

History of object and scene representations



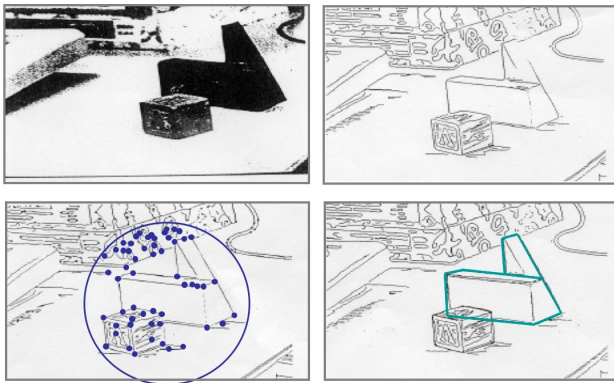
https://en.wikipedia.org/wiki/Landscape_with_the_Fall_of_Icarus

Outline

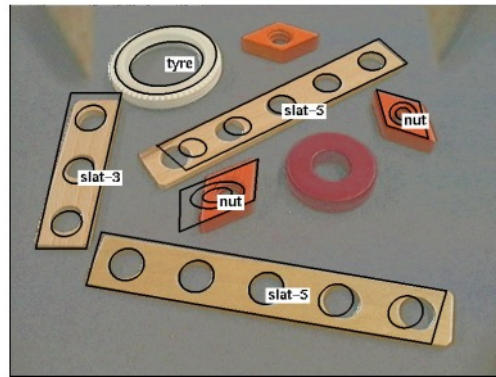
- Object representations
 - 3D shape
 - 3D primitives
 - 2D appearance-based models
 - 2D part-based models (deformable templates)
 - CNNs
- Scene representations
 - Structured representations
 - Appearance-based representations
 - Bottom-up and top-down perceptual organization
- Trends

3D object representations

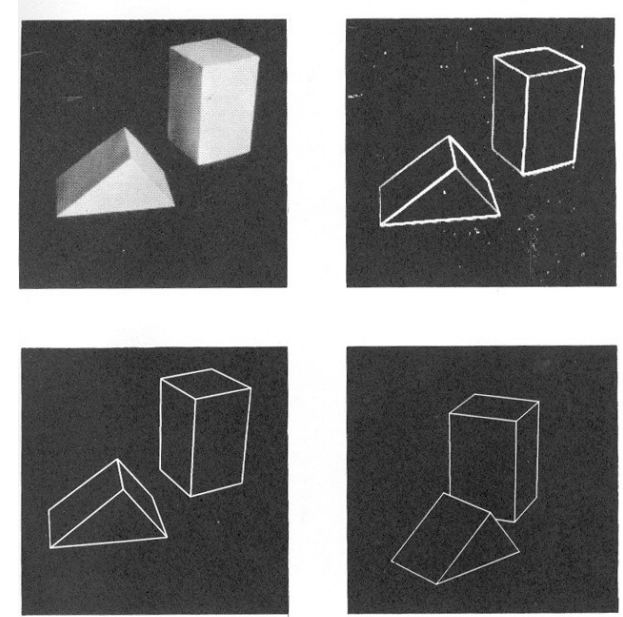
- Represent an object using its 3D model
- Recognition by *alignment* or *geometric invariants*



Alignment: Huttenlocher & Ullman (1987)

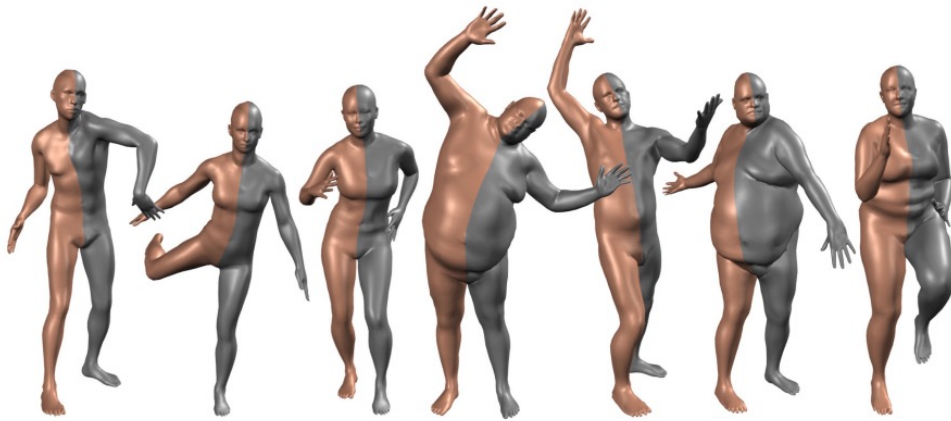


Invariants: [Graf \(2000\)](#)

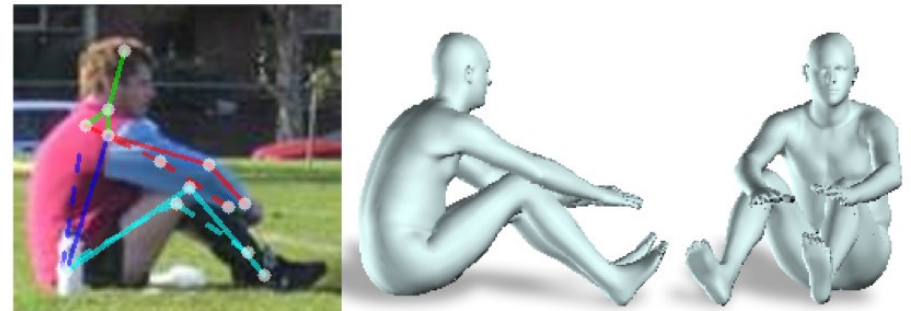


Roberts (1963)

Today: Category-specific 3D models

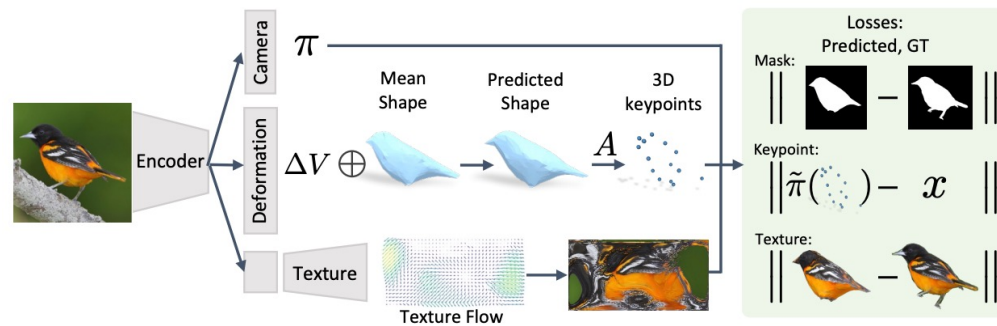


M. Loper et al. [SMPL: A Skinned Multi-Person Linear Model](#).
SIGGRAPH Asia, 2015



F. Bogo et al., [Keep it SMPL: Automatic Estimation of 3D Human Pose and Shape from a Single Image](#),
ECCV 2016

Today: Generic category-level 3D models



A. Kanazawa et al. [Learning Category-Specific Mesh Reconstruction from Image Collections](#). ECCV 2018

3D shape primitives: Generalized cylinders

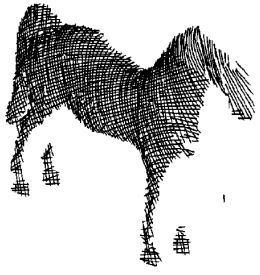


Figure 15
Horse

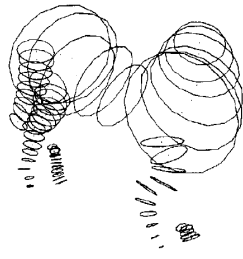


Figure 16
Analysis of Horse

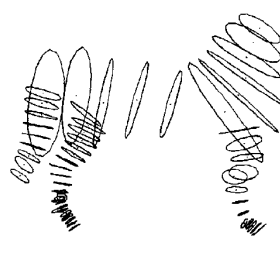
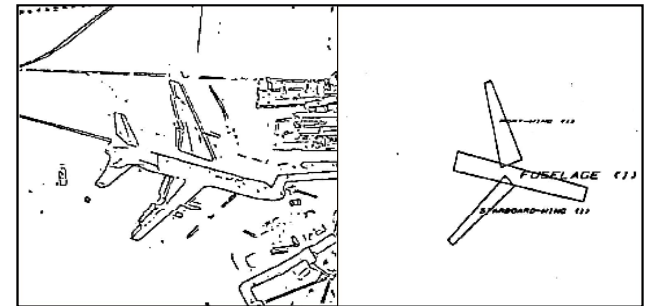
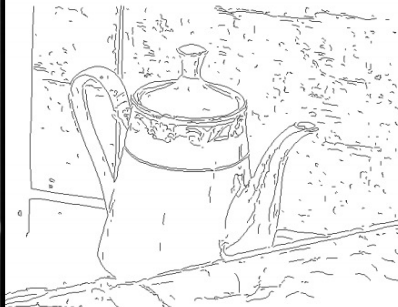


Figure 17
Analysis of Horse, Side View

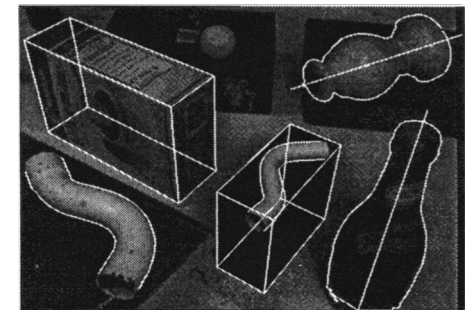
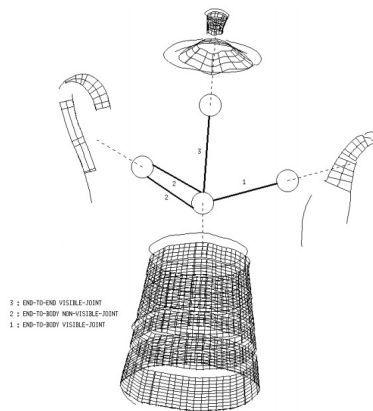
Binford (1971), Agin & Binford (1973)



Brooks (1981)



Zerroug & Nevatia (1994)



Zisserman et al. (1995)

Marr's 3D object representation

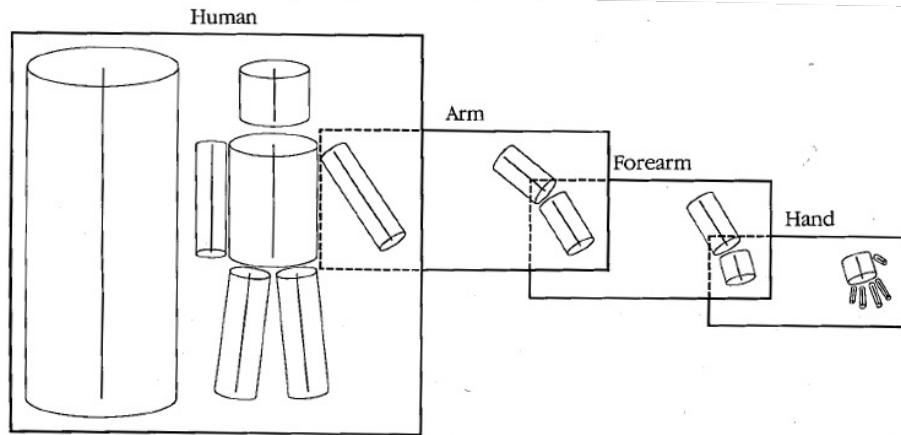
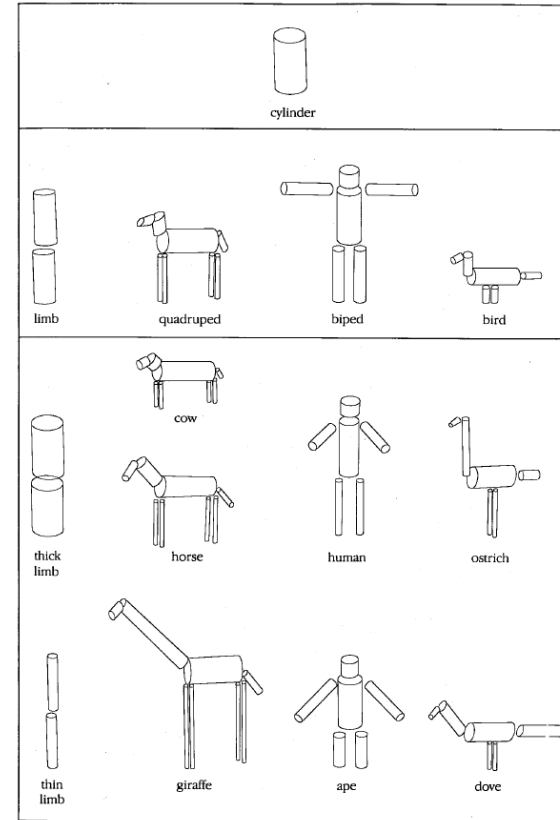
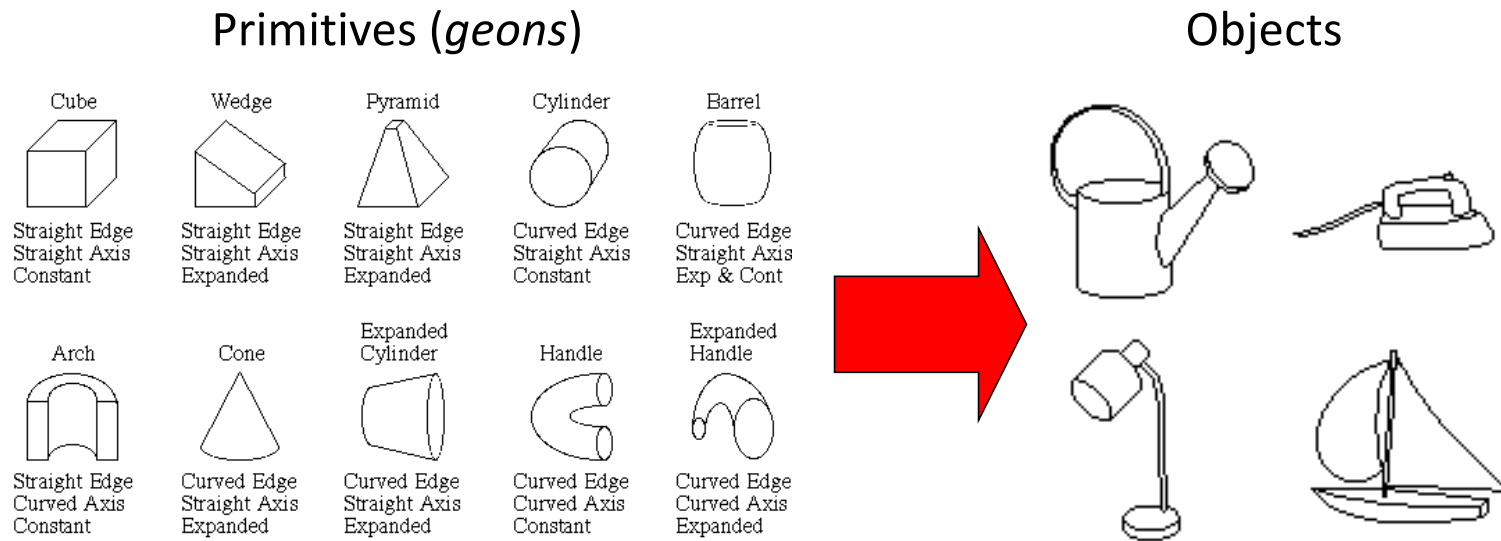


Figure 5-3. This diagram illustrates the organization of shape information in a 3-D model description. Each box corresponds to a 3-D model, with its model axis on the left side of the box and the arrangement of its component axes on the right. In addition, some component axes have 3-D models associated with them, as indicated by the way the boxes overlap. The relative arrangement of each model's component axes, however, is shown improperly, since it should be in an object-centered system rather than the viewer-centered projection used here (a more correct 3-D model is given by the table shown in Figure 5-5c). The important characteristics of this type of organization are: (1) Each 3-D model is a self-contained unit of shape information and has a limited complexity; (2) information appears in shape contexts appropriate for recognition (the disposition of a finger is most stable when specified relative to the hand that contains it); and (3) the representation can be manipulated flexibly. This approach limits the representation's scope, however, since it is only useful for shapes that have well-defined 3-D model decompositions. (Reprinted by permission from D. Marr and H. K. Nishihara, "Representation and recognition of the spatial organization of three-dimensional shapes," *Proc. R. Soc. Lond. B* 200, 269-294.)



Marr & Nishihara (1978)

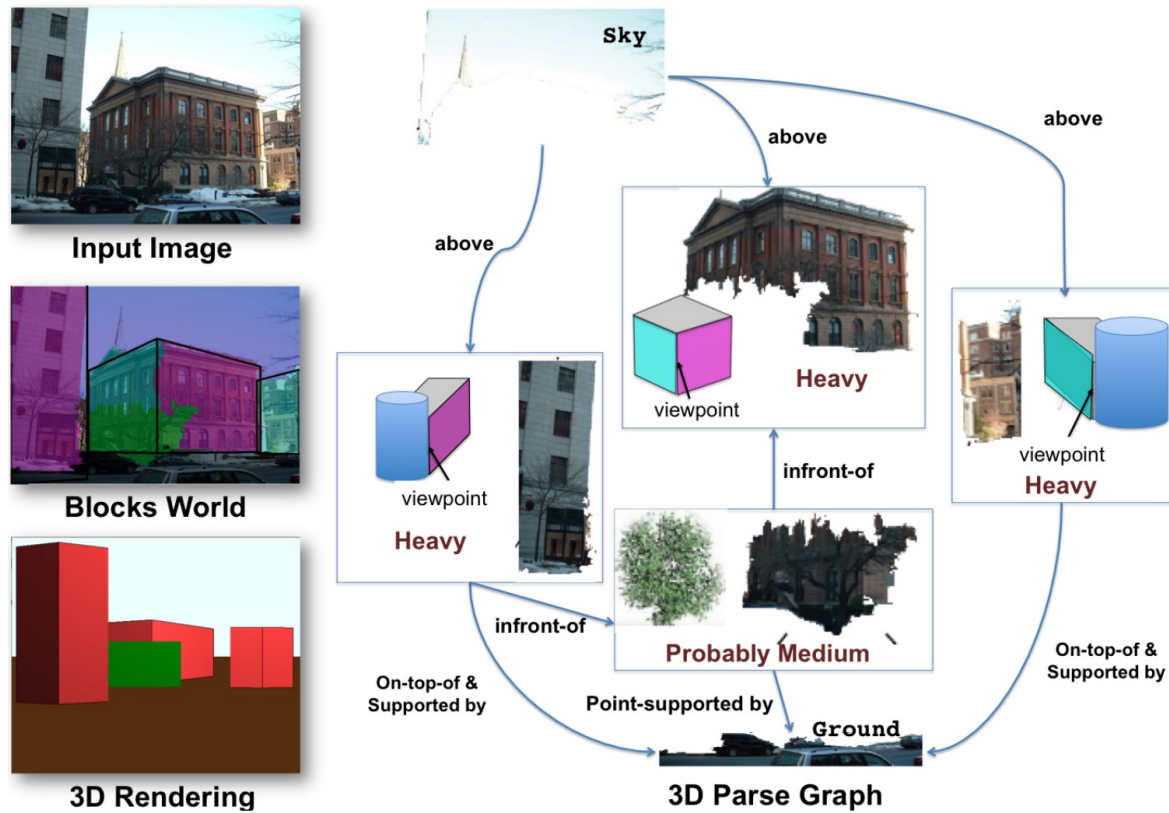
Psychological theory: Recognition by components



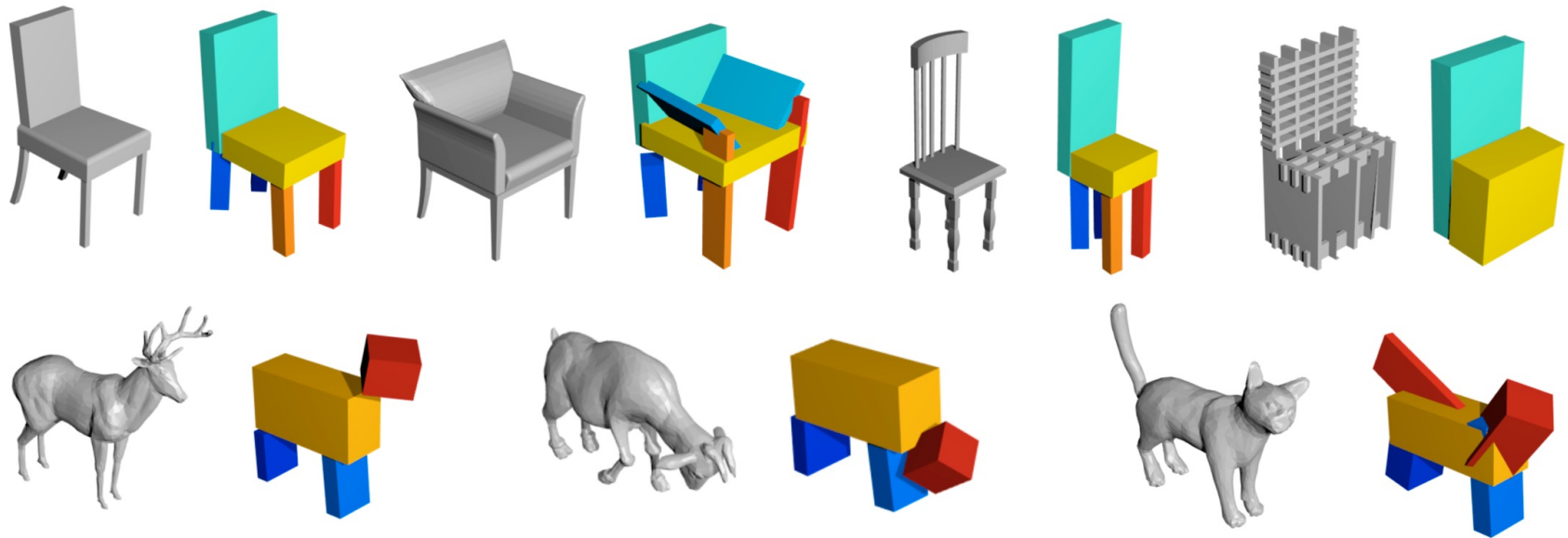
[Biederman \(1987\)](#)

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

Today: Revival of 3D primitives



Today: Revival of 3D primitives



“Here we do not wish to reprise the classic debates on the value of volumetric primitives – while they were oversold in the 70s and 80s, they suffer from complete neglect now, and we hope that this demonstration of feasibility of learning how to assemble an object from volumetric primitives will reignite interest.”

S. Tulsiani et al. [Learning Shape Abstractions by Assembling Volumetric Primitives](#). CVPR 2017

Today: Revival of 3D primitives

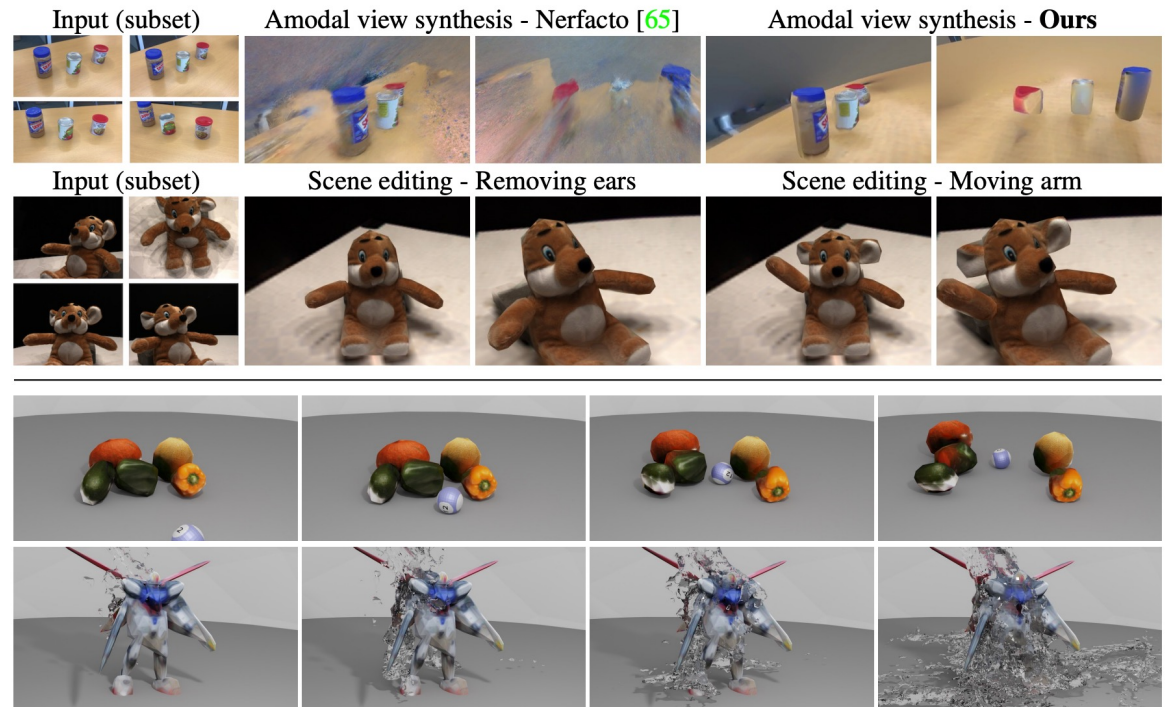


Figure 5: **Applications.** Amodal completion (1st row), scene editing (2nd) and physical simulations (bottom).

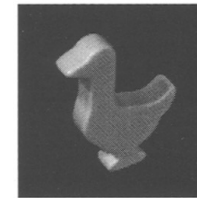
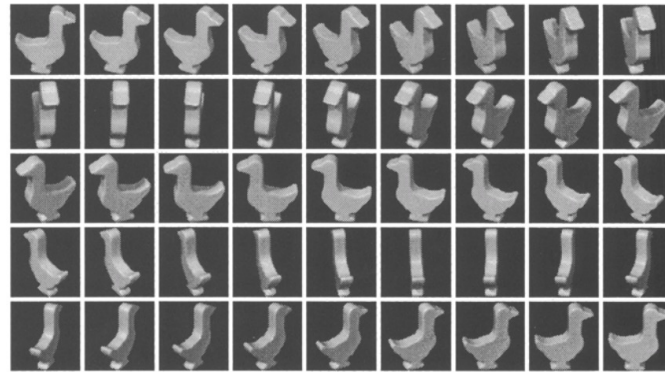
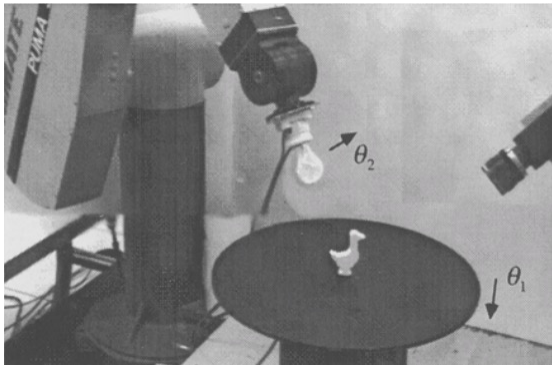
W. Liu et al. [Marching-Primitives: Shape Abstraction from Signed Distance Function](#). CVPR 2023

T. Monnier et al. [Differentiable Blocks World: Qualitative 3D Decomposition by Rendering Primitives](#). arXiv 2023

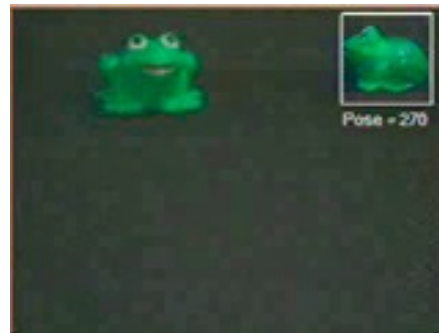
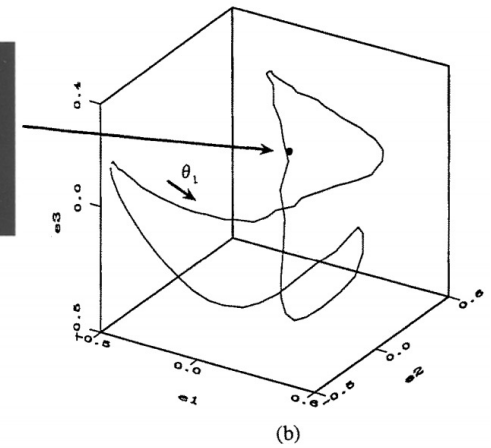
Outline

- Object representations
 - 3D shape
 - 3D primitives
 - 2D appearance-based models
 - 2D part-based models (deformable templates)
 - CNNs

Global appearance representations



(a)



[Recognition demo movie](#)

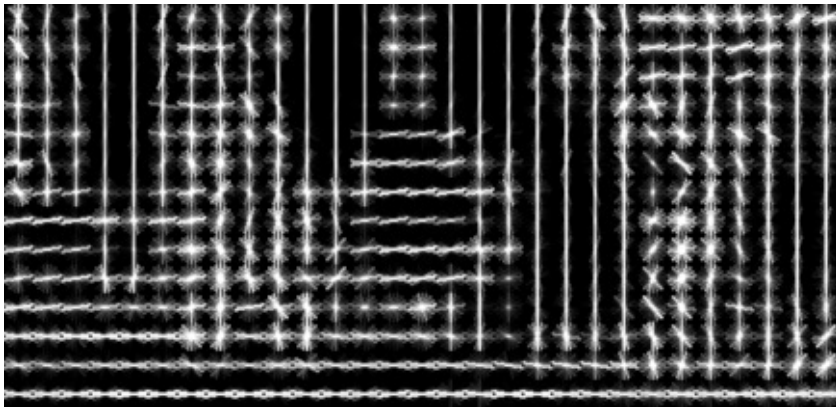
H. Murase and S. Nayar, [Visual learning and recognition of 3-d objects from appearance](#), IJCV 1995

J. Mundy et al., [An Experimental Comparison of Appearance and Geometric Model Based Recognition](#), 1996

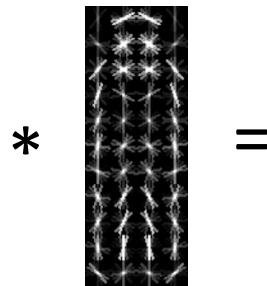
Global appearance representation: HOG template



HOG feature map



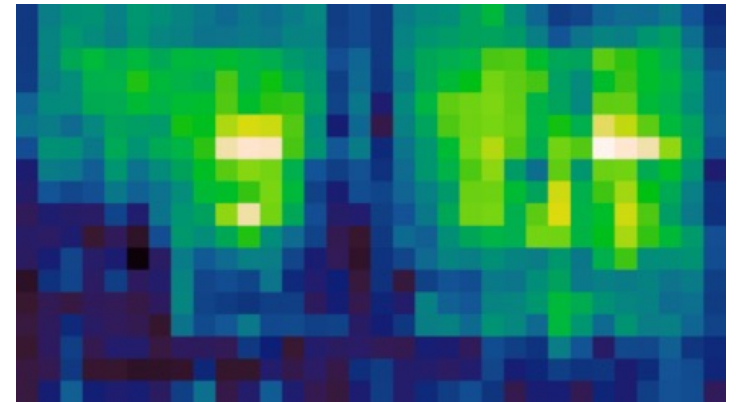
Template



*

=

Detector response map



N. Dalal and B. Triggs, [Histograms of Oriented Gradients for Human Detection](#), CVPR 2005

2D part-based representations: Local appearance + deformable shape

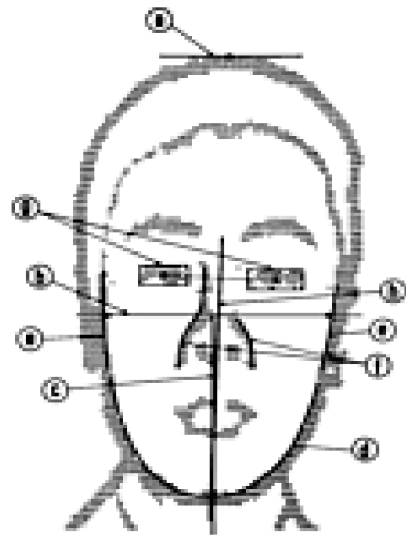
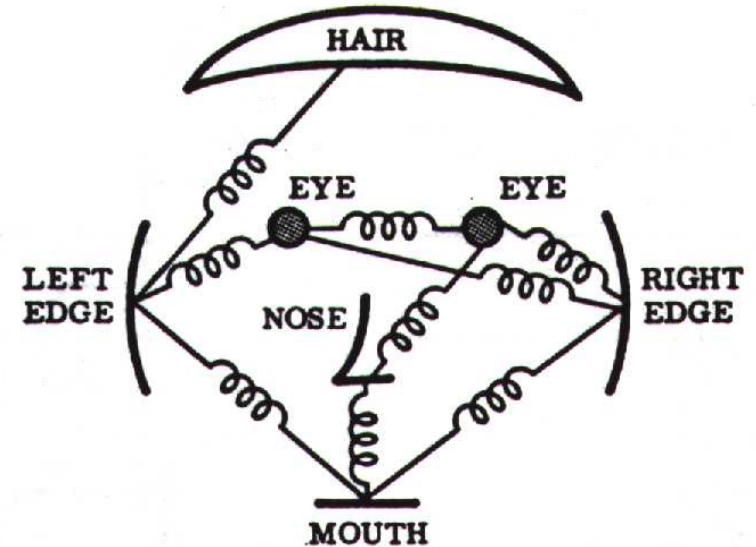


Figure 3-3

Typical sequence of the analysis steps.

- (a) top of head
- (b) cheeks and sides of face
- (c) nose, mouth, and chin
- (d) chin contour
- (e) face-side lines
- (f) nose lines
- (g) eyes
- (h) face axis

[Kanade \(1973\)](#)



[Fischler and Elschlager \(1973\)](#)

A procedural part-based recognition system

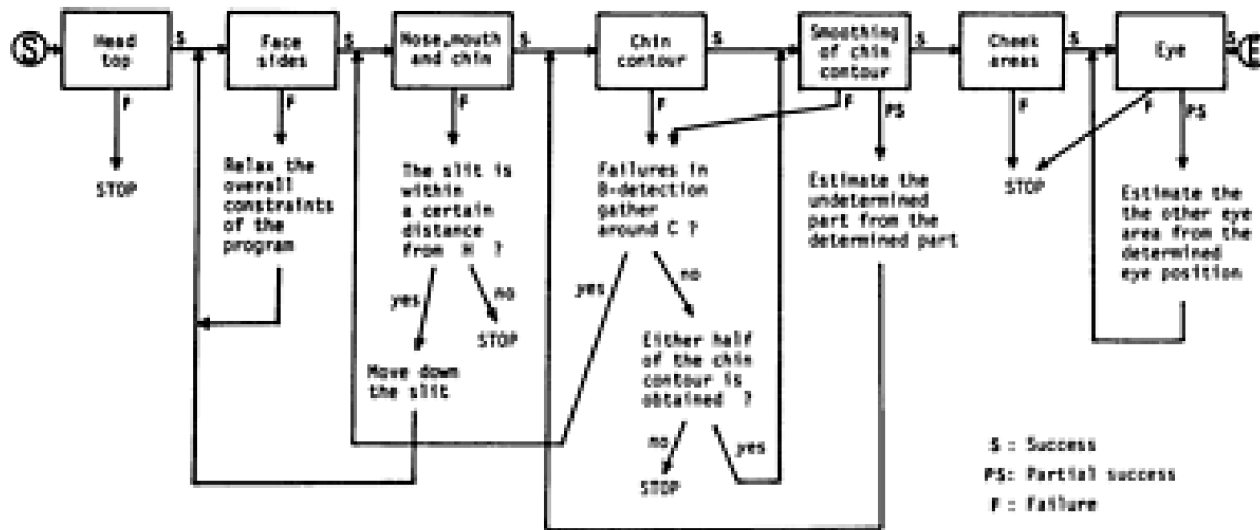


Figure 3-12 General flow of analysis program.

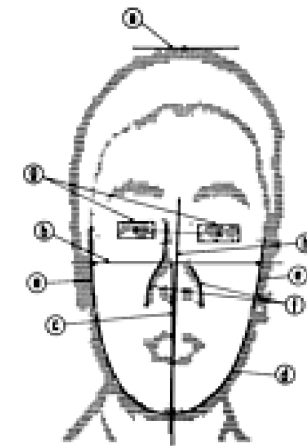


Figure 3-3

Typical sequence of the analysis steps.

- (a) top of head
- (b) cheeks and sides of face
- (c) nose, mouth, and chin
- (d) chin contour
- (e) face-side lines
- (f) nose lines
- (g) eyes
- (h) face axis

T. Kanade. [Picture Processing System by Computer Complex and Recognition of Human Faces](#). Ph.D. dissertation 1973

A procedural part-based recognition system

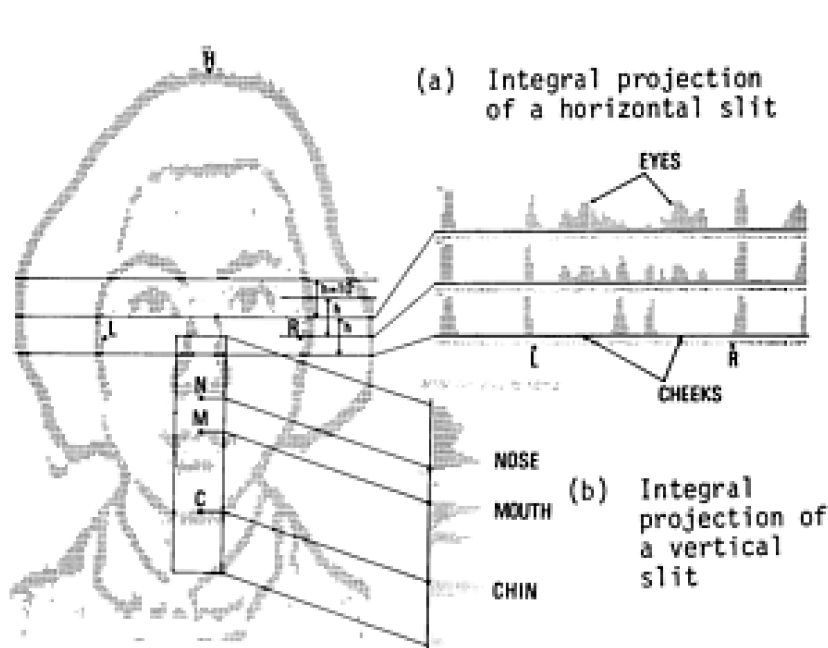


Figure 3-14 Detection of the face sides, nose, mouth, and chin by the application of slits.

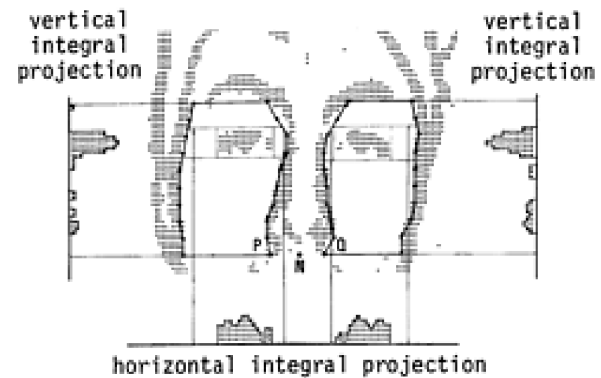


Figure 3-18 Cheek areas and rectangles which contain the eye.

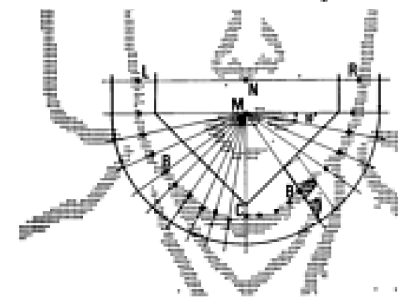


Figure 3-17

Extraction of chin contour. The search area is established and a slit is placed along each radial line.

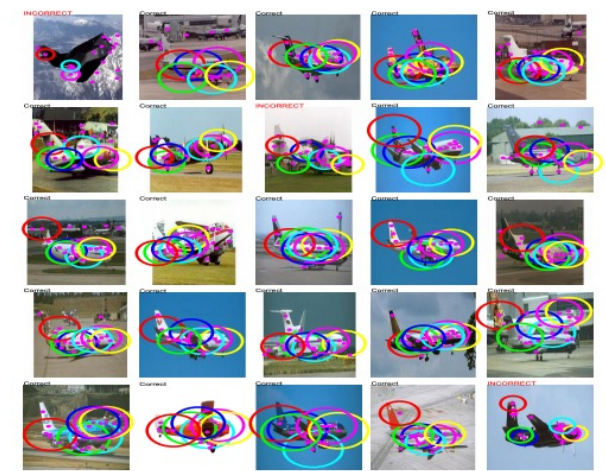
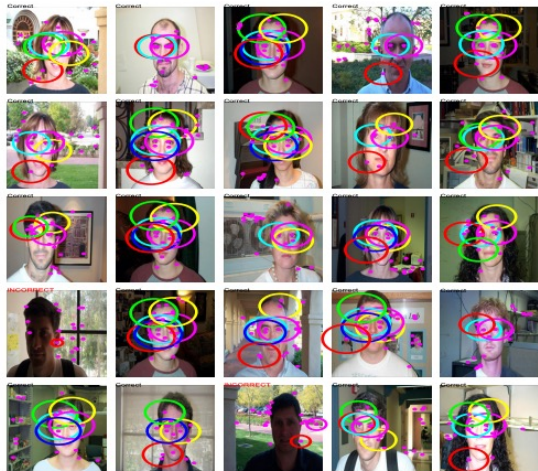
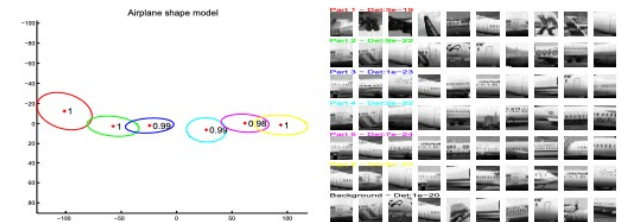
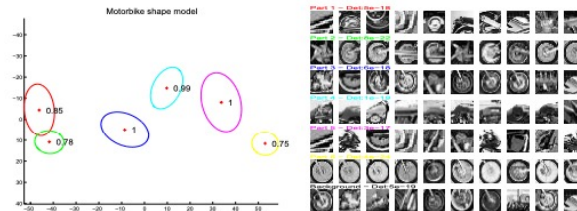
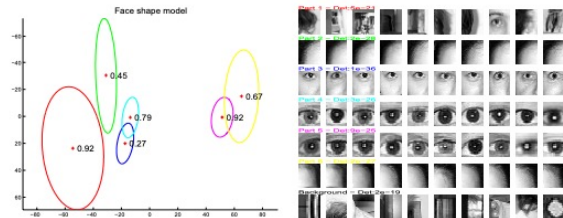
A procedural part-based recognition system

category of faces	number of faces	correct results	error or unrecovered failure	step in which the error or unrecovered failure occurred				
				face sides	nose mouth chin	chin contour	eyes	nose width
full face with no glasses or beard	670	608	62	5	14	17	19	7
full face with glasses	77	2	75	4	4	2	65	
face with turn or tilt	79	63	16	4	3	3	5	1
face with beard, and others	27		27		12	4		11

Table 3-1 Summary of results of analysis

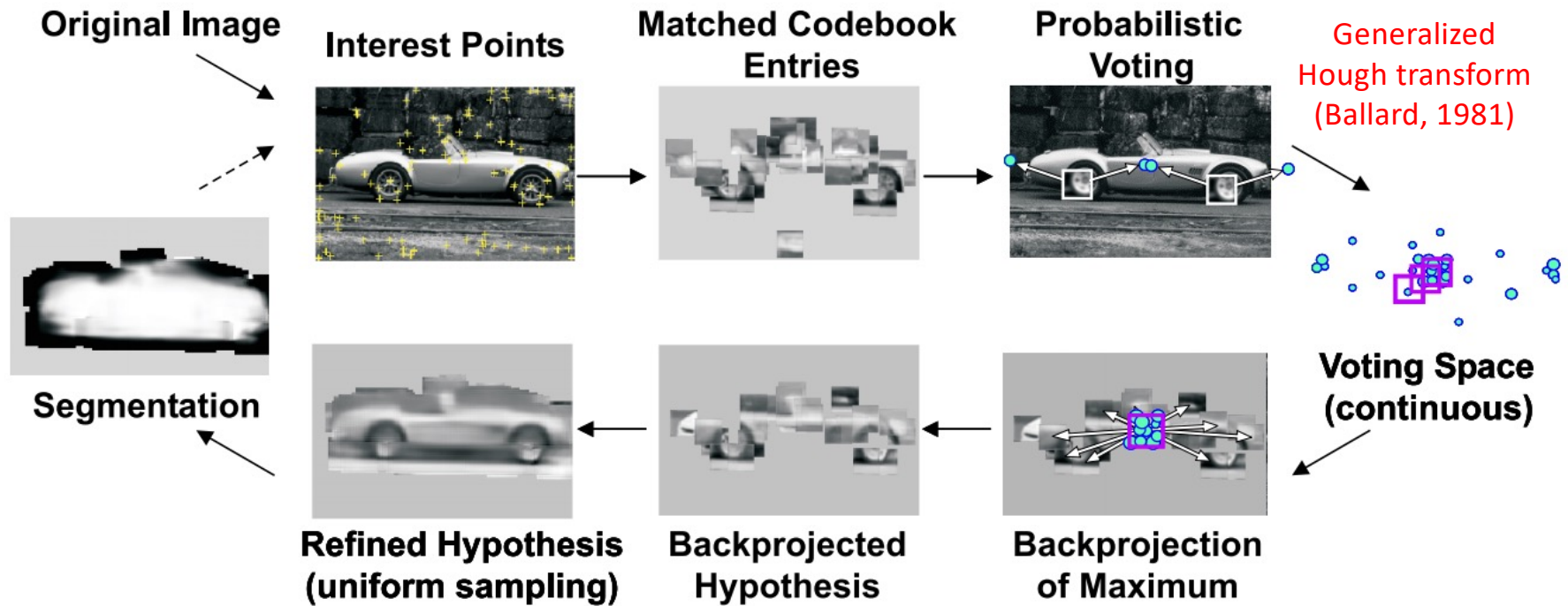
T. Kanade. [Picture Processing System by Computer Complex and Recognition of Human Faces](#). Ph.D. dissertation 1973

Constellation models



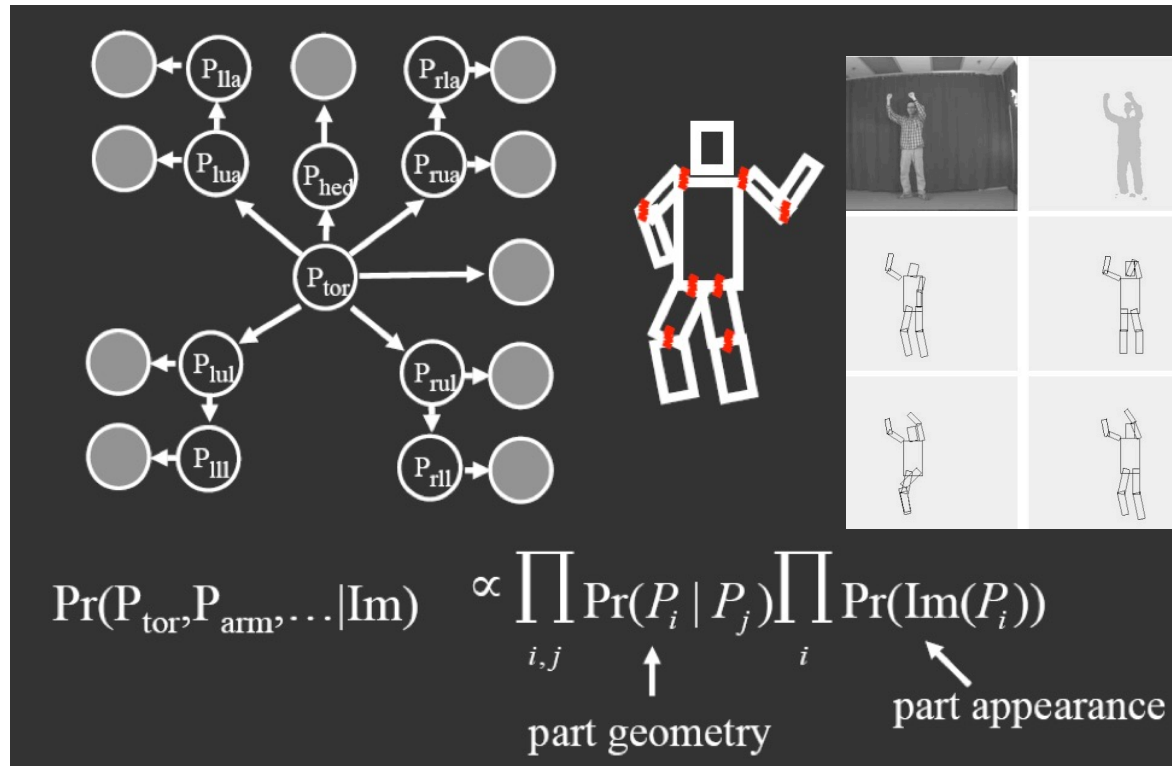
Burl, Weber, & Perona (1998); Weber, Welling & Perona (2000); Fergus, Perona & Zisserman (2003)

Implicit shape models



B. Leibe, A. Leonardis, and B. Schiele, [Combined Object Categorization and Segmentation with an Implicit Shape Model](#), ECCV Workshop on Statistical Learning in Computer Vision, 2004

Pictorial structures revived



P. Felzenszwalb and D. Huttenlocher, [Efficient matching of pictorial structures](#), CVPR 2000

P. Felzenszwalb and D. Huttenlocher, [Pictorial structures for object recognition](#), IJCV 2005

Last time: Object and scene representations

- Object representations
 - 3D shape
 - 3D primitives
 - 2D appearance-based models
 - 2D part-based models (deformable templates)
 - CNNs
- Scene representations
 - Structured representations
 - Appearance-based representations
 - Bottom-up and top-down perceptual organization
- Trends

Finding people used to be *really hard!*

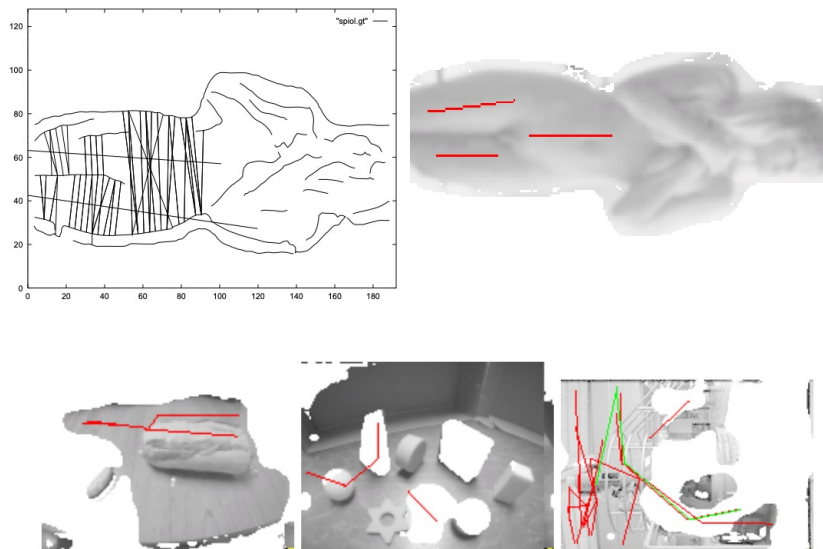


Fig. 6. Typical control images wrongly classified as containing naked people. These images contain people or skin-colored material (animal skin, wood, bread, off-white walls) and structures which the geometric grouper mistakes for spines or girdles. The grouper is frequently confused by groups of parallel edges, as in the industrial image.

[Fleck, Forsyth & Bregler \(1996\)](#)

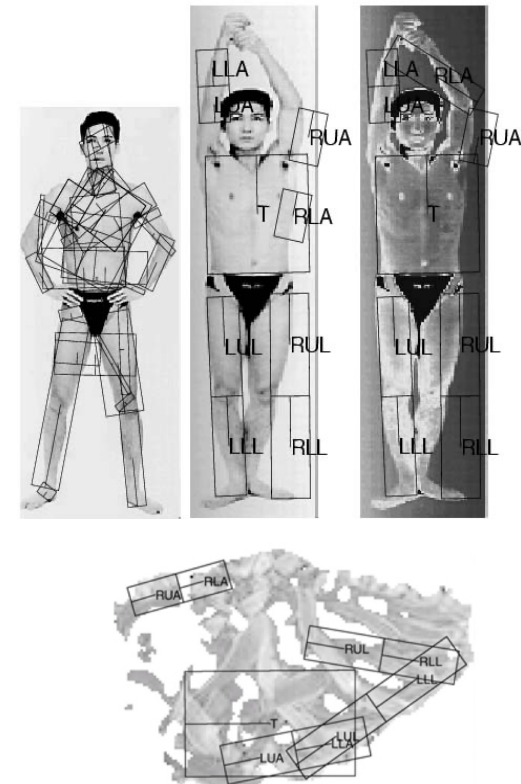
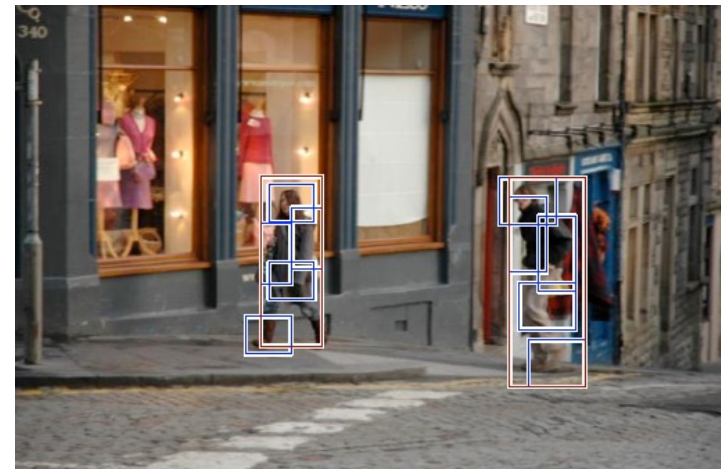
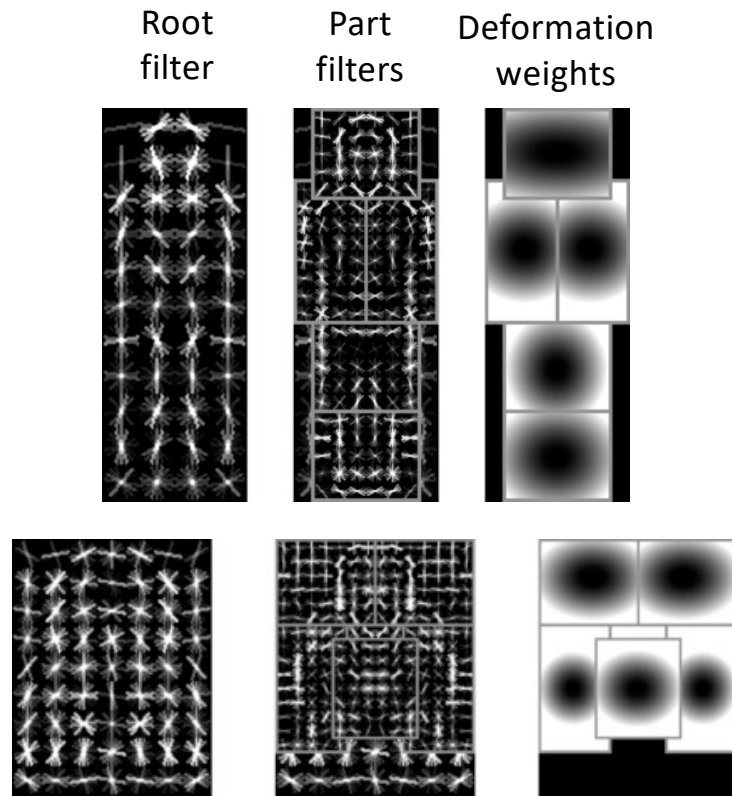


Figure 17. A negative image for which a human assembly was found. The assembly indeed looks like a configuration of a person. A better segment finder would not produce these segments and thus a person would not be detected. The white regions in the image are the pixels that have been masked out because they could not belong to a person due to their color.

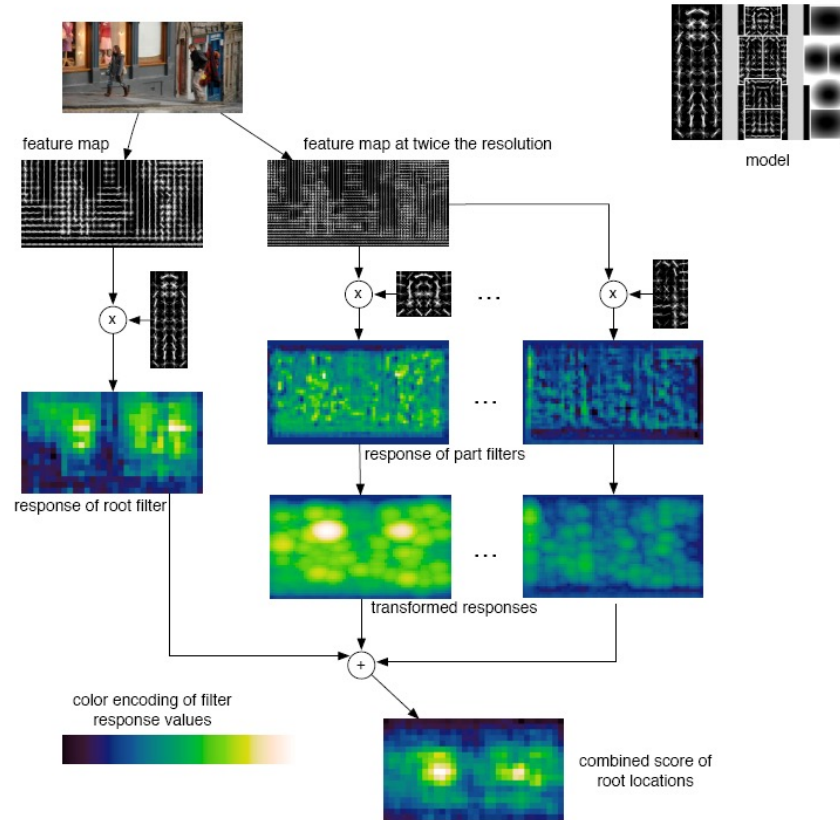
[Ioffe & Forsyth \(2001\)](#)

Discriminative deformable part models



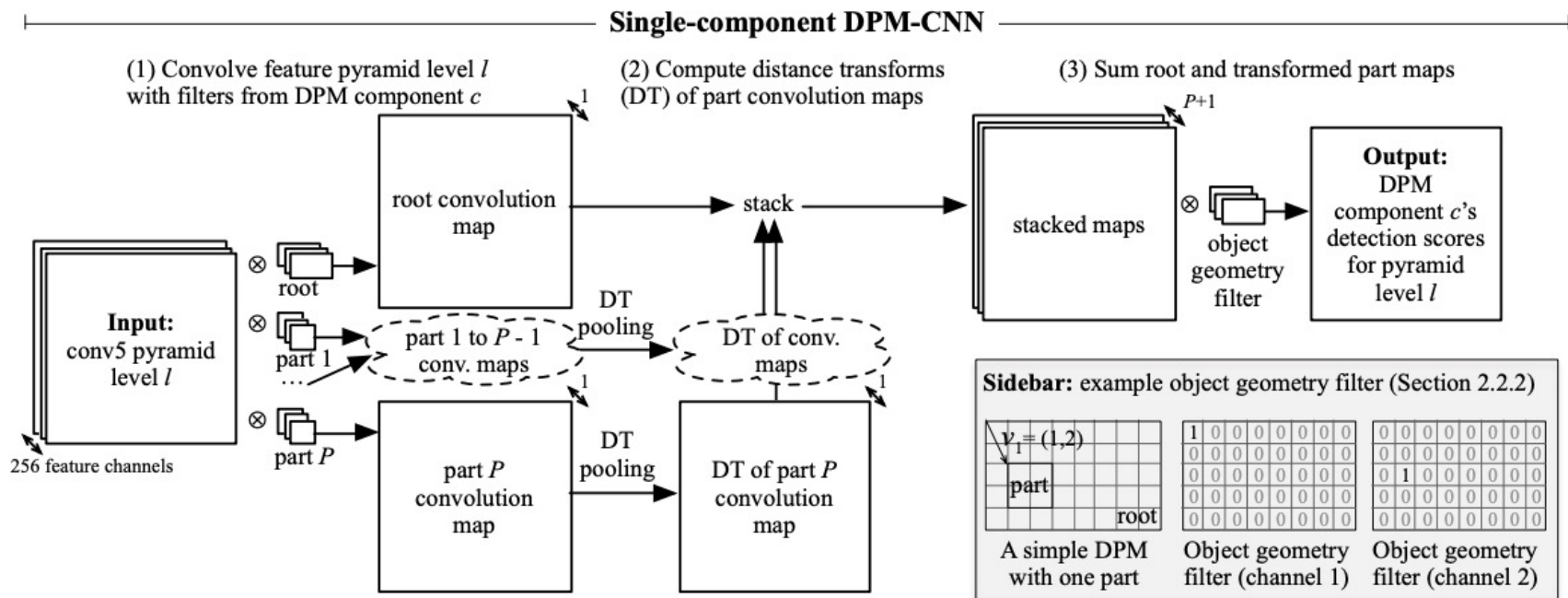
P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, [Object Detection with Discriminatively Trained Part-Based Models](#), PAMI 2009

Discriminative deformable part models



P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, [Object Detection with Discriminatively Trained Part-Based Models](#), PAMI 2009

Deformable part models as CNNs



R. Girshick, F. Iandola, T. Darrell, and J. Malik, [Deformable Part Models are Convolutional Neural Networks](#), CVPR 2015

CNNs as deformable part models

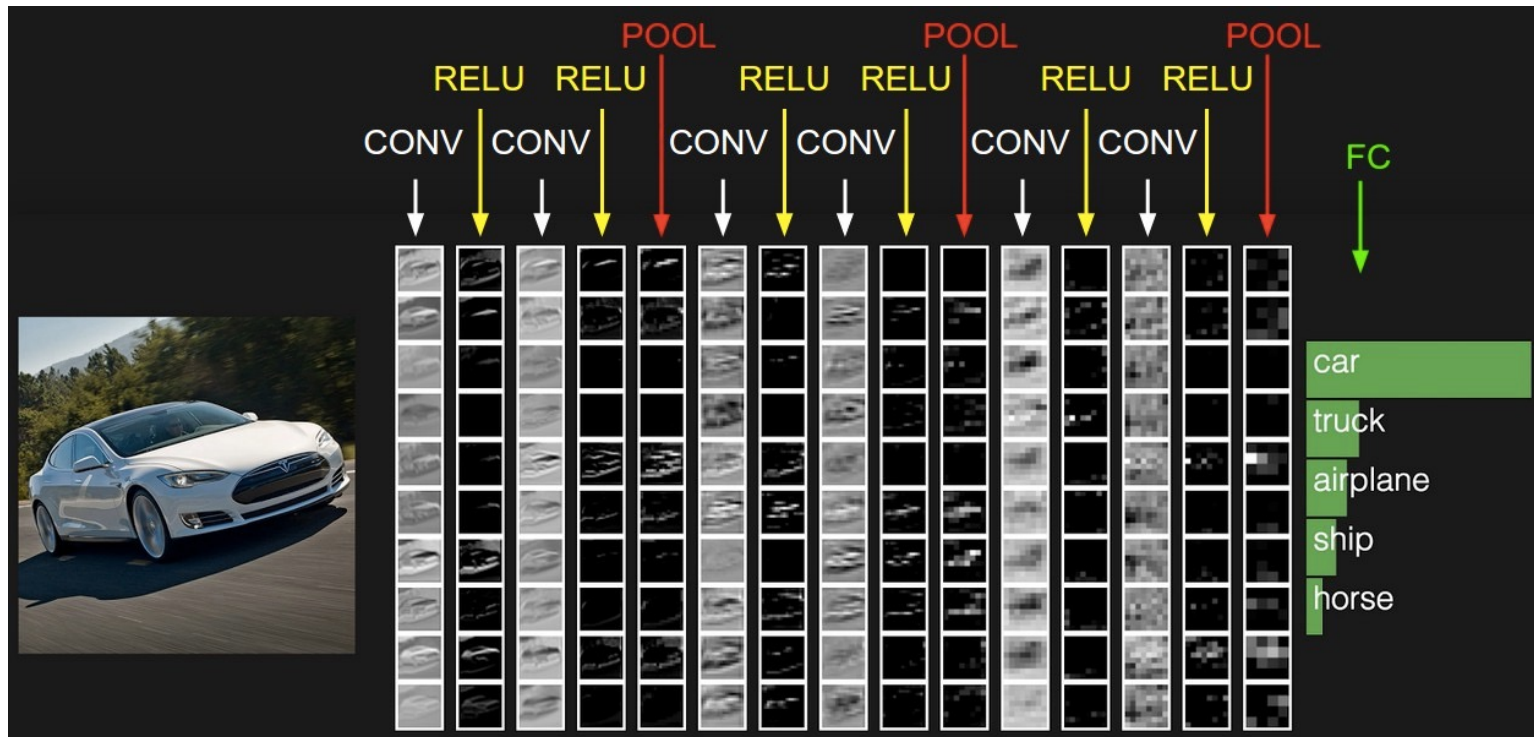
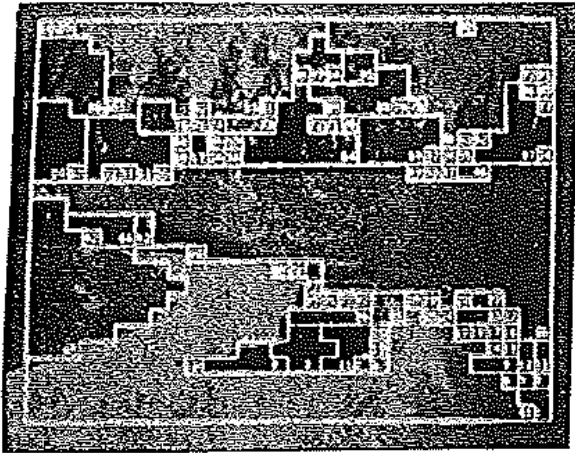


Image source: A. Karpathy

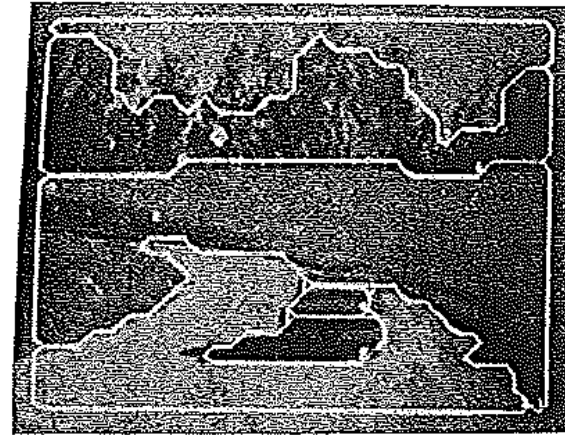
Outline

- Object representations
 - 3D shape
 - 3D primitives
 - 2D appearance-based models
 - 2D part-based models (deformable templates)
 - CNNs
- Scene representations
 - Structured representations
 - Appearance-based representations
 - Bottom-up and top-down perceptual organization

Structured scene representations



(B-2) Output of the non-semantic weakest boundary melted first region grower.

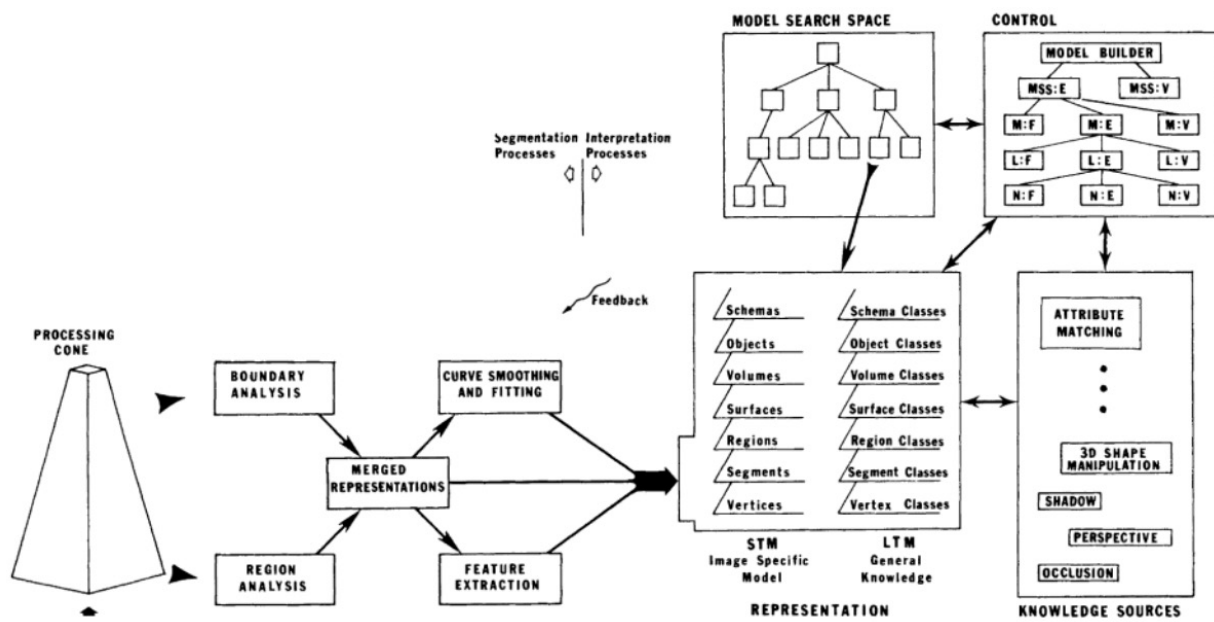


(B-5) Grouping regions by their assigned meaning, all regions considered mergable.

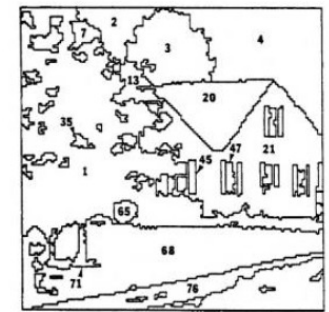
Slide credit:
A. Efros

- Approach has everything: color, super-pixels, bottom-up segmentation, top-down parsing, inter- and intra-region reasoning, Bayesian formulation!

Structured scene representations



(a)



(b)

	Summary of Identification Accuracy					
	All Regions		Regions in Which at Least One 5x4 Window Will Fit		Large Regions (area ≥ 65 pixels)	
	All 5 Objects	4 Objects After Collapsing Bush-Tree	All 5 Objects	4 Objects	5 Objects	4 Objects
Total Number of Regions	209		83		45	
Number of Target Regions	99		50		25	
Number of Non-Target Regions	110		33		20	
Number of Target Regions Correctly Identified	76	91	40	45	19	23
Number of Target Regions Incorrectly Identified	23	8	10	5	6	2
% Target Regions Correct	76.7	91.9	80.	90.	76.	92.

A. Hanson and E. Riseman, VISIONS: A computer system for interpreting scenes, Computer Vision Systems, 1978

Structured scene representations

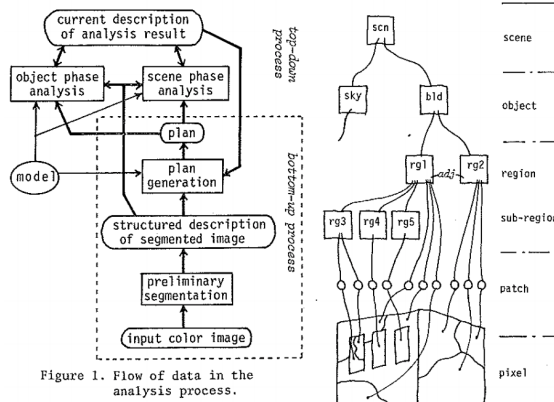


Figure 1. Flow of data in the analysis process.

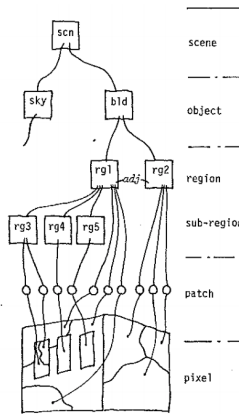


Figure 3. Structure of description built by the system.

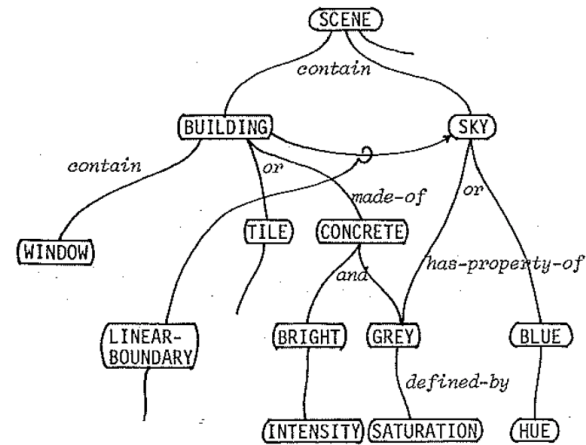


Figure 2. Semantic network for knowledge representation.

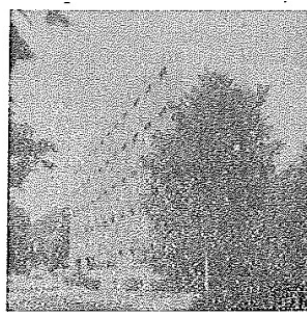
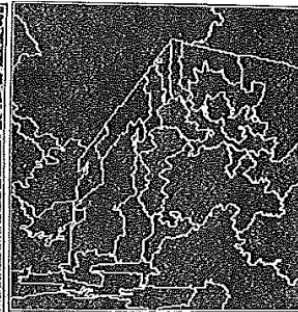


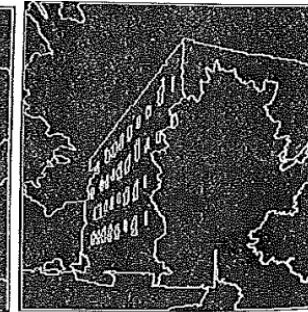
Figure 5-a. Digitized color scene.



5-b. Result of preliminary segmentation.

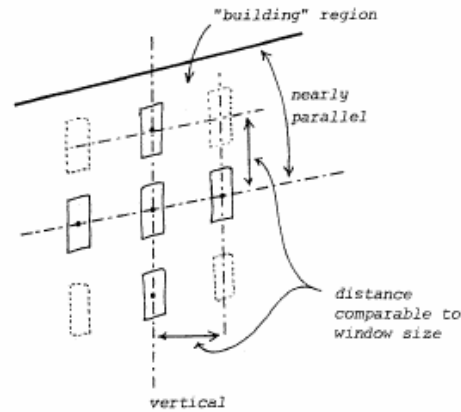


5-c. Plan image.



5-d. Result of semantic segmentation.

What went wrong?



(a) "windows" and "building"

```
[(ACT (IF (AND (IS-PLAN *PCH *MRGN) ..... (1)
              (*VERTICALLY-LONG *PCH))
          (THEN (GET-SET *PLSET (PLAN *MRGN) PATCHES) ..... (2)
                (AND (ALL-FETCH *WLIKE *PLSET ..... (3)
                      (AND (IS (LABEL *WLIKE) NIL)
                           (*VERTICALLY-LONG *WLIKE)))
                      (ALL-FETCH *WIND *WLIKE ..... (4)
                      (THERE-IS *WK *WLIKE
                       (*W-RELATION *WIND *WK))))))
          (THEN (CONCLUDE P-LABEL B-WINDOW)
                (FOR-EACH *WIND (AND (MUST-BE *WIND P-LABEL B-WINDOW)
                                     (DONE-FOR *WIND)))
                (SCORE-IS (ADD 2.1 (DIV (NUMBER-OF *WIND) 100.0))))
                (*PCH *MRGN])
```

(b) listing of the to-do rule for "windows" detection

Slide credit:
A. Efros

Ohta & Kanade (1978)

What went wrong?

Appendix-B Complete Listing of the Model

```

aSCENE knowledge-block-of-scene
(
  OBJECTS (aSKY aTREE aBUILDING aROAD aWINDOW)
  SUB-OBJECTS (aB-WINDOW aCAR aSHADOW)
  KEY-PATCH-IS (GREATERP (AREA aPOCH 300) aPOCH)
  PLAN-PHASE-GENERATION (aTV (BOUNDARY-LENGTH aPOCH aPOCH)
    (MULT (aG-B-DIFFERENCE aPOCH aPOCH)
      (BOUNDARY-CONTRAST aPOCH aPOCH)))
    aPOCH aPOCH)
  IF-PLAN-IS-MODIFIED (IF-DONE (
    rule-for-iso-detection
    (
      (FACT (IF (IS (OF HORIZON SCENE) NIL)
        (ALL-FETCH aPOCH aPLAN-REGIONS
          (IF (AND (NOT (PROBABLY ROAD aPOCH))
            (NOT (TOUCHING aPOCH LSH-SIDE))
            (ALL-FETCH aPOCH aPLAN-REGIONS
              (IF (AND (MAY-BE ROAD aPOCH)
                (MAYBE aPOCH aPOCH)
                (NOT (aG-COLOR aPOCH aPOCH))
                (FACING HORIZONTALLY aPOCH aPOCH))
              (MULT (SUB (FACING HORIZONTALLY aPOCH aPOCH) 0.5)
                (SUB (MAY-VALUE ROAD aPOCH) 0.5)
                (MAY-VALUE ROAD aPOCH))))))
            (VALUE aPOCH aPOCH)))
          (THEN (MAY (SCENE) ROAD-ZONE
            (WITH (MAY-LSH-SIDE aPOCH) 0.5 1.25))
            (MAY (SCENE) HORIZON (MAY-LSH-SIDE aPOCH))
            (EXECUTE PLAN-EVALUATOR)))
          ))
        )
    P-SELECT (TO-DO (
      rule-for-init-a-star
      (
        (FACT (AND (PROBABLY BUILDING aPOCH (NOTFOUND ROAD)))
          rule-for-tree-occlusion
          (
            (FACT (AND (aCAR aPOCH (MAYPER aPOCH)
              (OR (TOUCHING aPOCH UP-SIDE) (TOUCHING aPOCH SIDE))
              (THERE-IS aTR aREGIONS
                (AND (IS (LABEL aTR) TREE)
                  (MAYBE aPOCH aTR)
                  (TOUCHING aTR SIDE)
                  (aWIDTH-INZ aPOCH (Y-ZONE aTR))))))
              (THEN (CONCLUDE P-LABEL TREE)
                (CONCLUDE D-PERSE (WITH (OCCLUDE aTR FRAME))
                  (SCORE-IS 1.0))) aPOCH))
            rule-for-tree-garbage
            (
              (FACT (PROBABLY TREE aPOCH)
                (THEN (CONCLUDE P-LABEL TREE)
                  (SCORE-IS (MAY-VALUE TREE aPOCH))) aPOCH)))
            P-LABEL (IF-DONE (
              IF-some-rule-to-be-act-ivated-when-keypatch-is-labeled
              (
                (FACT (AND (IS (OF PLAN aPOCH) NIL))
                  (THEN (EXECUTE PLAN-EVALUATOR))) aPOCH)))
            )
          )
        aSKY knowledge-block-of-sky
        (PROPERTY-RULES (
          (GEN (NOT (aCLOUDER aPOCH)) (1.0 , 0.5)) aPOCH)
          (GEN (aSHINING aPOCH)) (1.0 , 0.2) aPOCH)
          (GEN (OR (aBLUE aPOCH (aCOPPER aPOCH)) (1.0 , 0.2)) aPOCH)
          (GEN (NOT (aTEXTURAL aPOCH)) (1.0 , 0.7)) aPOCH)
          (STR (TOUCHING aPOCH UP-SIDE) (0.7 , 0.2)) aPOCH))
        RELATION-RULES (
          (STR (AND (aLINEAR-BOUNDARY aPOCH aPOCH)
            (IF (aLINEAR-BOUNDARY (POSITION (GM aPOCH aPOCH))
              (0.0 , 0.5) FOR SKY) aPOCH aPOCH))
            (STR (IF (NOT (IS (OF BUILDING-ZONE SCENE)) NIL)
              (aSUZZY (a-RATIO aPOCH (OF BUILDING-ZONE SCENE)) (0.5 0.3)
                (0.0 , 0.5) FOR SCENE) aPOCH)))
        )
      )
    )
  )

```

```

(THEN (CONCLUDE P-LABEL BUILDING)
  (SCORE-IS (ADD 4.0 (CONFIDENCE-VALUE aPOCH)))) aPOCH)
(FACT (AND (PROBABLY ROAD aPOCH (NOTFOUND ROAD))
  (THEN (CONCLUDE P-LABEL ROAD)
    (SCORE-IS (ADD 4.0 (CONFIDENCE-VALUE aPOCH)))) aPOCH)
(FACT (AND (PROBABLY SKY aPOCH (NOTFOUND SKY))
  (THEN (CONCLUDE P-LABEL SKY)
    (SCORE-IS (ADD 4.0 (CONFIDENCE-VALUE aPOCH)))) aPOCH)
(FACT (AND (PROBABLY TREE aPOCH)
  (NOT (THERE-IS aTR aREGIONS
    (AND (IS (LABEL aTR) TREE)
      (OR (TOUCHING (PLAN aPOCH (PLAN aTR))
        (aWIDTH-INZ (PLAN aPOCH)
          (Y-ZONE aTR))))))))))
  (THEN (CONCLUDE P-LABEL TREE)
    (SCORE-IS (ADD 4.0 (CONFIDENCE-VALUE aPOCH)))) aPOCH)
rule-for-adjacent-wall-of-building
(FACT (AND (MAY-BE BUILDING aPOCH)
  (THERE-IS aWL aREGIONS
    (AND (IS (LABEL aWL) BUILDING)
      (NOT (IS (OF SHARP YEN (OBJECT aWL)))))
      (IS (OF ADJACENT (OBJECT aWL)))))
  (DIFFERENT-ZONE aPOCH aWL)))
(THEN (CONCLUDE P-LABEL BUILDING)
  (CONCLUDE D-PERSE (WITH ADJACENT aWL)
    (SCORE-IS (ADD 5.0 (MAY-VALUE BUILDING aPOCH)))) aPOCH)
rule-for-building-occlusion
(FACT (AND (MAY-BE BUILDING aPOCH)
  (THERE-IS aWL aREGIONS
    (AND (IS (LABEL aWL) BUILDING)
      (aSAME-ZONE aPOCH aWL)
      (aG-COLOR aPOCH aWL)
      (THERE-IS aTR aPATCHES
        (OR (IS (LABEL aTR) TREE)
          (AND (IS (LABEL aTR) BUILDING)
            (NOT (IS (OBJECT aWL)
              (OBJECT aTR))))))))))
  (THEN (CONCLUDE P-LABEL BUILDING)
    (CONCLUDE D-PERSE (WITH (OCCLUDE aWL (REGION aTR)))
      (SCORE-IS (ADD 1.0 (MAY-VALUE BUILDING aPOCH)))) aPOCH)
P-SELECT (
  TO-DO (
    (FACT (MAY-BE SKY aPOCH)
      (THEN (SCORE-IS (ADD 2.0 (MAY-VALUE SKY aPOCH)))) aPOCH)
    (FACT (AND (IS-PLAN aPOCH aPOCH) (aRIGHT aPOCH)
      (THEN (SCORE-IS 2.0)) aPOCH aPOCH)
    (FACT (aRIGHT aPOCH (THEN (SCORE-IS 0.5))) aPOCH))
  IF-DONE (
    (FACT aTA (THEN (CONCLUDE P-LABEL SKY)
      (CONCLUDE R-PERSE (MASTER aPOCH))) aPOCH))
  aPRIORI-VALUE-IS 0.1
  aTREE knowledge-block-of-tree
  (MADE-OF (aLEAVES)
    PROPERTY-RULES (
      (GEN (aTIDULE aPOCH) (0.0 , 0.3)) aPOCH)
      (STR (aHEAVY-TEXTURE aPOCH) (0.0 , 0.2)) aPOCH))
  P-SELECT (
    TO-DO (
      (FACT (MAY-BE TREE aPOCH)
        (THEN (SCORE-IS (ADD 2.0 (MAY-VALUE TREE aPOCH)))) aPOCH)
      (FACT (AND (IS-PLAN aPOCH aPOCH) (NOT (aSHINING aPOCH))
        (THEN (SCORE-IS 3.0)) aPOCH aPOCH)
      IF-DONE (
        (FACT aTA (THEN (CONCLUDE P-LABEL TREE)
          (CONCLUDE P-MERGE (MASTER aPOCH))) aPOCH))
      aPRIORI-VALUE-IS 0.2
  aROAD knowledge-block-of-road
  (MADE-OF (OR (aPAVT aCONCRETE)
    SUB-OBJECTS (OR aSHADOW)
    PROPERTY-RULES (
      (GEN (aCLOUDER aPOCH) (0.0 , 0.4)) aPOCH)
      (GEN (aHORIZONTALLY-LONG aPOCH) (0.7 , 0.2)) aPOCH)
      (STR (TOUCHING aPOCH LSH-SIDE) (0.0 , 0.2)) aPOCH))
  RELATION-RULES (
    (STR (AND (aG-COLOR aPOCH aPOCH) (TOUCHING aPOCH aPOCH)
      (0.0 , 0.2) FOR ROAD) aPOCH)
    (GEN (IF (NOT (IS (OF HORIZON SCENE)) NIL)
      (a-RATIO aPOCH (OF ROAD-ZONE SCENE))
        (0.0 , 0.3) FOR SCENE) aPOCH))
  )

```

```

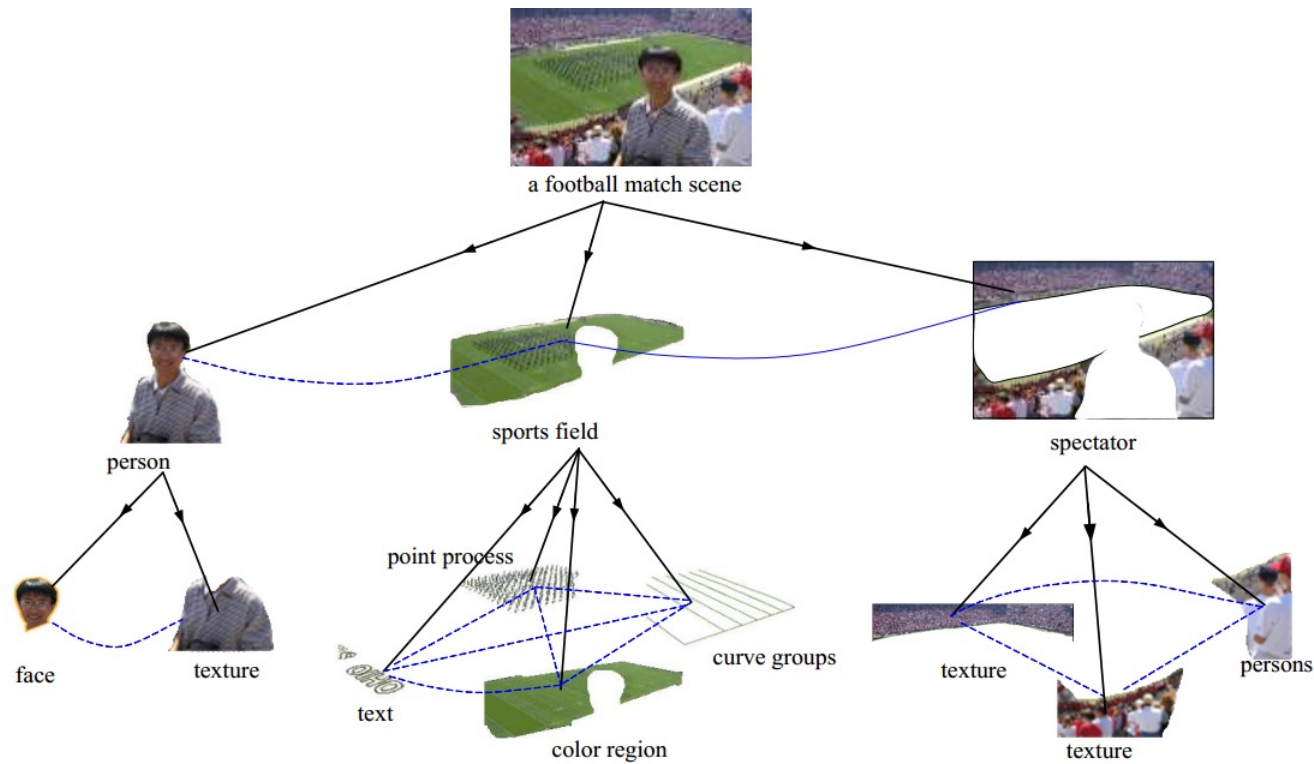
aBUILDING knowledge-block-of-building
(MADE-OF (OR (aCONCRETE aTILE aBRICK)
  SUB-OBJECTS (aSHADOW)
  PROPERTY-RULES (
    (GEN (aTIDULE aPOCH) (0.0 , 0.3)) aPOCH)
    (STR (aPANNHOLE aPOCH) (0.0 , 0.2)) aPOCH)
    (STR (aPANNHOLE aPOCH) (0.0 , 0.2)) aPOCH)
    (GEN (aHORIZON aPOCH) (0.0 , 0.3)) aPOCH))
  RELATION-RULES (
    (GEN (AND (aLINEAR-BOUNDARY aPOCH aPOCH)
      (IF (aLINEAR-BOUNDARY (NOT (POSITION (OF aPOCH aPOCH))
        (0.0 , 0.4) FOR SKY) aPOCH aPOCH))
        (STR (IF (NOT (IS (OF BUILDING-ZONE SCENE)) NIL)
          (AND (a-RATIO aPOCH (OF BUILDING-ZONE SCENE))
            (aPANNHOLE aPOCH))
            (0.0 , 0.3) FOR SCENE) aPOCH))
        )
    )
  P-SELECT (
    TO-DO (
      (FACT (AND (MAY-BE BUILDING aPOCH (aSAME-ZONE aPOCH aPOCH))
        (THEN (CONCLUDE P-LABEL BUILDING)
          (CONCLUDE R-PERSE aPOCH)
          (SCORE-IS (ADD 2.0 (MAY-VALUE BUILDING aPOCH))))
        (FACT (AND (NOT (IS-PLAN aPOCH aPOCH)) (aSAME-ZONE aPOCH aPOCH)
          (MAY-BE BUILDING (PLAN aPOCH))
          (THEN (CONCLUDE P-LABEL BUILDING)
            (CONCLUDE R-PERSE aPOCH)
            (SCORE-IS (ADD 1.0 (MAY-VALUE BUILDING (PLAN aPOCH))))
            (aPOCH aPOCH))
          rule-for-init-a-star
          (
            (FACT (IF (AND (IS-PLAN aPOCH aPOCH) (aSAME-ZONE aPOCH aPOCH)
              (aVERTICALLY-LONG aPOCH (aCONTRACT aPOCH (PLAN aPOCH)))
              (THEN (GET-SET (aSET (PLAN aPOCH) (PATCHES)
                (AND (ALL-FETCH aWL aLSET
                  (AND (IS (LABEL aWL) NIL)
                    (aSAME-ZONE aWL aPOCH)
                    (aVERTICALLY-LONG aWL aPOCH)
                    (aCONTRACT aWL aPLAN aPOCH))))
                )
            )
          )
        )
      )
    )
  )

```

Slide credit:
A. Efros

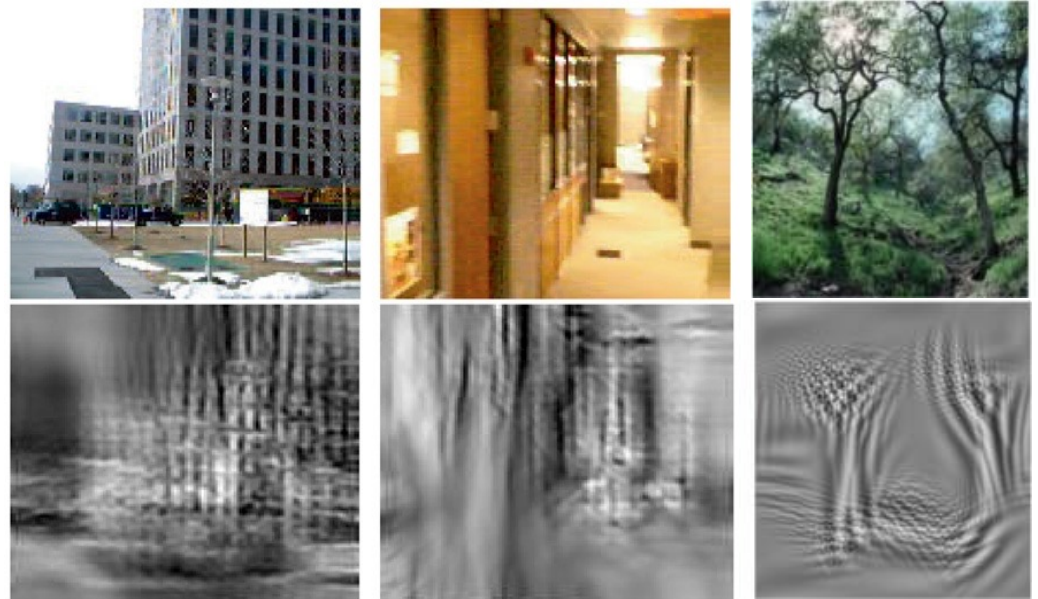
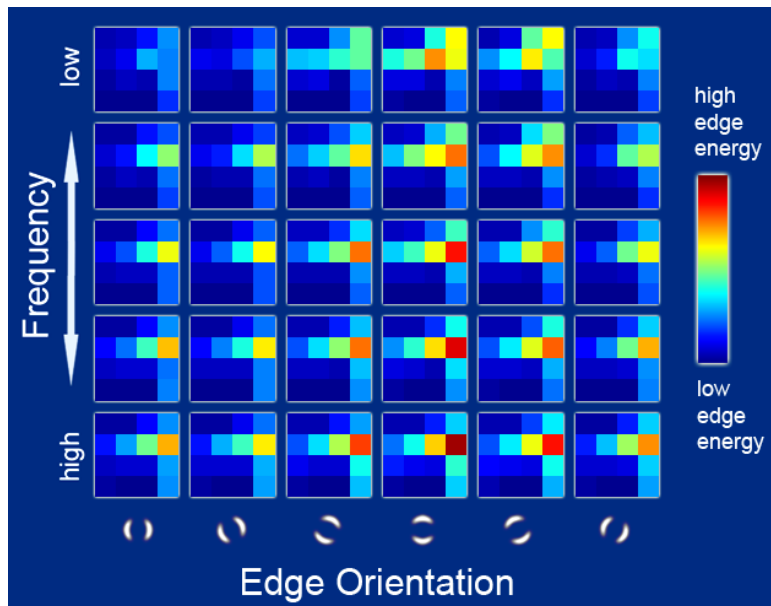
Ohta & Kanade (1978)

Structured scene representations revisited



Z. Tu et al. [Image Parsing: Unifying Segmentation, Detection, and Object Recognition](#), ICCV 2003, IJCV 2005

Appearance-based scene representation: GIST



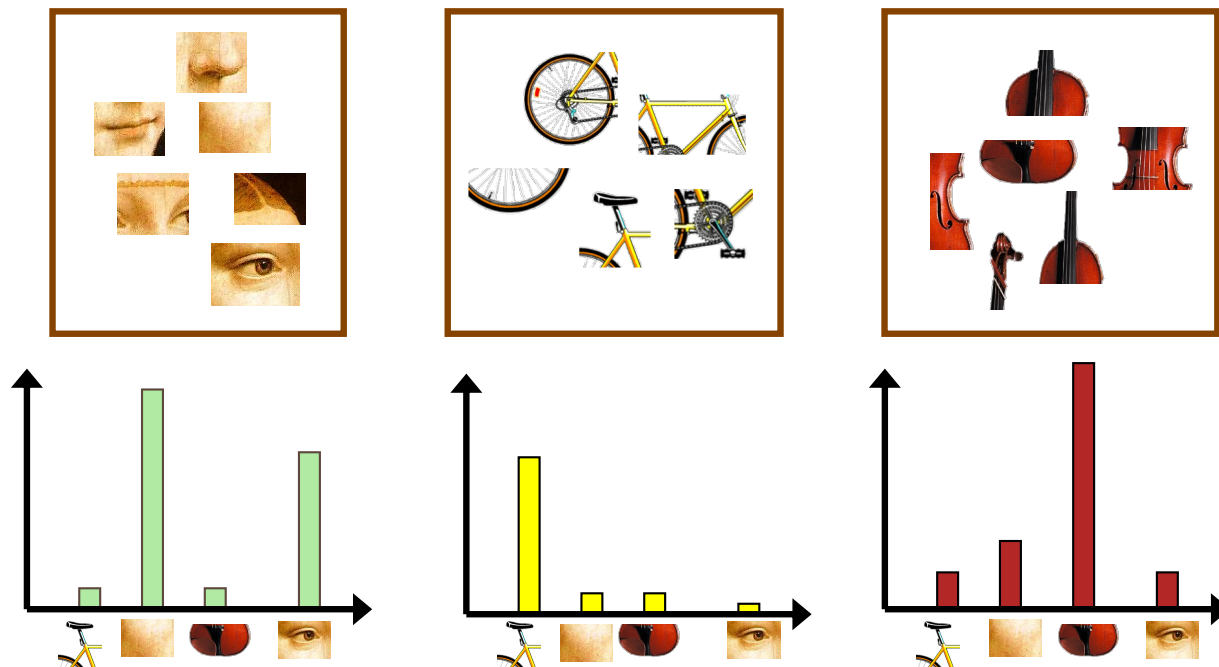
A. Oliva and A. Torralba. [Modeling the shape of the scene: A holistic representation of the spatial envelope](#). IJCV 2001

Appearance-based scene representation: GIST

- Matching of scenes based on GIST works surprisingly well – *given a large enough dataset*

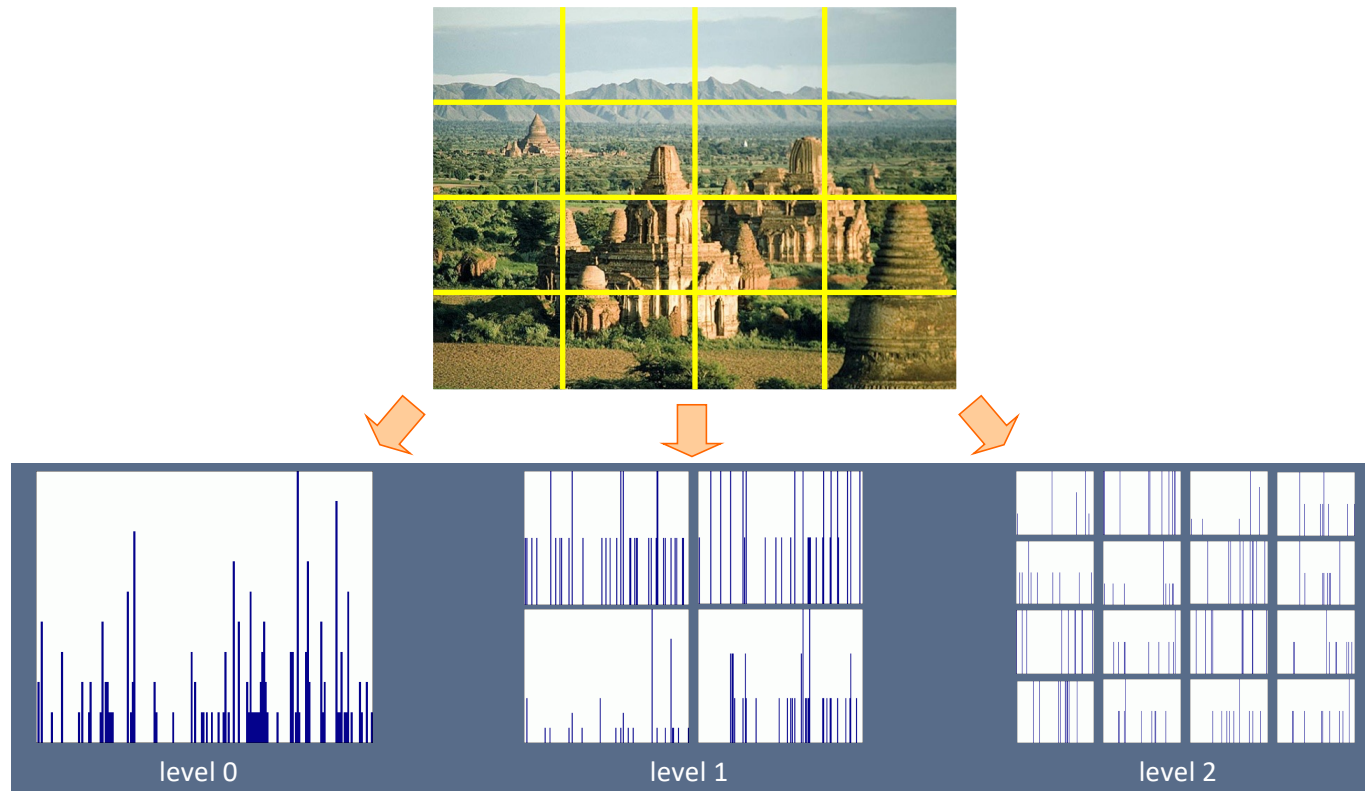


Appearance-based scene representation: Bag of features



Csurka et al. (2004), Willamowski et al. (2005), Grauman & Darrell (2005), Sivic et al. (2003, 2005)

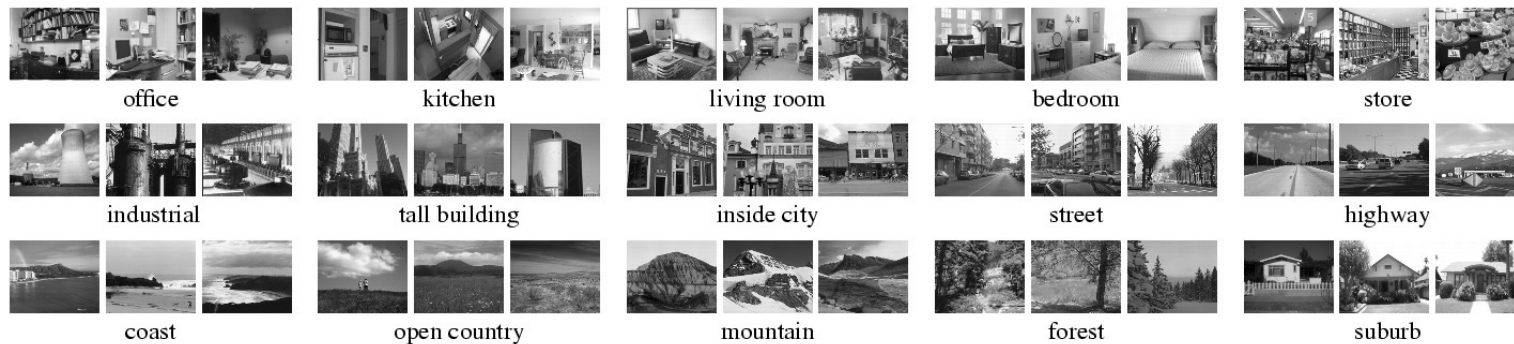
Spatial pyramids



Lazebnik, Schmid & Ponce (2006)

Spatial pyramids

15-category scene dataset

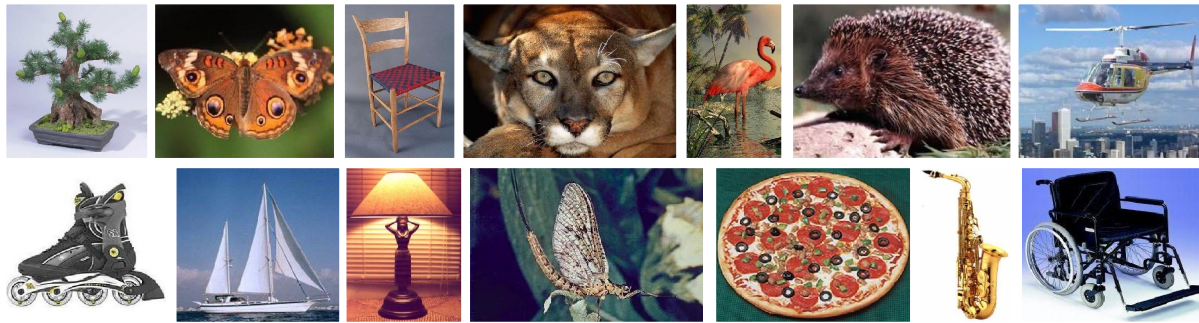


Multi-class classification results
(100 training images per class)

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

Spatial pyramids

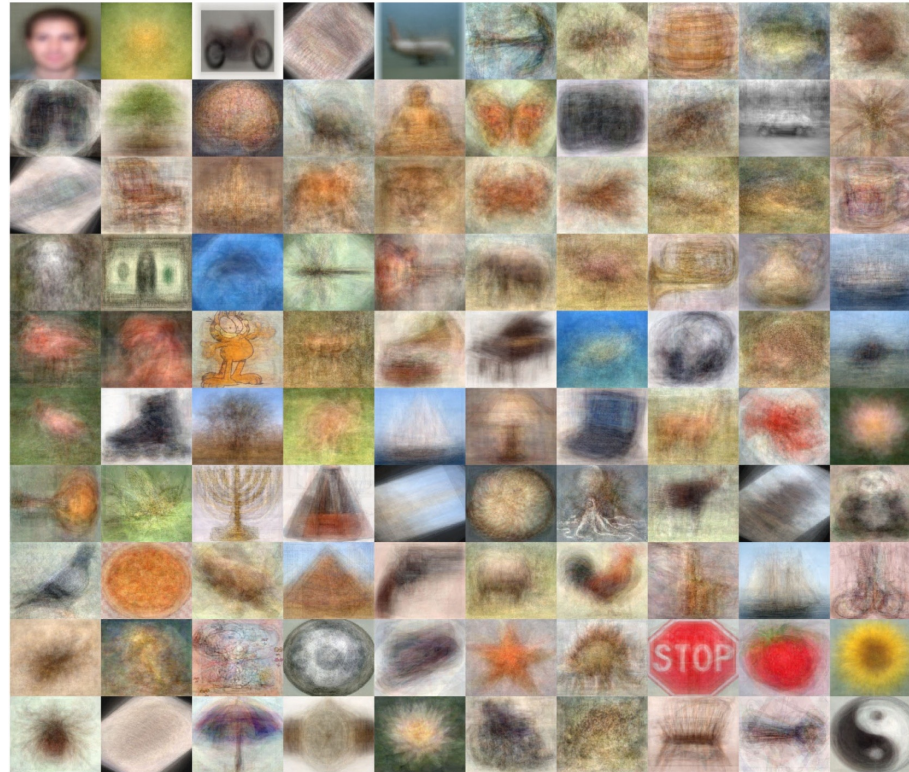
Multi-class classification results (30 training images per class)



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

	Weak features (16)		Strong features (200)	
Level	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

Caltech-101 dataset

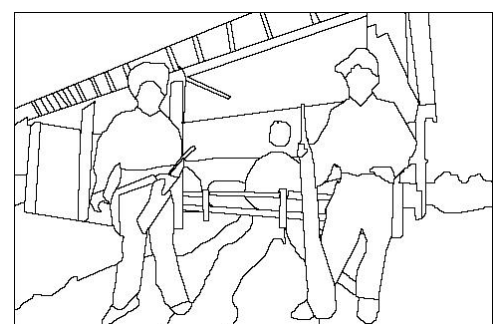
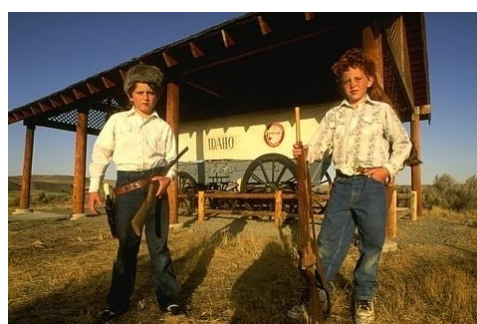
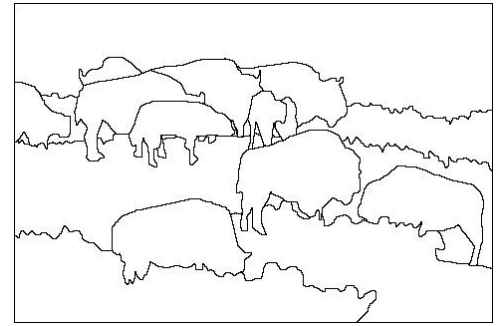


L. Fei-Fei et al. [Learning generative visual models from few training examples: An incremental Bayesian approach tested on 101 object categories](http://www.vision.caltech.edu/Image_Datasets/Caltech101/), CVPR 2004 Workshop on Generative-Model Based Vision
http://www.vision.caltech.edu/Image_Datasets/Caltech101/

Outline

- Object representations
 - 3D shape
 - 3D primitives
 - 2D appearance-based models
 - 2D part-based models (deformable templates)
 - CNNs
- Scene representations
 - Structured representations
 - Appearance-based representations
 - Bottom-up and top-down perceptual organization

Bottom-up perceptual organization



J. Malik et al. [Contour and Texture Analysis for Image Segmentation](#). IJCV 2001

D. Martin et al. [A Database of Human Segmented Natural Images and its Application to Evaluating Segmentation Algorithms and Measuring Ecological Statistics](#). ICCV 2001

Top-down perceptual organization: Semantic segmentation

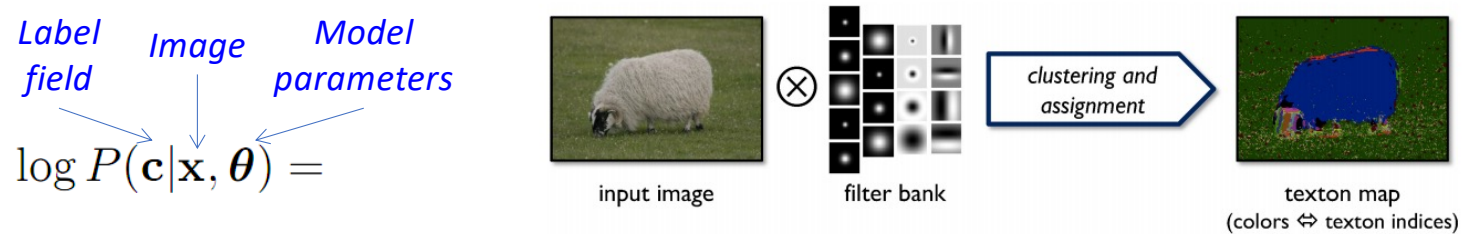
MSRC Dataset (2006)



object classes	building	grass	tree	cow	sheep	sky	airplane	water	face	car
bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat

J. Shotton et al. [TextonBoost: Joint Appearance, Shape And Context Modeling For Multi-class Object Recognition And Segmentation](#). ECCV 2006

Top-down perceptual organization: Semantic segmentation

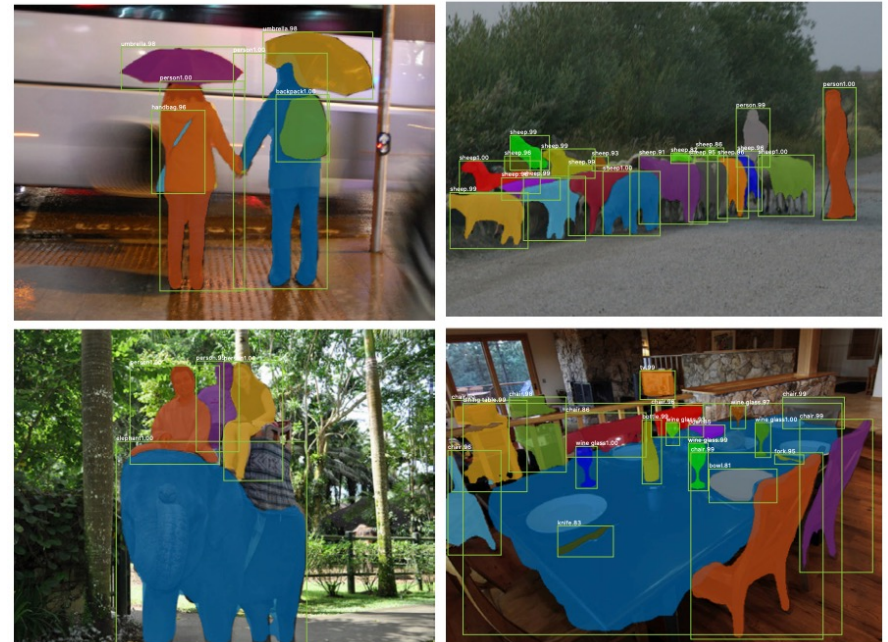
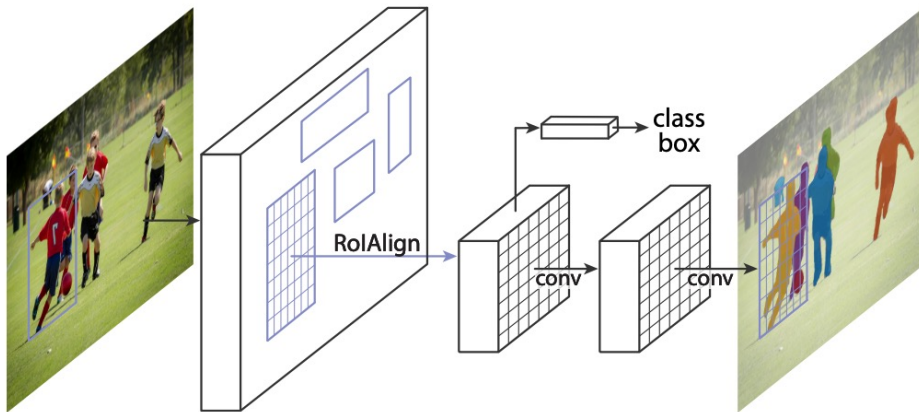


$$\sum_i \overbrace{\psi_i(c_i, \mathbf{x}; \boldsymbol{\theta}_\psi)}^{\text{texture-layout}} + \overbrace{\pi(c_i, x_i; \boldsymbol{\theta}_\pi)}^{\text{color}} + \overbrace{\lambda(c_i, i; \boldsymbol{\theta}_\lambda)}^{\text{location}} \quad \textit{Local data term}$$

$$+ \sum_{(i,j) \in \mathcal{E}} \overbrace{\phi(c_i, c_j, \mathbf{g}_{ij}(\mathbf{x}); \boldsymbol{\theta}_\phi)}^{\text{edge}} \quad \textit{Smoothing term}$$

J. Shotton et al. [TextonBoost: Joint Appearance, Shape And Context Modeling For Multi-class Object Recognition And Segmentation](#). ECCV 2006

Semantic segmentation today: Mask R-CNN

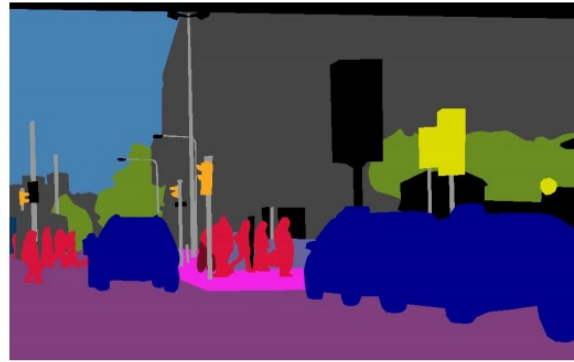


K. He, G. Gkioxari, P. Dollár, and R. Girshick, [Mask R-CNN](#), ICCV 2017 (Best Paper Award)

Panoptic segmentation



(a) image



(b) semantic segmentation



(c) instance segmentation



(d) panoptic segmentation

A. Kirillov et al. [Panoptic segmentation](#). CVPR 2019

Segment Anything Model

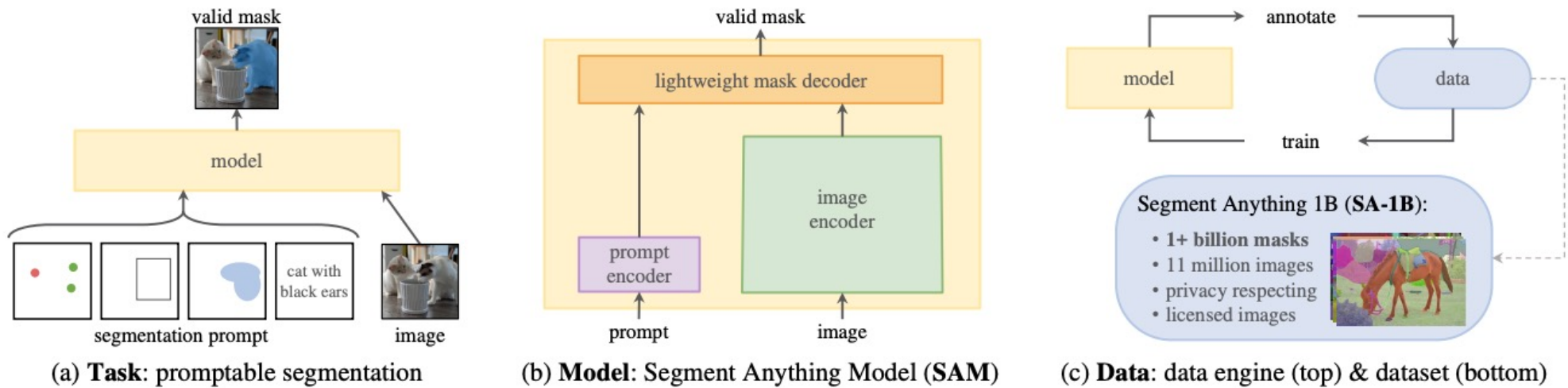


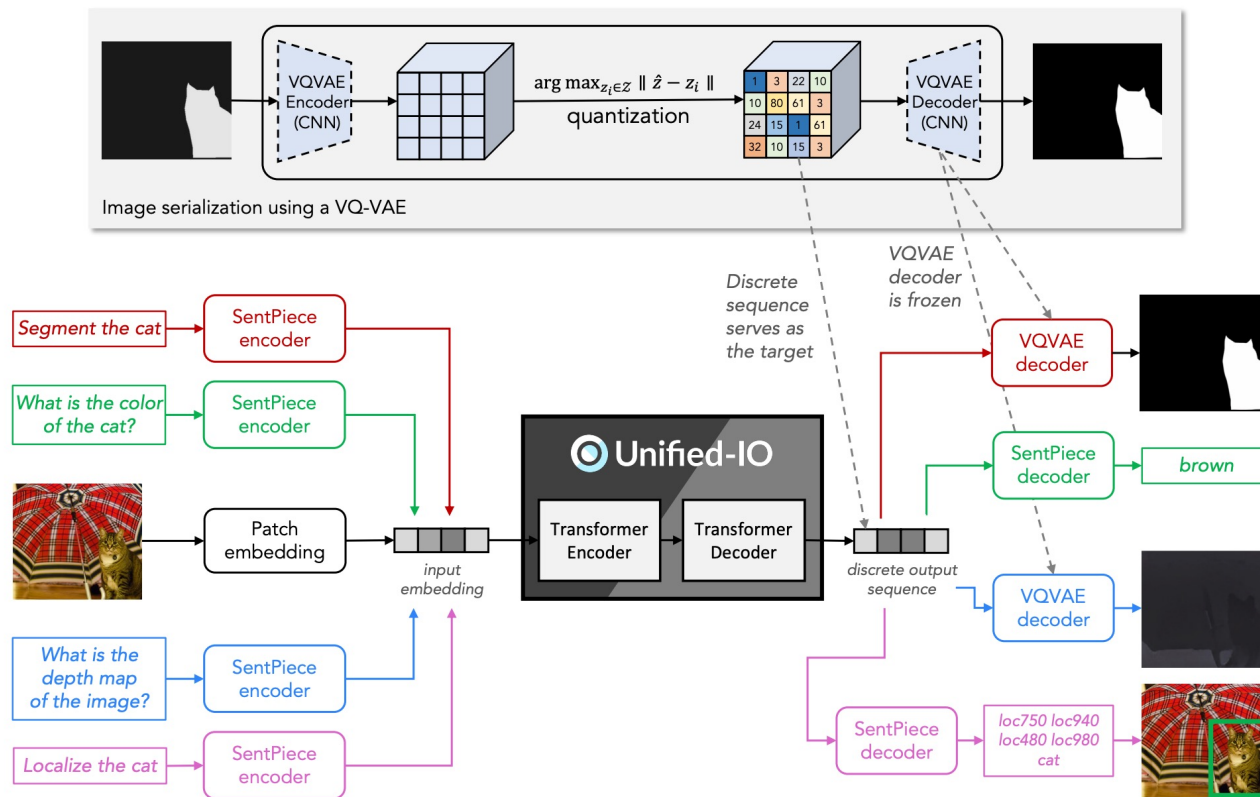
Figure 1: We aim to build a foundation model for segmentation by introducing three interconnected components: a promptable segmentation *task*, a segmentation *model* (SAM) that powers data annotation and enables zero-shot transfer to a range of tasks via prompt engineering, and a *data* engine for collecting SA-1B, our dataset of over 1 billion masks.

A. Kirillov et al. [Segment Anything](https://arxiv.org/abs/2304.02643). arXiv 2023
<https://segment-anything.com/>

Outline

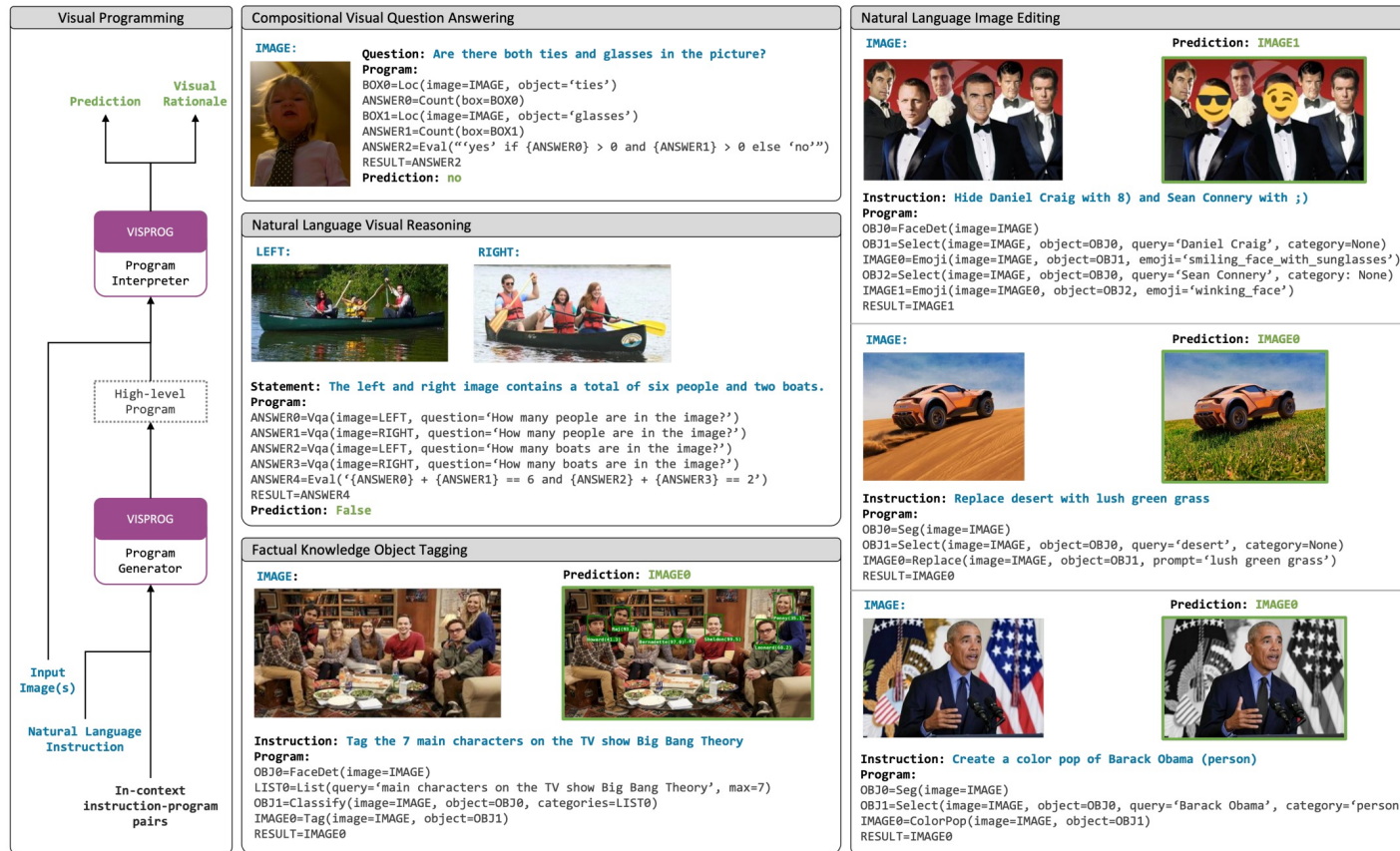
- Object representations
 - 3D shape
 - 3D primitives
 - 2D appearance-based models
 - 2D part-based models (deformable templates)
 - CNNs
- Scene representations
 - Structured representations
 - Appearance-based representations
 - Bottom-up and top-down perceptual organization
- Trends

Unified representations for many tasks



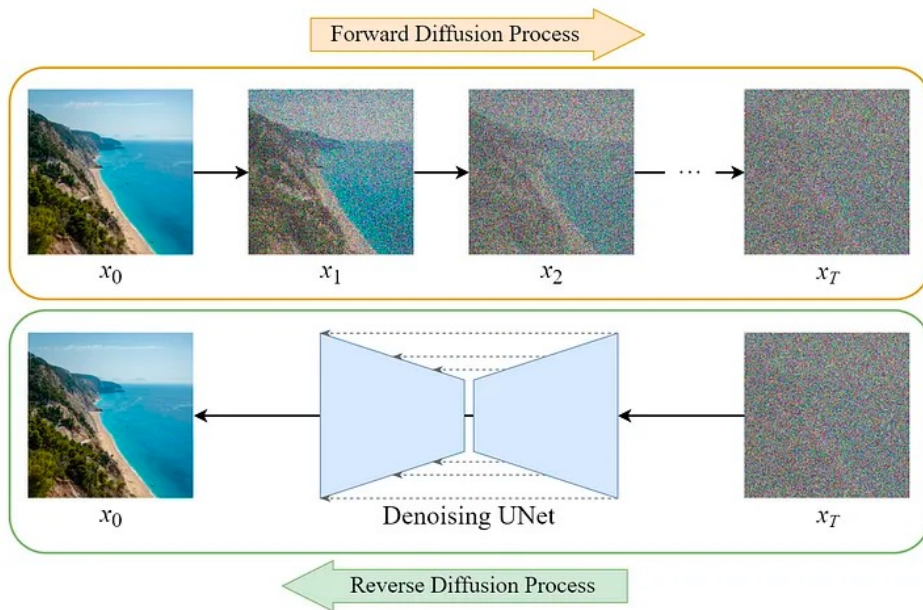
J. Lu et al. [Unified-I/O: A Unified Model for Vision, Language, and Multi-Modal Tasks](#). ICLR 2023

Unified representations for many tasks



Generative models

Diffusion models



[Figure source](#)



Sprouts in the shape of text 'Imagen' coming out of a fairytale book.



A photo of a Shiba Inu dog with a backpack riding a bike. It is wearing sunglasses and a beach hat.



A high contrast portrait of a very happy fuzzy panda dressed as a chef in a high end kitchen making dough. There is a painting of flowers on the wall behind him.



Teddy bears swimming at the Olympics 400m Butterfly event.



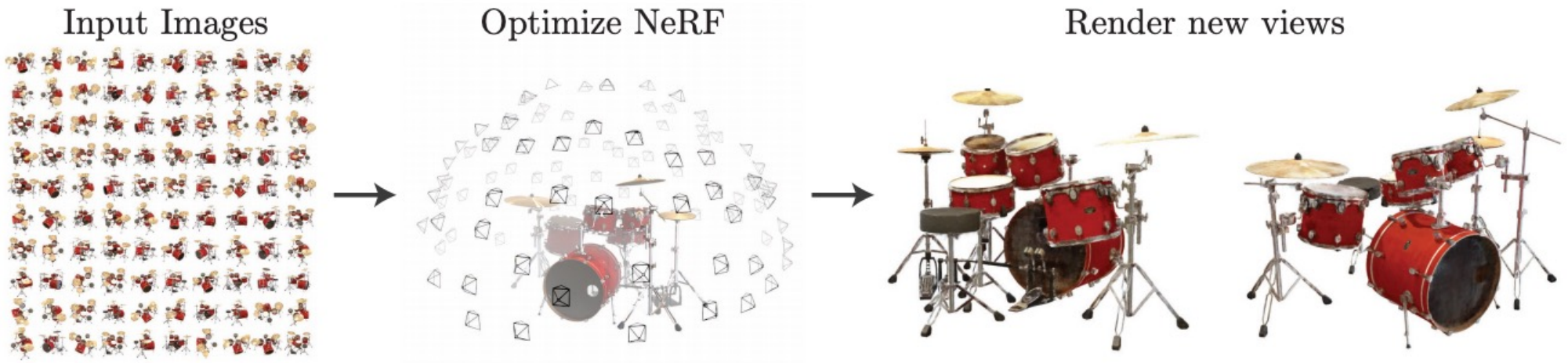
A cute corgi lives in a house made out of sushi.



A cute sloth holding a small treasure chest. A bright golden glow is coming from the chest.

C. Saharia et al. [Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding](#). NeurIPS 2022

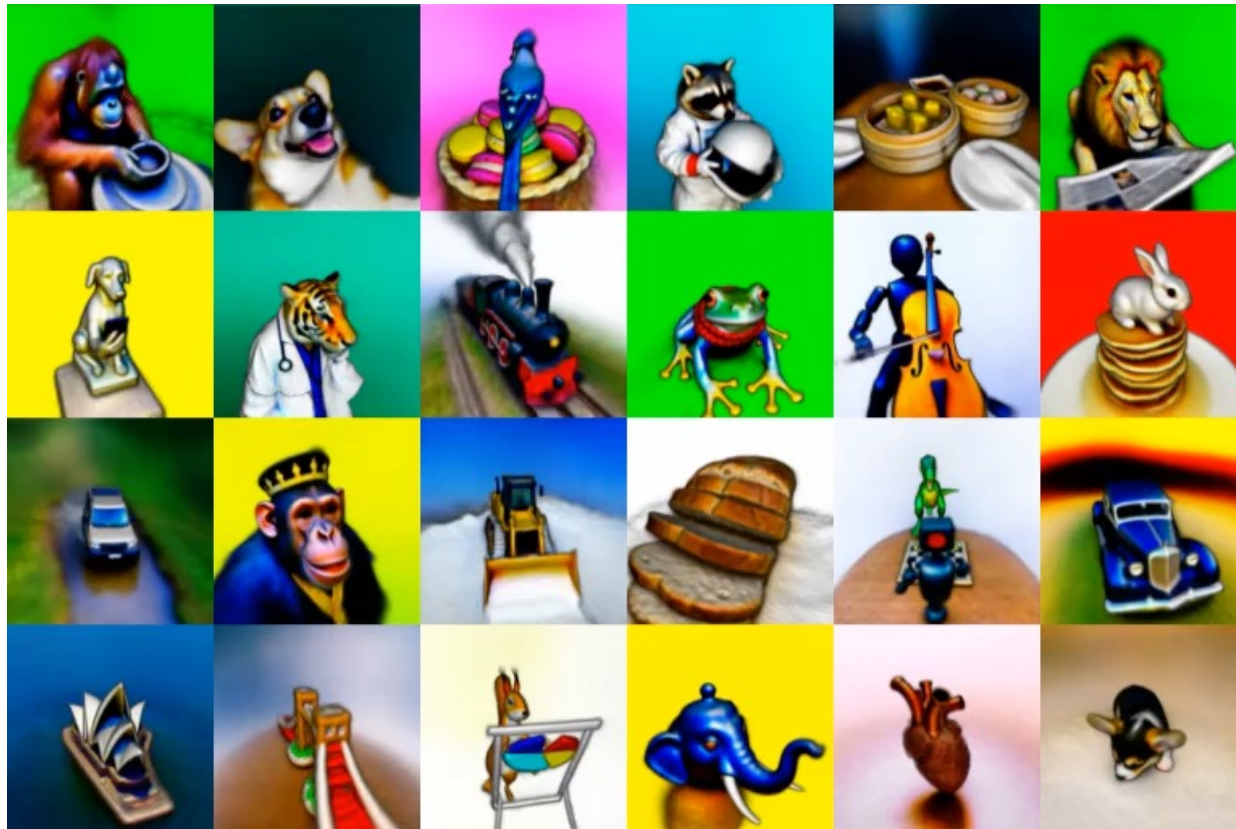
Neural 3D representations: NERFs



$$(x, y, z, \theta, \phi) \rightarrow \begin{matrix} \text{[Neural Network]} \\ F_{\Theta} \end{matrix} \rightarrow (RGB\sigma)$$

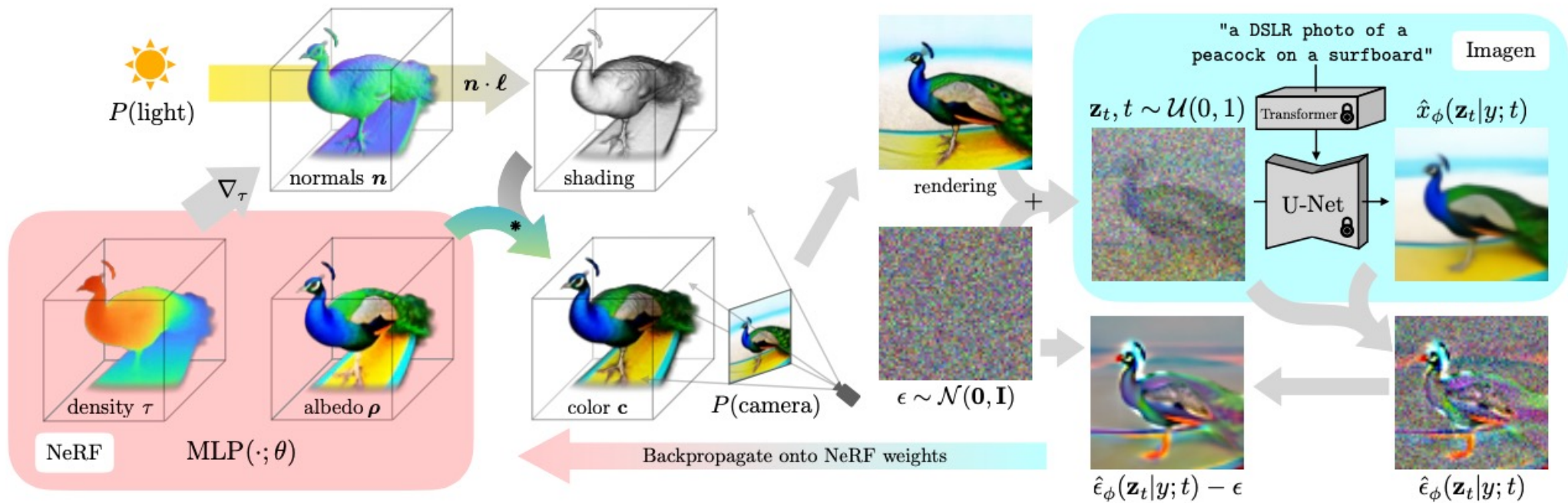
B. Mildenhall et al., [Representing Scenes as Neural Radiance Fields for View Synthesis](#), ECCV 2020

Connecting 2D to 3D: DreamFusion



B. Poole et al. [DreamFusion: Text-to-3D using 2D Diffusion](#). arXiv 2022

Connecting 2D to 3D: DreamFusion



B. Poole et al. [DreamFusion: Text-to-3D using 2D Diffusion](#). arXiv 2022