# Classical and Modern Recognition Techniques

# Today's outline

- We've covered Deep Convolutional Networks. But what did recognition techniques look like before AlexNet?
  - Bag of words models
  - Sliding window models
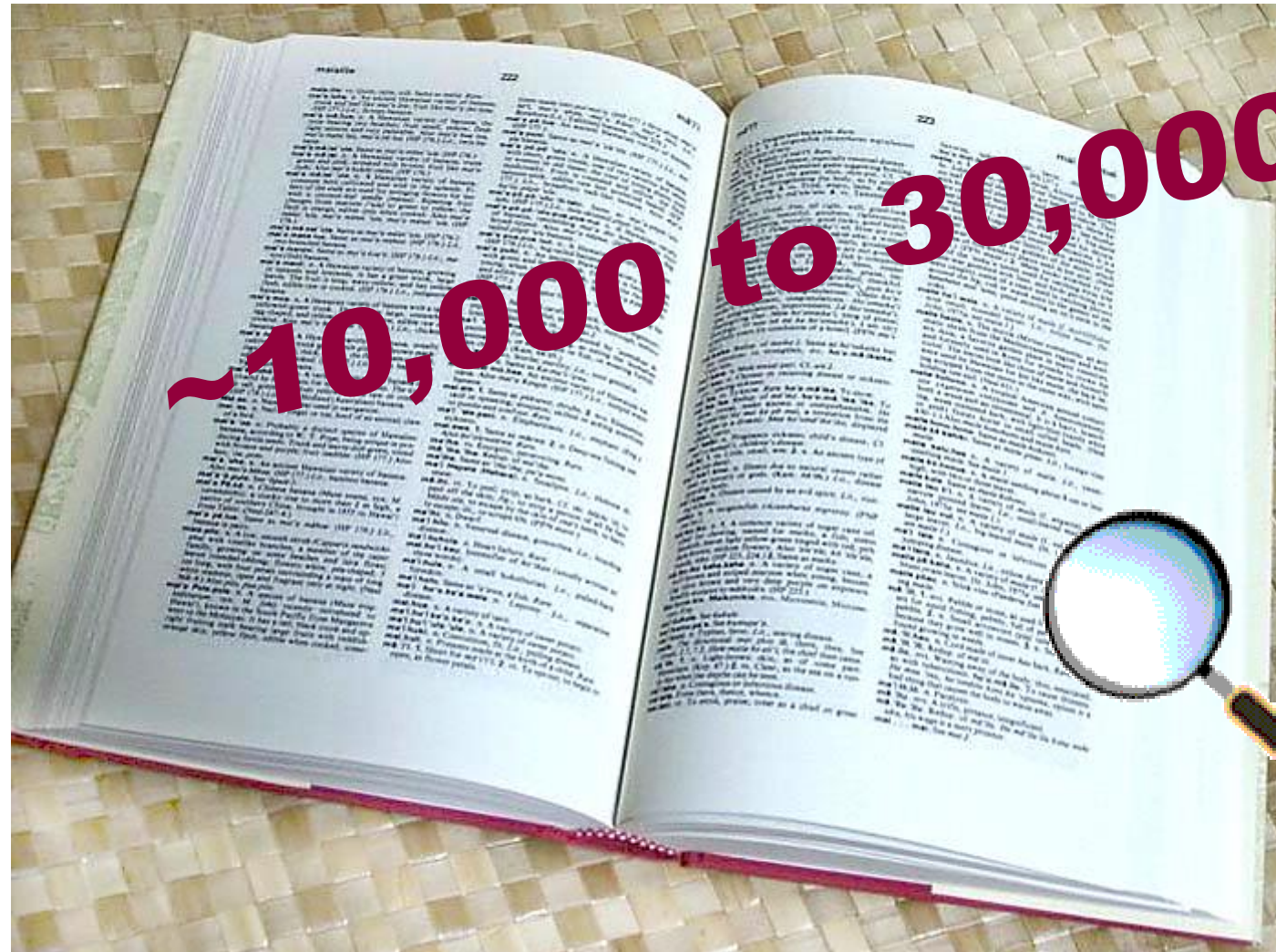- What do more recent deep learning architectures look like?

# Recognition: Overview and History



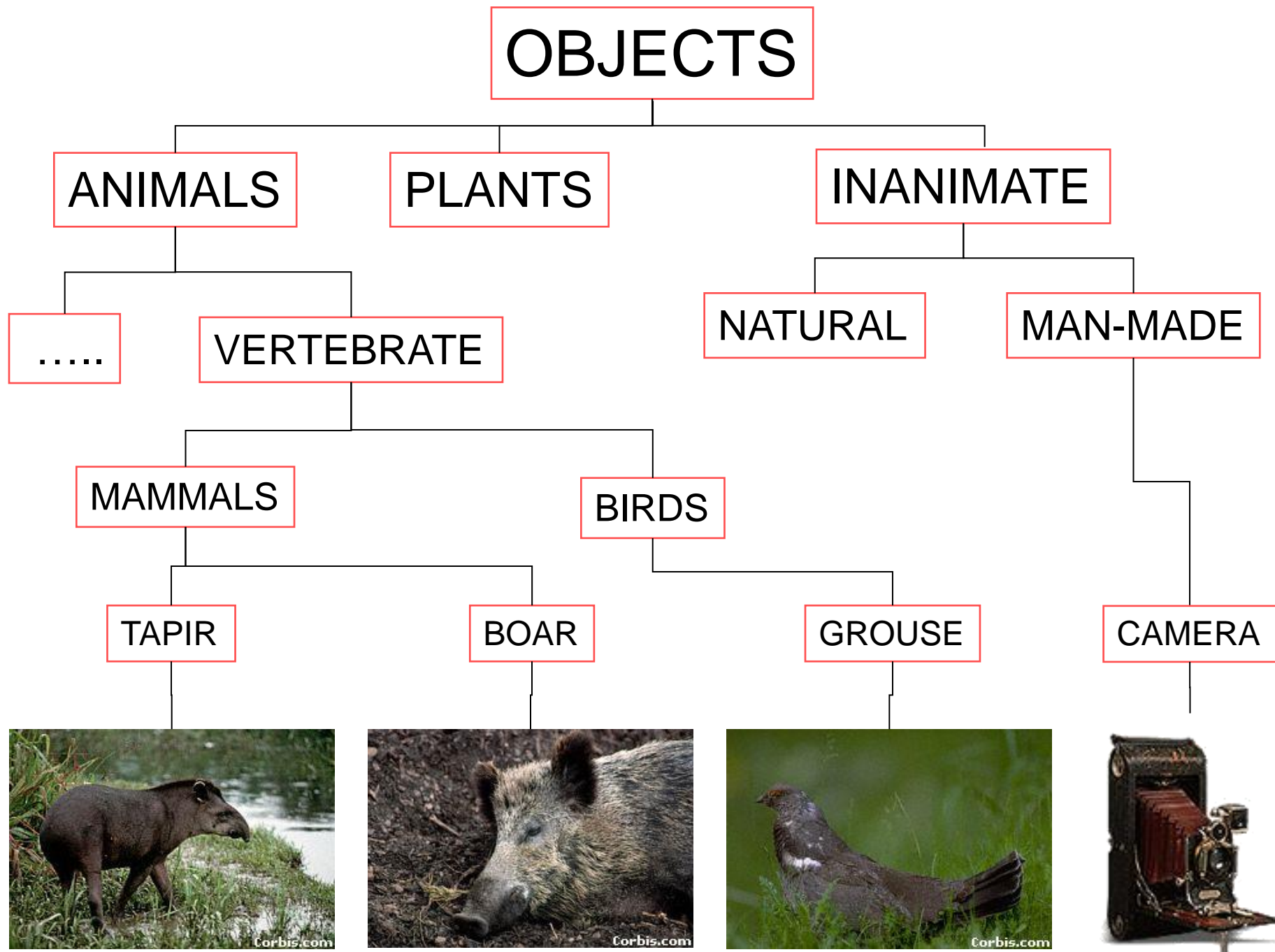Slides from Lana Lazebnik, Fei-Fei Li, Rob Fergus, Antonio Torralba, and Jean Ponce

# How many visual object categories are there?



~10,000 to 30,000

Biederman 1987

~10,000 to 30,000

# Specific recognition tasks



Svetlana Lazebnik

# Scene categorization or classification



- **outdoor/indoor**
- **city/forest/factory/etc.**

Svetlana Lazebnik

# Image annotation / tagging / attributes



- **street**
- **people**
- **building**
- **mountain**
- **tourism**
- **cloudy**
- **brick**
- **…**

Svetlana Lazebnik

# Object detection



**• find pedestrians**

Svetlana Lazebnik

# Image parsing / semantic segmentation



Svetlana Lazebnik

# Scene understanding?



Svetlana Lazebnik

# Recognition is all about modeling variability

Set of

Images

Variability: Camera position
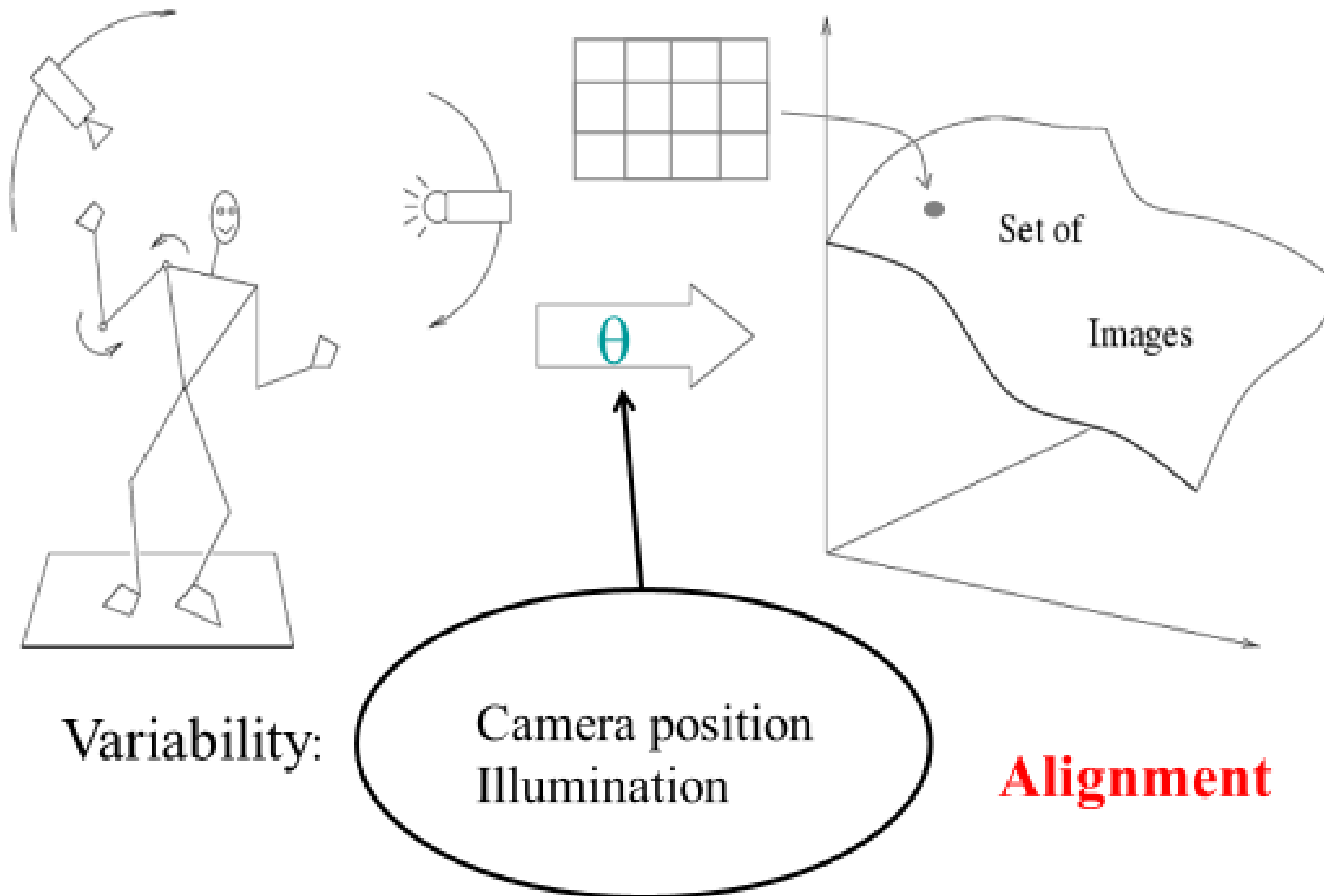
Illumination

Shape parameters

Within-class variations?

Svetlana Lazebnik

# Within-class variations



Svetlana Lazebnik

# History of ideas in recognition

- 1960s – early 1990s: the geometric era

Svetlana Lazebnik

Variability:

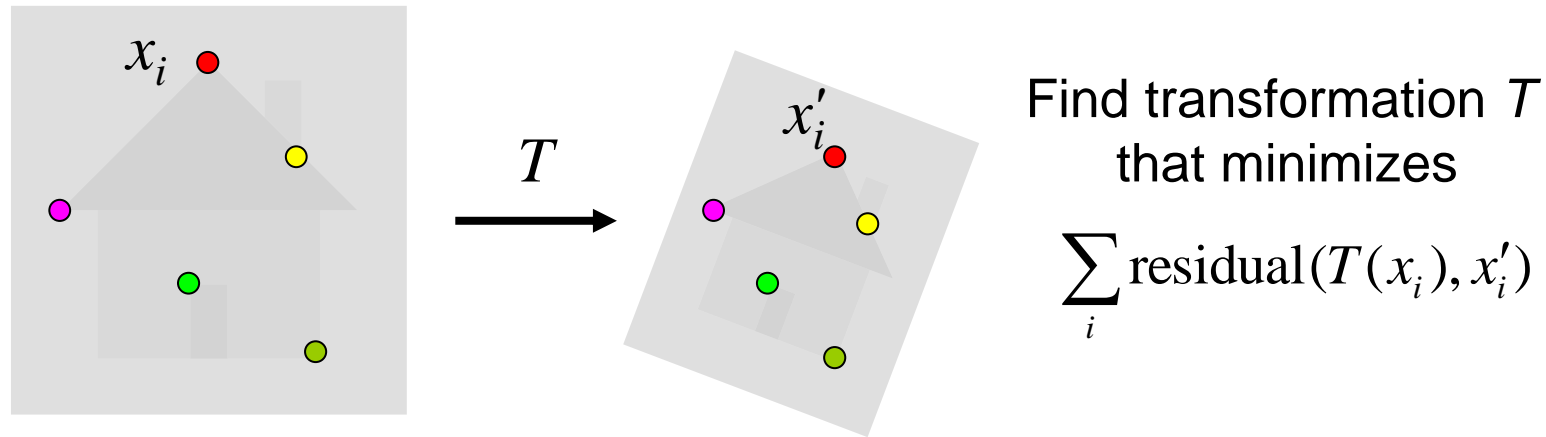Camera position
Illumination

**Alignment**

$\theta$

Set of

Images

Shape: assumed known

Roberts (1965); Lowe (1987); Faugeras & Hebert (1986); Grimson & Lozano-Perez (1986); Huttenlocher & Ullman (1987)

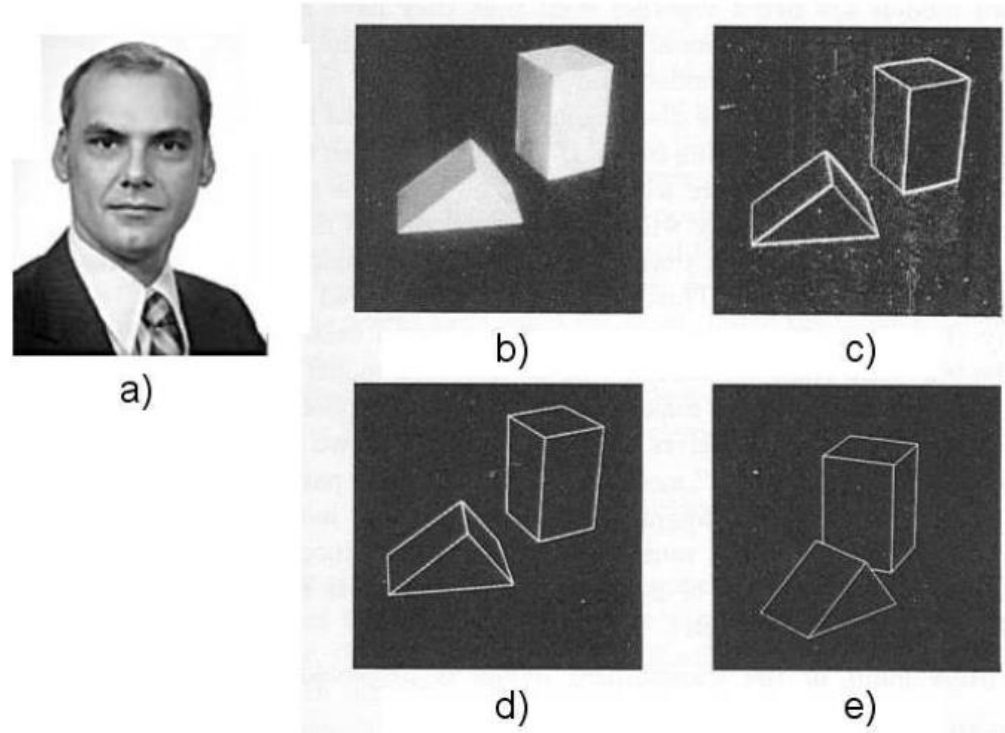Svetlana Lazebnik

# Recall: Alignment

- Alignment: fitting a model to a transformation between pairs of features (*matches*) in two images



$x_i$

$\xrightarrow{\;\;T\;\;}$

$x_i'$

Find transformation $T$ that minimizes
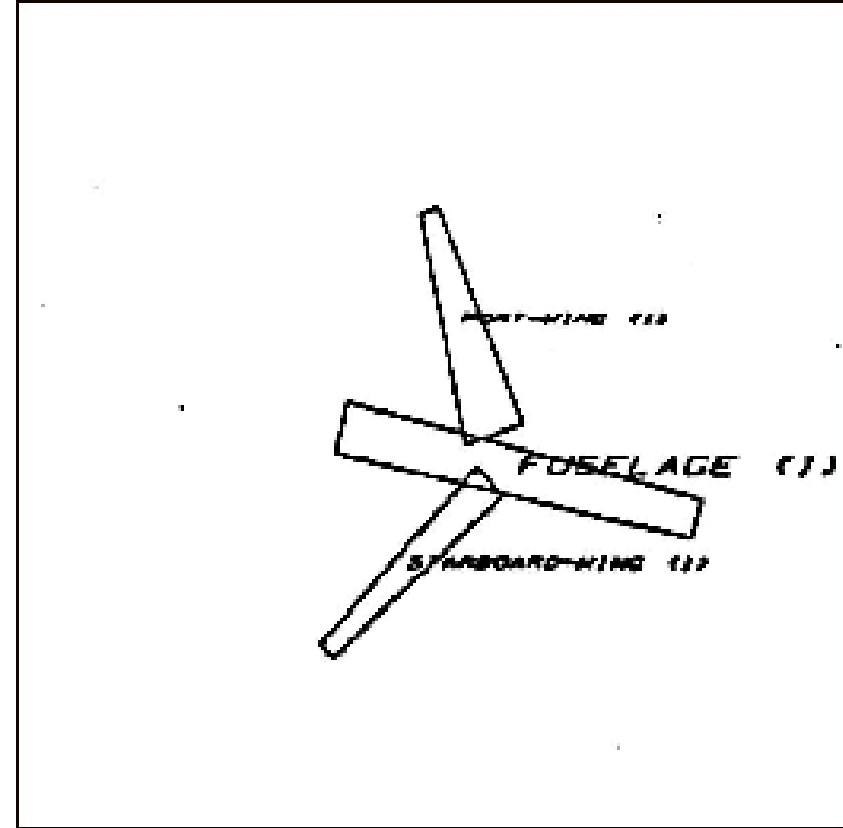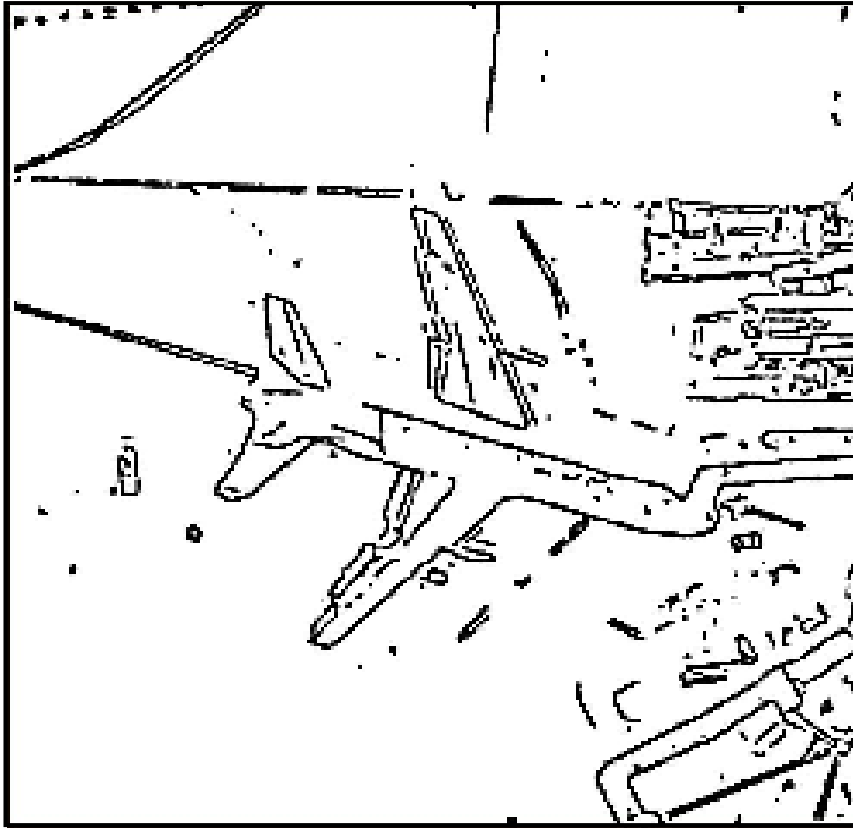
$$\sum_i \mathrm{residual}(T(x_i), x_i')$$

# Recognition as an alignment problem: Block world



L. G. Roberts, *Machine Perception of Three Dimensional Solids*, Ph.D. thesis, MIT Department of Electrical Engineering, 1963.

**Fig. 1.** A system for recognizing 3-d polyhedral scenes. a) L.G. Roberts. b)A blocks world scene. c)Detected edges using a 2x2 gradient operator. d) A 3-d polyhedral description of the scene, formed automatically from the single image. e) The 3-d scene displayed with a viewpoint different from the original image to demonstrate its accuracy and completeness. (b) - e) are taken from [64] with permission MIT Press.)

J. Mundy, Object Recognition in the Geometric Era: a Retrospective, 2006

# Representing and recognizing object categories is harder...


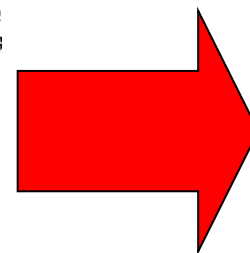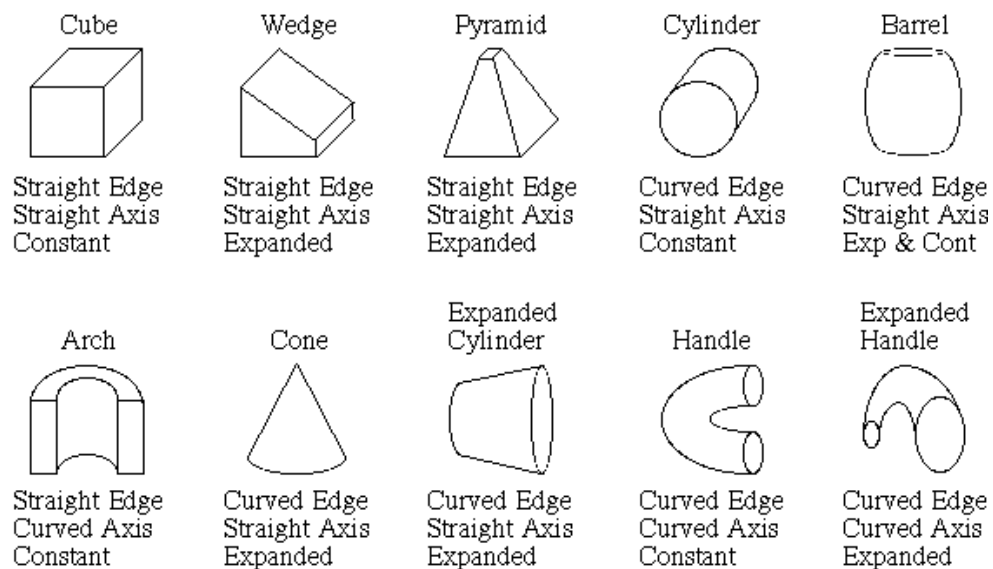
ACRONYM (Brooks and Binford, 1981)

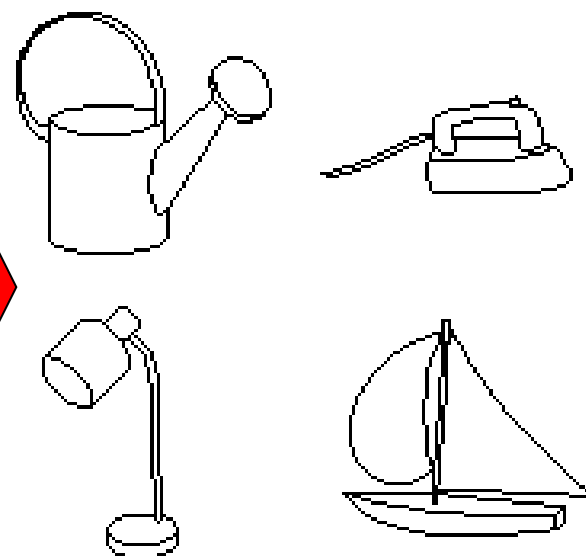Binford (1971), Nevatia & Binford (1972), Marr & Nishihara (1978)

# Recognition by components

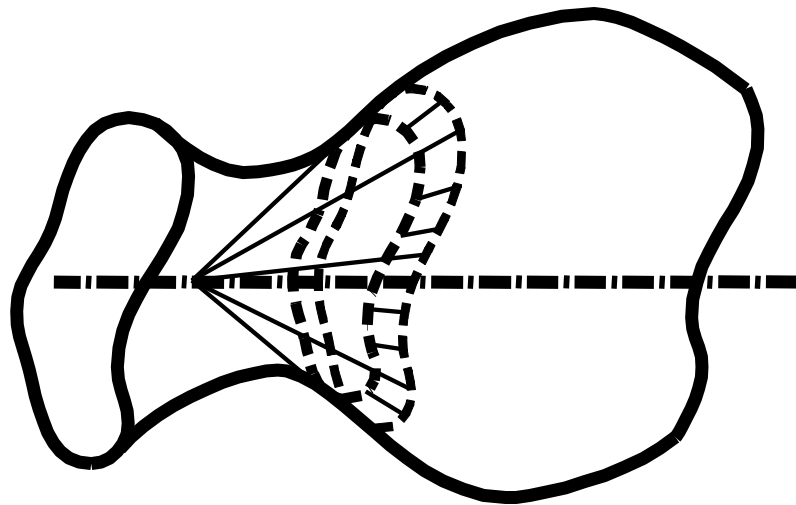Biederman (1987)

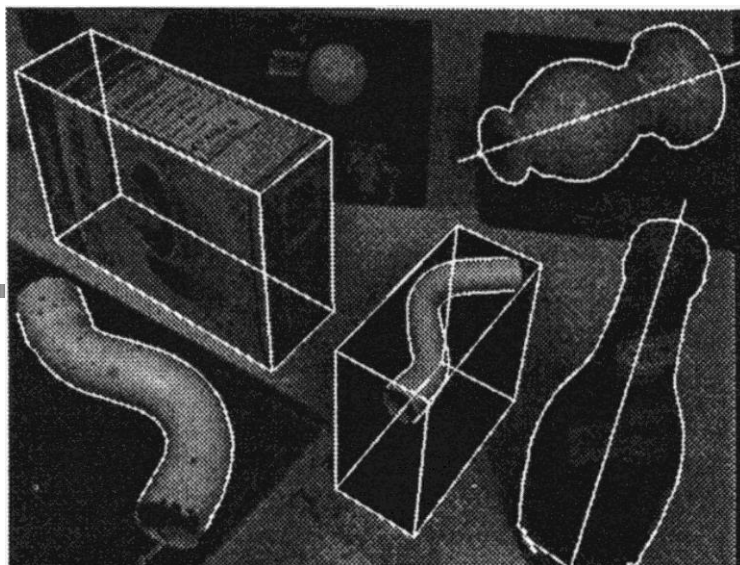Primitives (geons)                                      Objects

Svetlana Lazebnik

# General shape primitives?
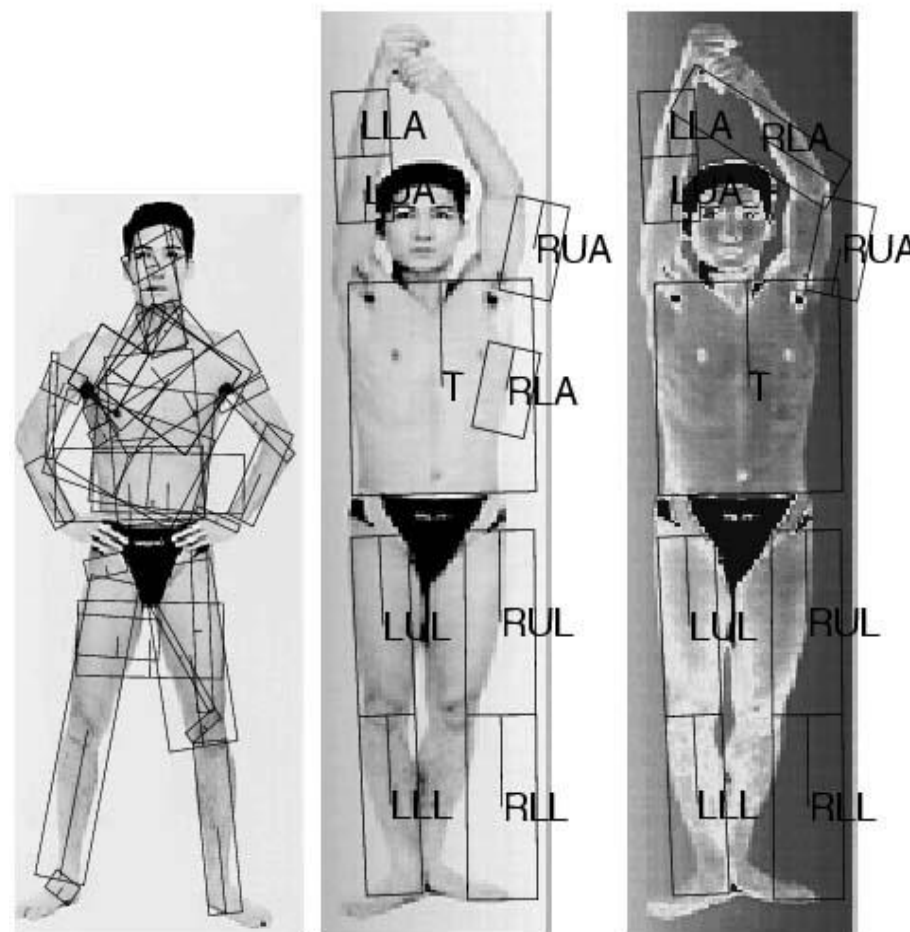
Generalized cylinders
Ponce et al. (1989)

Zisserman et al. (1995)
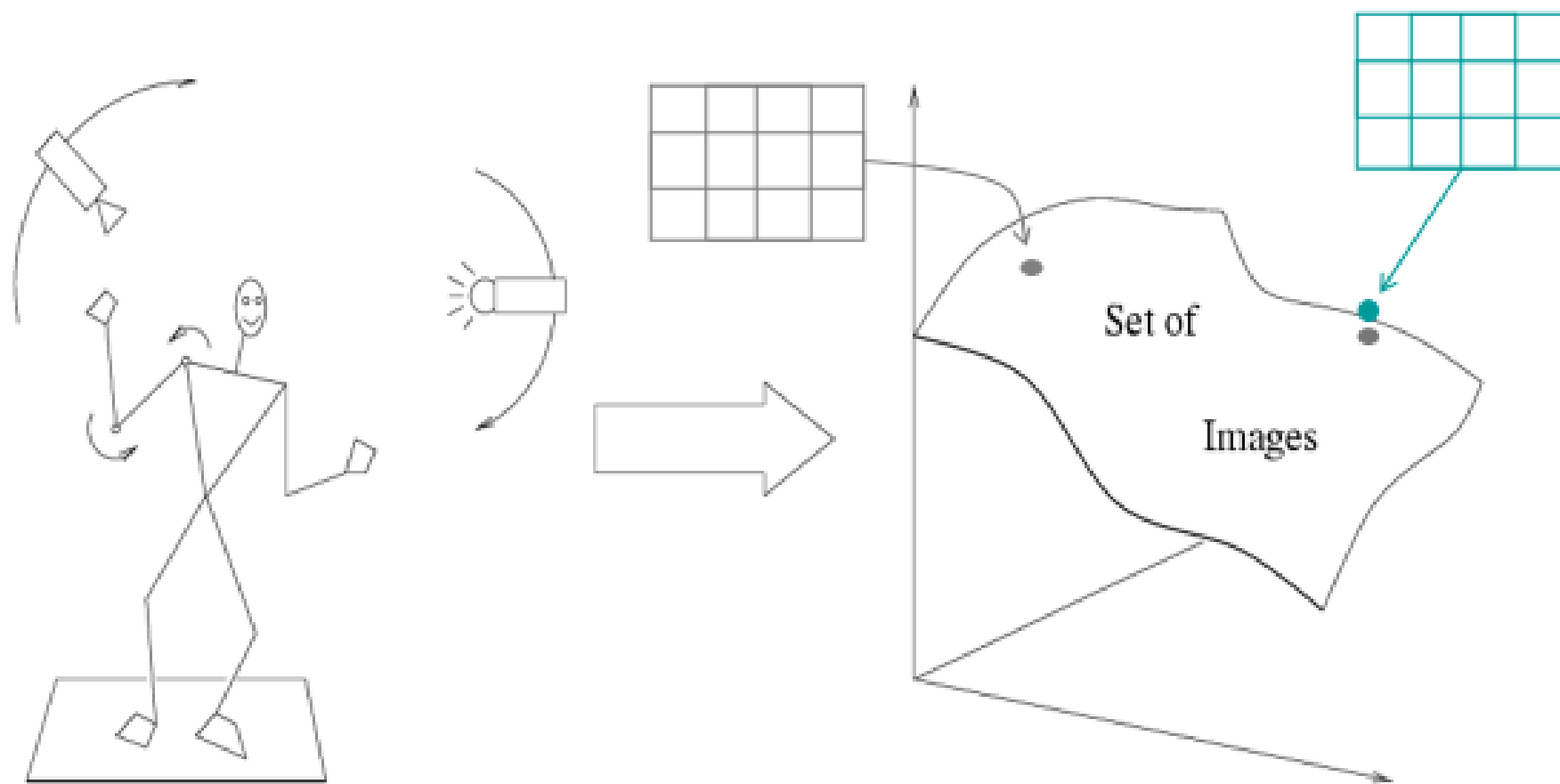
Forsyth (2000)

# History of ideas in recognition

- 1960s – early 1990s: the geometric era
- 1990s: appearance-based models

Svetlana Lazebnik

Empirical models of image variability

**Appearance-based techniques**

Turk & Pentland (1991); Murase & Nayar (1995); etc.

# Eigenfaces (Turk & Pentland, 1991)



| Experimental | Correct/Unknown Recognition Percentage | | |
|---|---|---|---|
| Condition | Lighting | Orientation | Scale |
| Forced classification | 96/0 | 85/0 | 64/0 |
| Forced 100% accuracy | 100/19 | 100/39 | 100/60 |
| Forced 20% unknown rate | 100/20 | 94/20 | 74/20 |

# Color Histograms



Swain and Ballard, Color Indexing, IJCV 1991.

# History of ideas in recognition

- 1960s – early 1990s: the geometric era
- 1990s: appearance-based models
- 1990s – present: sliding window approaches

Svetlana Lazebnik

# Sliding window approaches

# Sliding window approaches



- Turk and Pentland, 1991
- Belhumeur, Hespanha, & Kriegman, 1997
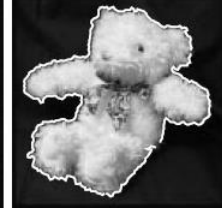- Schneiderman & Kanade 2004
- Viola and Jones, 2000



- Schneiderman & Kanade, 2004
- Argawal and Roth, 2002
- Poggio et al. 1993

# History of ideas in recognition

- 1960s – early 1990s: the geometric era
- 1990s: appearance-based models
- Mid-1990s: sliding window approaches
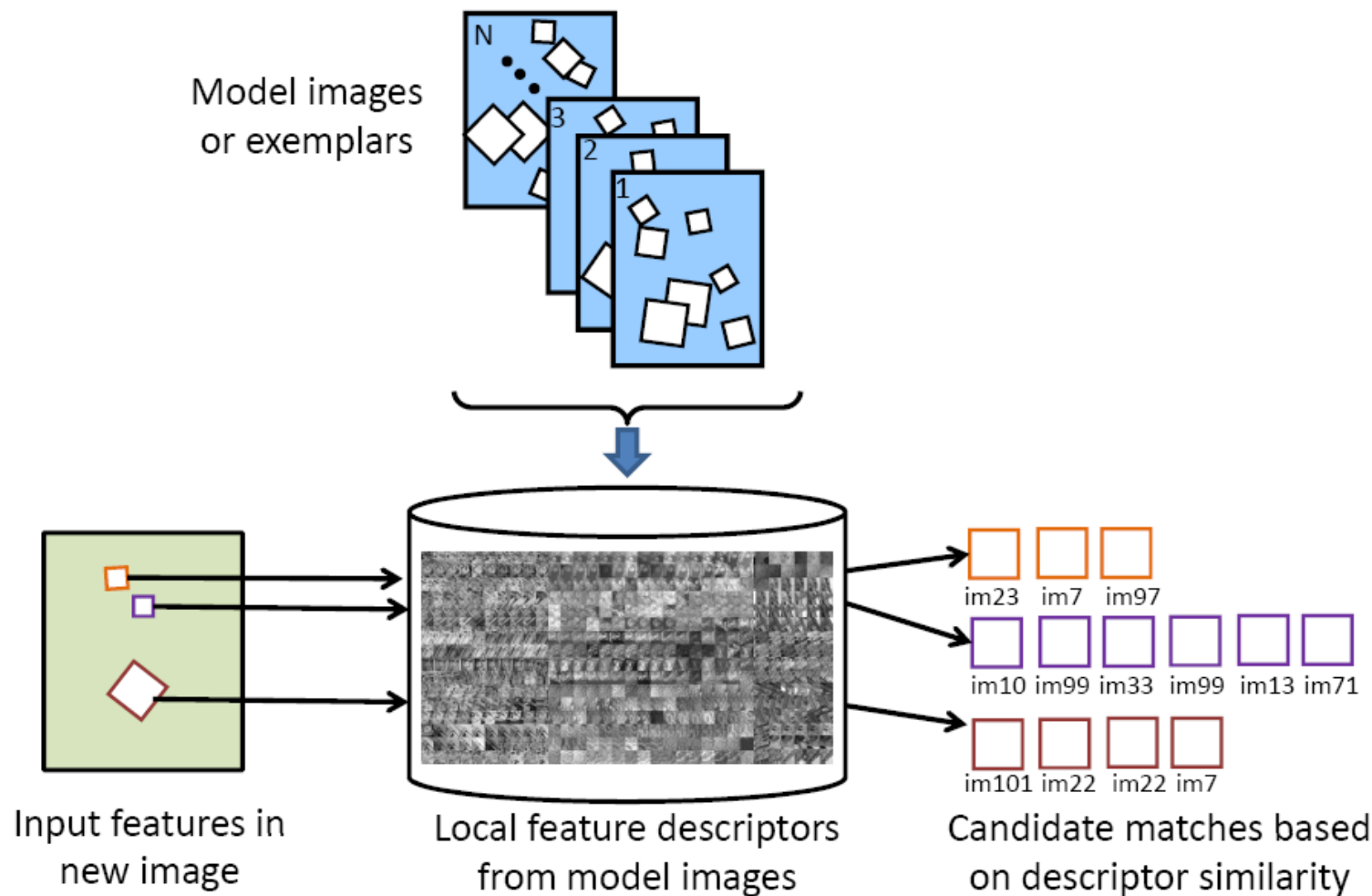- Late 1990s: local features
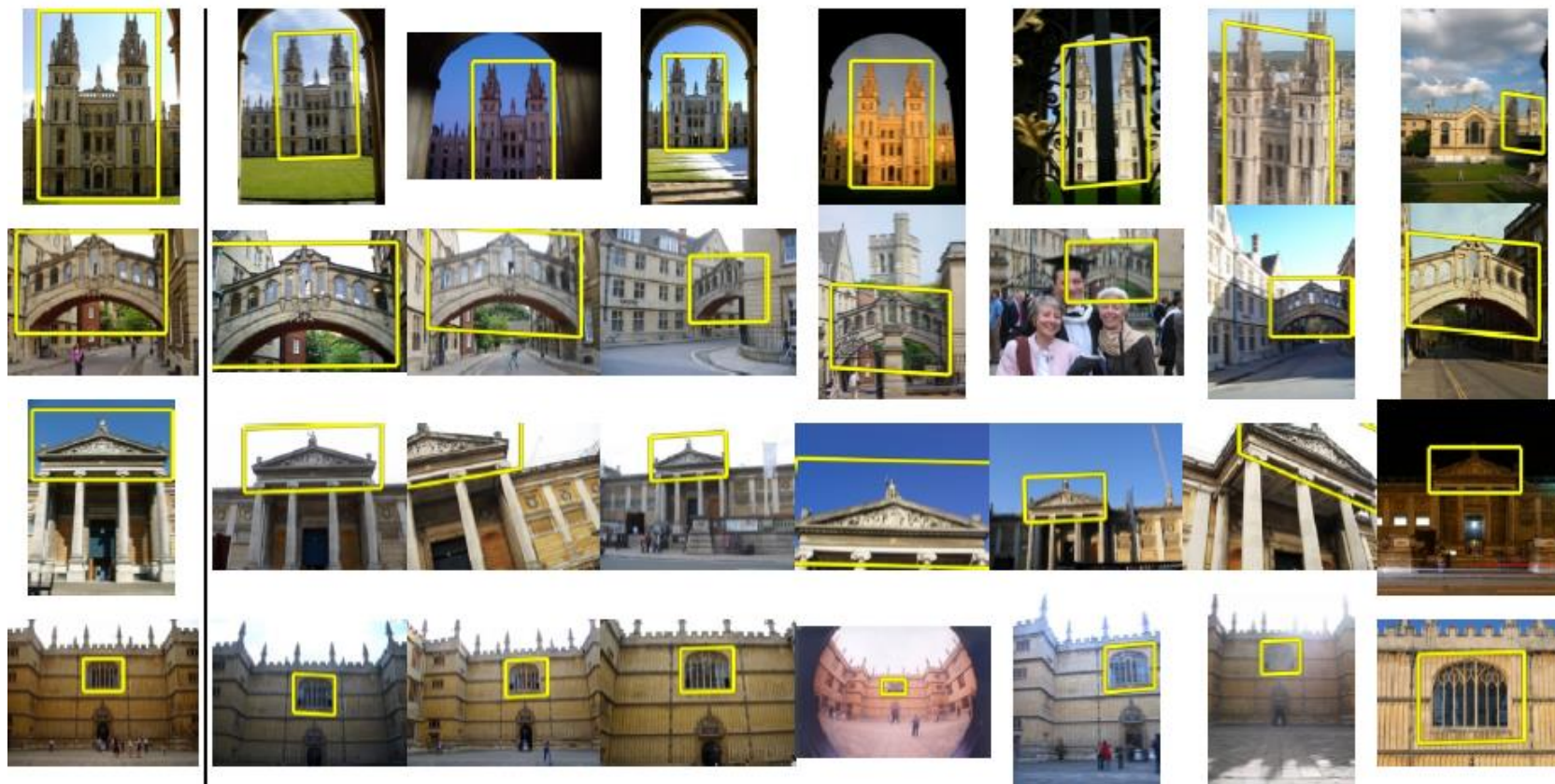
Svetlana Lazebnik

# Local features for object instance recognition



D. Lowe (1999, 2004)

# Large-scale image search

Combining local features, indexing, and spatial constraints



Model images or exemplars

Input features in new image

Local feature descriptors from model images

Candidate matches based on descriptor similarity

im23  im7  im97

im10  im99  im33  im99  im13  im71

im101  im22  im22  im7

Image credit: K. Grauman and B. Leibe

# Large-scale image search

Combining local features, indexing, and spatial constraints



Philbin et al. '07

# Large-scale image search

Combining local features, indexing, and spatial constraints
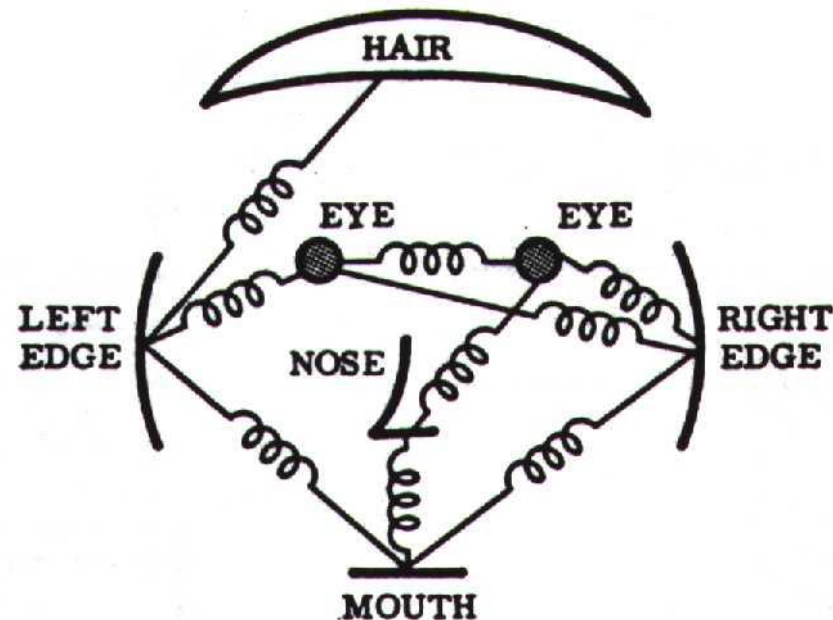


Svetlana Lazebnik

# History of ideas in recognition

- 1960s – early 1990s: the geometric era
- 1990s: appearance-based models
- Mid-1990s: sliding window approaches
- Late 1990s: local features
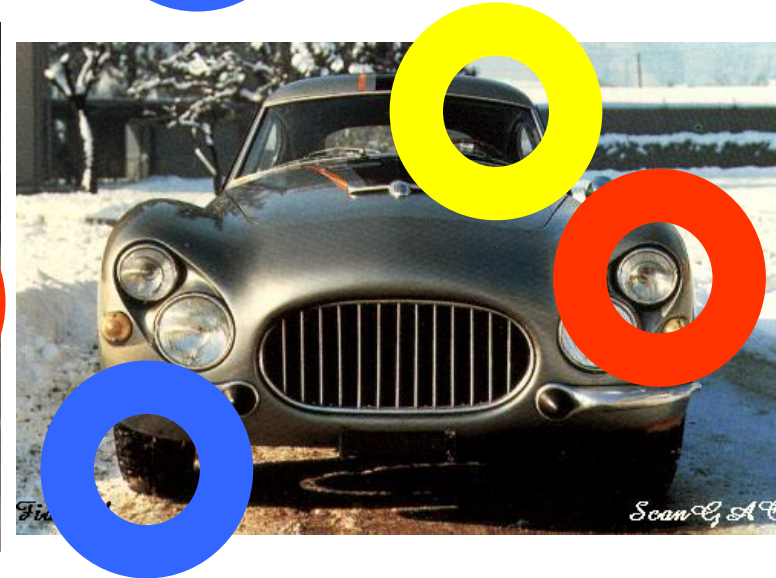- Early 2000s: parts-and-shape models

# Parts-and-shape models

- Model:
  - Object as a set of parts
  - Relative locations between parts
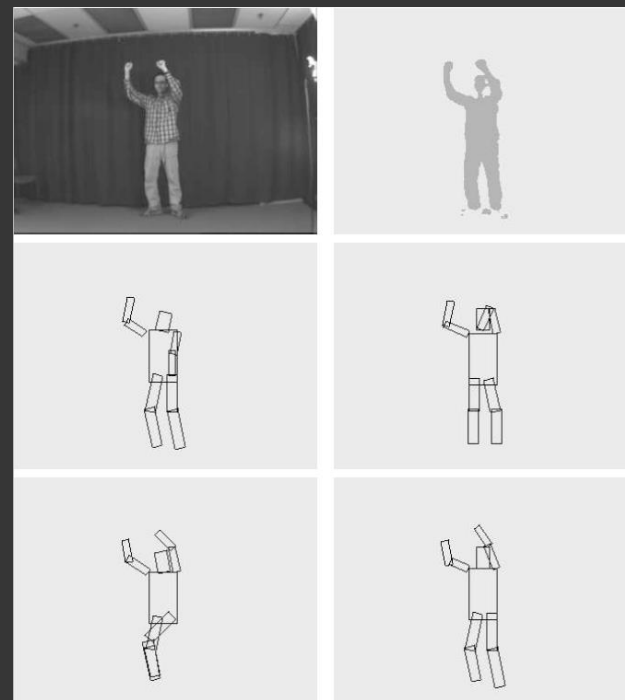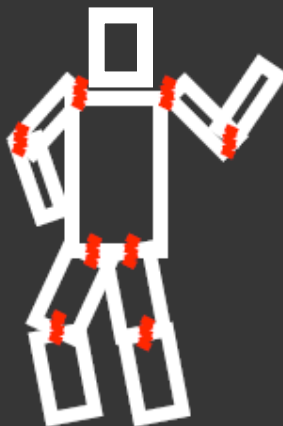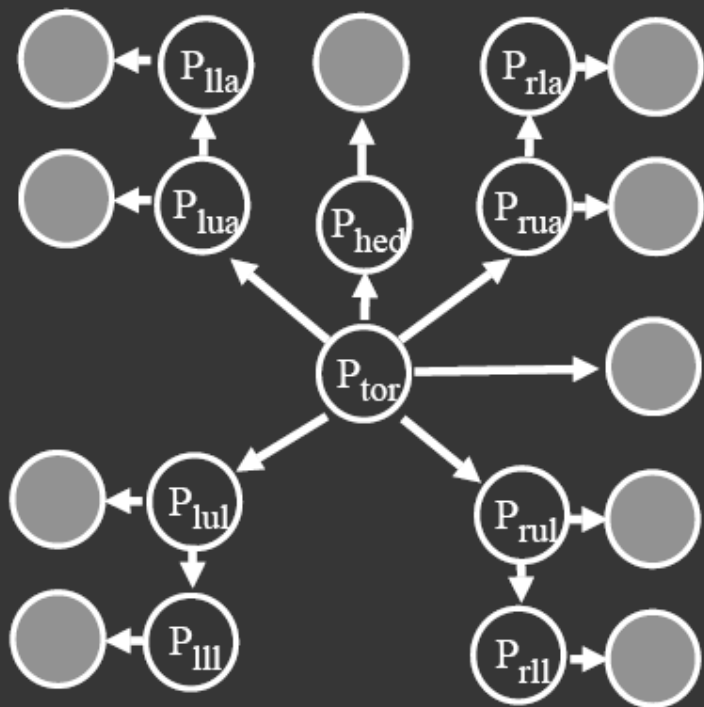  - Appearance of part

# Constellation models



Weber, Welling & Perona (2000), Fergus, Perona & Zisserman (2003)

# Pictorial structure model

Fischler and Elschlager(73), Felzenszwalb and Huttenlocher(00)
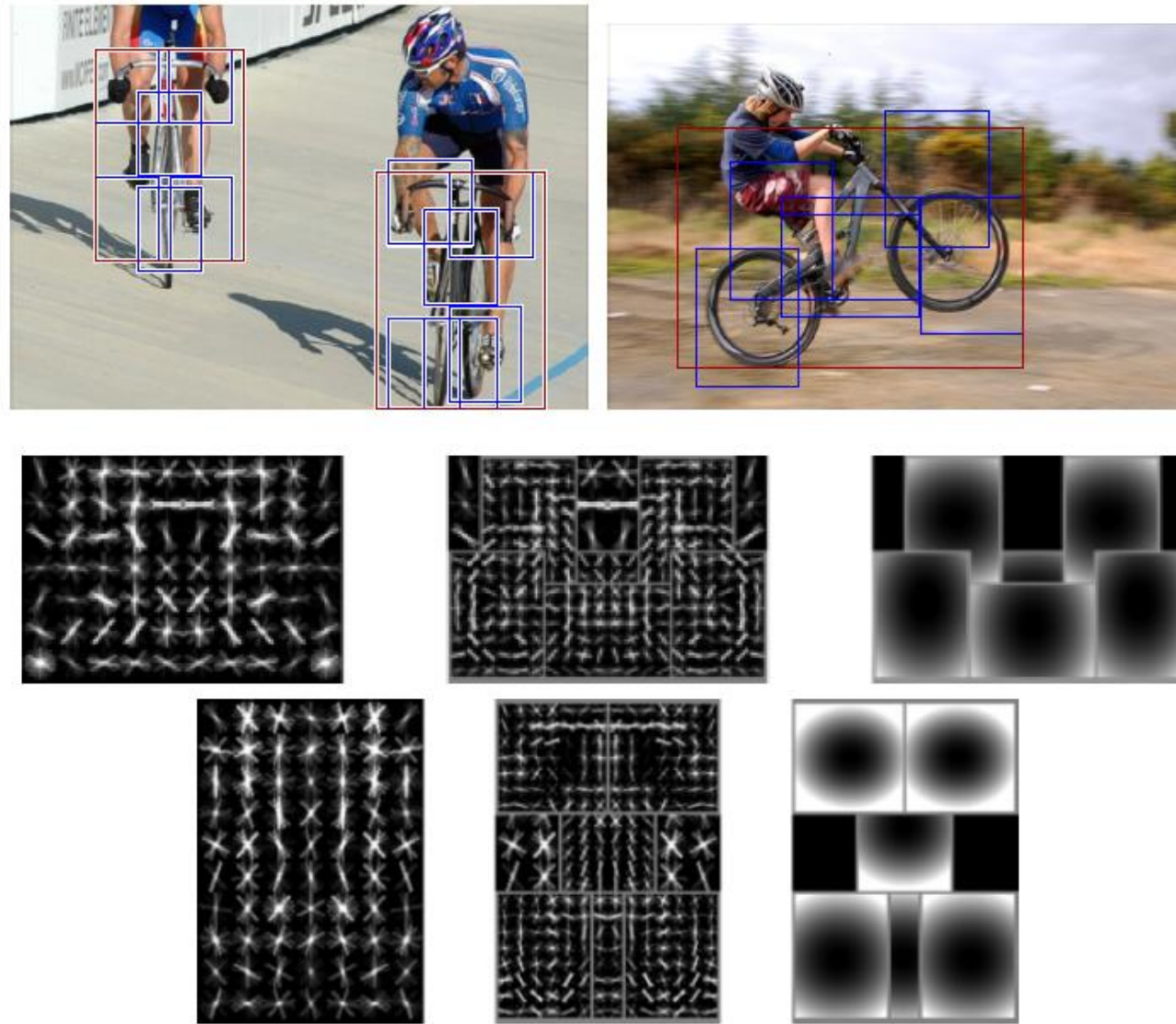


$$\text{Pr}(\text{P}_{\text{tor}}, \text{P}_{\text{arm}}, \dots | \text{Im}) \propto \prod_{i,j} \text{Pr}(P_i | P_j) \prod_i \text{Pr}(\text{Im}(P_i))$$

part geometry    part appearance

# Discriminatively trained part-based models



P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, "Object Detection with Discriminatively Trained Part-Based Models," PAMI 2009
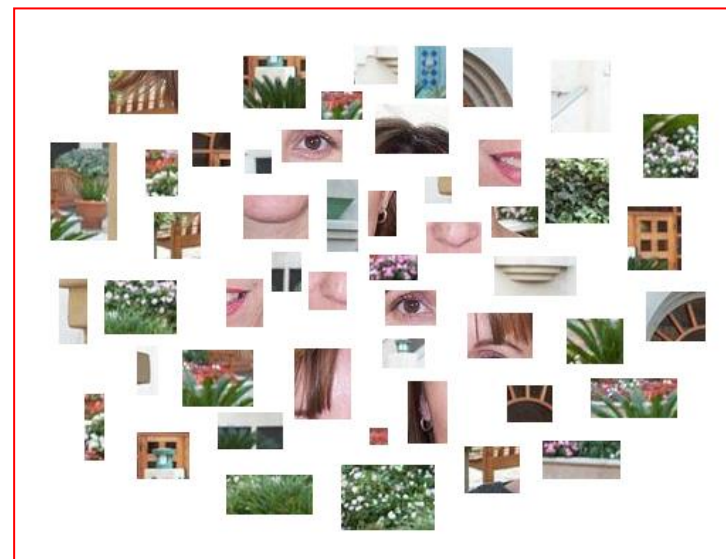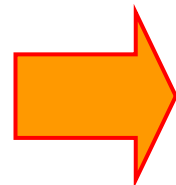
# History of ideas in recognition

- 1960s – early 1990s: the geometric era
- 1990s: appearance-based models
- Mid-1990s: sliding window approaches
- Late 1990s: local features
- Early 2000s: parts-and-shape models
- Mid-2000s: bags of features

# Bag-of-features models

# Bag-of-features models



Object → Bag of 'words'

Svetlana Lazebnik

# Objects as texture

- All of these are treated as being the same



- No distinction between foreground and background: scene recognition?

Svetlana Lazebnik

# Origin 1: Texture recognition

- Texture is characterized by the repetition of basic elements or *textons*

- For stochastic textures, it is the identity of the textons, not their spatial arrangement, that matters

Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

# Origin 1: Texture recognition



histogram

Universal texton dictionary

Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

# Origin 2: Bag-of-words models

- Orderless document representation: frequencies of words from a dictionary  Salton & McGill (1983)

# Origin 2: Bag-of-words models

- Orderless document representation: frequencies of words from a dictionary  Salton & McGill (1983)



US Presidential Speeches Tag Cloud
http://chir.ag/phernalia/preztags/

# Origin 2: Bag-of-words models

- Orderless document representation: frequencies of words from a dictionary   Salton & McGill (1983)



US Presidential Speeches Tag Cloud
http://chir.ag/phernalia/preztags/

# Origin 2: Bag-of-words models

- Orderless document representation: frequencies of words from a dictionary  Salton & McGill (1983)



US Presidential Speeches Tag Cloud
http://chir.ag/phernalia/preztags/

# Bag-of-features steps

1. Extract features
2. Learn "visual vocabulary"
3. Quantize features using visual vocabulary
4. Represent images by frequencies of "visual words"

# 1. Feature extraction

- Regular grid or interest regions

# 1. Feature extraction



**Compute descriptor**

**Normalize patch**

Detect patches

# 1. Feature extraction

# 2. Learning the visual vocabulary

# 2. Learning the visual vocabulary



Clustering

# 2. Learning the visual vocabulary



Visual vocabulary

Clustering

# Clustering and vector quantization

- Clustering is a common method for learning a visual vocabulary or codebook

  - Unsupervised learning process
  - Each cluster center produced by k-means becomes a codevector
  - Codebook can be learned on separate training set
  - Provided the training set is sufficiently representative, the codebook will be "universal"

- The codebook is used for quantizing features

  - A *vector quantizer* takes a feature vector and maps it to the index of the nearest codevector in a codebook
  - Codebook = visual vocabulary
  - Codevector = visual word

# Example codebook



**Appearance codebook**

...

# Visual vocabularies: Issues

- How to choose vocabulary size?
  - Too small: visual words not representative of all patches
  - Too large: quantization artifacts, overfitting

- Computational efficiency
  - Vocabulary trees
    (Nister & Stewenius, 2006)

# But what about layout?



All of these images have the same color histogram

# Spatial pyramid



Compute histogram in each spatial bin

# Spatial pyramid representation

- Extension of a bag of features
- Locally orderless representation at several levels of resolution



level 0

# Spatial pyramid representation

- Extension of a bag of features
- Locally orderless representation at several levels of resolution



level 0                           level 1

Lazebnik, Schmid & Ponce (CVPR 2006)

# Spatial pyramid representation

- Extension of a bag of features
- Locally orderless representation at several levels of resolution



level 0                    level 1                    level 2

# Scene category dataset



office    kitchen    living room    bedroom    store

industrial    tall building    inside city    street    highway

coast    open country    mountain    forest    suburb

## Multi-class classification results
### (100 training images per class)

| Level | Weak features (vocabulary size: 16) | | Strong features (vocabulary size: 200) | |
|---|---|---|---|---|
| | Single-level | Pyramid | Single-level | Pyramid |
| 0 (1 × 1) | 45.3 ±0.5 | | 72.2 ±0.6 | |
| 1 (2 × 2) | 53.6 ±0.3 | 56.2 ±0.6 | 77.9 ±0.6 | 79.0 ±0.5 |
| 2 (4 × 4) | 61.7 ±0.6 | 64.7 ±0.7 | 79.4 ±0.3 | **81.1** ±0.3 |
| 3 (8 × 8) | 63.3 ±0.8 | **66.8** ±0.6 | 77.2 ±0.4 | 80.7 ±0.3 |

# Caltech101 dataset

**Multi-class classification results (30 training images per class)**

| Level | Weak features (16) | | Strong features (200) | |
|-------|--------------|---------|--------------|---------|
|       | Single-level | Pyramid | Single-level | Pyramid |
| 0 | 15.5 ±0.9 |          | 41.2 ±1.2 |          |
| 1 | 31.4 ±1.2 | 32.8 ±1.3 | 55.9 ±0.9 | 57.0 ±0.8 |
| 2 | 47.2 ±1.1 | 49.3 ±1.4 | 63.6 ±0.9 | **64.6** ±0.8 |
| 3 | 52.2 ±0.8 | **54.0** ±1.1 | 60.3 ±0.9 | 64.6 ±0.7 |

# History of ideas in recognition

- 1960s – early 1990s: the geometric era
- 1990s: appearance-based models
- Mid-1990s: sliding window approaches
- Late 1990s: local features
- Early 2000s: parts-and-shape models
- Mid-2000s: bags of features
- Present trends: combination of local and global methods, context, *deep learning*

Svetlana Lazebnik

# Beyond AlexNet

# VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE-SCALE IMAGE RECOGNITION

**Karen Simonyan & Andrew Zisserman 2015**

**These are the "VGG" networks.**
**"Perceptual Loss" in generative deep learning refers to these networks**

| ConvNet Configuration | | | | | |
|---|---|---|---|---|---|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| input (224 × 224 RGB image) | | | | | |
| conv3-64 | conv3-64 **LRN** | conv3-64 **conv3-64** | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 |
| maxpool | | | | | |
| conv3-128 | conv3-128 | conv3-128 **conv3-128** | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 |
| maxpool | | | | | |
| conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 **conv1-256** | conv3-256 conv3-256 **conv3-256** | conv3-256 conv3-256 conv3-256 **conv3-256** |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 **conv1-512** | conv3-512 conv3-512 **conv3-512** | conv3-512 conv3-512 conv3-512 **conv3-512** |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 **conv1-512** | conv3-512 conv3-512 **conv3-512** | conv3-512 conv3-512 conv3-512 **conv3-512** |
| maxpool | | | | | |
| FC-4096 | | | | | |
| FC-4096 | | | | | |
| FC-1000 | | | | | |
| soft-max | | | | | |

Table 2: **Number of parameters** (in millions).

| Network | A,A-LRN | B | C | D | E |
|---|---|---|---|---|---|
| Number of parameters | 133 | 133 | 134 | 138 | 144 |

Table 4: **ConvNet performance at multiple test scales.**

| ConvNet config. (Table 1) | smallest image side | | top-1 val. error (%) | top-5 val. error (%) |
|---|---|---|---|---|
| | train ($S$) | test ($Q$) | | |
| B | 256 | 224,256,288 | 28.2 | 9.6 |
| C | 256 | 224,256,288 | 27.7 | 9.2 |
| | 384 | 352,384,416 | 27.8 | 9.2 |
| | [256; 512] | 256,384,512 | 26.3 | 8.2 |
| D | 256 | 224,256,288 | 26.6 | 8.6 |
| | 384 | 352,384,416 | 26.5 | 8.6 |
| | [256; 512] | 256,384,512 | **24.8** | **7.5** |
| E | 256 | 224,256,288 | 26.9 | 8.7 |
| | 384 | 352,384,416 | 26.7 | 8.6 |
| | [256; 512] | 256,384,512 | **24.8** | **7.5** |

# Going Deeper with Convolutions

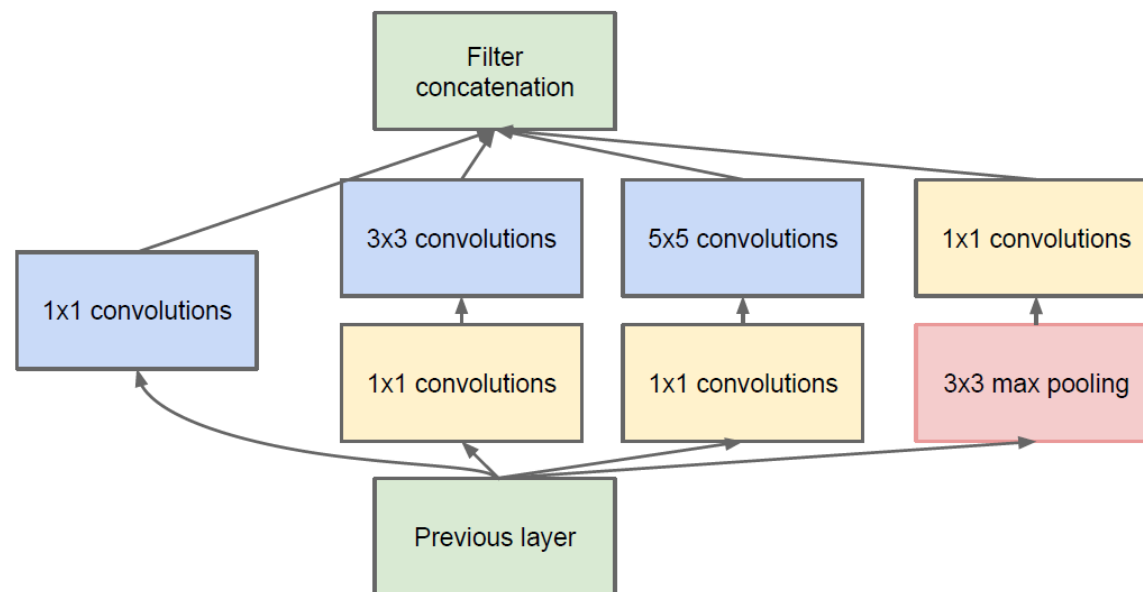**Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, Andrew Rabinovich**
**2015**

**This is the "Inception" architecture or "GoogLeNet"**

**\*The architecture blocks are called "Inception" modules
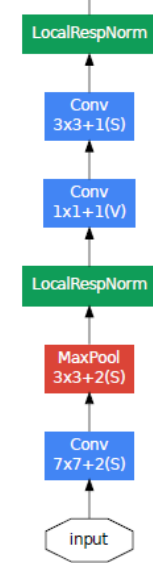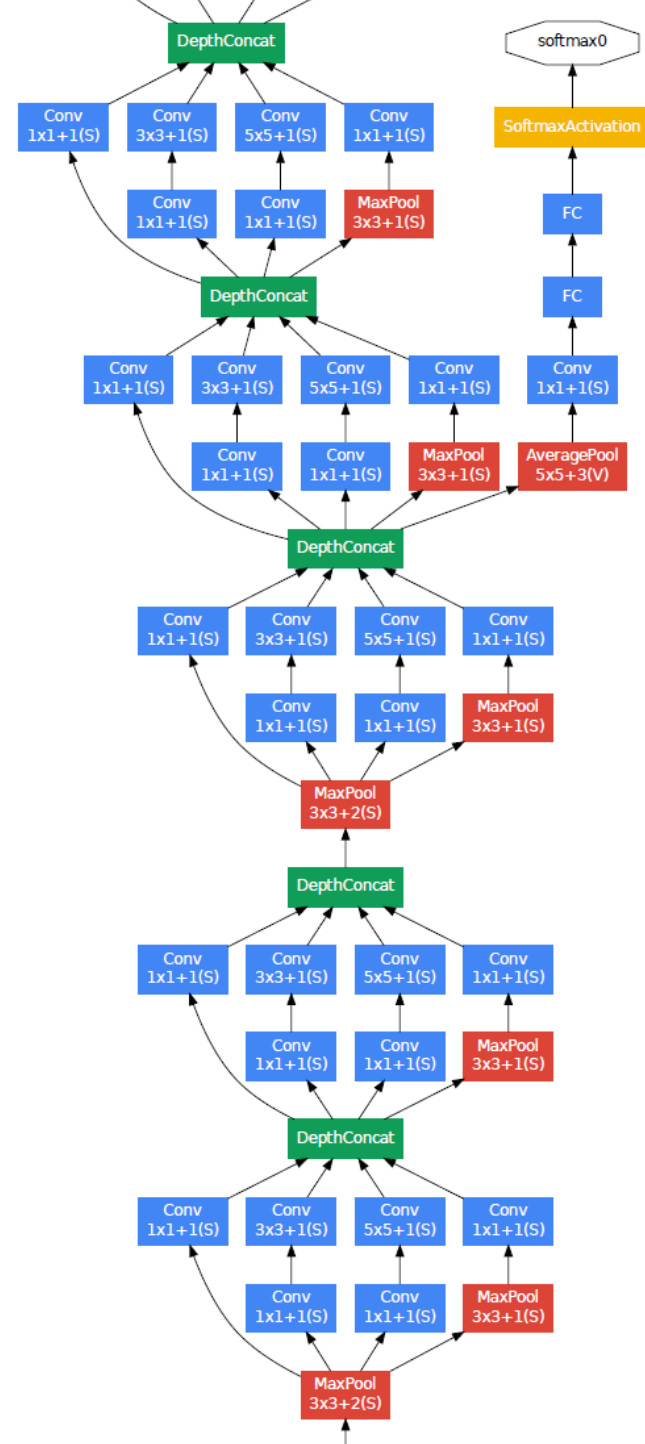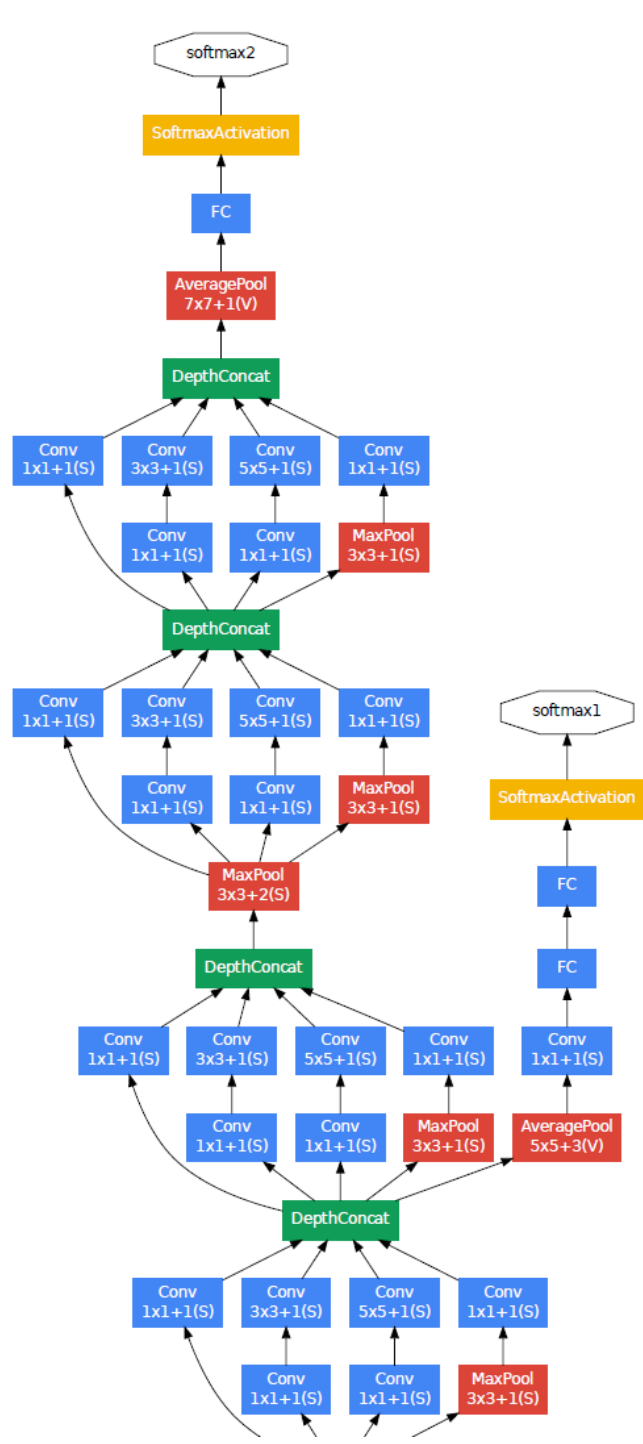and the collection of them into a particular net is "GoogLeNet"**

(a) Inception module, naïve version



(b) Inception module with dimensionality reduction

| type | patch size/ stride | output size | depth | #1×1 | #3×3 reduce | #3×3 | #5×5 reduce | #5×5 | pool proj | params | ops |
|---|---|---|---|---|---|---|---|---|---|---|---|
| convolution | 7×7/2 | 112×112×64 | 1 | | | | | | | 2.7K | 34M |
| max pool | 3×3/2 | 56×56×64 | 0 | | | | | | | | |
| convolution | 3×3/1 | 56×56×192 | 2 | | 64 | 192 | | | | 112K | 360M |
| max pool | 3×3/2 | 28×28×192 | 0 | | | | | | | | |
| inception (3a) | | 28×28×256 | 2 | 64 | 96 | 128 | 16 | 32 | 32 | 159K | 128M |
| inception (3b) | | 28×28×480 | 2 | 128 | 128 | 192 | 32 | 96 | 64 | 380K | 304M |
| max pool | 3×3/2 | 14×14×480 | 0 | | | | | | | | |
| inception (4a) | | 14×14×512 | 2 | 192 | 96 | 208 | 16 | 48 | 64 | 364K | 73M |
| inception (4b) | | 14×14×512 | 2 | 160 | 112 | 224 | 24 | 64 | 64 | 437K | 88M |
| inception (4c) | | 14×14×512 | 2 | 128 | 128 | 256 | 24 | 64 | 64 | 463K | 100M |
| inception (4d) | | 14×14×528 | 2 | 112 | 144 | 288 | 32 | 64 | 64 | 580K | 119M |
| inception (4e) | | 14×14×832 | 2 | 256 | 160 | 320 | 32 | 128 | 128 | 840K | 170M |
| max pool | 3×3/2 | 7×7×832 | 0 | | | | | | | | |
| inception (5a) | | 7×7×832 | 2 | 256 | 160 | 320 | 32 | 128 | 128 | 1072K | 54M |
| inception (5b) | | 7×7×1024 | 2 | 384 | 192 | 384 | 48 | 128 | 128 | 1388K | 71M |
| avg pool | 7×7/1 | 1×1×1024 | 0 | | | | | | | | |
| dropout (40%) | | 1×1×1024 | 0 | | | | | | | | |
| linear | | 1×1×1000 | 1 | | | | | | | 1000K | 1M |
| softmax | | 1×1×1000 | 0 | | | | | | | | |

Only 6.8 million parameters. AlexNet ~60 million, VGG up to 138 million

| Team | Year | Place | Error (top-5) | Uses external data |
|---|---|---|---|---|
| SuperVision | 2012 | 1st | 16.4% | no |
| SuperVision | 2012 | 1st | 15.3% | Imagenet 22k |
| Clarifai | 2013 | 1st | 11.7% | no |
| Clarifai | 2013 | 1st | 11.2% | Imagenet 22k |
| MSRA | 2014 | 3rd | 7.35% | no |
| VGG | 2014 | 2nd | 7.32% | no |
| GoogLeNet | 2014 | 1st | 6.67% | no |

Table 2: Classification performance.

| Number of models | Number of Crops | Cost | Top-5 error | compared to base |
|---|---|---|---|---|
| 1 | 1 | 1 | 10.07% | base |
| 1 | 10 | 10 | 9.15% | -0.92% |
| 1 | 144 | 144 | 7.89% | -2.18% |
| 7 | 1 | 7 | 8.09% | -1.98% |
| 7 | 10 | 70 | 7.62% | -2.45% |
| 7 | 144 | 1008 | 6.67% | -3.45% |

# Surely it would be ridiculous to go any deeper...

- To be continued with ResNet