

Deep Learning in Computer Vision – Convolutional Neural Network

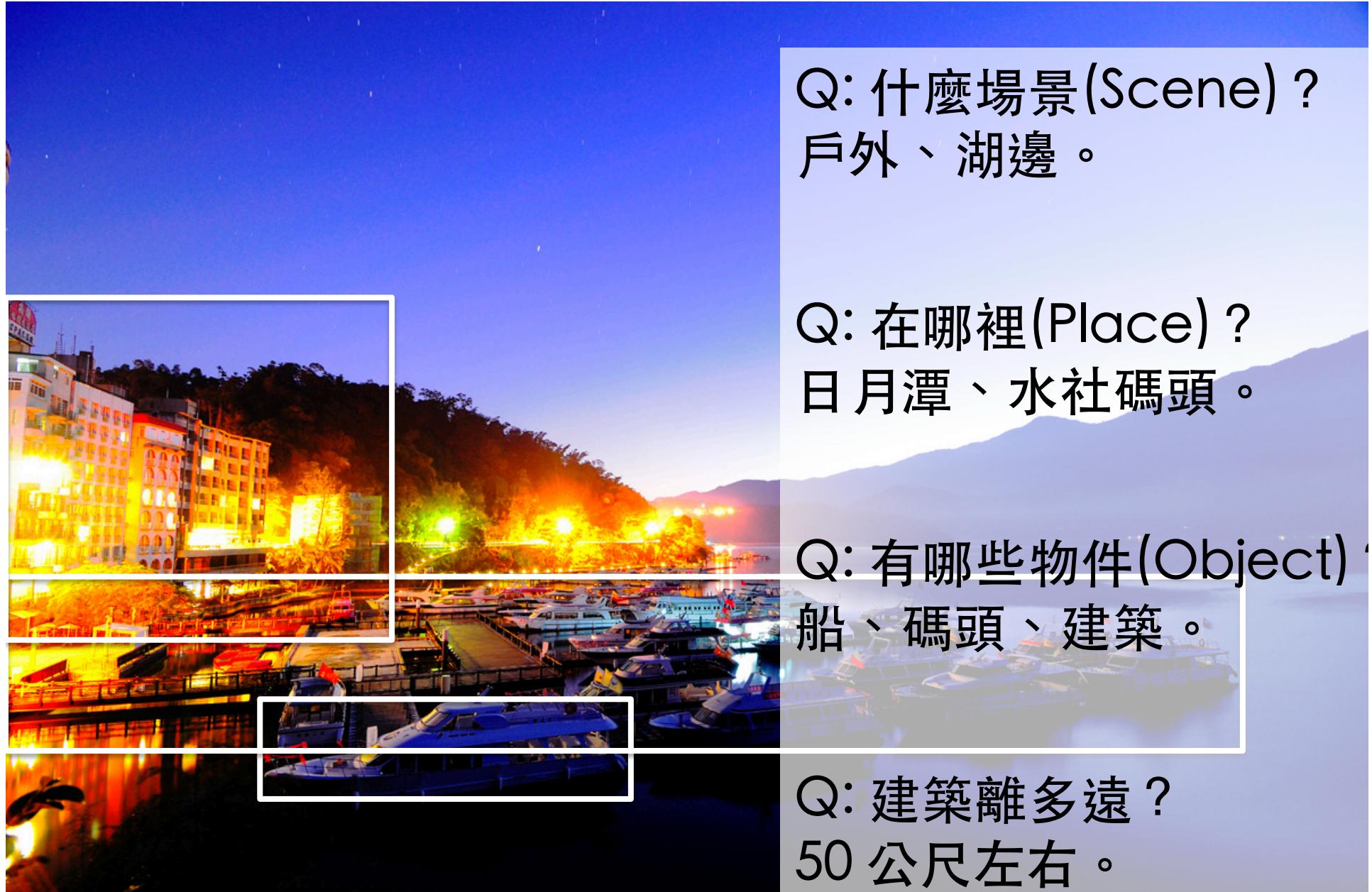
國立清華大學

孫民教授



VSLab

電腦視覺 Computer Vision



電腦視覺 Computer Vision



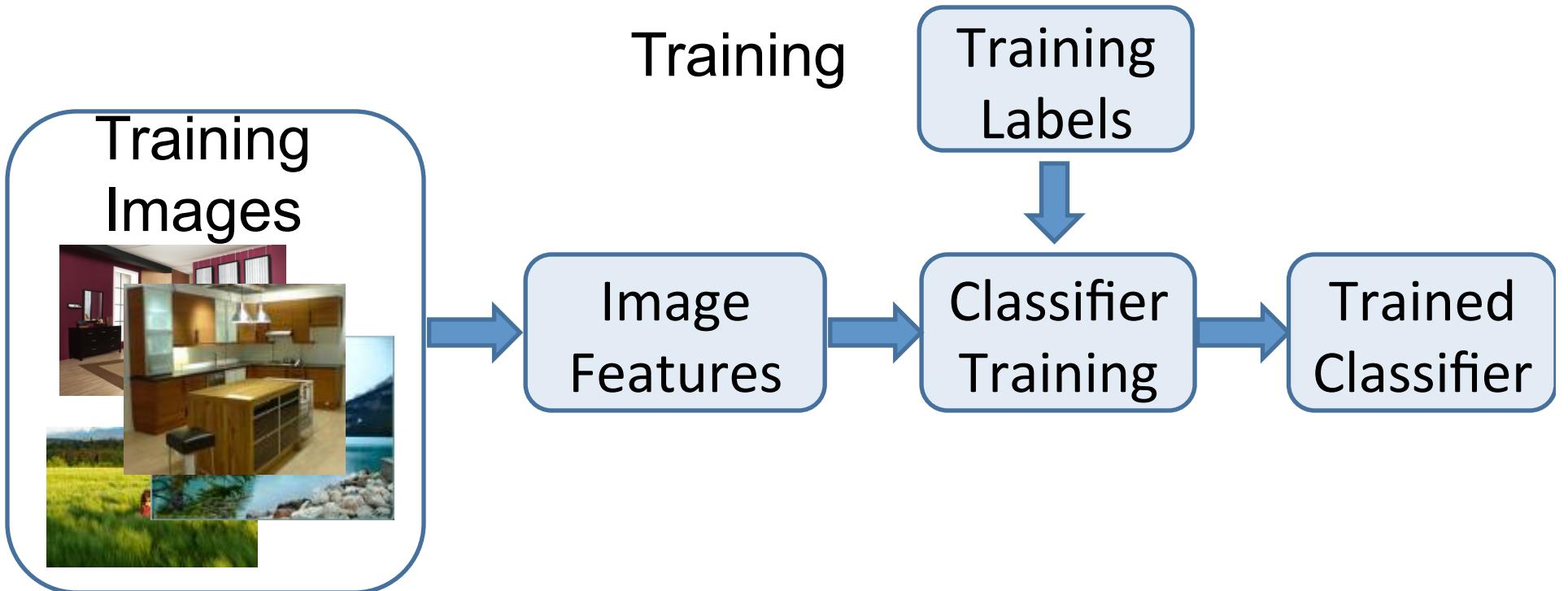
Q: 什麼場景(Scene) ?
戶外、湖邊。
(Classification)

Q: 在哪裡(Place) ?
日月潭、水社碼頭。
(Classification)

Q: 有哪些物件(Object) ?
船、碼頭、建築。
(Classification)

Q: 建築離多遠 ?
50 公尺左右。

Image Classification



Slides: Jame Hays,
Isabelle Guyon,
Erik Sudderth,
Mark Johnson,
Derek Hoiem

Image Classification

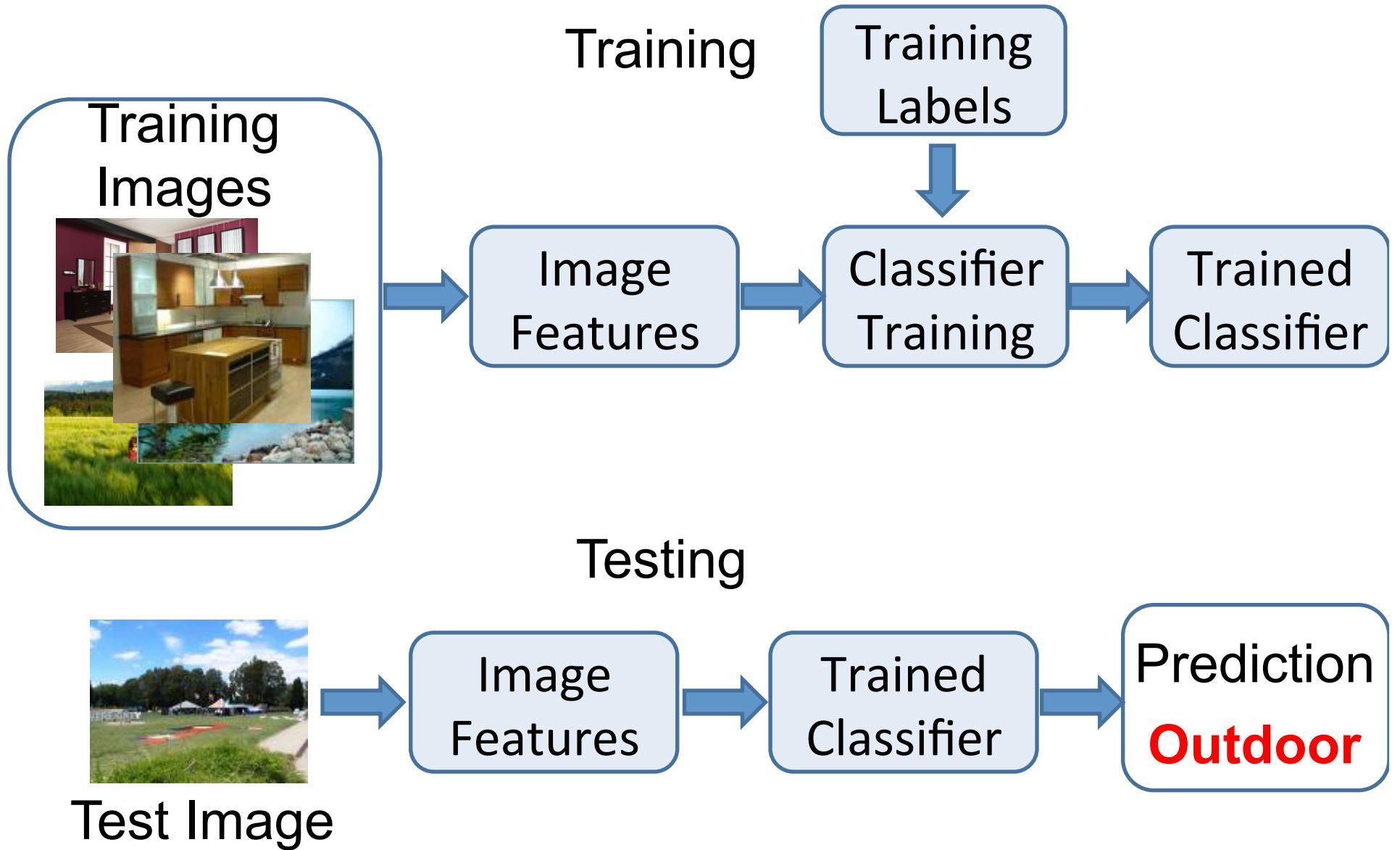
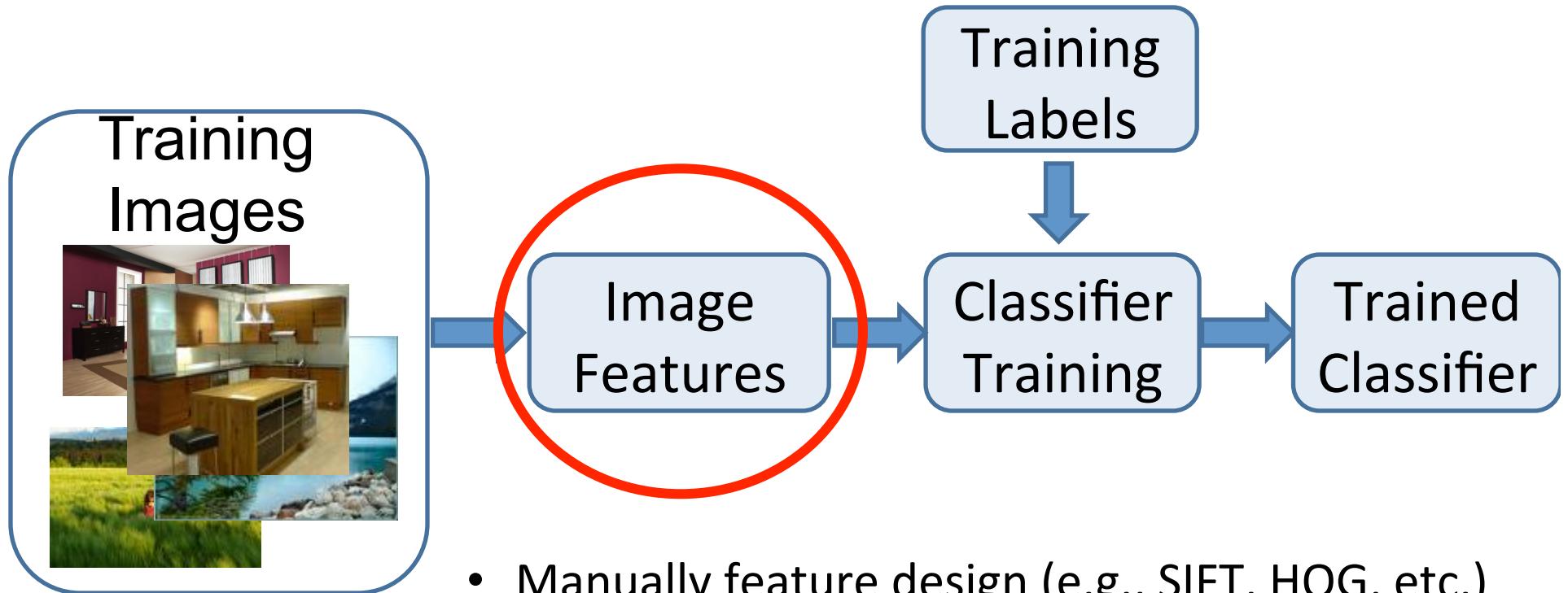


Image Classification

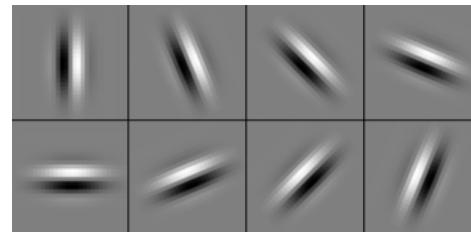


- Manually feature design (e.g., SIFT, HOG, etc.)
- Multiple separately learned layers (e.g., K-means, Pyramid, etc.)

SIFT Descriptor

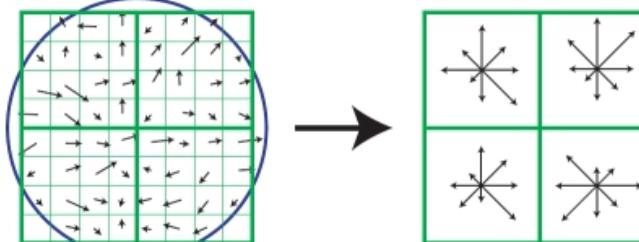
Image
Pixels

Apply
oriented filters

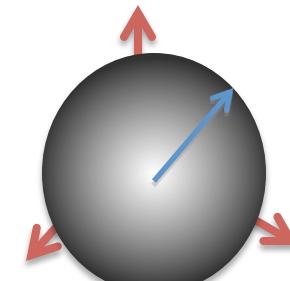


Lowe [IJCV 2004]

Spatial pool
(Sum)



Normalize to
unit length



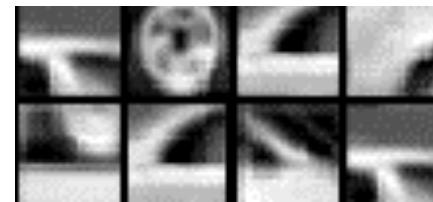
Feature
Vector

slide credit: R. Fergus

Spatial Pyramid Matching

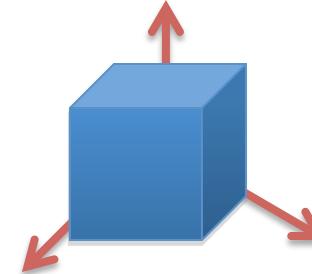
SIFT
Features

Filter with
Visual Words

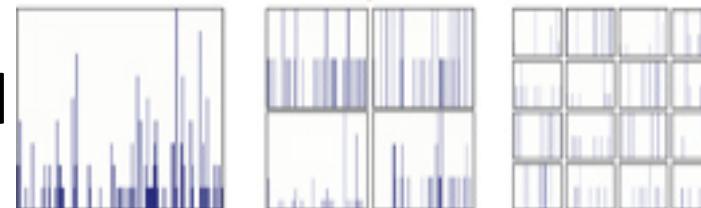


Lazebnik,
Schmid,
Ponce
[CVPR 2006]

Max



Multi-scale
spatial pool
(Sum)

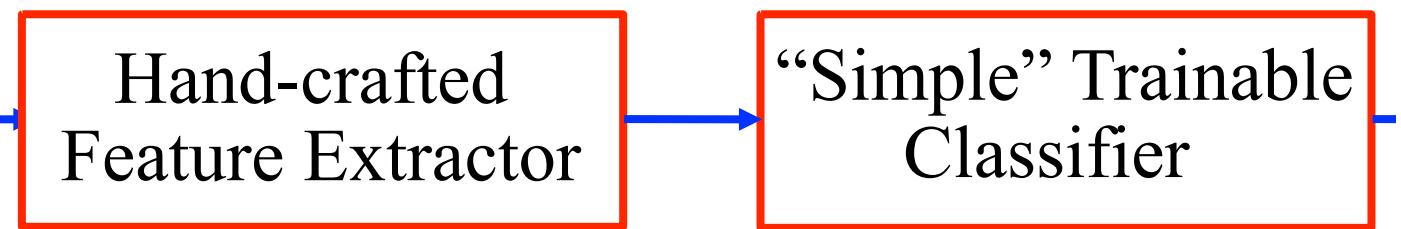


Classifier

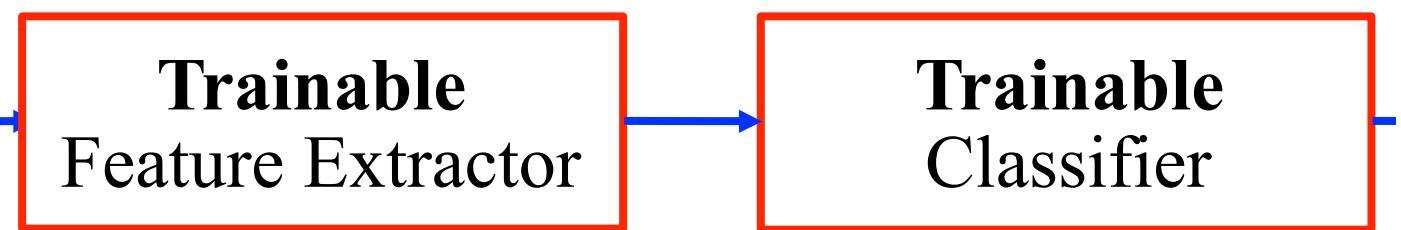
slide credit: R. Fergus

Comparison

- Classical CV

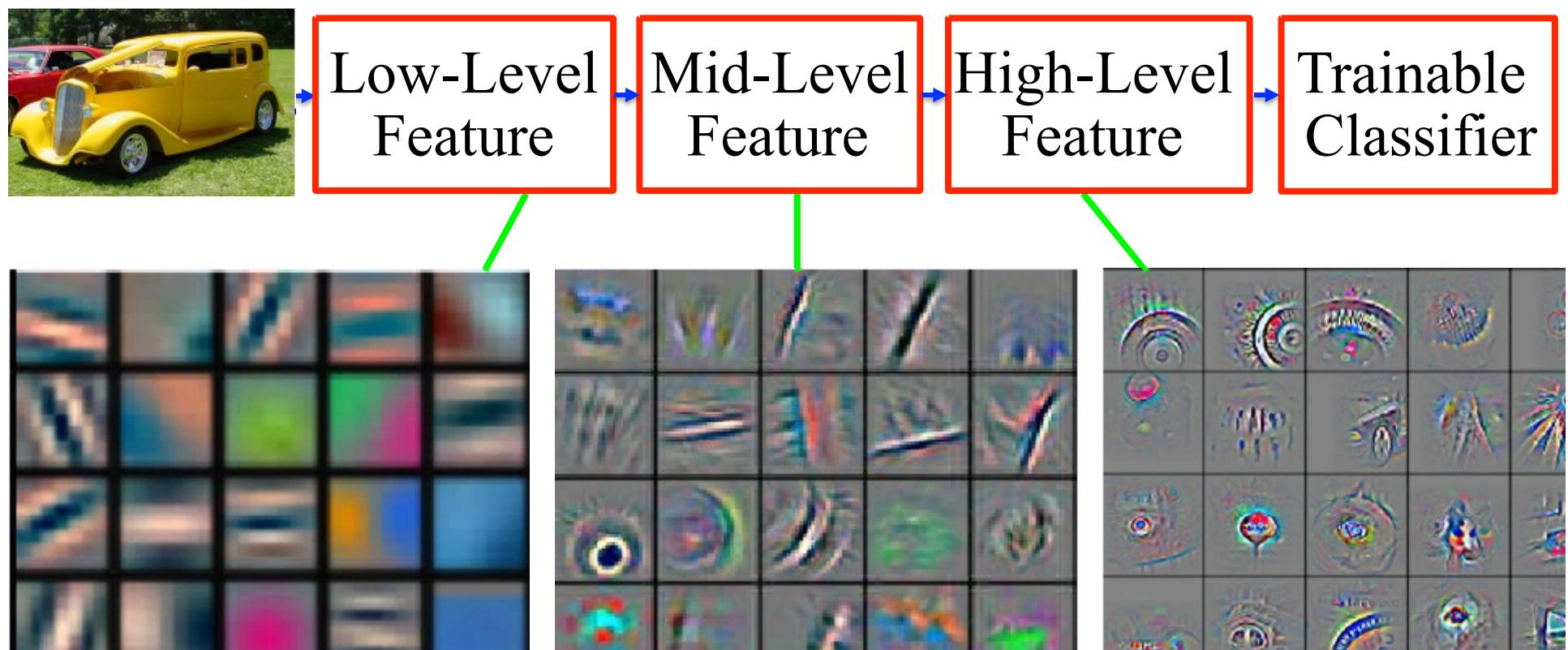
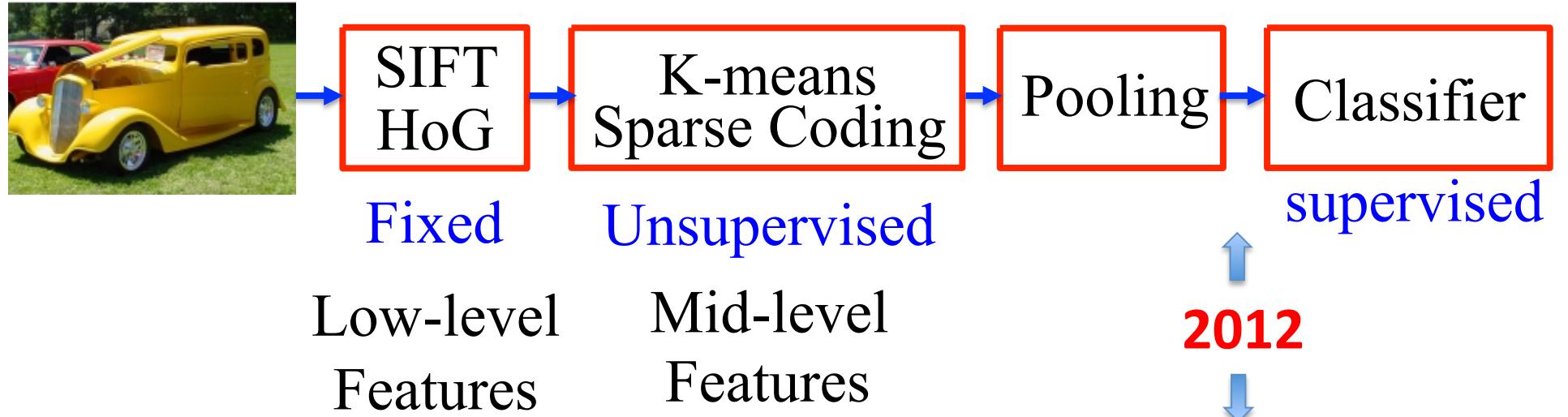


- Deep Learning



End-to-End Learning

Slides: Y LeCun
MA Ranzato



1K Image Classification

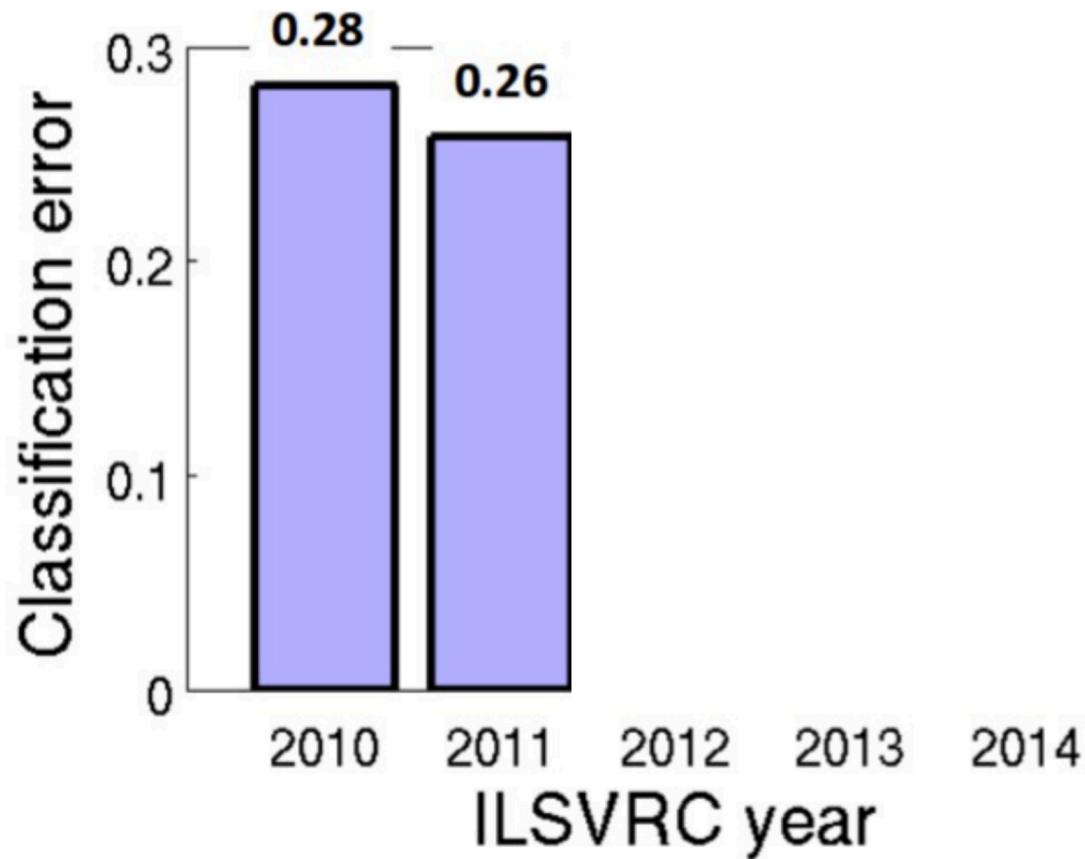


Figure from Olga Russakovsky ECCV'14 workshop

Deep
Learning



1K Image Classification

Convolution Neuron Networks
(CNN)

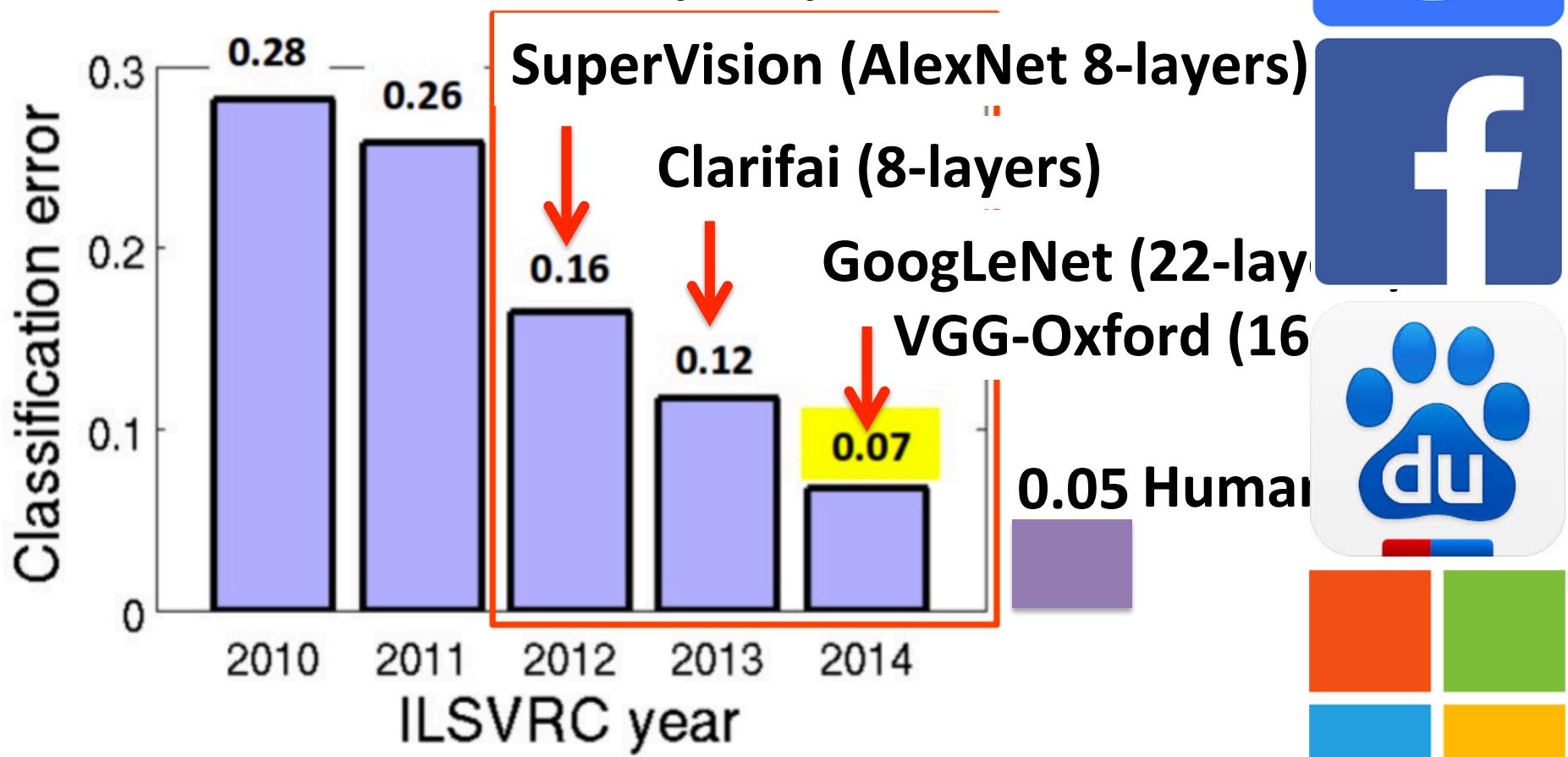
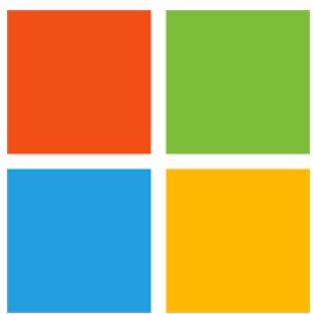
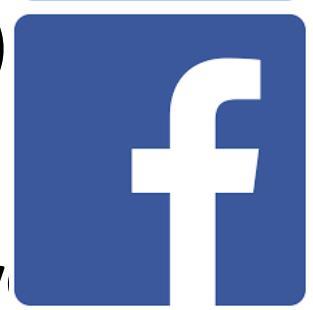
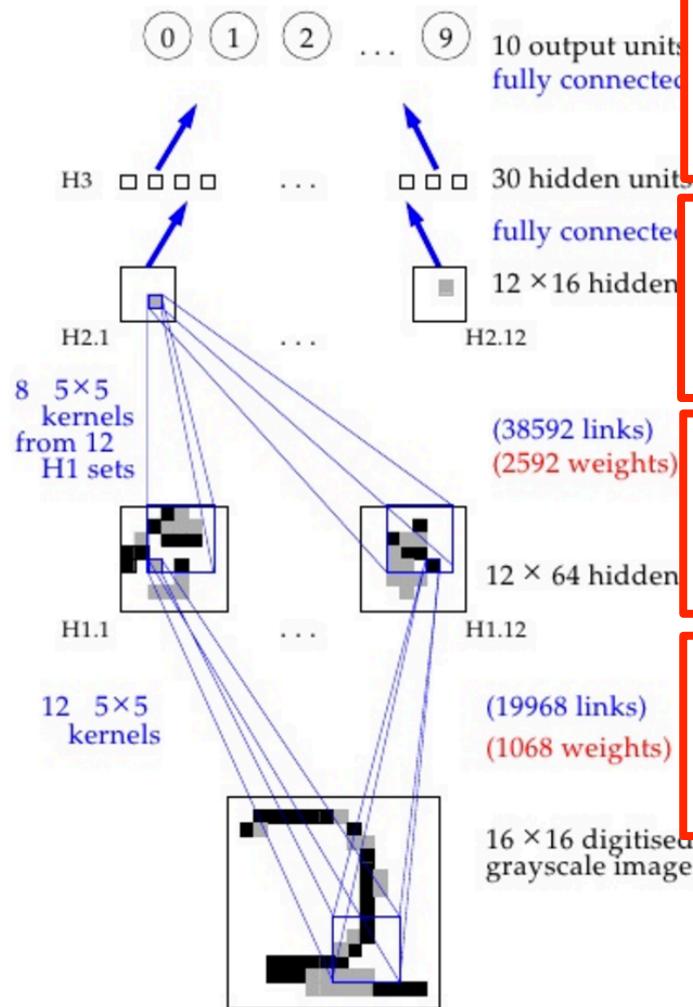


Figure from Olga Russakovsky ECCV'14 workshop



Classic CNN: LeNet

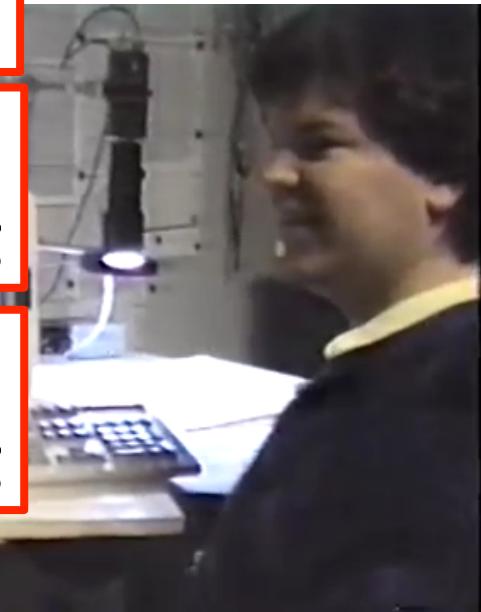


Fully Connected
+ Non-linearity

Fully Connected
+ Non-linearity

Convolution
+Spatial Pooling

Convolution
+Spatial Pooling



LeNet-1, 1993

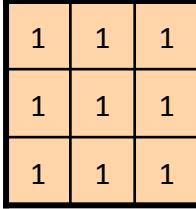
Before weight sharing 64660 links

After weight sharing 9760 weights

Backpropagation Applied to Handwritten Zip Code. Yann Lecun, 1989

What is 2D Convolution?

Image filtering

$$g[\cdot, \cdot] \frac{1}{9}$$


$$f[., .]$$

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

$$h[., .]$$

		0						

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

Image filtering

$f[.,.]$

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

$h[.,.]$

0	10									

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

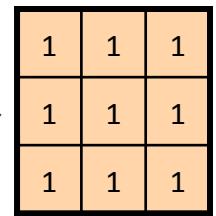


Image filtering

$f[.,.]$

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

$h[.,.]$

	0	10	20							

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

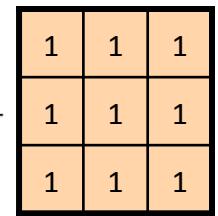


Image filtering

$f[.,.]$

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

$h[.,.]$

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

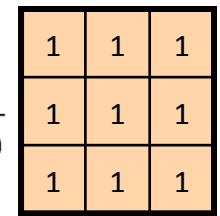


Image filtering

$f[.,.]$

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

$h[.,.]$

			0	10	20	30	30				

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

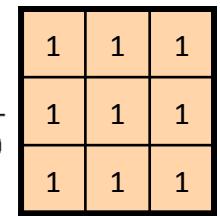


Image filtering

$$g[\cdot, \cdot] \begin{matrix} 1 \\ 9 \end{matrix} \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ \hline \end{array}$$

$$f[.,.]$$

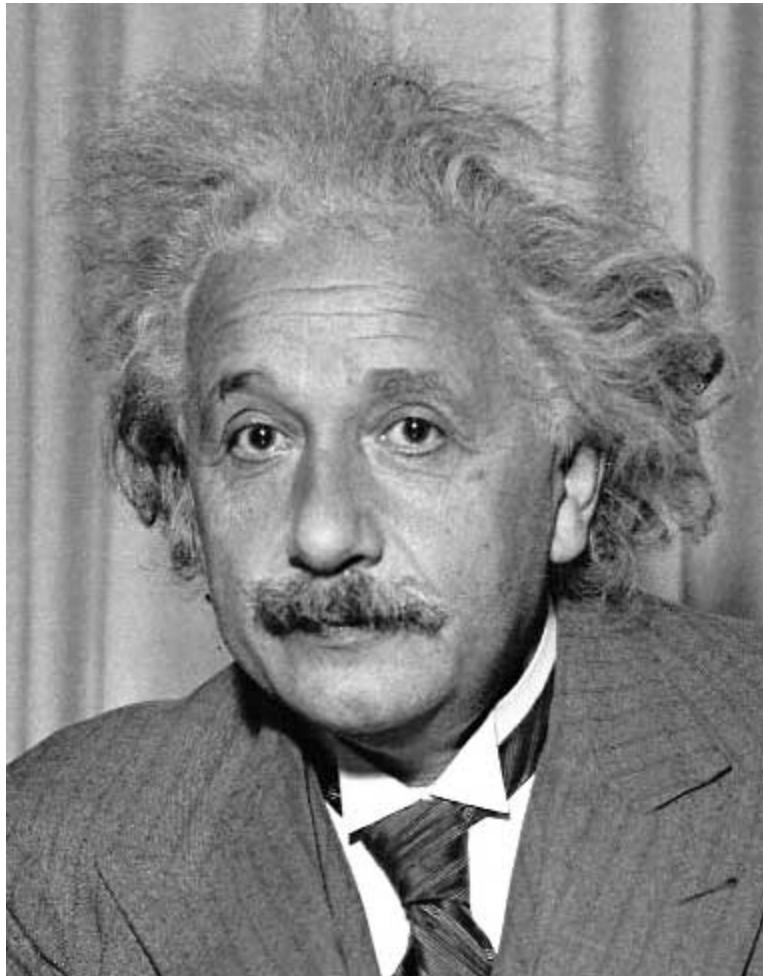
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$h[.,.]$$

	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	

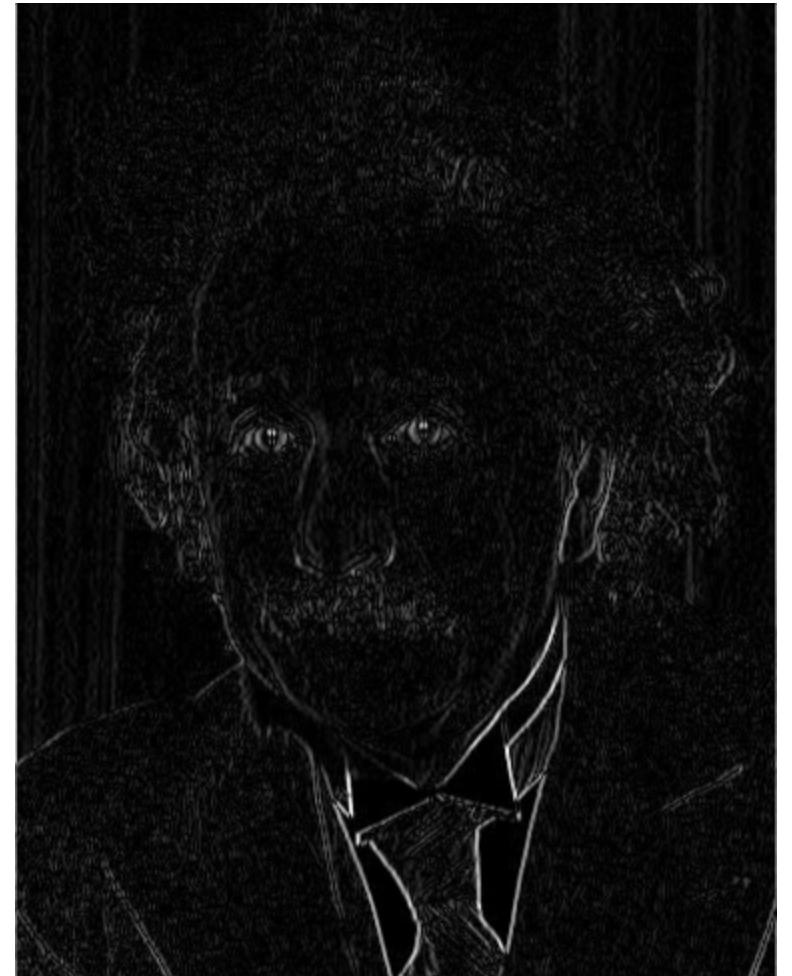
$$h[m,n] = \sum_{k,l} g[k,l] f[m+k, n+l]$$

Image filtering



1	0	-1
2	0	-2
1	0	-1

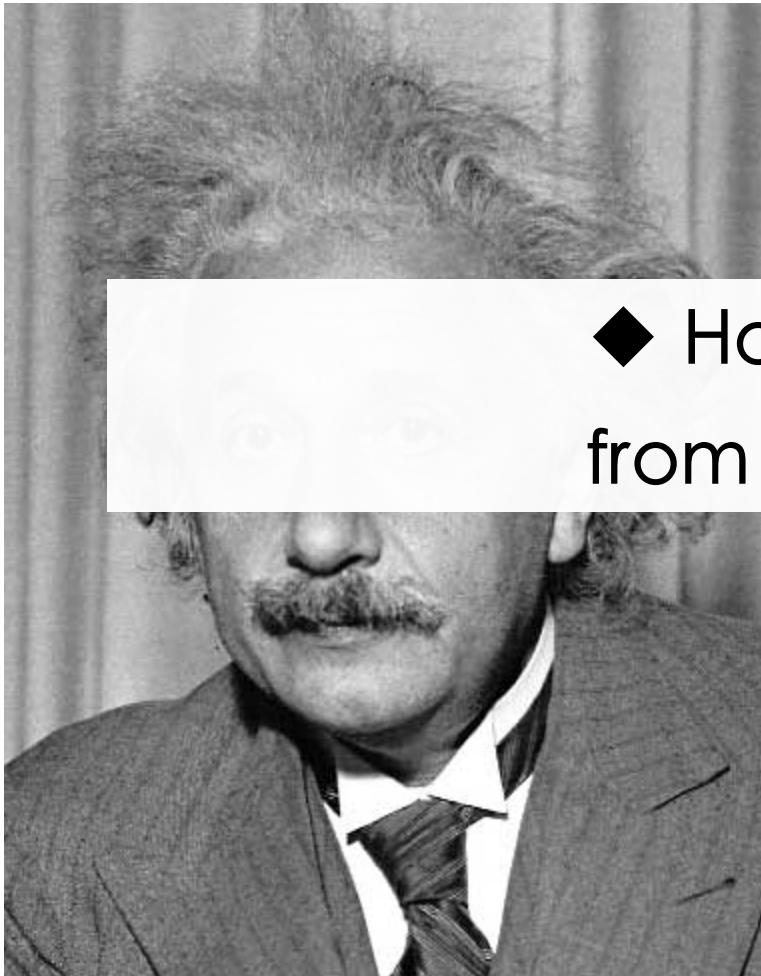
Sobel



Vertical Edge
(absolute value) 23

http://en.wikipedia.org/wiki/Sobel_operator

Image filtering



◆ How to learn filters
from **Big Visual** data?

0	0	0
-1	-2	-1

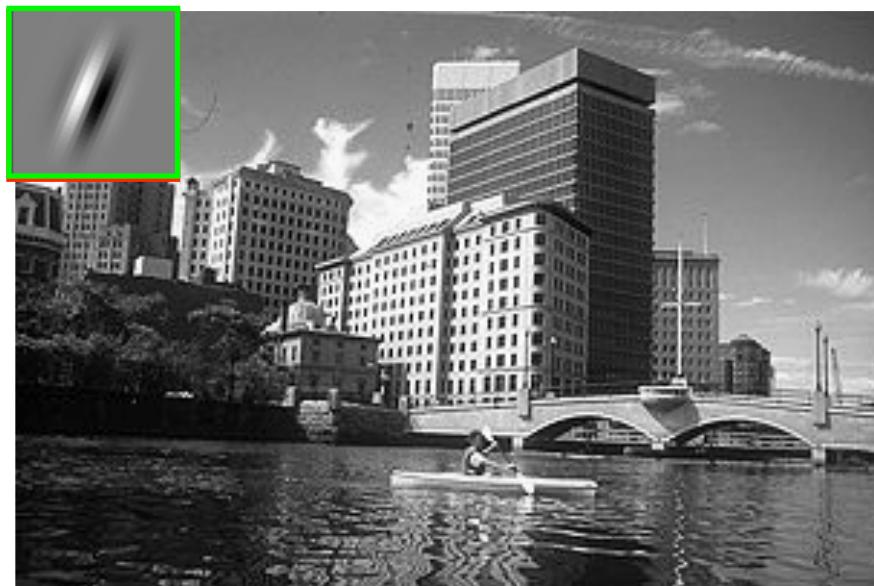
Sobel



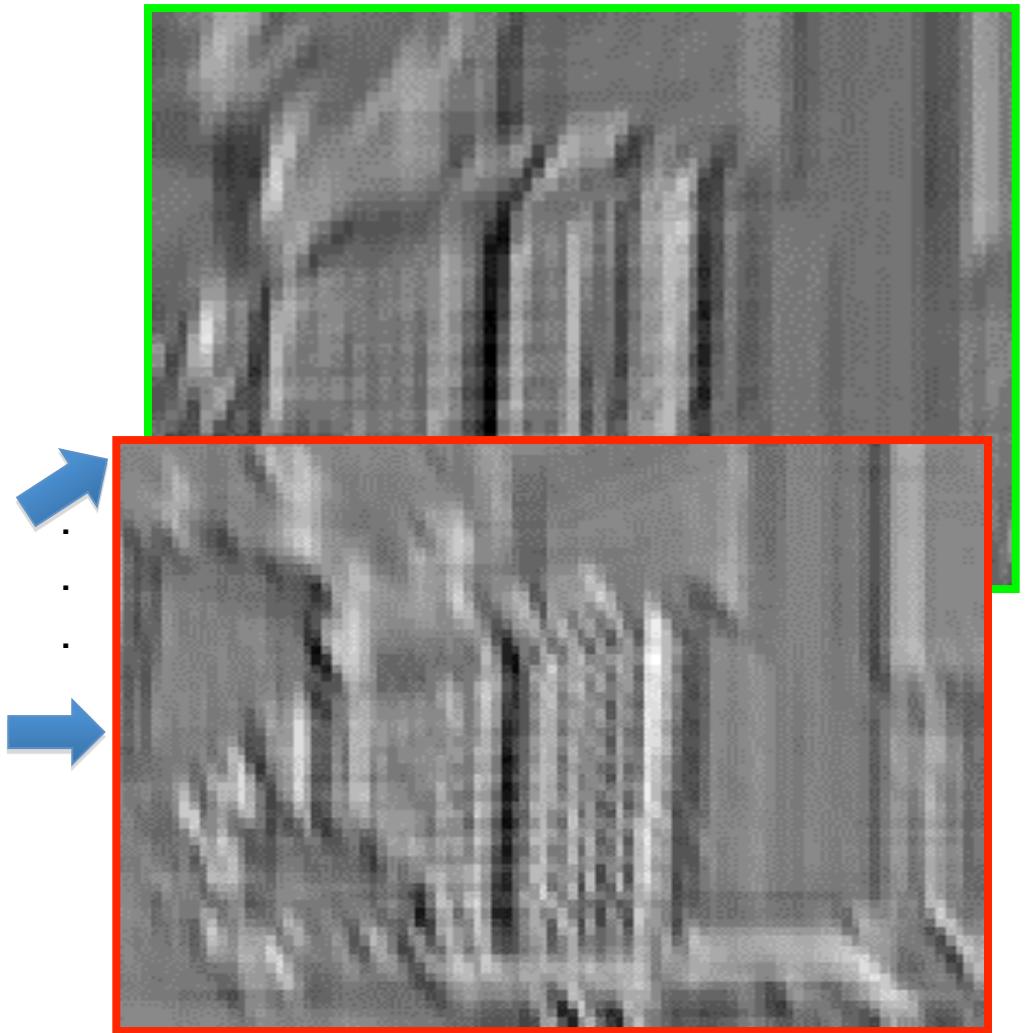
Horizontal Edge
(absolute value) 24

What is 2D Convolution?

- Weighted moving sum



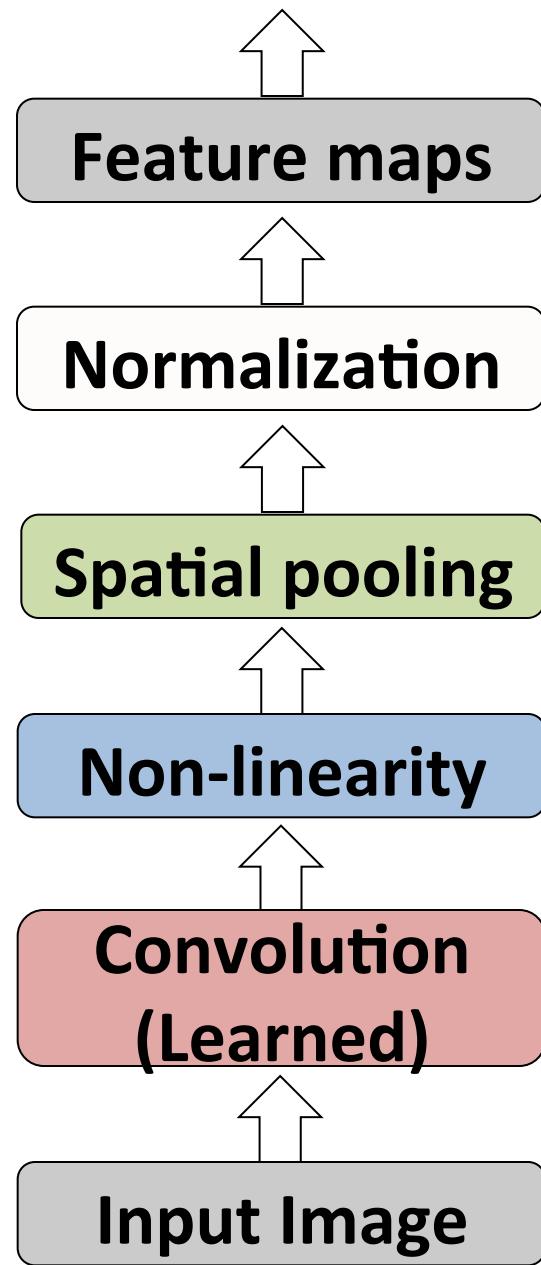
Input



Feature Activation Map

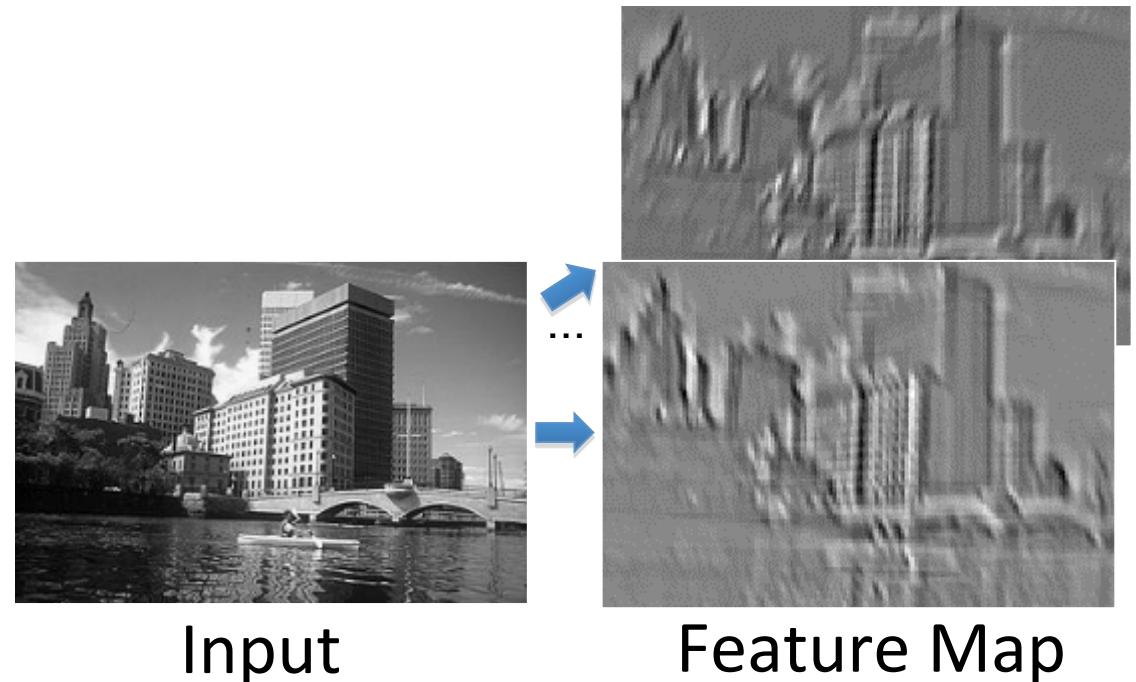
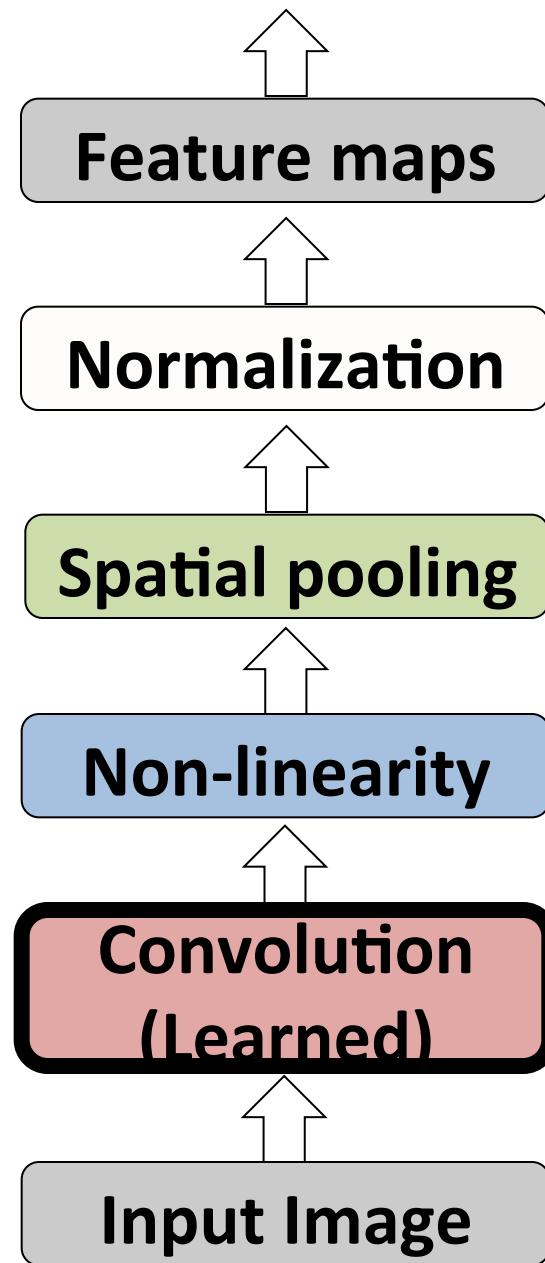
slide credit: S. Lazebnik

Convolutional Neural Networks



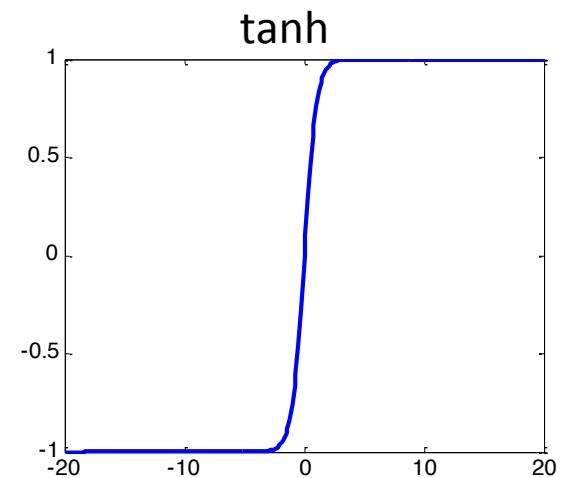
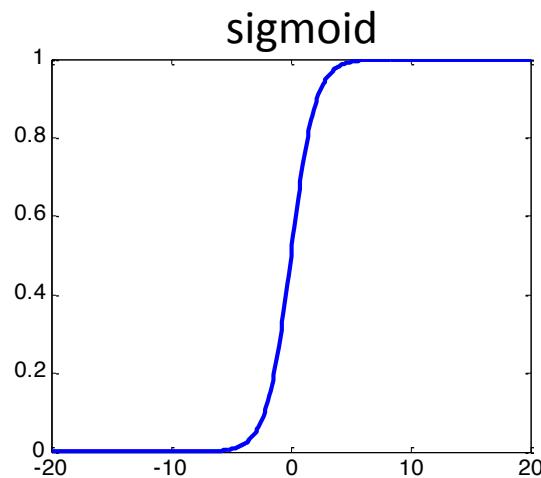
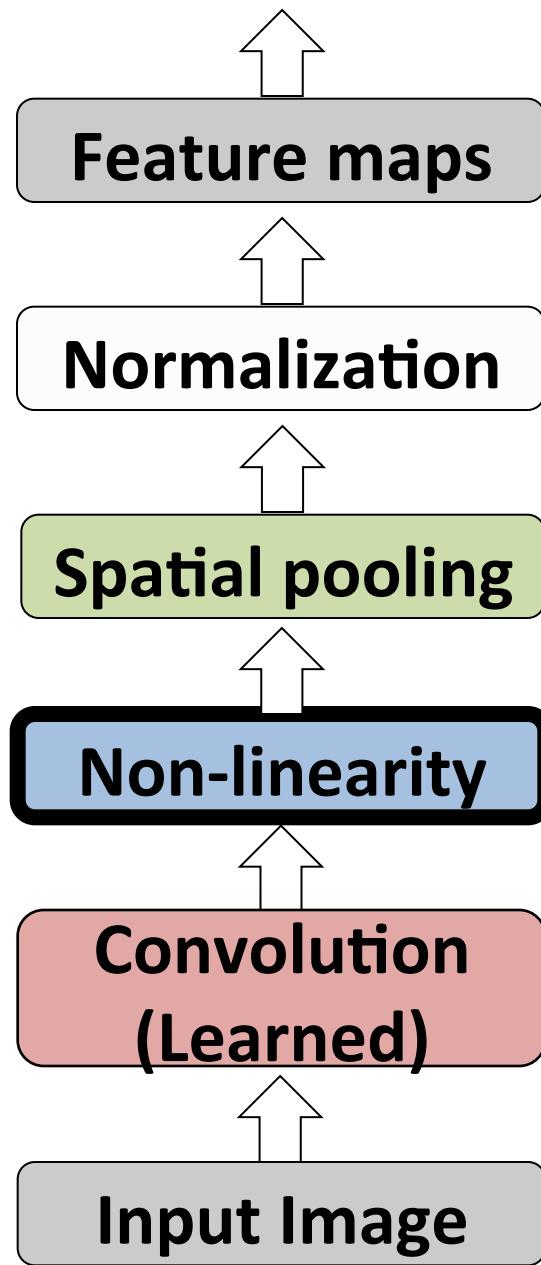
slide credit: S. Lazebnik

Convolutional Neural Networks

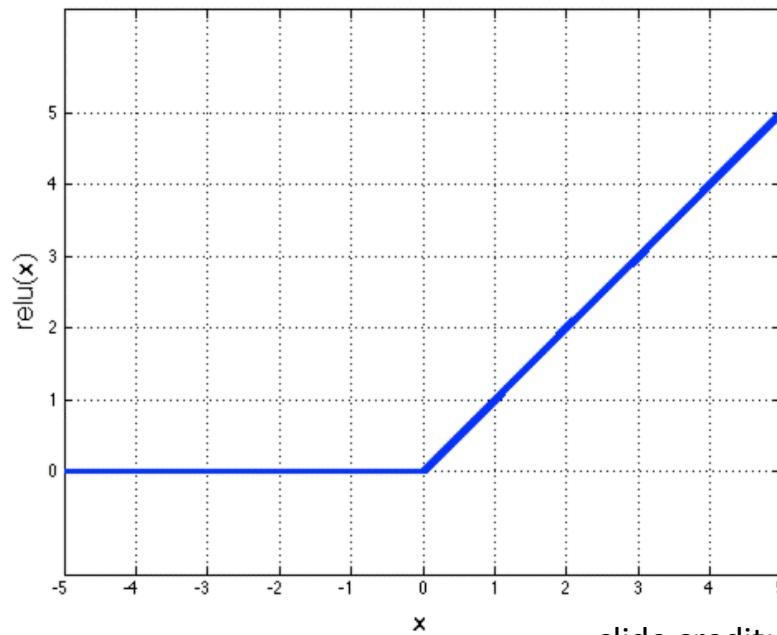


slide credit: S. Lazebnik

Convolutional Neural Networks

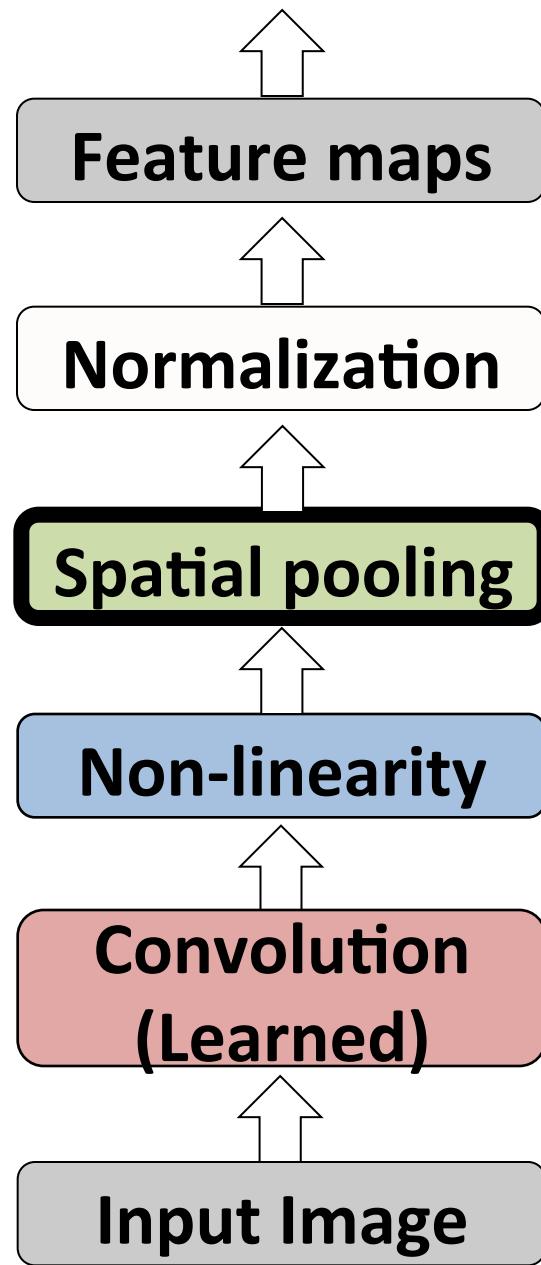


Rectified Linear Unit (ReLU)

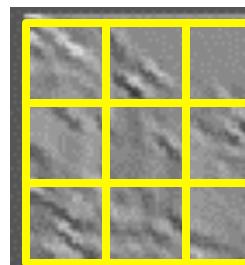


slide credit: S. Lazebnik

Convolutional Neural Networks



KernelSize



Stride

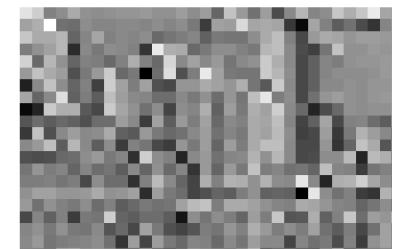


Max-pooling:

a **non-linear** down-sampling

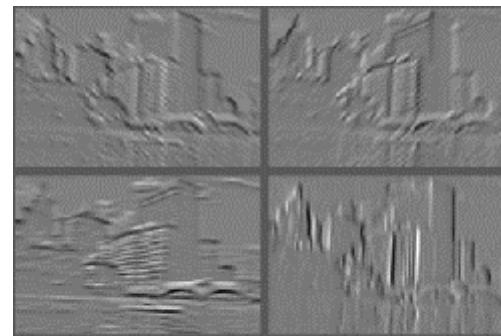
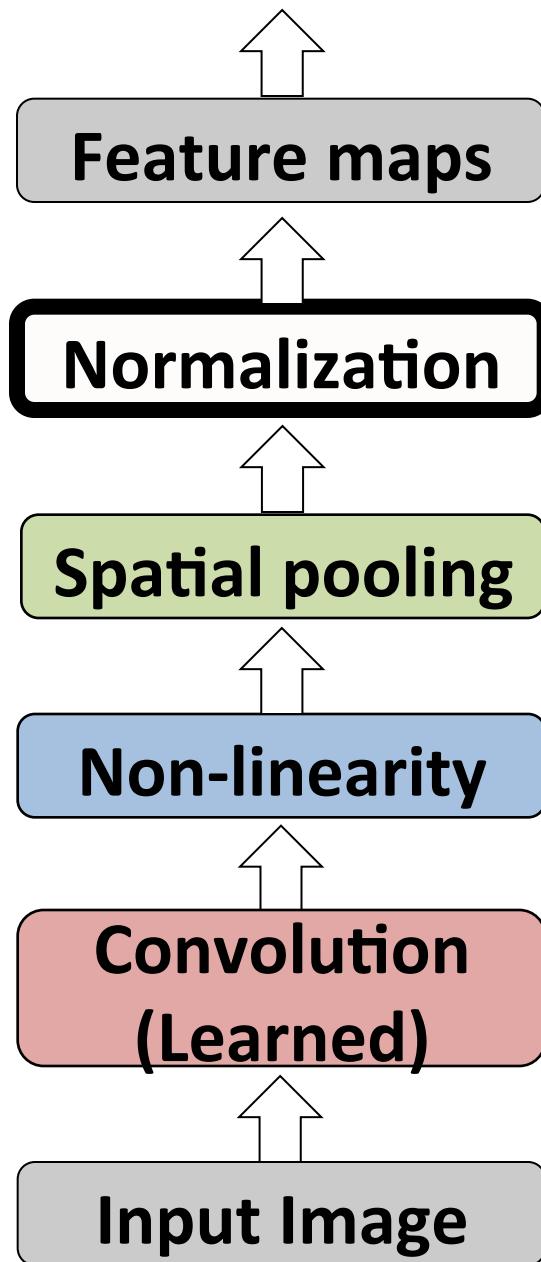
Provide *translation invariance*

Max pooling

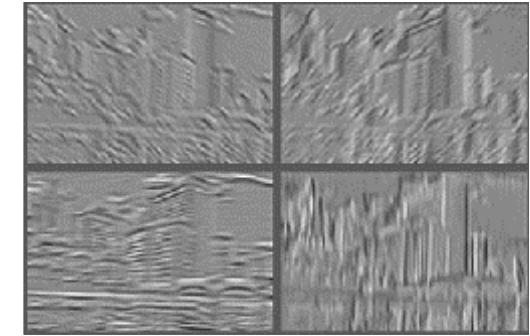


slide credit: S. Lazebnik

Convolutional Neural Networks



Feature Maps

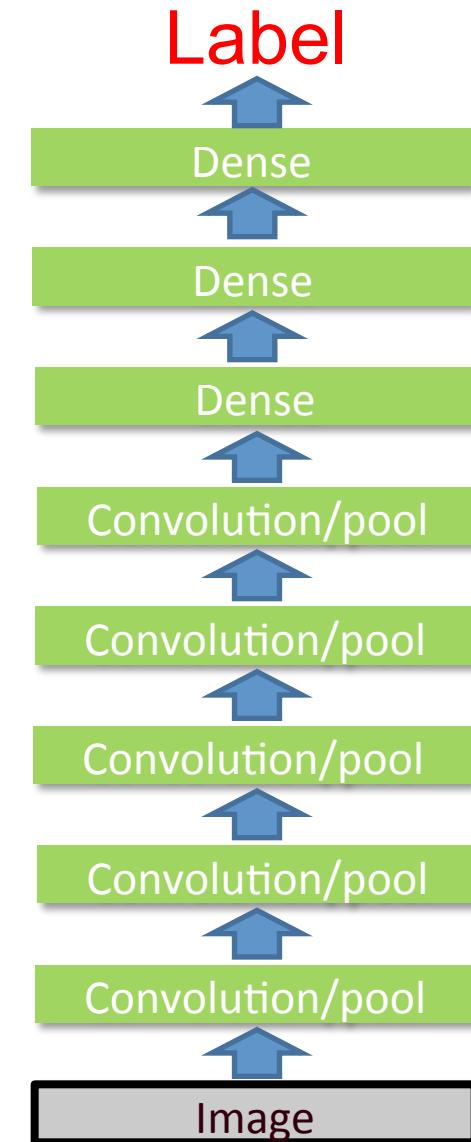
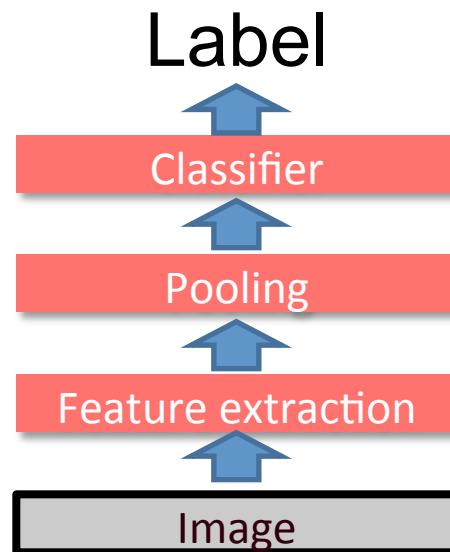


Feature Maps
After Contrast
Normalization

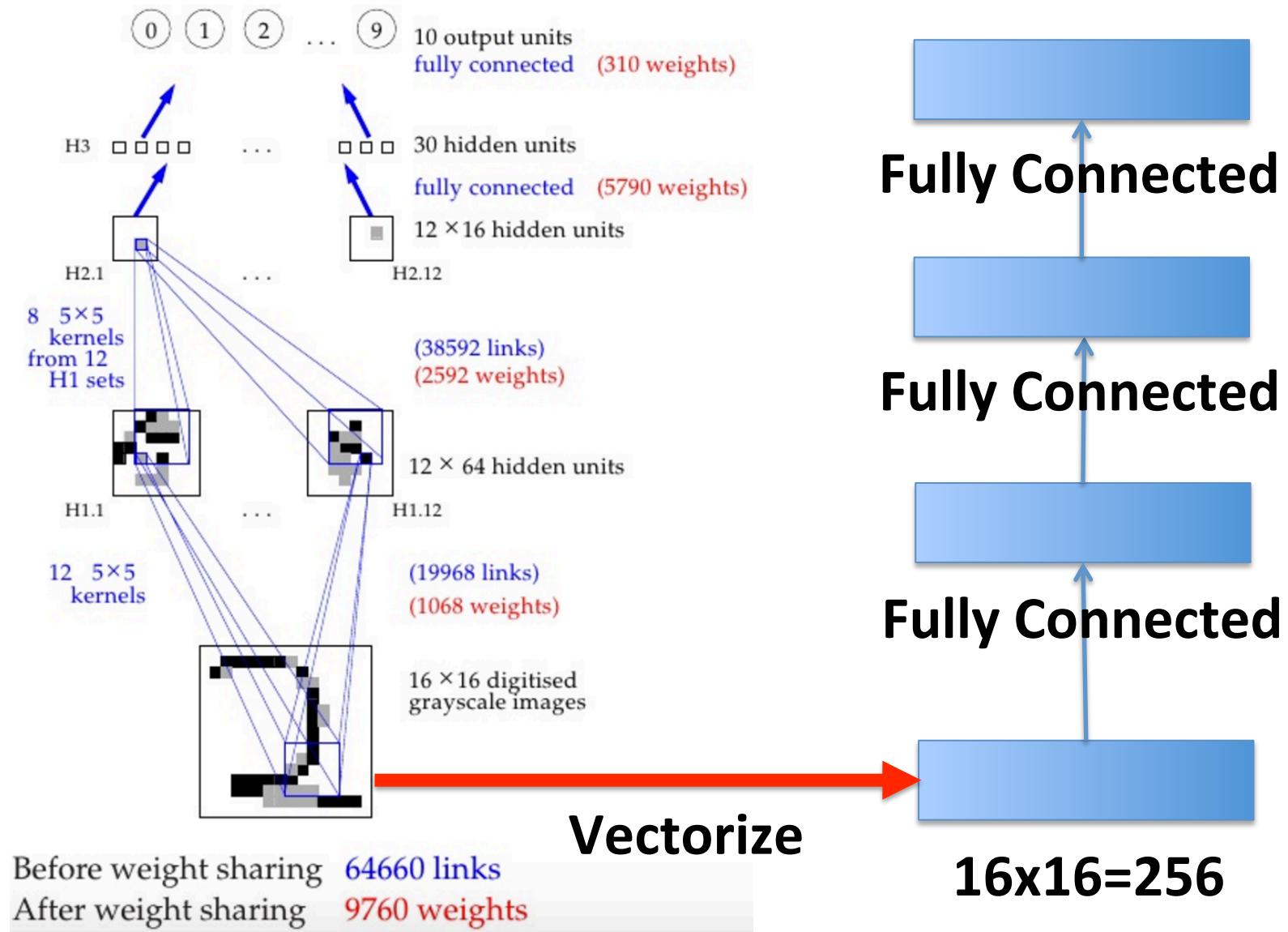
Equalizes the feature maps locally:
weak signals are boosted, whereas
strong signals are suppressed.

Engineered vs. learned features

Convolutional filters are trained in a supervised manner by **back-propagating** classification error

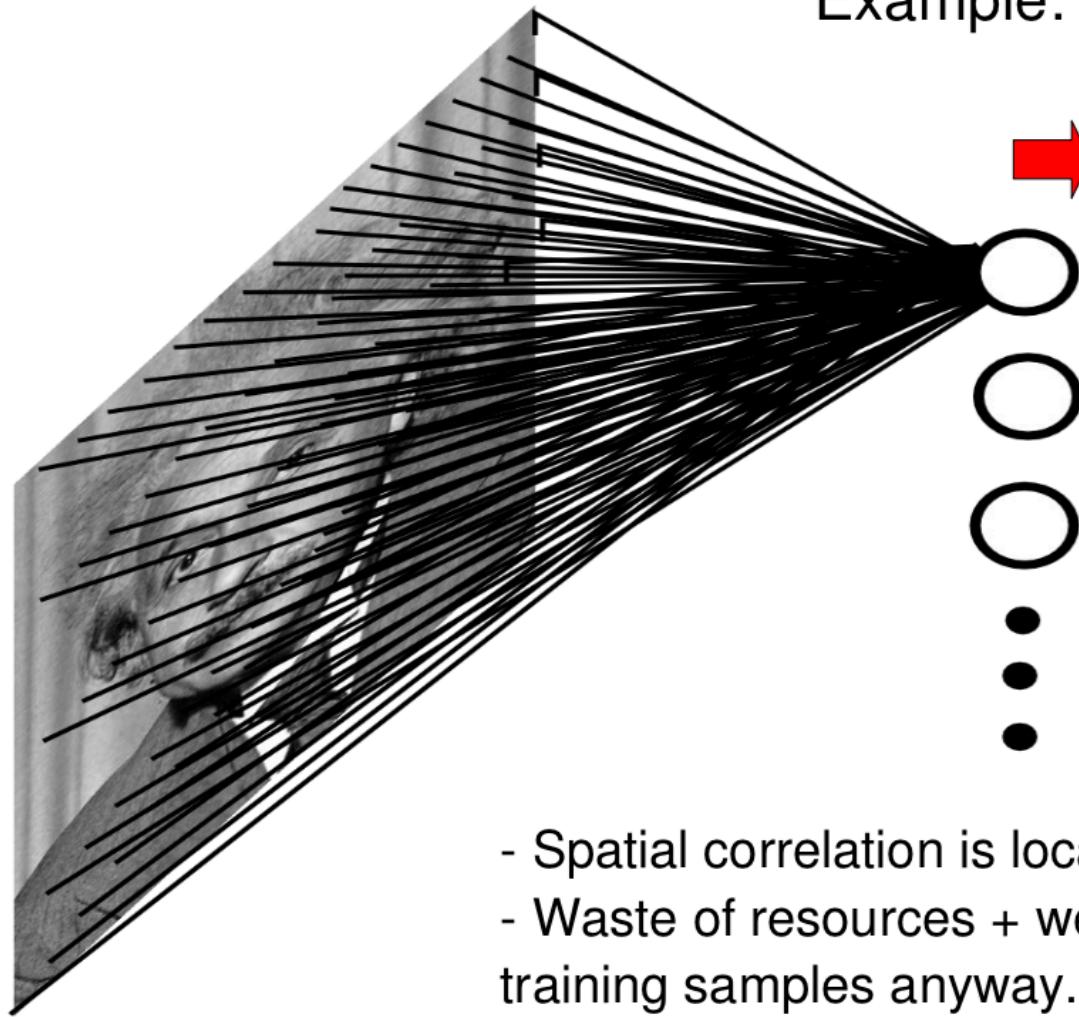


Why Convolution?



Yann Lecun, 1989

Fully Connected Layer



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

- Spatial correlation is local
- Waste of resources + we have not enough training samples anyway..

33