# Object Recognition

- REPRESENTATION PROBLEM
  2D & 3D OBJECT MODELING
  SHAPE, COLOR, TEXTURE, ...
- MODEL ACQUISITION PROBLEM
- SEARCH PROBLEM
  HOW TO SELECT IMAGE
  FEATURES & INDEX
  INTO DB of ALL OBJECTS
  TO FIND BEST MATCH ?

  ⇒ SELECTION/SEGMENTATION ⟨ FEATURE DETECTION, PERCEPTUAL GROUPING
  INDEXING
  MATCHING (SIMILARITY)
  CORRESPONDENCE
  VERIFICATION



Pictures in Our Heads

Lighting affects appearance



The "Margaret Thatcher Illusion" by Peter Thompson



The "Margaret Thatcher Illusion" by Peter Thompson



a    b    c    d

Figure 29. Four configurations of a nonrigid object.

## Recognition Problems

- What is it?
  - Object detection
- Who is it?
  - Recognizing identity
    - Object recognition
    - Category recognition
- What are they doing?
  - Activity recognition
- All of these are **classification** problems
  - Choose one class from a list of possible candidates

## Face Detection



## Face Recognition



P. Sinha and T. Poggio, Last but not least, *Perception* **31**, 2002, 133.

# What Makes Recognition Hard?

- Intrinsic variability within each class
- Pose variability
- Illumination variability
- Background variability
- Segmentation problem
  - What region within an image contains the object?
- Feature selection problem
  - What features describe shape and appearance?

---

## UNKNOWNS

**DATA SELECTION:**
WHAT SUBSET OF DATA CORRESPONDS TO A SINGLE OBJECT?
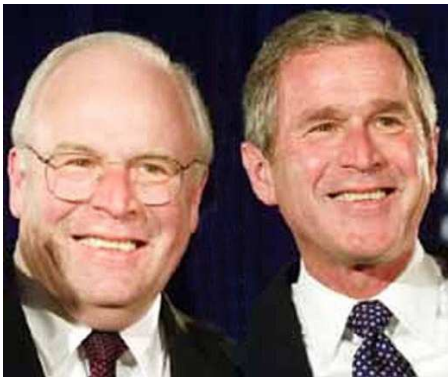
**OBJECT IDENTIFICATION:**
WHICH OBJECT MODEL CORRESPONDS TO DATA SUBSET

**OBJECT INSTANCE:**
FOR NON-RIGID OBJECTS OR OBJECT CLASSES, SPECIFICATION OF THE PARAMETERS THAT COMPLETELY DESCRIBE A GIVEN LEGAL INSTANCE OF OBJECT

**FEATURE CORRESPONDENCE:**
WHICH INDIVIDUAL MODEL FEATURES CORRESPOND TO EACH DATA FEATURE?

**POSE:** POSITION & ORIENTATION OF VIEWER WRT OBJECT IN SCENE

---

## RECOGNITION INVARIANTS

METHOD SHOULD BE **INVARIANT** UNDER

VIEWPOINT INVARIANTS {
- TRANSLATION
- ROTATION
- SCALE
- PARTIAL OCCLUSION (SELF & FROM OTHER OBJECTS)
}

SENSOR INVARIANTS
- SENSOR NOISE
- ILLUMINATION / SHADOWS
- "LOCAL" ERRORS IN EARLY PROCESSING MODULES (E.G., EDGE DETECTION)

OBJECT INVARIANCE
- INTRINSIC SHAPE "DISTORTIONS" (E.G. ARTICULATED OBJECTS)

---

## SEARCH SPACES

**DATA SPACE:** $2^n$, $n = \#$ image features

**OBJECT-MODEL SPACE:** $M$, $\#$ models

**MODEL-PARAMETERIZATION SPACE:** $d^k$
where $k = \#$ parameters
$d = \#$ values per parameter

**CORRESPONDENCE SPACE:** $(m+1)^n$
where $m = \#$ model features
$n = \#$ image features
(assuming all $n$ features part of 1 object)

**POSE SPACE:** (TRANSFORMATION SPACE)
2D objects: 3D or 4D
(2: translation, 1 rotation, 1 scale)

3D Objects:
Orthography: 5D
(2: viewing direction, 1: rotation in image plane, 2: translation)

Perspective: 6D
(3: translation, 3: rotation)

---

4

## Slide 1

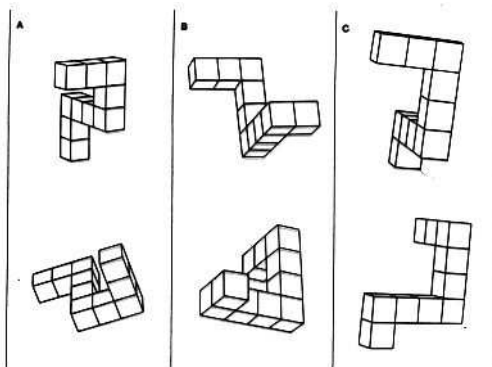MATCHING MODELS w/ IMAGES

2D - 2D MATCH
- SMALL # "CHARACTERISTIC VIEWS"
  - MAY BE USEFUL FOR FAMILIAR OBJECTS W/ FEW "TYPICAL" VIEWS
- VIEWER-CENTERED MODELS

2½D/3D - 3D MATCH
- RECOVER BOTTOM-UP FROM IMAGE(S) 2½D OR 3D DATA
  - E.G. STEREO
    SHAPE-FROM X
    GEONS          [Biederman]
    GENERALIZED CYLINDERS [Marr]

2D - 3D MATCH
- OBJECT-CENTERED MODELS
- VIEWPOINT-INVARIANT FEATURES
  OR VIEWPOINT-VARIANT FEATURES

## Slide 2

METHOD 0: DIRECT APPROACH
- FOR EACH MODEL OBJECT, STORE A SUFFICIENTLY LARGE NUMBER OF DIFFERENT VIEWS
- DEFINE A SIMILARITY MEASURE TO FIND BEST MATCH BETWEEN IMAGE DATA AND DIFFERENT VIEWS OF ALL OBJECT MODELS
  E.G. CROSS-CORRELATION W/ A TEMPLATE
- DISADVANTAGES:
  - VERY LARGE SPACE OF POSSIBLE VIEWS (TRANSFORMATIONS)
  - BRUTE-FORCE SEARCH OF "POSE SPACE"
  - OBJECTS SHOULD BE RECOGNIZED IN NOVEL VIEWS

Ex. 2D position & 2D orientation invariant to 1 pixel and 1° ⟹
$512 \times 512 \times 360 \approx 90{,}000{,}000$ poses

## Slide 3



## Slide 4

METHOD 1: INVARIANT GLOBAL PROPERTIES
- SELECT A SET OF GLOBAL, INVARIANT PROPERTIES — EXPLOITS REGULARITIES ACROSS VIEWPOINTS — VIEWPOINT-INVARIANT FEATURES
  E.G. COMPACTNESS $(P^2/A)$
  AREA (A), MOMENTS $(\sum\sum x^i y^j f(x,y))$
- FEATURES MUST BE RELATIVELY SIMPLE TO COMPUTE ⟹ NEED PROCEDURE FOR DETECTING/COMPUTING EACH SUCH PROPERTY
- IF OBJECT'S PROPERTIES HAVE CHARACTERISTIC RANGES OF VALUES ⟹ POLYTOPE IN n-D FEATURE SPACE
- PATTERN RECOGNITION
  GIBSON
  SRI VISION MODULE
- ARE THERE INVARIANT FEATURES FOR REAL, 3D OBJECTS??

## Slide 1

Given n-D feature vector, $\langle f_1, f_2, \ldots, f_n \rangle$, match w/ N possible objects:



Decision Tree Classifier

```
        f_2 < 2.4  --Yes--> Nut
            |
            No
            |
        f_1 < 3.7  --Yes--> Wrench
            |
            No
            |
        Screwdriver
```

Nearest-Neighbor Classifier

## Slide 2

### A Geometric Invariant

CROSS RATIO

Given 4 points $(x_i, y_i, z_i)$ that are collinear in scene

let $(x_{i_K}, y_{i_K}, z_{i_K})$ be the perspective projection of those 4 points into an image.

$$\text{Cross Ratio} = \frac{(x_3 - x_1)(x_4 - x_2)}{(x_3 - x_2)(x_4 - x_1)}$$

$$= \frac{(x_{i_3} - x_{i_1})(x_{i_4} - x_{i_2})}{(x_{i_3} - x_{i_2})(x_{i_4} - x_{i_1})}$$

and similarly for y's.

$\Rightarrow$ If we know coords of 3 collinear points, we can compute coords of any other point on same line from its image coords

## Slide 3

# Approximate Invariants

- Over a limited range of viewpoint variation
- Parallelism
- Collinearity
- Angle between a pair of lines
- Co-termination

## Slide 4

### MATCHING USING OCCLUDING CONTOURS
Kriegman & Ponce, 1990, PAMI

- Given: Smooth, 3D shape

- Calculate (offline) the implicit equation of the occluding contour
$$F(u, v, V) = 0$$

- (For each aspect)
Solve for pose by minimizing average distance between model contour and image edges

$$\min \sum_i F^2(u_i, v_i, V)$$

where $(u_i, v_i)$ are image edge pts.

(Nonlinear least-squares minimization of an expression for the distance b/w contour & image $\Rightarrow$ Iterative techniques)

Figure 3. The result of minimizing the mean square error of the implicit contour equation: (a). The contours at each iteration of the minimization are shown overlaid on the edge points used in the minimization; the Canny edges are drawn as little circles. (b). The result of recognition for the white torus.

To conclude, let us give an example. Consider a torus. The expression $C$ is too complex to be given here. With the following substitutions: $\hat{x} = \tilde{x}\sin\beta$, $\hat{y} = \tilde{y}\sin\beta$, $c = \cos^2\beta$, the expression $C_1$ is:

$$0 = C_1(\hat{x}, \hat{y}, c) =$$

---

## Method 2: Object Decomposition into Parts
### (+ Correspondence Space Search)

- Decompose object into constituent **parts**, where each part is simple, generic
  - E.g., 1D contours — lines, corners, "codons"
  - 2D surface patches — holes
  - 3D volumes — generalized cylinders "geons"

- Decomposition done independent of viewpoint

- 1. Segment image into parts
  2. Classify parts
  3. Describe object in terms of parts
     - set of features
     - spatial relations between parts

- Examples:
  Perceptron
  Bolles' local-feature focus method

---

# Pose consistency

- Correspondences between image features and model features are not independent
- A small number of correspondences yields a camera --- the others must be consistent with this

- Strategy:
  - Generate hypotheses using small numbers of correspondences (e.g. triples of points for a calibrated perspective camera, etc., etc.)
  - Backproject and verify
- Notice that the main issue here is camera calibration
- Appropriate groups are "frame groups"

---

## Bolles' Local Feature Focus Method

- 2D local features: holes, convex corners, concave corners
- Describe each object as a graph
- For each occurrence of each feature type in each model, select subgraph "centered" at node which is sufficient to distinguish it from all others
- Compile a table for each feature type which says where to look for other nearby features in order to confirm it.

## Ex. FOCUS-FEATURE TABLE

"FOCUS" FEATURE: HOLE (i.e., CIRCLE)

POSSIBLE OBJECT FEATURES: HOLE1 in OBJECT 3
HOLE 4 in OBJECT 5
HOLE 6 in OBJECT 5
⋮

NEARBY FEATURES:

| | | |
|---|---|---|
| CORNER | DISTANCE, ORIENTATION | CORNER 3 in OBJ 2 |
| CORNER | DISTANCE, ORIENTATION | CORNER 1 in OBJ 5 |
| HOLE | DISTANCE | HOLE 1 in OBJ 2 OR HOLE 5 in OBJ 4 OR ... |

---

- BUILD A GRAPH FROM NODE-NODE MATCHES —

  NODE = (MODEL NODE, IMAGE FEATURE) PAIR

  ARC = CONNECTS PAIRS OF NODES WHICH ARE NEIGHBORS IN MODEL *AND* IMAGE FEATURES ARE DISTINCT
  (⇒ LOCALLY CONSISTENT (PAIRWISE) ASSIGNMENT)

- FIND MAXIMAL CLIQUE —
  LARGEST CONNECTED SUBGRAPH
  = LARGEST CLUSTER OF MUTUALLY-CONSISTENT MATCHES

- VERIFY MATCH BY COMPARING COMPLETE MODEL w/ IMAGE

---



Graph of pairwise-consistent feature matches

- MAXIMAL CLIQUE IN RED.
- EXAMPLE OF SEARCHING THE "CORRESPONDENCE SPACE"

---



Hypothesized Solution:

Corner 6
Hole 1 (focus feature)
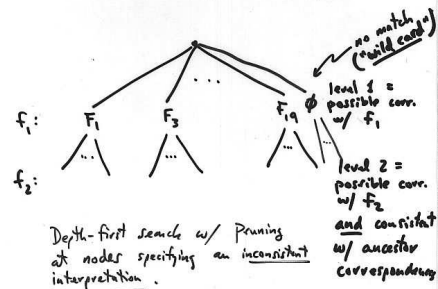Hole 2
Corner 2

Verification:

1. Compute pose
2. Project model to image space
3. Cross-correlation

## Interpretation Trees

Using geometric constraints to limit search of Correspondence Space

Given: Image features $f_1, ..., f_n$
Model features $F_1, ..., F_m$



no match
("wild card")

level 1 = possible corr. w/ $f_1$

level 2 = possible corr. w/ $f_2$ and consistent w/ ancestor correspondences

Depth-first search w/ Pruning at nodes specifying an inconsistent interpretation.

---

Size of Correspondence Space =

$$(m+1)^n$$

Height of Interpretation Tree = $n$

Goal: find path from root to leaf such that all $n$ correspondences are consistent w/ a single model

Approach: Use 1st and 2nd order geometric constraints to prune interpretation tree as we traverse it.

Possible Constraints: Unary Constraints: properties of $f_i$ and $F_j$ must be similar (i.e. match)

Ex: length of edge
corner angle
color
texture

---

Binary Constraints: consistency constraints between $(f_{i_1}, F_{j_1})$ and $(f_{i_2}, F_{j_2})$ pairings.

Ex: Relative angle between 2 edges
Range of distances between 2 edges
Range of directions from $f_{i_1}$ to $f_{i_2}$
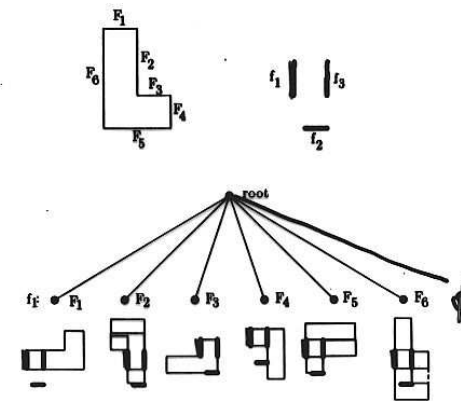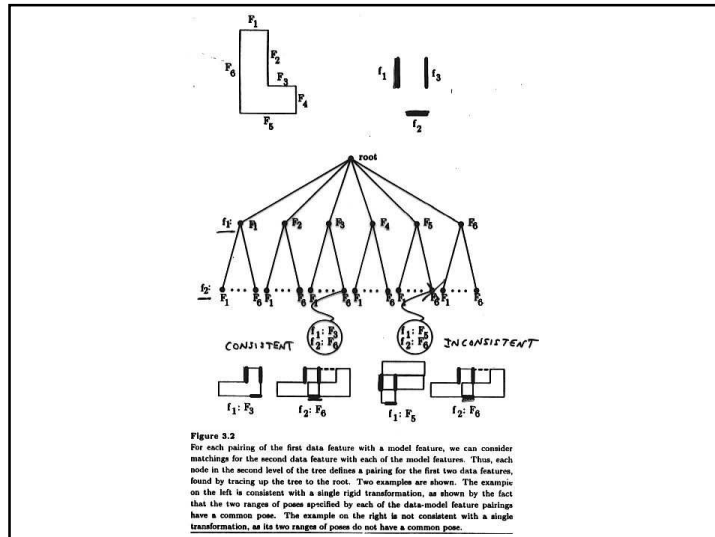
etc.

---



**Figure 3.1**
We can build a tree of possible interpretations, by first considering all the ways of matching the first data feature, $f_1$, to each of the model features, $F_j, j = 1, ..., m$. In the bottom part of the figure, we show an example of these pairings for the model and data shown at the top. In some cases, due to occlusion of the data features, a range of possible poses is given.

**Figure 3.2**
For each pairing of the first data feature with a model feature, we can consider matchings for the second data feature with each of the model features. Thus, each node in the second level of the tree defines a pairing for the first two data features, found by tracing up the tree to the root. Two examples are shown. The example on the left is consistent with a single rigid transformation, as shown by the fact that the two ranges of poses specified by each of the data-model feature pairings have a common pose. The example on the right is not consistent with a single transformation, as its two ranges of poses do not have a common pose.

---

Interpretation Tree Search Algorithm

- Depth-first search with pruning (cut-offs)

- At each node, test all unary & binary constraints. If all are satisfied, continue. Otherwise, backtrack

- At leaf node have an hypothesis for a feasible interpretation. Verify by solving for pose given correspondences, and compare projected model w/ data.
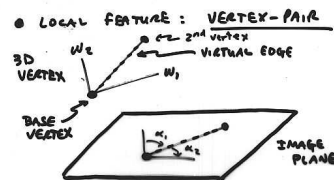


---

Method 3: Combining Local Evidence + Pose Space Search

- First match, then find best viewing transformation

- Define 6D parameter space — pose space of all possible (discretely sampled) viewing transformations

- 1. For each matching (image feature, model feat.) pair, "vote" in parameter space for all transformations they describe.
  2. Find "peak" in parameter space.

- Examples:
  Generalized Hough Transform
  Thompson & Mundy, 1987
  Hinton, 1981

---

# Voting on Pose

- Each model leads to many correct sets of correspondences, each of which has the same pose
  - Vote on pose, in an accumulator array
  - This is a Hough transform, with all its issues

THOMPSON & MUNDY'S ALGORITHM

• LOCAL FEATURE: VERTEX-PAIR

Each vertex = intersection of 2 edges.

Given image w/ n vertices,
$\Rightarrow 2n^2$ vertex-pairs

• GIVEN A CORRESPONDENCE BETWEEN MODEL & IMAGE VERTICES & EDGES, THERE IS A UNIQUE TRANSFORMATION (AFFINE) WHICH MAPS IMAGE VERTEX-PAIR TO MODEL VERTEX-PAIR
($\Rightarrow$ MODELS WEAK PERSPECTIVE PROJECTION)



ALGORITHM

FOREACH IMAGE VERTEX-PAIR DO
FOREACH MODEL VERTEX-PAIR DO

COMPUTE 6-ELEMENT TRANSFORMATION T THAT MAPS ONE TO OTHER

VOTE (INCREMENT "BIN") FOR 1 POINT IN 6D TRANSFORM PARAM. SPACE CORRESPONDING TO T.

FIND "CLUSTERS" OF VOTES (PEAK FINDING)
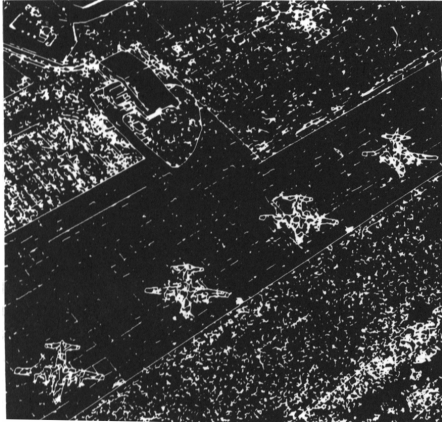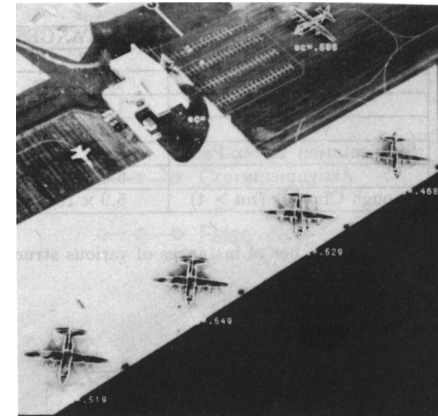
VERIFY PROJECTED MODEL W/ IMAGE



Figure from "The evolution and testing of a model-based object recognition system", J.L. Mundy and A. Heller, Proc. Int. Conf. Computer Vision, 1990 copyright 1990 IEEE



From "The evolution and testing of a model-based object recognition system", J.L. Mundy and A. Heller, Proc. ICCV, 1990

11

Figure from "The evolution and testing of a model-based object recognition system", J.L. Mundy and A. Heller, Proc. Int. Conf. Computer Vision, 1990 copyright 1990 IEEE


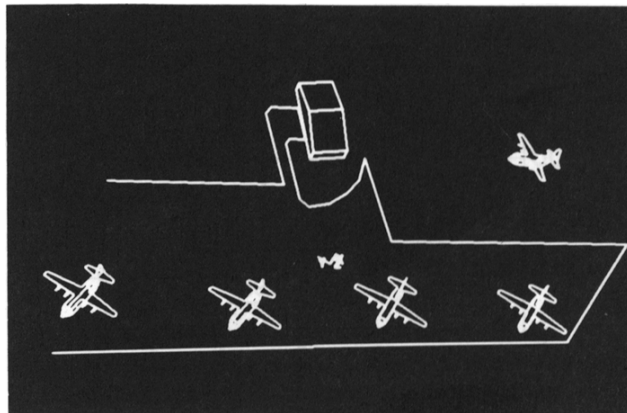
Figure from "The evolution and testing of a model-based object recognition system", J.L. Mundy and A. Heller, Proc. Int. Conf. Computer Vision, 1990 copyright 1990 IEEE



Figure from "The evolution and testing of a model-based object recognition system", J.L. Mundy and A. Heller, Proc. Int. Conf. Computer Vision, 1990 copyright 1990 IEEE

GEOMETRIC HASHING

- Wolfson, Lamdan et al. 1988

- PLANAR, RIGID objects

- Affine camera model

- Property: An affine transform of a planar rigid body is uniquely defined by the transformation of 3 ordered, non-collinear pts [e.g., klein, 1925]

$\Rightarrow$ Pick 3 non-collinear model pts $B_0, B_1, B_2$ as affine basis

Any pt $P$ represented by the affine coords $(\alpha, \beta)$ where

$$P = \alpha (B_1 - B_0) + \beta (B_2 - B_0) + B_0$$

$(\alpha, \beta)$ invariant under affine transf. $T$

$$TP = \alpha (TB_1 - TB_0) + \beta (TB_2 - TB_0) + TB_0$$

12

## Invariants

- There are geometric properties that are invariant to camera transformations
- Easiest case: view a plane object in weak perspective
- Assume we have three base points P_i on the object
  - then any other point on the object can be written as

$$P_k = P_1 + \mu_{ka}(P_2 - P_1) + \mu_{kb}(P_3 - P_1)$$

- Now image points are obtained by multiplying by a plane affine transformation, so

$$p_k = AP_k$$
$$= A\big(P_1 + \mu_{ka}(P_2 - P_1) + \mu_{kb}(P_3 - P_1)\big)$$
$$= p_1 + \mu_{ka}(p_2 - p_1) + \mu_{kb}(p_3 - p_1)$$

## Invariants

- This means that, if I know the base points in the image, I can read off the μ values for the object
  - they're the same in object and in image --- **invariant**

- Suggests a strategy rather like the Hough transform
  - search correspondences, form μ's and vote

## Geometric Hashing

- Vote on identity and correspondence using invariants
  - Take hypotheses with large enough votes
- Fill up a table, indexed by μ's, with
  - the base points and fourth point that yield those μ's
  - the object identity

Offline Modeling/Learning Phase
1. Select m feature pts on planar model(s).
2. foreach ordered, non-collinear triple of model pts do
   Compute coords of other (m-3) pts
   foreach coord (α,β) store
   Hash (α,β) = (object, basis)

$O(m^4)$ time and space.

**Algorithm 18.3:** Geometric hashing: voting on identity and point labels

For all groups of three image points $T(I)$
   For every other image point $p$
      Compute the $\mu$'s from $p$ and $T(I)$
      Obtain the table entry at these values
         if there is one, it will label the three points in $T(I)$
         with the name of the object
         and the names of these particular points.
      Cluster these labels;
         if there are enough labels, backproject and verify
      end
   end
end

---

**Online Matching Phase**

1. **Feature Detection :**
   Extract $n$ interest pts from image

2. **Choose Basis :**
   Pick 3 non-collinear pts.
   Compute affine coords of other
   $(n-3)$ pts.

3. For each affine coord,
   vote for all (model, basis) pairs
   in entry in hash table.

4. Find **Peaks** in object-Basis space
   If none, goto step 2.

5. From (model pt, image pt) pairs in Peaks
   compute best-fit affine transform
   (least squares)

6. Verify complete transformed model
   w/ image

7. If no peak verified, goto step 2
   $O(n^4)$ time

---

**ALIGNMENT METHOD**

- HUTTENLOCHER AND ULLMAN, 1987
- ARBITRARY (NON-PLANAR) objects
  defined by points and edge contours

- Weak-perspective camera model
  $$x' = \Pi(sRx + b)$$
       ↑ 2D     ↑ 3D

  where $\Pi$ = orthographic proj
  $s$ = scale
  $R$ = 3D rotation  } 6 params.
  $b$ = 2D translation

- Type of Affine transformation
  $x' = Lx + b$ where $L$ is $2 \times 2$

- OK when object far from camera
  and object depth small relative to distance
  from camera

---

**Main Idea:**
Separate problem of finding best model
from problem of finding best transformation
(3D → 2D determined by viewpoint)
of model to image

1. Hypothesize viewpoint for each model
   ⇒ solve for viewing transformation
   "normalization" step
2. Select best model by matching
   "normalized" model w/ image

Ex. CHARACTER RECOGNITION

- For each letter, store its description in some canonical position, orientation & size

- Given "viewed" letter, "undo" shift, rotation and scale:

  + Shift center-of-mass to fixed location (origin)

  + Scale convex hull to fixed size

  + Rotate major axis to fixed orientation (or detect orientations of key features & orient these — e.g., $\underset{R}{D}$, $\underset{R}{L}$)

- Match "normalized" input w/ each model and select best match.

# Recognizing Human Actions

- Movement and posture change
  - run, walk, crawl, jump, hop, swim, skate, sit, stand, kneel, lie, dance (various), …
- Object manipulation
  - pick, carry, hold, lift, throw, catch, push, pull, write, type, touch, hit, press, stroke, shake, stir, turn, eat, drink, cut, stab, kick, point, drive, bike, insert, extract, juggle, play musical instrument (various)…
- Conversational gesture
  - point, …
- Sign Language

# Activities and Situation Assessment

- Example: Withdrawing money from an ATM
- Activities constructed by composing actions. Partial order plans may be a good model.
- Activities may involve multiple agents
- Detecting unusual situations or activity patterns is facilitated by the video activity transform

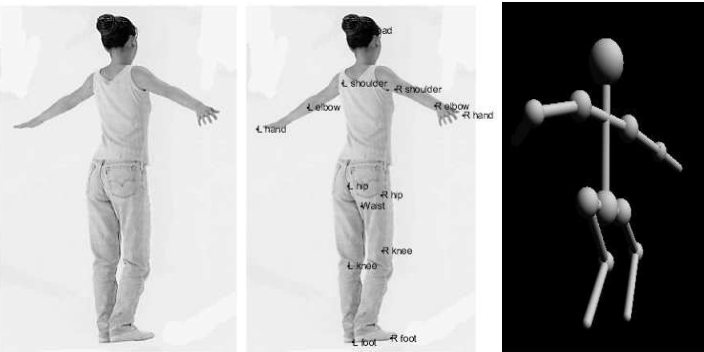| Objects in Space | Actions in Space-Time |
|---|---|
| • Segment/Region-of-interest | • Segment/volume-of-interest |
| • Features (points, curves, wavelet coefficients..) | • Features (points, curves, wavelets, motion vectors..) |
| • Correspondence and deform into alignment | • Correspondence and deform into alignment |
| • Recover parameters of generative model | • Recover parameters of generative model |
| • Discriminative classifier | • Discriminative classifier |

## Key Cues for Action Recognition

- "Morpho-kinesics" of action (shape and movement of the body)
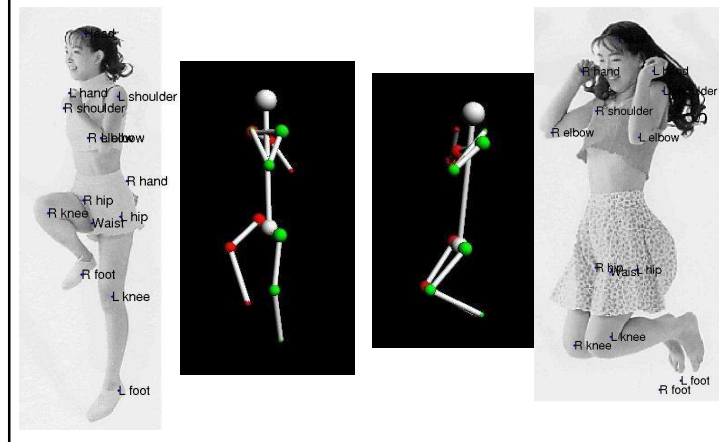- Identity of the object/s
- Activity context

## Image/Video    Stick figure    Action

- Stick figures can be specified in a variety of ways or at various resolutions (degrees of freedom)
  - 2D joint positions
  - 3D joint positions
  - Joint angles
- Complete representation
- Evidence that it is effectively computable

## Human Body Configurations



## Human Body Configurations

# Mathematical Challenges

- Modeling shape variation
- Nearest neighbor search in high dimensions
- Combining statistical optimality with computational efficiency
- Reconstruction algorithms for novel sensing modalities