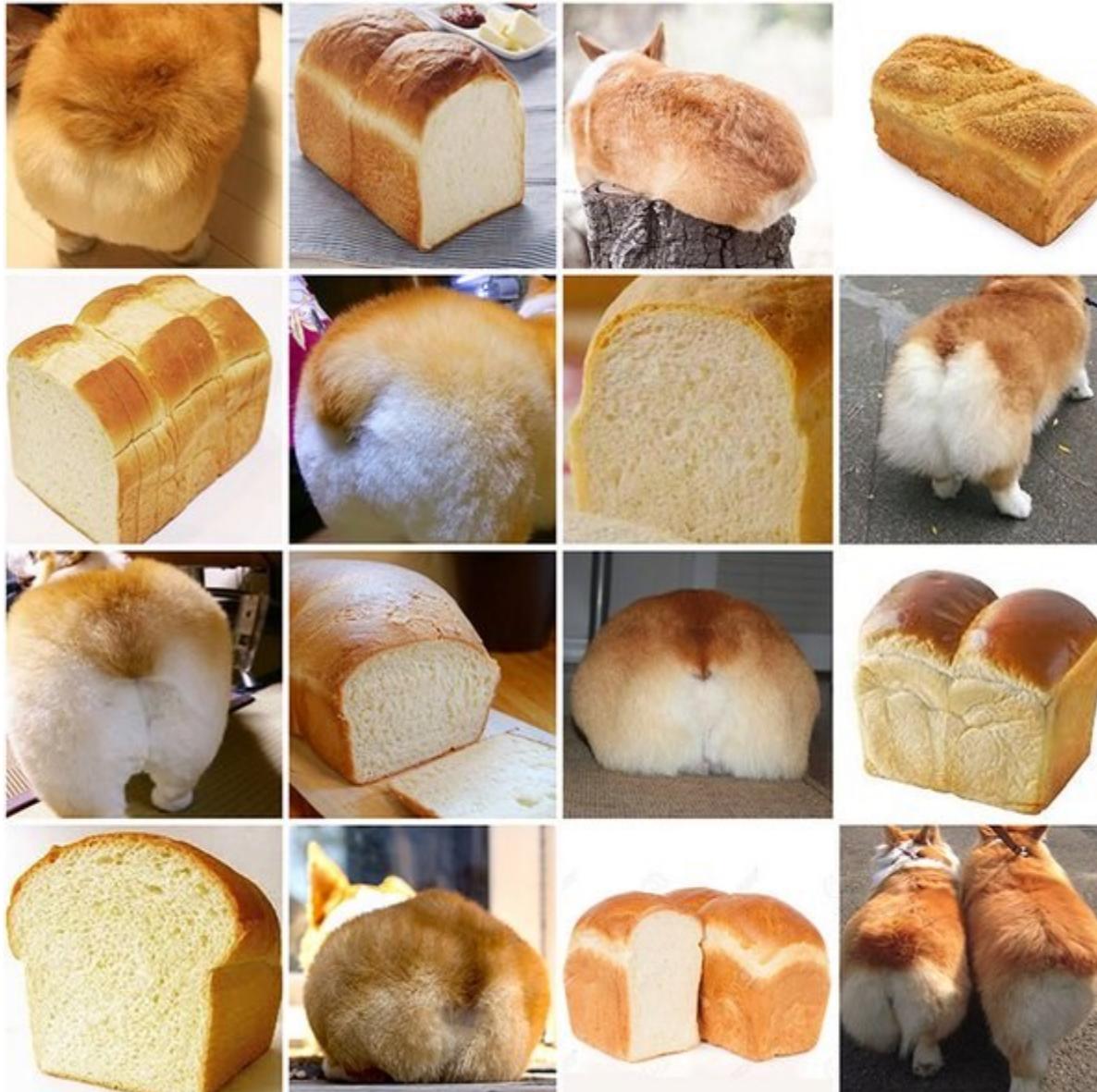


Advanced Computer Vision



FUTURE VISION

Convolutional Neural Networks







Wow



false positives

what class



so misclassified



@teenybiscuit

no good filtr

cool kernel

Goals

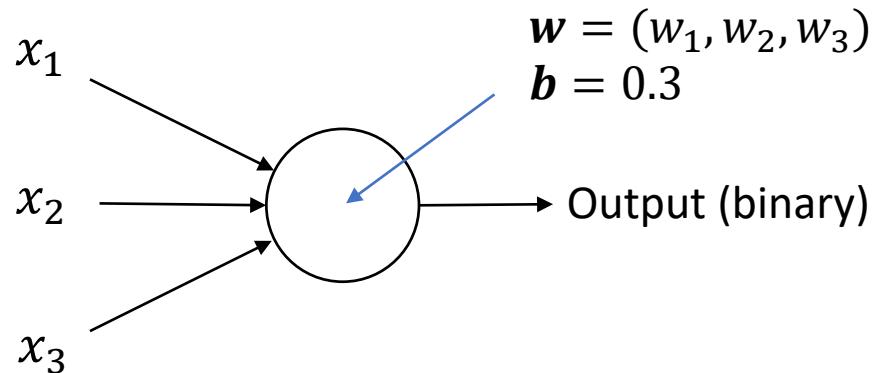
Build a classifier which is more powerful at representing complex functions *and* more suited to the learning problem.

What does this mean?

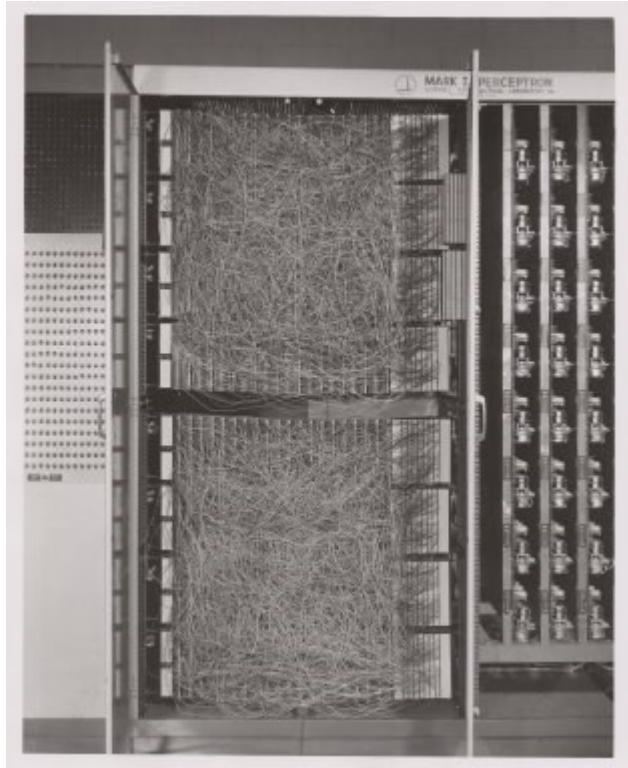
1. Assume that the *underlying data generating function* relies on a composition of factors.
2. Learn a feature representation that is specific to the dataset.

Neural Networks

- Basic building block for composition is a *perceptron* (Rosenblatt c.1960)
- Linear classifier – vector of weights w and a ‘bias’ b



$$\text{output} = \begin{cases} 0 & \text{if } w \cdot x + b \leq 0 \\ 1 & \text{if } w \cdot x + b > 0 \end{cases} \quad w \cdot x \equiv \sum_j w_j x_j$$



Mark 1 Perceptron
c.1960

20x20 pixel
camera feed

Universality

A single-layer network can learn any function:

- So long as it is differentiable
- To some approximation;
More perceptrons = a better approximation

Visual proof (Michael Nielson):

<http://neuralnetworksanddeeplearning.com/chap4.html>

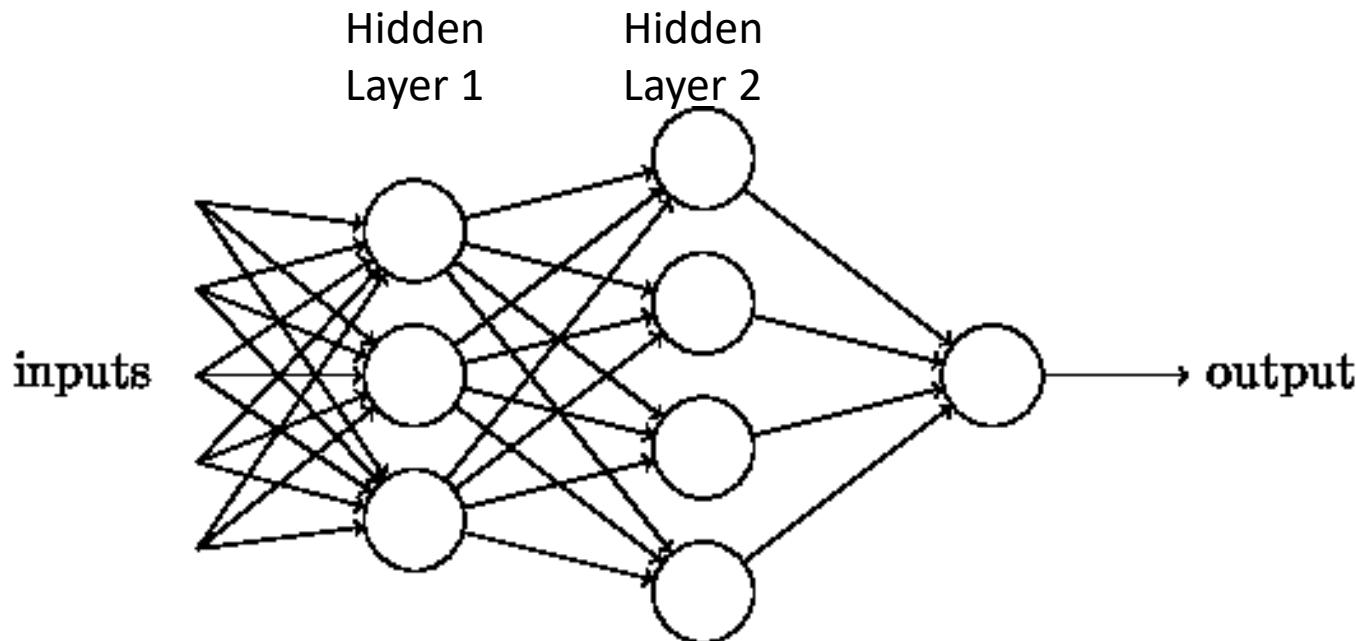
*If a single-layer network can learn any function...
...given enough parameters...*

...then why do we go deeper?

Intuitively, composition is efficient because it allows *reuse*.

Empirically, deep networks do a better job than shallow networks at learning such hierarchies of knowledge.

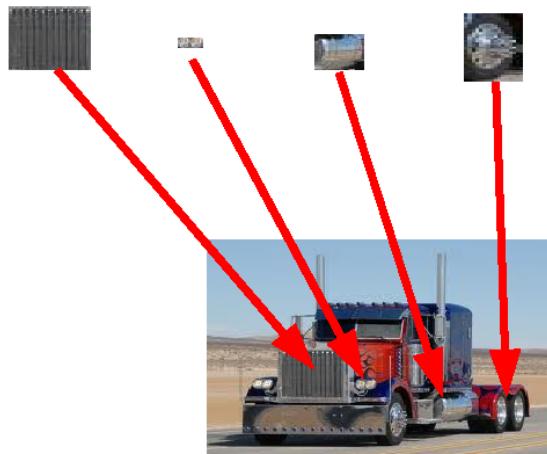
Composition



Layers that are in between the input and the output are called *hidden layers*, because we are going to *learn* their weights via an optimization process.

Interpretation of many layers

[0 0 1 0 0 0 0 1 0 0 1 1 0 0 1 0 ...] truck feature



Exponentially more efficient than a
1-of-N representation (a la k-means)

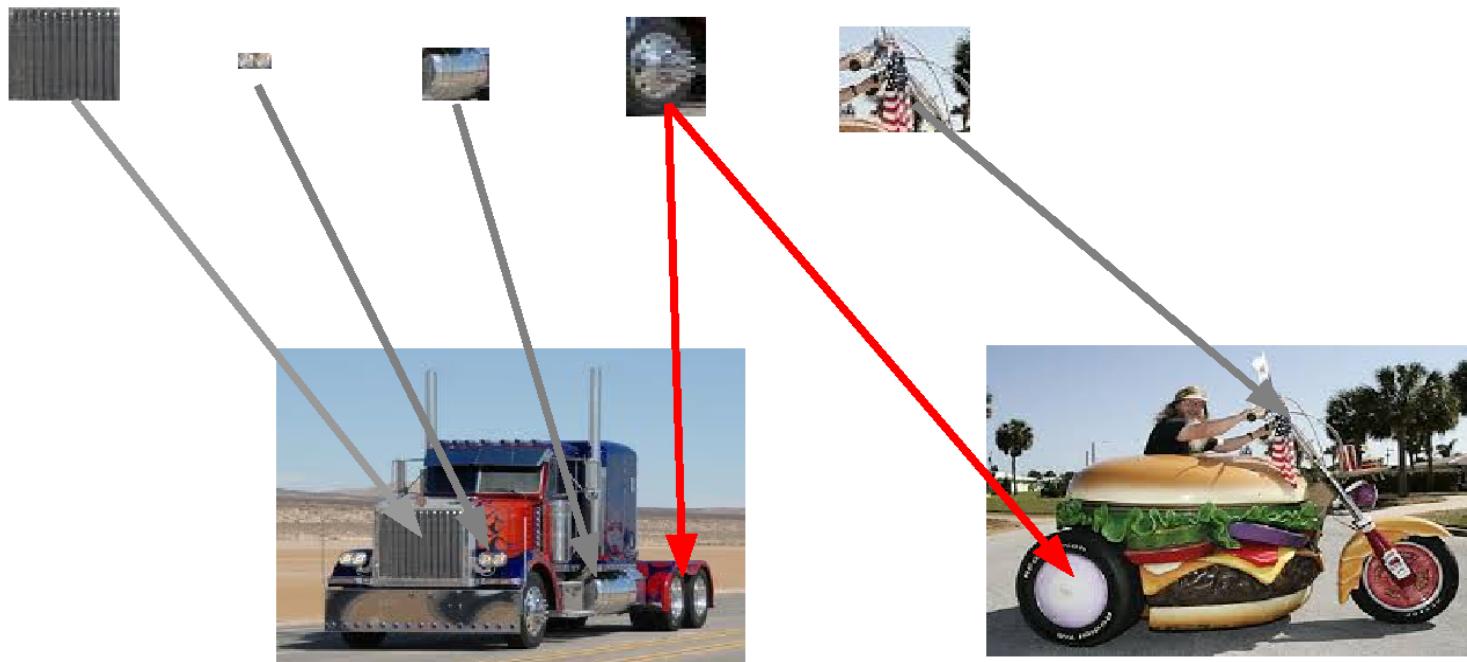
14

Ranzato 

Interpretation

[1 1 0 0 0 1 0 1 0 0 0 0 1 1 0 1 ...] motorbike

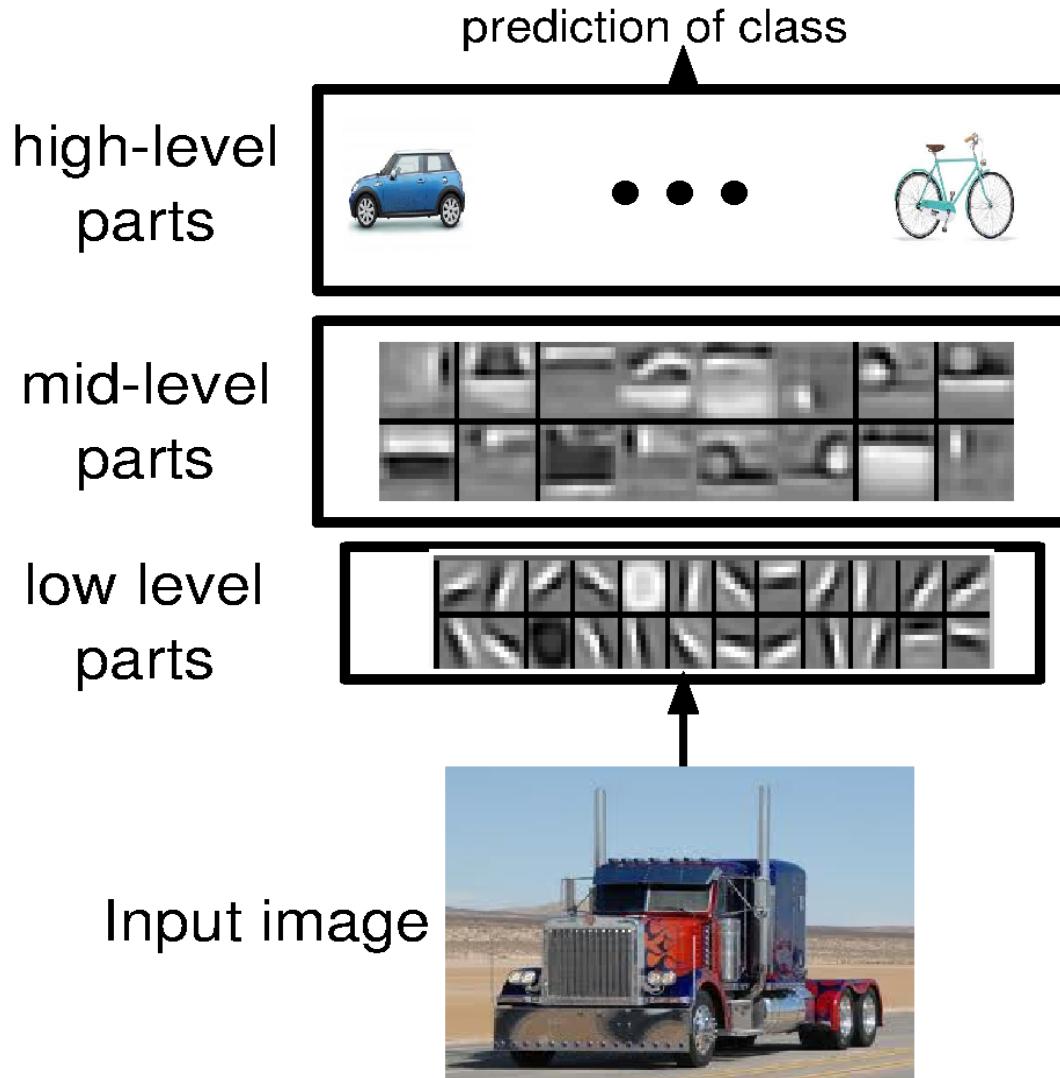
[0 0 1 0 0 0 0 1 0 0 1 1 0 0 1 0 ...] truck



15

Ranzato

Interpretation



- distributed representations
- feature sharing
- compositionality

NOTE: Not actually the weights; a demonstrative visualization!

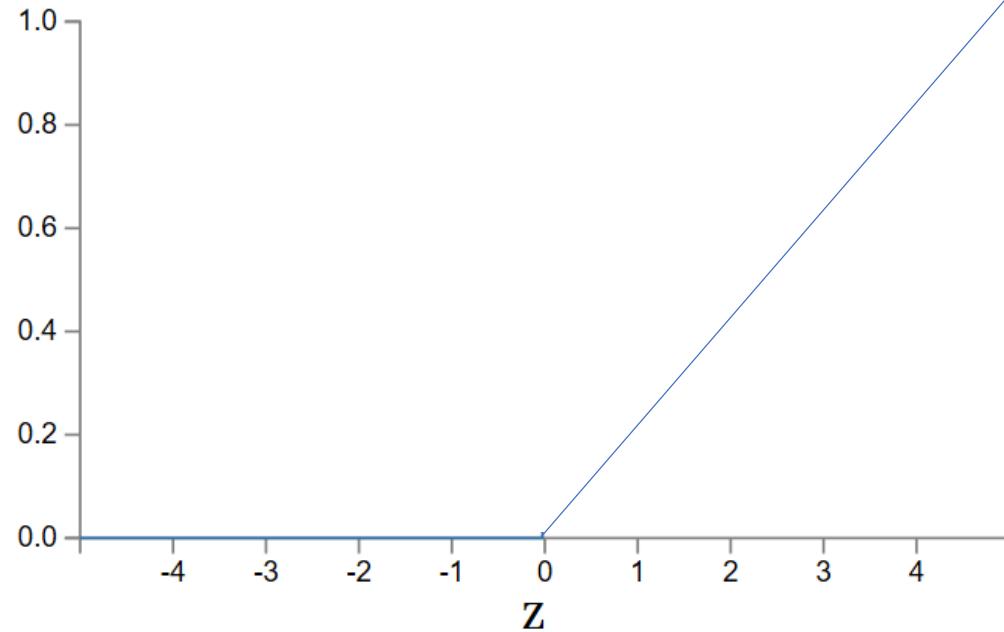
16

Ranzato

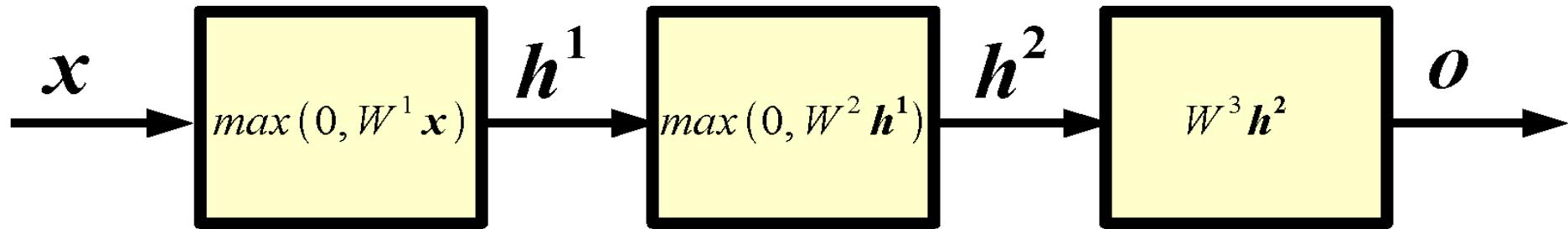
Lee et al. "Convolutional DBN's ..." ICML 2009

Activation functions: Rectified Linear Unit

- ReLU $f(x) = \max(0, x)$



Neural Networks: example



x input

h^1 1-st layer hidden units

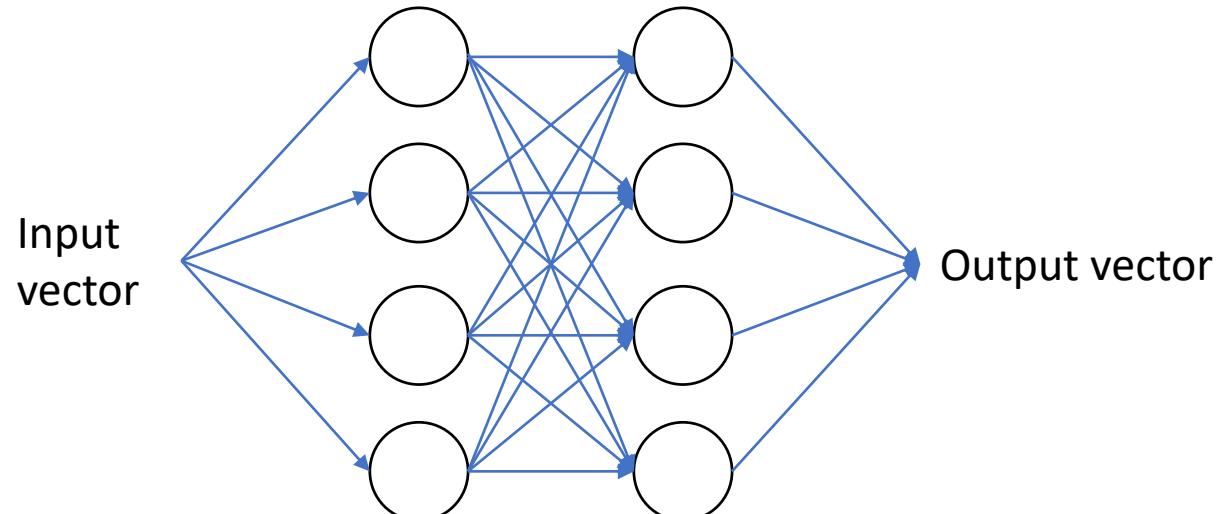
h^2 2-nd layer hidden units

o output

Example of a 2 hidden layer neural network (or 4 layer network, counting also input and output).

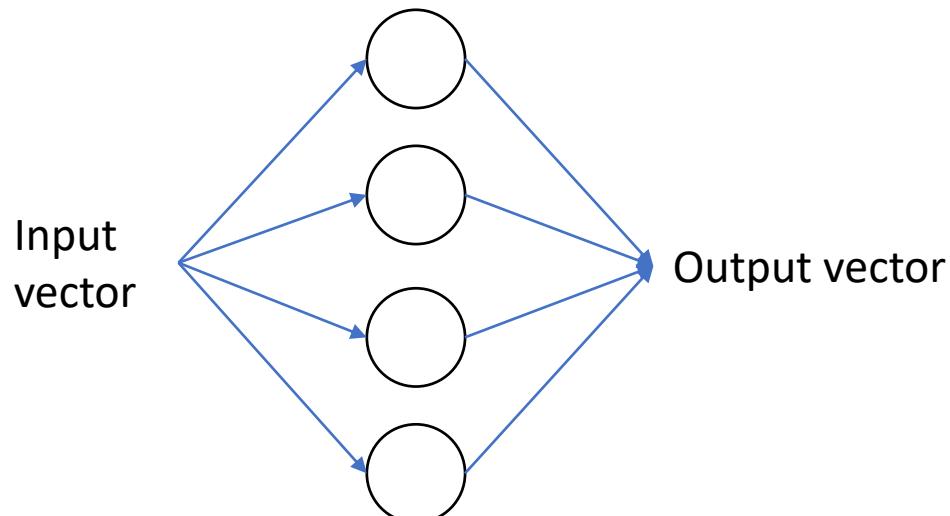
Does anyone pass along the weight without an activation function?

No – this is linear chaining.



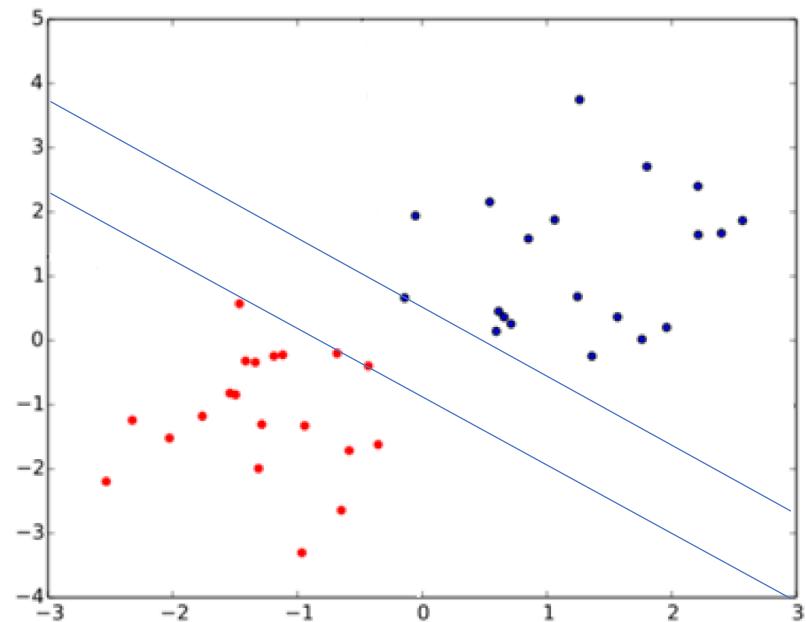
Does anyone pass along the weight without an activation function?

No – this is linear chaining.



What is the relationship between SVMs and perceptrons?

SVMs attempt to learn the support vectors which maximize the margin between classes.



What is the relationship between SVMs and perceptrons?

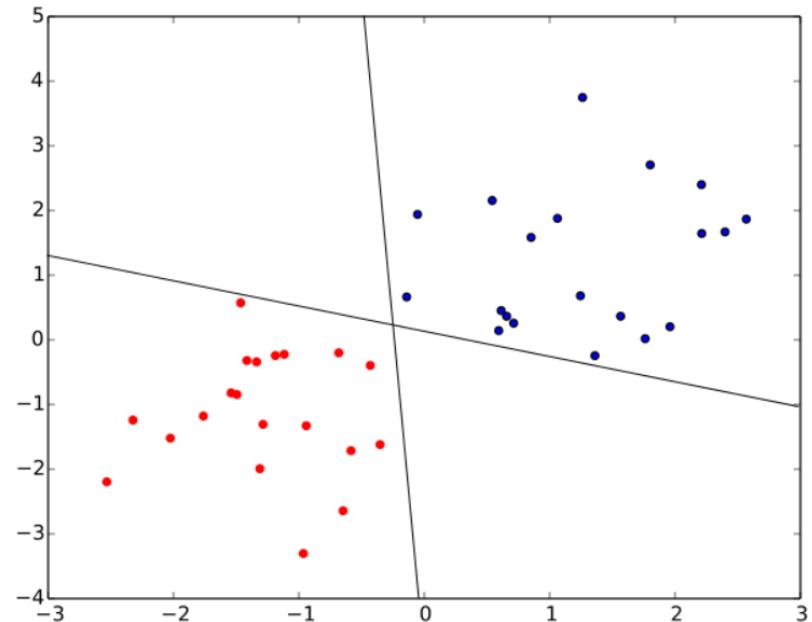
SVMs attempt to learn the support vectors which maximize the margin between classes.

A perceptron does not.

Both of these perceptron classifiers are equivalent.

‘Perceptron of optimal stability’ is used in SVM:

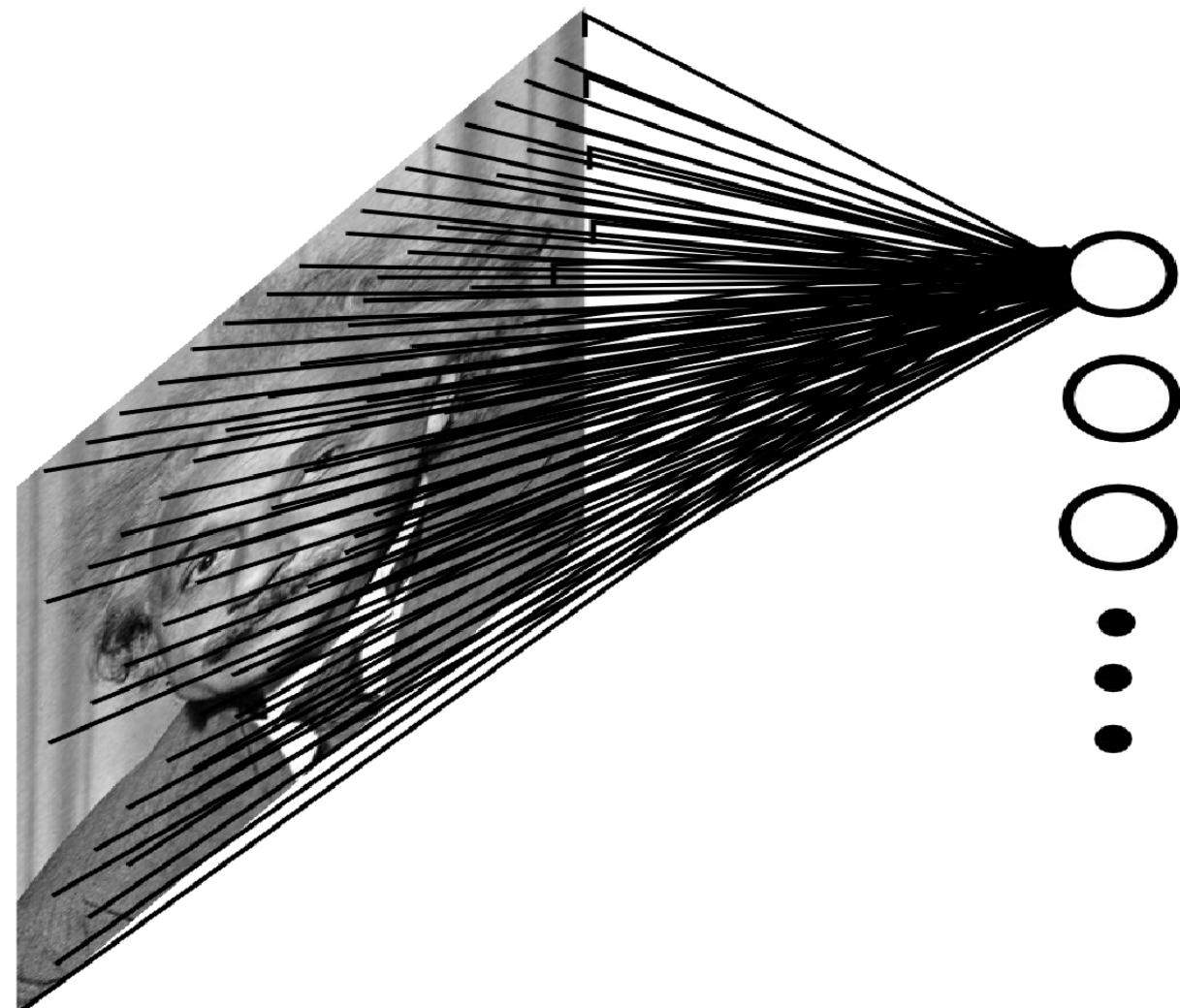
Perceptron
+ optimal stability
+ kernel trick
= *foundations of SVM*



Outline

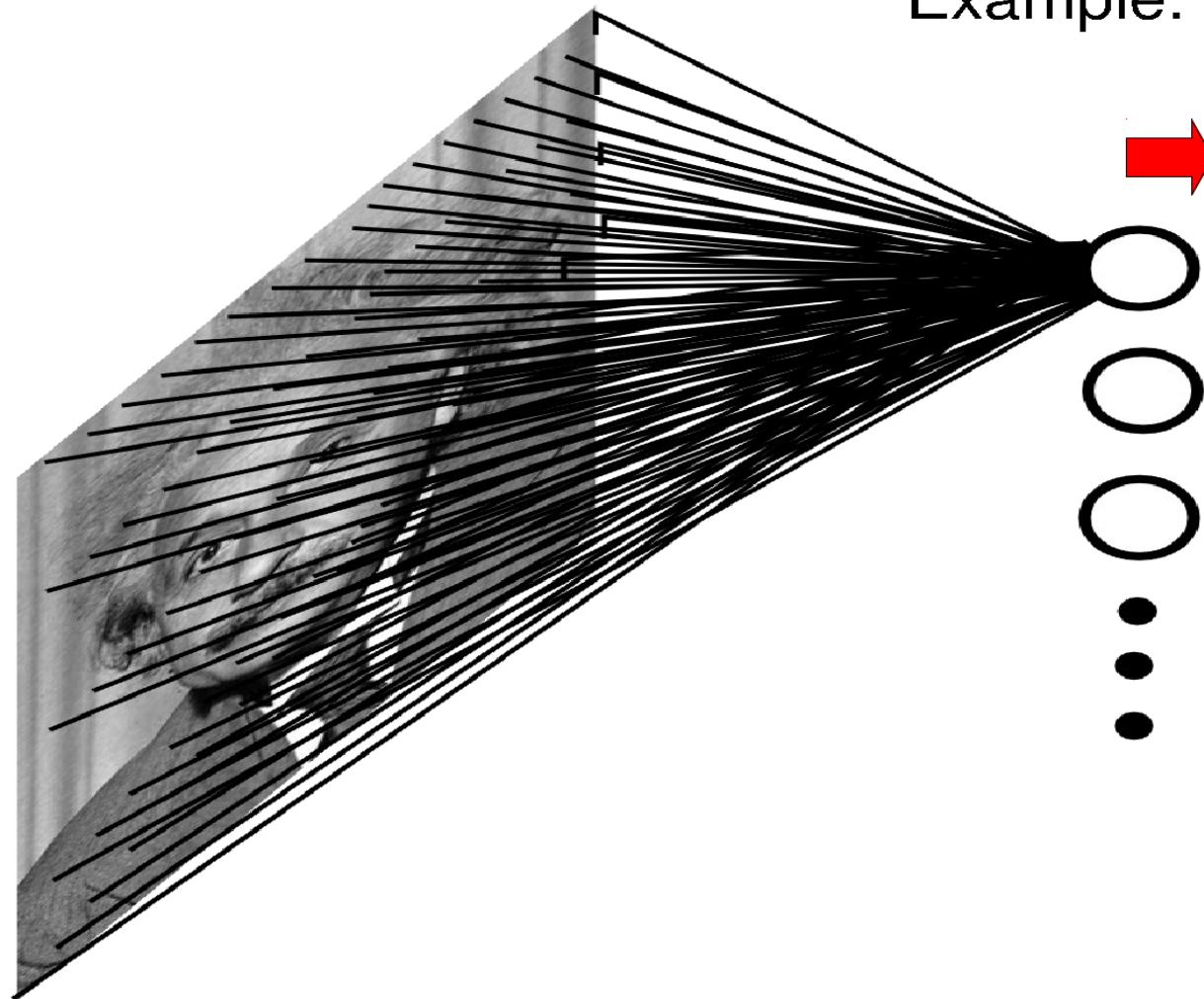
- Supervised Neural Networks
- Convolutional Neural Networks
- Examples
- Tips

Images as input to neural networks

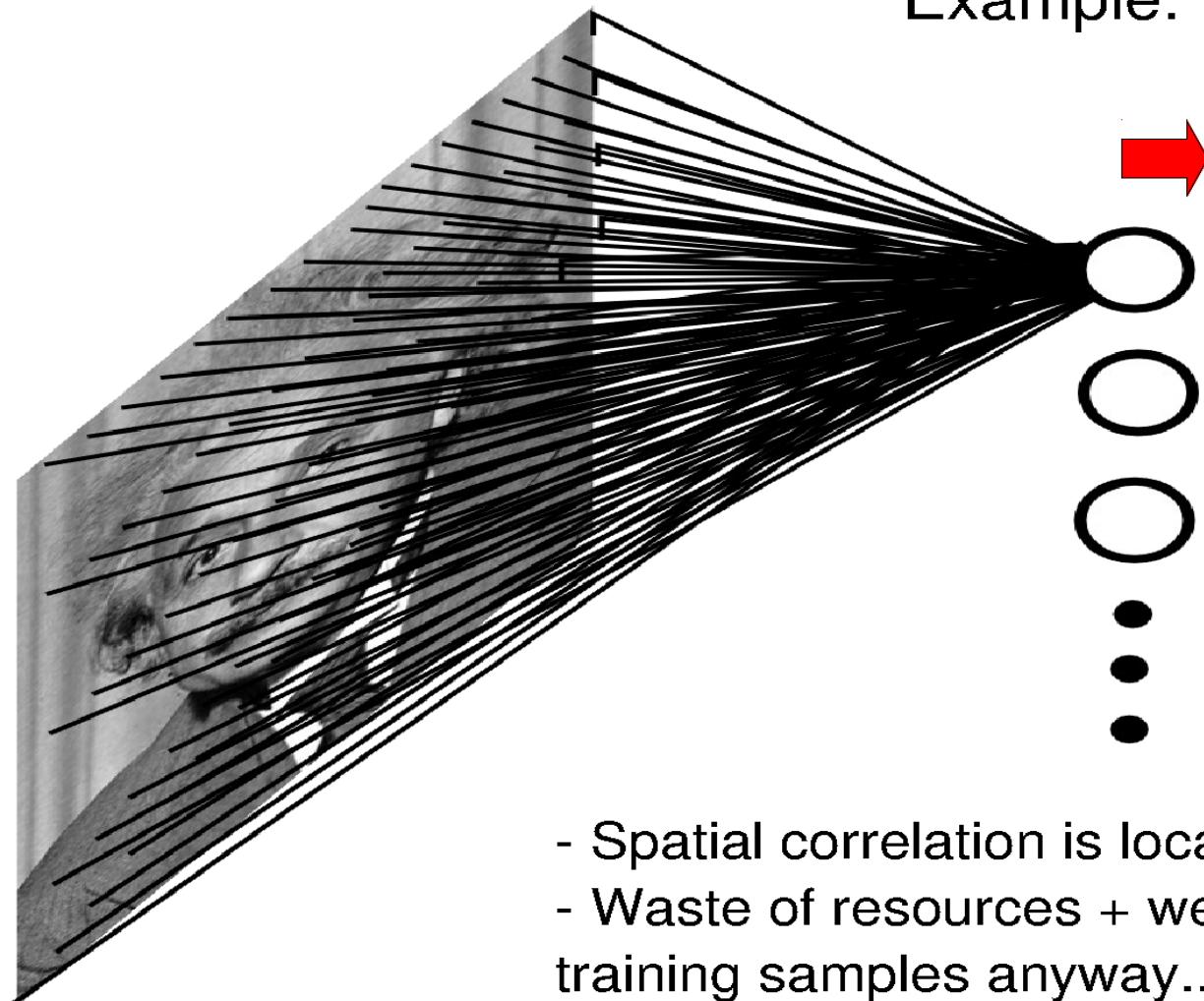


Images as input to neural networks

Example: 200x200 image
40K hidden units
→ ~2B parameters!!!



Images as input to neural networks

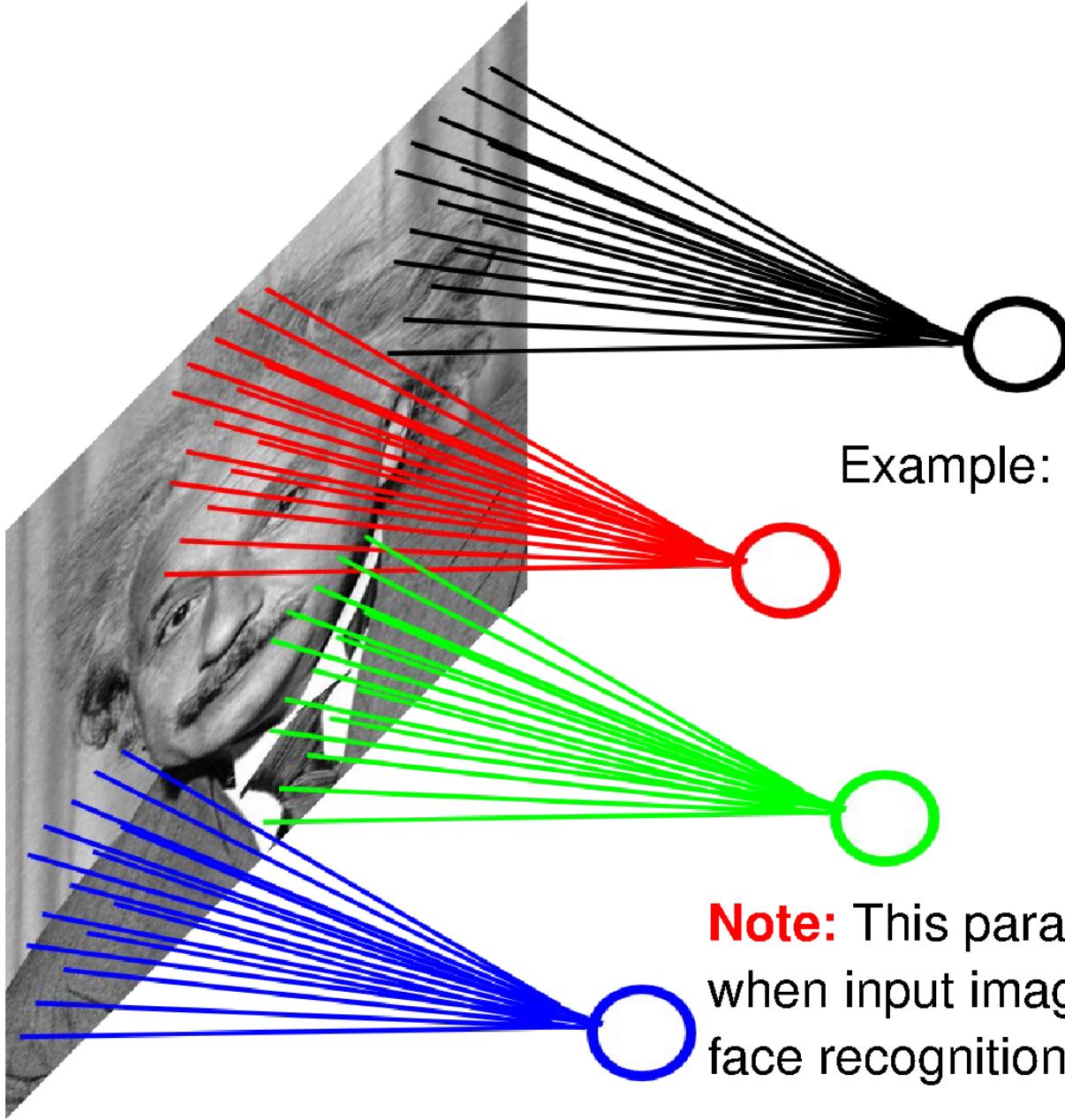


33

Ranzato

Motivation

- Sparse interactions – *receptive fields*
 - Assume that in an image, we care about ‘local neighborhoods’ only for a given neural network layer.
 - Composition of layers will expand local -> global.



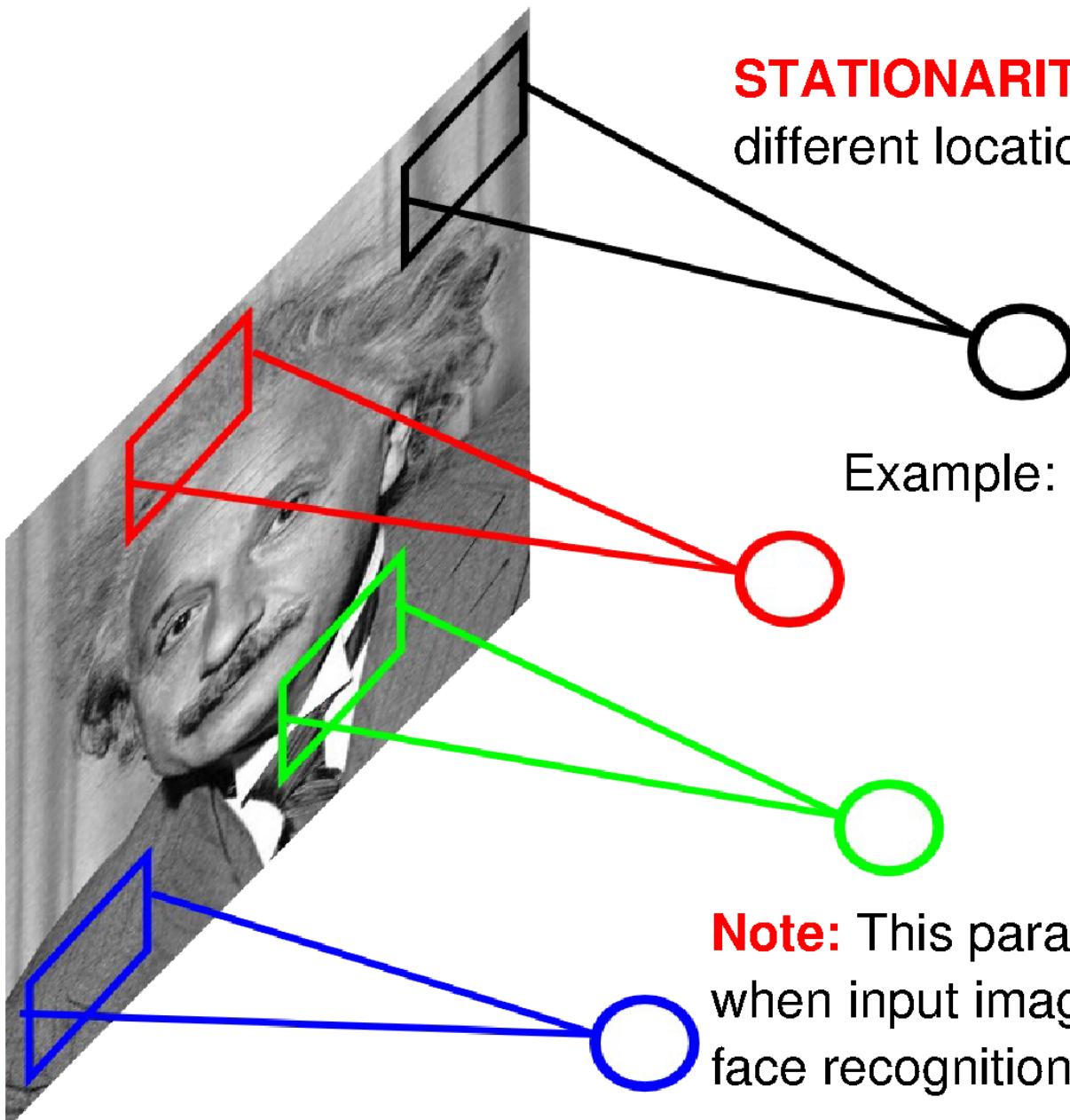
Example: 200x200 image
40K hidden units
Filter size: 10x10
4M parameters

Note: This parameterization is good
when input image is registered (e.g.,
face recognition).

34

Ranzato 

STATIONARITY? Statistics is similar at different locations



Example:
200x200 image
40K hidden units
Filter size: 10x10
4M parameters

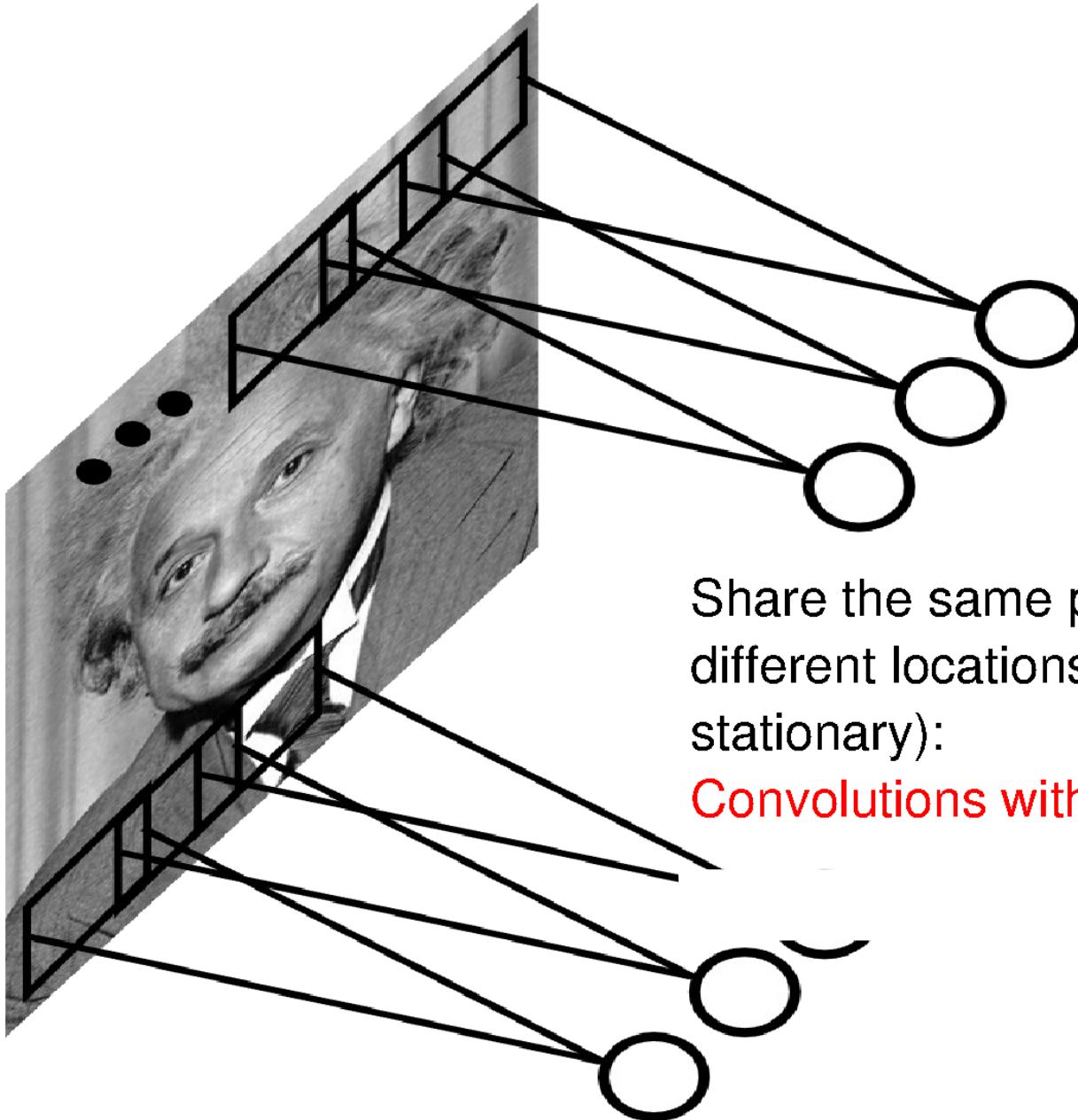
Note: This parameterization is good when input image is registered (e.g.,
face recognition).

35

Ranzato

Motivation

- Sparse interactions – *receptive fields*
 - Assume that in an image, we care about ‘local neighborhoods’ only for a given neural network layer.
 - Composition of layers will expand local -> global.
- Parameter sharing
 - ‘Tied weights’ – use same weights for more than one perceptron in the neural network.
 - Leads to *equivariant representation*
 - If input changes (e.g., translates), then output changes similarly



Share the same parameters across
different locations (assuming input is
stationary):

Convolutions with learned kernels

Filtering reminder: Correlation (rotated convolution)

$$f[\cdot, \cdot] \quad \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$$I[., .]$$

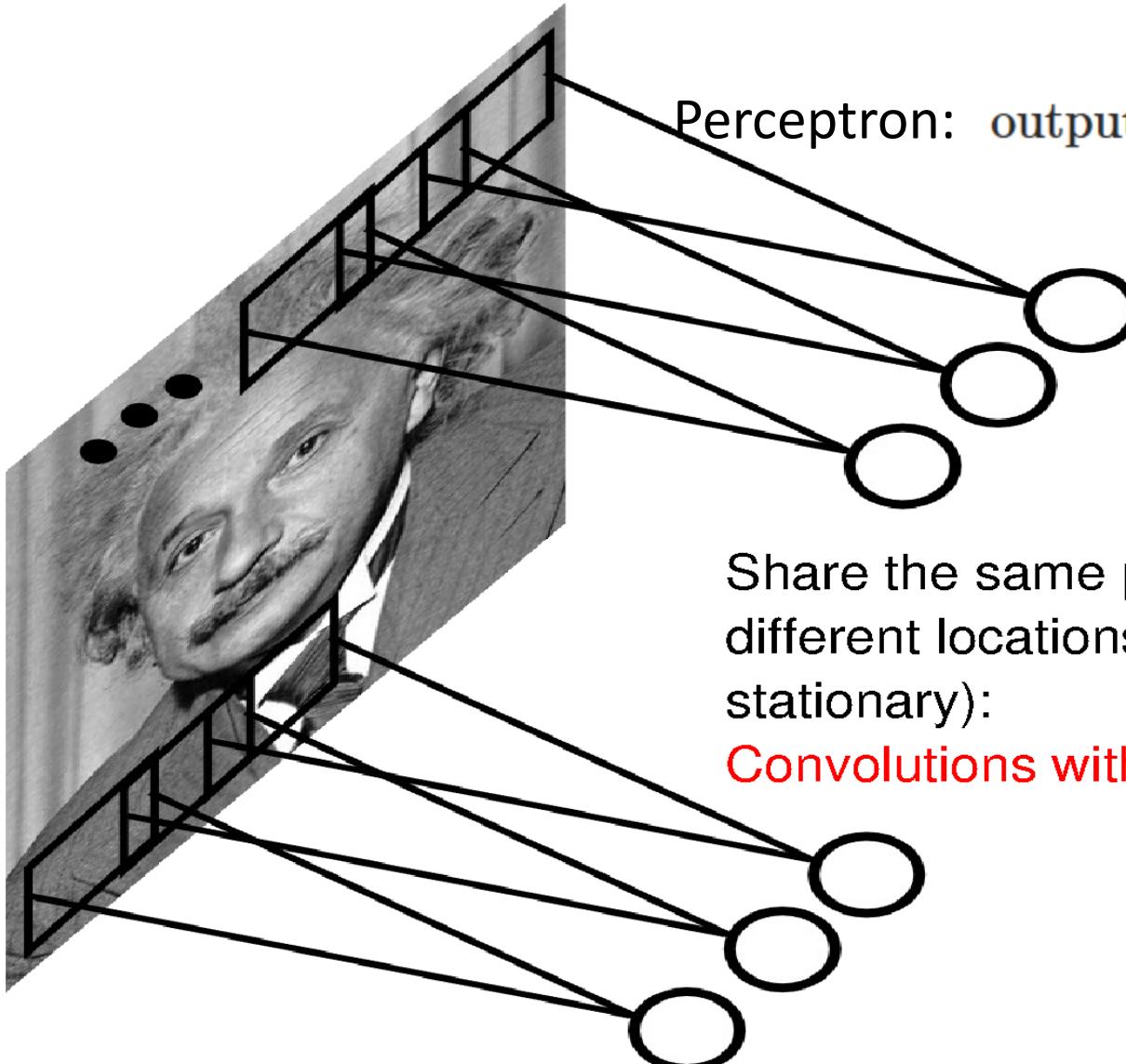
$$h[., .]$$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

	0	10	20	30	30	30	20	10	
	0	20	40	60	60	60	40	20	
	0	30	60	90	90	90	60	30	
	0	30	50	80	80	90	60	30	
	0	30	50	80	80	90	60	30	
	0	20	30	50	50	60	40	20	
	10	20	30	30	30	30	20	10	
	10	10	10	0	0	0	0	0	

$$h[m, n] = \sum_{k,l} f[k, l] I[m + k, n + l]$$

Convolutional Layer



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

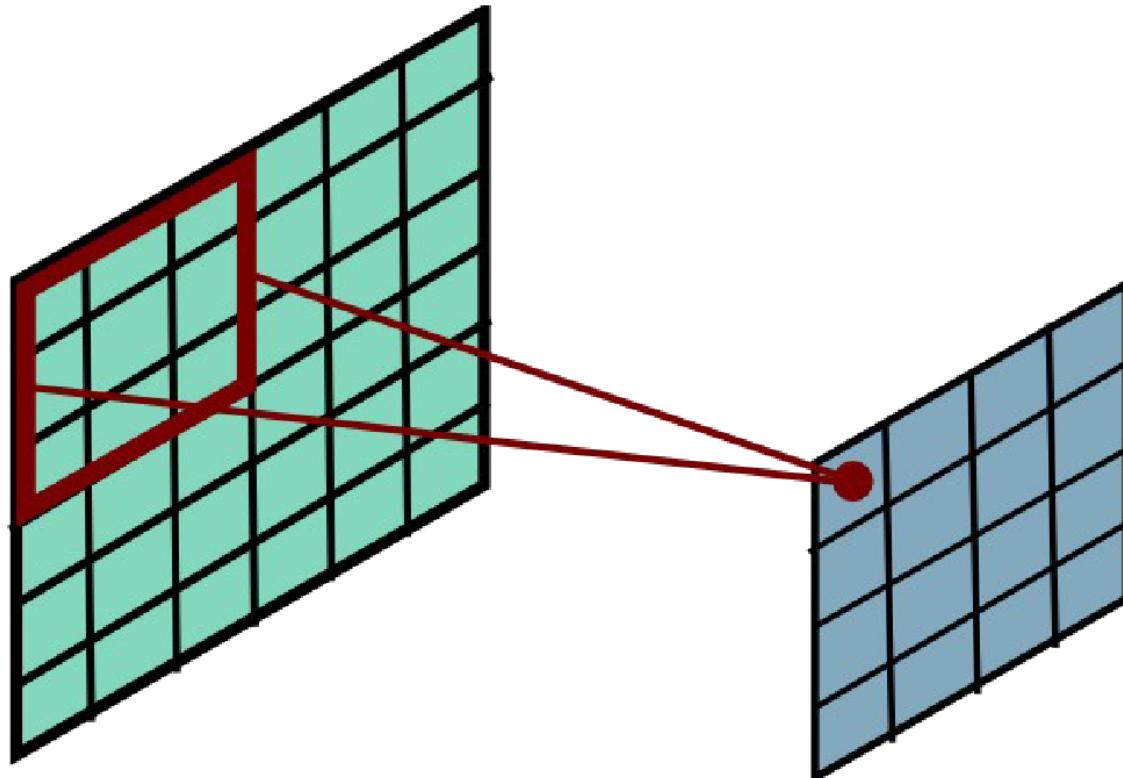
$$w \cdot x \equiv \sum_j w_j x_j,$$

This is convolution!

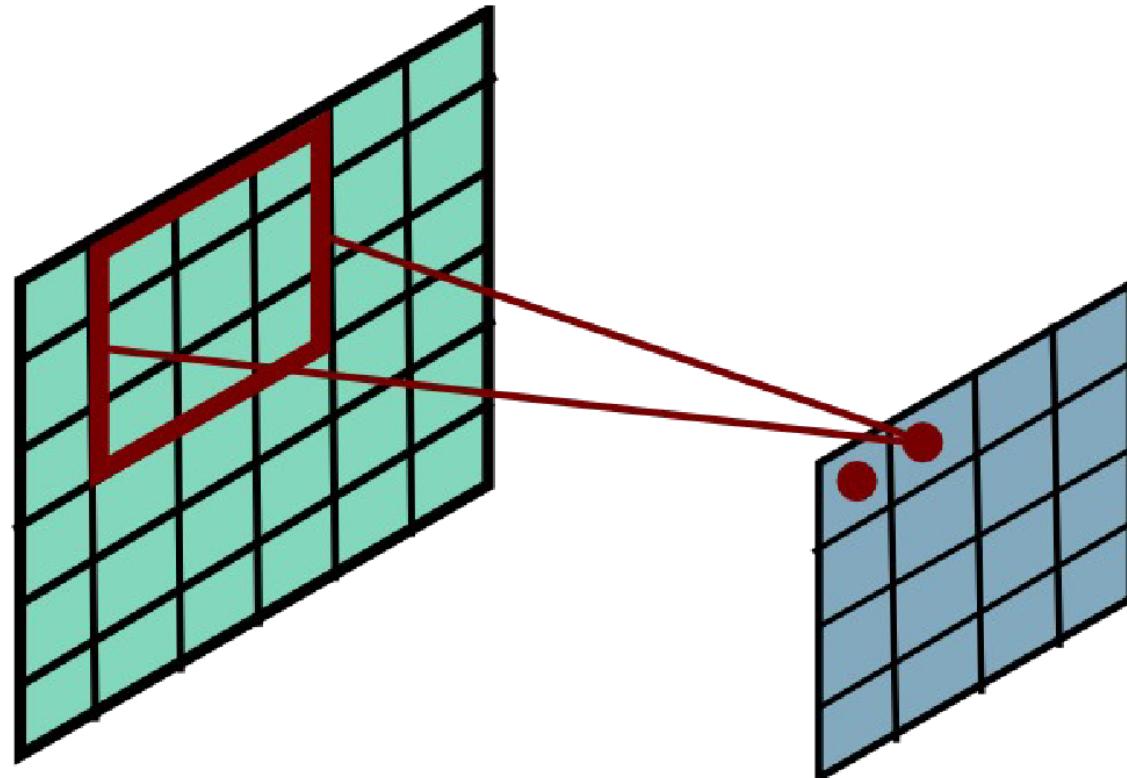
Share the same parameters across different locations (assuming input is stationary):

Convolutions with learned kernels

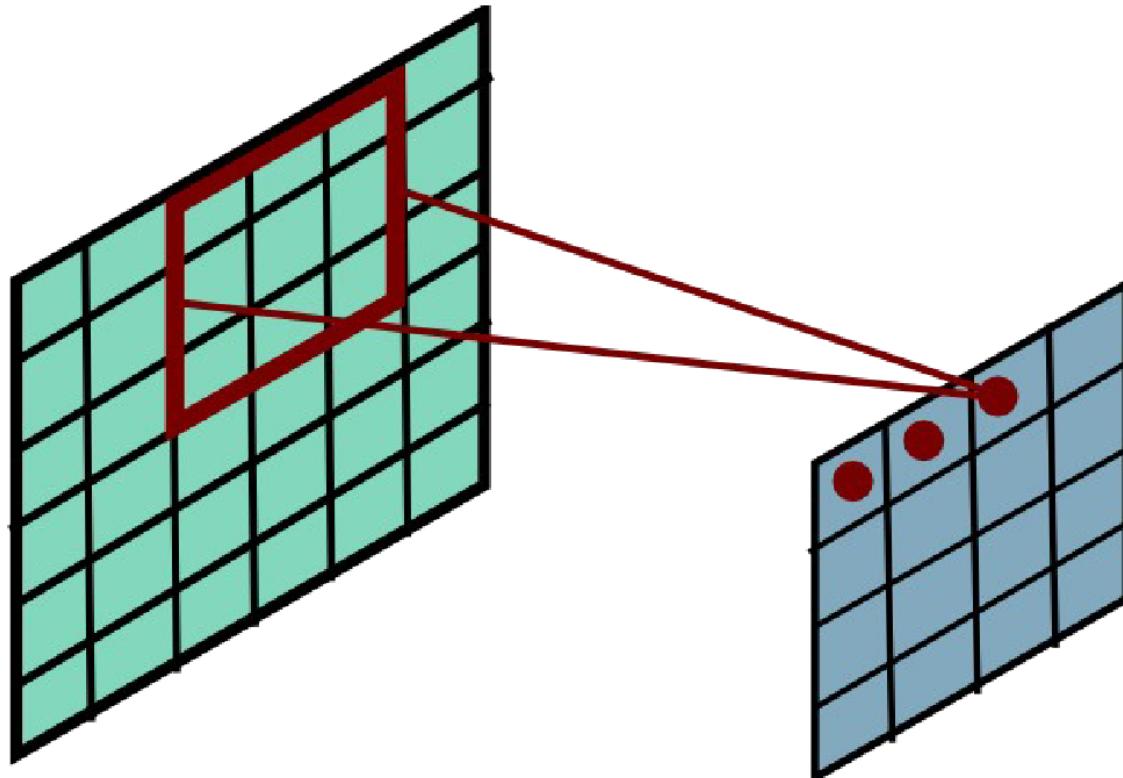
Convolutional Layer



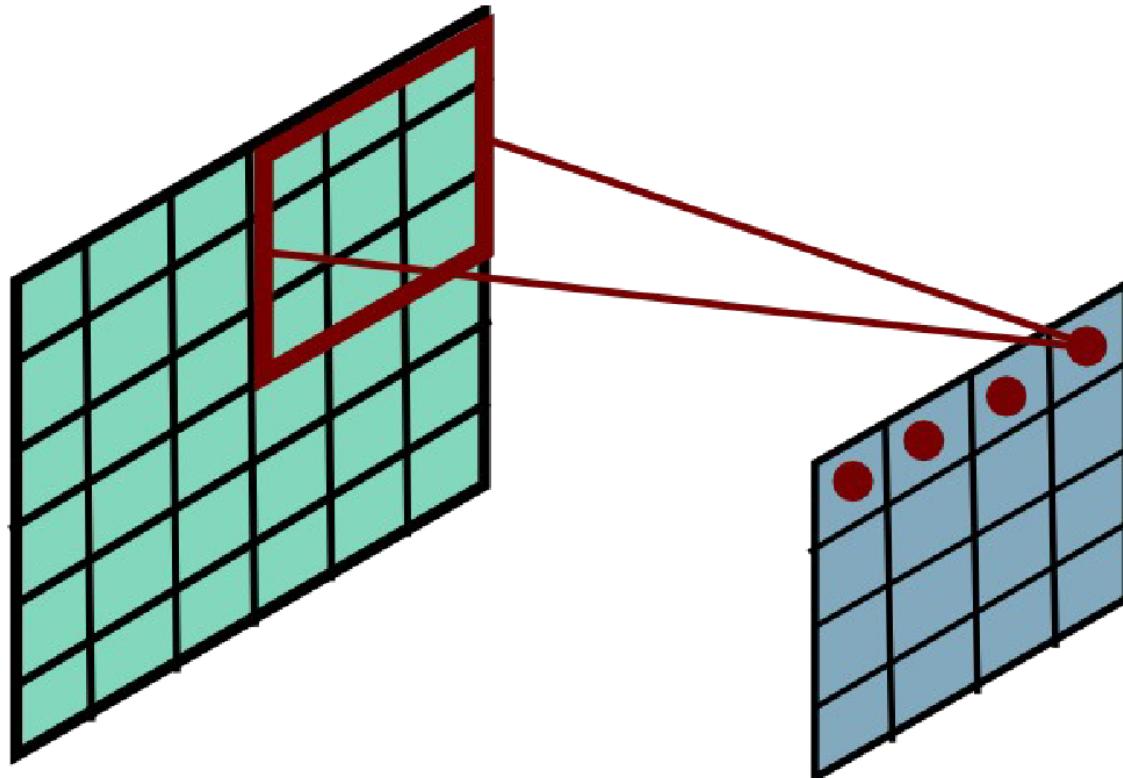
Convolutional Layer



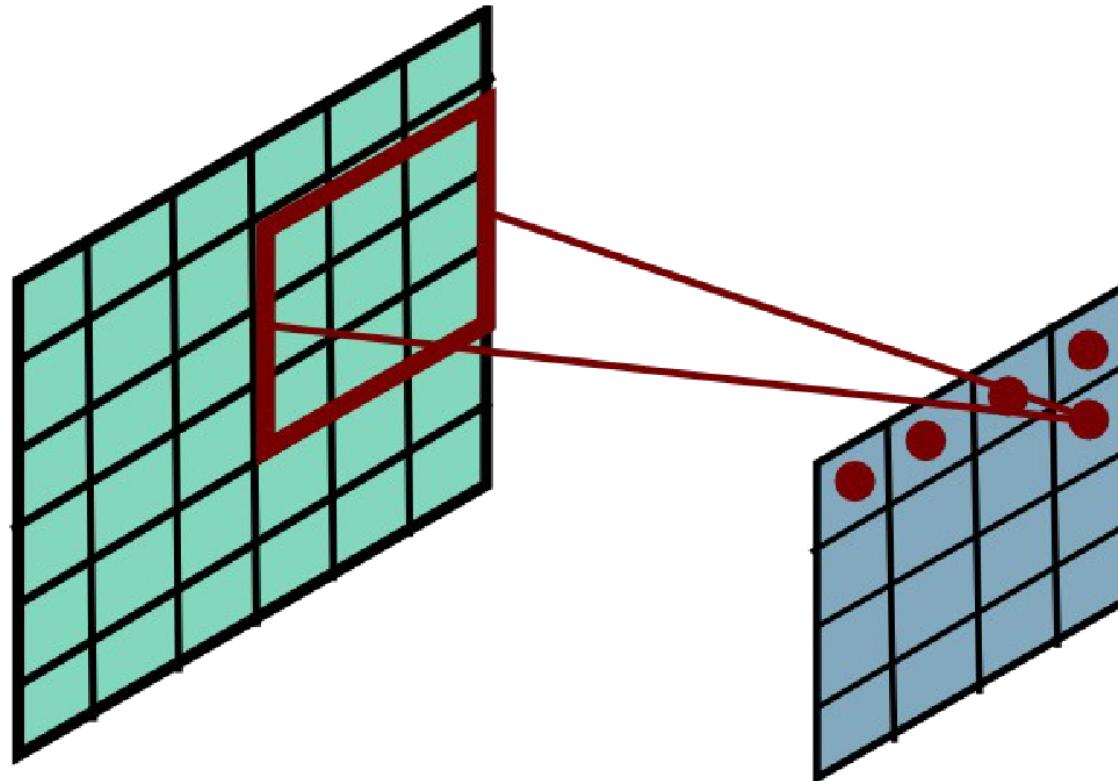
Convolutional Layer



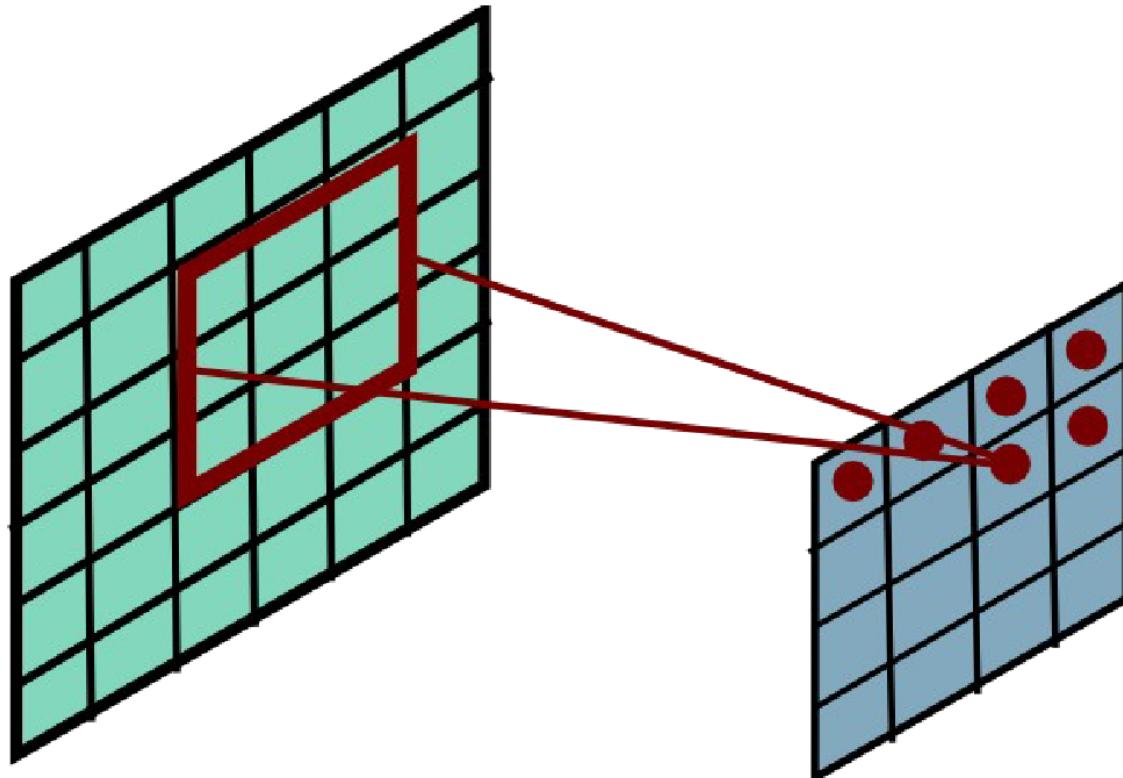
Convolutional Layer



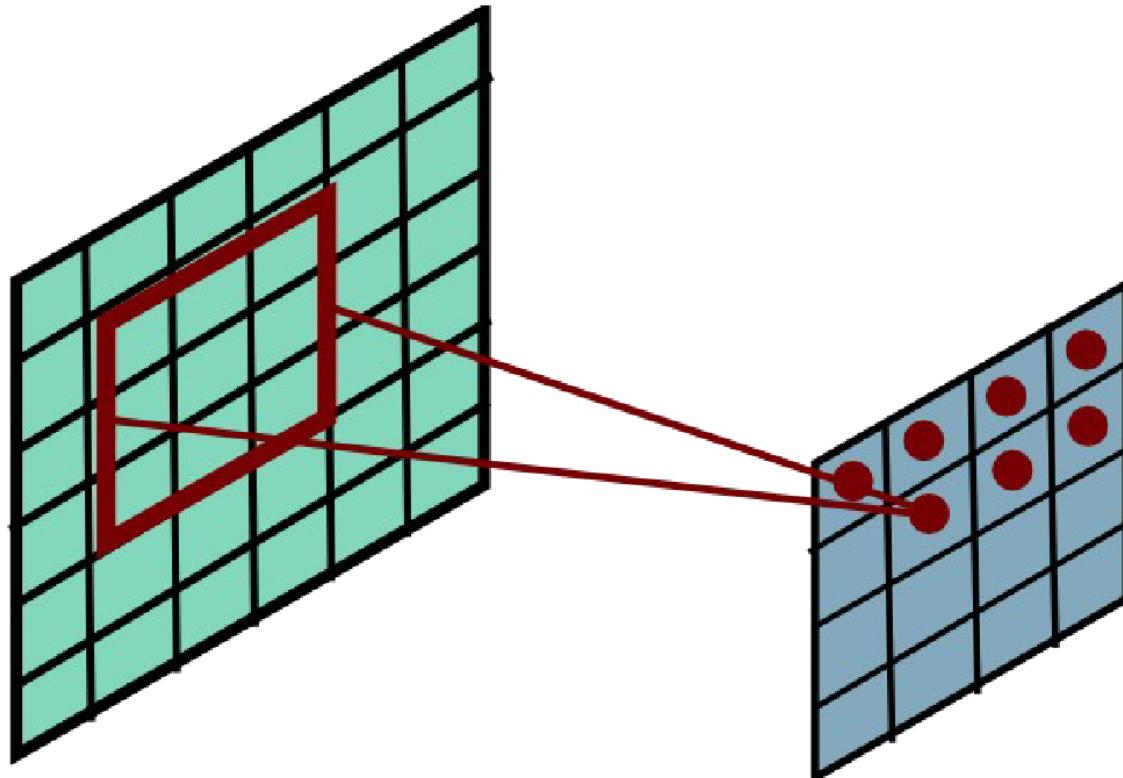
Convolutional Layer



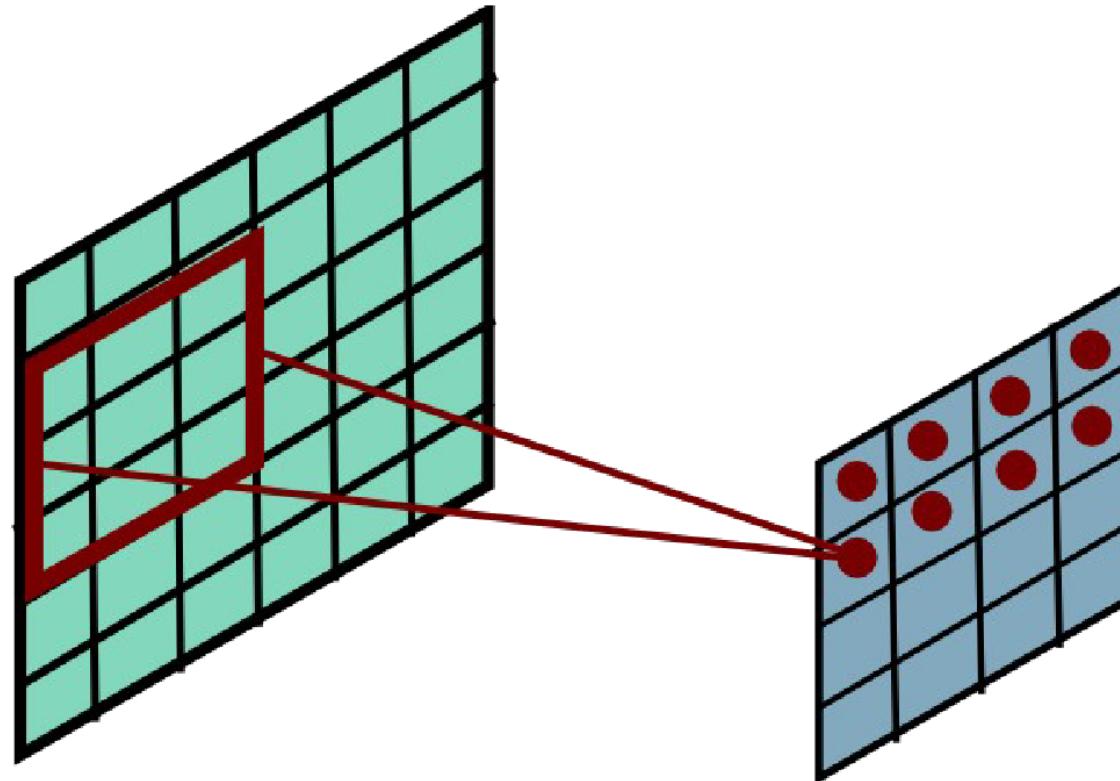
Convolutional Layer



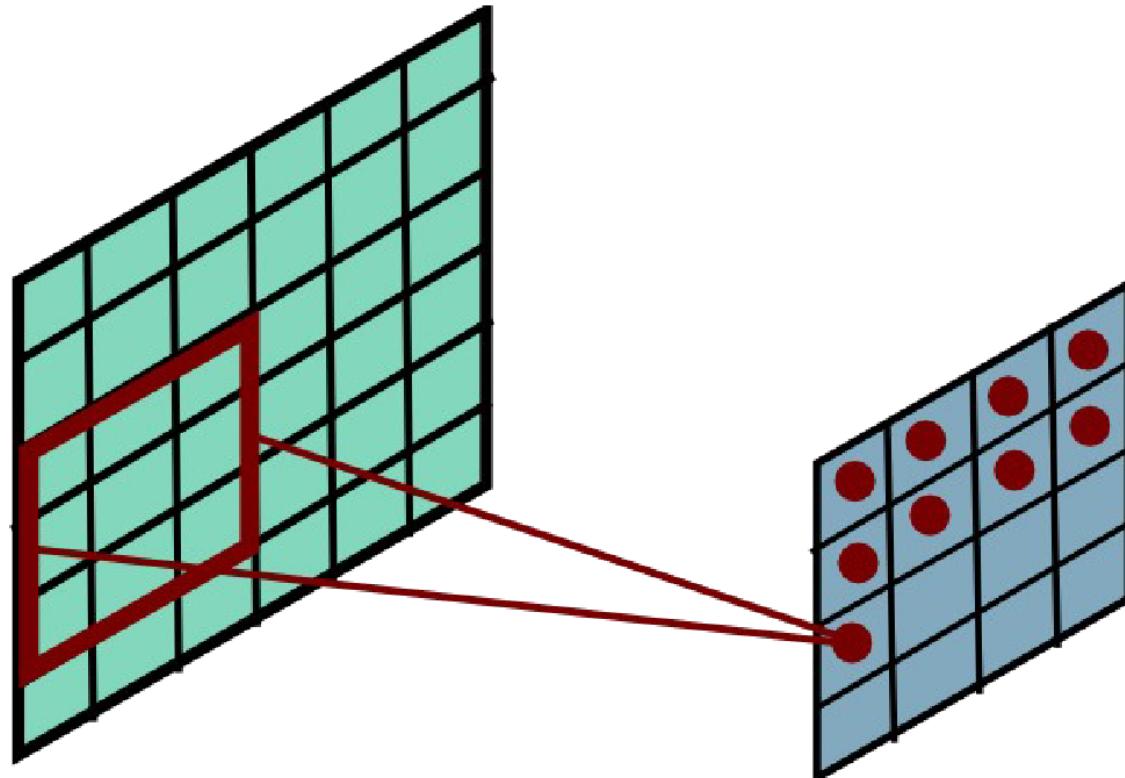
Convolutional Layer



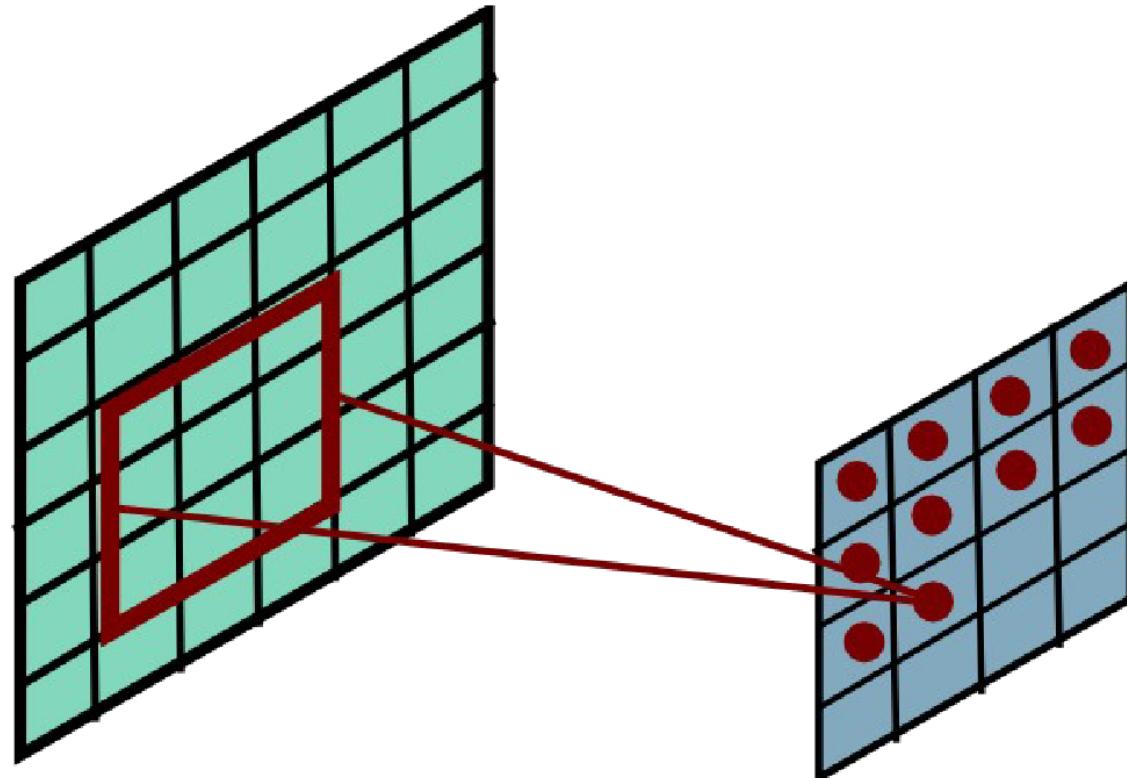
Convolutional Layer



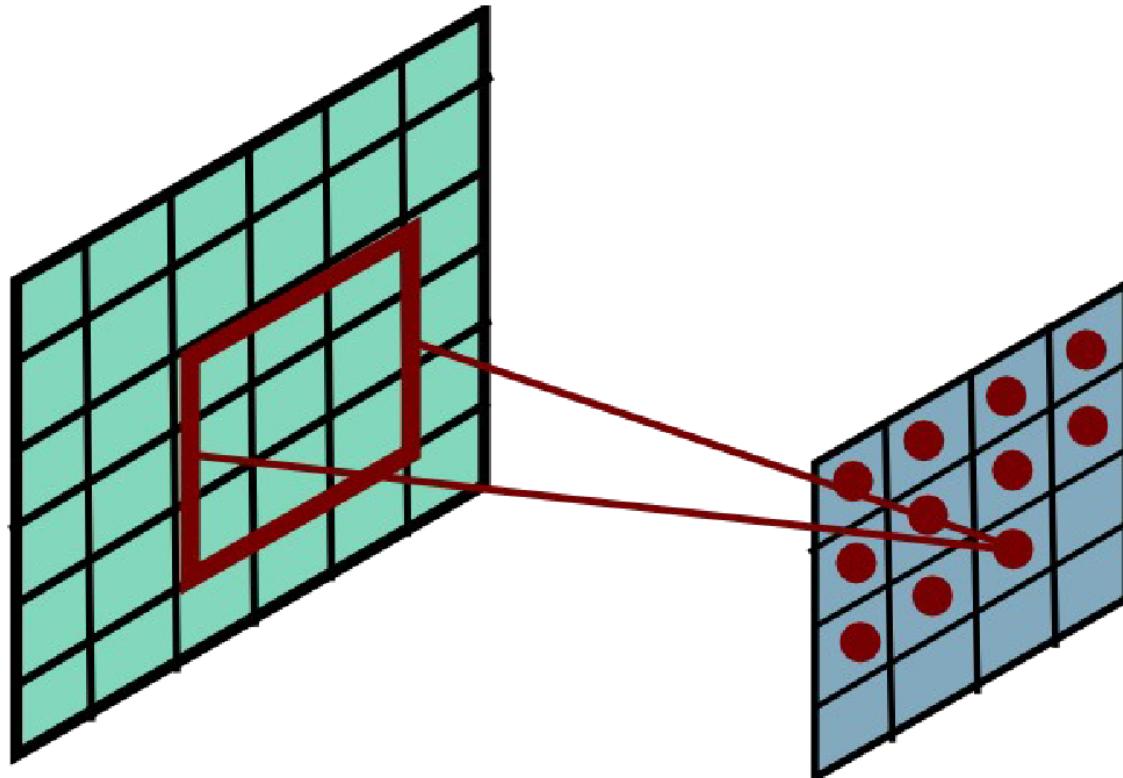
Convolutional Layer



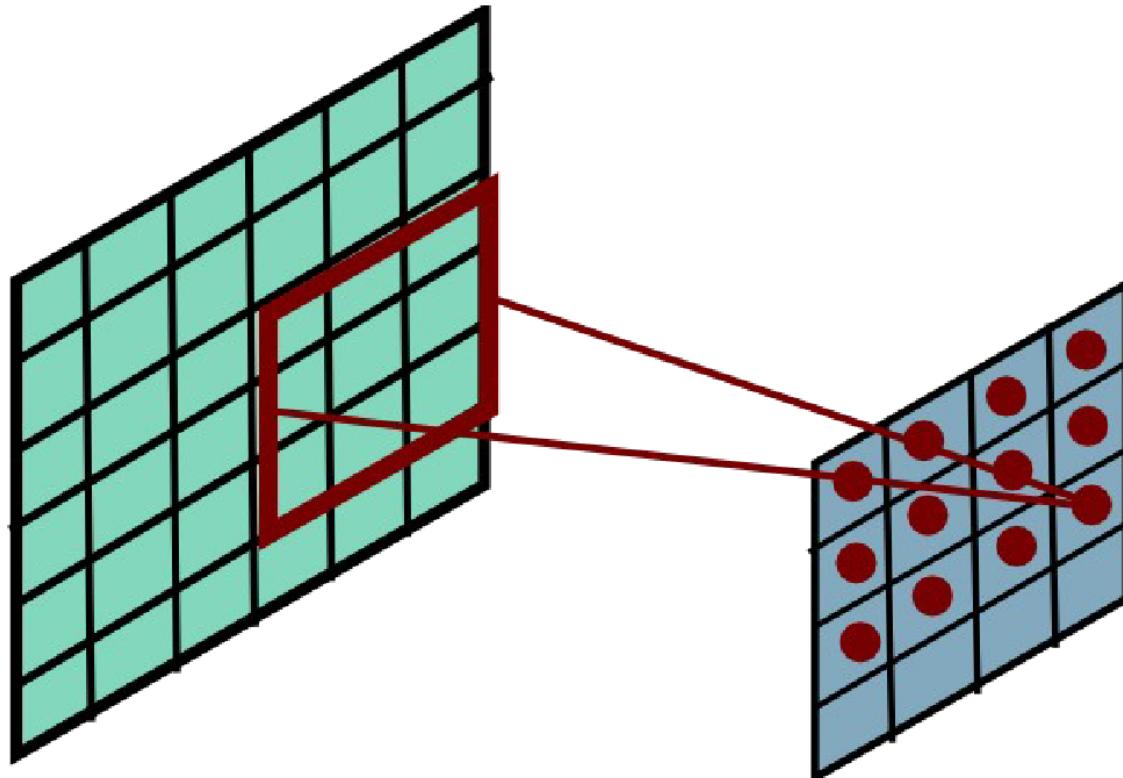
Convolutional Layer



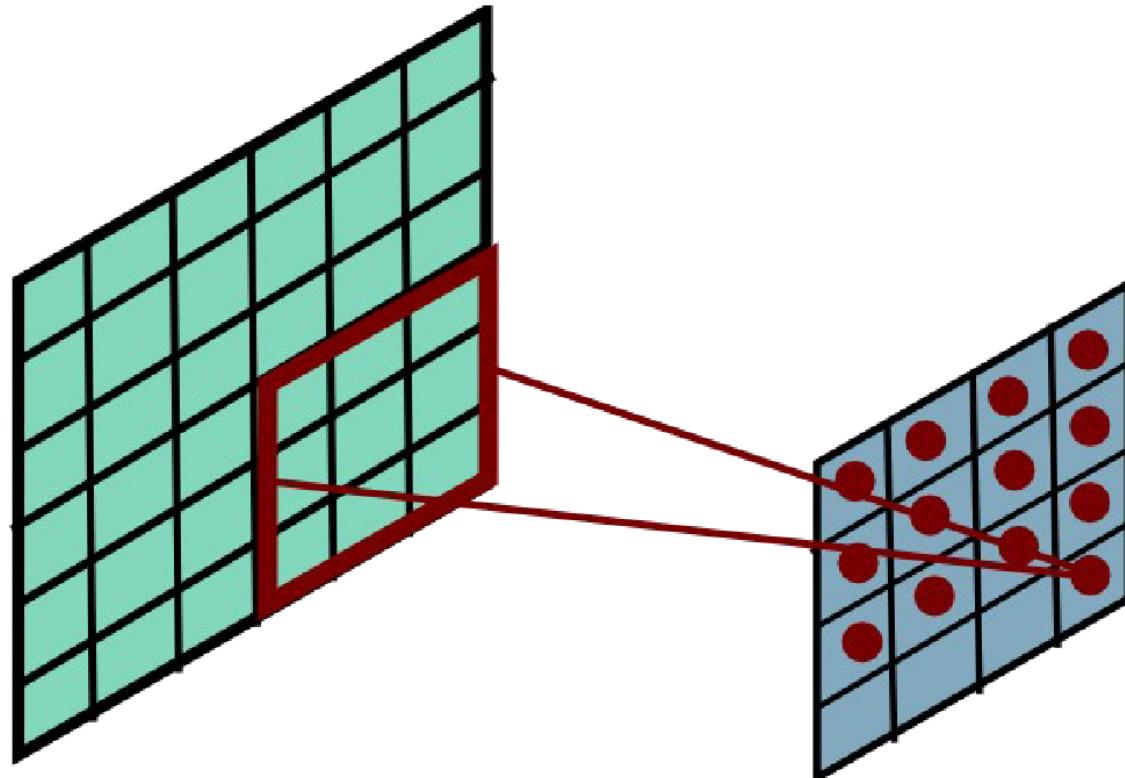
Convolutional Layer



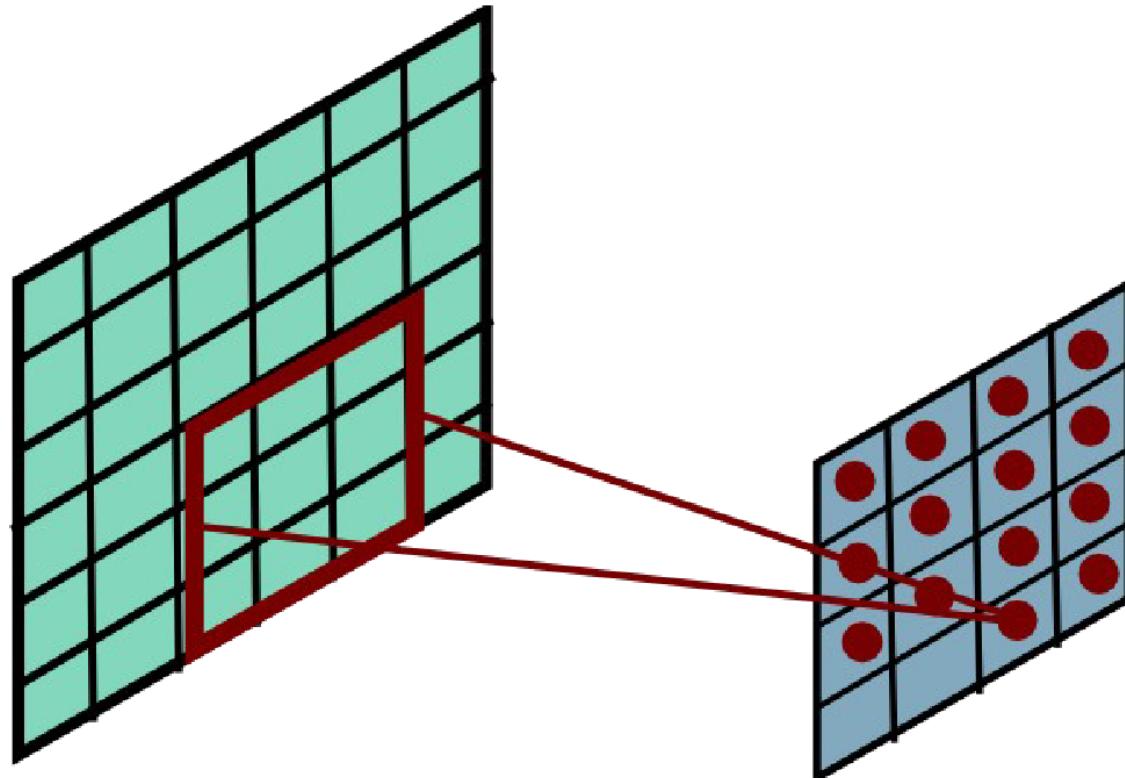
Convolutional Layer



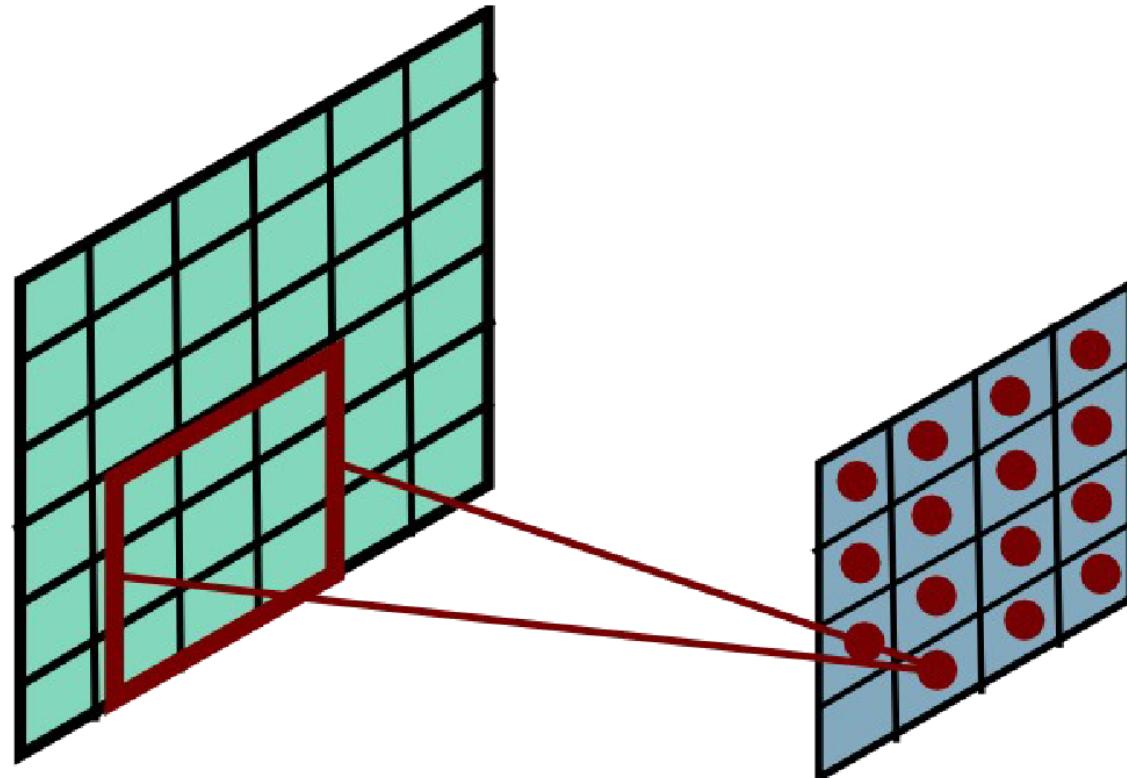
Convolutional Layer



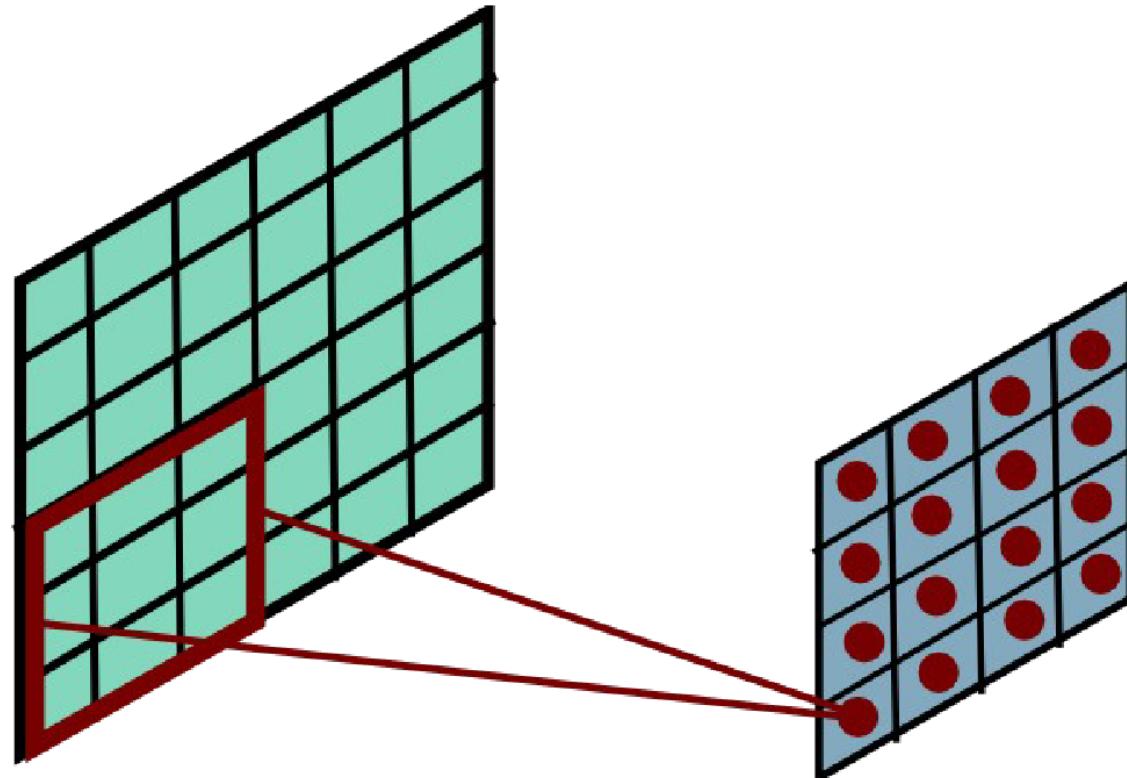
Convolutional Layer



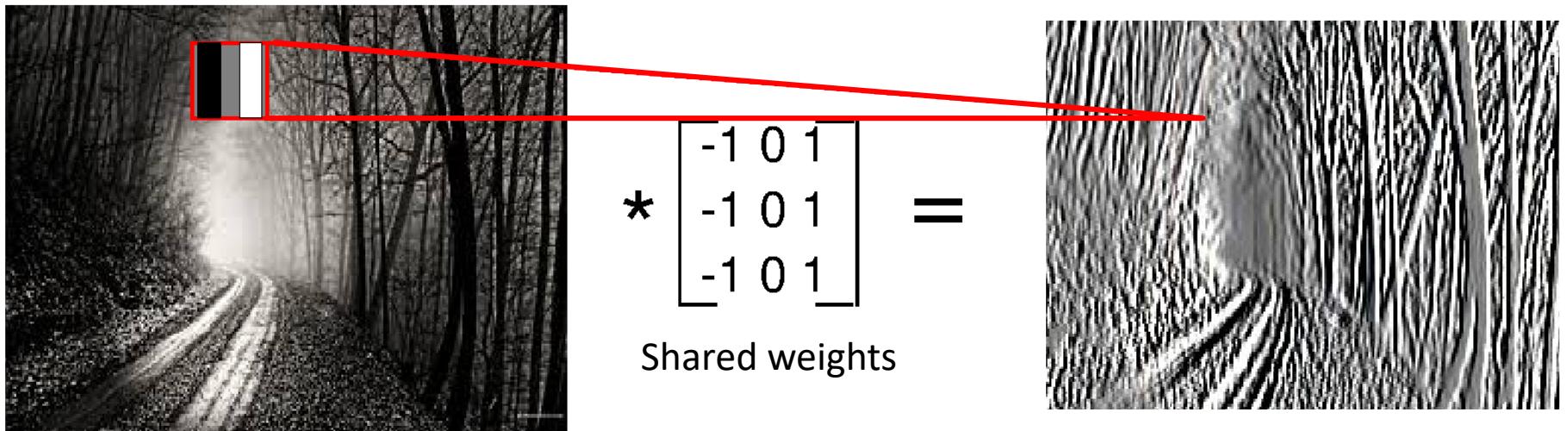
Convolutional Layer



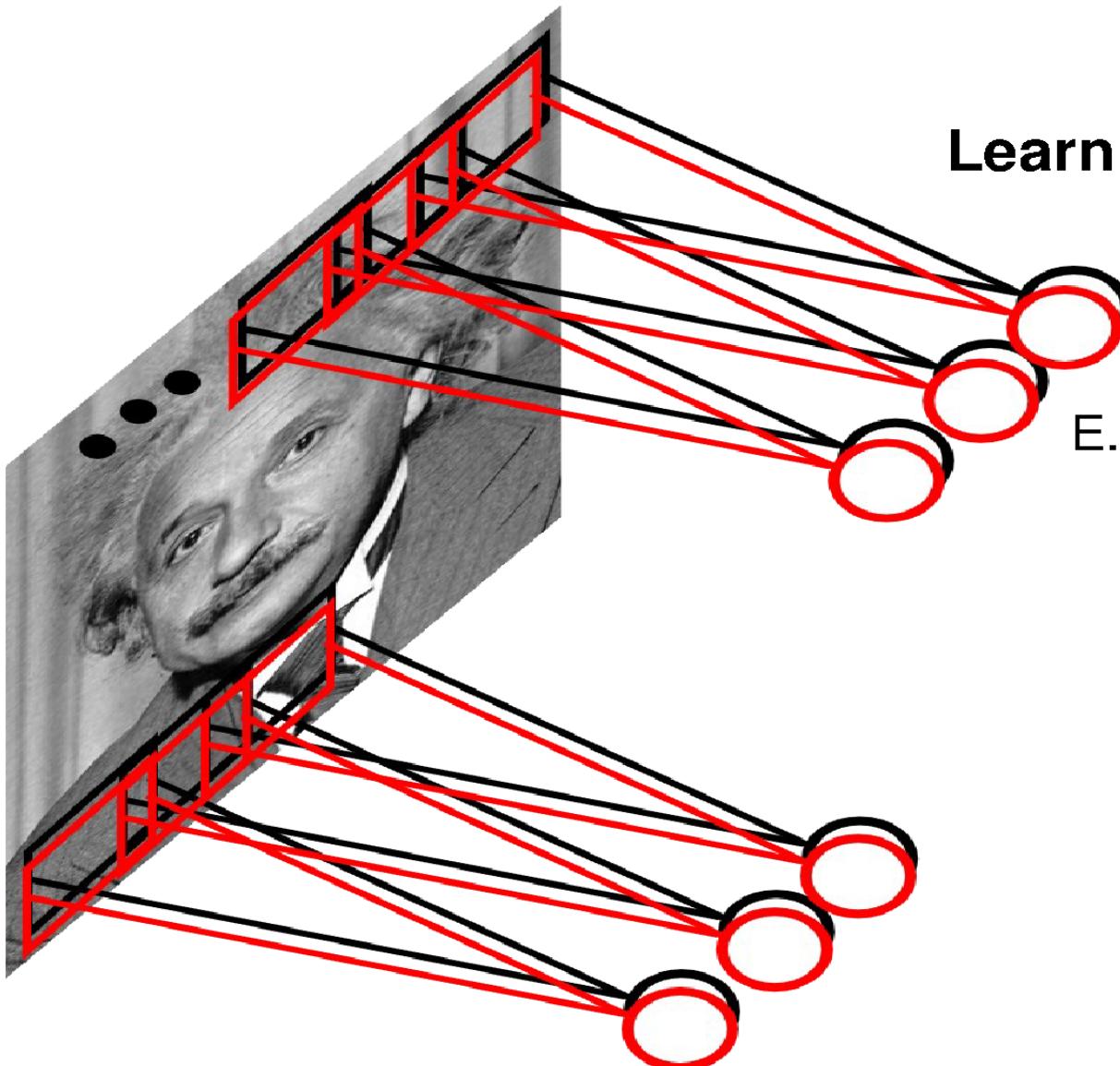
Convolutional Layer



Convolutional Layer



Convolutional Layer

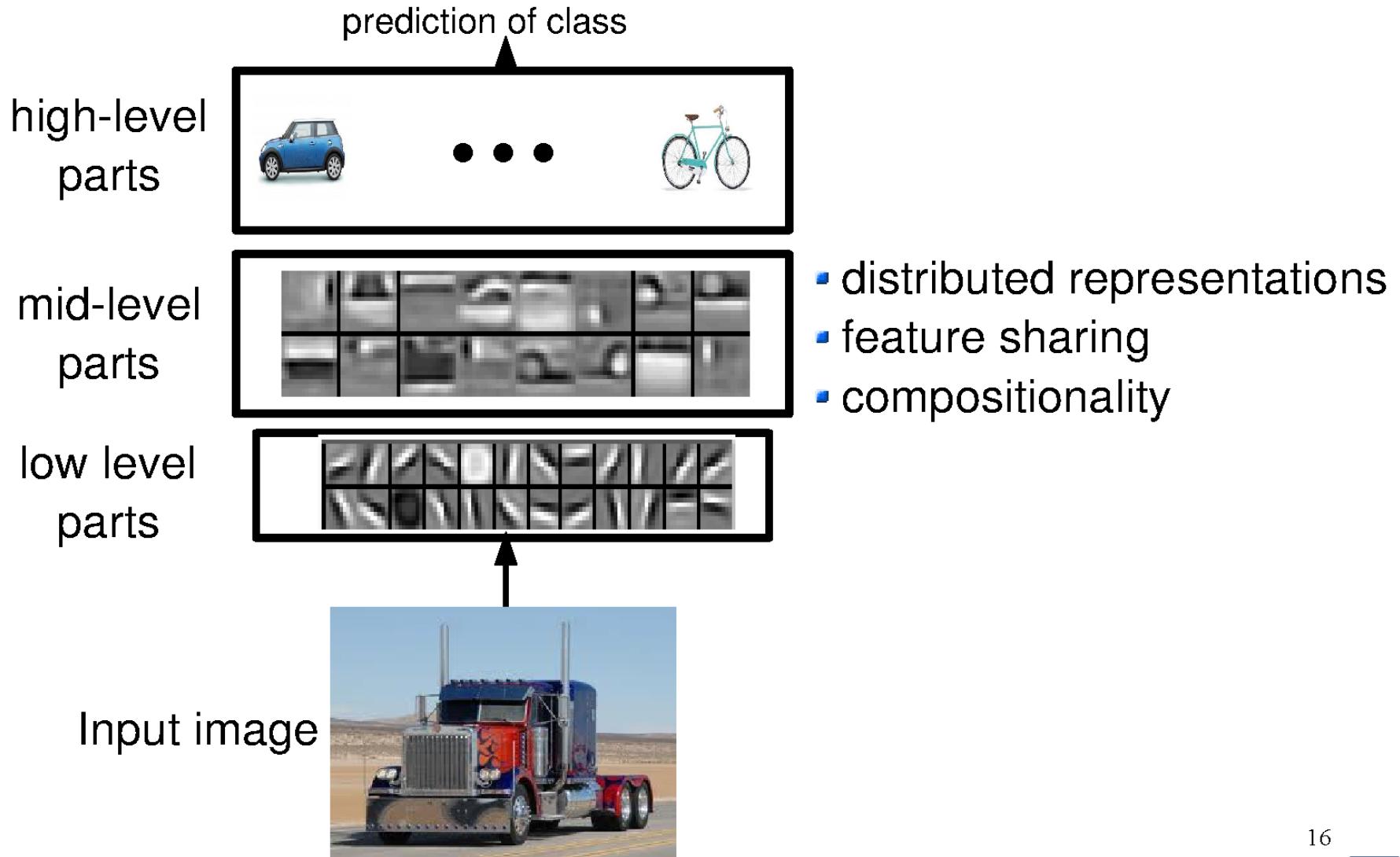


Learn multiple filters.

Filter = 'local' perceptron.
Also called *kernel*.

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

Interpretation



Lee et al. "Convolutional DBN's ..." ICML 2009

16

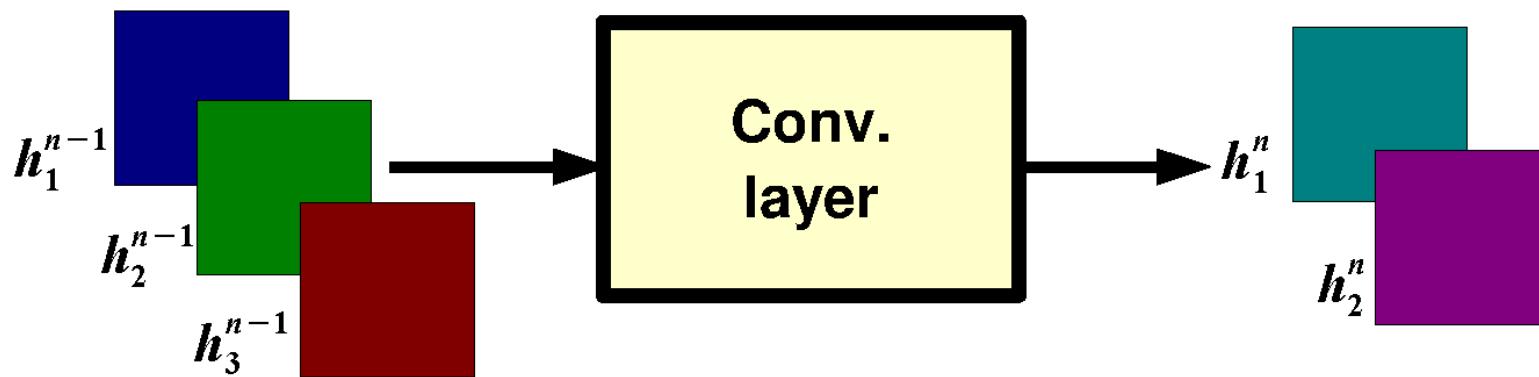
Ranzato

Convolutional Layer

$$h_j^n = \max \left(0, \sum_{k=1}^K h_k^{n-1} * w_{kj}^n \right)$$

output feature map input feature map kernel

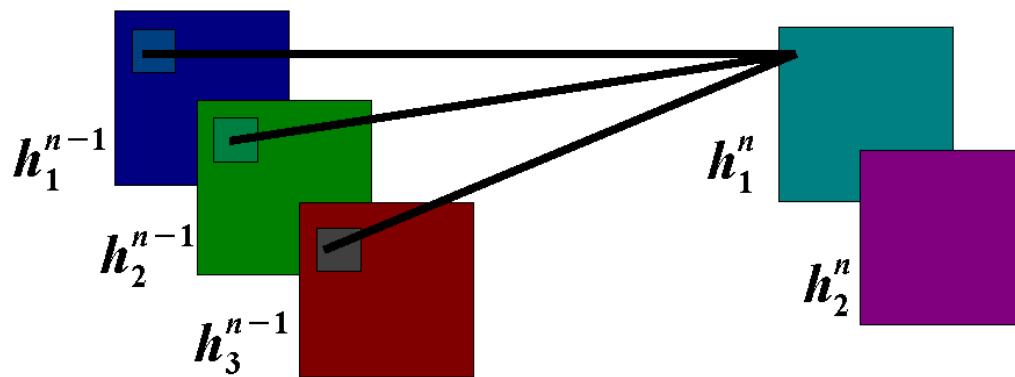
*n = layer number
K = kernel size
j = # channels (input) or # filters (depth)*



Convolutional Layer

$$h_j^n = \max \left(0, \sum_{k=1}^K h_k^{n-1} * w_{kj}^n \right)$$

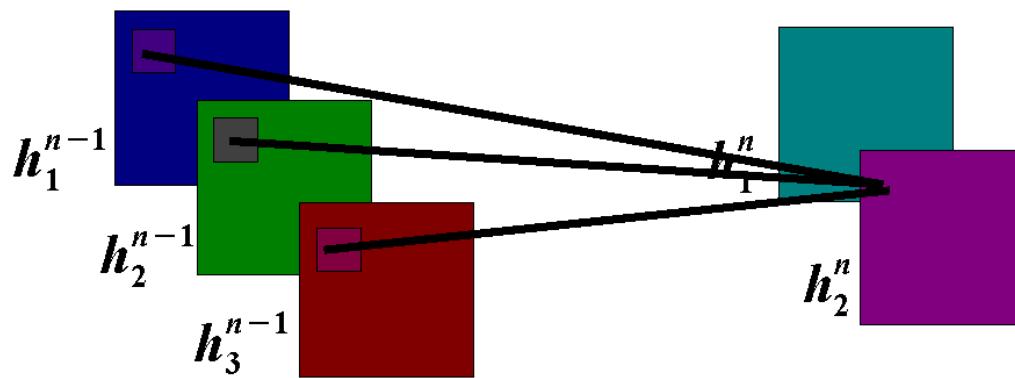
output feature map **input feature map** **kernel**



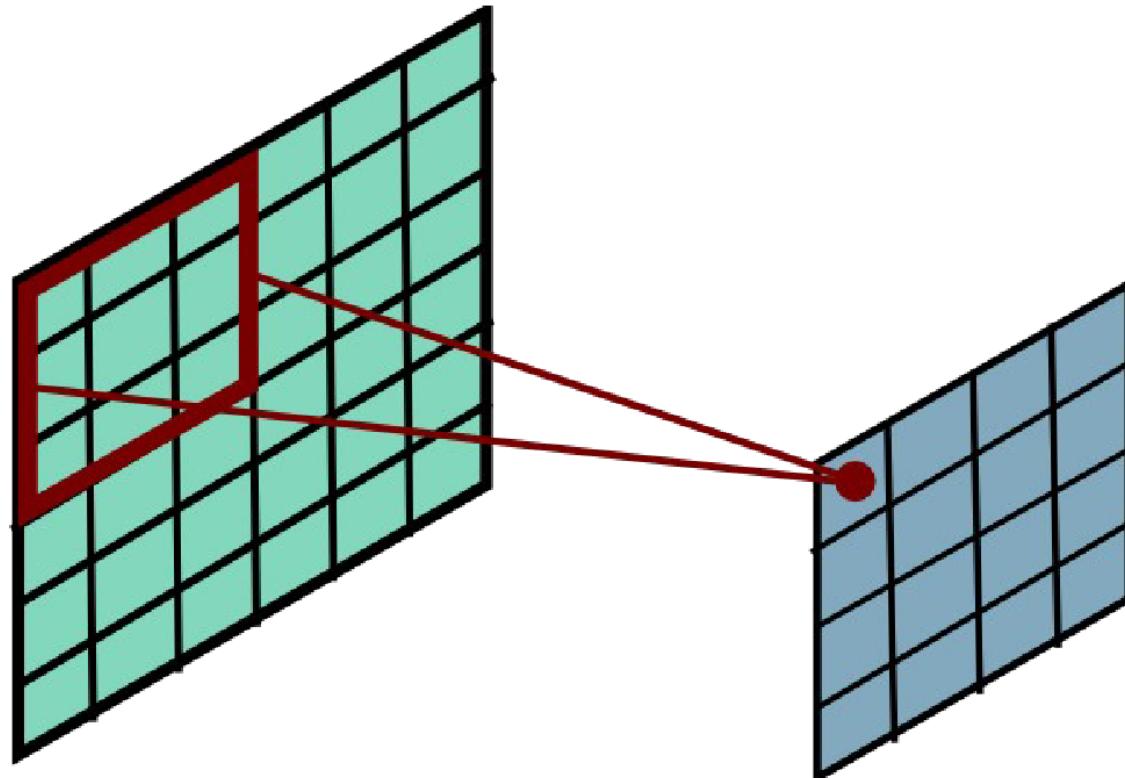
Convolutional Layer

$$h_j^n = \max \left(0, \sum_{k=1}^K h_k^{n-1} * w_{kj}^n \right)$$

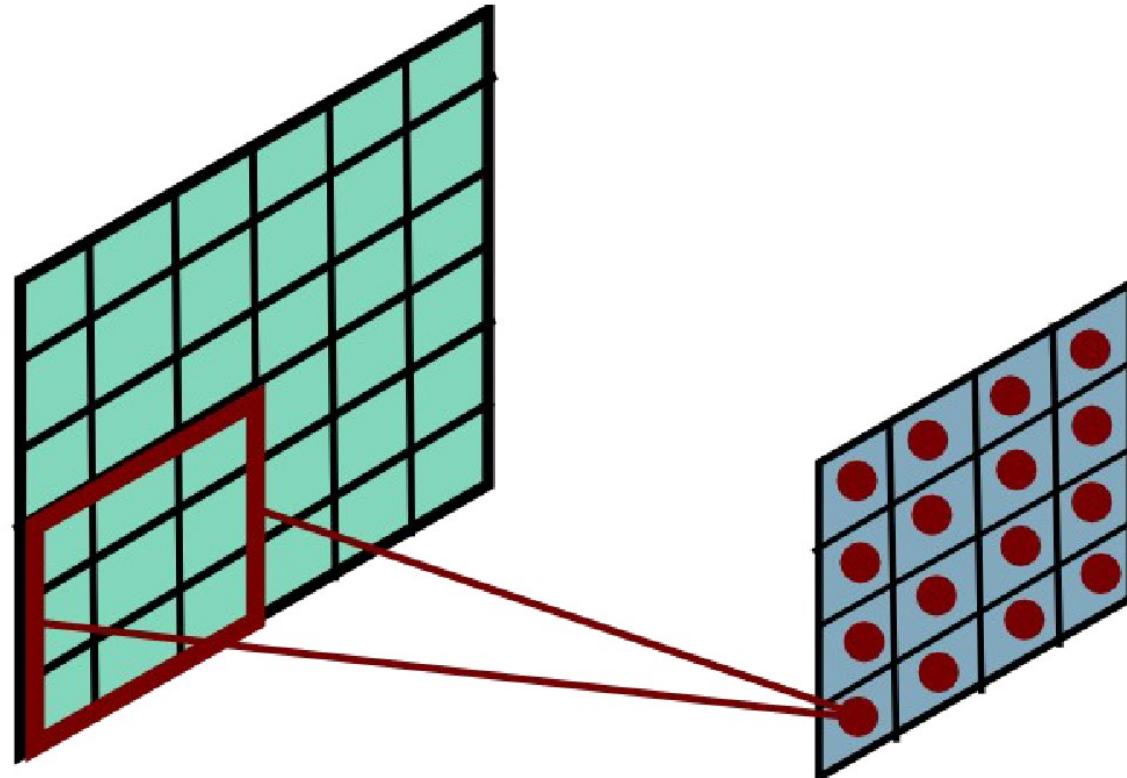
output feature map **input feature map** **kernel**



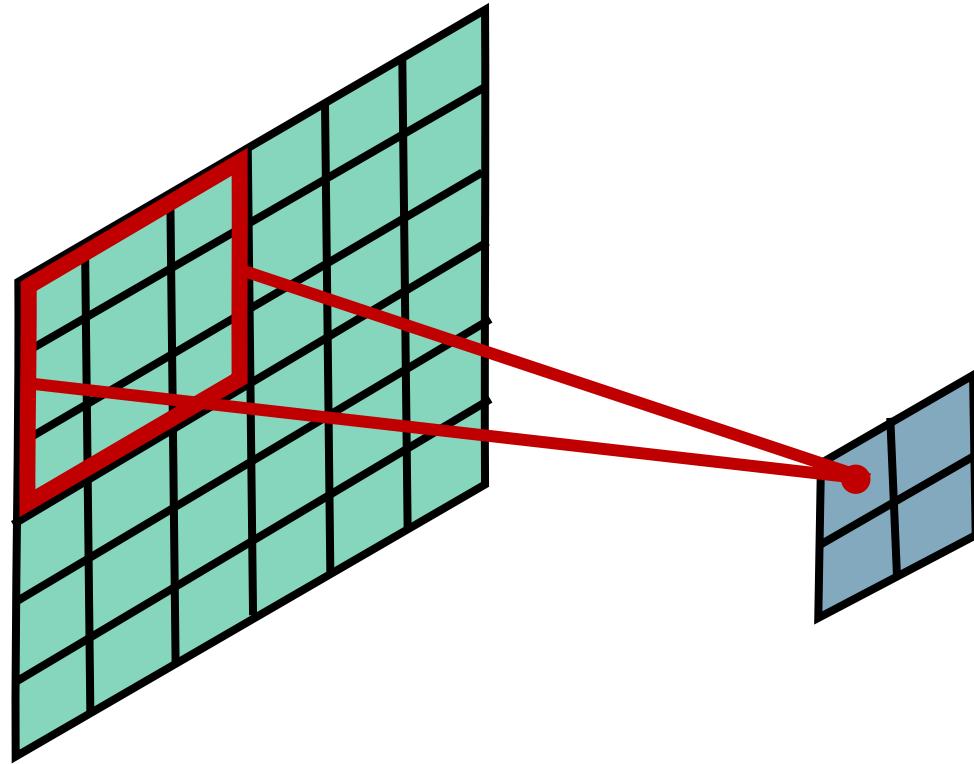
Stride = 1



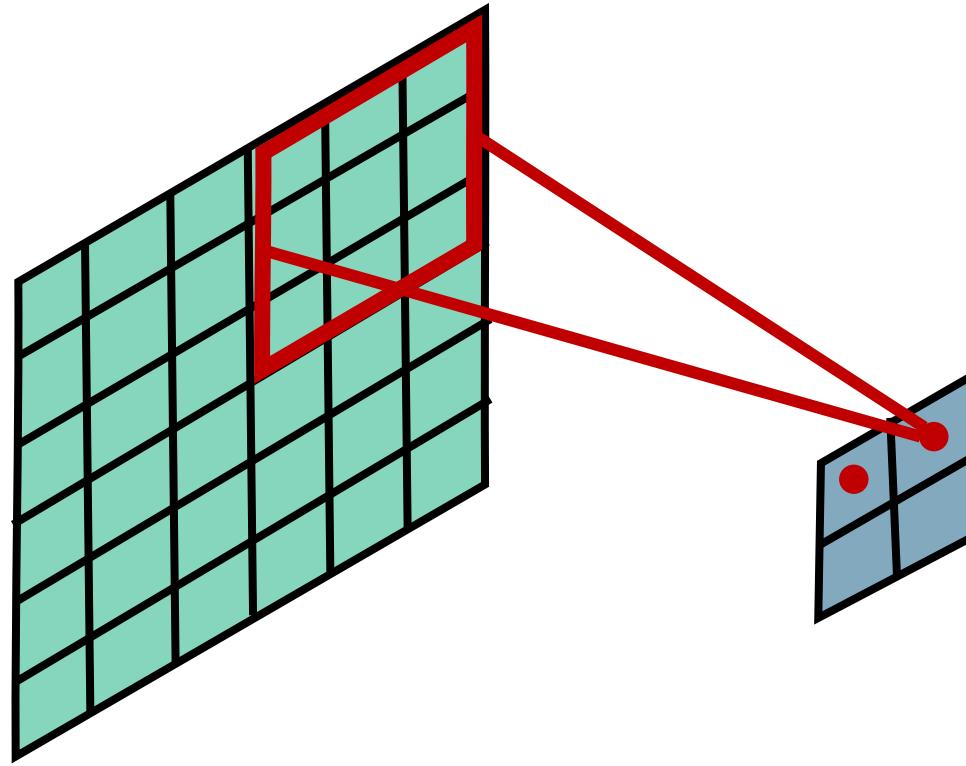
Stride = 1



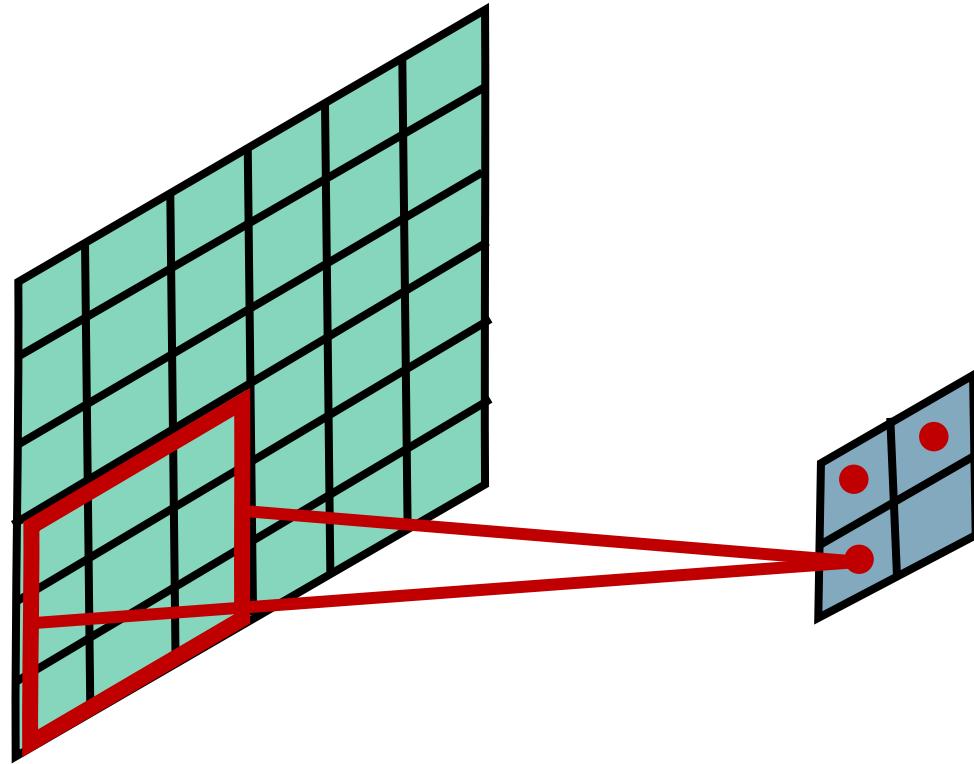
Stride = 3



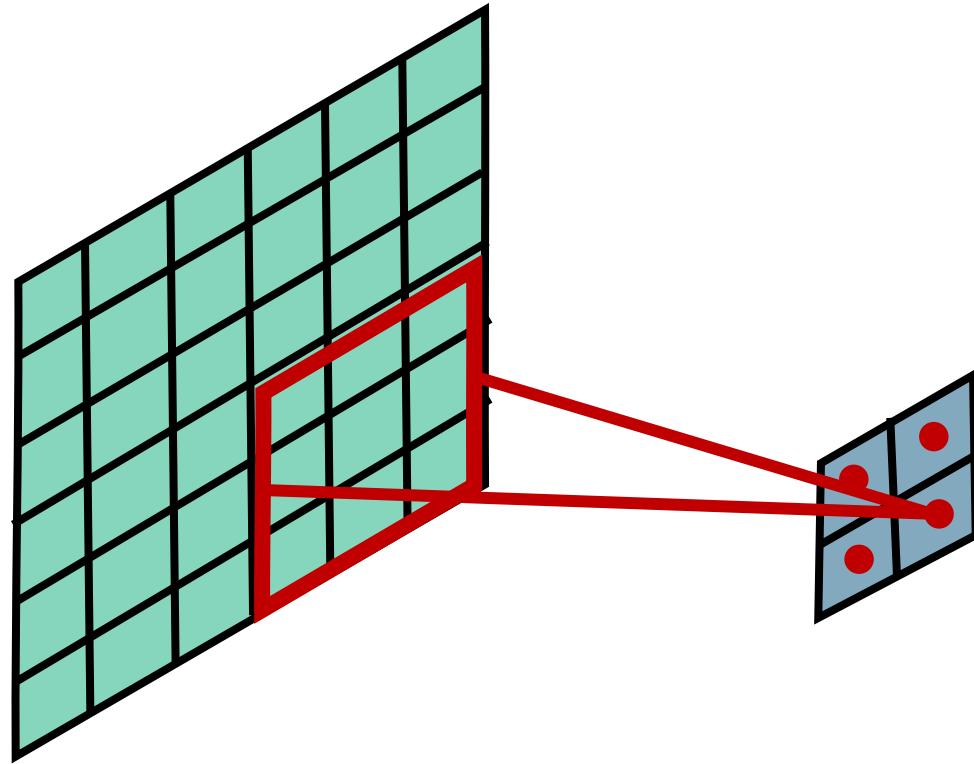
Stride = 3



Stride = 3



Stride = 3



Convolutional Layer

Question: What is the size of the output? What's the computational cost?

Answer: It is proportional to the number of filters and depends on the stride. If kernels have size $K \times K$, input has size $D \times D$, stride is 1, and there are M input feature maps and N output feature maps then:

- the input has size $M @ D \times D$
- the output has size $N @ (D - K + 1) \times (D - K + 1)$
- the kernels have $M \times N \times K \times K$ coefficients (which have to be learned)
- cost: $M * K * K * N * (D - K + 1) * (D - K + 1)$

Convolutional Layer

Question: What is the size of the output? What's the computational cost?

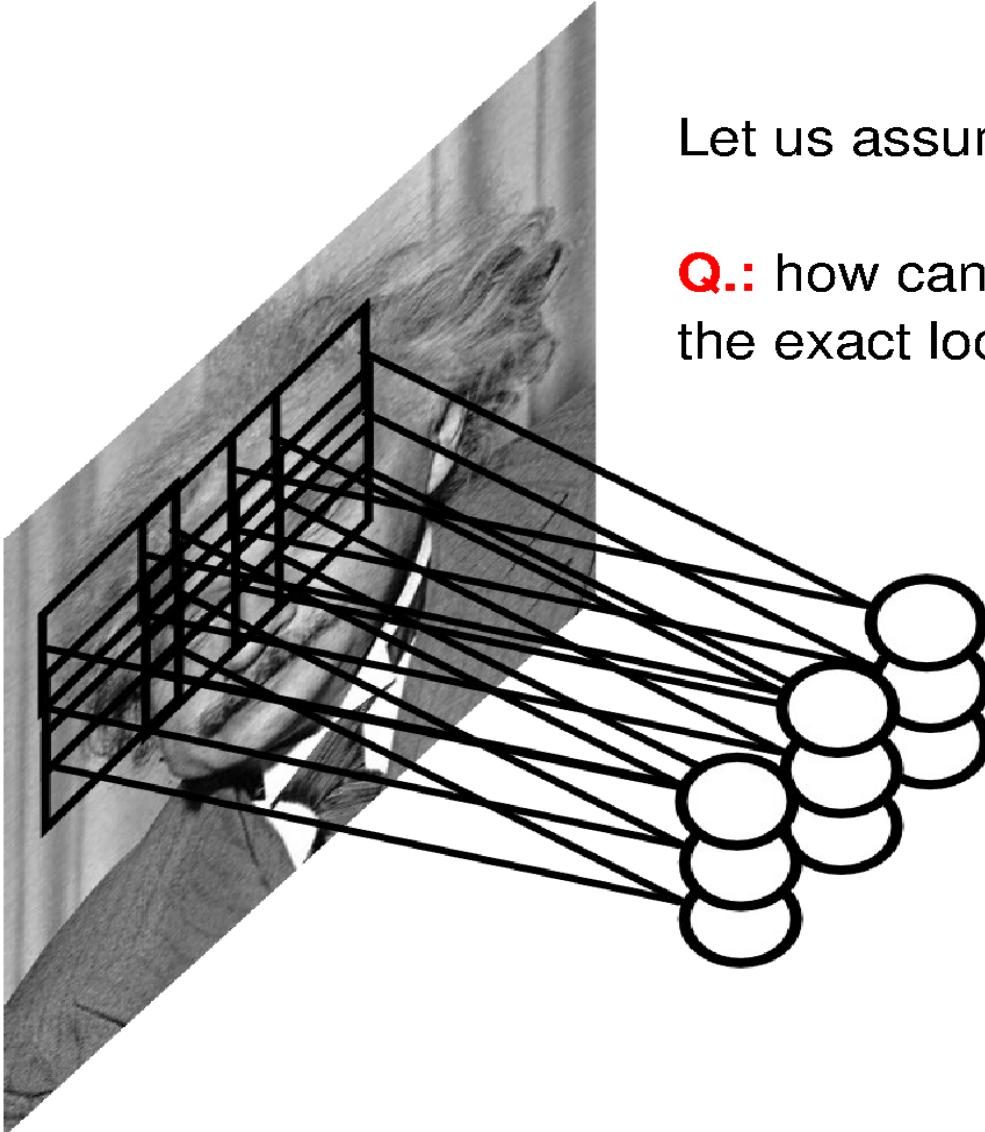
Answer: It is proportional to the number of filters and depends on the stride. If kernels have size $K \times K$, input has size $D \times D$, stride is 1, and there are M input feature maps and N output feature maps then:

- the input has size $M @ D \times D$
- the output has size $N @ (D - K + 1) \times (D - K + 1)$
- the kernels have $M \times N \times K \times K$ coefficients (which have to be learned)
- cost: $M \cdot K^2 \cdot N \cdot (D - K + 1)^2$

Question: How many feature maps? What's the size of the filters?

Answer: Usually, there are more output feature maps than input feature maps. Convolutional layers can increase the number of hidden units by big factors (and are expensive to compute). The size of the filters has to match the size/scale of the patterns we want to detect (task dependent).

Pooling Layer



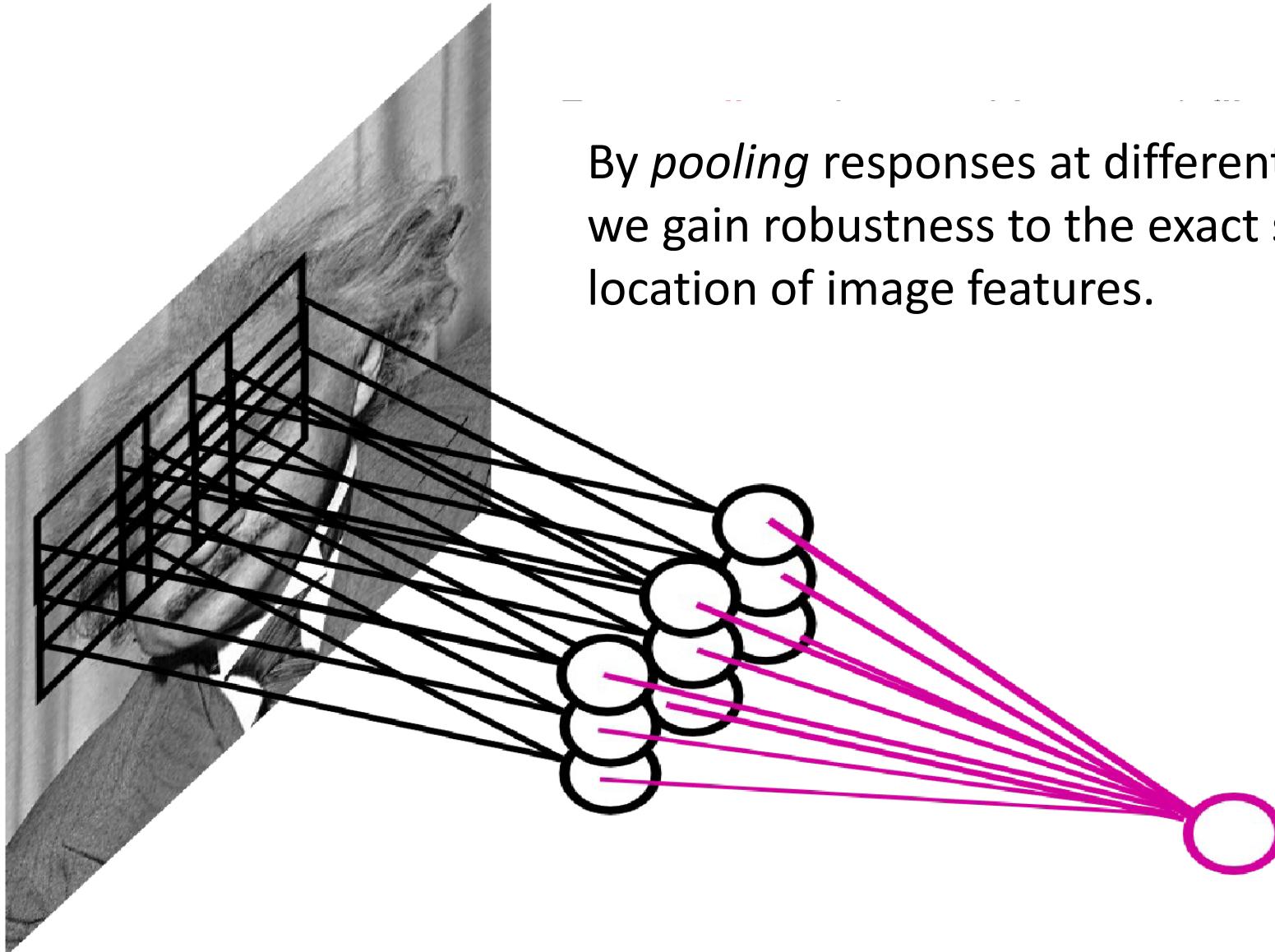
Let us assume filter is an “eye” detector.

Q.: how can we make the detection robust to the exact location of the eye?

60

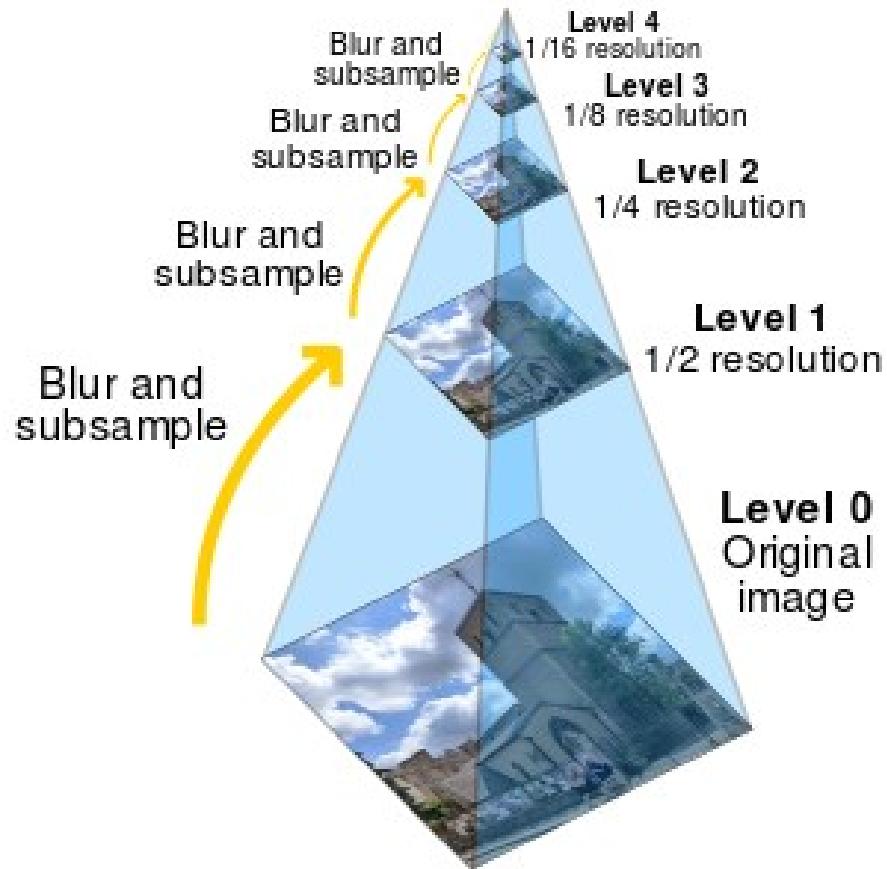
Ranzato 

Pooling Layer



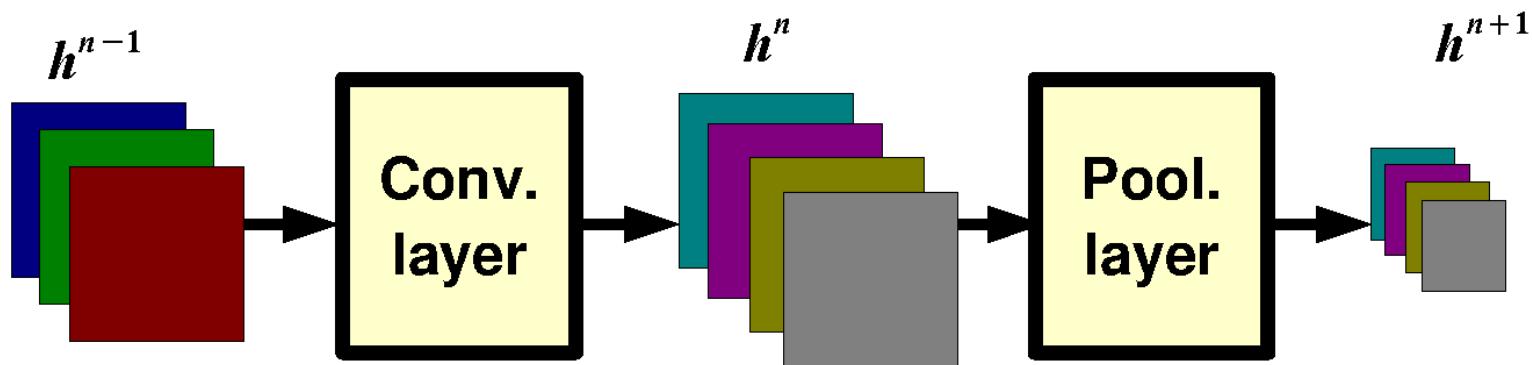
By *pooling* responses at different locations,
we gain robustness to the exact spatial
location of image features.

Pooling is similar to downsampling



...except sometimes we don't want to blur,
as other functions might be better for classification.

Pooling Layer: Receptive Field Size



Pooling Layer: Examples

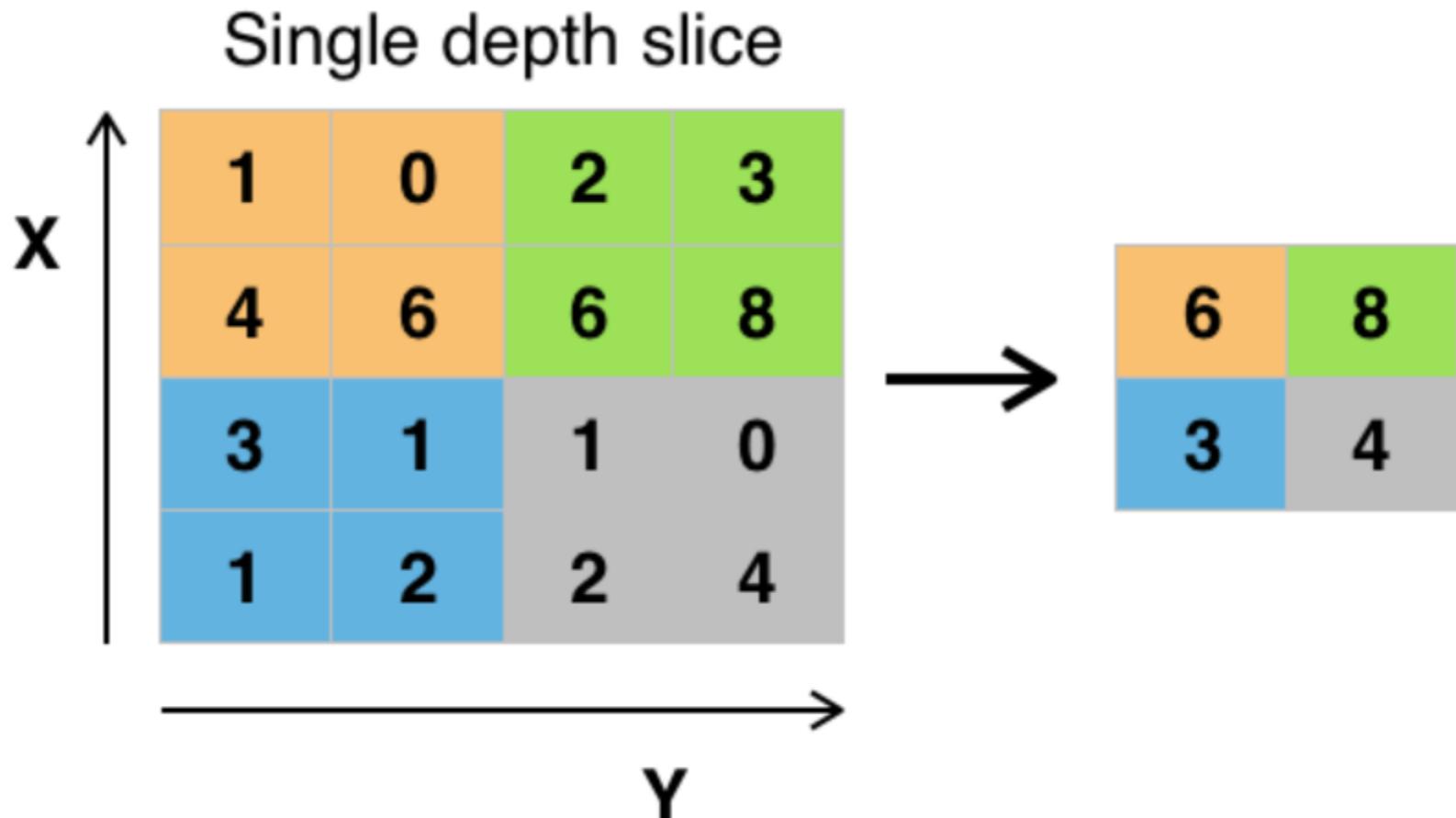
Max-pooling:

$$h_j^n(x, y) = \max_{\bar{x} \in N(x), \bar{y} \in N(y)} h_j^{n-1}(\bar{x}, \bar{y})$$

Average-pooling:

$$h_j^n(x, y) = 1/K \sum_{\bar{x} \in N(x), \bar{y} \in N(y)} h_j^{n-1}(\bar{x}, \bar{y})$$

Max pooling



Pooling Layer: Examples

Max-pooling:

$$h_j^n(x, y) = \max_{\bar{x} \in N(x), \bar{y} \in N(y)} h_j^{n-1}(\bar{x}, \bar{y})$$

Average-pooling:

$$h_j^n(x, y) = 1/K \sum_{\bar{x} \in N(x), \bar{y} \in N(y)} h_j^{n-1}(\bar{x}, \bar{y})$$

L2-pooling:

$$h_j^n(x, y) = \sqrt{\sum_{\bar{x} \in N(x), \bar{y} \in N(y)} h_j^{n-1}(\bar{x}, \bar{y})^2}$$

L2-pooling over features:

$$h_j^n(x, y) = \sqrt{\sum_{k \in N(j)} h_k^{n-1}(x, y)^2}$$

62

Pooling Layer

Question: What is the size of the output? What's the computational cost?

Answer: The size of the output depends on the stride between the pools. For instance, if pools do not overlap and have size $K \times K$, and the input has size $D \times D$ with M input feature maps, then:

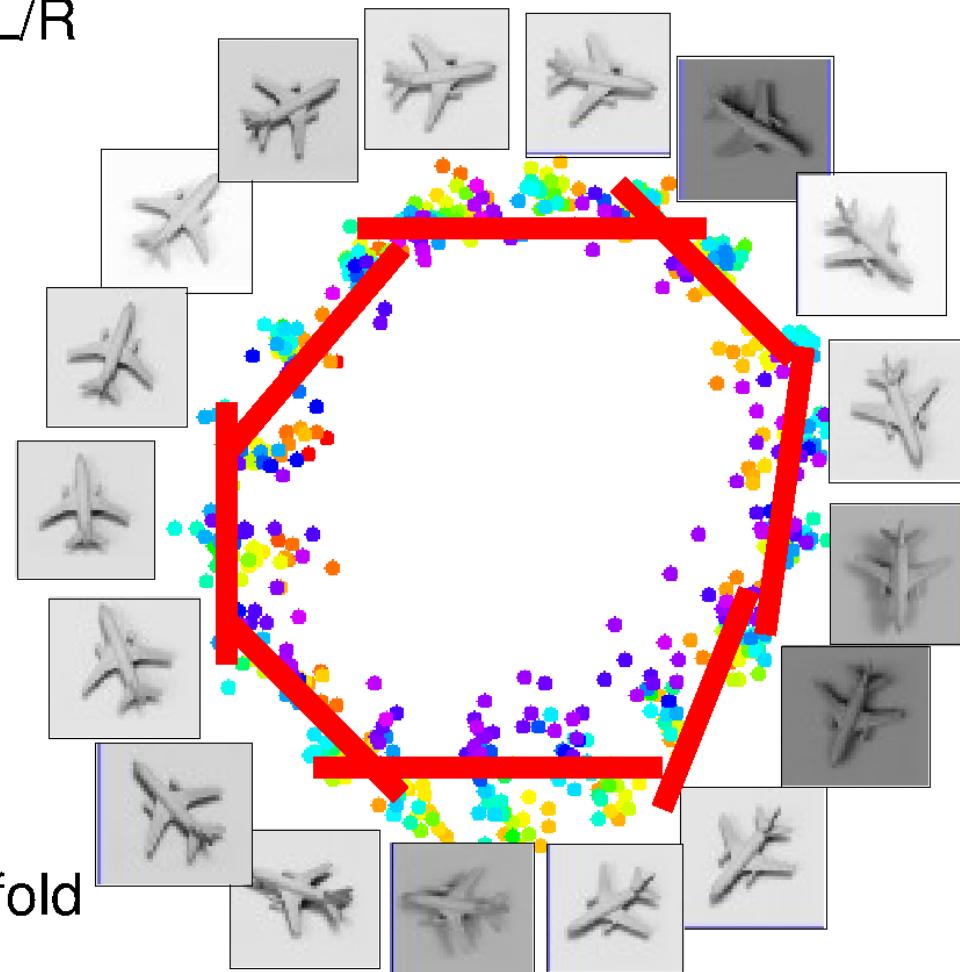
- output is $M @ (D/K) \times (D/K)$
- the computational cost is proportional to the size of the input (negligible compared to a convolutional layer)

Question: How should I set the size of the pools?

Answer: It depends on how much “invariant” or robust to distortions we want the representation to be. It is best to pool slowly (via a few stacks of conv-pooling layers).

Pooling Layer: Interpretation

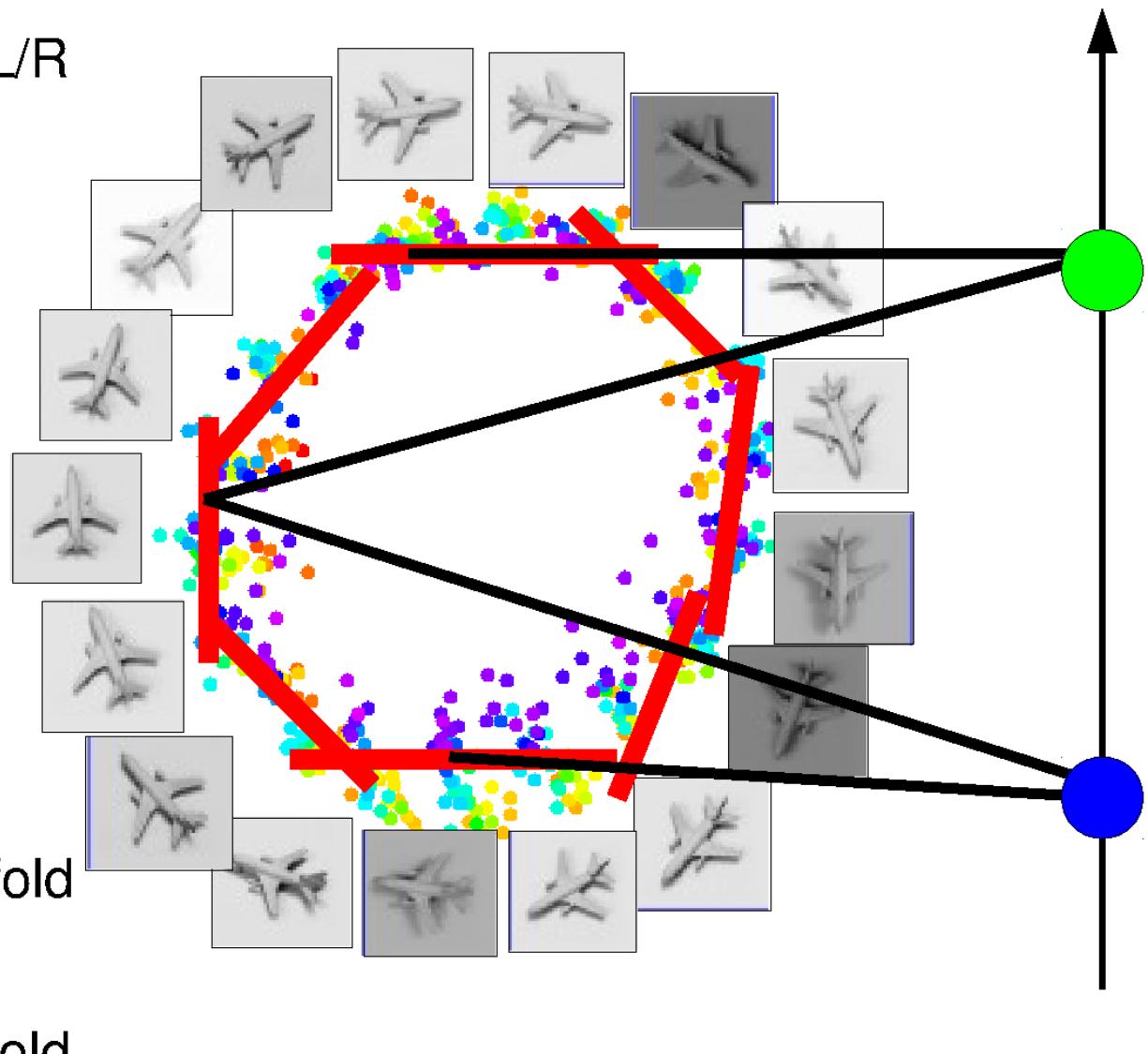
Task: detect orientation L/R



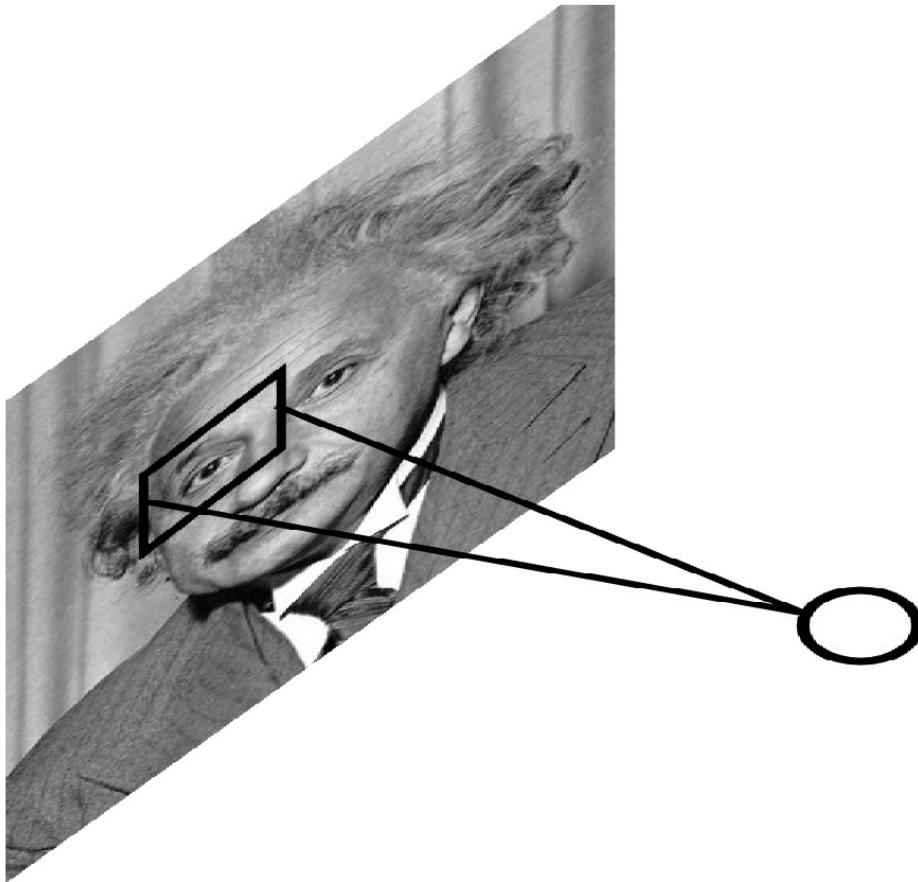
Conv layer:
linearizes manifold

Pooling Layer: Interpretation

Task: detect orientation L/R



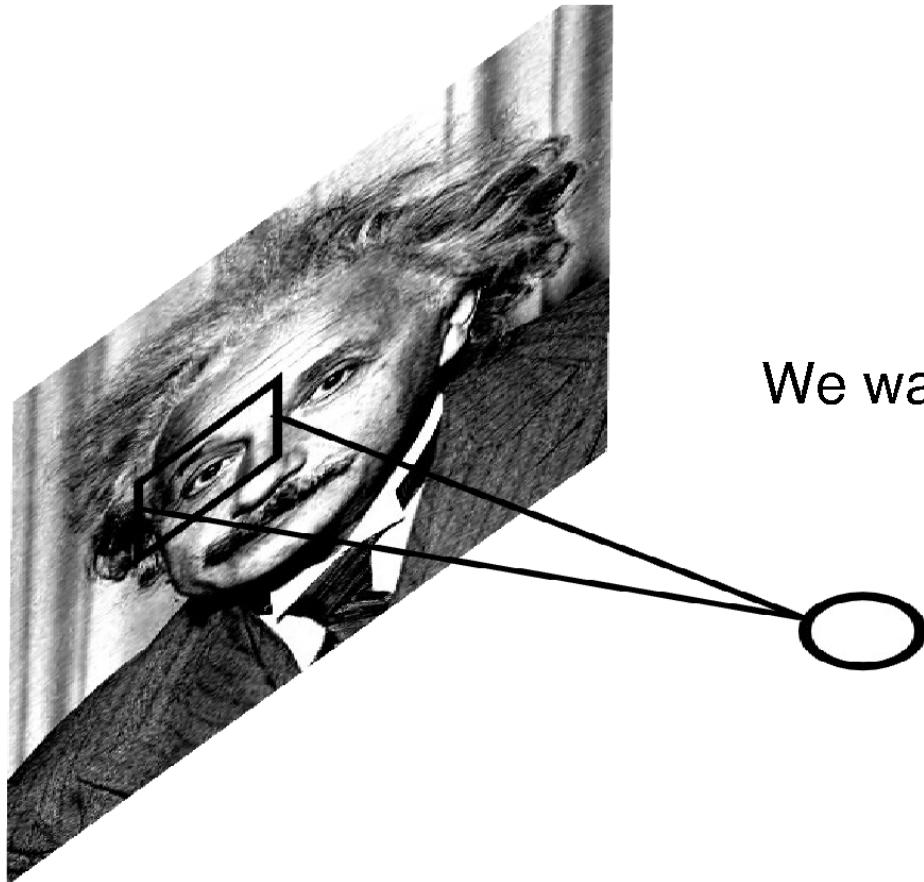
Local Contrast Normalization



68

Ranzato 

Local Contrast Normalization



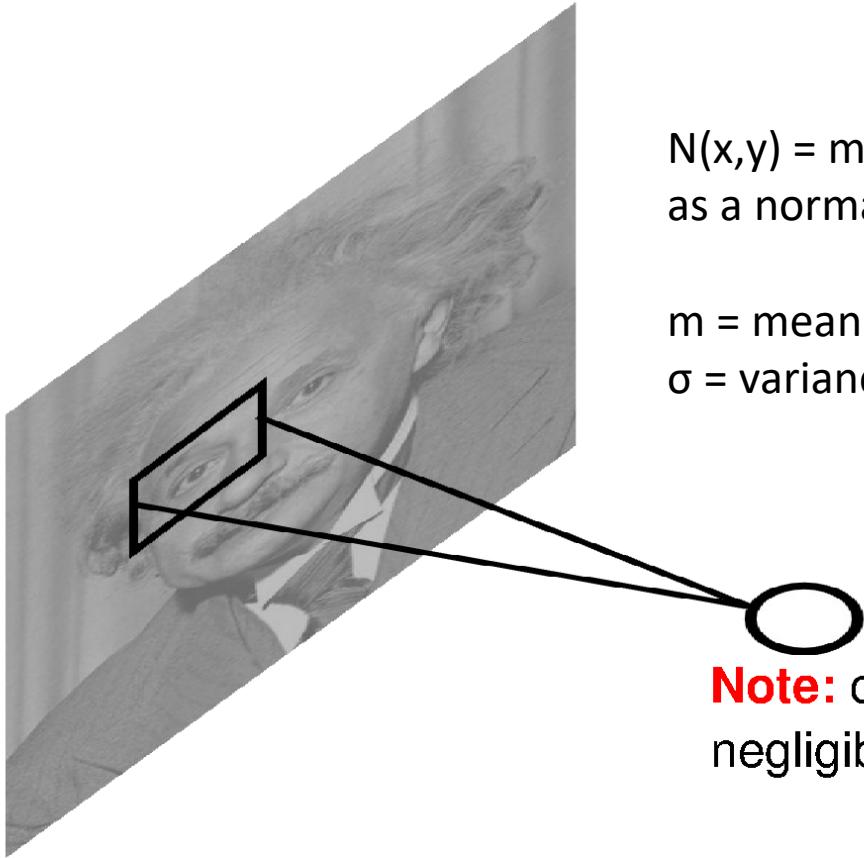
We want the same response.

69

Ranzato 

Local Contrast Normalization

$$h^{i+1}(x, y) = \frac{h^i(x, y) - m^i(N(x, y))}{\sigma^i(N(x, y))}$$



$N(x,y)$ = model pixel values in window
as a normal distribution

m = mean

σ = variance

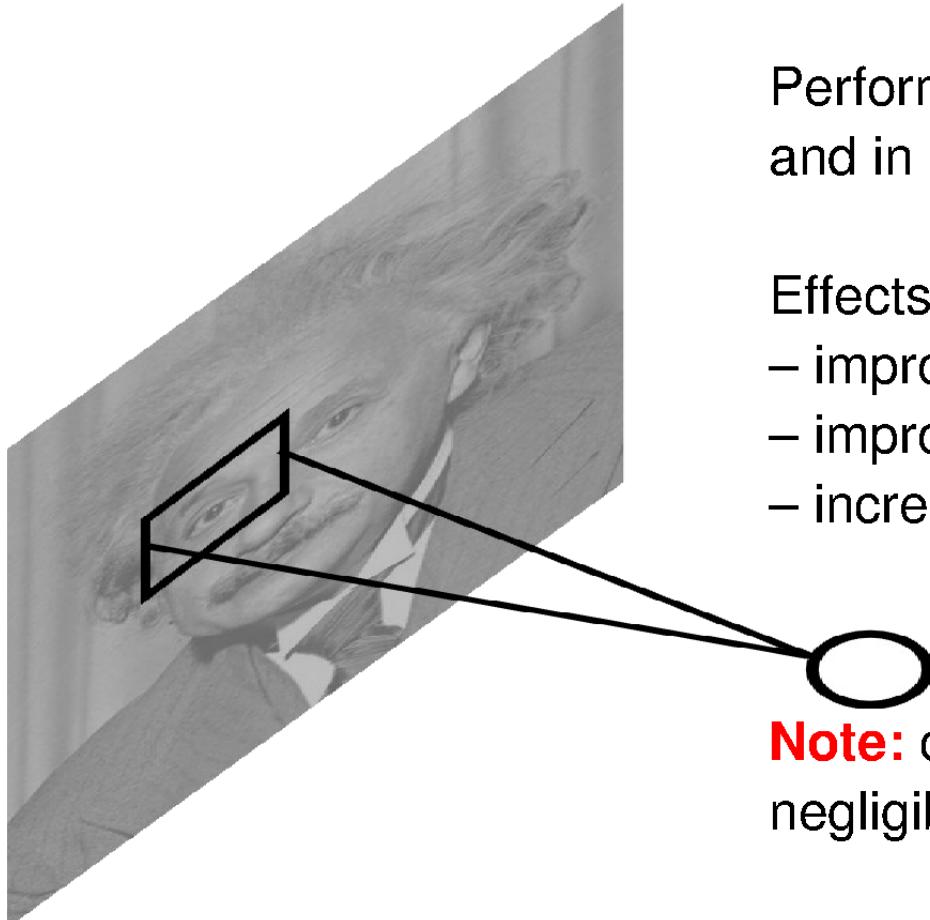
Note: computational cost is
negligible w.r.t. conv. layer.

70

Ranzato

Local Contrast Normalization

$$h^{i+1}(x, y) = \frac{h^i(x, y) - m^i(N(x, y))}{\sigma^i(N(x, y))}$$



Performed also across features
and in the higher layers..

Effects:

- improves invariance
- improves optimization
- increases sparsity

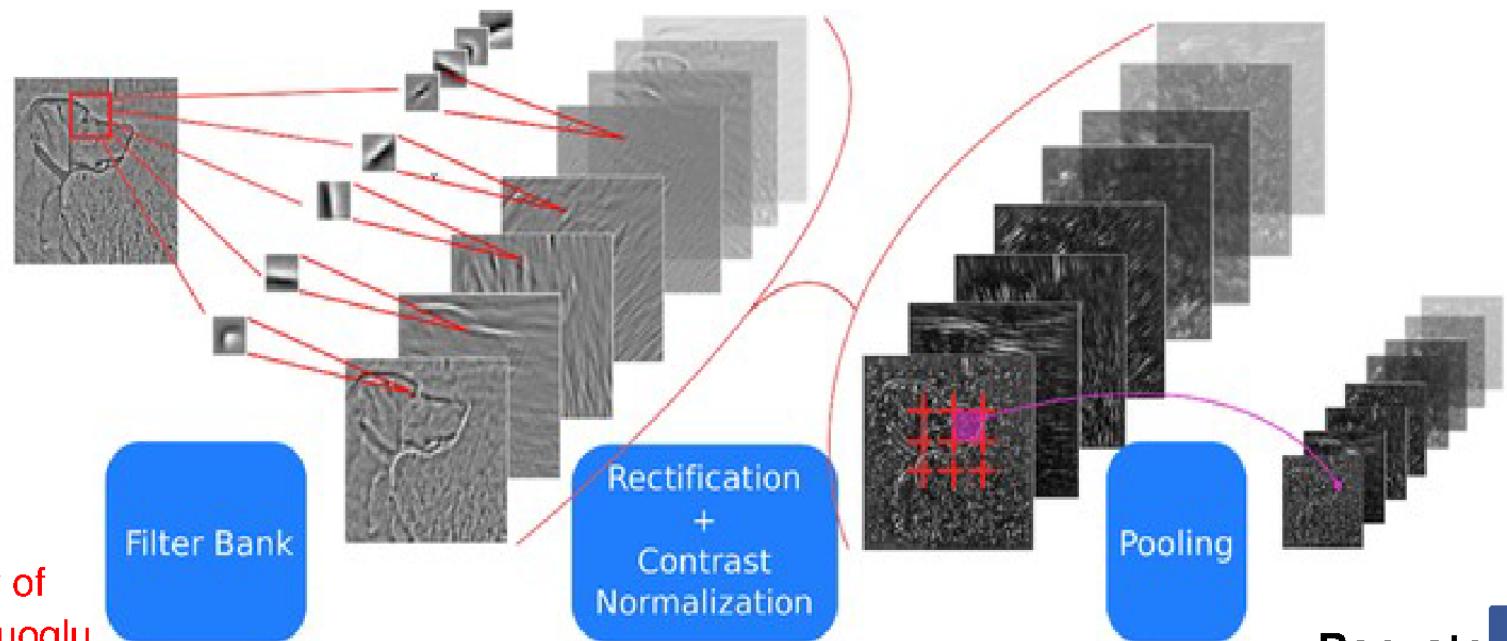
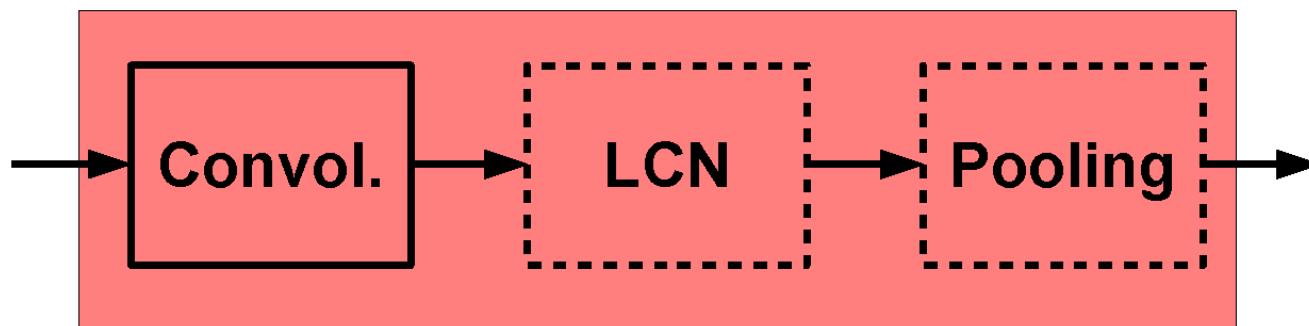
Note: computational cost is
negligible w.r.t. conv. layer.

70

Ranzato 

ConvNets: Typical Stage

One stage (zoom)

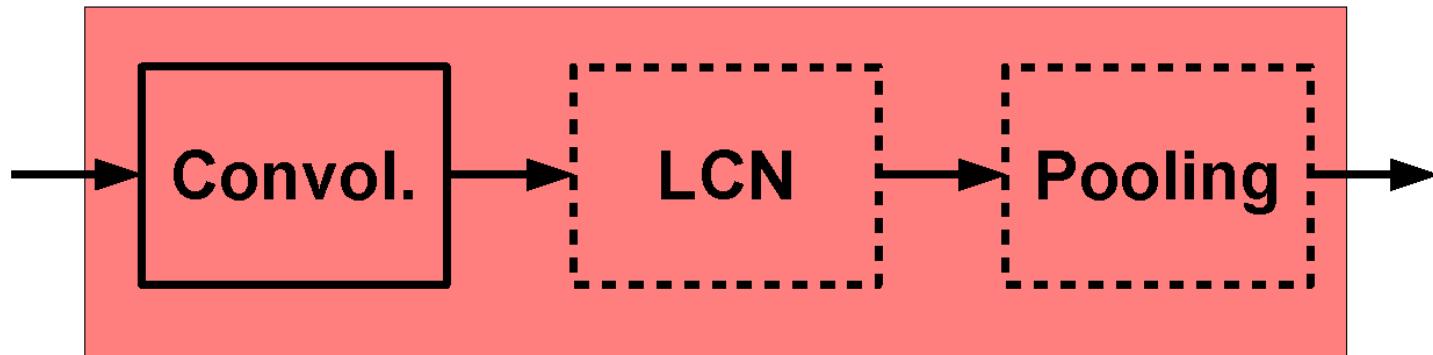


courtesy of
K. Kavukcuoglu

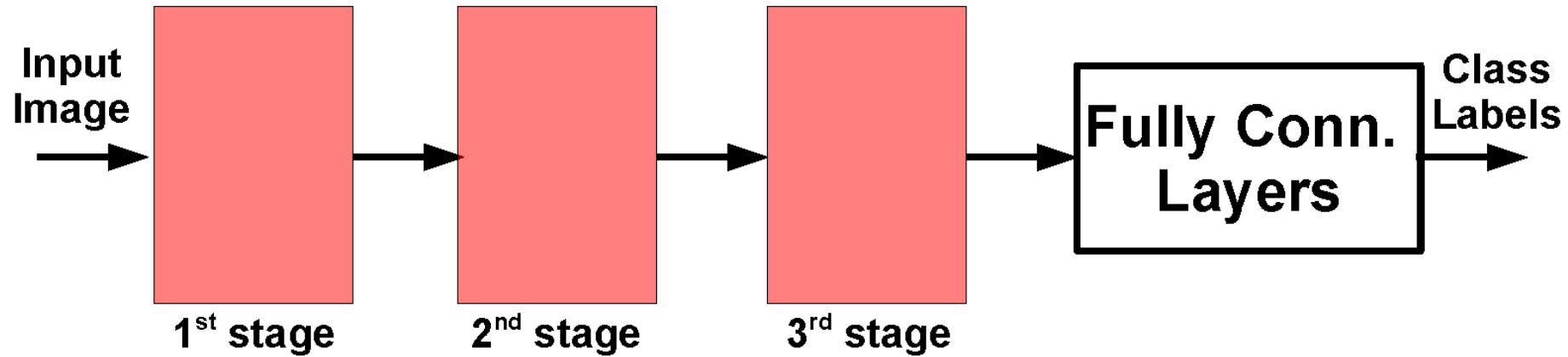
Ranzato

ConvNets: Typical Architecture

One stage (zoom)

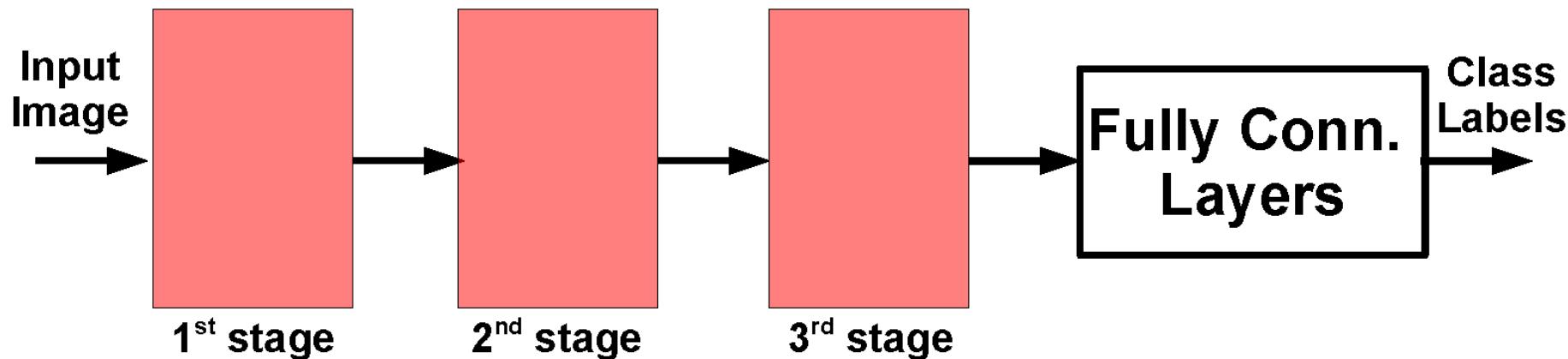


Whole system



ConvNets: Typical Architecture

Whole system



Conceptually similar to:

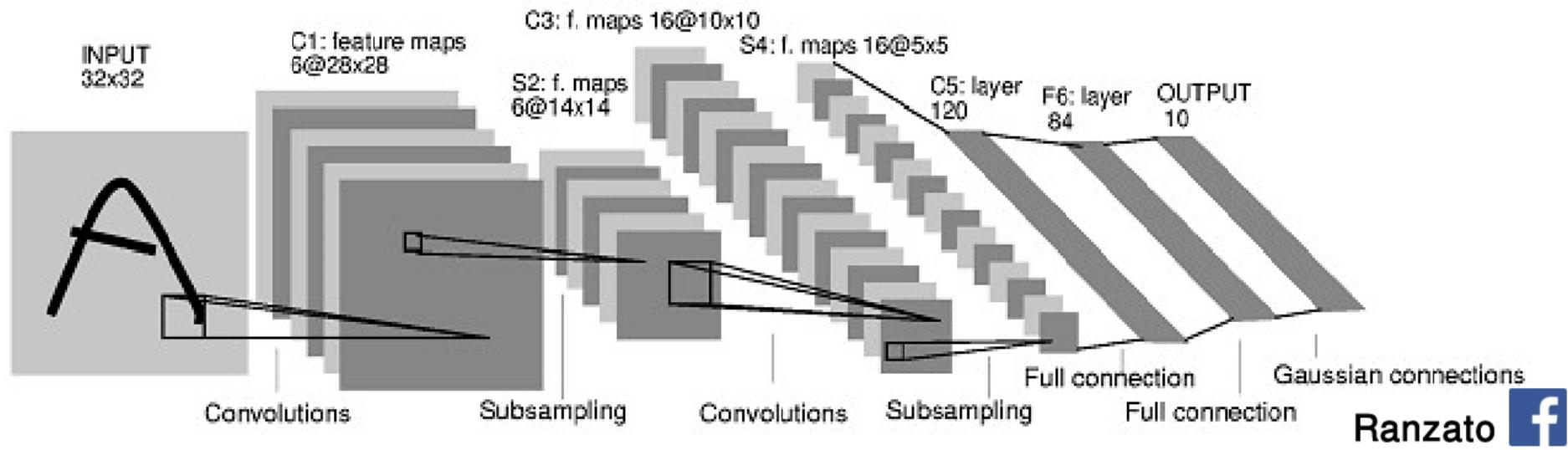
SIFT → K-Means → Pyramid Pooling → SVM

Lazebnik et al. "...Spatial Pyramid Matching..." CVPR 2006

SIFT → Fisher Vect. → Pooling → SVM

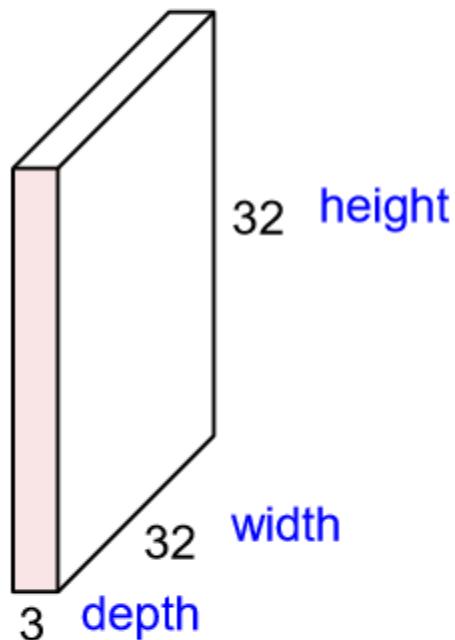
Sanchez et al. "Image classification with F.V.: Theory and practice" IJCV 2012

Yann LeCun's MNIST CNN architecture



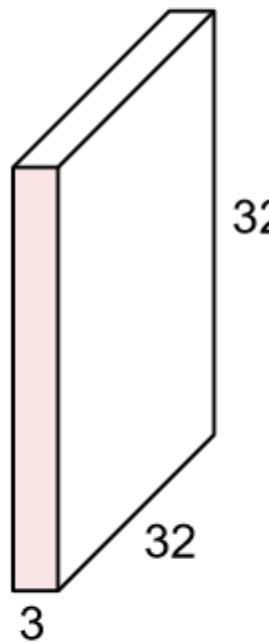
Convolutions: More detail

32x32x3 image



Convolutions: More detail

32x32x3 image

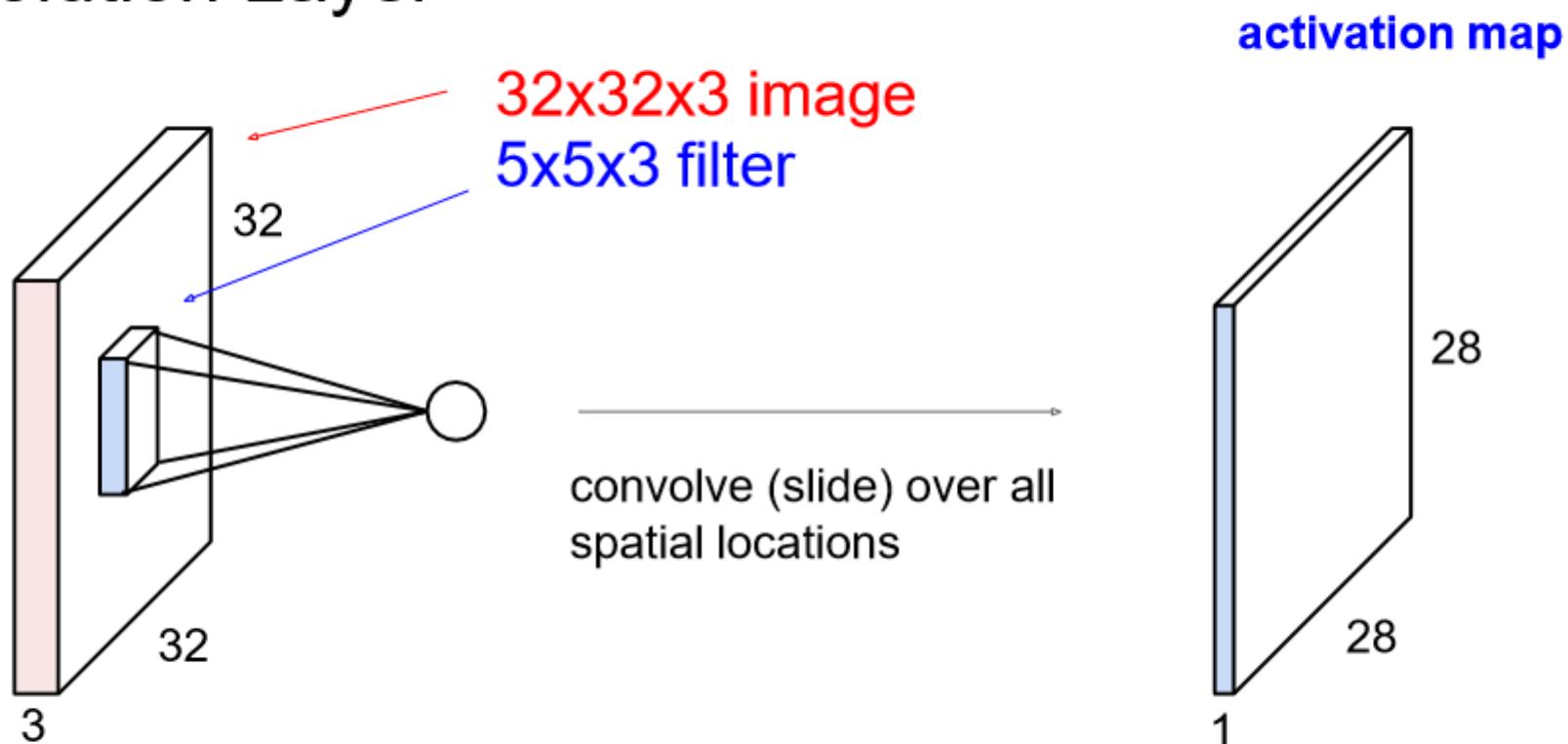


5x5x3 filter



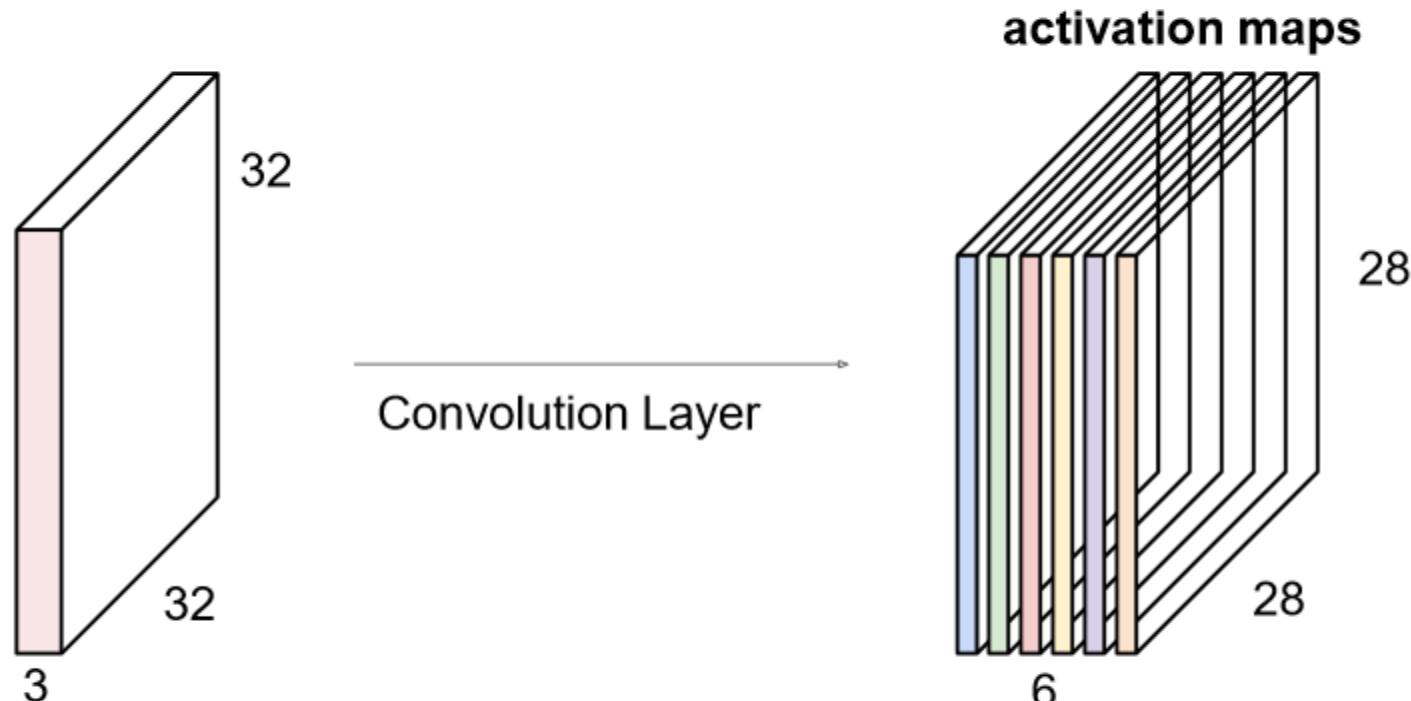
Convolutions: More detail

Convolution Layer



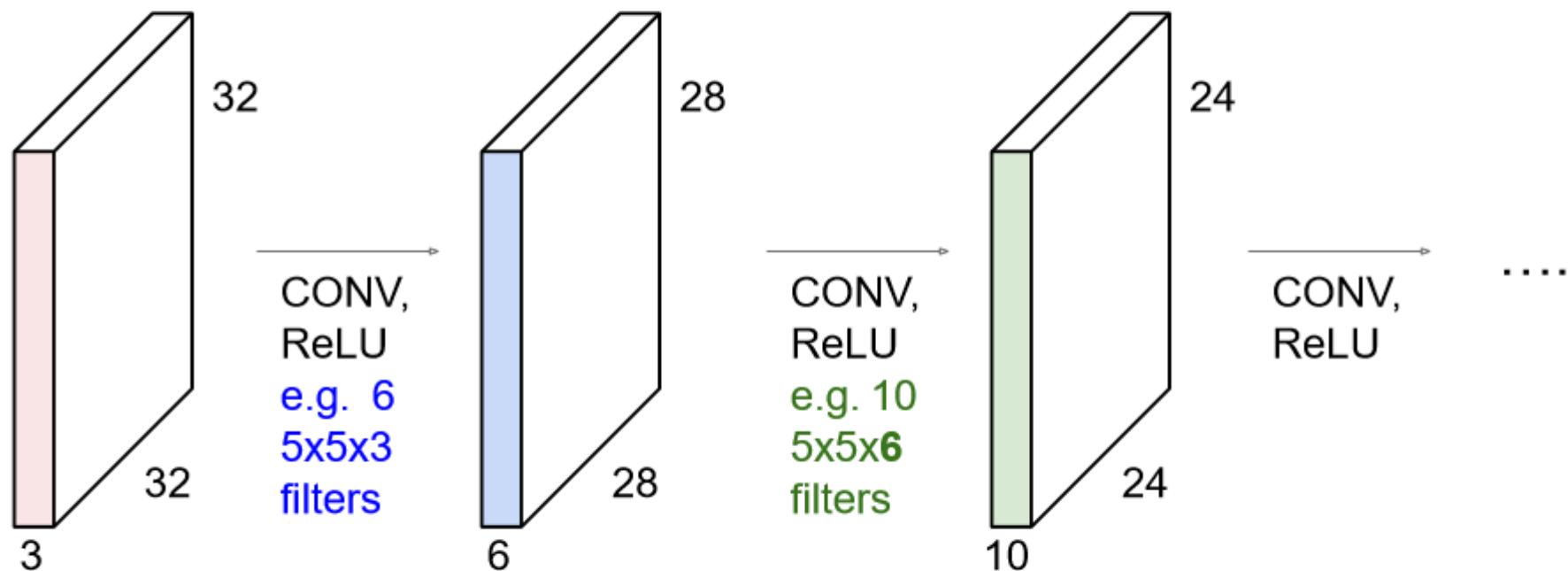
Convolutions: More detail

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

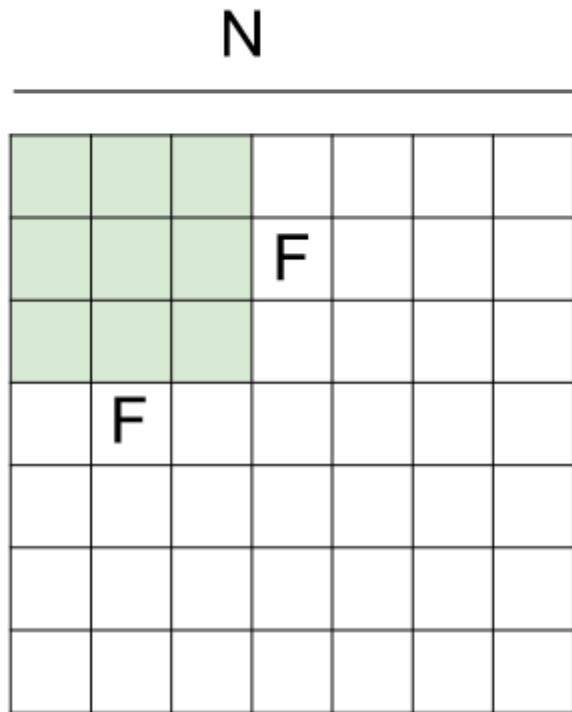


We stack these up to get a “new image” of size $28 \times 28 \times 6$!

Convolutions: More detail



Convolutions: More detail



Output size:
 $(N - F) / \text{stride} + 1$

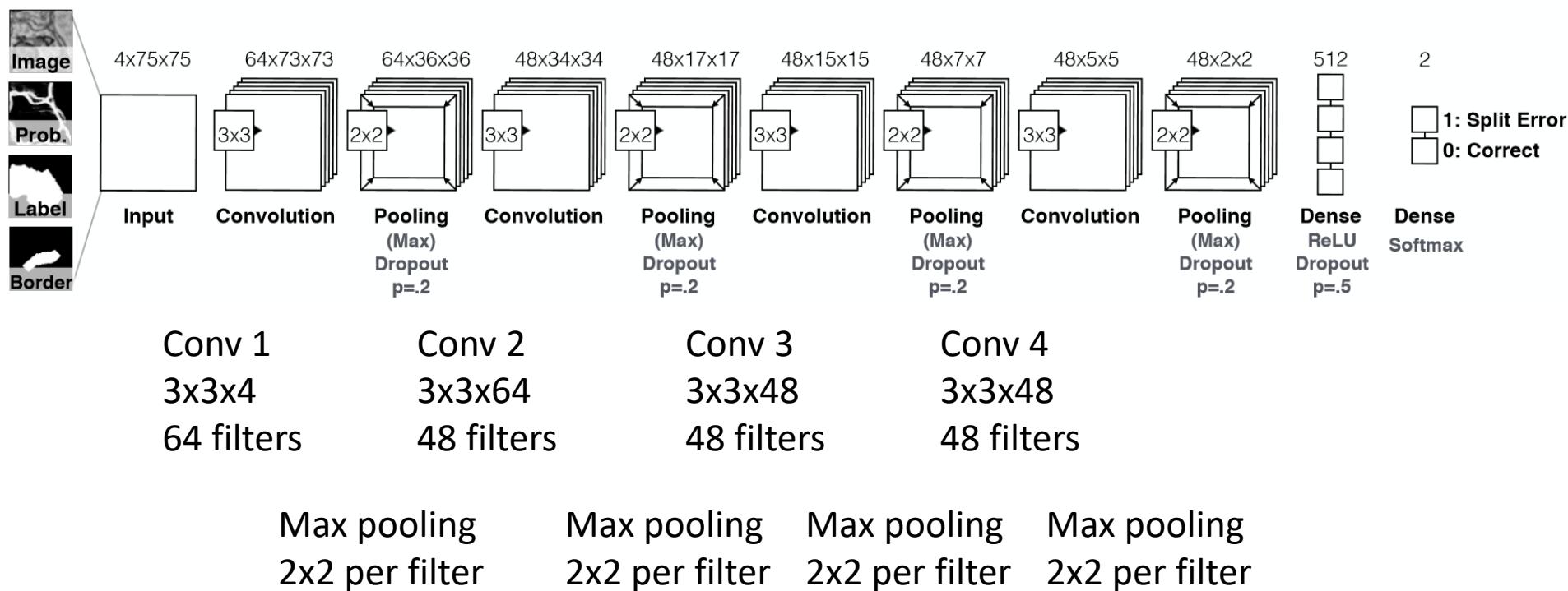
N

Our connectomics diagram

Auto-generated from network declaration by *nolearn* (for Lasagne / Theano)

Input

75x75x4



Reading architecture diagrams

Layers

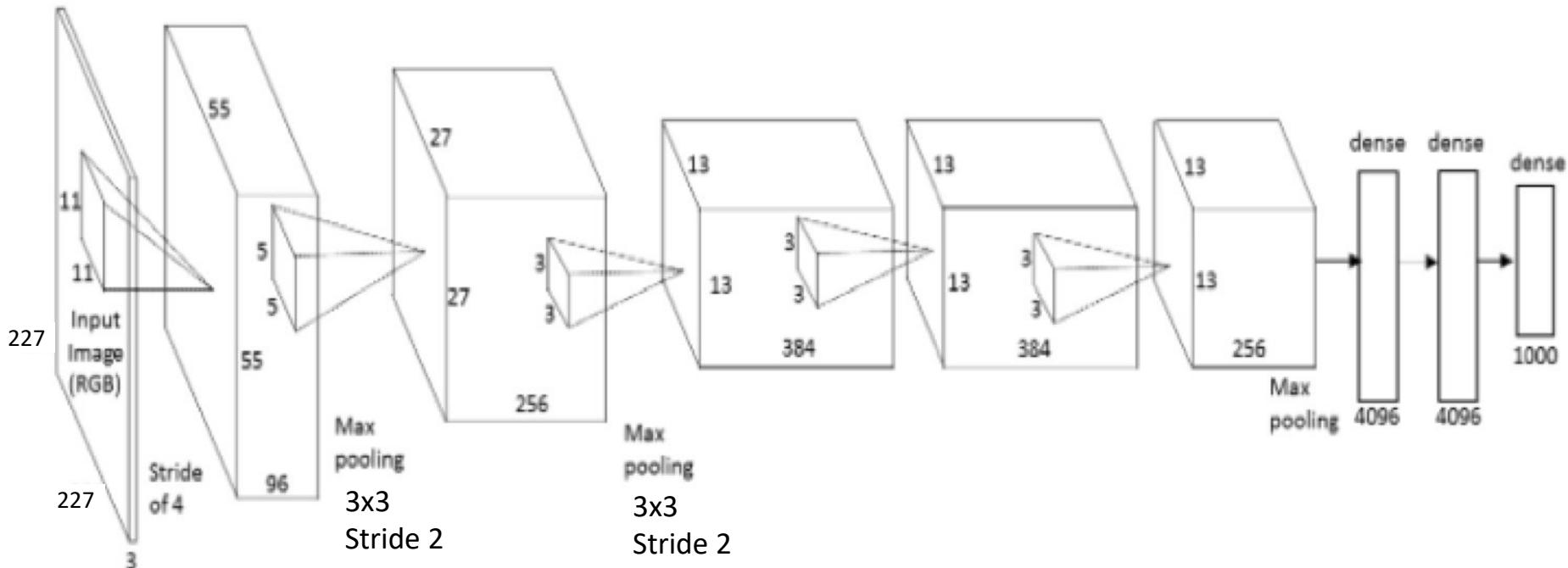
- Kernel sizes
- Strides
- # channels
- # kernels
- Max pooling

params	AlexNet	FLOPs
4M	FC 1000	4M
16M	FC 4096 / ReLU	16M
37M	FC 4096 / ReLU	37M
	Max Pool 3x3s2	
442K	Conv 3x3s1, 256 / ReLU	74M
1.3M	Conv 3x3s1, 384 / ReLU	112M
884K	Conv 3x3s1, 384 / ReLU	149M
	Max Pool 3x3s2	
	Local Response Norm	
307K	Conv 5x5s1, 256 / ReLU	223M
	Max Pool 3x3s2	
	Local Response Norm	
35K	Conv 11x11s4, 96 / ReLU	105M

AlexNet diagram (simplified)

Input size

227 x 227 x 3



Conv 1

$11 \times 11 \times 3$

Stride 4

96 filters

Conv 2

$5 \times 5 \times 48$

Stride 1

256 filters

Conv 3

$3 \times 3 \times 256$

Stride 1

384 filters

Conv 4

$3 \times 3 \times 192$

Stride 1

384 filters

Conv 4

$3 \times 3 \times 192$

Stride 1

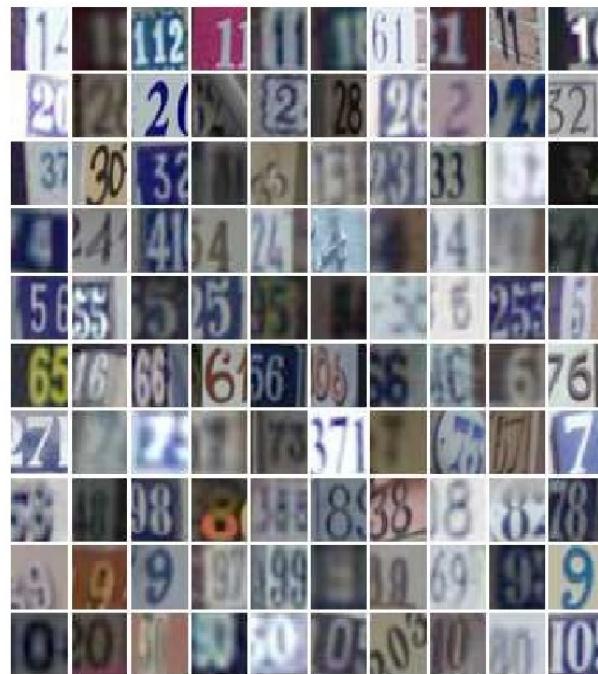
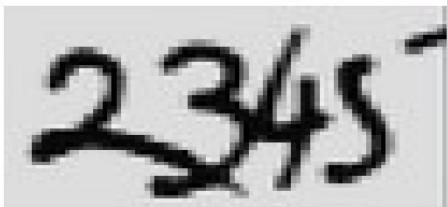
256 filters

Outline

- Supervised Neural Networks
- Convolutional Neural Networks
- Examples
- Tips

CONV NETS: EXAMPLES

- OCR / House number & Traffic sign classification



Ciresan et al. "MCDNN for image classification" CVPR 2012

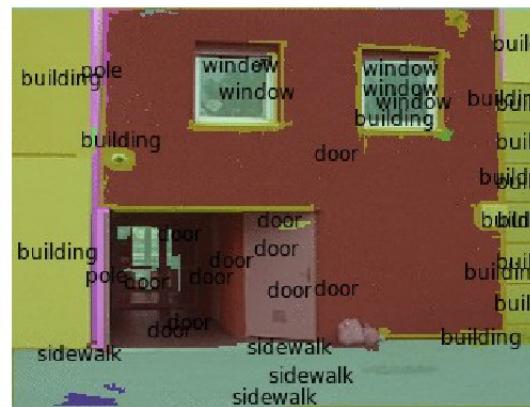
Wan et al. "Regularization of neural networks using dropconnect" ICML 2013

82

Jaderberg et al. "Synthetic data and ANN for natural scene text recognition" arXiv 2014

CONV NETS: EXAMPLES

- Scene Parsing



Farabet et al. "Learning hierarchical features for scene labeling" PAMI 2013

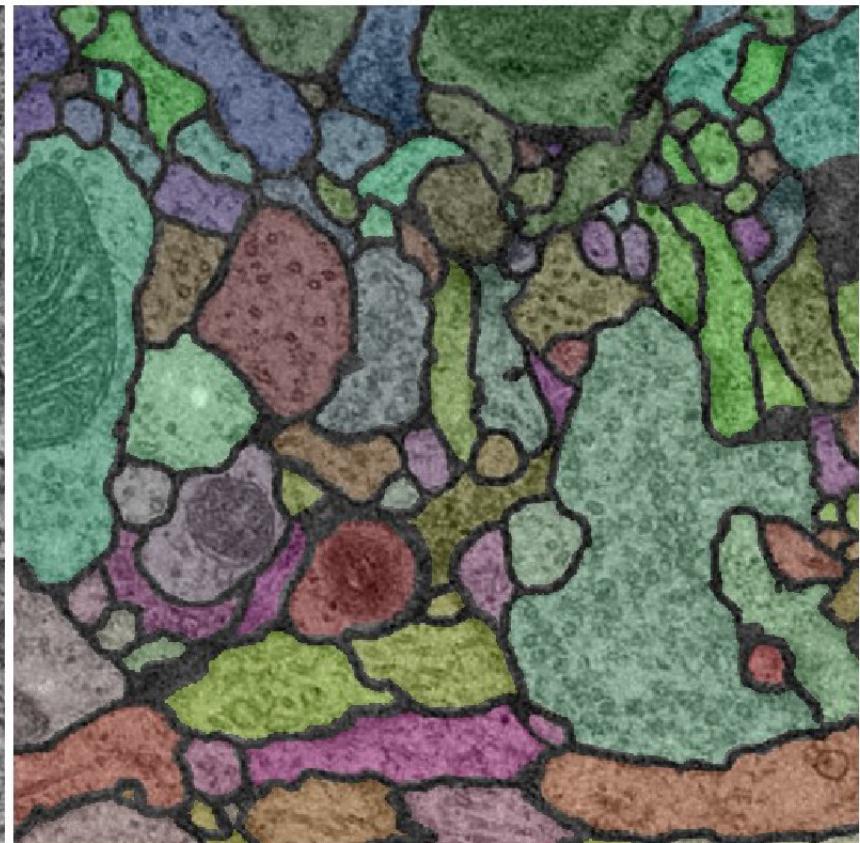
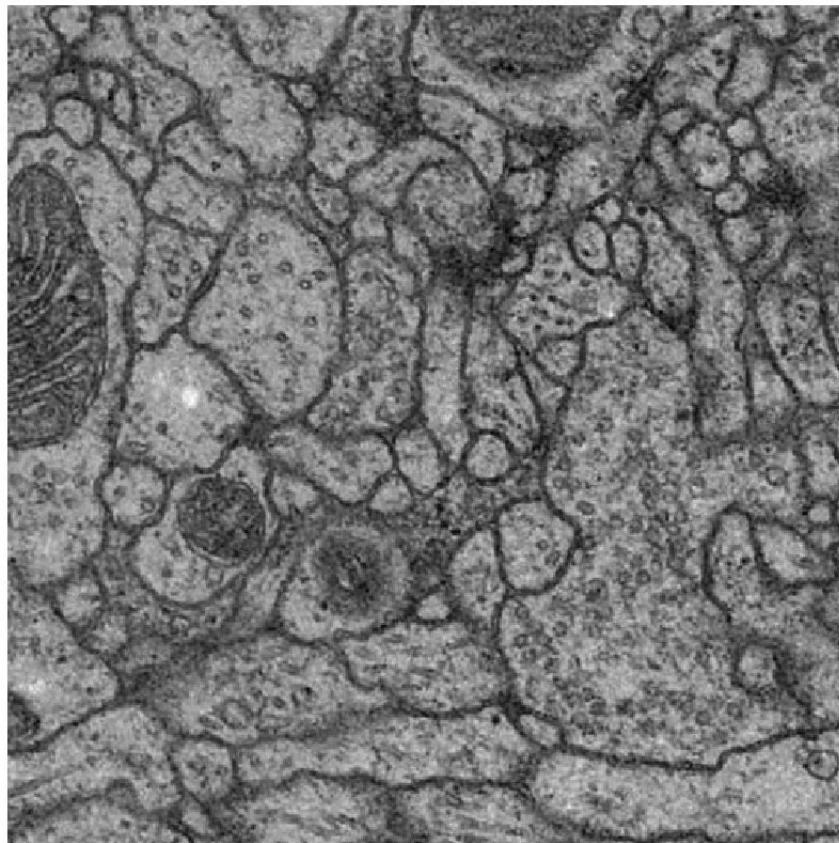
Pinheiro et al. "Recurrent CNN for scene parsing" arxiv 2013

85

Ranzato

CONV NETS: EXAMPLES

- Segmentation 3D volumetric images



Ciresan et al. "DNN segment neuronal membranes..." NIPS 2012

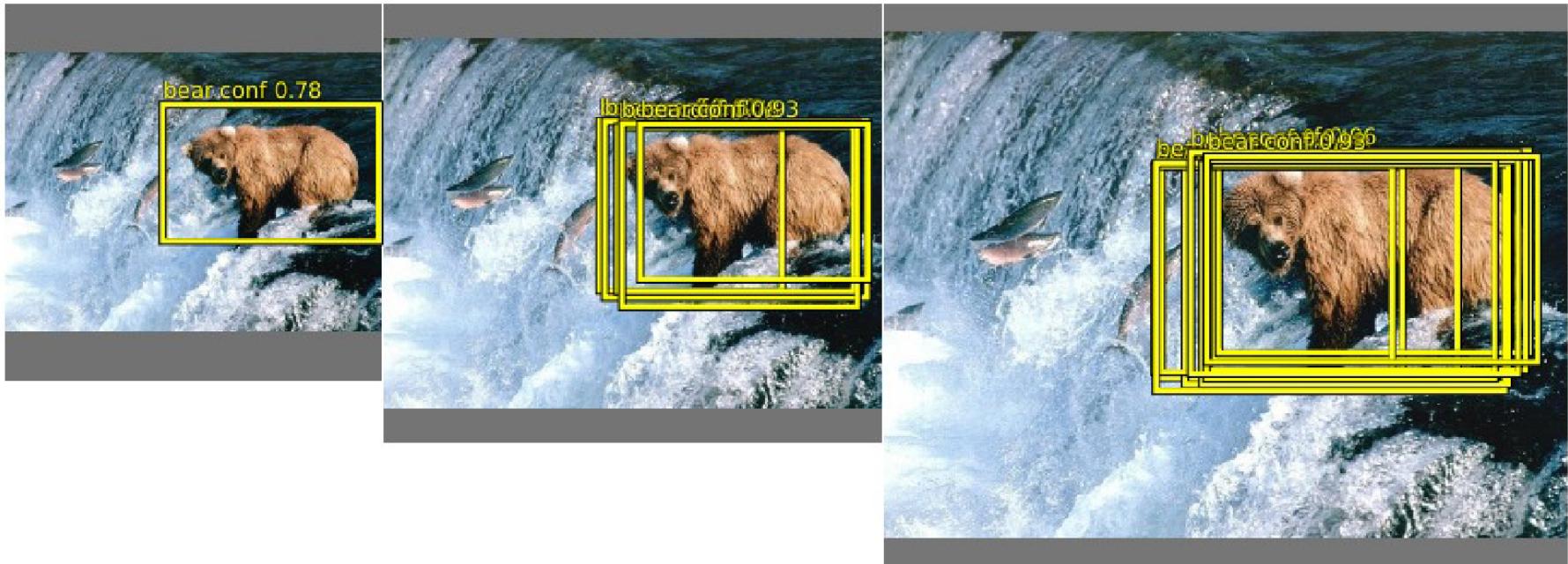
Turaga et al. "Maximin learning of image segmentation" NIPS 2009

86

Ranzato 

CONV NETS: EXAMPLES

- Object detection



Sermanet et al. “OverFeat: Integrated recognition, localization, ...” arxiv 2013

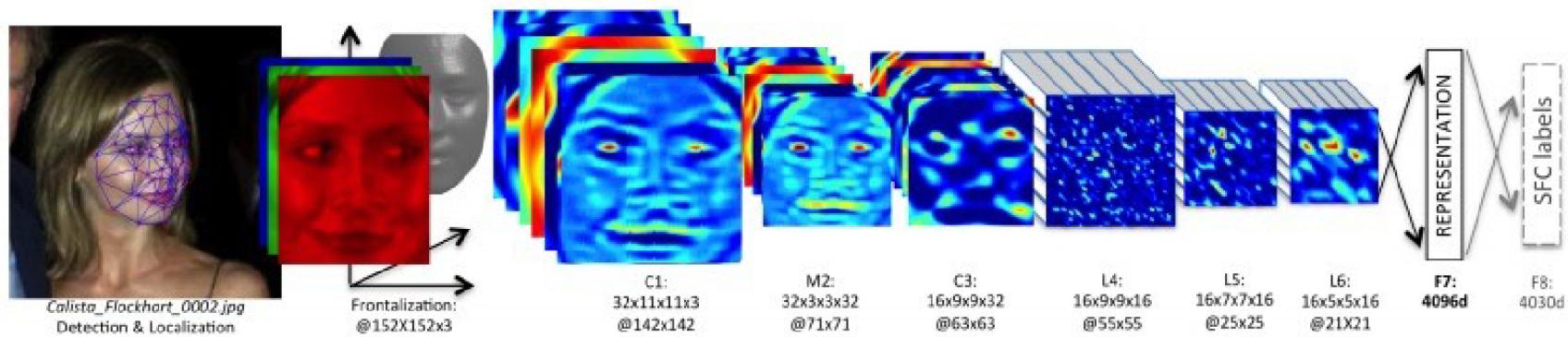
Girshick et al. “Rich feature hierarchies for accurate object detection...” arxiv 2013 91

Szegedy et al. “DNN for object detection” NIPS 2013

Ranzato

CONV NETS: EXAMPLES

- Face Verification & Identification

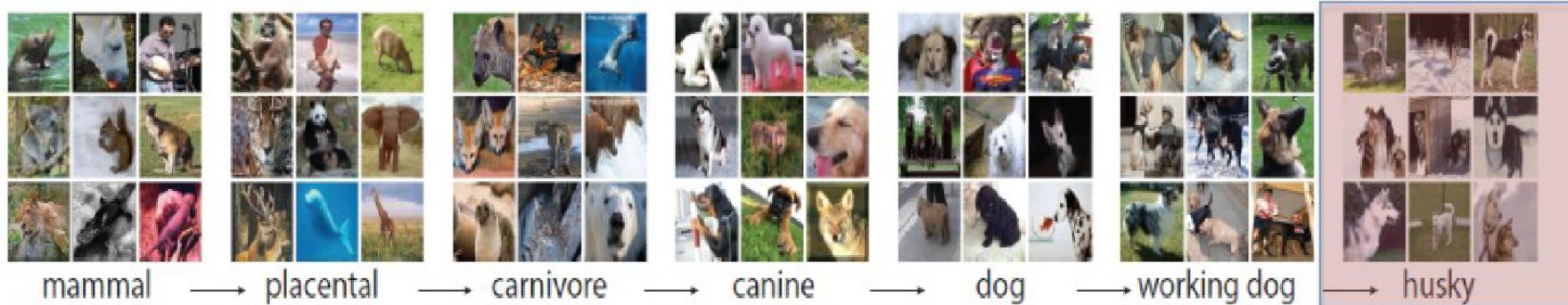


92

Taigman et al. "DeepFace..." CVPR 2014

Ranzato

Dataset: ImageNet 2012



- [S: \(n\) Eskimo dog](#), husky (breed of heavy-coated Arctic sled dog)
 - *direct hypernym / inherited hypernym / sister term*
- [S: \(n\) working dog](#) (any of several breeds of usually large powerful dogs bred to work as draft animals and guard and guide dogs)
 - [S: \(n\) dog, domestic dog, *Canis familiaris*](#) (a member of the genus *Canis* (probably descended from the common wolf) that has been domesticated by man since prehistoric times; occurs in many breeds) "the dog barked all night"
 - [S: \(n\) canine, canid](#) (any of various fissiped mammals with nonretractile claws and typically long muzzles)
 - [S: \(n\) carnivore](#) (a terrestrial or aquatic flesh-eating mammal) "terrestrial carnivores have four or five clawed digits on each limb"
 - [S: \(n\) placental, placental mammal, eutherian, eutherian mammal](#) (mammals having a placenta; all mammals except monotremes and marsupials)
 - [S: \(n\) mammal, mammalian](#) (any warm-blooded vertebrate having the skin more or less covered with hair; young are born alive except for the small subclass of monotremes and nourished with milk)
 - [S: \(n\) vertebrate, craniate](#) (animals having a bony or cartilaginous skeleton with a segmented spinal column and a large brain enclosed in a skull or cranium)
 - [S: \(n\) chordate](#) (any animal of the phylum Chordata having a notochord or spinal column)
 - [S: \(n\) animal, animate being, beast, brute, creature, fauna](#) (a living organism characterized by voluntary movement)
 - [S: \(n\) organism, being](#) (a living thing that has (or can develop) the ability to act or function independently)
 - [S: \(n\) living thing, animate thing](#) (a living (or once living) entity)
 - [S: \(n\) whole, unit](#) (an assemblage of parts that is regarded as a single entity) "how big is that part compared to the whole?", "the team is a unit"
 - [S: \(n\) object, physical object](#) (a tangible and visible entity; an entity that can cast a shadow) "it was full of rackets, balls and other objects"
 - [S: \(n\) physical entity](#) (an entity that has physical existence)
 - [S: \(n\) entity](#) (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

Deng et al. "Imagenet: a large scale hierarchical image database" CVPR 2009

**mite****container ship****motor scooter****leopard**

mite	container ship	motor scooter	leopard
black widow	lifeboat	go-kart	jaguar
cockroach	amphibian	moped	cheetah
tick	fireboat	bumper car	snow leopard
starfish	drilling platform	golfcart	Egyptian cat

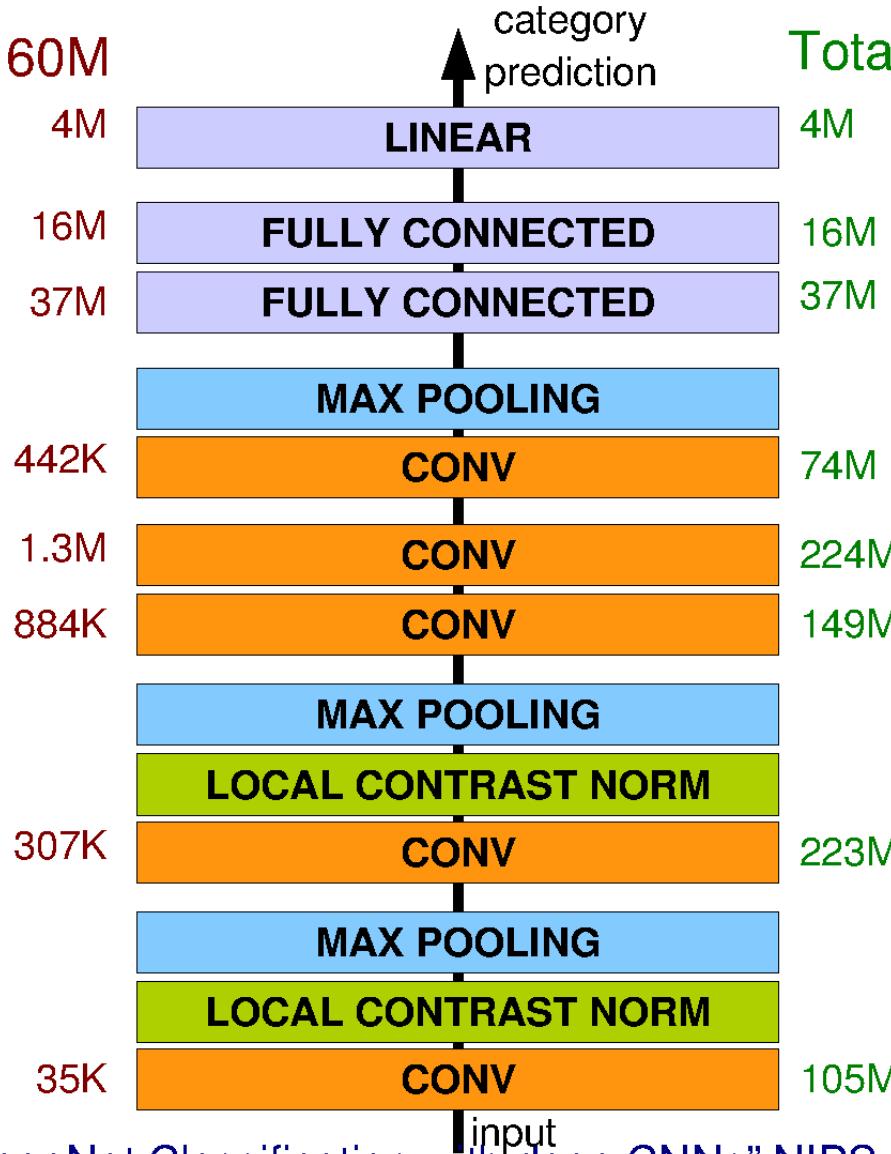
**grille****mushroom****cherry****Madagascar cat**

convertible	agaric	dalmatian	squirrel monkey
grille	mushroom	grape	spider monkey
pickup	jelly fungus	elderberry	titi
beach wagon	gill fungus	ffordshire bullterrier	indri
fire engine	dead-man's-fingers	currant	howler monkey

Architecture for Classification

Total nr. params: 60M

Total nr. flops: 832M

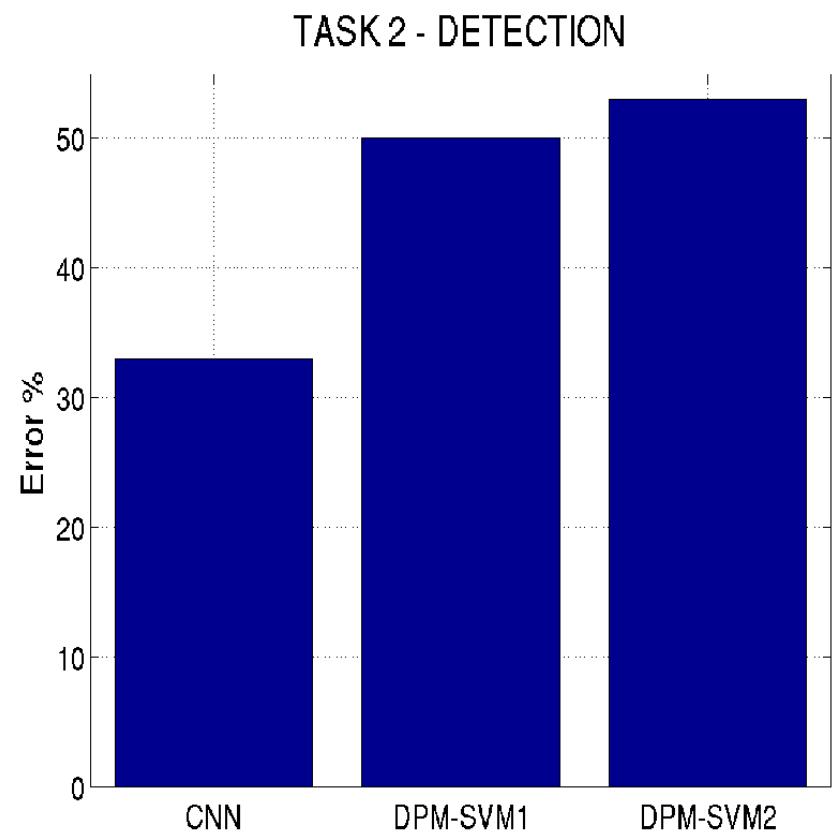
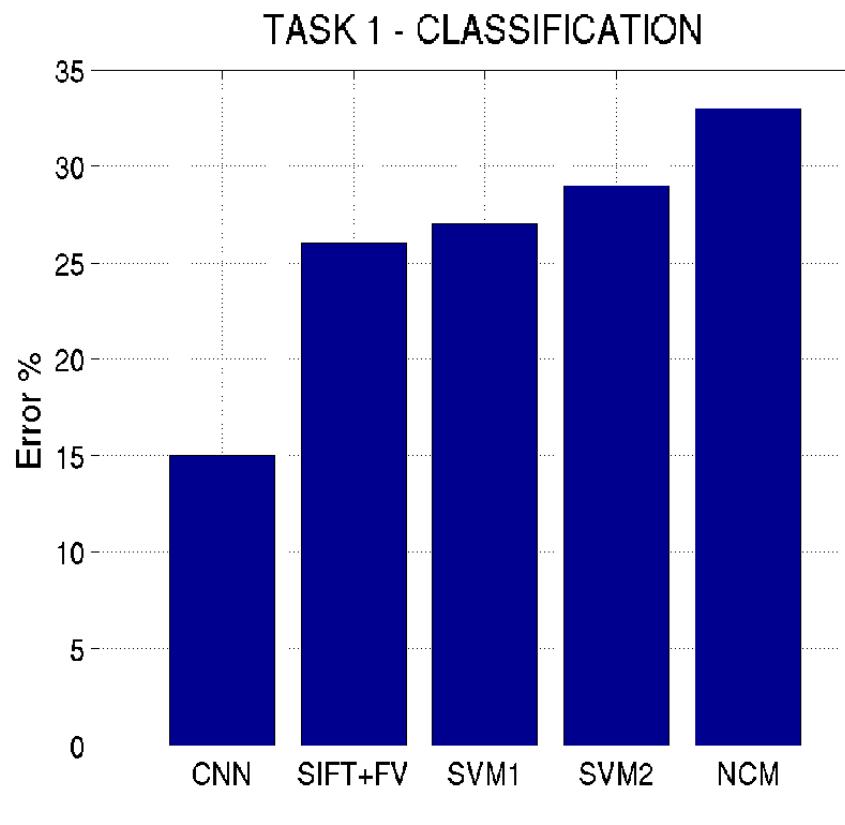


Krizhevsky et al. "ImageNet Classification with deep CNNs" NIPS 2012

96

Ranzato

Results: ILSVRC 2012



98

Krizhevsky et al. "ImageNet Classification with deep CNNs" NIPS 2012

Ranzato