# Category Recognition



Jia-Bin Huang

Virginia Tech
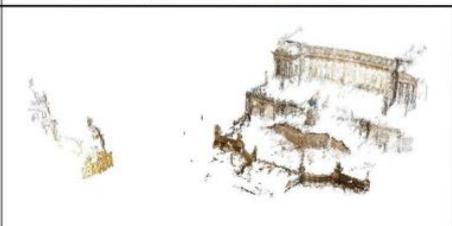
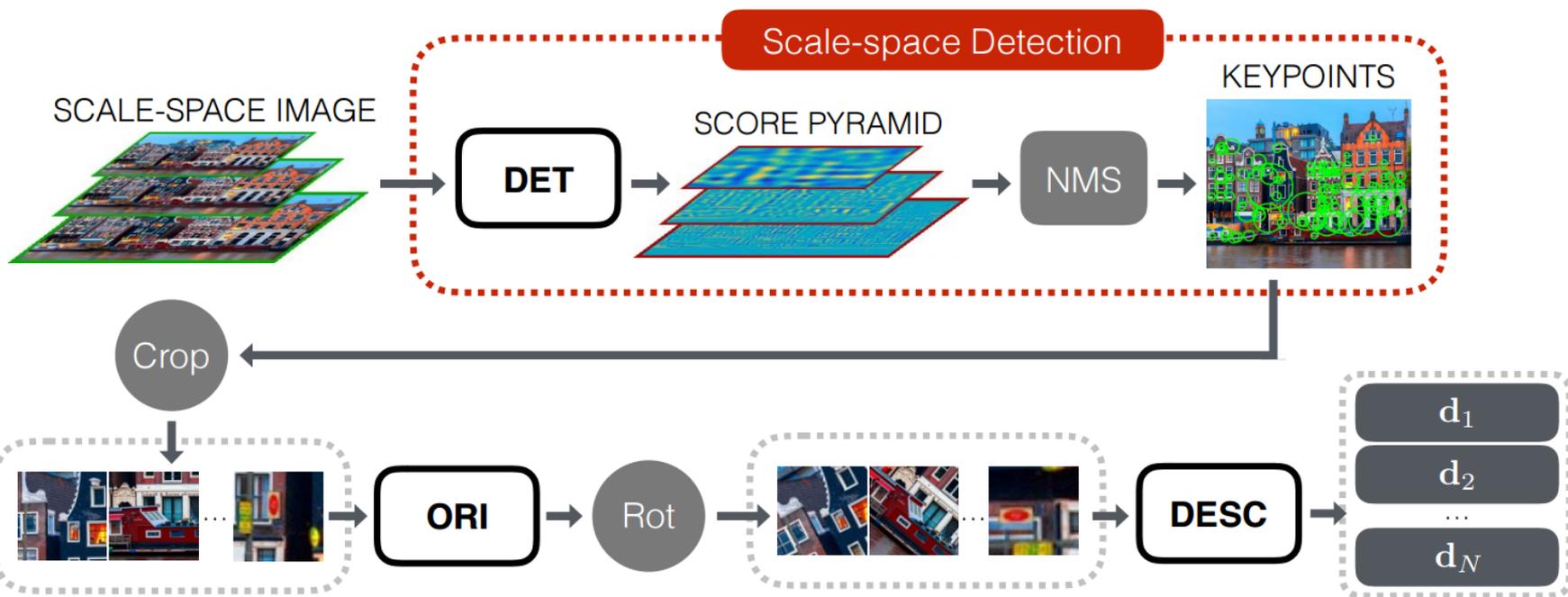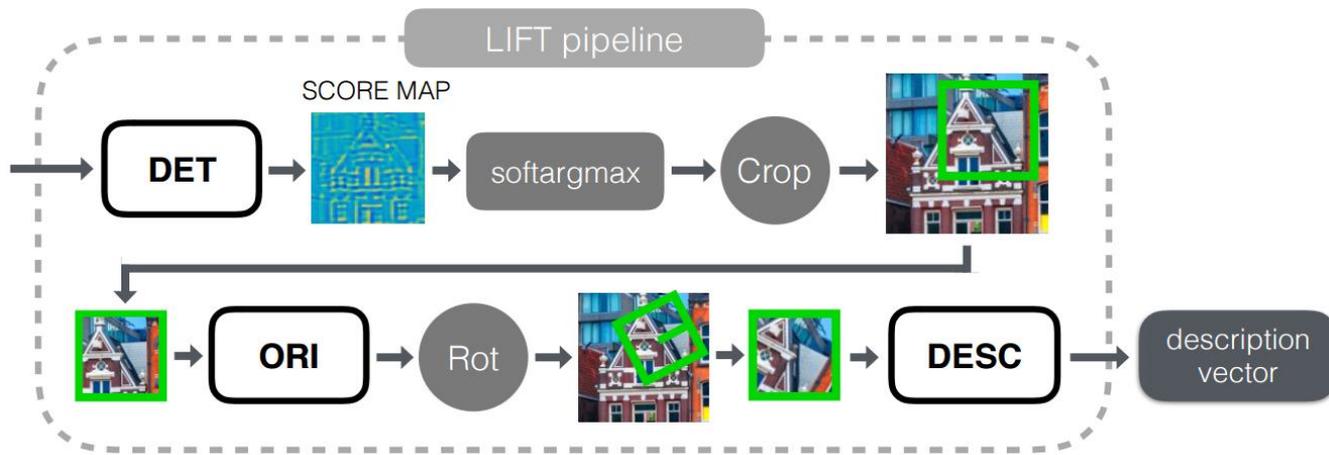ECE 6554 Advanced Computer Vision

# Administrative stuffs

- Presentation and discussion leads assigned
  - https://docs.google.com/spreadsheets/d/1P5pfyCio5flq3QCy4Mo1XS66I6d14jqDxE2Tny4efVs/edit#gid=0


- Questions?

# Today's class

- Finish instance recognition

- Category recognition

- Convolutional neural network

| Day Image | Day Model | Night Model | Fused Night Model | Night Image |
|---|---|---|---|---|



From Dusk till Dawn: Modeling in the Dark, CVPR 2016

LIFT pipeline

SCORE MAP

DET → softargmax → Crop

ORI → Rot → DESC → description vector

Scale-space Detection

SCALE-SPACE IMAGE

DET → SCORE PYRAMID → NMS → KEYPOINTS

Crop

ORI → Rot → DESC → $\mathbf{d}_1$ $\mathbf{d}_2$ ... $\mathbf{d}_N$

Lift: Learned invariant feature transform, ECCV 2016
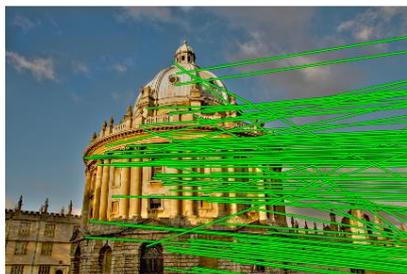
# Instance recognition

- Motivation – visual search
- Visual words
  - quantization, index, bags of words
- Spatial verification
  - affine; RANSAC, Hough
- Other text retrieval tools
  - tf-idf, query expansion
- Example applications

# Instance recognition: remaining issues

- How to summarize the content of an entire image? And gauge overall similarity?

- How large should the vocabulary be?  How to perform quantization efficiently?

- Is having the same set of visual words enough to identify the object/scene?  How to verify spatial agreement?

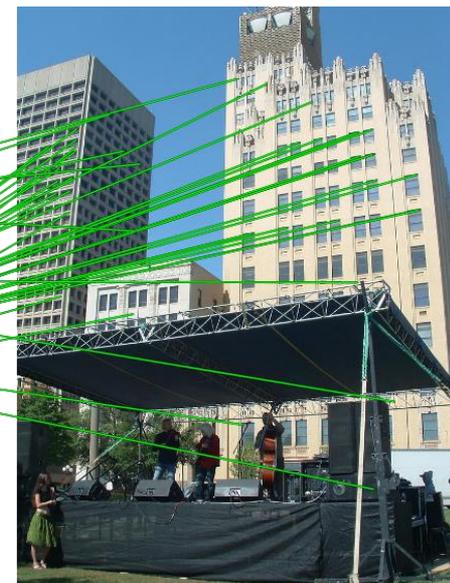- How to score the retrieval results?
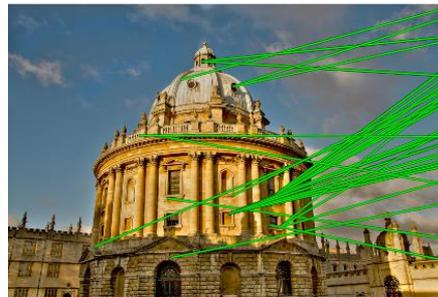
Kristen Grauman

# Spatial Verification

Query



DB image with high BoW similarity

Query



DB image with high BoW similarity

Both image pairs have many visual words in common.
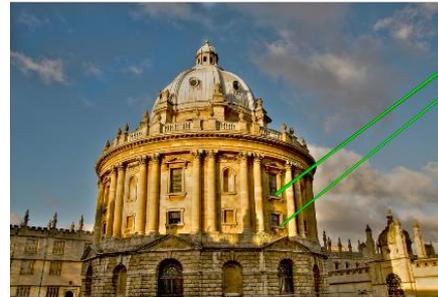
Slide credit: Ondrej Chum

# Spatial Verification



Query

DB image with high BoW similarity

Query

DB image with high BoW similarity

Only some of the matches are mutually consistent
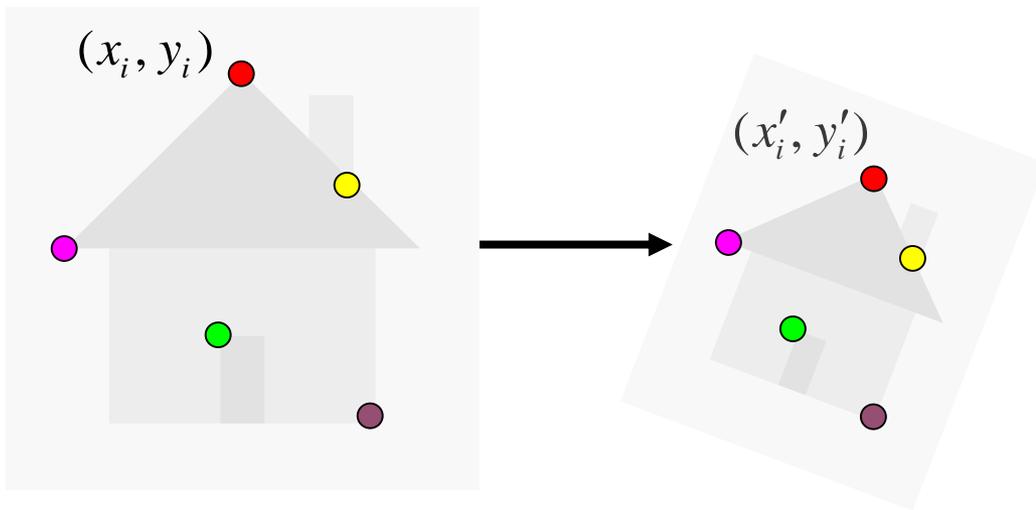
Slide credit: Ondrej Chum

# Spatial Verification: two basic strategies

- RANSAC
  - Typically sort by BoW similarity as initial filter
  - Verify by checking support (inliers) for possible transformations
    - e.g., "success" if find a transformation with > N inlier correspondences

- Generalized Hough Transform
  - Let each matched feature cast a vote on location, scale, orientation of the model object
  - Verify parameters with enough votes

Kristen Grauman

# RANSAC verification

# Recall: Fitting an affine transformation
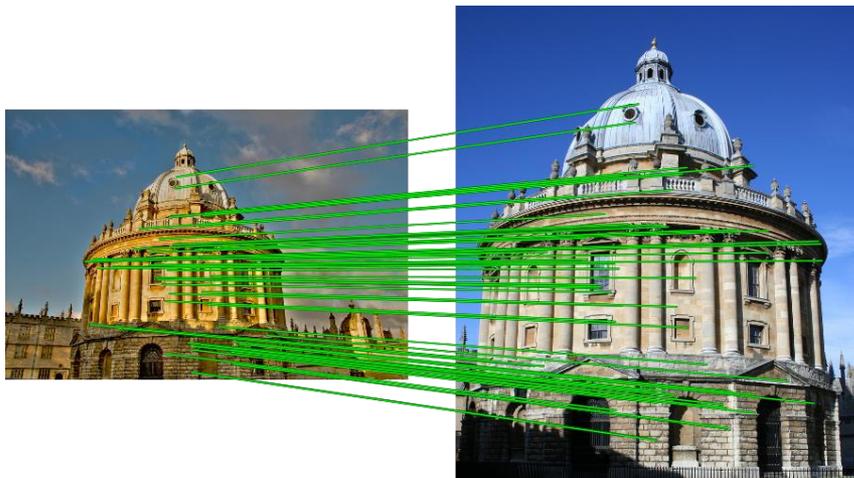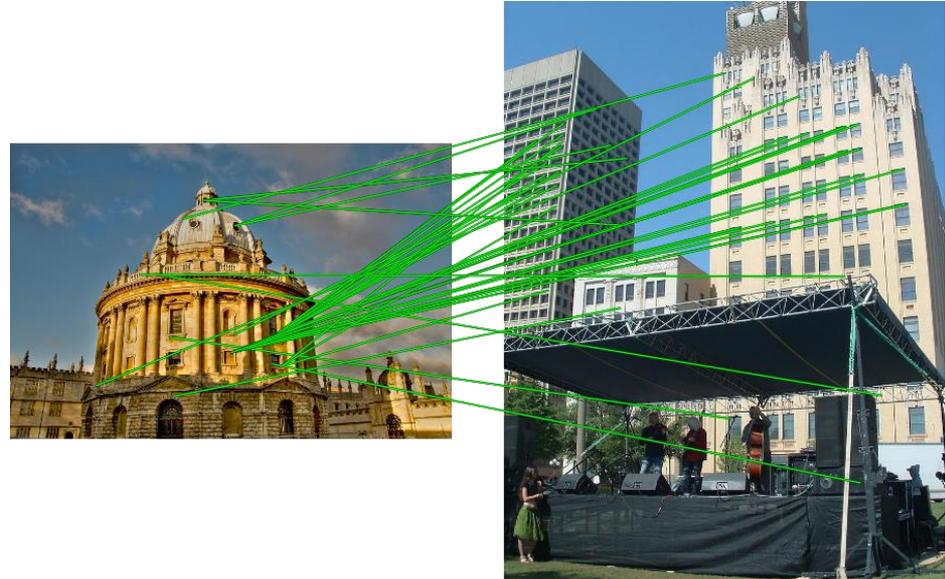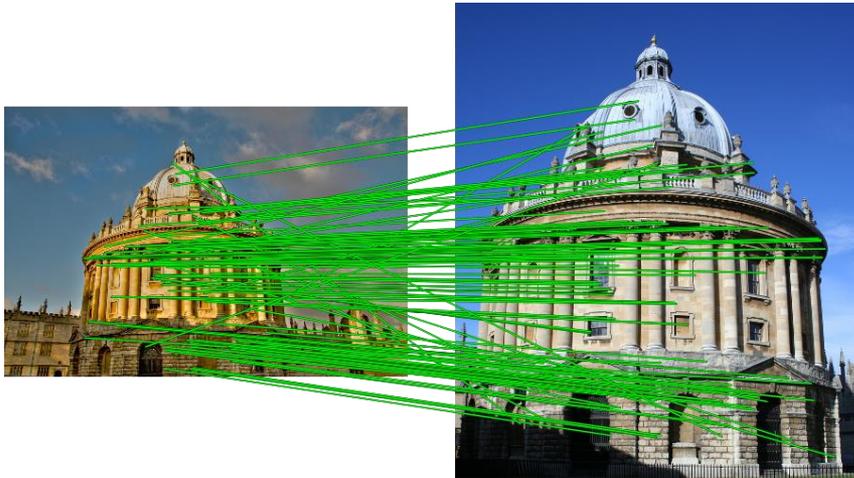
$(x_i, y_i)$

$(x_i', y_i')$

Approximates viewpoint changes for roughly planar objects and roughly orthographic cameras.

$$\begin{bmatrix} x_i' \\ y_i' \end{bmatrix} = \begin{bmatrix} m_1 & m_2 \\ m_3 & m_4 \end{bmatrix} \begin{bmatrix} x_i \\ y_i \end{bmatrix} + \begin{bmatrix} t_1 \\ t_2 \end{bmatrix}$$

$$\begin{bmatrix} & & \cdots & & & \\ x_i & y_i & 0 & 0 & 1 & 0 \\ 0 & 0 & x_i & y_i & 0 & 1 \\ & & \cdots & & & \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \\ m_3 \\ m_4 \\ t_1 \\ t_2 \end{bmatrix} = \begin{bmatrix} \cdots \\ x_i' \\ y_i' \\ \cdots \end{bmatrix}$$

12

# RANSAC verification

# Video Google System

1. Collect all words within query region
2. Inverted file index to find relevant frames
3. Compare word counts
4. Spatial verification

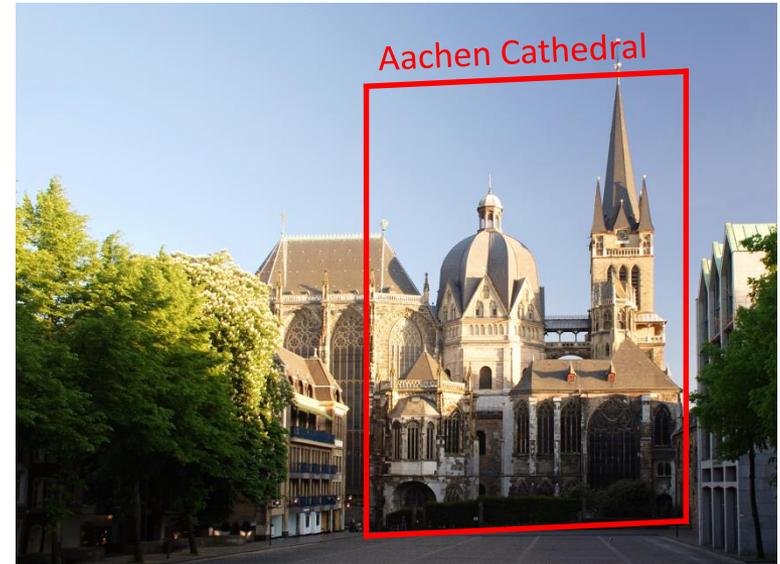Sivic & Zisserman, ICCV 2003

- Demo online at :
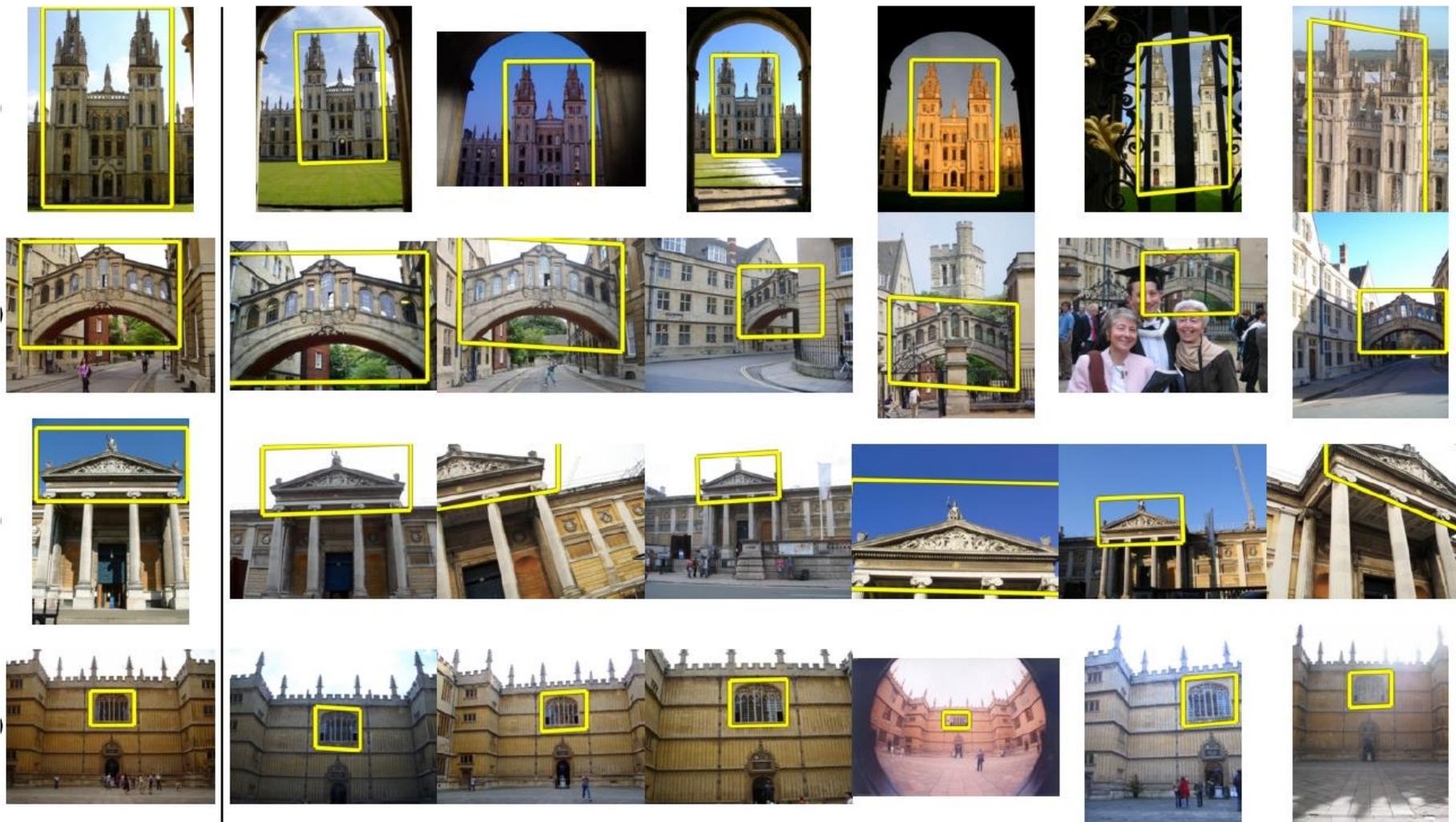http://www.robots.ox.ac.uk/~vgg/research/vgoogle/index.html



Query region

Retrieved frames

Kristen Grauman

# Example Applications



Aachen Cathedral

**Mobile tourist guide**
- **Self-localization**
- **Object/building recognition**
- **Photo/video augmentation**

[Quack, Leibe, Van Gool, CIVR'08]

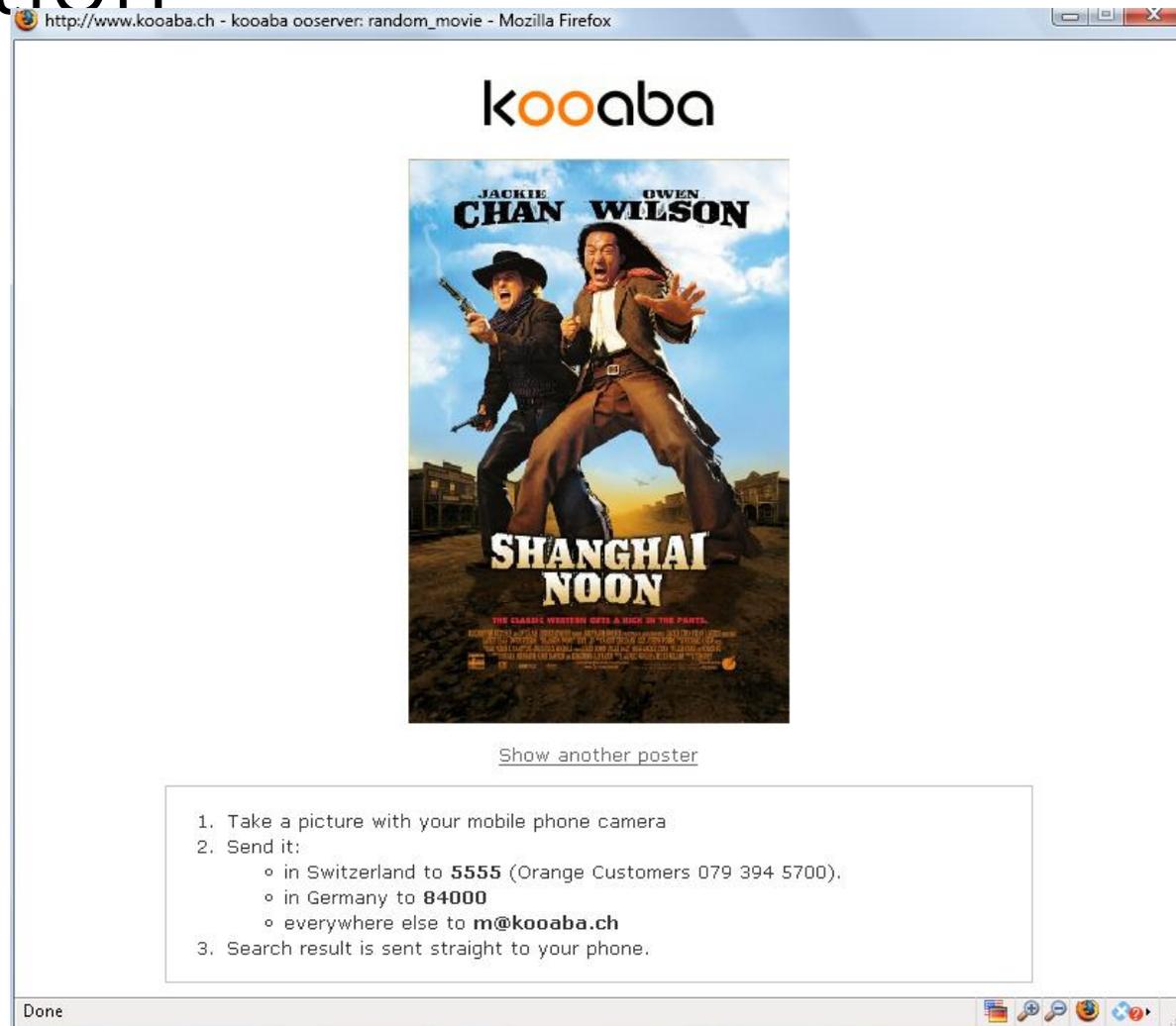# Application: Large-Scale Retrieval



Query        Results from 5k Flickr images (demo available for 100k set)
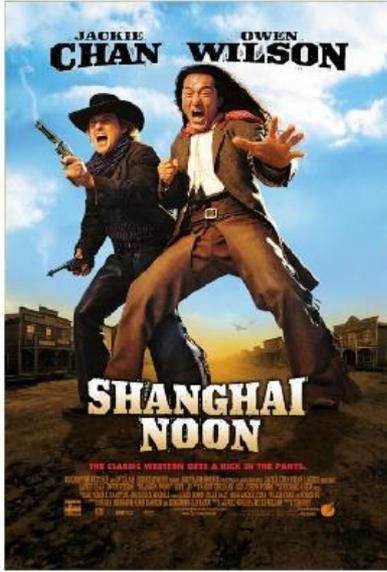
[Philbin CVPR'07]

# Web Demo: Movie Poster Recognition

50'000 movie posters indexed

Query-by-image from mobile phone available in Switzer-land



http://www.kooaba.com/en/products_engine.html#

# Google Goggles

Use pictures to search the web. ▷ Watch a video

**Get Google Goggles**

**Android (1.6+ required)**
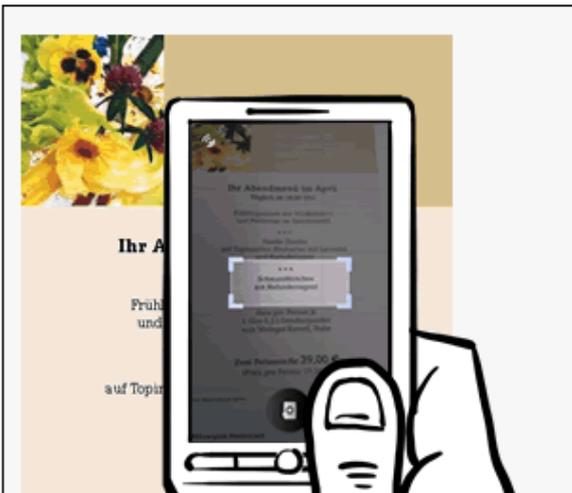Download from Android Market.

**Send Goggles to Android phone**

New! **iPhone (iOS 4.0 required)**
Download from the App Store.

**Send Goggles to iPhone**

New!

Menu
Crêpes-8
œufs-7

Text

Landmarks

GLOBAL HISTORY A SUMMARY

Books

Schlomo Inc
555-1212
Sinc@gmail.com

Contact Info

Artwork

Wine

Logos

Google goggles labs 11:03

Lammkoteletts vom Biobauern mit Schalotten, Tomatencoulis und Basilikum-Gnocchi

German (auto) » English

Lamb chops from the farmers with the shallots, tomato sauce and basil gnocchi
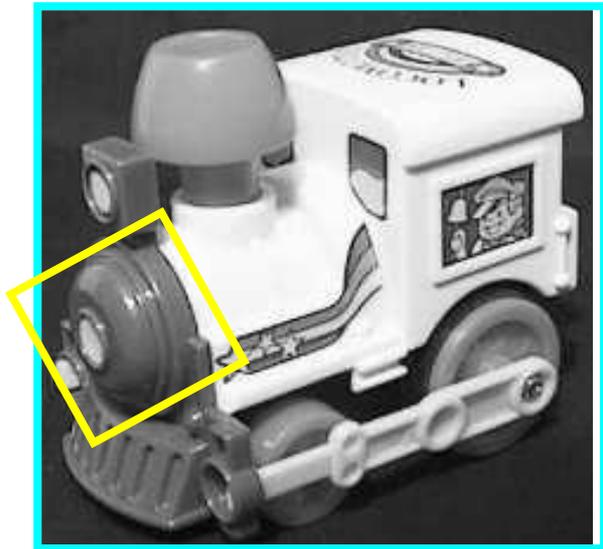
18

# Spatial Verification: two basic strategies

- RANSAC
    - Typically sort by BoW similarity as initial filter
    - Verify by checking support (inliers) for possible transformations
        - e.g., "success" if find a transformation with > N inlier correspondences

- Generalized Hough Transform
    - Let each matched feature cast a vote on location, scale, orientation of the model object
    - Verify parameters with enough votes

Kristen Grauman

# Voting: Generalized Hough Transform

- If we use scale, rotation, and translation invariant local features, then each feature match gives an alignment hypothesis (for scale, translation, and orientation of model in image).
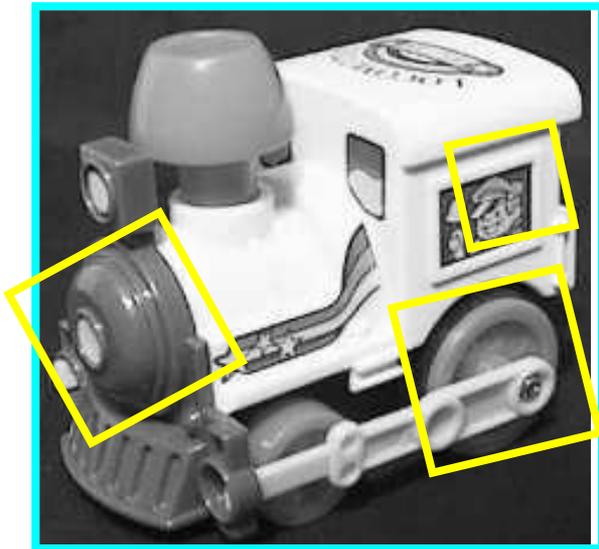


Model



Novel image

20

Adapted from Lana Lazebnik

# Voting: Generalized Hough Transform

- A hypothesis generated by a single match may be unreliable,
- So let each match **vote** for a hypothesis in Hough space



Model

Novel image

# Gen Hough Transform details (Lowe's system)

- **Training phase:** For each model feature, record 2D location, scale, and orientation of model (relative to normalized feature frame)

- **Test phase:** Let each match btwn a test SIFT feature and a model feature vote in a 4D Hough space
  - Use broad bin sizes of 30 degrees for orientation, a factor of 2 for scale, and 0.25 times image size for location
  - Vote for two closest bins in each dimension

- Find all bins with at least three votes and perform geometric verification
  - Estimate least squares *affine* transformation
  - Search for additional features that agree with the alignment

David G. Lowe. **"Distinctive image features from scale-invariant keypoints."** *IJCV* 60 (2), pp. 91-110, 2004.

# Example result



Background subtract for model boundaries

Objects recognized,

Recognition in spite of occlusion

[Lowe]

# Recall: difficulties of voting

- Noise/clutter can lead to as many votes as true target

- Bin size for the accumulator array must be chosen carefully

- In practice, good idea to make broad bins and spread votes to nearby bins, since verification stage can prune bad vote peaks.

# Gen Hough vs RANSAC

**GHT**

- Single correspondence -> vote for all consistent parameters
- Represents uncertainty in the model parameter space
- Linear complexity in number of correspondences and number of voting cells; beyond 4D vote space impractical
- Can handle high outlier ratio

**RANSAC**

- Minimal subset of correspondences to estimate model -> count inliers
- Represents uncertainty in image space
- Must search all data points to check for inliers each iteration
- Scales better to high-d parameter spaces

Kristen Grauman

# What else can we borrow from text retrieval?

Index

China is forecasting a trade surplus of $90bn (£51bn) to $100bn this year, a threefold increase on 2004's $32bn. The Commerce Ministry said the surplus would ... dicted 30% jump in expor... h a 18% rise in imp... ly to further ... hat China's ... deliber... the sur... one fac... Xiaochua... more to bo... stayed withi... value of the yua... % in July and permitted it ... w band, but the US wants the yuan to be ... d to trade freely. However, Beijing has made ... that it will take its time and tread careful... allowing the yuan to rise further in value.

**China, trade, surplus, commerce, exports, imports, US, yuan, bank, domestic, foreign, increase, trade, value**

# *tf-idf* weighting

- **T**erm **f**requency – **i**nverse **d**ocument **f**requency
- Describe frame by frequency of each word within it, downweight words that appear often in the database
- (Standard weighting for text retrieval)

Number of occurrences of word i in document d

Total number of documents in database

$$t_i = \frac{n_{id}}{n_d} \log \frac{N}{n_i}$$

Number of words in document d

Number of documents word i occurs in, in whole database

Kristen Grauman

# Query Expansion



Results

Query image

Spatial verification

New query

New results

Chum, Philbin, Sivic, Isard, Zisserman: Total Recall..., ICCV 2007

# Recognition via alignment

**Pros**:
- Effective when we are able to find reliable features within clutter
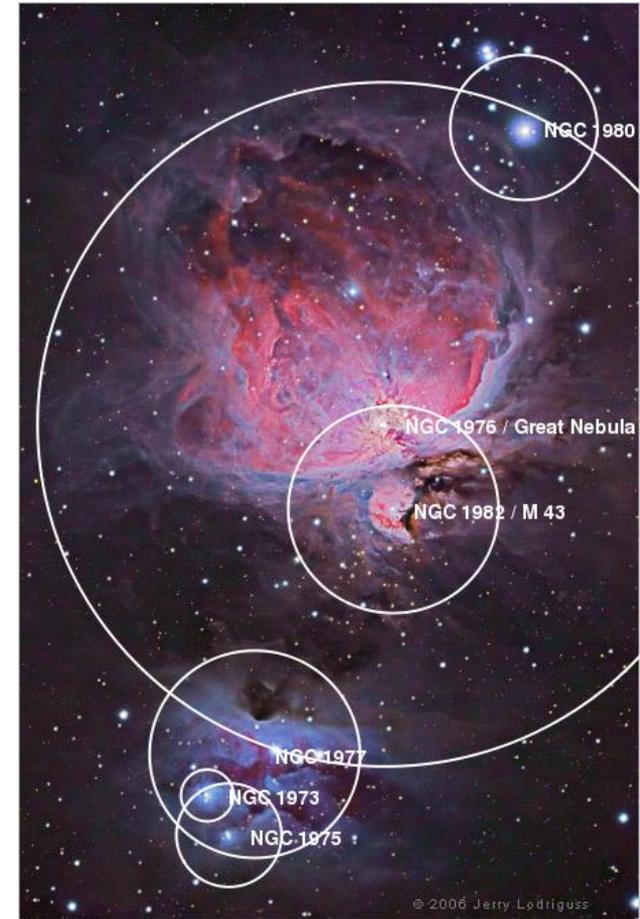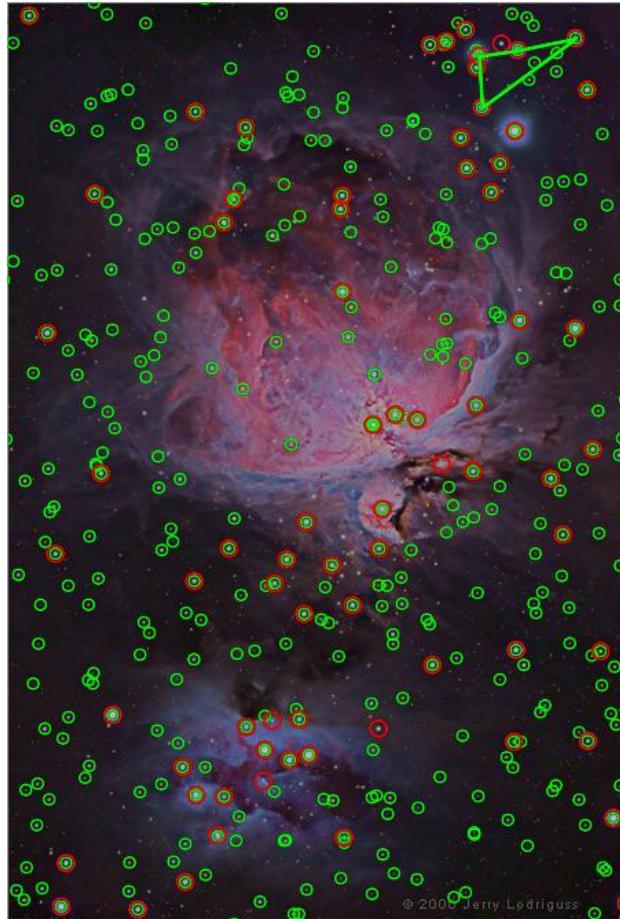- Great results for matching specific instances

**Cons**:
- Scaling with number of models
- Spatial verification as post-processing – not seamless, expensive for large-scale problems
- Not suited for category recognition.

Kristen Grauman

# Making the Sky Searchable:
# Fast Geometric Hashing for Automated Astrometry

Sam Roweis, Dustin Lang & Keir Mierle
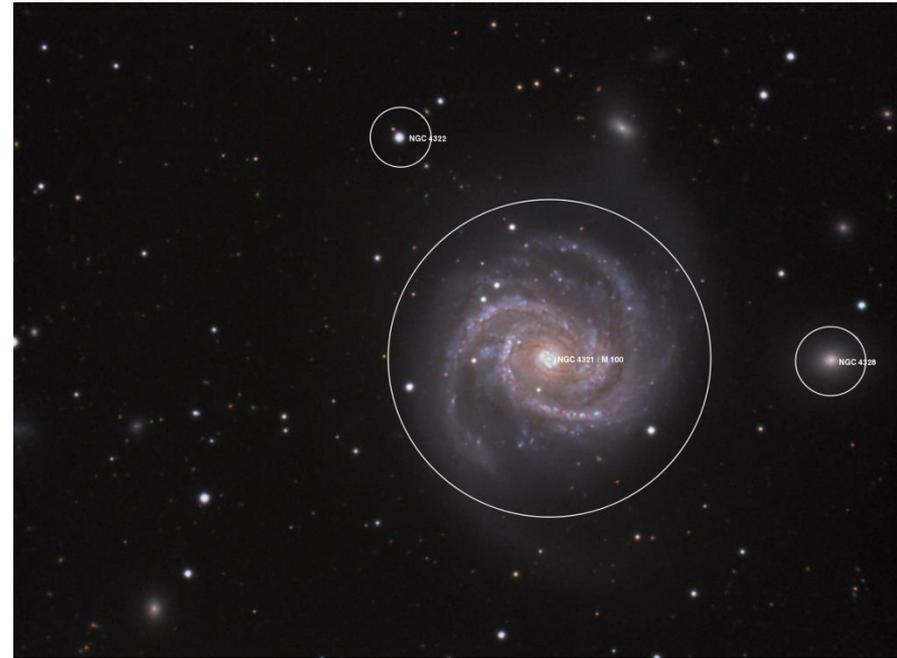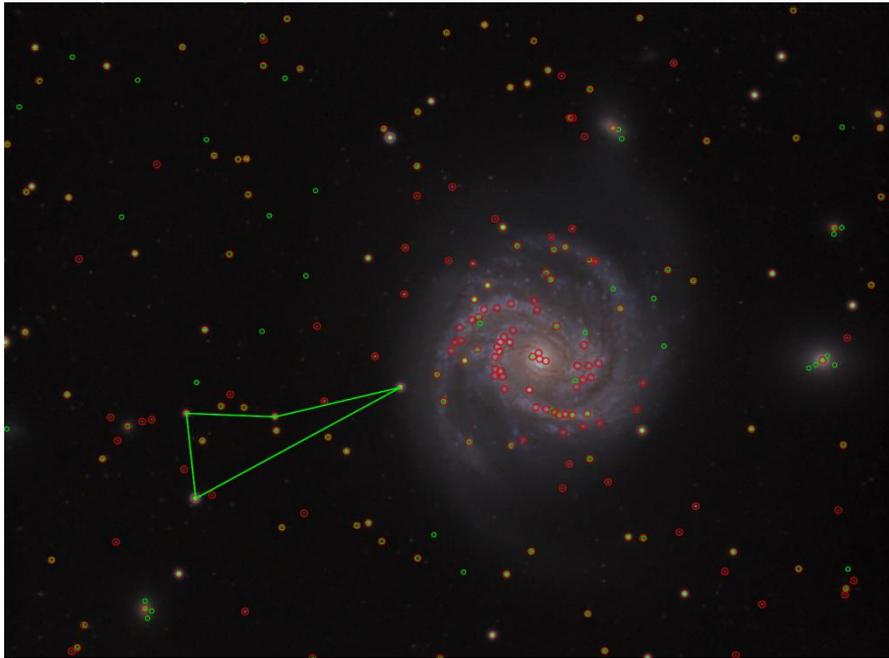University of Toronto

David Hogg & Michael Blanton
New York University

# Example



A shot of the Great Nebula, by Jerry Lodriguss (c.2006), from astropix.com
http://astrometry.net/gallery.html

32

# Example





An amateur shot of M100, by Filippo Ciferri (c.2007) from flickr.com
http://astrometry.net/gallery.html

33

# Example



A beautiful image of Bode's nebula (c.2007) by Peter Bresseler, from starlightfriend.de
http://astrometry.net/gallery.html

# Things to remember

- **Matching local invariant features**

  - Useful not only to provide matches for multi-view geometry, but also to find objects and scenes.

- **Bag of words** representation: quantize feature space to make discrete set of visual words

  - Summarize image by distribution of words
  - Index individual words

- **Inverted index**: pre-compute index to enable faster search at query time

- **Recognition of instances via alignment:** matching local features followed by spatial verification

  - Robust fitting : RANSAC, GHT

# Discussion – Think-pair-share

- Find a person you don't know

- Discuss
  - strength,
  - weakness, and
  - potential extension

- Share with class
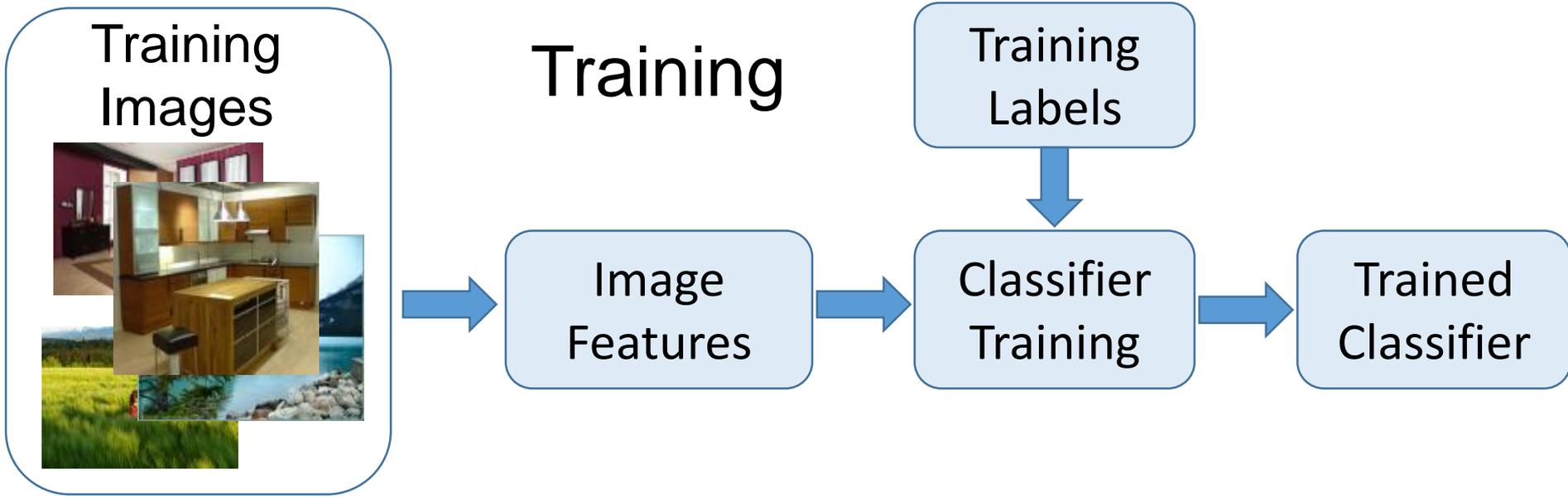
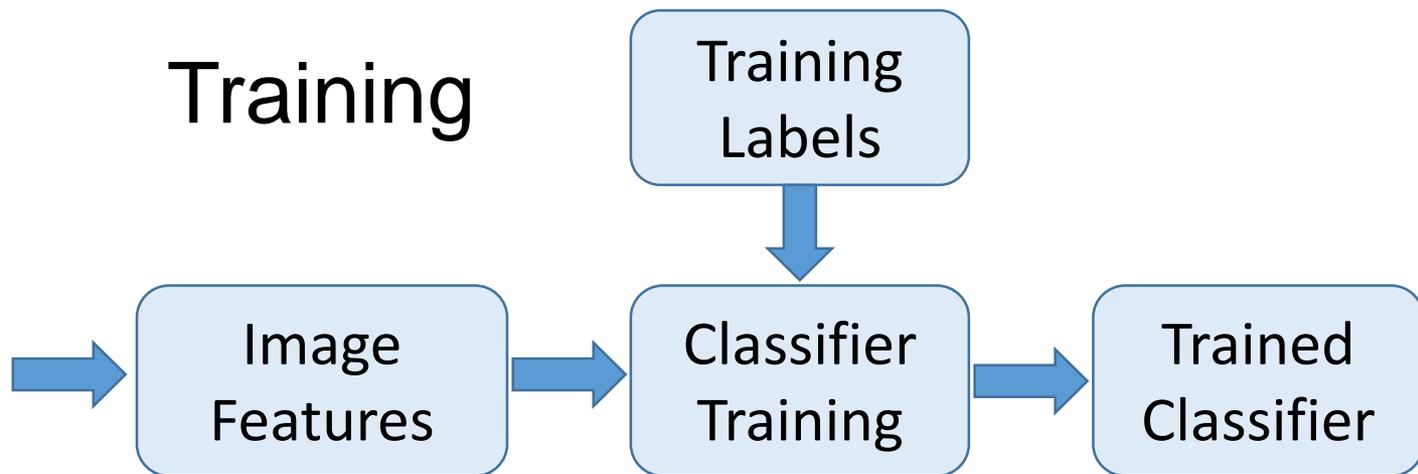# Image Categorization: Training phase



Training Images

Training

Training Labels

Image Features

Classifier Training

Trained Classifier

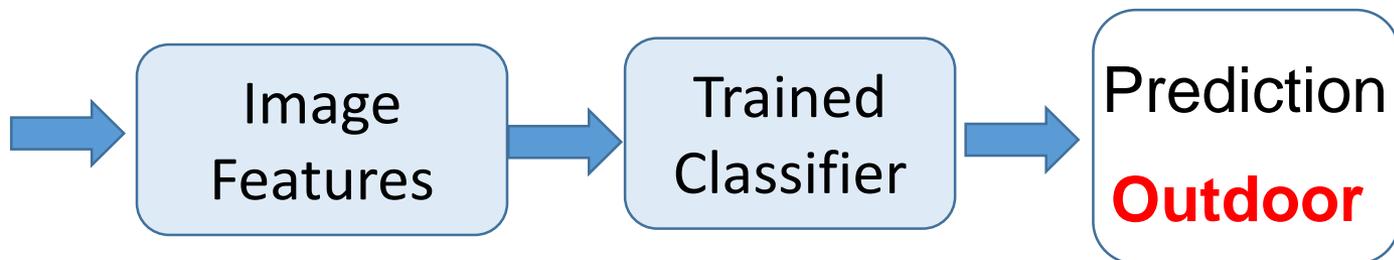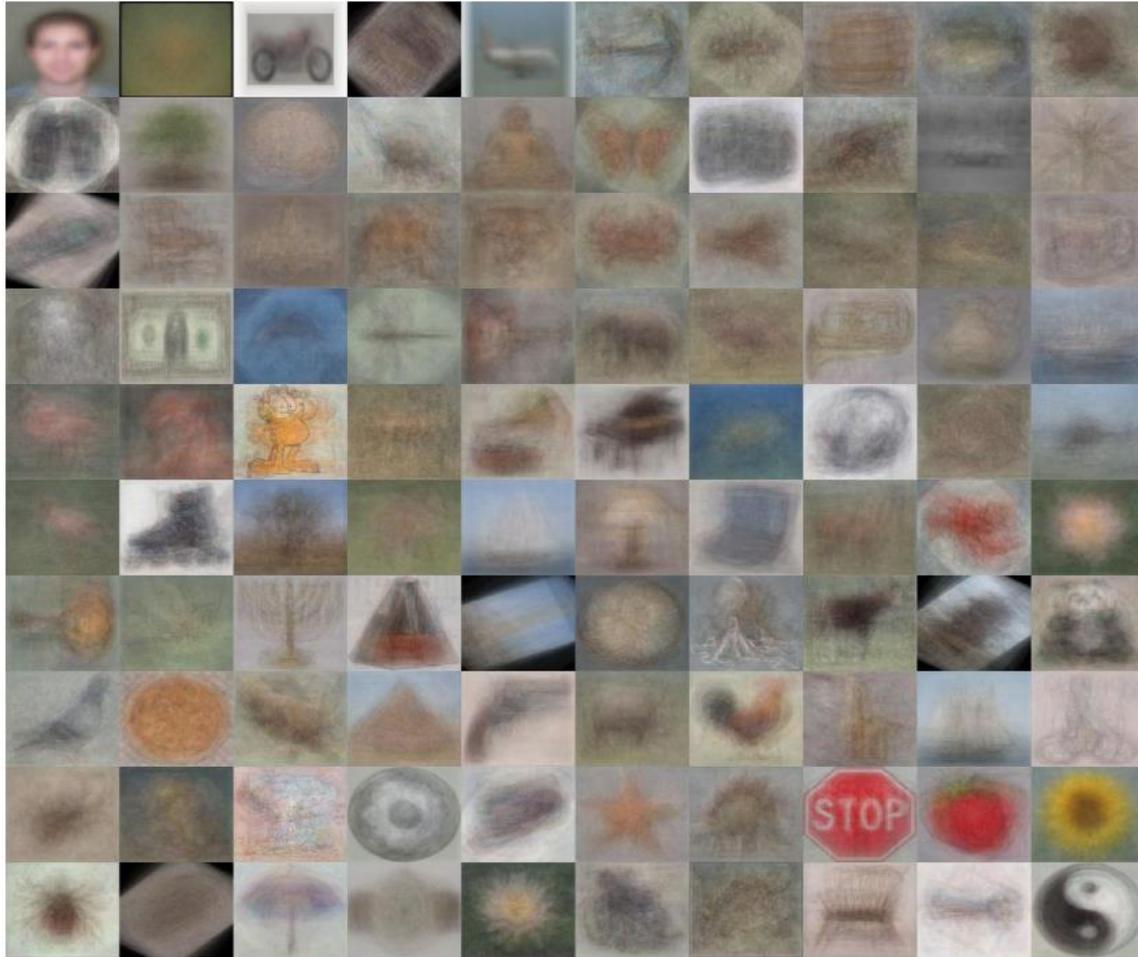# Image Categorization: Testing phase

# Image categorization

- Cat vs Dog

# Image categorization
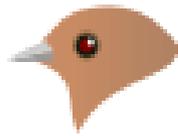
- Object recognition



Caltech 101 Average Object Images

# Image categorization

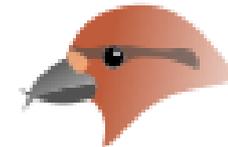- Fine-grained recognition



Generalist    Insect catching    Grain eating    Coniferous-seed eating    Nectar feeding

Chiseling    Dip netting    Surface skimming    Scything    Probing

Aerial fishing    Pursuit fishing    Scavenging    Raptorial    Filter feeding

[Visipedia Project](Visipedia Project)

# Image categorization

- Place recognition



spare bedroom · teenage bedroom · romantic bedroom

wooded kitchen · messy kitchen · stylish kitchen

darkest forest path · wintering forest path · greener forest path

rocky coast · misty coast · sunny coast

Places Database [Zhou et al. NIPS 2014]

# Image categorization

- Visual font recognition



Top Ranked Fonts (Space Coast):
- Adobe Caslon Pro Bold
- Rotation LT Std Bold
- Gazette LT Std Bold
- Baskerville Cyr LT Std Bold
- Joanna MT Std Bold

Top Ranked Fonts (Saturday):
- Hypatia Sans Pro Black
- Gill Sans Std Bold
- Montara Bold Gothic
- IT Ckabel Std Bold
- Myriad Arabic Black

[Chen et al. CVPR 2014]

# Image categorization

- Dating historical photos



| 1940 | 1953 | 1966 | 1977 |

[Palermo et al. ECCV 2012]

# Image categorization

- Image style recognition



HDR    Macro    Baroque    Roccoco

Vintage    Noir    Northern Renaissance    Cubism

Minimal    Hazy    Impressionism    Post-Impressionism

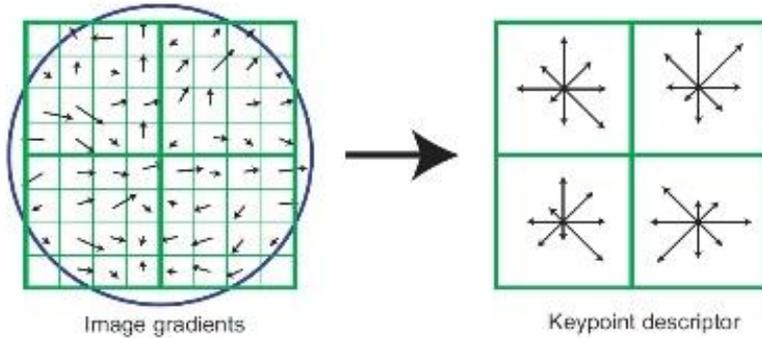Long Exposure    Romantic    Abs. Expressionism    Color Field Painting
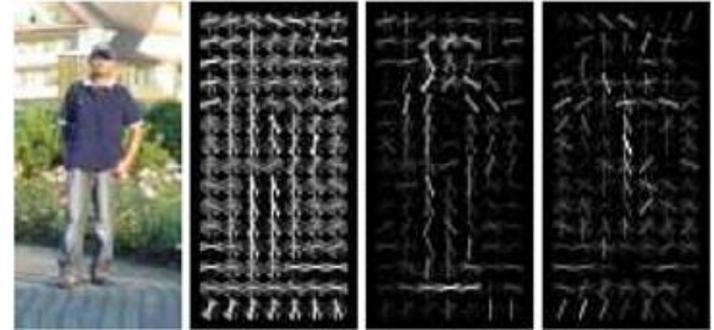
Flickr Style: 80K images covering 20 styles.    Wikipaintings: 85K images for 25 art genres.
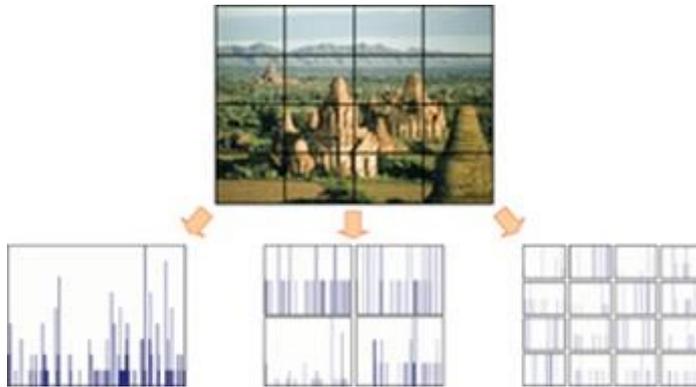
[Karayev et al. BMVC 2014]

# Features are the Keys
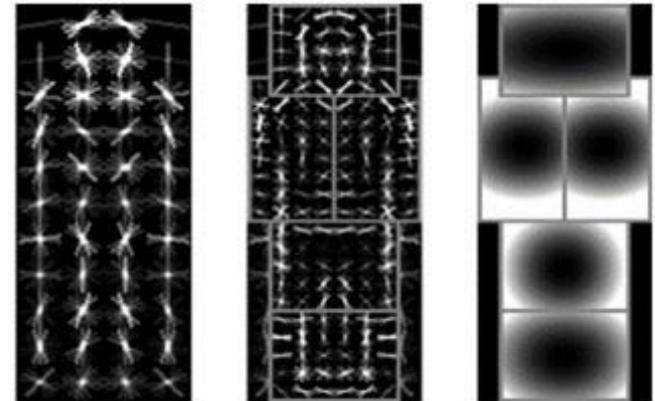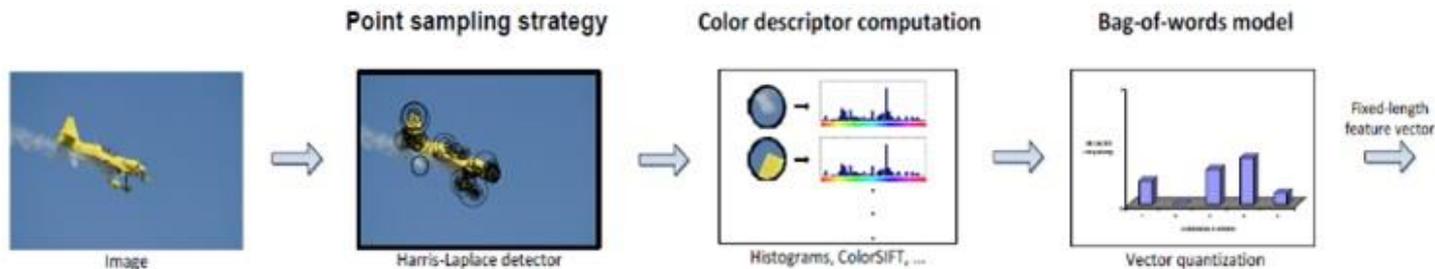


SIFT [Loewe IJCV 04]



HOG [Dalal and Triggs CVPR 05]



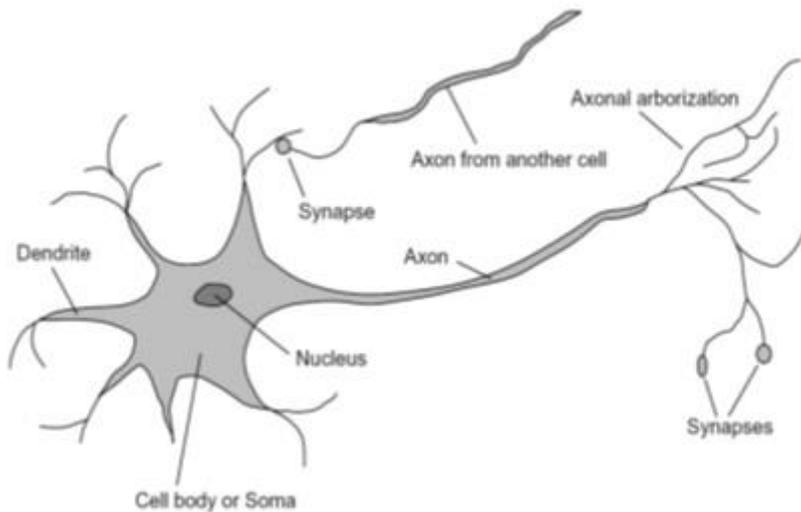SPM [Lazebnik et al. CVPR 06]



DPM [Felzenszwalb et al. PAMI 10]



Color Descriptor [Van De Sande et al. PAMI 10]

# Learning a Hierarchy of Feature Extractors

- Each layer of hierarchy extracts features from output of previous layer

- All the way from pixels → classifier

- Layers have the (nearly) same structure

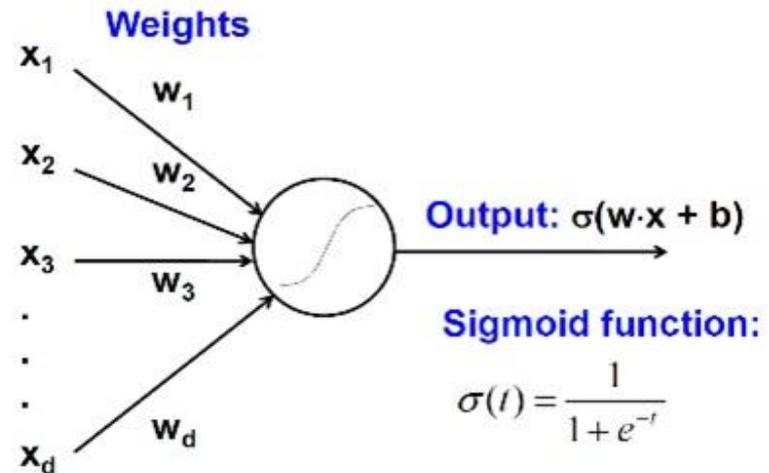Image/video → **Layer 1** → **Layer 2** → **Layer 3** → Labels
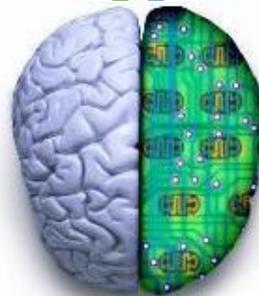
# Biological neuron and Perceptrons



A biological neuron



An artificial neuron (Perceptron)
- a linear classifier

# Simple, Complex and Hypercomplex cells
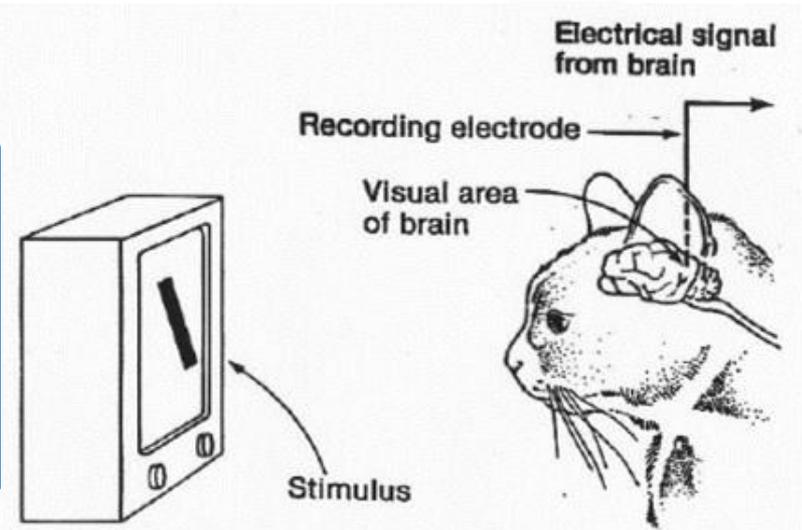




David H. Hubel and Torsten Wiesel



Electrical signal from brain

Recording electrode

Visual area of brain
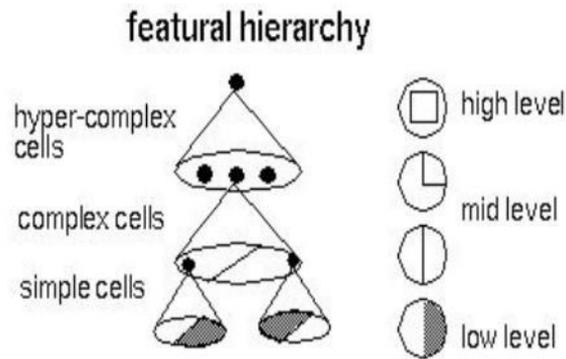
Stimulus

Suggested a **hierarchy** of **feature detectors** in the visual cortex, with higher level features responding to patterns of activation in lower level cells, and propagating activation upwards to still higher level cells.
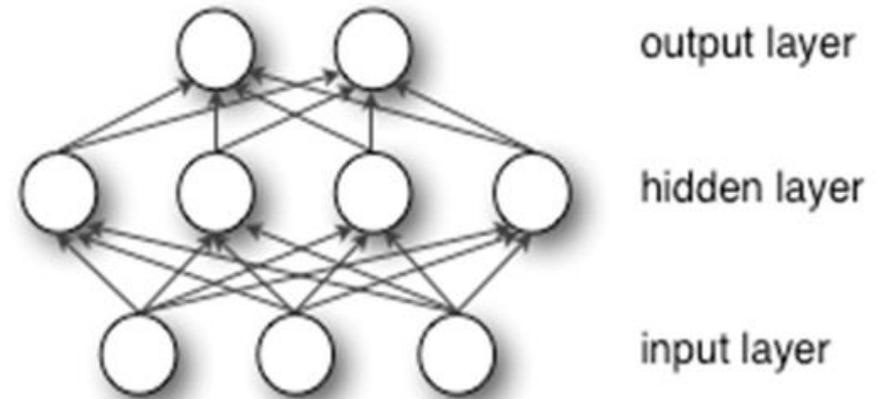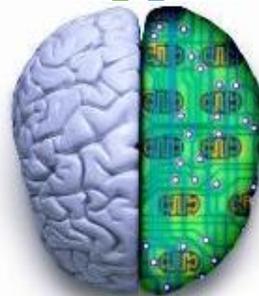
**David Hubel's** Eye, Brain, and Vision

# Hubel/Wiesel Architecture and Multi-layer Neural Network



Hubel and Weisel's architecture



Multi-layer Neural Network
- A *non-linear* classifier

# Multi-layer Neural Network

- A non-linear classifier

- **Training:** find network weights **w** to minimize the error between true training labels $y_i$ and estimated labels $f_w(x_i)$

$$E(\mathbf{w}) = \sum_{i=1}^{N} \left( y_i - f_\mathbf{w}(\mathbf{x}_i) \right)^2$$

- Minimization can be done by gradient descent provided $f$ is differentiable

- This training method is called **back-propagation**

output layer

hidden layer

input layer

# Convolutional Neural Networks

- Also known as CNN, ConvNet, DCN

- CNN = a multi-layer neural network with
  1. Local connectivity
  2. Weight sharing

# CNN: Local Connectivity



**Global** connectivity       **Local** connectivity

- \# input units (neurons): 7

- \# hidden units: 3

- Number of parameters
  - Global connectivity: 3 x 7 = 21
  - Local connectivity:   3 x 3 = 9

# CNN: Weight Sharing
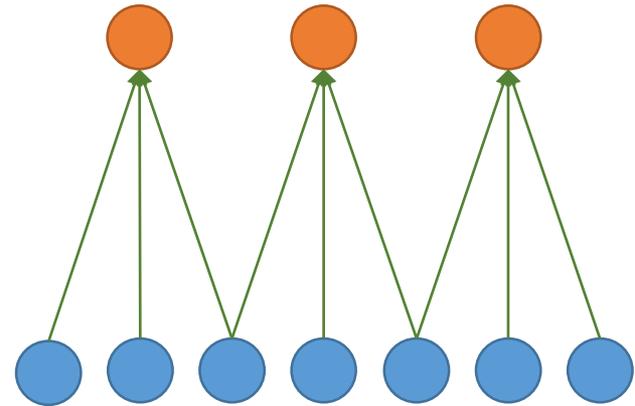


Hidden layer

Input layer

$w_1$  $w_3$  $w_5$  $w_7$  $w_9$
$w_2$  $w_4$  $w_6$  $w_8$

$w_1$  $w_3$  $w_2$  $w_1$  $w_3$
$w_2$  $w_1$  $w_3$  $w_2$

**Without** weight sharing

**With** weight sharing

- # input units (neurons): 7

- # hidden units: 3

- Number of parameters
  - Without weight sharing: 3 x 3 = 9
  - With weight sharing :    3 x 1 = 3

# CNN with multiple input channels



Hidden layer

Input layer    Channel 1

Channel 2

**Single** input channel

**Multiple** input channels

Filter weights

Filter weights

# CNN with multiple output maps

Hidden layer

Map 1

Map 2

Input layer

**Single** output map

Filter weights

**Multiple** output maps

Filter 1

Filter 2

Filter weights

# Putting them together

- Local connectivity
- Weight sharing
- Handling multiple input channels
- Handling multiple output maps

Weight sharing

Local connectivity

# input channels

# output (activation) maps

Image credit: A. Karpathy

# Neocognitron [Fukushima, Biological Cybernetics 1980]



Fig. 1. Correspondence between the hierarchy model by Hubel and Wiesel, and the neural network of the neocognitron



Deformation-Resistant Recognition

S-cells: (simple)
 - extract local features

C-cells: (complex)
 - allow for positional errors

# LeNet [LeCun et al. 1998]



C3: f. maps 16@10x10

S4: f. maps 16@5x5

INPUT 32x32

C1: feature maps 6@28x28

S2: f. maps 6@14x14

C5: layer 120

F6: layer 84

OUTPUT 10

Convolutions    Subsampling    Convolutions    Subsampling    Full connection    Gaussian connections    Full connection



LeNet-1 from 1993

Gradient-based learning applied to document recognition [LeCun, Bottou, Bengio, Haffner 1998]

# What is a Convolution?

- Weighted moving sum



Input

Feature Activation Map

slide credit: S. Lazebnik

# Convolutional Neural Networks

Feature maps

↑

Normalization

↑

Spatial pooling

↑

Non-linearity

↑

Convolution
(Learned)

↑

Input Image

# Convolutional Neural Networks



Input

Feature Map

# Convolutional Neural Networks

Feature maps

↑

Normalization

↑

Spatial pooling

↑

**Non-linearity**

↑

Convolution
(Learned)

↑

Input Image

## Rectified Linear Unit (ReLU)

# Convolutional Neural Networks

Feature maps

Normalization

**Spatial pooling**

Non-linearity

Convolution
(Learned)

Input Image

Max pooling

Max-pooling: a non-linear down-sampling

Provide *translation invariance*

# Convolutional Neural Networks

Feature maps

Normalization

Spatial pooling

Non-linearity

Convolution
(Learned)

Input Image



Feature Maps



Feature Maps
After Contrast
Normalization

# Convolutional Neural Networks



slide credit: S. Lazebnik

# Engineered vs. learned features

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

**Label**

| Dense |
| Dense |
| Dense |
| Convolution/pool |
| Convolution/pool |
| Convolution/pool |
| Convolution/pool |
| Convolution/pool |
| Convolution/pool |
| **Image** |

**Label**

| Classifier |
| Pooling |
| Feature extraction |
| **Image** |

INPUT 32x32

C1: feature maps 6@28x28

C3: f. maps 16@10x10

S2: f. maps 6@14x14

S4: f. maps 16@5x5

C5: layer 120

F6: layer 84

OUTPUT 10

Convolutions

Subsampling

Convolutions

Subsampling

Full connection

Full connection

Gaussian connections

**Gradient-Based Learning Applied to Document Recognition**, LeCun, Bottou, Bengio and Haffner, Proc. of the IEEE, **1998**



**Imagenet Classification with Deep Convolutional Neural Networks**, Krizhevsky, Sutskever, and Hinton, NIPS **2012**

Slide Credit: L. Zitnick

INPUT
32x32

C1: feature maps
6@28x28

C3: f. maps 16@10x10

S2: f. maps
6@14x14

S4: f. maps 16@5x5

C5: layer
120

F6: layer
84

OUTPUT
10

Convolutions   Subsampling   Convolutions   Subsampling   Full connection   Full connection   Gaussian connections

**Gradient-Based Learning Applied to Document Recognition**, LeCun, Bottou, Bengio and Haffner, Proc. of the IEEE, 1998

224   5   5   13   3   3   13   dense   dense
11   55   27   3   3   3   13   1000
11   128   3   192   128   2048   2048
224   48
3

**Imagenet Class... Networks**, Kriz...

GPUs + Data*

* Rectified activations and dropout

Slide Credit: L. Zitnick

# SIFT Descriptor

Image Pixels → Apply gradient filters

Spatial pool (Sum)

Normalize to unit length

→ Feature Vector

# SIFT Descriptor

Image Pixels ⟹

Apply oriented filters

Spatial pool (Sum)

Normalize to unit length

⟹ Feature Vector

# Spatial Pyramid Matching



SIFT Features

Filter with Visual Words

Max

Multi-scale spatial pool (Sum)

Classifier

Lazebnik, Schmid, Ponce [CVPR 2006]

# Deformable Part Model



## DeepPyramid DPM

(1) Color image pyramid

3

level $L$

(2) Truncated SuperVision CNN

For each pyramid level $l$

(output layer is conv5)

3

image pyramid level 1

(3) Conv5 feature pyramid

256

level $L$

(4) DPM-CNN

For each pyramid level $l$

256

conv5 pyramid level 1

(1/16th spatial resolution of the image)

(5) DPM score pyramid

level $L$

DPM score pyramid level 1

Deformable Part Models are Convolutional Neural Networks [Girshick et al. CVPR 15]

# AlexNet

- Similar framework to LeCun'98 but:
  - Bigger model (7 hidden layers, 650,000 units, 60,000,000 params)
  - More data ($10^6$ vs. $10^3$ images)
  - GPU implementation (50x speedup over CPU)
    - Trained on two GPUs for a week



A. Krizhevsky, I. Sutskever, and G. Hinton,
ImageNet Classification with Deep Convolutional Neural Networks, NIPS 2012

# Using CNN for Image Classification



Fully connected layer Fc7
d = 4096

AlexNet

Fixed input size:
224x224x3

Averaging

d = 4096

Softmax Layer

"Jia-Bin"

# Progress on ImageNet



ImageNet Image Classification Top5 Error

| | | | | | |
|---|---|---|---|---|---|
| 16.4 | 11.7 | 7.3 | 6.7 | 3.57 | 3.08 |
| 2012 AlexNet | 2013 ZF | 2014 VGG | 2014 GoogLeNet | 2015 ResNet | 2016 GoogLeNet-v4 |



I WAS WINNING IMAGENET

UNTIL A DEEPER MODEL CAME ALONG

# VGG-Net

- The deeper, the better

- Key design choices:
  - 3x3 conv. Kernels
    - very small
  - conv. stride 1
    - no loss of information

- Other details:
  - Rectification (ReLU) non-linearity
  - 5 max-pool layers (x2 reduction)
  - no normalization
  - 3 fully-connected (FC) layers

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096
FC-1000
softmax

# VGG-Net

- Why 3x3 layers?
  - Stacked conv. layers have a large receptive field
  - two 3x3 layers – 5x5 receptive field
  - three 3x3 layers – 7x7 receptive field

- More non-linearity
  - Less parameters to learn
  - ~140M per net

5

5

1st 3x3 conv. layer

2nd 3x3 conv. layer

# ResNet

- Can we just increase the #layer?



- How can we train very deep network?
  - Residual learning



| method | top-5 err. (test) |
|---|---|
| VGG [41] (ILSVRC'14) | 7.32 |
| GoogLeNet [44] (ILSVRC'14) | 6.66 |
| VGG [41] (v5) | 6.8 |
| PReLU-net [13] | 4.94 |
| BN-inception [16] | 4.82 |
| **ResNet (ILSVRC'15)** | **3.57** |

# DenseNet

- Shorter connections (like ResNet) help
- Why not just connect them all?

# Training Convolutional Neural Networks

- Backpropagation + stochastic gradient descent with momentum
  - [Neural Networks: Tricks of the Trade](#)
- Dropout
- Data augmentation
- Batch normalization
- Initialization
  - Transfer learning

# Training CNN with gradient descent

- A CNN as composition of functions

$$f_{\boldsymbol{w}}(\boldsymbol{x}) = f_L(\dots (f_2(f_1(\boldsymbol{x}; \boldsymbol{w_1}); \boldsymbol{w_2}) \dots ; \boldsymbol{w_L})$$

- Parameters

$$\boldsymbol{w} = (\boldsymbol{w_1}, \boldsymbol{w_2}, \dots \boldsymbol{w_L})$$

- Empirical loss function

$$L(\boldsymbol{w}) = \frac{1}{n} \sum_i l(z_i, f_{\boldsymbol{w}}(\boldsymbol{x_i}))$$

- Gradient descent

$$\boldsymbol{w^{t+1}} = \boldsymbol{w^t} - \eta_t \frac{\partial \boldsymbol{f}}{\partial \boldsymbol{w}}(\boldsymbol{w^t})$$

New weight

Old weight

Learning rate

Gradient

# An Illustrative example

$$f(x, y) = xy, \qquad \frac{\partial f}{\partial x} = y, \frac{\partial f}{\partial y} = x$$

Example: $x = 4, y = -3 \Rightarrow f(x, y) = -12$

Partial derivatives
$$\frac{\partial f}{\partial x} = -3, \qquad \frac{\partial f}{\partial y} = 4$$

Gradient
$$\nabla f = [\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}]$$

Example credit: Andrej Karpathy

$$f(x, y, z) = (x + y)z = qz$$

$$q = x + y$$
$$\frac{\partial q}{\partial x} = 1, \qquad \frac{\partial q}{\partial y} = 1$$

$$f = qz$$
$$\frac{\partial f}{\partial q} = z, \qquad \frac{\partial f}{\partial z} = q$$

Goal: compute the gradient
$$\nabla f = [\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}]$$

$$f(x, y, z) = (x + y)z = qz$$

$$q = x + y$$
$$\frac{\partial q}{\partial x} = 1, \qquad \frac{\partial q}{\partial y} = 1$$

$$f = qz$$
$$\frac{\partial f}{\partial q} = z, \qquad \frac{\partial f}{\partial z} = q$$

**Chain rule:**
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q}\frac{\partial q}{\partial x}$$

```
# set some inputs
x = -2; y = 5; z = -4

# perform the forward pass
q = x + y # q becomes 3
f = q * z # f becomes -12

# perform the backward pass (backpropagation) in reverse order:
# first backprop through f = q * z
dfdz = q # df/dz = q, so gradient on z becomes 3
dfdq = z # df/dq = z, so gradient on q becomes -4
# now backprop through q = x + y
dfdx = 1.0 * dfdq # dq/dx = 1. And the multiplication here is the chain rule!
dfdy = 1.0 * dfdq # dq/dy = 1
```
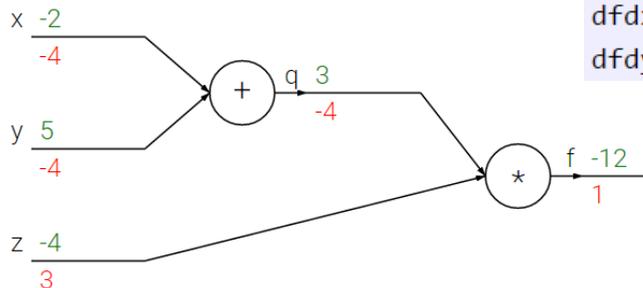
x  -2
   -4

y  5
   -4

q  3
   -4

z  -4
   3

+

*

f  -12
   1

Example credit: Andrej Karpathy

# Backpropagation (recursive chain rule)



$w_1$

$w_2$

$w_n$

$q$

$\frac{\partial f}{\partial q}$

$$\frac{\partial f}{\partial w_i} = \frac{\partial q}{\partial w_i} \frac{\partial f}{\partial q}$$

Local gradient

Gate gradient

Can be computed during forward pass

The gate receives this during backprop

# Dropout



(a) Standard Neural Net

(b) After applying dropout.

Intuition: successful conspiracies
- 50 people planning a conspiracy

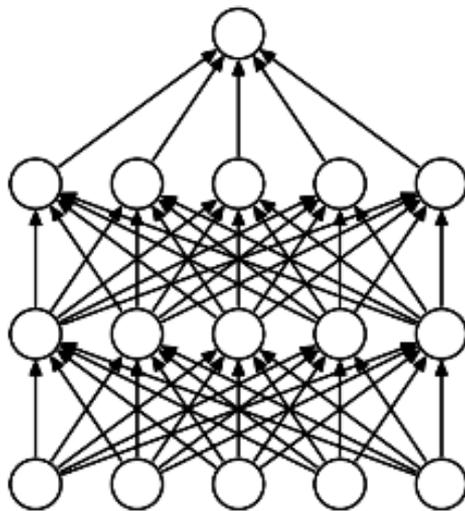- Strategy A: plan a big conspiracy involving 50 people
  - Likely to fail. 50 people need to play their parts correctly.

- Strategy B: plan 10 conspiracies each involving 5 people
  - Likely to succeed!

Dropout: A simple way to prevent neural networks from overfitting [Srivastava JMLR 2014]

# Dropout



(a) Standard Neural Net

(b) After applying dropout.



Without dropout

With dropout



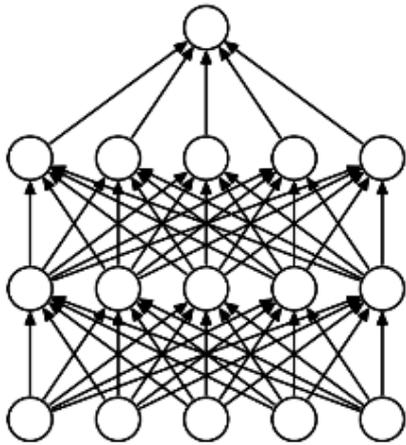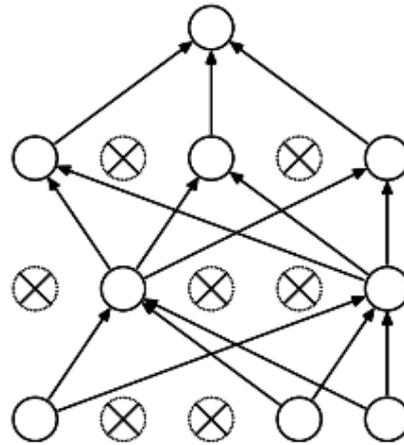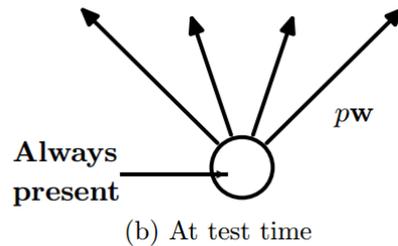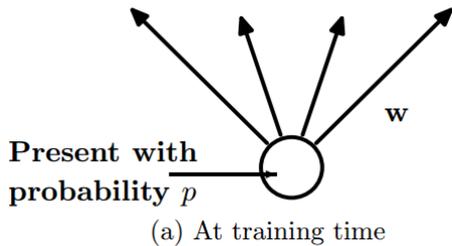**Present with probability $p$**   w

(a) At training time

**Always present**   $pw$

(b) At test time

**Main Idea**: approximately combining exponentially many different neural network architectures efficiently

| Model | Top-1 (val) | Top-5 (val) | Top-5 (test) |
|---|---|---|---|
| SVM on Fisher Vectors of Dense SIFT and Color Statistics | - | - | 27.3 |
| Avg of classifiers over FVs of SIFT, LBP, GIST and CSIFT | - | - | 26.2 |
| Conv Net + dropout (Krizhevsky et al., 2012) | 40.7 | 18.2 | - |
| Avg of 5 Conv Nets + dropout (Krizhevsky et al., 2012) | 38.1 | 16.4 | 16.4 |

Table 6: Results on the ILSVRC-2012 validation/test set.

Dropout: A simple way to prevent neural networks from overfitting [Srivastava JMLR 2014]

# Data Augmentation (Jittering)

- Create *virtual* training s
  - Horizontal flip
  - Random crop
  - Color casting
  - Geometric distortion

Deep Image [Wu et al. 2015]

# Parametric Rectified Linear Unit



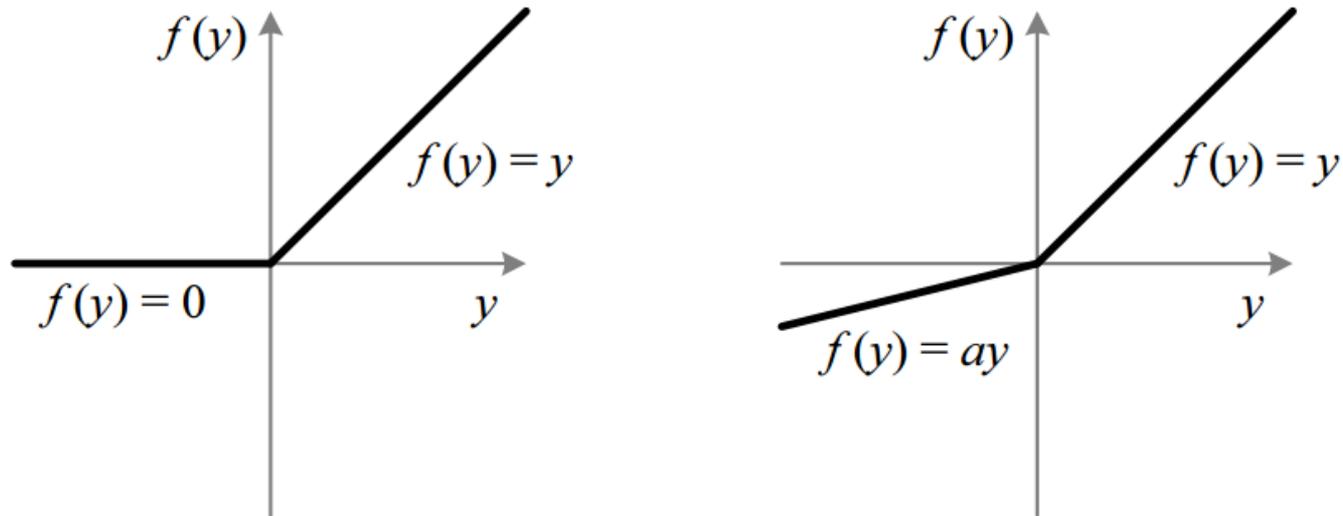| | team | top-5 (**test**) |
|---|---|---|
| in competition ILSVRC 14 | MSRA, SPP-nets [11] | 8.06 |
| | VGG [25] | 7.32 |
| | GoogLeNet [29] | 6.66 |
| post-competition | VGG [25] (arXiv v5) | 6.8 |
| | Baidu [32] | 5.98 |
| | **MSRA, PReLU-nets** | **4.94** |

Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification [He et al. 2015]

# Batch Normalization

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
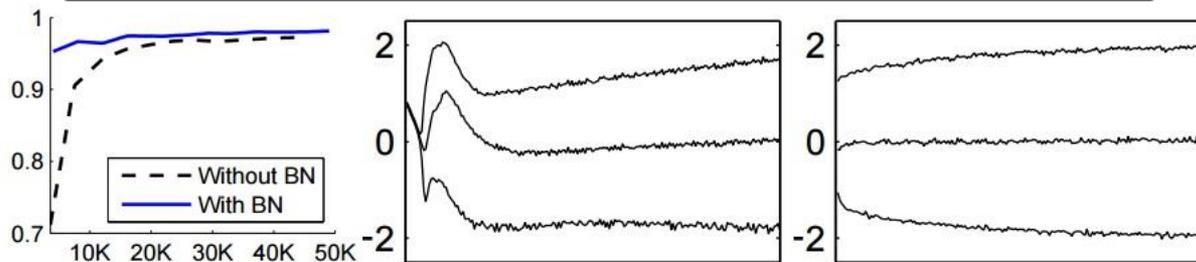Parameters to be learned: $\gamma, \beta$

**Output:** $\{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i \qquad \text{// mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{\mathcal{B}})^2 \qquad \text{// mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}$$

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}$$



(a)　(b) Without BN　(c) With BN

Legend: — — Without BN　—— With BN

Batch Normalization: Accelerating Deep Network Training by
Reducing Internal Covariate Shift [Ioffe and Szegedy 2015]

# Things to remember

- Visual categorization help transfer knowledge


- Convolutional neural networks
  - A cascade of conv + ReLU + pool
  - Representation learning
  - Advanced architectures
  - Tricks for training CNN