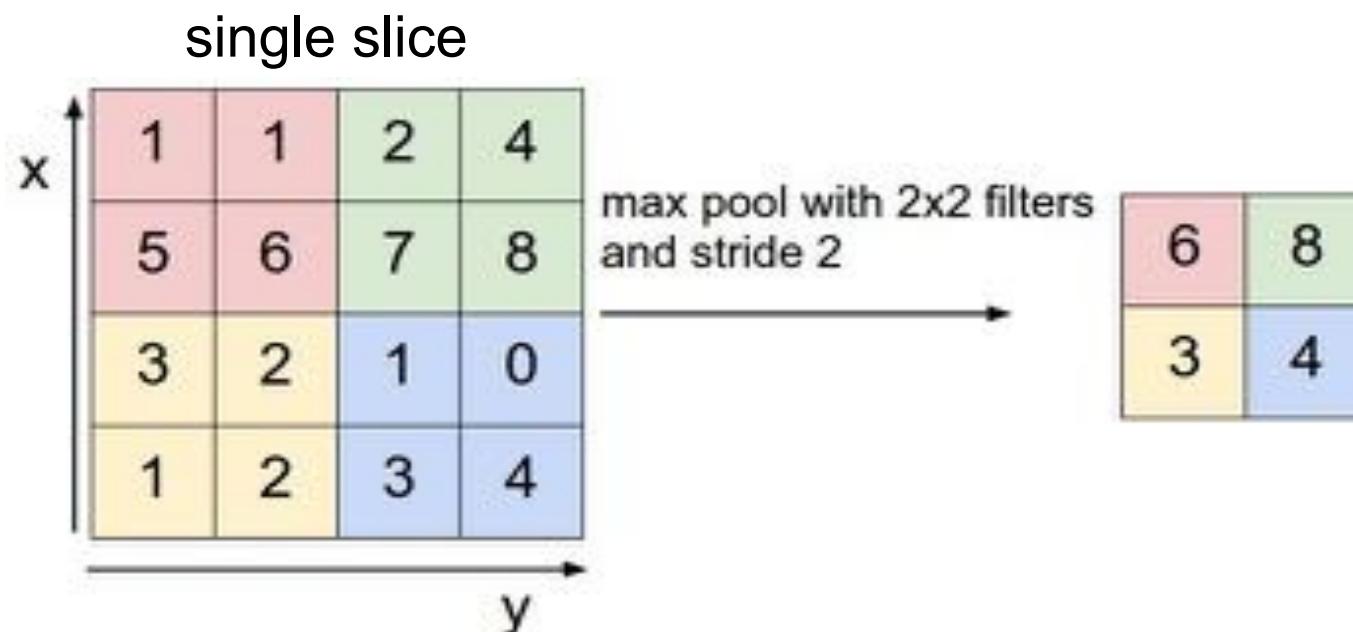
Max Pooling

- Similar to filtering, but output the maximum entry in the patch, instead of a weighted sum



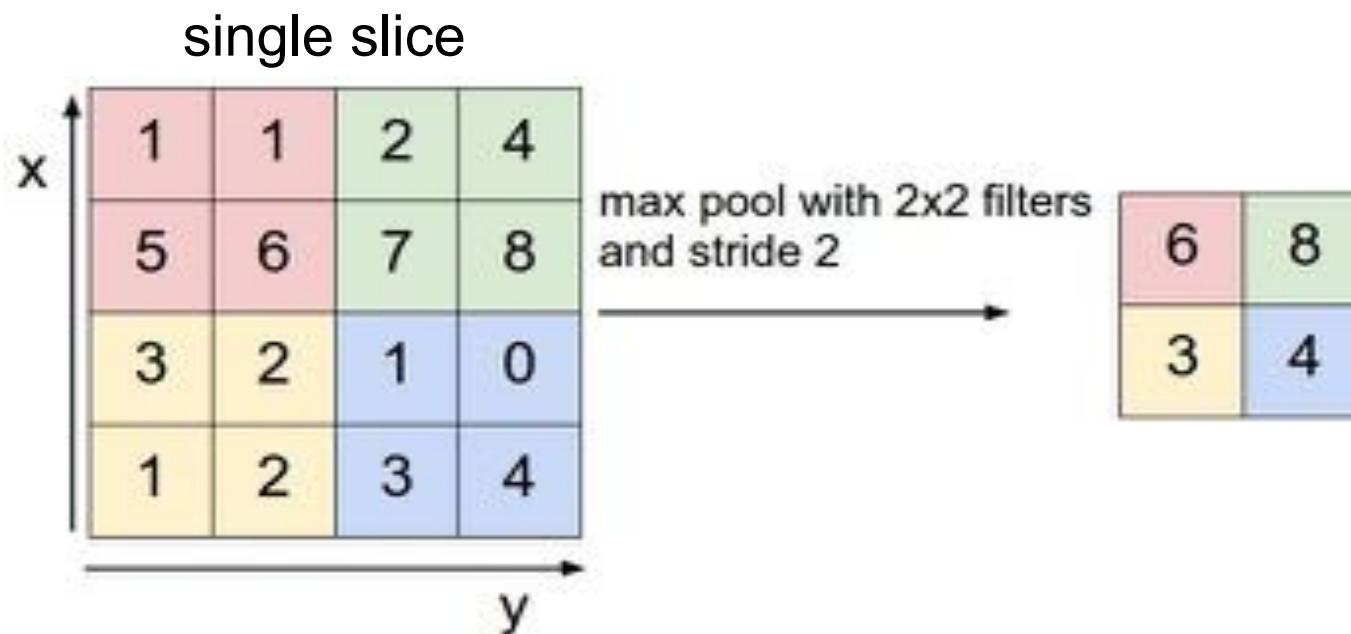
Max Pooling

- Similar to filtering, but output the maximum entry in the patch, instead of a weighted sum
- Typically output image is smaller than input → filters look at successively larger areas in the input image



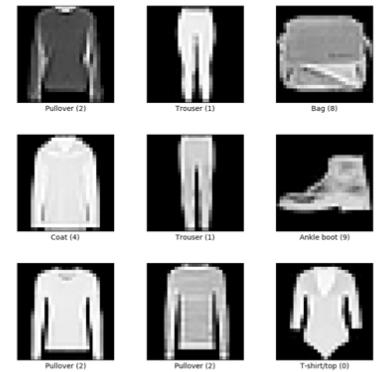
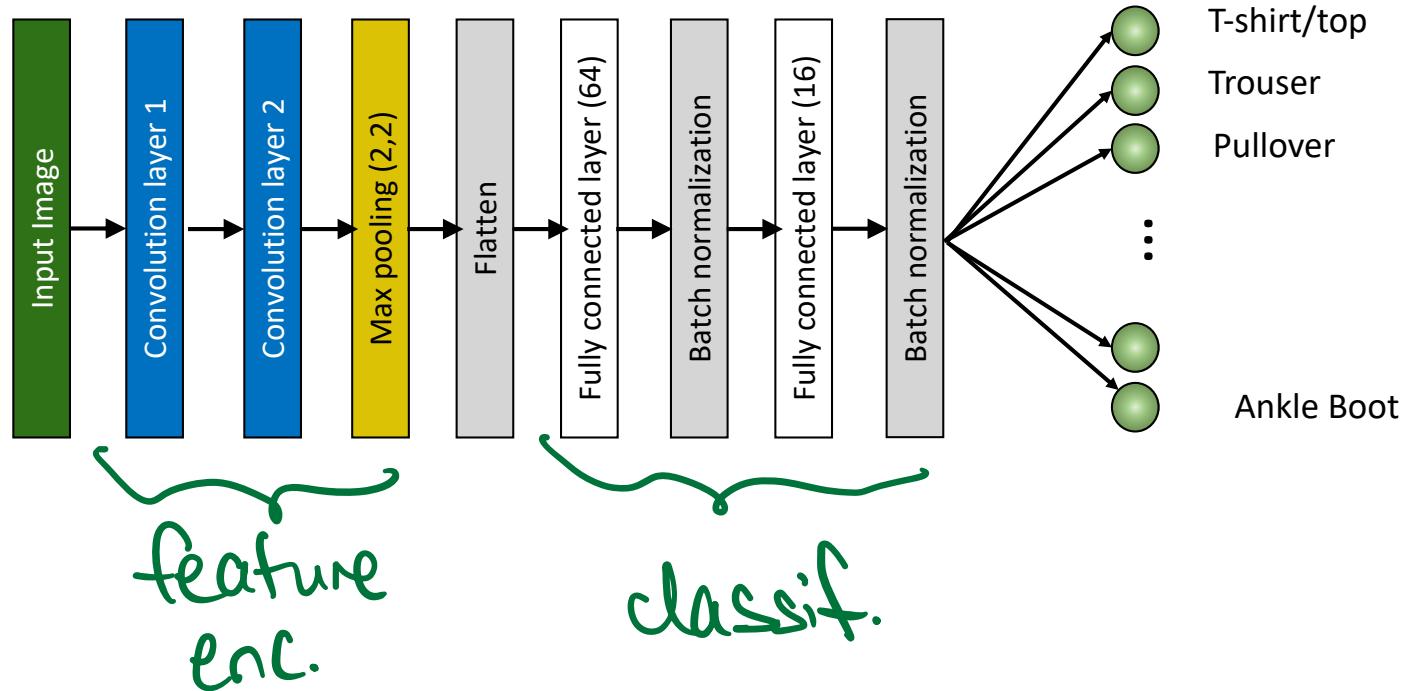
Max Pooling

- Similar to filtering, but output the maximum entry in the patch, instead of a weighted sum
- Typically output image is smaller than input → filters look at successively larger areas in the input image
- Usually applied per channel, not across channel



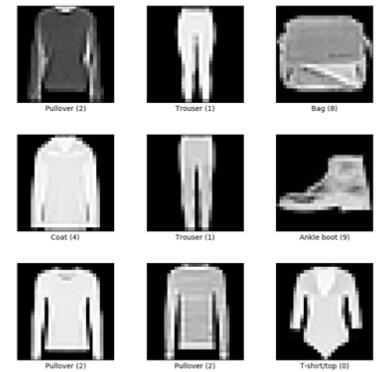
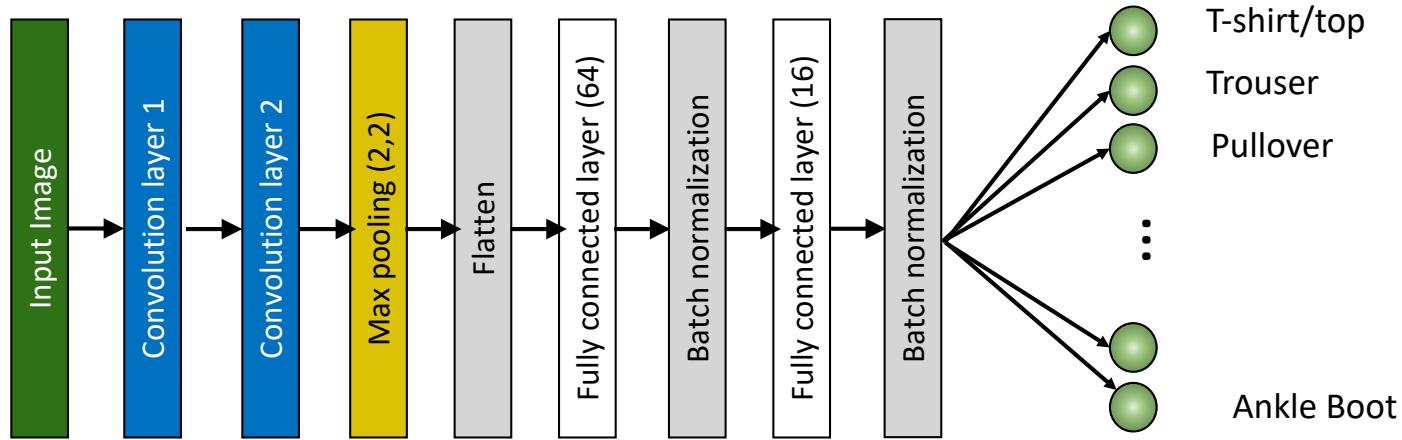
Case study: CNN for Fashion MNIST

- 2 convolution layers with 16 filters of size (3,3)



Case study: CNN for Fashion MNIST

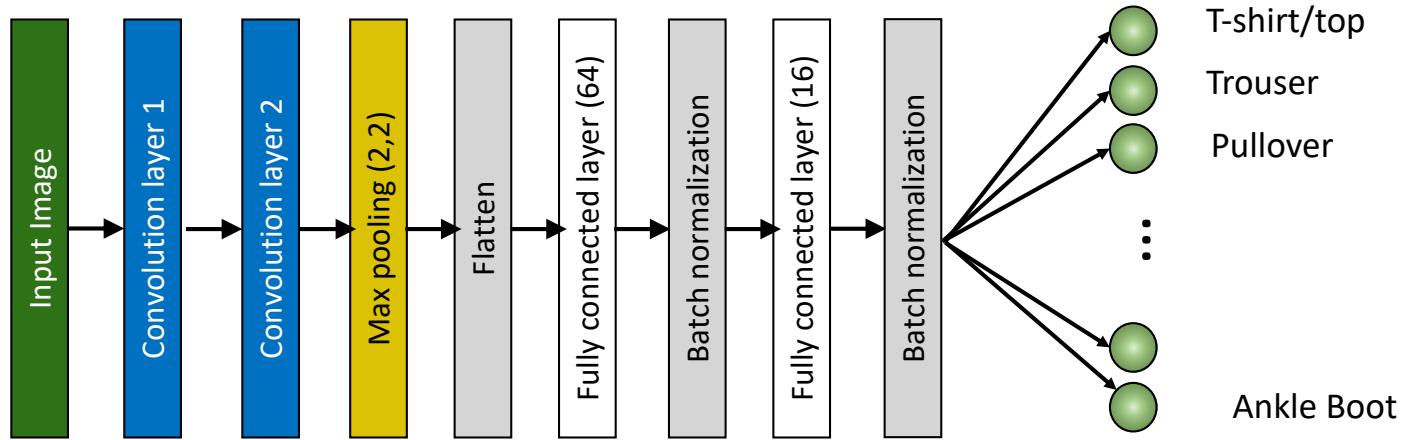
- 2 convolution layers with 16 filters of size (3,3)



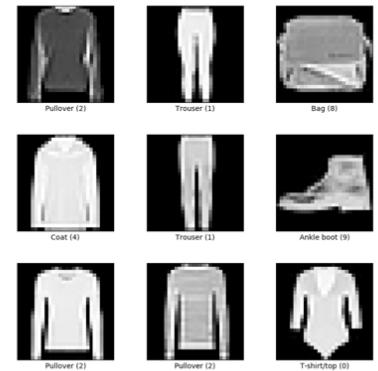
- Train with Adam for 10 epochs, mini-batch of 128

Case study: CNN for Fashion MNIST

- 2 convolution layers with 16 filters of size (3,3)

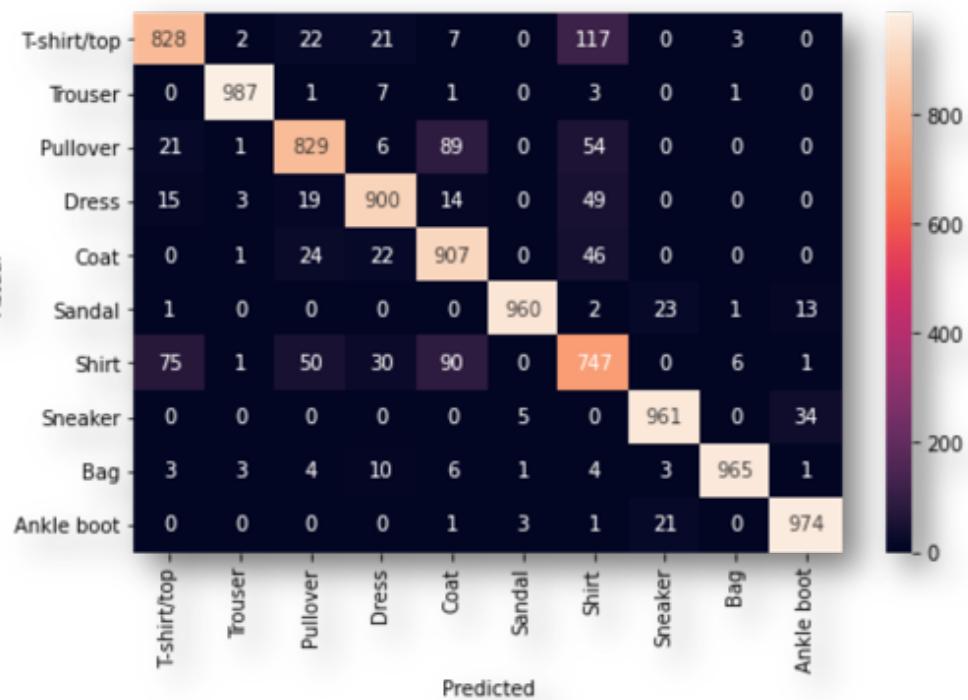


- Train with Adam for 10 epochs, mini-batch of 128
- Accuracy: 0.91 (vs 0.87 for fully connected)

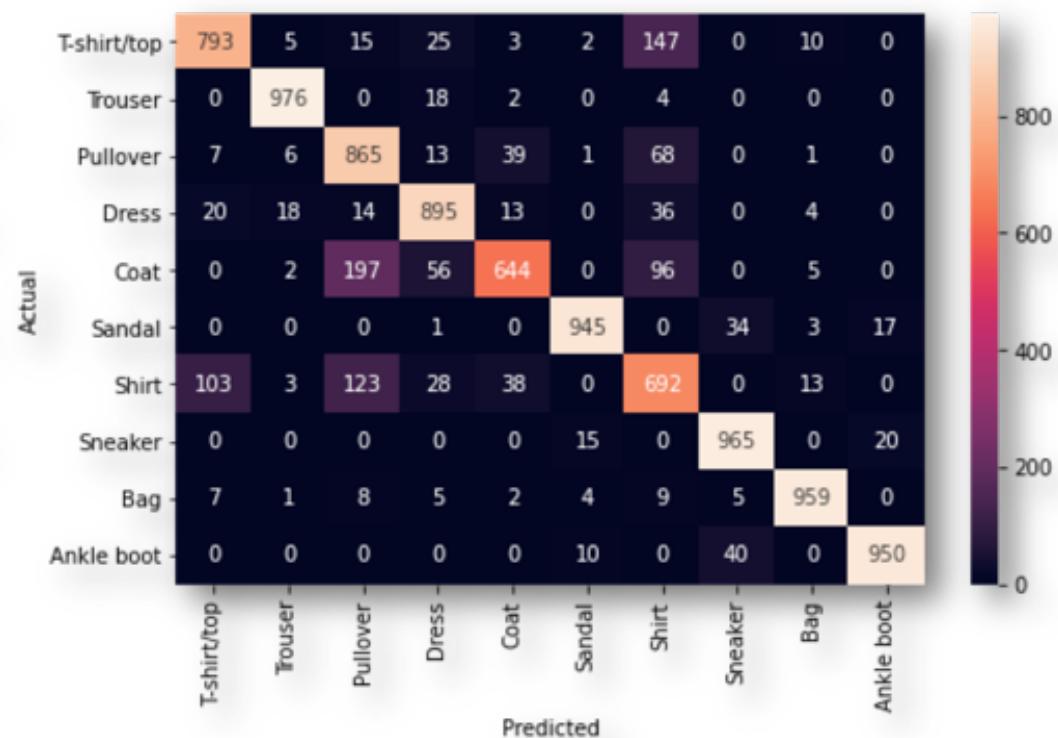


Confusion matrices: CNN vs. fully connected

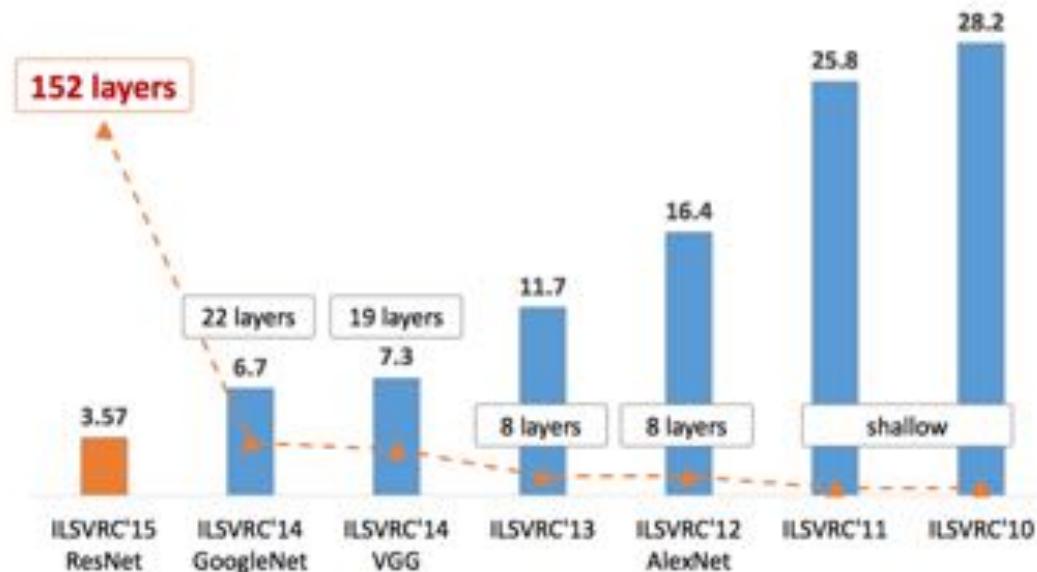
CNN



FNN

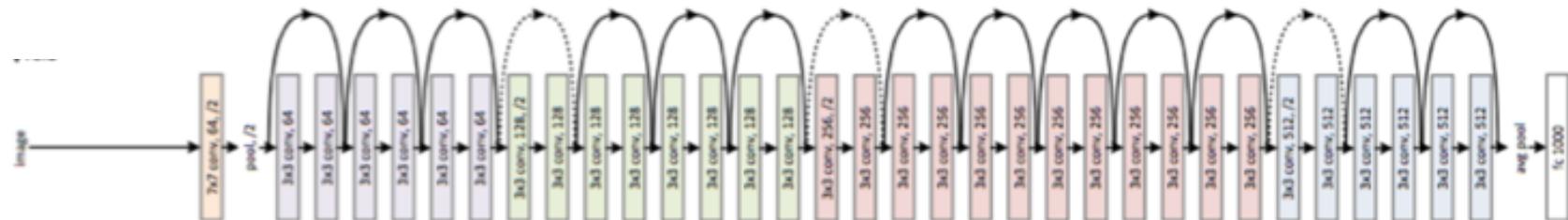
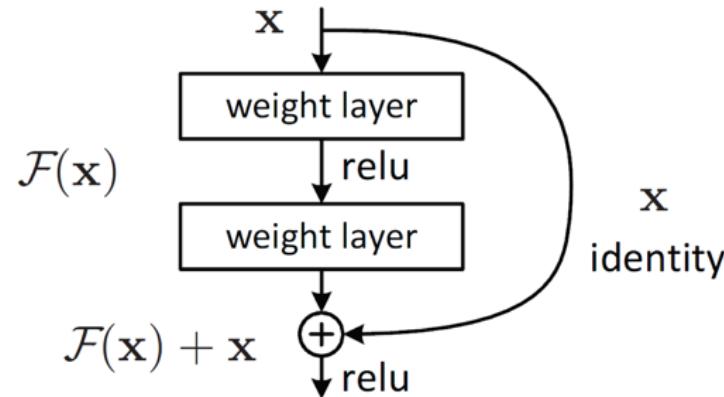


Deep CNNs



Optional: ResNet

- Skip connections improve training for very deep network
- Idea: add a layer's input to its output (shortcuts)



Summary: CNN model

- **Convolution:** “local detectors”
spatial locality
- **Weight sharing:**
apply same detector to all image patches
efficiency: much fewer parameters!, translation invariance
- **Pooling:** “did I see this pattern anywhere at all in this region?”
abstract away locality
- Trained with SGD and backpropagation

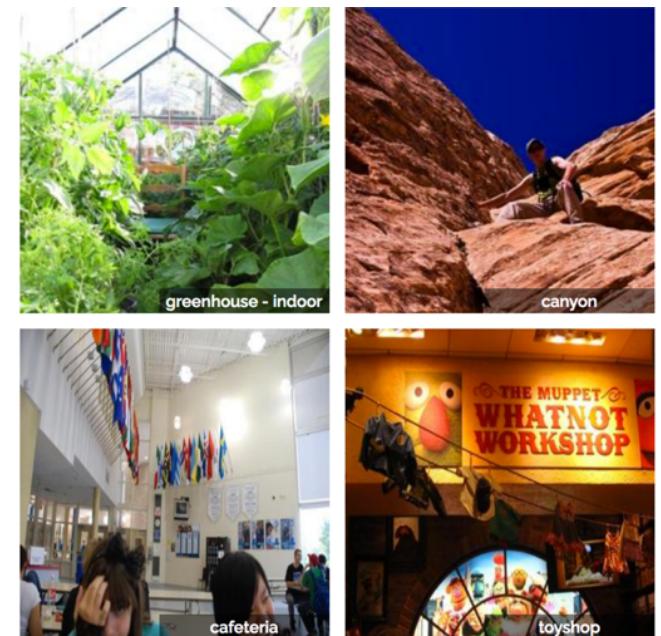
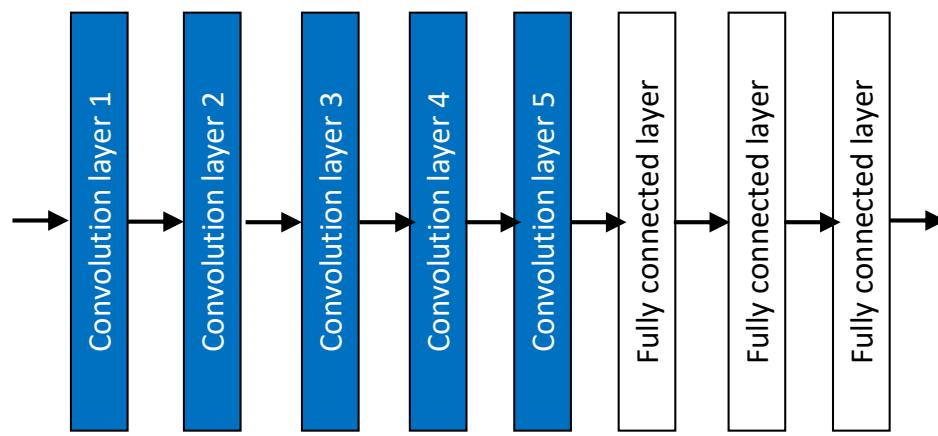
Outline: CNNs

- Spatial locality & invariance
- Convolutional Neural Networks
 - Convolution and filters
 - Max Pooling
- Example: Fashion MNIST
- CNN training
 - Illustrations: what do CNNs (not) learn?

What do neurons in a CNNs learn?

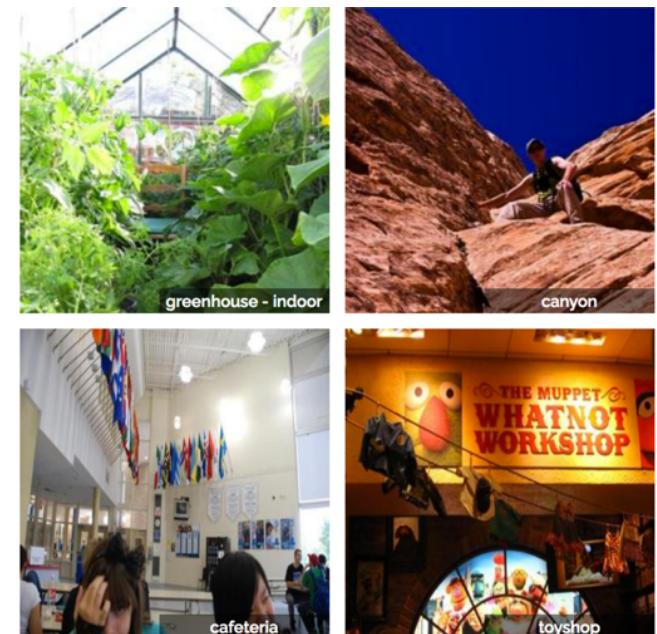
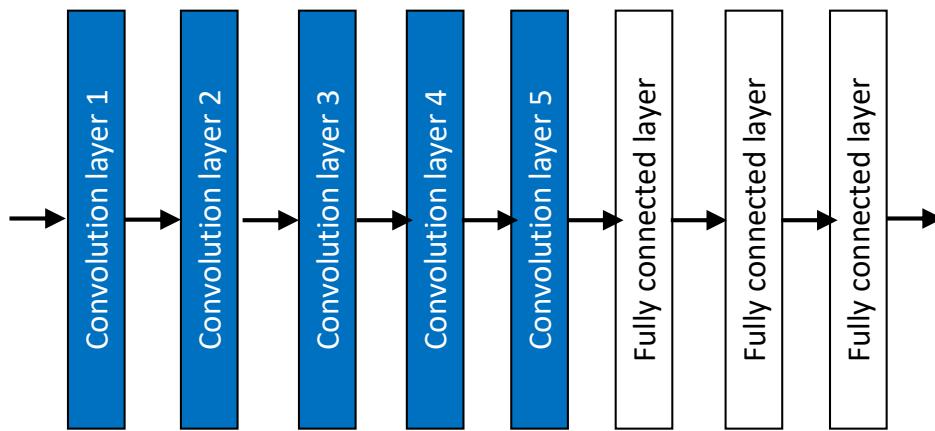
- **Visualization:**

View neurons as pattern detectors -- find patches that activate a neuron, and see whether they classify a “concept”



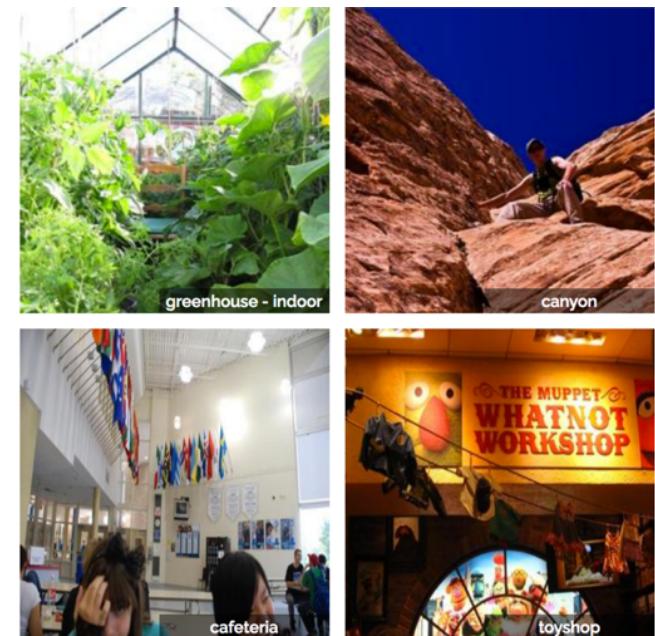
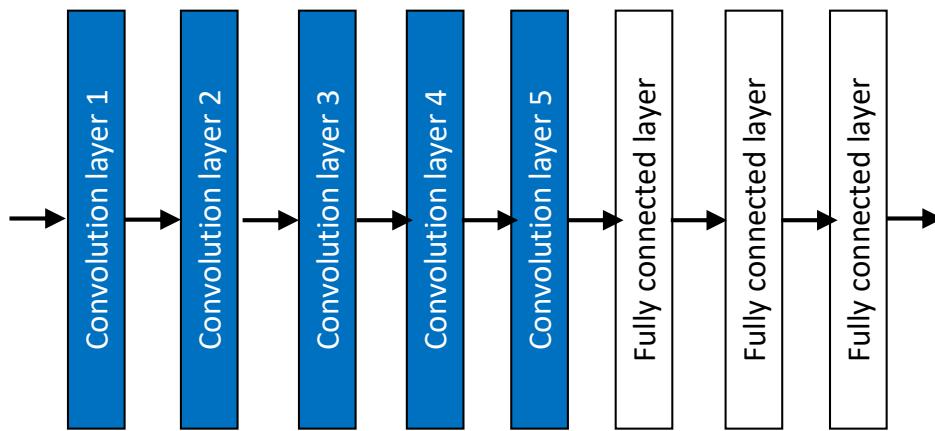
What do neurons in a CNNs learn?

- **Visualization:**
View neurons as pattern detectors -- find patches that activate a neuron, and see whether they classify a “concept”
- What would you expect? Hierarchical representations
edges -> shapes/textures -> parts -> objects

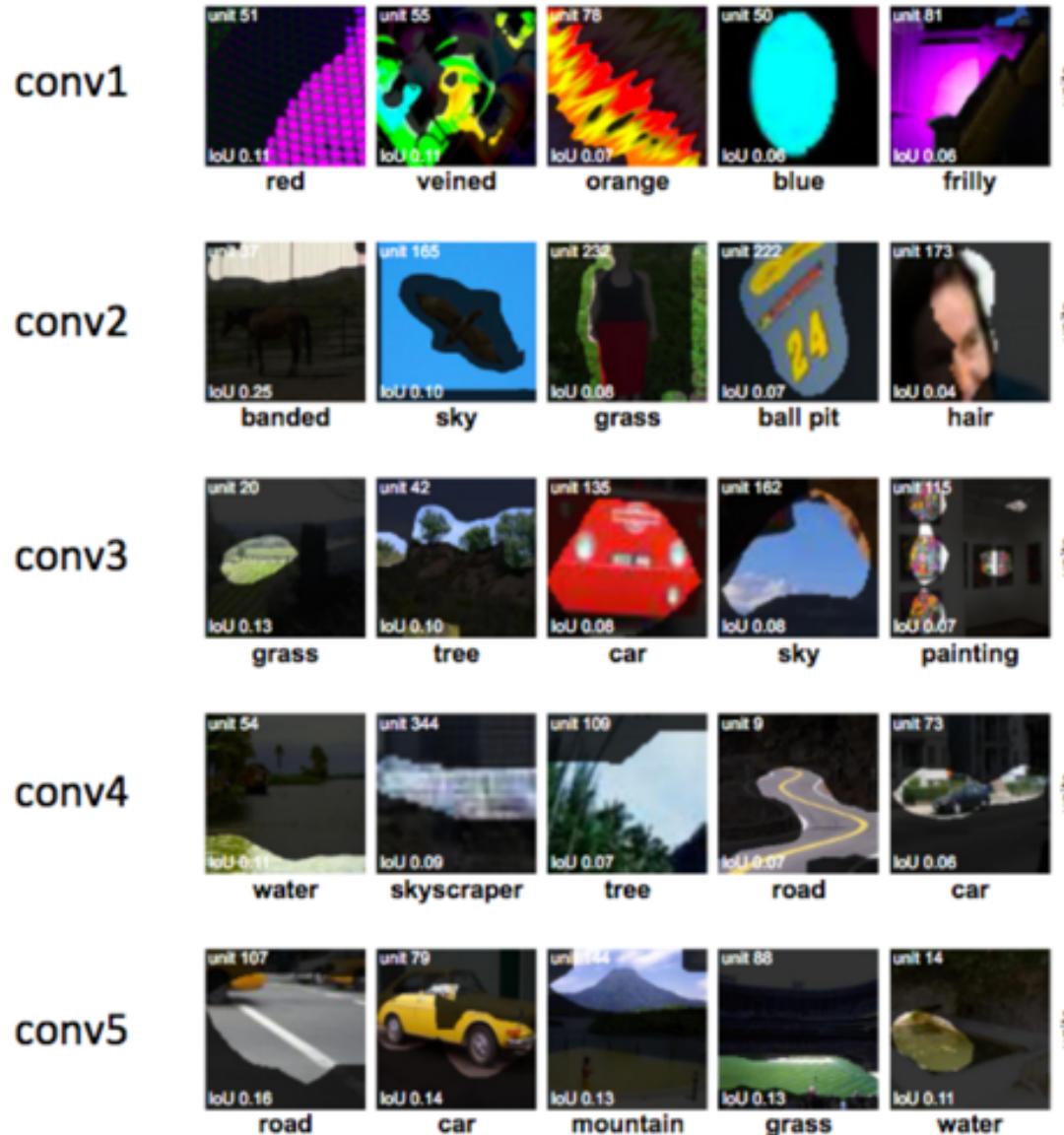


What do neurons in a CNNs learn?

- **Visualization:**
View neurons as pattern detectors -- find patches that activate a neuron, and see whether they classify a “concept”
- What would you expect? Hierarchical representations
edges -> shapes/textures -> parts -> objects
- Example: train AlexNet to classify scenes



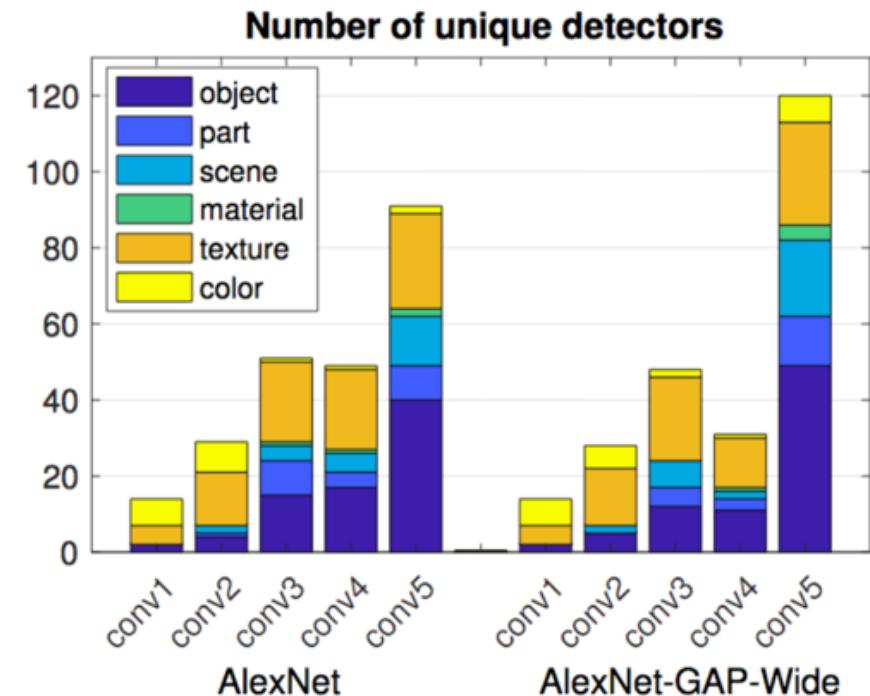
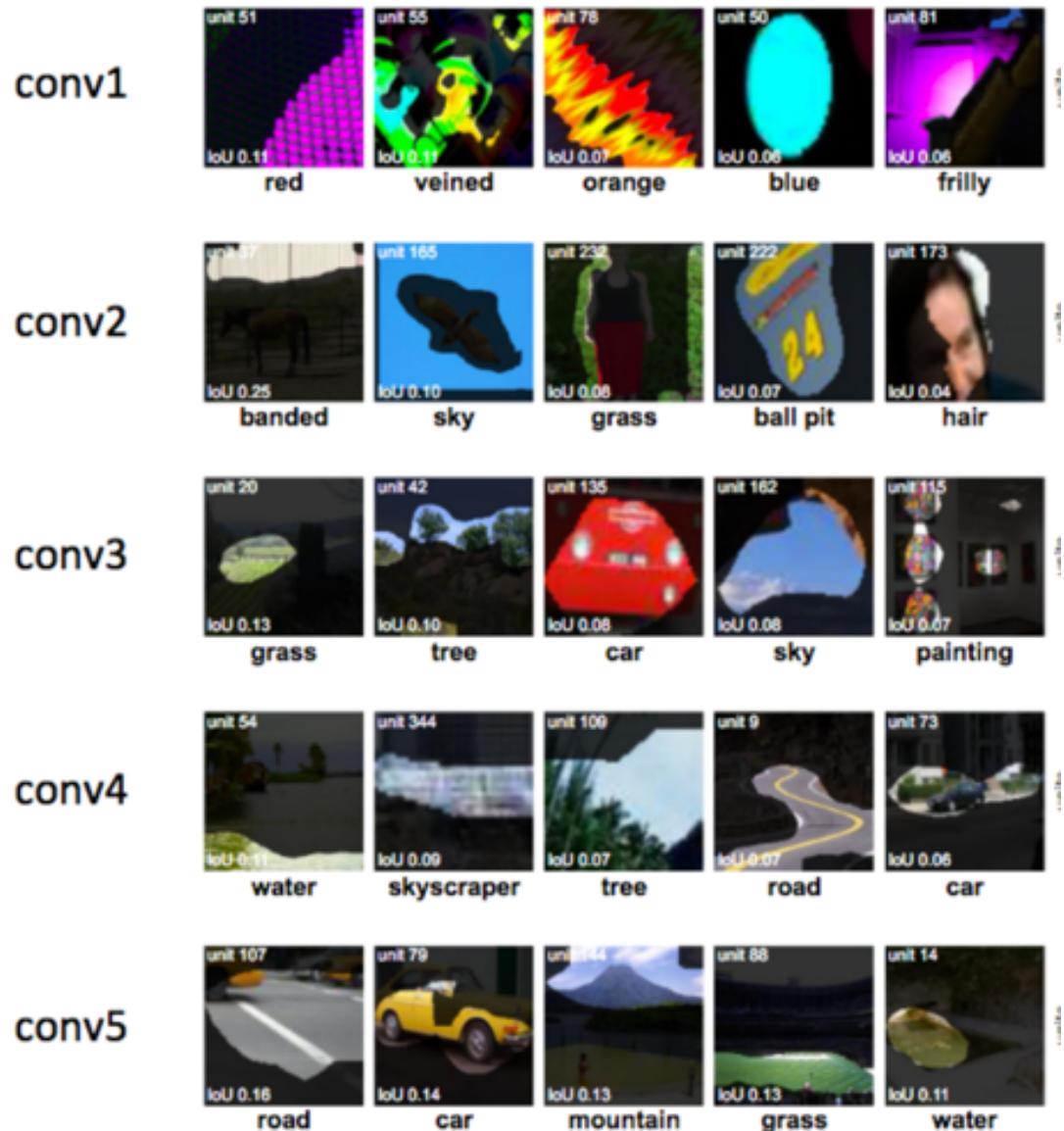
What do neurons in a CNNs learn?



First layer(s): detectors for colors, edges, textures
Later layers: detectors for parts, objects, scenes

<http://netdissect.csail.mit.edu/>

What do neurons in a CNNs learn?

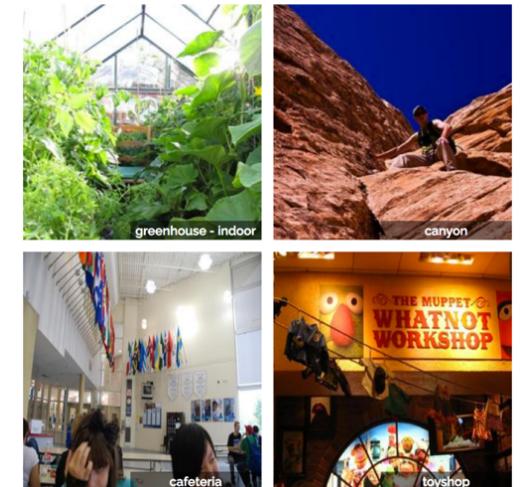


First layer(s): detectors for colors, edges, textures
Later layers: detectors for parts, objects, scenes

<http://netdissect.csail.mit.edu/>

Example: scene detection

When you learn to classify scenes ...



... the network learns object detectors!

House

res5c unit 1410



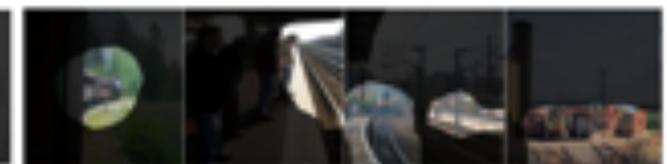
Dog

IoU=0.142 res5c unit 1573



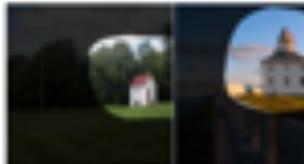
Train

IoU=0.216 res5c unit 924



ResNet-152

res5c unit 301



IoU=0.087



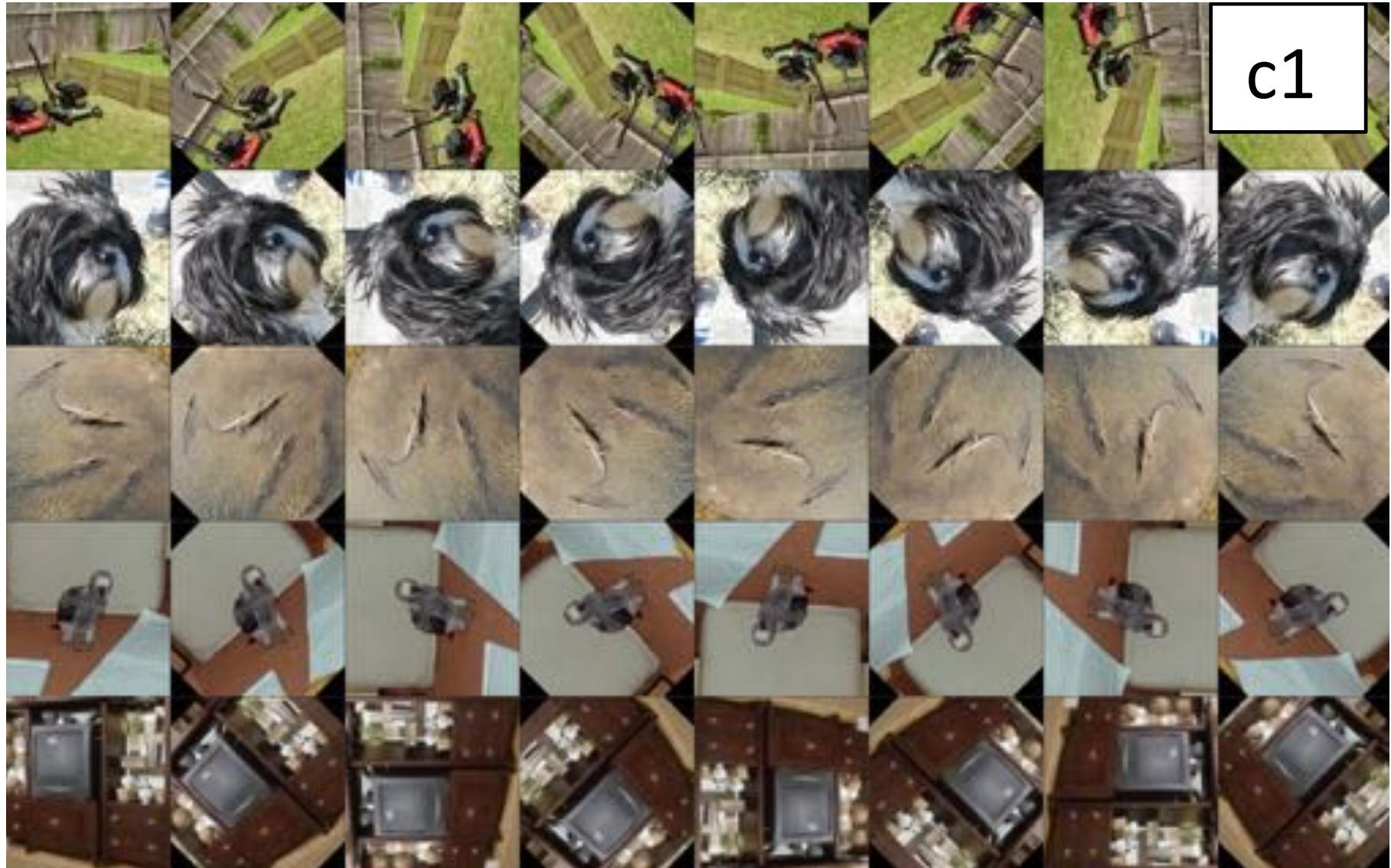
IoU=0.193

res5c unit 2001

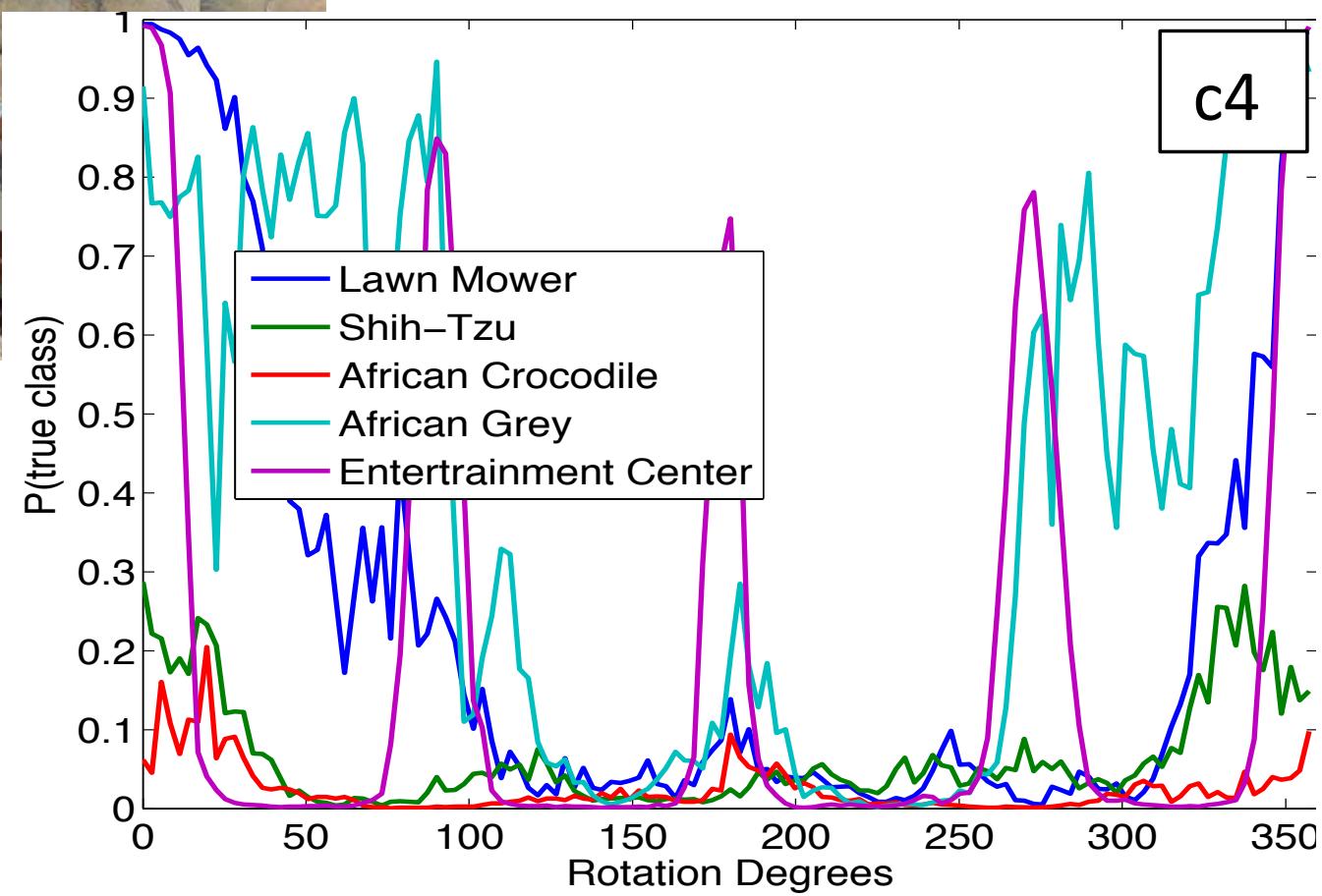
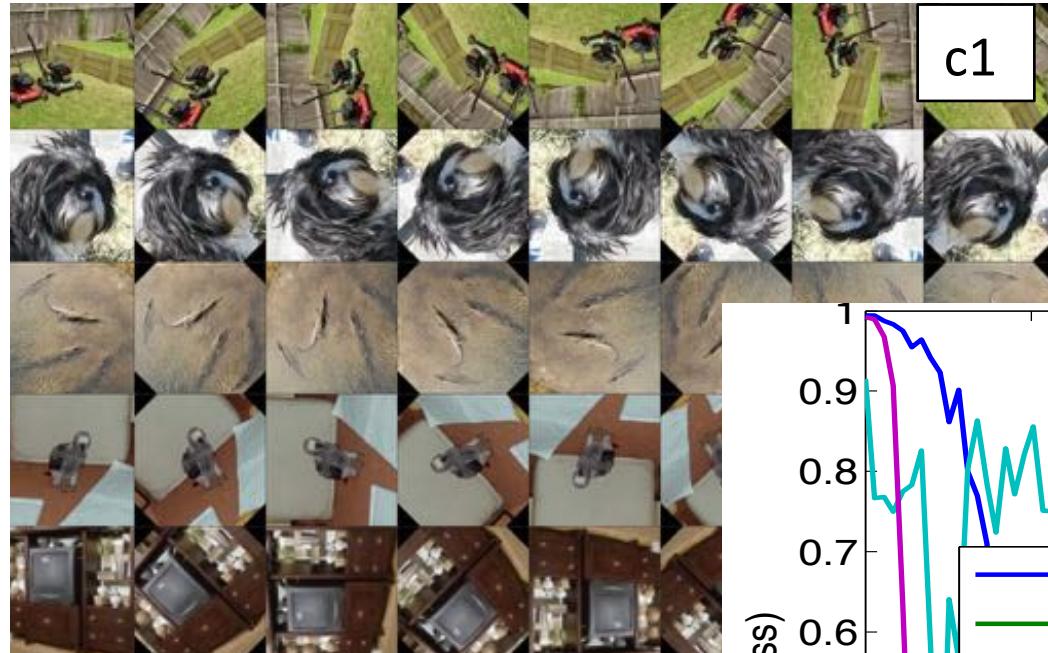
IoU=0.255



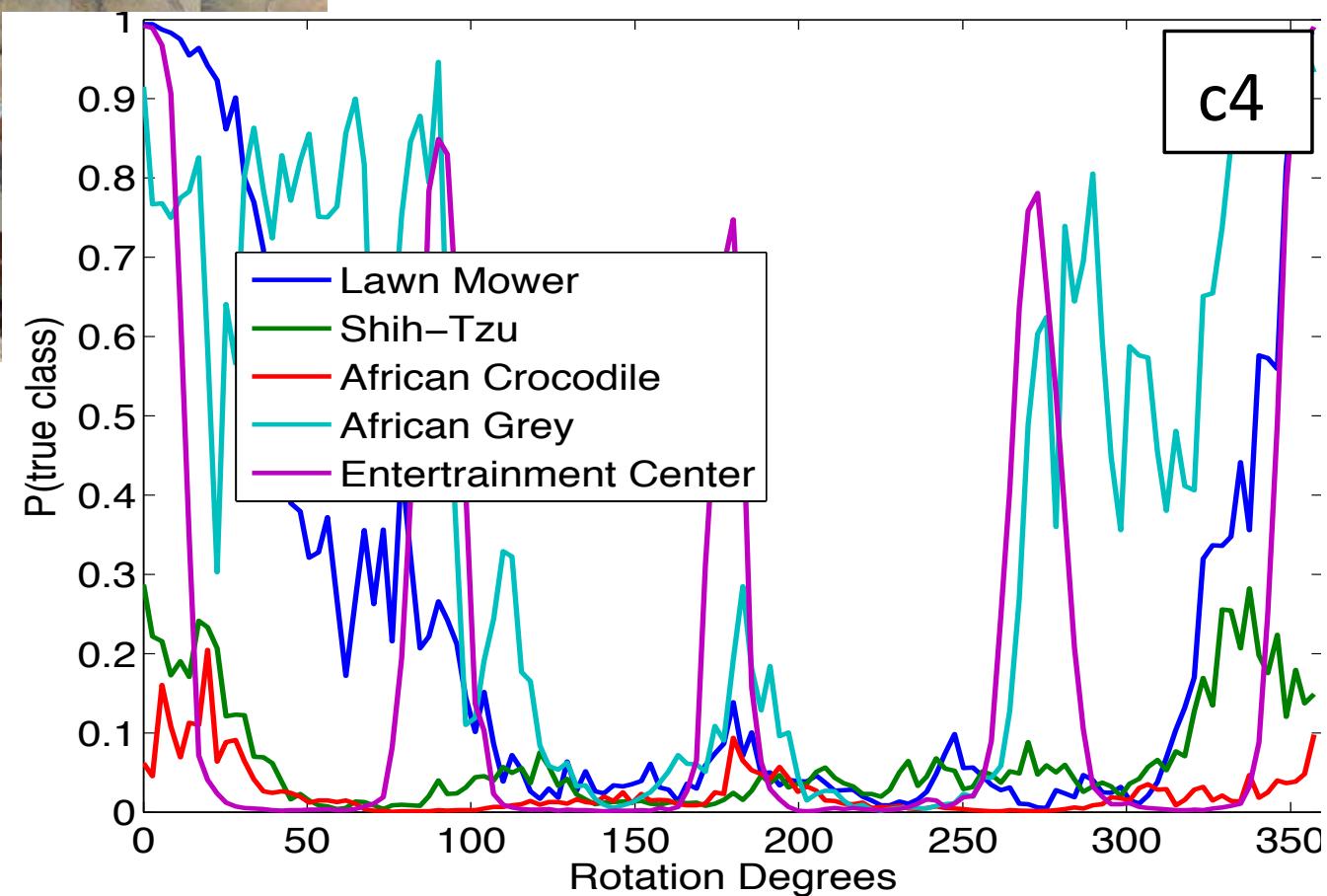
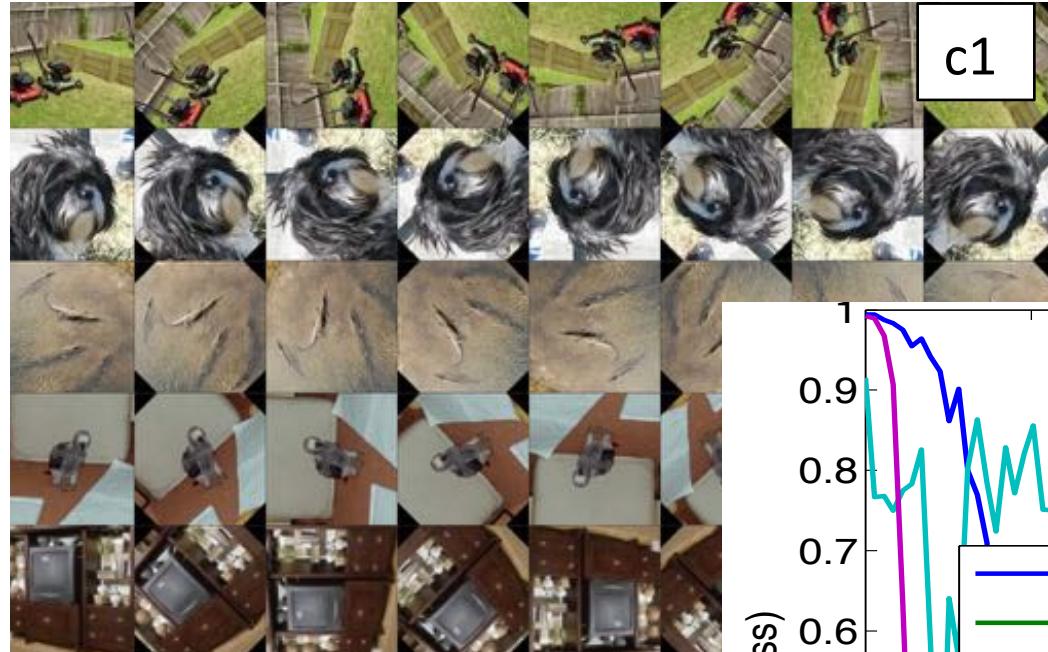
Input rotations



Input rotations



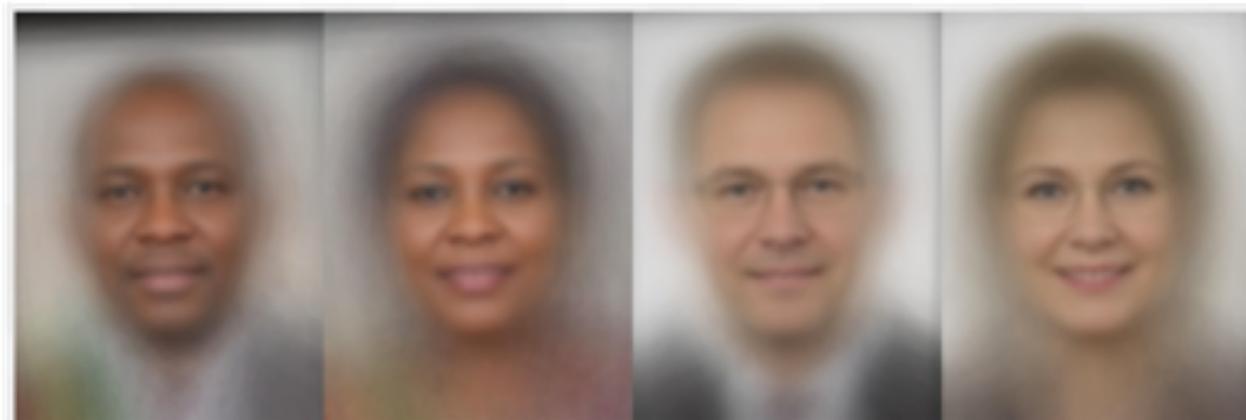
Input rotations



CNNs can fail when
the input is rotated.

Bias

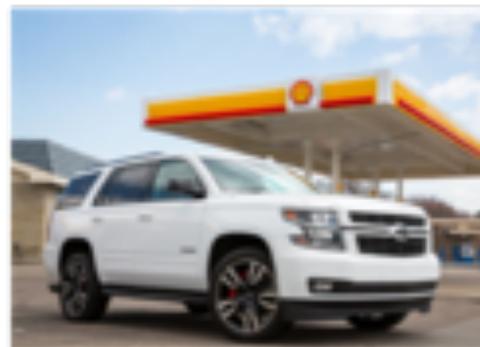
Gender Classifier	Darker Male	Darker Female	Lighter Male	Lighter Female	Largest Gap
Microsoft	94.0%	79.2%	100%	98.3%	20.8%
FACE++	99.3%	65.5%	99.2%	94.0%	33.8%
IBM	88.0%	65.3%	99.7%	92.9%	34.4%



Errors can be
very different for
different classes.

Are we using the right features?

“car”

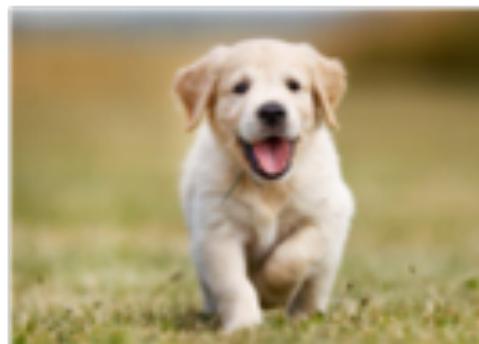


Are we using the right features?

“car”

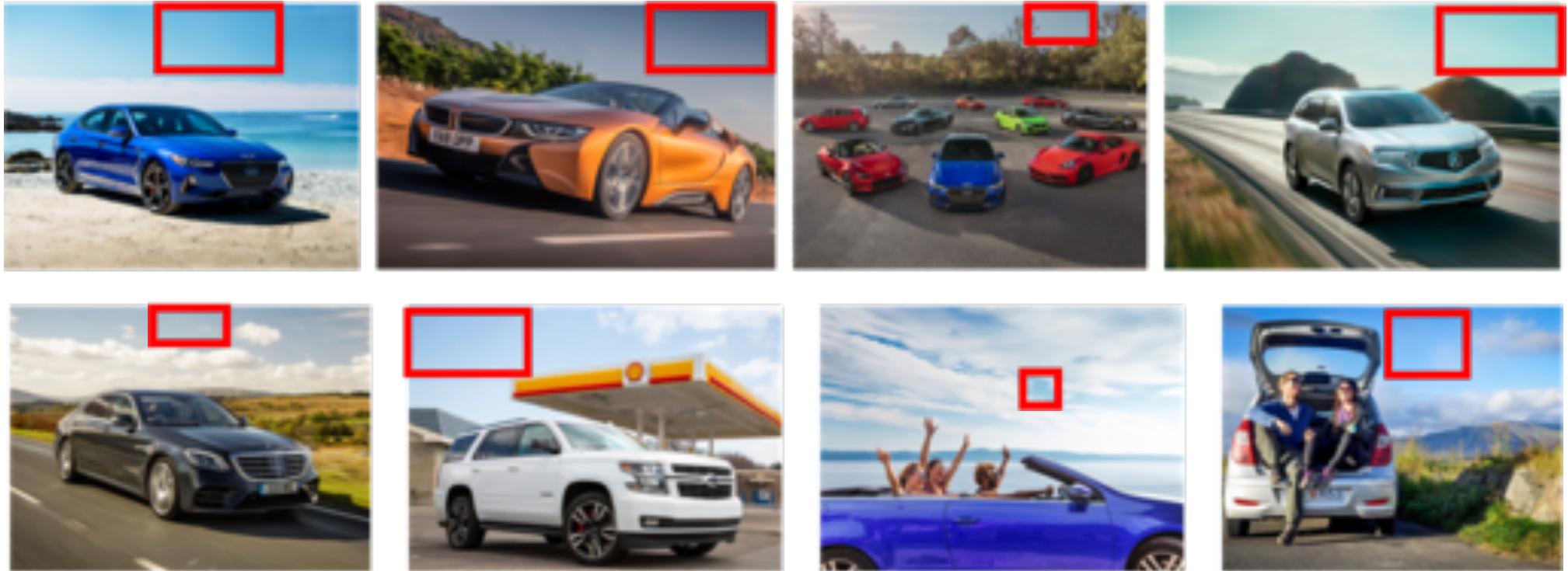


“not a car”



Are we using the right features?

“car”

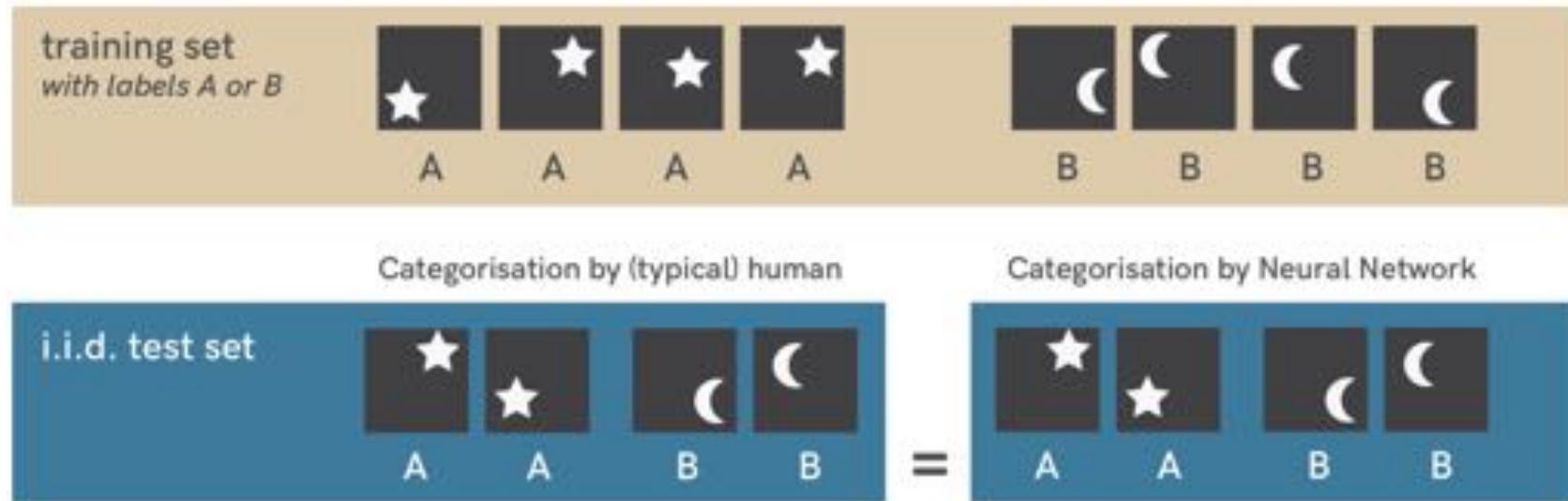


“not a car”



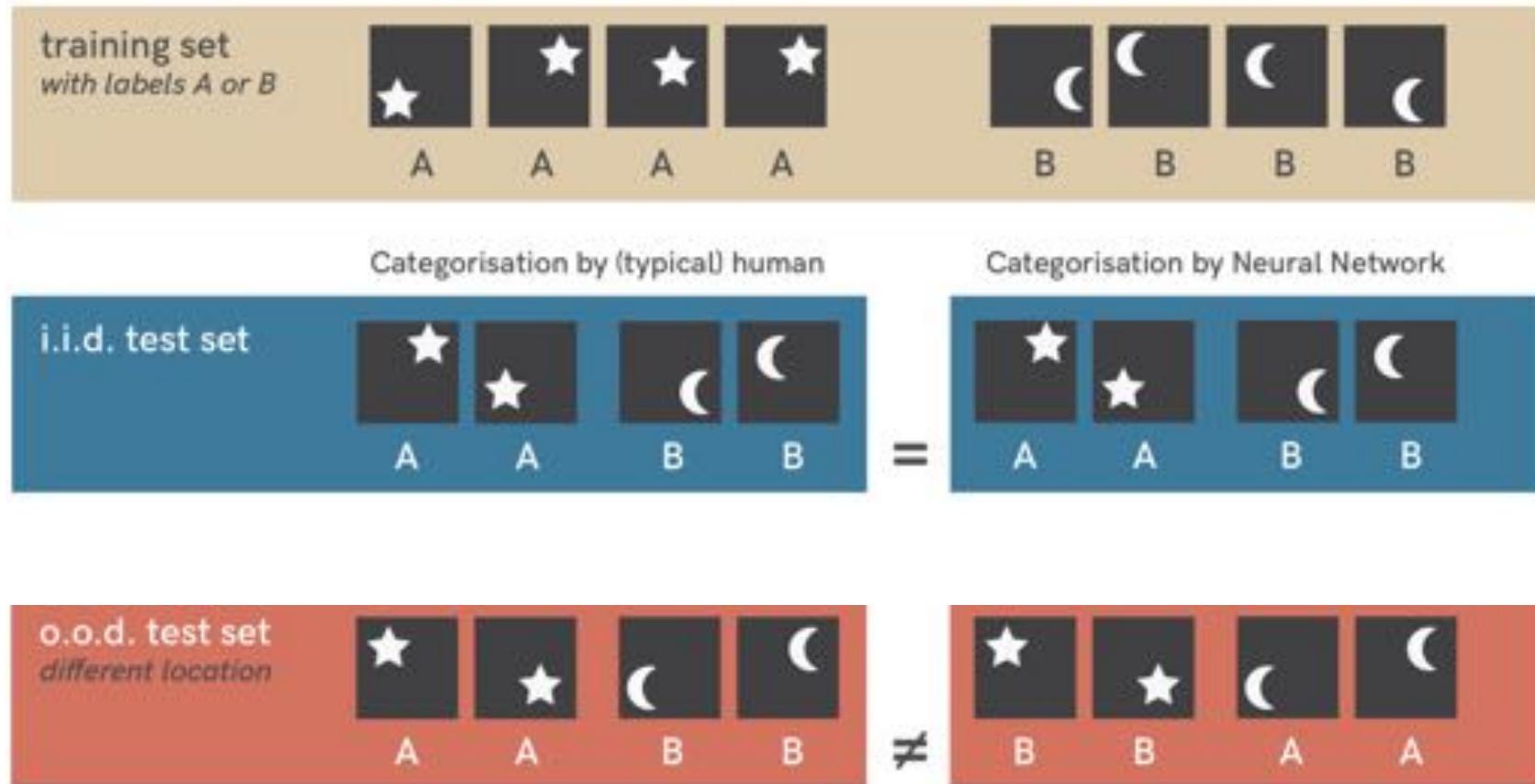
“Shortcut solutions”

Good test performance within distribution, but bad under distribution shifts



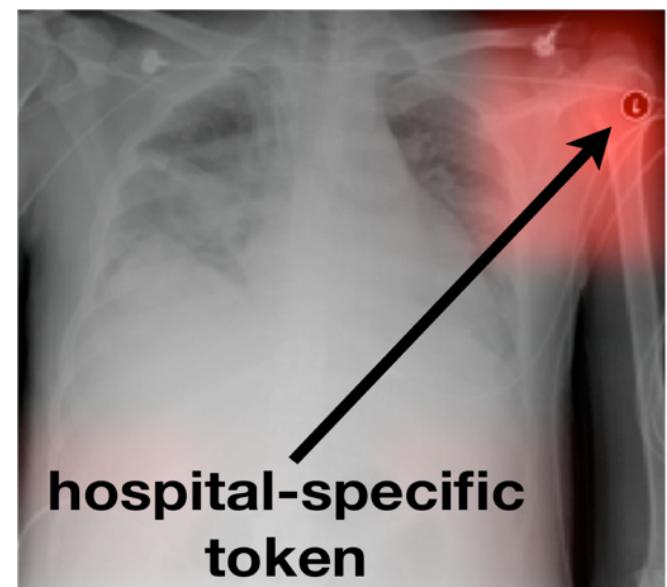
“Shortcut solutions”

Good test performance within distribution, but bad under distribution shifts



“Shortcuts” more generally

Detecting pneumonia from x-rays

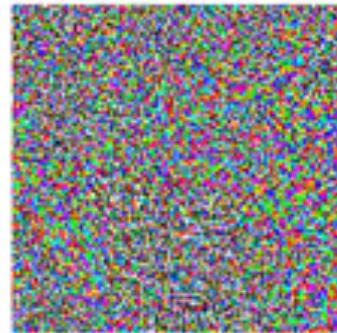


“Shortcuts” more generally



x
“panda”
57.7% confidence

$+ .007 \times$



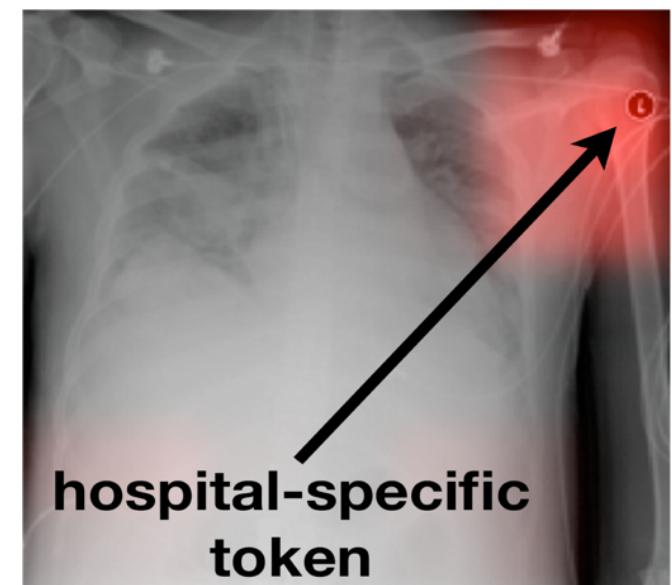
$\text{sign}(\nabla_x J(\theta, x, y))$
“nematode”
8.2% confidence



$x +$
 $\epsilon \text{sign}(\nabla_x J(\theta, x, y))$
“gibbon”
99.3 % confidence

(Goodfellow et al 2014,...)

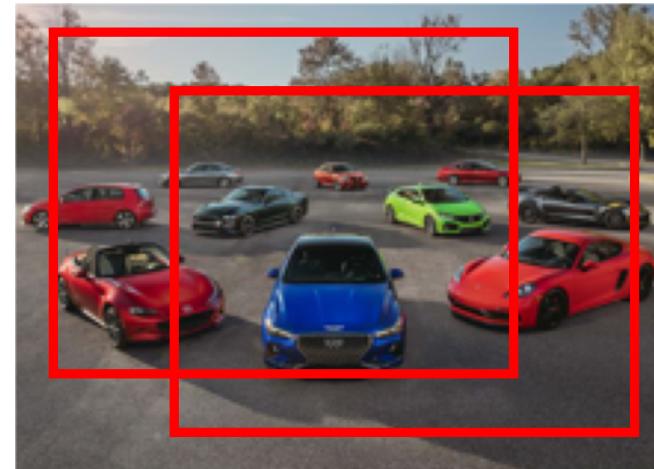
Detecting pneumonia from x-rays



Are we using the right features?

What can we do?

- collect more data
- data augmentation (cropping, rotation, ...)
- pre-training
- ...



Examples of data transformations



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate ($90^\circ, 180^\circ, 270^\circ$)



(g) Cutout



(h) Gaussian noise



(i) Gaussian blur



(j) Sobel filtering

(Chen, Kornblith, Norouzi, Hinton 2020)

Selection of tools: e.g. Tensorflow, Keras preprocessing layers

Summary: CNNs

- Neural Network for images
- **Convolution:** “local detectors”
spatial locality, efficiency via weight sharing
- **Pooling:** abstract away locality
- Trained with SGD and backpropagation
- Modern CNNs are deep
- Data augmentation to improve performance

regul: $\|\theta\|_2^2 = \sum_{j=1}^d \theta_j^2$ $\|\theta\|_1 = \sum_{j=1}^d |\theta_j|$