



Pictorial Structures for Object Recognition

PEDRO F. FELZENSZWALB

Artificial Intelligence Lab, Massachusetts Institute of Technology

pff@ai.mit.edu

DANIEL P. HUTTENLOCHER

Computer Science Department, Cornell University

dph@cs.cornell.edu

Received November 1, 2002; Revised February 12, 2004; Accepted March 25, 2004

First online version published in September, 2004

Au: Mismatch
file disk
followed.

Abstract. In this paper we present a computationally efficient framework for part-based modeling and recognition of objects. Our work is motivated by the pictorial structure models introduced by Fischler and Elschlager. The basic idea is to represent an object by a collection of parts arranged in a deformable configuration. The appearance of each part is modeled separately, and the deformable configuration is represented by spring-like connections between pairs of parts. These models allow for qualitative descriptions of visual appearance, and are suitable for generic recognition problems. We address the problem of using pictorial structure models to find instances of an object in an image as well as the problem of learning an object model from training examples, presenting efficient algorithms in both cases. We demonstrate the techniques by learning models that represent faces and human bodies and using the resulting models to locate the corresponding objects in novel images.

Keywords: part-based object recognition, statistical models, energy minimization

1. Introduction

Research in object recognition is increasingly concerned with the ability to recognize generic classes of objects rather than just specific instances. In this paper, we consider both the problem of recognizing objects using generic part-based models and the problem of learning such models from example images. Our work is motivated by the pictorial structure representation introduced by Fischler and Elschlager (1973) thirty years ago, where an object is modeled by a collection of parts arranged in a deformable configuration. Each part encodes local visual properties of the object, and the deformable configuration is characterized by spring-like connections between certain pairs of parts. The best match of such a model to an image is found by mini-

mizing an energy function that measures both a match cost for each part and a deformation cost for each pair of connected parts.

While the pictorial structure formulation is appealing in its simplicity and generality, several shortcomings have limited its use: (i) the resulting energy minimization problem is hard to solve efficiently, (ii) the model has many parameters, and (iii) it is often desirable to find more than a single best (minimum energy) match. In this paper we address these limitations, providing techniques that are practical for a broad range of object recognition problems. We illustrate the method for two quite different generic recognition tasks, finding faces and finding people. For faces, the parts are features such as the eyes, nose and mouth, and the spring-like connections allow for variation in the relative

35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50



Figure 1. Sample results for detection of a face (a); and a human body (b). Each image shows the globally best location for the corresponding object, as computed by our algorithms. The object models were learned from training examples.

51 locations of these features. For people, the parts are the
 52 limbs, torso and head, and the spring-like connections
 53 allow for articulation at the joints. Matching results
 54 with these two models are illustrated in Fig. 1

55 The main contributions of this paper are three-fold.
 56 First, we provide an efficient algorithm for the classi-
 57 cal pictorial structure energy minimization problem de-
 58 scribed in Fischler and Elschlager (1973), for the case
 59 where the connections between parts do not form any
 60 cycles and are of a particular (but quite general) type.
 61 Many objects, including faces, people and animals can
 62 be represented by such acyclic multi-part models. Sec-
 63 ond, we introduce a method for learning these mod-
 64 els from training examples. This method learns all the
 65 model parameters, including the structure of connec-
 66 tions between parts. Third, we develop techniques for
 67 finding multiple good hypotheses for the location of an
 68 object in an image rather than just a single best solu-
 69 tion. Finding multiple hypotheses is important for tasks
 70 where there may be several instances of an object in
 71 an image, as well as for cases where imprecision in the
 72 model may result in the desired match not being the one
 73 with the minimum energy. We address the problems of
 74 learning models from examples and of hypothesizing
 75 multiple matches by expressing the pictorial structure
 76 framework in a statistical setting.

77 1.1. Pictorial Structures

78 A pictorial structure model for an object is given by
 79 a collection of parts with connections between cer-
 80 tain pairs of parts. The framework is quite general,
 81 in the sense that it is independent of the specific

scheme used to model the appearance of each part 82
 as well as the type of connections between parts. A 83
 natural way to express such a model is in terms of 84
 an undirected graph $G = (V, E)$, where the vertices 85
 $V = \{v_1, \dots, v_n\}$ correspond to the n parts, and there 86
 is an edge $(v_i, v_j) \in E$ for each pair of connected parts 87
 v_i and v_j . An instance of the object is given by a con- 88
 figuration $L = (l_1, \dots, l_n)$, where each l_i specifies the 89
 location of part v_i . Sometimes we refer to L simply 90
 as the object location, but “configuration” emphasizes 91
 the part-based representation. The location of each part 92
 can simply specify its position in the image, but more 93
 complex parameterizations are also possible. For ex- 94
 ample, for the person model in Section 6 the location 95
 of a part specifies a position, orientation and an amount 96
 of foreshortening. 97

In Fischler and Elschlager (1973) the problem of 98
 matching a pictorial structure to an image is defined 99
 in terms of an energy function to be minimized. The 100
 cost or energy of a particular configuration depends 101
 both on how well each part matches the image data 102
 at its location, and how well the relative locations of 103
 the parts agree with the deformable model. Given an 104
 image, let $m_i(l_i)$ be a function measuring the degree of 105
 mismatch when part v_i is placed at location l_i in the 106
 image. For a given pair of connected parts let $d_{ij}(l_i, l_j)$ 107
 be a function measuring the degree of deformation of 108
 the model when part v_i is placed at location l_i and part 109
 v_j is placed at location l_j . Then an optimal match of 110
 the model to the image is naturally defined as 111

$$L^* = \arg \min_L \left(\sum_{i=1}^n m_i(l_i) + \sum_{(v_i, v_j) \in E} d_{ij}(l_i, l_j) \right), \quad (1)$$

112 which is a configuration minimizing the sum of the
 113 match costs m_i for each part and the deformation costs
 114 d_{ij} for connected pairs of parts. Generally the defor-
 115 mation costs are only a function of the relative position
 116 of one part with respect to another, making the model
 117 invariant to certain global transformations. Note that
 118 matching a pictorial structure model to an image does
 119 not involve making any initial decisions about locations
 120 of individual parts, rather an overall decision is made
 121 based on both the part match costs and the deformation
 122 costs together.

123 This energy function is simple and makes intuitive
 124 sense. However, previous methods have used heuristics
 125 or local search techniques that do not find an optimal
 126 solution and depend on having good initialization. In
 127 contrast we present an efficient algorithm that can find
 128 a global minimum of the energy function without any
 129 initialization.

130 Pictorial structures can be used to represent quite
 131 generic objects. For example, the appearance models
 132 for the individual parts can be a blob of some color
 133 and orientation, or capture the response of local ori-
 134 ented filters. The connections between parts can en-
 135 code generic relationships such as “close to”, “to the
 136 left of”, or more precise geometrical constraints such
 137 as ideal joint angles. Since both the part models and
 138 the relationships between parts can be generic, picto-
 139 rial structures provide a powerful framework. Suppose
 140 we want to model the appearance of the human body. It
 141 makes sense to represent the body as an articulated ob-
 142 ject, with joints connecting different body parts. With
 143 pictorial structures we can use a coarse model, consist-
 144 ing of a small number of parts connected by flexible
 145 joints. The combination of simple appearance models
 146 for the parts and structural relations between parts pro-
 147 vides sufficient context to find the human body as a
 148 whole, even when it would be difficult to find generic
 149 parts such as “lower-leg” or “upper-arm” on their own.

150 1.2. Efficient Algorithms

151 Our primary goal is to take the pictorial structure frame-
 152 work, and use it to efficiently solve object recognition
 153 and model learning problems. We consider a natural
 154 class of pictorial structure models and present efficient
 155 algorithms both for matching such models to images
 156 and for learning models from examples. These efficient
 157 algorithms are based on two restrictions on the form of
 158 the pictorial structure models. First, our methods re-
 159 quire that the graph G be acyclic (i.e., form a tree).

160 Second the methods require that the relationships be-
 161 tween connected pairs of parts be expressed in a par-
 162 ticular form.

163 Restricting the connections between parts to a tree
 164 structure is natural for many classes of objects. For ex-
 165 ample, the connections between parts of many animate
 166 objects can form a tree corresponding to their skeletal
 167 structure. Many other kinds of objects can be repre-
 168 sented using a tree such as a star-graph, where there
 169 is one central part to which all the other parts are con-
 170 nected. When the graph G is a tree it is possible to com-
 171 pute the best match of the model to an image in poly-
 172 nomial time. This is done using a generalization of the
 173 Viterbi algorithm (Rabiner and Juang, 1993). Related
 174 methods are known in the Bayesian Network commu-
 175 nity as belief propagation algorithms (Pearl, 1988). The
 176 fastest such polynomial time algorithms run in $O(h^2n)$
 177 time, where n is the number of object parts, and h is
 178 a discrete number of possible locations for each part.
 179 Unfortunately this is too slow in most cases because the
 180 number of possible locations for a single part is usually
 181 quite large – in the hundreds of thousands or millions.

182 The restriction that we impose on the form of con-
 183 nections between parts enables an improvement in the
 184 running time of the matching algorithms so that it be-
 185 comes essentially linear rather than quadratic in the
 186 number of possible locations for each part. We require
 187 that $d_{ij}(l_i, l_j)$ be a Mahalanobis distance between trans-
 188 formed locations,

$$d_{ij}(l_i, l_j) = (T_{ij}(l_i) - T_{ji}(l_j))^T M_{ij}^{-1} (T_{ij}(l_i) - T_{ji}(l_j)), \quad (2)$$

189 The matrix M_{ij} should be diagonal, and for simplicity
 190 we will assume that T_{ij} , and T_{ji} are one-to-one. We
 191 further require that it be possible to represent the set
 192 of possible transformed locations $T_{ij}(l_i)$ and $T_{ji}(l_j)$ as
 193 positions in a grid. These functions capture the ideal
 194 relative locations for parts v_i and v_j . The distance be-
 195 tween the transformed locations, weighted by M_{ij}^{-1} ,
 196 measures the deformation of a “spring” connecting the
 197 two parts. This special form for the deformation costs
 198 allows for matching algorithms that run in time linear in
 199 the number of grid positions of the transformed space.
 200 Often this is the same as the number of possible loca-
 201 tions for each part, but sometimes it may be slightly
 202 larger. As we will see, a broad class of interesting re-
 203 lationships can be represented in this form, including
 204 those illustrated in Sections 5 and 6.

205 The asymptotic running time of the matching algo-
 206 rithms that we develop is thus nearly optimal, in

207 the sense that the methods run in essentially the same
 208 asymptotic time as simply matching each part to the
 209 image separately, without accounting for the connec-
 210 tions between them. In practice, the algorithms are also
 211 quite fast, finding the globally best match of a picto-
 212 rial structure to an image in just a few seconds using a
 213 desktop computer.

214 *1.3. Statistical Formulation*

215 In their original work, Fischler and Elschlager only
 216 considered the problem of finding the best match of
 217 a pictorial structure model to an image. As discussed
 218 above, they characterized this problem using the en-
 219 ergy function in Eq. (1). While this energy function
 220 intuitively makes sense, it has many free paramet-
 221 ers. For each different object, one has to construct
 222 a model, which includes picking appearance param-
 223 eters for each part, a set of edges connecting pairs of
 224 parts and the characteristics of the connections. We
 225 are interested in automatically learning these param-
 226 eters from examples. Moreover, the energy minimiza-
 227 tion formulation only characterizes the problem of find-
 228 ing the best match of a model to an image, whereas
 229 it is often desirable to find multiple good potential
 230 matches.

231 These questions are naturally addressed using a sta-
 232 tistical framework for pictorial structure models which
 233 we describe in Section 2. In this framework, the en-
 234 ergy minimization problem introduced by Fischler and
 235 Elschlager is equivalent to finding the maximum a pos-
 236 teriori estimate of the object configuration given an ob-
 237 served image. The statistical formulation can be used to
 238 learn the parameters of a model from examples. In fact,
 239 *all* model parameters can be learned from a few training
 240 examples using maximum likelihood estimation. This
 241 is of practical as well as theoretical interest, since it is
 242 generally not possible to find the best parameters for a
 243 deformable model by trial and error.

244 The statistical framework also provides a natural way
 245 of finding several good matches of a model to an im-
 246 age rather than finding just the best one. The idea is
 247 to consider primarily good matches without consider-
 248 ing many bad ones. We can achieve this by sampling
 249 object configurations from their posterior probability
 250 distribution given an observed image. Sampling makes
 251 it possible to find many locations for which the pos-
 252 terior is high, and to subsequently select one or more
 253 of those using an independent method. This procedure
 254 lets us use imprecise models for generating hypotheses

and can be seen as a mechanism for visual selection 255
 (see Amit and Geman, 1999). 256

1.4. Related Work 257

Research in object recognition has been dominated 258
 by approaches that separate processing into distinct 259
 stages of feature extraction and matching. In the first 260
 stage, discrete primitives, or “features” are detected. In 261
 the second stage, stored models are matched against 262
 those features. For instance, in the pioneering work of 263
 Roberts (1965) children’s blocks were recognized by 264
 first extracting edges and corners from images and then 265
 matching these features to polyhedral models of the 266
 blocks. The model-based recognition paradigm of the 267
 1980’s similarly followed this approach. These meth- 268
 ods focus largely on the problem of efficiently search- 269
 ing for correspondences between features that have 270
 been extracted from an image, and features of a stored 271
 model. Examples include interpretation tree search 272
 (Ayache and Faugeras, 1986; Grimson and Lozano- 273
 Perez, 1987), the alignment method (Huttenlocher and 274
 Ullman, 1990), RANSAC (Fischler and Bolles, 1981) 275
 and geometric hashing (Lamdan et al., 1990). 276

Limitations of the simple features used by most 277
 earlier model-based recognition techniques led to a 278
 quite different class of recognition methods, devel- 279
 oped in the 1990’s, which operate directly on images 280
 rather than first extracting discrete features. These in- 281
 clude both appearance-based methods (e.g., Turk and 282
 Pentland, 1991; Murase and Nayar, 1995) and 283
 template-based methods such as Hausdorff matching 284
 (Huttenlocher et al., 1993). Such approaches treat im- 285
 ages as the entities to be recognized, rather than having 286
 more abstract models based on features or other primi- 287
 tives. One or more training images of an object are used 288
 to form a “template” that is used as a model. This model 289
 is then compared to new images to determine whether 290
 or not the target is present, generally by explicitly con- 291
 sidering possible transformations of the template. 292

The matching of pictorial structures is an alternative 293
 approach that in many ways combines the appearance- 294
 based and geometric techniques. The energy minimiza- 295
 tion problem associated with these models as defined 296
 in Eq. (1) incorporates match costs for the individual 297
 parts and deformation costs for the geometric configu- 298
 ration into a single overall problem. Thus the approach 299
 provides a means of simultaneously using appearance 300
 and geometry, rather than first making binary decisions 301
 about the possible locations of parts or features. The 302

303 main drawback of the pictorial structures approach has
 304 been the computational difficulty of the energy mini-
 305 mization problem, which we address here for a class of
 306 models.

307 There have been other part-based recognition meth-
 308 ods, which like the pictorial structures approach are
 309 based on separately modeling the appearance of in-
 310 dividual parts and the geometric relations between
 311 them. However most of these part-based methods make
 312 binary decisions about potential part locations (e.g.,
 313 Pentland, 1987; Dickinson et al., 1993; Rivlin et al.,
 314 1995; Burl and Perona, 1996). Moreover, most part-
 315 based methods use some kind of search heuristics, such
 316 as first matching a particular “distinctive” part and then
 317 searching for other parts given that initial match, in or-
 318 der to avoid the combinatorial explosion of the con-
 319 figuration space. Such heuristics make it difficult to
 320 handle occlusion, particularly for those parts that are
 321 considered first in the search.

322 In Burl et al. (1998) models similar to pictorial struc-
 323 tures were used to represent objects in terms of a con-
 324 stellations of local features. In these models, rather than
 325 there being connections between pairs of parts, all the
 326 parts are constrained with respect to a central coordi-
 327 nate system using a Gaussian distribution. Like the
 328 pictorial structures formulation, no binary decisions are
 329 made about part or feature locations. These models,
 330 however, are not well suited for representing articu-
 331 lated objects, as a joint Gaussian distribution cannot
 332 capture multiple articulation points. Moreover, in Burl
 333 et al. (1998) the matching algorithms use heuristics that
 334 don’t necessarily find the optimal match of a model to
 335 an image.

336 The problem of finding people in images using
 337 coarse part-based two-dimensional models was con-
 338 sidered in Ioffe and Forsyth (2001). This is one of two
 339 domains that we use to illustrate the pictorial struc-
 340 tures approach. Two different methods are reported in
 341 Ioffe and Forsyth (2001). The first method makes bi-
 342 nary decisions about the possible locations for indi-
 343 vidual parts and subsequently searches for groups of
 344 parts that match the overall model. The second method
 345 uses sequential importance sampling (particle filtering)
 346 to generate increasingly larger configurations of parts.
 347 We also describe a sampling-based technique, however
 348 rather than employing approximate distributions ob-
 349 tained via sequential importance sampling, our method
 350 is based on efficiently computing the *exact* (discrete)
 351 posterior distribution for the object configuration and
 352 then sampling from that posterior.

In illustrating the pictorial structures approach us- 353
 ing the problem of finding people in images we 354
 employ simple part models based on binary images 355
 obtained by background subtraction. This suggests 356
 comparisons with silhouette-based deformable match- 357
 ing techniques (e.g., Gdalyahu and Weinshall, 1999; 358
 Sebastian et al., 2001). These approaches are quite dif- 359
 ferent, however. First of all, silhouette-based methods 360
 generally operate using boundary contours, requiring 361
 good segmentation of the object from the background. 362
 In contrast, the models we use are not based on a bound- 363
 ary representation and operate directly on binary im- 364
 ages. For example, a single part could match a region 365
 of the image that has several disconnected components. 366
 Secondly, deformable matching methods are generally 367
 based on two-dimensional shape representations rather 368
 than highly parameterized models. Thus they do not ap- 369
 ply to cases such as an articulated body where in some 370
 configurations the parts can cross one another yielding 371
 vastly different shapes. 372

373 Finally we note that models similar to pictorial struc-
 374 tures have recently been used for tracking people by
 375 matching models at each frame (Ramanan and Forsyth,
 376 2003). In contrast, most work on tracking highly articu-
 377 lated objects such as people relies heavily on motion
 378 information (Bregler and Malik, 1998; Ju et al., 1996)
 379 and only performs incremental updates in the object
 380 configuration. In such approaches, some other method
 381 is used to find an initial match of the model to the image,
 382 and then tracking commences from that initial condi-
 383 tion. Pictorial structures can be used to solve this track
 384 initialization problem, or as demonstrated in Ramanan
 385 and Forsyth (2003) can be used as a tracking method
 386 on their own.

2. Statistical Framework 387

As noted in the introduction, the pictorial structure en- 388
 ergy minimization problem can be viewed in terms 389
 of statistical estimation. The statistical framework de- 390
 scribed here is useful for addressing two of the three 391
 questions that we consider in this paper, that of learn- 392
 ing pictorial structure models from examples and that 393
 of finding multiple good matches of a model to an im- 394
 age. For the third question, that of efficiently minimiz- 395
 ing the energy in Eq. (1), the statistical formulation 396
 provides relatively little insight, however it unifies the 397
 three questions in a common framework. 398

A standard way of approaching object recognition 399
 in a statistical setting is as follows. Let θ be a set of 400

401 parameters that define an object model, I denote an
 402 image, and as before let L denote a configuration of
 403 the object (a location for each part). The distribution
 404 $p(I | L, \theta)$ captures the imaging process, and measures
 405 the likelihood of seeing a particular image given that
 406 an object is at some location. The distribution $p(L | \theta)$
 407 measures the prior probability that an object is at a
 408 particular location. Finally, the posterior distribution,
 409 $p(L | I, \theta)$, characterizes the probability that the object
 410 configuration is L given the model θ and the image I .
 411 Using Bayes' rule the posterior can be written as,

$$p(L | I, \theta) \propto p(I | L, \theta)p(L | \theta). \quad (3)$$

412 A common drawback of the Bayesian formulation
 413 is the difficulty of determining a prior distribution,
 414 $p(L | \theta)$, that is both informative and generally appli-
 415 cable. For instance, a uniform prior is general but pro-
 416 vides no information. On the other hand a prior which
 417 says that the object is in the lower left corner of the
 418 image is highly informative but of little use in general.
 419 For pictorial structures, the prior over configurations
 420 encodes information about the *relative* positions of the
 421 parts, which can be both informative and general. For
 422 instance, for a human body model such a prior can
 423 capture which are likely relative orientations of two
 424 connected limbs.

425 A number of interesting problems can be character-
 426 ized in terms of this statistical framework,

- 427 • MAP estimation—this is the problem of finding a
 428 location L with maximum posterior probability. In
 429 some sense, the MAP estimate is our best guess for
 430 the location of the object. In our framework this will
 431 be equivalent to the energy minimization problem
 432 defined by Eq. (1).
- 433 • Sampling from the posterior—sampling provides a
 434 natural way to hypothesize many good potential
 435 matches of a model to an image, rather than just
 436 finding the best one. This is useful to detect multiple
 437 instances of an object in an image and to find possible
 438 locations of an object with an imprecise model.
- 439 • Model estimation—this is the problem of finding θ
 440 which specifies a good model for a particular ob-
 441 ject. The statistical framework allows us to learn
 442 the model parameters from training examples using
 443 maximum likelihood estimation.

444 Our pictorial structure models are parametrized by
 445 $\theta = (u, E, c)$, where $u = \{u_1, \dots, u_n\}$ are appear-
 446 ance parameters, the set of edges E indicates which

parts are connected, and $c = \{c_{ij} | (v_i, v_j) \in E\}$ are
 447 connection parameters. There is a separate appearance
 448 model for each part, but the exact method used to model
 449 the appearance of parts is not important at this point.
 450 In Section 5 we model appearance using image deriva-
 451 tives around a point, to represent local features of a face
 452 such as the tip of the nose or the corners of the mouth.
 453 In Section 6 we model appearance using rectangular
 454 shapes, to represent individual body parts. In practice,
 455 the appearance modeling scheme just needs to provide
 456 a distribution $p(I | l_i, u_i)$ up to a normalizing constant,
 457 which measures the likelihood of seeing a particular
 458 image, given that a part with appearance parameters u_i
 459 is at location l_i . This distribution does not have to be
 460 a precise generative model, an approximate measure is
 461 good enough in practice.

We model the likelihood of seeing an image given
 463 that the object is at some configuration by the product
 464 of the individual likelihoods,
 465

$$p(I | L, \theta) = p(I | L, u) \propto \prod_{i=1}^n p(I | l_i, u_i). \quad (4)$$

This approximation is good if the parts do not overlap,
 466 as in this case they generate different portions of the
 467 image. But the approximation can be bad if one part
 468 occludes another. For the iconic models described in
 469 Section 5 the prior distribution over configurations en-
 470 forces that the parts do not overlap (the probability of
 471 a configuration with overlap is very small). For the ar-
 472 ticulated models described in Section 6 there is much
 473 less constraint on the locations of parts, and parts can
 474 easily overlap. In this case we demonstrate that a good
 475 estimate of the object configuration can be found by ob-
 476 taining multiple samples from the posterior distribution
 477 and then selecting one of them using an independent
 478 method. This shows that sampling from the posterior
 479 can be useful for handling modeling error.
 480

The prior distribution over object configurations is
 481 captured by a tree-structured Markov random field with
 482 edge set E . In general, the joint distribution for a tree-
 483 structured prior can be expressed as,
 484

$$p(L | \theta) = \frac{\prod_{(v_i, v_j) \in E} p(l_i, l_j | \theta)}{\prod_{v_i \in V} p(l_i | \theta)^{\deg v_i - 1}},$$

where $\deg v_i$ is the degree of vertex v_i in the graph de-
 485 fined by E . We do not model any preference over the
 486 absolute location of each part, only over their relative
 487 configuration. This means that $p(l_i | \theta)$ is constant, and
 488 we let it equal one for simplicity. The joint distributions
 489

490 for pairs of parts connected by edges are characterized
 491 by the parameters $c = \{c_{ij} \mid (v_i, v_j) \in E\}$. Since we let
 492 $p(l_i \mid \theta) = 1$, the prior distribution over object configura-
 493 tions is given by,

$$p(L \mid \theta) = p(L \mid E, c) = \prod_{(v_i, v_j) \in E} p(l_i, l_j \mid c_{ij}). \quad (5)$$

494 Note that both $p(l_i, l_j \mid c_{ij})$ and $p(L \mid E, c)$ are im-
 495 proper priors (see Berger, 1985). This is a consequence
 496 of using an uninformative prior over absolute locations
 497 for each part.

498 In Eq. (4) we defined the form of $p(I \mid L, \theta)$, the
 499 likelihood of seeing an image given that the object is
 500 at a some configuration, and in Eq. (5) we defined the
 501 form of $p(L \mid \theta)$, the prior probability that the object
 502 would assume a particular configuration. These can be
 503 substituted into Eq. (3) yielding,

$$P(L \mid I, \theta) \propto \left(\prod_{i=1}^n p(I \mid l_i, u_i) \prod_{(v_i, v_j) \in E} p(l_i, l_j \mid c_{ij}) \right).$$

504 Taking the negative logarithm of this equation yields
 505 the same energy function that is being minimized in
 506 Eq. (1), where $m_i(l_i) = -\log p(I \mid l_i, u_i)$ is a match
 507 cost measuring how well part v_i matches the image
 508 data at location l_i , and $d_{ij}(l_i, l_j) = -\log p(l_i, l_j \mid c_{ij})$
 509 is a deformation cost measuring how well the relative
 510 locations for v_i and v_j agree with the prior model. Thus
 511 we see that the MAP estimation problem for the statisti-
 512 cal models introduced in this section is equivalent to
 513 the original energy minimization problem for pictorial
 514 structures described in Fischler and Elschlager (1973).

515 As discussed in the introduction our efficient algo-
 516 rithms require that the deformation costs be expressed
 517 in a particular form as shown in Eq. (2). This require-
 518 ment has a natural interpretation in terms of the statisti-
 519 cal models. Since $d_{ij}(l_i, l_j) = -\log p(l_i, l_j \mid c_{ij})$, it is
 520 equivalent to assume that the joint prior distribution for
 521 the locations of a pair of connected parts is given by a
 522 Gaussian over the displacement between transformed
 523 locations,

$$p(l_i, l_j \mid c_{ij}) \propto \mathcal{N}(T_{ij}(l_i) - T_{ji}(l_j), 0, D_{ij}), \quad (6)$$

524 where T_{ij} , T_{ji} , and D_{ij} are the connection parameters
 525 encoded by c_{ij} . These parameters correspond to the
 526 ones in Eq. (2) where $D_{ij} = M_{ij}/2$ is a diagonal co-
 527 variance matrix.

3. Learning Model Parameters 528

Suppose we are given a set of example images 529
 $\{I^1, \dots, I^m\}$ and corresponding object configurations 530
 $\{L^1, \dots, L^m\}$ for each image. We want to use the train- 531
 ing examples to obtain estimates for the model param- 532
 eters $\theta = (u, E, c)$, where $u = \{u_1, \dots, u_n\}$ are the 533
 appearance parameters for each part, E is the set of con- 534
 nections between parts, and $c = \{c_{ij} \mid (v_i, v_j) \in E\}$ 535
 are the connection parameters. The maximum likeli- 536
 hood (ML) estimate of θ is, by definition, the value θ^* 537
 that maximizes 538

$$p(I^1, \dots, I^m, L^1, \dots, L^m \mid \theta) = \prod_{k=1}^m p(I^k, L^k \mid \theta),$$

where the right hand side is obtained by assuming 539
 that each example was generated independently. Since 540
 $p(I, L \mid \theta) = p(I \mid L, \theta)p(L \mid \theta)$, the ML estimate is 541

$$\theta^* = \arg \max_{\theta} \prod_{k=1}^m p(I^k \mid L^k, \theta) \prod_{k=1}^m p(L^k \mid \theta). \quad (7)$$

The first term in this equation depends only on the ap- 542
 pearance of the parts, while the second term depends 543
 only on the set of connections and connection param- 544
 eters. Below we show that one can independently solve 545
 for the appearance models of the individual parts and 546
 the structural model given by the connections and their 547
 parameters. As a consequence, any kind of part models 548
 can be used in this framework as long as there is a max- 549
 imum likelihood estimation procedure for learning the 550
 model parameters for a single part from examples. We 551
 use quite simple part models in this paper because our 552
 focus is on developing a general framework and provid- 553
 ing efficient algorithms that can be used with many 554
 different modeling schemes. 555

3.1. Estimating the Appearance Parameters 556

From Eq. (7) we get 557

$$u^* = \arg \max_u \prod_{k=1}^m p(I^k \mid L^k, u).$$

The likelihood of seeing image I^k , given the configura- 558
 tion L^k for the object is given by Eq. (4). Thus, 559

$$\begin{aligned} u^* &= \arg \max_u \prod_{k=1}^m \prod_{i=1}^n p(I^k \mid l_i^k, u_i) \\ &= \arg \max_u \prod_{i=1}^n \prod_{k=1}^m p(I^k \mid l_i^k, u_i). \end{aligned}$$

560 Looking at the right hand side we see that to find u^* we
 561 can independently solve for the u_i^* ,

$$u_i^* = \arg \max_{u_i} \prod_{k=1}^m p(I^k | l_i^k, u_i).$$

562 This is exactly the ML estimate of the appearance
 563 parameters for part v_i , given independent examples
 564 $\{(I^1, l_i^1), \dots, (I^m, l_i^m)\}$. Solving for u_i^* depends on
 565 picking a specific modeling scheme for the parts, and
 566 we return to this in Sections 5 and 6.

567 3.2. Estimating the Dependencies

568 From Eq. (7) we get

$$E^*, c^* = \arg \max_{E, c} \prod_{k=1}^m p(L^k | E, c). \quad (8)$$

569 We need to pick a set of edges that form a tree and
 570 the connection parameters for each edge. This can be
 571 done in a similar way to the algorithm of Chow and Liu
 572 (1968), which estimates a tree distribution for discrete
 573 random variables. Eq. (5) defines the prior probability
 574 of the object assuming configuration L^k as,

$$p(L^k | E, c) = \prod_{(v_i, v_j) \in E} p(l_i^k, l_j^k | c_{ij}).$$

575 Plugging this into Eq. (8) and re-ordering the factors
 576 we get,

$$E^*, c^* = \arg \max_{E, c} \prod_{(v_i, v_j) \in E} \prod_{k=1}^m p(l_i^k, l_j^k | c_{ij}). \quad (9)$$

577 We can estimate the parameters for each possible con-
 578 nection independently, even before we know which
 579 connections will actually be in E as,

$$c_{ij}^* = \arg \max_{c_{ij}} \prod_{k=1}^m p(l_i^k, l_j^k | c_{ij}).$$

580 This is the ML estimate for the joint distribution
 581 of l_i and l_j , given independent examples $\{(l_i^1, l_j^1),$
 582 $\dots, (l_i^m, l_j^m)\}$. Solving for c_{ij}^* depends on picking a
 583 specific representation for the joint distributions. In-
 584 dependent of the exact form of $p(l_i, l_j | c_{ij})$, and how
 585 to compute c_{ij}^* (which we consider later, as it varies
 586 with different modeling schemes), we can characterize
 587 the “quality” of a connection between two parts as the

probability of the examples under the ML estimate for
 their joint distribution, 588 589

$$q(v_i, v_j) = \prod_{k=1}^m p(l_i^k, l_j^k | c_{ij}^*).$$

Intuitively, the quality of a connection between two
 parts measures the extent to which their locations are
 related. These quantities can be used to estimate the
 connection set E^* as follows. We know that E^* should
 form a tree, and according to Eq. (9) we let, 590 591 592 593 594

$$E^* = \arg \max_E \prod_{(v_i, v_j) \in E} q(v_i, v_j) \\ = \arg \min_E \sum_{(v_i, v_j) \in E} -\log q(v_i, v_j). \quad (10)$$

The right hand side is obtained by taking the nega-
 tive logarithm of the quantity being maximized (and
 thus finding the argument minimizing the value, in-
 stead of maximizing it). Solving for E^* is equivalent to
 the problem of computing the minimum spanning tree
 (MST) of a graph. We build a complete graph on the
 vertices V , and associate a weight $-\log q(v_i, v_j)$ with
 each edge (v_i, v_j) . The MST of this graph is the tree
 with minimum total weight, which is exactly the set of
 edges defined by Eq. (10). The MST problem is well
 known (see Cormen et al., 1996) and can be solved ef-
 ficiently. Kruskal’s algorithm can be used to compute
 the MST in $O(n^2 \log n)$ time, since we have a complete
 graph with n nodes. 595 596 597 598 599 600 601 602 603 604 605 606 607 608

4. Matching Algorithms 609

In this section we present two efficient algorithms for
 matching tree-structured models to images with con-
 nections of the form in Eqs. (2) and (6). The first al-
 gorithm solves the energy minimization problem in
 Eq. (1), which in the statistical framework is equiva-
 lent to finding the MAP estimate of the object location
 given an observed image. The second algorithm sam-
 ples configurations from the posterior distribution. In
 Felzenszwalb and Huttenlocher (2000) we described
 a version of the energy minimization algorithm that
 uses a different restriction on the form of connections
 between parts. That form did not allow for efficient
 sampling from the posterior distribution. 610 611 612 613 614 615 616 617 618 619 620 621 622

623 4.1. Energy Minimization or MAP Estimate

624 As discussed in Section 1.1, the problem of finding the
625 best match of a pictorial structure model to an image
626 is defined by the following equation,

$$L^* = \arg \min_L \left(\sum_{i=1}^n m_i(l_i) + \sum_{(v_i, v_j) \in E} d_{ij}(l_i, l_j) \right).$$

627 The form of this minimization is quite general, and it
628 appears in a number of problems in computer vision,
629 including MAP estimation of Markov random fields for
630 low-level vision such as image restoration and stereo
631 and optimization of active contour models (snakes).
632 While the form of the minimization is shared with these
633 other problems, the structure of the graph and space of
634 possible solutions differ substantially. This changes the
635 computational nature of the problem.

636 Solving this minimization for arbitrary graphs and
637 arbitrary functions m_i , d_{ij} is an NP-hard problem
638 (see Boykov et al., 2001). However, when the graph
639 $G = (V, E)$ has a restricted form, the problem can be
640 solved more efficiently. For instance, with first-order
641 snakes the graph is simply a chain, which enables a dynamic
642 programming solution that takes $O(h^2n)$ time
643 (see Amini et al., 1990), where as before we use n to
644 denote the number of parts in the model and h is a discrete
645 number of possible locations for each part. Moreover, with
646 snakes the minimization is done over a small number of
647 locations for each vertex (e.g., the current location plus
648 the 8 neighbors on the image grid). This minimization is
649 then iterated until the change in energy is small. The key
650 to an efficient algorithm for snakes is that the number of
651 possible locations for each part, h , is small in each
652 iteration, as the dynamic programming solution is quadratic
653 in this value. Another source of efficient algorithms has
654 been in restricting d_{ij} to a particular form. This approach
655 has been particularly fruitful in some recent work on
656 MRFs for low-level vision (Boykov et al., 2001; Ishikawa
657 and Geiger, 1998). Here we use constraints on both the
658 structure of the graph and the form of d_{ij} .

659 By restricting the graphs to trees, a similar kind of
660 dynamic programming can be applied as is done for
661 chains, making the minimization problem polynomial
662 rather than exponential time. The precise technique is
663 described in Section 4.1.1. However, this $O(h^2n)$ algorithm
664 is not practical in most cases, because for pictorial
665 structures the number of possible locations for each part
666 is usually huge.

668 Recall our restricted form for d_{ij} shown in Eq. (2) in
669 terms of a Mahalanobis distance between transformed
670 locations,

$$d_{ij}(l_i, l_j) = (T_{ij}(l_i) - T_{ji}(l_j))^T M_{ij}^{-1} (T_{ij}(l_i) - T_{ji}(l_j)).$$

671 We will show how this restriction can be used to obtain a
672 minimization algorithm that runs in $O(h'n)$ rather than
673 $O(h^2n)$ time, where h' is the number of grid locations
674 in a discretization of the space of transformed locations
675 given by T_{ij} and T_{ji} . The relationship between h'
676 and h depends on the particular transformations being
677 used, but in most cases the two quantities have similar
678 value. This makes it quite practical to compute a
679 *globally optimal match* of a pictorial structure model to an
680 image, up to the discretization of the possible locations.
681 We first discuss the overall minimization problem for
682 tree-structured models and then turn to the method that
683 exploits the form of d_{ij} .

684 **4.1.1. Efficient Minimization.** In this section, we describe
685 an algorithm for finding a configuration $L^* = (l_1^*, \dots, l_n^*)$
686 that minimizes Eq. (1) when the graph G is a tree, which
687 is based on the well known Viterbi recurrence. Given $G = (V, E)$,
688 let $v_r \in V$ be an arbitrarily chosen root vertex (this choice
689 does not affect the results). From this root, each vertex
690 $v_i \in V$ has a depth d_i which is the number of edges between
691 it and v_r (and the depth of v_r is 0). The children, C_i ,
692 of vertex v_i are those neighboring vertices, if any, of
693 depth $(d_i + 1)$. Every vertex v_i other than the root has a
694 unique parent, which is the neighboring vertex of depth
695 $(d_i - 1)$.

697 For any vertex v_j with no children (i.e., any leaf of the
698 rooted tree), the best location l_j^* for that vertex can be
699 computed as a function of the location of just its parent,
700 v_i . The only edge incident on v_j is (v_i, v_j) , thus the
701 only contribution of l_j to the energy in (1) is $m_j(l_j) +$
702 $d_{ij}(l_i, l_j)$. The quality of the best location for v_j
703 given location l_i for v_i is

$$B_j(l_i) = \min_{l_j} (m_j(l_j) + d_{ij}(l_i, l_j)), \quad (11)$$

704 and the best location for v_j as a function of l_i can be
705 obtained by replacing the min in the equation above with
706 $\arg \min$.

707 For any vertex v_j other than the root, assume that the
708 function $B_c(l_j)$ is known for each child $v_c \in C_j$. That is,
709 the quality of the best location for each child is known with
710 respect to the location of v_j . Then the

711 quality of the best location for v_j given a location for
712 its parent v_i is

$$B_j(l_i) = \min_{l_j} \left(m_j(l_j) + d_{ij}(l_i, l_j) + \sum_{v_c \in C_j} B_c(l_j) \right). \quad (12)$$

713 Again, the best location for v_j as a function of l_i can
714 be obtained by replacing the min in the equation above
715 with arg min. This equation subsumes (11) because for
716 a leaf node the sum over its children is simply empty.
717 Finally, for the root v_r , if $B_c(l_r)$ is known for each child
718 $v_c \in C_r$ then the best location for the root is

$$l_r^* = \arg \min_{l_r} \left(m_r(l_r) + \sum_{v_c \in C_r} B_c(l_r) \right).$$

719 That is, the minimization in (1) can be expressed re-
720 cursively in terms of the $(n - 1)$ functions $B_j(l_i)$ for
721 each vertex $v_j \in V$ (other than the root). These re-
722 cursive equations suggest a simple algorithm. Let d be
723 the maximum depth in the tree. For each node v_j with
724 depth d , compute $B_j(l_i)$, where v_i is the parent of v_j .
725 These are all leaf nodes, so clearly $B_j(l_i)$ can be com-
726 puted as in (11). Next, for each node v_j with depth
727 $(d - 1)$ compute $B_j(l_i)$, where again v_i is the parent of
728 v_j . Clearly, $B_c(l_j)$ has been computed for every child
729 v_c of v_j , because the children have depth d . Thus $B_j(l_i)$
730 can be computed as in (12). Continue in this manner,
731 decreasing the depth until reaching the root at depth
732 zero. Besides computing each B_j we also compute B'_j ,
733 which indicates the best location of v_j as a function of
734 its parent location (obtained by replacing the min in B_j
735 with arg min). At this point, we compute the optimal
736 location l_r^* for the root. The optimal location L^* for
737 all the parts can be computed by tracing back from the
738 root to each leaf. We know the optimal location of each
739 node given the location of its parent, and the optimal
740 location of each parent is now known starting from the
741 root.

742 The overall running time of this algorithm is $O(Hn)$,
743 where H reflects the time required to compute each
744 $B_j(l_i)$ and $B'_j(l_i)$. In the general case this takes $O(h^2)$
745 time as it is necessary to consider every location of a
746 child node for each possible location of the parent. In
747 the next section, we show how to compute each $B_j(l_i)$
748 and $B'_j(l_i)$ more efficiently when d_{ij} is restricted to be
749 in the form of Eq. (2).

4.1.2. Generalized Distance Transforms. Traditional distance transforms are defined for sets of points on a grid. Suppose we have a grid \mathcal{G} , and $\rho(x, y)$ is some measure of distance between points on the grid. Given a point set $B \subseteq \mathcal{G}$, the distance transform of B specifies for each location in the grid, the distance to the closest point in the set,

$$\mathcal{D}_B(x) = \min_{y \in B} \rho(x, y).$$

In particular, \mathcal{D}_B is zero at any point in B , and is small at nearby locations. The distance transform is commonly used for matching edge based models (see Borgefors, 1988; Huttenlocher et al., 1993). The trivial way to compute this function takes $O(k|B|)$ time, where k is the number of locations in the grid. On the other hand, efficient algorithms exist to compute the distance transform in $O(k)$ time, independent of the number of points in B (see Borgefors, 1986; Karzanov, 1992). These algorithms have small constants and are very fast in practice. In order to compute the distance transform, it is commonly expressed as

$$\mathcal{D}_B(x) = \min_{y \in \mathcal{G}} (\rho(x, y) + 1_B(y)),$$

where $1_B(y)$ is an indicator function for membership in the set B , that has value 0 when $y \in B$ and ∞ otherwise. This suggests a generalization of distance transforms where the indicator function is replaced with some arbitrary function over the grid \mathcal{G} ,

$$\mathcal{D}_f(x) = \min_{y \in \mathcal{G}} (\rho(x, y) + f(y)).$$

Intuitively, for each grid location x , the transform finds a location y that is close to x and for which $f(y)$ is small. Note that if there is a location where $f(x)$ has a small value, \mathcal{D}_f will have small value at x and nearby locations.

With the restricted form of d_{ij} in Eq. (2), the functions $B_j(l_i)$ that must be computed by the dynamic programming algorithm can be rewritten as generalized distance transforms, where the distance in the grid, $\rho(x, y)$, is given by the Mahalanobis distance defined by M_{ij} ,

$$B_j(l_i) = \mathcal{D}_f(T_{ij}(l_i)),$$

750
751
752
753
754
755
756

757
758
759
760
761
762
763
764
765
766
767
768

769
770
771
772
773

774
775
776
777
778
779
780
781
782
783
784

785 where

$$f(y) = \begin{cases} m_j(T_{ji}^{-1}(y)) + \sum_{v_c \in C_j} B_c(T_{ji}^{-1}(y)) & \text{if } y \in \text{range}(T_{ji}) \\ \infty & \text{otherwise} \end{cases}$$

786 The grid \mathcal{G} specifies a discrete set of possible values
 787 for $T_{ji}(l_j)$ that are considered during the minimization.
 788 This in turn specifies a discrete set of locations l_j . There
 789 is an approximation being made, since the set of discrete
 790 values for $T_{ji}(l_j)$ (the locations in the grid) might
 791 not match the set of discrete values for $T_{ij}(l_i)$ (where
 792 we need the value of \mathcal{D}_f). We can simply define the
 793 value of the distance transform at a non-grid position
 794 to be the value of the closest grid point. The error introduced
 795 by this approximation is small (as the transform by definition
 796 changes slowly).

797 The same algorithms that efficiently compute the
 798 classical distance transform can be used to compute
 799 the generalized distance transform under different distances,
 800 by replacing the indicator function $1_B(x)$ with
 801 an arbitrary function $f(x)$. In particular we use the
 802 method of Karzanov (originally in Karzanov, 1992, but
 803 see Rucklidge, 1996) for a better description) to compute
 804 the transform of a function under a Mahalanobis distance
 805 with diagonal covariance matrix. This algorithm can also
 806 compute $B'_j(l_i)$, the best location for
 807 v_j as a function of its parent location, as it computes
 808 $B_j(l_i)$.

809 4.2. Sampling from the Posterior

810 We now turn to the problem of sampling from the posterior
 811 distribution of object configurations. The sampling
 812 problem can be solved with a very similar algorithm
 813 to the one described in the previous section. The relationship
 814 between the two cases is analogous to the relationship
 815 between the forward-backward and the Viterbi algorithms
 816 for hidden Markov models. Basically the sampling algorithm
 817 works directly with the probability distributions instead of
 818 their negative logarithms, and the maximizations in the
 819 recursive equations are replaced by summations.

820 As we saw in Section 2 the posterior distribution for
 821 our models is given by

$$p(L | I, \theta) \propto \left(\prod_{i=1}^n p(I | l_i, u_i) \prod_{(v_i, v_j) \in E} p(l_i, l_j | c_{ij}) \right).$$

Like before, let $v_r \in V$ be an arbitrarily chosen root
 823 vertex, and the children of v_i be C_i . The algorithm
 824 works by first computing $p(l_r | I, \theta)$. We then sample
 825 a location for the root from that distribution. Next we
 826 sample a location for each child, v_c , of the root from
 827 $p(l_c | l_r, I, \theta)$. We can continue in this manner until we
 828 have sampled a location for each part. The marginal
 829 distribution for the root location is,
 830

$$p(l_r | I, \theta) \propto \sum_{l_1} \cdots \sum_{l_{r-1}} \sum_{l_{r+1}} \cdots \sum_{l_n} \left(\prod_{i=1}^n p(I | l_i, u_i) \prod_{(v_i, v_j) \in E} p(l_i, l_j | c_{ij}) \right).$$

Computing the distribution in this form would take exponential
 831 time. But since the set of dependencies between parts form a tree,
 832 we can rewrite the distribution as,
 833
 834

$$p(l_r | I, \theta) \propto p(I | l_r, u_r) \prod_{v_c \in C_r} S_c(l_r).$$

The functions $S_j(l_i)$ are similar to the $B_j(l_i)$ we used
 835 for the energy minimization algorithm,
 836

$$S_j(l_i) \propto \sum_{l_j} \left(p(I | l_j, u_j) p(l_i, l_j | c_{ij}) \prod_{v_c \in C_j} S_c(l_j) \right). \quad (13)$$

These recursive functions already give a polynomial algorithm
 837 to compute $p(l_r | I, \theta)$ up to a normalizing constant.
 838 As in the energy minimization algorithm we can compute the
 839 S functions starting from the leaf vertices. The trivial way
 840 to compute each $S_j(l_i)$ takes $O(h^2)$ time. For each location
 841 l_i we evaluate the function by explicitly summing over all
 842 possible locations l_j . We will show how to compute each
 843 $S_j(l_i)$ more efficiently for the case where $p(l_i, l_j | c_{ij})$
 844 is in the special form given by Eq. (6). But first let's see
 845 what we need to do after we sample a location for the root
 846 from its marginal distribution. If we have a location for the
 847 parent v_i of v_j we can write,
 848
 849

$$p(l_j | l_i, I, \theta) \propto p(I | l_j, u_j) p(l_i, l_j | c_{ij}) \prod_{v_c \in C_j} S_c(l_j). \quad (14)$$

If we have already computed the S functions we can
 850 compute this distribution in $O(h)$ time. So once we have
 851 sampled a location for the root, we can sample a
 852

853 location for each of its children. Next we sample a loca-
 854 tion for the nodes at the third level of the tree, and so on
 855 until we sample a location for every part. Note that if we
 856 want to sample multiple times we only need to compute
 857 the S functions once. And when the location of a parent
 858 node is fixed, we only need to compute the distribution
 859 in (14) for locations of the children where $p(l_i, l_j | c_{ij})$
 860 is not too small. So sampling multiple times is not much
 861 more costly than sampling once.

862 **4.2.1. Computing the S Functions.** We want to ef-
 863 ficiently compute the function in Eq. (13). We will do
 864 this by writing the function as a Gaussian convolution in
 865 the transformed space of locations given by T_{ij} and T_{ji} .
 866 Using the special form of $p(l_i, l_j | c_{ij})$ we can write,

$$S_j(l_i) \propto \sum_{l_j} \left(\mathcal{N}(T_{ij}(l_i) - T_{ji}(l_j), 0, D_{ij}) p(I | l_j, u_j) \prod_{v_c \in C_j} S_c(l_j) \right).$$

867 This can be seen as a Gaussian convolution in the trans-
 868 formed space:

$$S_j(l_i) \propto (F \otimes f)(T_{ij}(l_i)),$$

869 where F is a Gaussian filter with covariance D_{ij} , \otimes is
 870 the convolution operator, and

$$f(y) = \begin{cases} p(I | T_{ji}^{-1}(y), u_j) \prod_{v_c \in C_j} S_c(T_{ji}^{-1}(y)) & \text{if } y \in \text{range}(T_{ji}) \\ 0 & \text{otherwise} \end{cases}$$

871 Just like when computing the generalized distance
 872 transform, the convolution is done over a discrete grid
 873 which specifies possible values for $T_{ji}(l_j)$. The Gaus-
 874 sian filter F is separable since the covariance matrix is
 875 diagonal. We can compute a good approximation for
 876 the convolution in time linear in h' , the set of grid loca-
 877 tions, using the techniques from Wells, III (1986). This
 878 gives an overall $O(h'n)$ time algorithm for sampling a
 879 configuration from the posterior distribution.

5. Iconic Models 880

The framework presented so far is general in the sense 881
 that it doesn't fully specify how objects are represented. 882
 A particular modeling scheme must define the pose 883
 space for the object parts, the form of the appearance 884
 model for each part, and the type of connections be- 885
 tween parts. In this section we describe models that rep- 886
 resent objects by the appearance of local image patches 887
 and spatial relationships between those patches. This 888
 type of model has been popular in the context of face 889
 detection (see Fischler and Elschlager, 1973; Burl et al., 890
 1998). We first describe how we model the appearance 891
 of a part, and later describe how we model spatial re- 892
 lationships between parts. Learning an iconic model 893
 involves picking labeled landmarks on a number of inst- 894
 ances of the target object. From these training exam- 895
 ples both the appearance models for each part and the 896
 spatial relationships between parts are automatically 897
 estimated, using the procedure described in Section 3. 898
 In Section 5.3 we show some experiments with face 899
 detection. 900

5.1. Parts 901

In this class of models the location of a part is specified 902
 by its (x, y) position in the image, so we have a two- 903
 dimensional pose space for each part. To model the 904
 appearance of each individual part we use the iconic 905
 representation introduced in Rao and Ballard (1995). 906
 The iconic representation is based on the response of 907
 Gaussian derivative filters of different orders, orienta- 908
 tions and scales. An image patch centered at some po- 909
 sition is represented by a high-dimensional vector that 910
 collects all the responses of a set of filters at that point. 911
 This vector is normalized and called the iconic index 912
 at that position. Figure 2 shows the nine filters used 913
 to build the iconic representation at a fixed scale. In prac- 914
 tice, we use three scales, given by $\sigma_1 = 1$, $\sigma_2 = 2$, and 915
 $\sigma_3 = 4$, the standard deviations of the Gaussian filters. 916
 So we get a 27 dimensional vector. The iconic index is 917
 fairly insensitive to changes in lighting conditions. For 918
 example, it is invariant to gain and bias. Invariance to 919
 bias is a consequence of using image derivative filters, 920

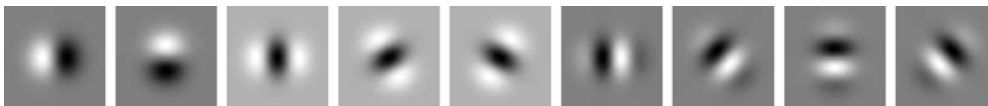


Figure 2. Gaussian derivative basis functions used in the iconic representation.

921 and the normalization provides the invariance to gain.
 922 Iconic indices are also relatively insensitive to small
 923 changes in scale and other image deformations. They
 924 can also be made invariant to image rotation, although
 925 we use an orientation-sensitive representation here.

926 The appearance of a part is modeled by a distribution
 927 over iconic indices. Specifically, we model the distri-
 928 bution of iconic indices at the location of a part as a
 929 Gaussian with diagonal covariance matrix. Using a di-
 930 agonal covariance matrix makes it possible to estimate
 931 the distribution with a small number of examples. If
 932 many examples are available, a full Gaussian or even
 933 more complex distributions such as a mixture of Gaus-
 934 sians, or a non-parametric estimate could be used. Un-
 935 der the Gaussian model, the appearance parameters for
 936 each part are $u_i = (\mu_i, \Sigma_i)$, a mean vector and a co-
 937 variance matrix. We have,

$$p(I | l_i, u_i) \propto \mathcal{N}(\alpha(l_i), \mu_i, \Sigma_i),$$

938 where $\alpha(l_i)$ is the iconic index at location l_i in the im-
 939 age. We can easily estimate the maximum likelihood
 940 parameters of this distribution, as required by the learn-
 941 ing technique in Section 3, using the mean and covari-
 942 ance of the iconic indices corresponding to the positive
 943 examples of a particular part.

944 Note that we could use other methods to repre-
 945 sent the appearance of image patches. In particular,
 946 we experimented with the eigenspace techniques from
 947 Moghaddam and Pentland (1997). With a small num-
 948 ber of training examples the eigenspace methods are no
 949 better than the iconic representation, and the iconic rep-
 950 resentation can be computed more efficiently. In fact,
 951 the iconic representation can be computed very fast by
 952 convolving each level of a Gaussian pyramid with small
 953 x - y separable filters (see Freeman and Adelson, 1991).

954 5.2. Spatial Relations

955 The spatial configuration of the parts is modeled by
 956 a collection of springs connecting pairs of parts. Each
 957 connection (v_i, v_j) is characterized by the ideal relative
 958 location of the two connected parts s_{ij} , and a full co-
 959 variance matrix Σ_{ij} which in some sense corresponds
 960 to the stiffness of the spring connecting the two parts.
 961 So the connection parameters are $c_{ij} = (s_{ij}, \Sigma_{ij})$. We
 962 model the distribution of the relative location of part
 963 v_i with respect to the location of part v_j as a Gaussian
 964 with mean s_{ij} and covariance Σ_{ij} ,

$$p(l_i, l_j | c_{ij}) = \mathcal{N}(l_i - l_j, s_{ij}, \Sigma_{ij}). \quad (15)$$

So, ideally the location of part v_i is the location of 965
 part v_j shifted by s_{ij} . Since the models are deformable, 966
 the location of v_i can vary by paying a cost that de- 967
 pends on the covariance matrix. This corresponds to 968
 stretching the spring. Because we have a full covari- 969
 ance matrix, stretching in different directions can have 970
 different costs. For example, two parts can be highly 971
 constrained to be at the same vertical position, while 972
 their relative horizontal position may be uncertain. As 973
 with the appearance models for the individual parts, the 974
 maximum likelihood parameters of these spatial distri- 975
 butions for pairs of parts can easily be estimated using 976
 training examples. 977

In practice, we need to write the joint distribution of 978
 l_i and l_j in the specific form required by our algorithms. 979
 It must be a Gaussian distribution with zero mean and 980
 diagonal covariance in a transformed space. To do this, 981
 we first compute the singular value decomposition of 982
 the covariance matrix $\Sigma_{ij} = U_{ij} D_{ij} U_{ij}^T$. Now the fol- 983
 lowing transformations can be defined, 984

$$T_{ij}(l_i) = U_{ij}^T(l_i - s_{ij}), \quad \text{and} \quad T_{ji}(l_j) = U_{ij}^T(l_j),$$

which allow us to write Eq. (15) in the correct form, 985

$$p(l_i, l_j | c_{ij}) = \mathcal{N}(T_{ij}(l_i) - T_{ji}(l_j), 0, D_{ij}).$$

5.3. Experiments 986

To test the iconic modes just described we used the ML 987
 estimation procedure from Section 3 to train a model of 988
 frontal faces, and the MAP estimation technique from 989
 Section 4.1 to detect faces in novel images. Our first 990
 model has five parts, corresponding to the eyes, nose, 991
 and corners of the mouth. To generate training exam- 992
 ples we labeled the location of each part in twenty dif- 993
 ferent images (from the Yale face database). More train- 994
 ing examples were automatically generated by scaling 995
 and rotating each training image by a small amount. 996
 This makes our model handle some variation in ori- 997
 entation and scale. Some of the training examples and 998
 the structure of the learned model are shown in Fig. 3. 999
 Remember that we never told the system which pairs 1000
 of parts should be connected together. Determining the 1001
 structure is part of the ML parameter estimation pro- 1002
 cedure. 1003

We tested the resulting model by matching it to novel 1004
 images using the energy minimization algorithm for 1005
 finding the MAP estimate of the object location. Note 1006



Figure 3. Three examples from the first training set showing the locations of the labeled features and the structure of the learned model.

1007 that *all* model parameters were automatically estimated
1008 with the maximum likelihood procedure. Thus, there
1009 are no “knobs” to tune in the matching algorithm. Some
1010 matching results are shown in Fig. 4. Both the learning

and matching algorithms are extremely fast. Using a 1011
desktop computer it took a few seconds to learn the 1012
model and less than a second to compute the MAP es- 1013
timate in each image. These experiments demonstrate 1014



Figure 4. Matching results.



Figure 5. Matching results on occluded faces. The top row shows some input images and the bottom row shows the corresponding matching results. The MAP estimate was a good match when the faces had up to two of five parts occluded and incorrect when three parts were occluded.

1015 that we can learn a useful model from training
1016 examples.

1017 Figure 5 illustrates matching results on images with
1018 partially occluded faces. The matching algorithm au-
1019 tomatically handles such partial occlusion in a robust
1020 way, finding a good configuration of all the parts when
1021 up to two of the five parts are occluded. The occluded
1022 parts are placed at reasonable locations because of the
1023 constraints between parts. Moreover, it does not matter
1024 which parts are occluded because our matching algo-

rithm finds the global minimum of the energy function, 1025
independent of the choice of root used by the dynamic 1026
programming approach. When three of the five parts 1027
are occluded the best match of the model to the image 1028
was incorrect. 1029

Figure 6 illustrates matching results on an image that 1030
contains multiple faces. Recall that our energy mini- 1031
mization algorithm computes the optimal location for 1032
the model as a function of the location of a root part. 1033
To detect multiple faces we first find the best overall 1034



Figure 6. Matching results on an image with multiple faces. See text for description.

1035 location for the root. We then exclude nearby locations
1036 and find the best remaining one and so on for addi-
1037 tional detections. Each root location yields an object
1038 configuration that is optimal with respect to that lo-
1039 cation of the root. In this example we simply found
1040 the best three locations for the model, alternatively a
1041 threshold could be used to find all matches above a cer-
1042 tain quality. Multiple detections could also have been
1043 generated with the sampling techniques together with
1044 a separate verification technique.

1045 We also learned a larger model, this one with nine
1046 parts. We now have three parts for each eye, one for
1047 the left corner, one for the right corner and one for the
1048 pupil. This is a useful model to detect gaze direction.
1049 Figure 7 shows one of the training examples and the
1050 learned model. Also, in Fig. 7, there is a detailed illus-
1051 tration of the connections to the left corner of the right
1052 eye (part 1). The ellipses illustrate the location uncer-
1053 tainty for the other parts, when this part is at some fixed
1054 location. They are level sets of the probability distri-
1055 bution for the location of parts 2, 3, and 4, given that
1056 part 1 is fixed. Note that the location of the pupil (part 2)
1057 is much more constrained with respect to the location
1058 of the eye corner than any other part, as would be ex-
1059 pected intuitively. Also note that the distributions are
1060 not spherically symmetric, as they reflect the typical
1061 variation in the relative locations of parts. We see that
1062 the algorithm both learned an interesting structure for
1063 the model, and automatically determined a rich set of
1064 constraints between the locations of different pairs of
1065 parts.

1066 6. Articulated Models

1067 In this section we present a scheme to model articulated
1068 objects. Our main motivation is to construct a system
1069 that can estimate the pose of human bodies. We concen-

trate on detecting objects in binary images such as those
1070 obtained by background subtraction. Figure 8 shows
1071 an example input and matching result. Binary images
1072 characterize well the problem of pose estimation for an
1073 articulated object. We want to find an object configu-
1074 ration that covers the foreground pixels and leaves the
1075 background pixels uncovered. Our method works with
1076 very noisy input, including substantial occlusion which
1077 we illustrate with examples. Note that in order to de-
1078 tect articulated bodies we use the sampling techniques
1079 in Section 4.2 instead of computing the MAP estimate
1080 for the object location. This is important because the
1081 models for articulated bodies are imprecise rather than
1082 being accurate generative models. 1083

6.1. Parts 1084

For simplicity, we assume that the image of an object
1085 is generated by a scaled orthographic projection, so
1086 that parallel features in the model remain parallel in
1087 the image. For images of human forms this is generally
1088 a reasonable assumption. We further assume that the
1089 scale factor of the projection is known. We can easily
1090 add an extra parameter to our search space in order to
1091 relax this latter assumption. 1092

Suppose that objects are composed of a number of
1093 rigid parts, connected by flexible joints. If a rigid part
1094 is more or less cylindrical, its projection can be ap-
1095 proximated by a rectangle. The width of the rectangle
1096 comes from the diameter of the cylinder and is fixed,
1097 while the length of the rectangle depends on the length
1098 of the cylinder but can vary due to foreshortening. We
1099 model the projection of a part as a rectangle paramete-
1100 rized by (x, y, s, θ) . The center of the rectangle is given
1101 in image coordinates (x, y) , the length is defined by the
1102 amount of foreshortening $s \in [0, 1]$, and the orienta-
1103 tion is given by θ . So we have a four-dimensional pose
1104 space for each part. 1105

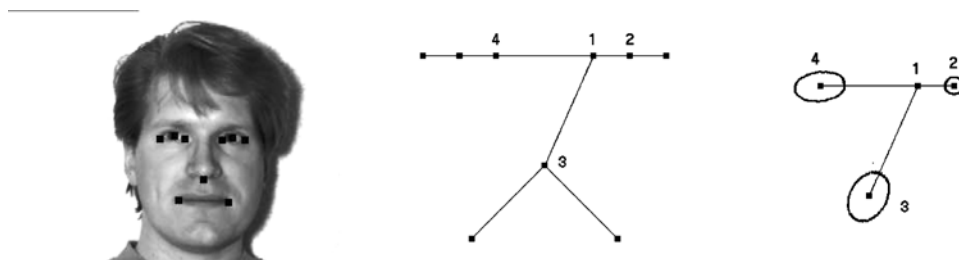


Figure 7. One example from the second training set, the structure of the learned model, and a pictorial illustration of the connections to one of the parts, showing the location uncertainty for parts 2, 3, and 4, when part 1 is at a fixed position.

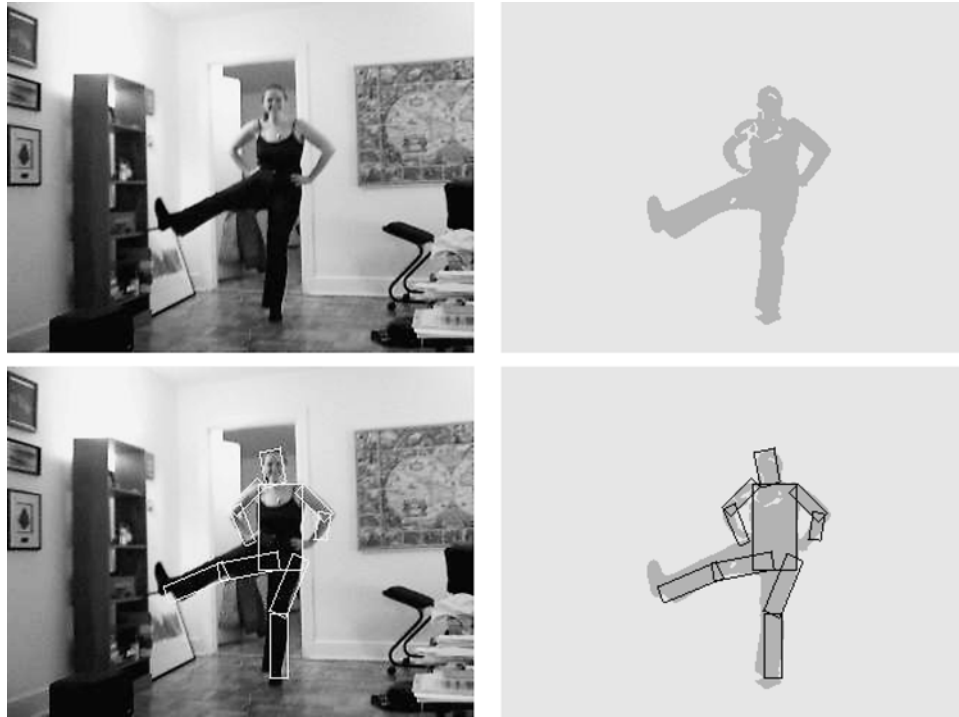


Figure 8. Input image, binary image obtained by background subtraction, and matching result superimposed on both images.

1106 We model the likelihood of observing an image given
 1107 a particular location for a part in the following way.
 1108 First, each pixel in the image is generated independ-
 1109 ently. Pixels inside a part are foreground pixels with
 1110 probability q_1 . Intuitively, q_1 should be close to one, ex-
 1111 pressing the idea that parts occlude the background. We
 1112 also model a border area around each part (see Fig. 9).
 1113 In this area, pixels belong to the foreground with prob-
 1114 ability q_2 . In practice, when we estimate q_2 from data
 1115 we see that pixels around a part tend to be background.
 1116 We assume that pixels outside both areas are equally
 1117 likely to be background or foreground pixels. Thus,

$$p(I | l_i, u_i) = q_1^{\text{count}_1} (1 - q_1)^{(\text{area}_1 - \text{count}_1)} q_2^{\text{count}_2} \times (1 - q_2)^{(\text{area}_2 - \text{count}_2)} 0.5^{(t - \text{area}_1 - \text{area}_2)},$$

1118 where count_1 is the number of foreground pixels inside
 1119 the rectangle, and area_1 is the area of the rectangle.
 1120 count_2 and area_2 are similar measures corresponding to
 1121 the border area, and t is the total number of pixels in the
 1122 image. So the appearance parameters are $u_i = (q_1, q_2)$,
 1123 and it is straightforward to estimate these parameters
 1124 from training examples.

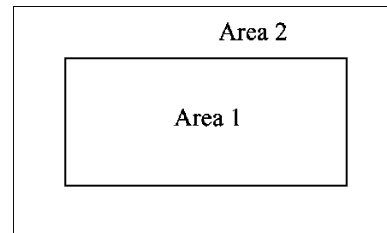


Figure 9. A rectangular part. area_1 is the area inside the part, and area_2 is the border area around it.

To make the probability measure robust we consider 1125
 a slightly dilated version of the foreground when com- 1126
 1127 puting count_1 , and to compute count_2 we erode the
 1128 foreground (in practice we dilate and erode the binary
 1129 images by two pixels). Computing the likelihood for
 1130 every possible location of a part can be done efficiently
 1131 by convolving the image with uniform filters. Each con-
 1132 volution counts the number of pixels inside a rectangle
 1133 (specified by the filter) at every possible translation.
 1134 Intuitively, our model of $p(I | l_i, u_i)$ is reasonable
 1135 for a single part. The likelihood favors large parts,
 1136 as they explain a larger area of the image. But re-
 1137 member that we model $p(I | L, u)$ as a product of the

1138 individual likelihoods for each part. For a configuration
 1139 with overlapping parts, this measure “over-counts” ev-
 1140 idence. Suppose we have an object with two parts. The
 1141 likelihood of an image is the same if the two parts are
 1142 arranged to explain different areas of the image, or if
 1143 the two parts are on top of each other and explain the
 1144 same area twice. Therefore, with this measure the MAP
 1145 estimate of an object configuration can be a bad guess
 1146 for its true position. This is not because the posterior
 1147 probability of the true configuration is low, but because
 1148 there are configurations which have high posterior and
 1149 are wrong. In our experiments, we obtain a number of
 1150 configurations which have high posterior probability
 1151 by sampling from that distribution. We then select one
 1152 of the samples by computing a quality measure that
 1153 does not over-count evidence.

1154 There is one more thing we have to take into account
 1155 for sampling to work. When $p(I | L, u)$ over-counts ev-
 1156 idence, it tends to create high peaks. This in turn creates
 1157 high peaks in the posterior. When a distribution has a
 1158 very strong peak, sampling from the distribution will
 1159 almost always obtain the location of the peak. To en-
 1160 sure that we get a number of different hypotheses from
 1161 sampling we use a smoothed version of the likelihood
 1162 function, defined as

$$p'(I | L, u) \propto p(I | L, u)^{1/T} \propto \prod_{i=1}^n p(I | l_i, u_i)^{1/T},$$

1163 where T controls the degree of smoothing. This is a
 1164 standard technique, borrowed from the principle of an-
 1165 nealing (see Geman and Geman, 1984). In all our ex-
 1166 periments we used $T = 10$.

1167 6.2. Spatial Relations

1168 For the articulated objects, pairs of parts are connected
 1169 by flexible joints. A pair of connected parts is illus-
 1170 trated in Fig. 10. The location of the joint is specified
 1171 by two points (x_{ij}, y_{ij}) and (x_{ji}, y_{ji}) , one in the co-
 1172 ordinate frame of each part, as indicated by circles in
 1173 Fig. 10(a). In an ideal configuration these points co-
 1174 incide, as illustrated in Fig. 10(b). The ideal relative
 1175 orientation is given by θ_{ij} , the difference between the
 1176 orientation of the two parts.

1177 Suppose $l_i = (x_i, y_i, s_i, \theta_i)$ and $l_j = (x_j, y_j, s_j, \theta_j)$
 1178 are the locations of two connected parts. The joint prob-
 1179 ability for the two locations is based on the deviation
 1180 between their ideal relative values and the observed

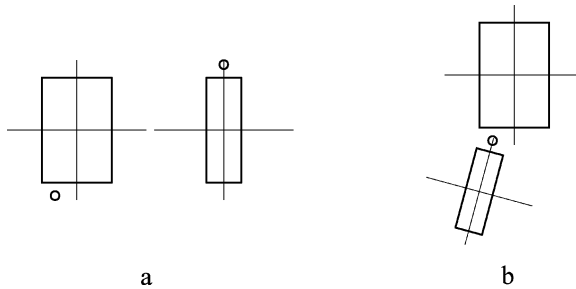


Figure 10. Two parts of an articulated object, (a) in their own coordinate system and (b) the ideal configuration of the pair.

ones, 1181

$$p(l_i, l_j | c_{ij}) = \mathcal{N}(x'_i - x'_j, 0, \sigma_x^2) \mathcal{N}(y'_i - y'_j, 0, \sigma_y^2) \mathcal{N}(s_i - s_j, 0, \sigma_s^2) \mathcal{M}(\theta_i - \theta_j, \theta_{ij}, k), \quad (16)$$

where (x'_i, y'_i) and (x'_j, y'_j) are the positions of the joints 1182 in image coordinates. Let R_θ be the matrix that per- 1183 forms a rotation of θ radians about the origin. Then, 1184

$$\begin{aligned} \begin{bmatrix} x'_i \\ y'_i \end{bmatrix} &= \begin{bmatrix} x_i \\ y_i \end{bmatrix} + s_i R_{\theta_i} \begin{bmatrix} x_{ij} \\ y_{ij} \end{bmatrix}, \quad \text{and} \quad \begin{bmatrix} x'_j \\ y'_j \end{bmatrix} \\ &= \begin{bmatrix} x_j \\ y_j \end{bmatrix} + s_j R_{\theta_j} \begin{bmatrix} x_{ji} \\ y_{ji} \end{bmatrix}. \end{aligned}$$

The distribution over angles, \mathcal{M} , is the von Mises dis- 1185 tribution (Gumbel et al., 1953), 1186

$$\mathcal{M}(\theta, \mu, k) \propto e^{k \cos(\theta - \mu)}.$$

The first two terms in the joint distribution measure the 1187 horizontal and vertical distances between the observed 1188 joint positions in the image. The third term measures 1189 the difference in foreshortening between the two parts. 1190 The last term measures the difference between the rel- 1191 ative angle of the two parts and the ideal relative angle. 1192 Usually σ_x and σ_y will be small so parts tend to be 1193 aligned at their joint. And if k is small the angle be- 1194 tween the two parts is fairly unconstrained, modeling 1195 a revolute joint. The connection parameters under this 1196 model are, 1197

$$c_{ij} = (x_{ij}, y_{ij}, x_{ji}, y_{ji}, \sigma_x^2, \sigma_y^2, \sigma_s^2, \theta_{ij}, k).$$

1198 Finding the maximum likelihood estimate of σ_s^2 is easy
 1199 since we just have a Gaussian distribution over $s_i -$
 1200 s_j . Similarly, there are known methods for finding the
 1201 ML parameters (θ_{ij}, k) of a von Mises distribution (see
 1202 Gumbel et al., 1953). The ML estimate of the joint
 1203 location in each part are the values $(x_{ij}, y_{ij}, x_{ji}, y_{ji})$
 1204 which minimize the sum of square distances between
 1205 (x'_i, y'_i) and (x'_j, y'_j) over the examples. We can compute
 1206 this as a linear least squares problem.

1207 We need to write the joint distribution of l_i and l_j in
 1208 the specific form required by our algorithms. It must
 1209 be a Gaussian distribution with zero mean and diagonal
 1210 covariance in a transformed space, as shown in Eq. (6).
 1211 First note that a von Mises distribution over angular
 1212 parameters can be specified in terms of a Gaussian over
 1213 the unit vector representation of the angles. Let $\vec{\alpha}$ and $\vec{\beta}$
 1214 be the unit vectors corresponding to two angles α and
 1215 β . That is, $\vec{\alpha} = [\cos(\alpha), \sin(\alpha)]^T$, and similarly for $\vec{\beta}$.
 1216 Then,

$$\cos(\alpha - \beta) = \vec{\alpha} \cdot \vec{\beta} = -\frac{\|\vec{\alpha} - \vec{\beta}\|^2 - 2}{2}.$$

1217 Now let

$$\begin{aligned} T_{ij}(l_i) &= (x'_i, y'_i, s_i, \cos(\theta_i + \theta_{ij}), \sin(\theta_i + \theta_{ij})), \\ T_{ji}(l_j) &= (x'_j, y'_j, s_j, \cos(\theta_j), \sin(\theta_j)), \\ D_{ij} &= \text{diag}(\sigma_x^2, \sigma_y^2, \sigma_s^2, 1/k, 1/k), \end{aligned}$$

1218 which allow us to write Eq. (16) in the right form,

$$p(l_i, l_j | c_{ij}) \propto \mathcal{N}(T_{ji}(l_j) - T_{ij}(l_i), 0, D_{ij}).$$

1219 For these models, the number of discrete locations h' in
 1220 the transformed space is a little larger than the number
 1221 of locations h for each part. This is because we repre-
 1222 sent the orientation of a part as a unit vector which lives
 1223 in a two-dimensional grid. In practice, we use 32 pos-
 1224 sible angles for each part, and represent them as points
 1225 in a 11×11 grid, which makes h' about four times h .

1226 6.3. Experiments

1227 We use a coarse articulated model to represent the hu-
 1228 man body. Our model has ten parts, corresponding to
 1229 the torso, head, two parts per arm and two parts per
 1230 leg. To generate training examples we labeled the lo-
 1231 cation of each part in ten different images (without too
 1232 much precision). The learned model is illustrated in
 1233 Fig. 11. The crosses indicate joints between parts. We

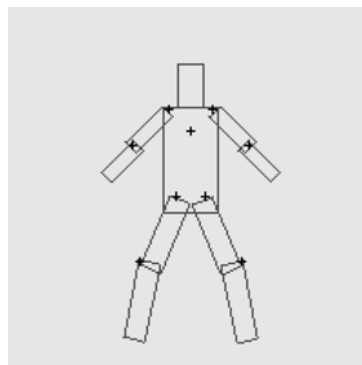


Figure 11. Human body model learned from example configurations.

never told the system which parts should be connected
 together, this is automatically learned during the ML
 learning procedure. Note that the correct structure was
 learned, and the joint locations agree with the human
 body anatomy (the joint in the middle of the torso con-
 nects to the head). The configuration of parts shown in
 Fig. 11 was obtained by fixing the position of the torso
 and placing all other parts in their optimal location with
 respect to each other.

We tested the model by matching it to novel im-
 ages. As described in Section 6.1, we sample config-
 urations from the posterior distribution to obtain mul-
 tiple hypotheses and rate each sample using a sepa-
 rate measure. For each sample we compute the Cham-
 fer distance between the shape of the object under the
 hypothesized configuration and the binary image ob-
 tained from the input. The Chamfer distance is a robust
 measure of binary correlation (Borgefors, 1988). The
 matching process is illustrated in Fig. 12. First, a binary
 image is obtained from the original image using back-
 ground subtraction. We use this binary image as input
 to the sampling algorithm to obtain a number of dif-
 ferent pose hypotheses. The best pose is then selected
 using the Chamfer measure.

More matching results are shown in Fig. 13. For
 each image, we sampled two-hundred object configu-
 rations from the posterior distribution and picked the
 best one under the Chamfer distance. Using a desk-
 top computer it took about one minute to process each
 image. The space of possible locations for each part
 was discretized into a $70 \times 70 \times 10 \times 32$ grid, corre-
 sponding to (x, y, s, θ) parameters. There are over 1.5
 million locations for each part, making any algorithm
 that considers locations for pairs of parts at a time im-
 practical.

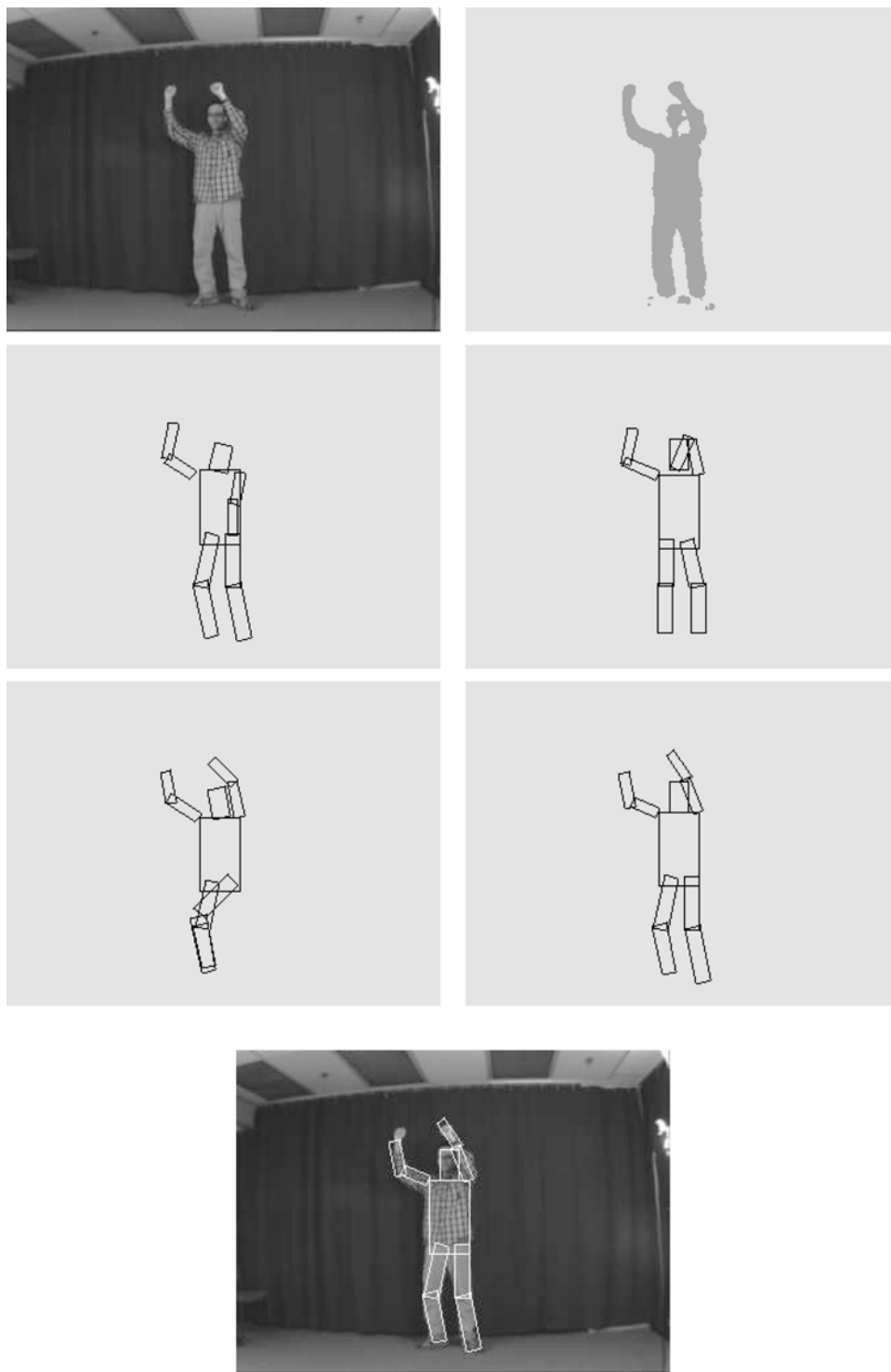


Figure 12. Input image, binary image, random samples from the posterior distribution of configurations, and best result selected using the Chamfer distance.

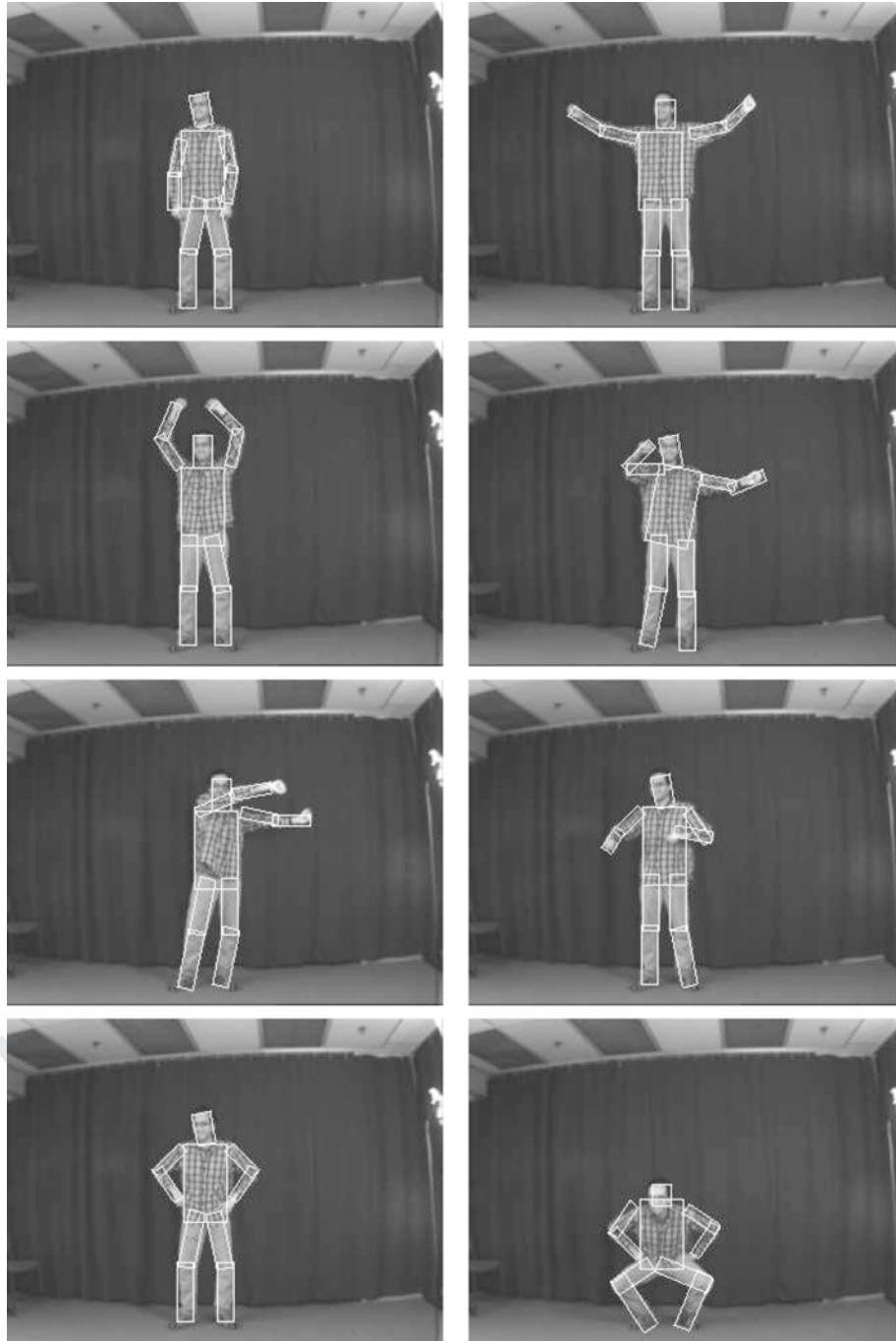


Figure 13. Matching results (sampling 200 times).

1269 Of course, sometimes the estimated pose is not cor-
 1270 rect. The most common source of error comes from
 1271 ambiguities in the binary images. Figure 14 shows an
 1272 example where the image does not provide enough in-

formation to estimate the position of one arm. Even in 1273
 that case we get a fairly good estimate. We can detect 1274
 when ambiguities happen because we obtain many dif- 1275
 ferent poses with equally good Chamfer score. Thus 1276

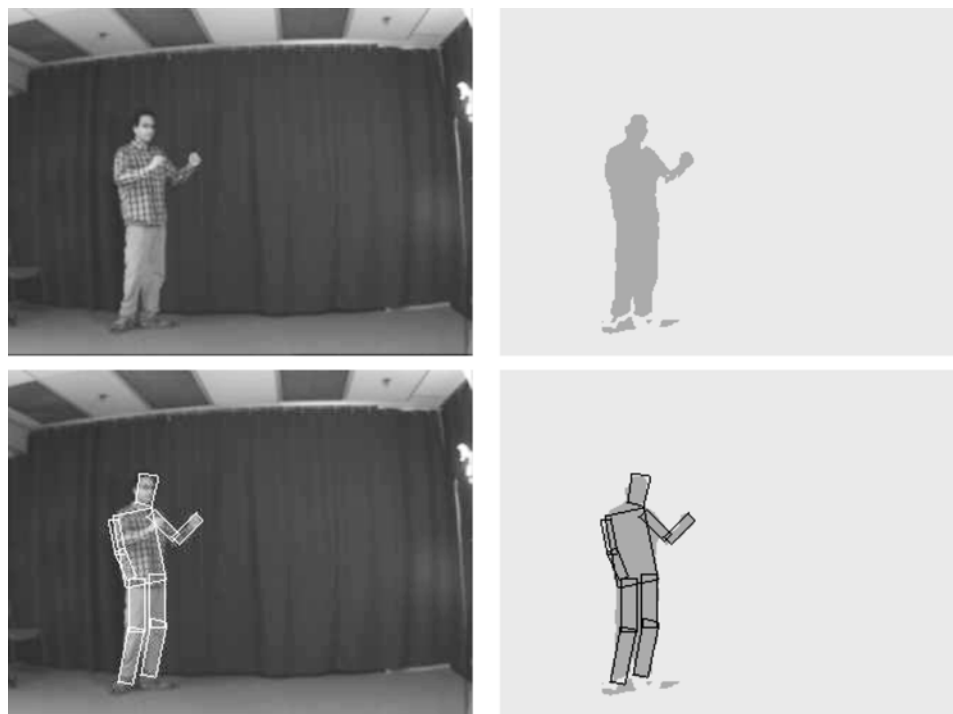


Figure 14. In this case, the binary image doesn't provide enough information to estimate the position of one arm.

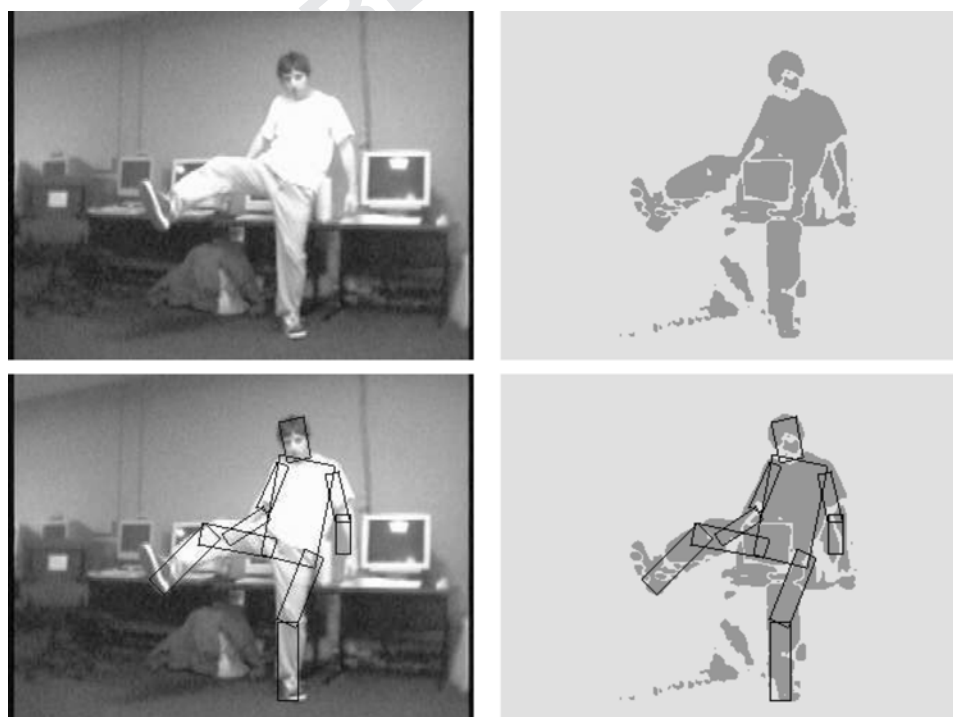


Figure 15. This example illustrates how our method works well with noisy images.

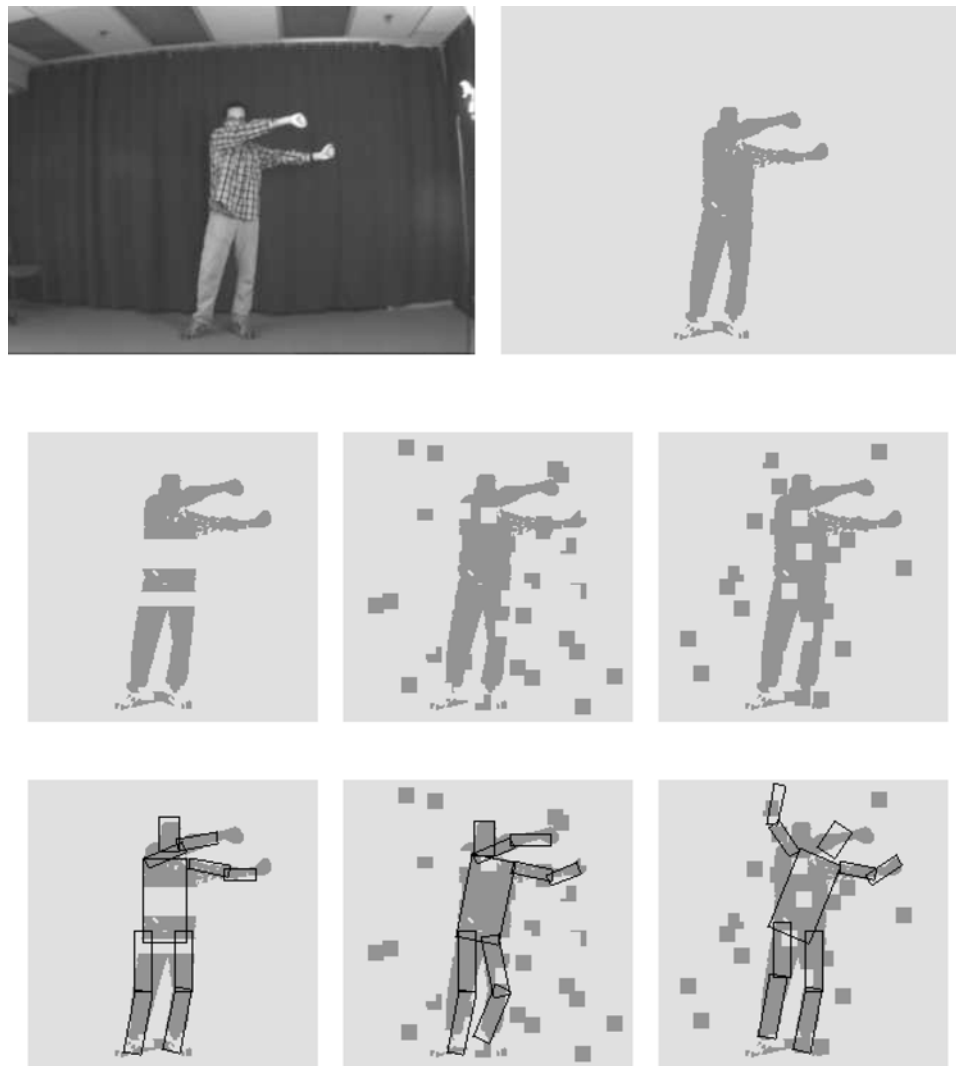


Figure 16. Matching results on corrupted images. The top row shows the original input, the middle row shows corrupted versions of the binary input and the last row shows the matching results. The first two cases demonstrate how the algorithm can handle good amounts of noise and occlusion. The third case shows an incorrect matching result.

1277 we know that there are different configurations that are
1278 equally good interpretations of the image.

1279 Figure 15 shows how our method works well on a
1280 noisy input. More examples of matching to noisy inputs
1281 are shown in Fig. 16, using corrupted binary images,
1282 including a case where large portions of the foreground
1283 are missing. These examples illustrate two of the main
1284 advantages of our approach. It would be difficult to de-
1285 tect body parts individually on inputs such as these, but
1286 the dependencies between parts provide sufficient con-
1287 text to detect the human body as a whole. Moreover, the
1288 presence of clutter and occlusion create difficulties for

heuristics or local search techniques, while our global
method can find the correct configuration in these cases.

7. Summary

1291

This paper describes a statistical framework for rep-
resenting the visual appearance of objects composed
of rigid parts arranged in a deformable configuration.
The models are based on the pictorial structure repre-
sentation introduced in Fischler and Elschlager (1973),
which allows for qualitative descriptions of appearance

1298 and is suitable for generic recognition problems. There
 1299 are three main contributions in the paper. First, we in-
 1300 troduce efficient algorithms for finding the best *global*
 1301 match of a large class of pictorial structure models to
 1302 an image. In contrast, prior work use heuristics or lo-
 1303 cal search techniques that must be somehow initial-
 1304 ized near the right answer. Second, we introduce the
 1305 use of statistical sampling techniques to identify mul-
 1306 tiple good matches of a model to an image. Third, our
 1307 use of a statistical formulation provides a natural way
 1308 of learning pictorial structure models from labeled ex-
 1309 ample images. Most of the prior work uses manually
 1310 constructed models, which are difficult to create and to
 1311 validate.

1312 One of the difficulties in representing generic objects
 1313 is the large variation in shape and photometric infor-
 1314 mation in each object class. Pictorial structure models
 1315 represent the appearance of each part separately and
 1316 explicitly capture the spatial configuration of the parts
 1317 independently of their appearances. This framework is
 1318 general, in the sense that it is independent of the spe-
 1319 cific method used to represent the appearance of parts,
 1320 and the type of the geometric relationships between
 1321 the parts. By using a general framework we have pro-
 1322 vided a set of computational mechanisms that can be
 1323 used for many different modeling schemes. We have de-
 1324 scribed two quite different modeling schemes, one was
 1325 used to model faces and the other to model articulated
 1326 bodies.

1327 Acknowledgments

1328 We would like to thank Romer Rosales for providing a
 1329 database of human body images with a variety of poses.

1330 References

- 1331 Amini, A.A., Weymouth, T.E., and Jain, R.C. 1990. Using dy-
 1332 namic programming for solving variational problems in vision.
 1333 *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
 1334 12(9):855–867.
- 1335 Amit, Y. and Geman, D. 1999. A computational model for visual
 1336 selection. *Neural Computation*, 11(7):1691–1715.
- 1337 Ayache, N.J. and Faugeras, O.D. 1986. Hyper: A new approach
 1338 for the recognition and positioning of two-dimensional objects.
 1339 *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
 1340 8(1):44–54.
- 1341 Berger, J.O. 1985. *Statistical Decision Theory and Bayesian Analy-*
 1342 *sis*. Springer-Verlag.
- 1343 Borgefors, G. 1986. Distance transformations in digital images. *Com-*
 1344 *puter Vision, Graphics, and Image Processing*, 34(3):344–371.
- Borgefors, G. 1988. Hierarchical chamfer matching: A parametric
 edge matching algorithm. *IEEE Transactions on Pattern Analysis*
 and *Machine Intelligence*, 10(6):849–865.
- Boykov, Y., Veksler, O., and Zabih, R. 2001. Fast approximate en-
 ergy minimization via graph cuts. *IEEE Transactions on Pattern*
Analysis and Machine Intelligence, 23(11):1222–1239.
- Bregler, C. and Malik, J. 1998. Tracking people with twists and
 exponential maps. In *IEEE Conference on Computer Vision and*
Pattern Recognition, pp. 8–15.
- Burl, M.C. and Perona, P. 1996. Recognition of planar object classes.
 In *IEEE Conference on Computer Vision and Pattern Recognition*,
 pp. 223–230.
- Burl, M.C., Weber, M., and Perona, P. 1998. A probabilistic approach
 to object recognition using local photometry and global geometry.
 In *European Conference on Computer Vision*, pp. II:628–641.
- Chow, C.K. and Liu, C.N. 1968. Approximating discrete probability
 distributions with dependence trees. *IEEE Transactions on Infor-*
mation Theory, 14(3):462–467.
- Cormen, T.H., Leiserson, C.E., and Rivest, R.L. 1996. *Introduction*
 to *Algorithms*. MIT Press and McGraw-Hill.
- Dickinson, S.J., Biederman, I., Pentland, A.P., Eklundh, J.O.,
 Bergevin, R., and Munck-Fairwood, R.C. 1993. The use of geons
 for generic 3-d object recognition. In *International Joint Confer-*
ence on Artificial Intelligence, pp. 1693–1699.
- Felzenszwalb, P.F. and Huttenlocher, D.P. 2000. Efficient matching
 of pictorial structures. In *IEEE Conference on Computer Vision*
 and *Pattern Recognition*, pp. II:66–73.
- Fischler, M.A. and Bolles, R.C. 1981. Random sample consensus: A
 paradigm for model fitting with applications to image analysis and
 automated cartography. *Communications of the ACM*, 24(6):381–
 395.
- Fischler, M.A. and Elschlager, R.A. 1973. The representation and
 matching of pictorial structures. *IEEE Transactions on Computer*,
 22(1):67–92.
- Freeman, W.T. and Adelson, E.H. 1991. The design and use of steer-
 able filters. *IEEE Transactions on Pattern Analysis and Machine*
Intelligence, 13(9):891–906.
- Gdalyahu, Y. and Weinshall, D. 1999. Flexible syntactic matching
 of curves and its application to automatic hierarchical classifica-
 tion of silhouettes. *IEEE Transactions on Pattern Analysis and*
Machine Intelligence, 21(12):1312–1328.
- Geman, S. and Geman, D. 1984. Stochastic relaxation, Gibbs dis-
 tributions, and the Bayesian restoration of images. *IEEE Trans-*
actions on Pattern Analysis and Machine Intelligence, 6(6):721–
 741.
- Grimson, W.E.L. and Lozano-Perez, T. 1987. Localizing overlapping
 parts by searching the interpretation tree. *IEEE Transactions on*
Pattern Analysis and Machine Intelligence, 9(4):469–482.
- Gumbel, E.J., Greenwood, J.A., and Durand, D. 1953. The circular
 normal distribution: Theory and tables. *Journal of the American*
Statistical Association, 48:131–152.
- Huttenlocher, D.P., Klanderman, G.A., and Rucklidge, W.J. 1993.
 Comparing images using the hausdorff distance. *IEEE Transac-*
tions on Pattern Analysis and Machine Intelligence, 15(9):850–
 863.
- Huttenlocher, D.P. and Ullman, S. 1990. Recognizing solid objects
 by alignment with an image. *International Journal of Computer*
Vision, 5(2):195–212.
- Ioffe, S. and Forsyth, D.A. 2001. Probabilistic methods for finding
 people. *International Journal of Computer Vision*, 43(1):45–68.

- 1405 Ishikawa, H. and Geiger, D. 1998. Segmentation by grouping junctions. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 125–131.
- 1406
- 1407
- 1408 Ju, S.X., Black, M.J., and Yacoob, Y. 1996. Cardboard people: A parameterized model of articulated motion. In *International Conference on Automatic Face and Gesture Recognition*, pp. 38–44.
- 1409
- 1410 Karzanov, A.V. 1992. Quick algorithm for determining the distances from the points of the given subset of an integer lattice to the points of its complement. *Cybernetics and System Analysis*, pp. 177–181.
- 1411
- 1412 Translation from the Russian by Julia Komissarchik.
- 1413
- 1414 Lamdan, Y., Schwartz, J.T., and Wolfson, H.J. 1990. Affine invariant model-based object recognition. *IEEE Transactions on Robotics and Automation*, 6(5):578–589.
- 1415
- 1416 Moghaddam, B. and Pentland, A.P. 1997. Probabilistic visual learning for object representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(7):696–710.
- 1417
- 1418 Murase, H. and Nayar, S.K. 1995. Visual learning and recognition of 3-d objects from appearance. *International Journal of Computer Vision*, 14(1):5–24.
- 1419
- 1420 Pearl, J. 1988. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann.
- 1421
- 1422 Pentland, A.P. 1987. Recognition by parts. In *IEEE International Conference on Computer Vision*, pp. 612–620.
- 1423
- 1424 Rabiner, L. and Juang, B. 1993. *Fundamentals of Speech Recognition*. Prentice Hall. 1428
- 1429
- 1430 Ramanan, D. and Forsyth, D.A. 2003. Finding and tracking people from the bottom up. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. II: 467–474. 1431
- 1432
- 1433 Rao, R.P.N. and Ballard, D.H. 1995. An active vision architecture based on iconic representations. *Artificial Intelligence*, 78(1/2):461–505. 1434
- 1435
- 1436 Rivlin, E., Dickinson, S.J., and Rosenfeld, A. Recognition by functional parts. *Computer Vision and Image Understanding*, 62(2):164–176, September 1995. 1437
- 1438
- 1439 Roberts, L.G. 1965. Machine perception of 3-d solids. In *Optical and Electro-optical Information Processing*, pp. 159–197. 1440
- 1441
- 1442 Rucklidge, W. 1996. *Efficient Visual Recognition Using the Hausdorff Distance*. Springer-Verlag, LNCS 1173. 1443
- 1444
- 1445 Sebastian, T.B., Klein, P.N., and Kimia, B.B. 2001. Recognition of shapes by editing shock graphs. In *IEEE International Conference on Computer Vision*, pp. I:755–762. 1446
- 1447
- 1448 Turk, M. and Pentland, A.P. 1991. Eigenfaces for recognition. *Journal of Cognitive Neuroscience*, 3(1):71–96. 1449
- 1450
- 1451 Wells, W.M. III 1986. Efficient synthesis of Gaussian filters by cascaded uniform filters. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 8(2):234–239. 1452

Copyright of International Journal of Computer Vision is the property of Kluwer Academic Publishing / Academic and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.