



# Articulated Body Motion Capture by Stochastic Search

JONATHAN DEUTSCHER AND IAN REID

*Department of Engineering Science, University of Oxford, Oxford, OX13PJ, United Kingdom*

jdeutsch@robots.ox.ac.uk

ian.reid@eng.ox.ac.uk

*Received August 19, 2003; Revised April 9, 2004; Accepted April 9, 2004*

*First online version published in October, 2004*

**Abstract.** We develop a modified particle filter which is shown to be effective at searching the high-dimensional configuration spaces (c. 30 + dimensions) encountered in visual tracking of articulated body motion. The algorithm uses a continuation principle, based on annealing, to introduce the influence of narrow peaks in the fitness function, gradually. The new algorithm, termed annealed particle filtering, is shown to be capable of recovering full articulated body motion efficiently. A mechanism for achieving a soft partitioning of the search space is described and implemented, and shown to improve the algorithm's performance. Likewise, the introduction of a crossover operator is shown to improve the effectiveness of the search for kinematic trees (such as a human body). Results are given for a variety of agile motions such as walking, running and jumping.

**Keywords:** human motion capture, visual tracking, particle filtering, genetic algorithms

## 1. Introduction

A popular form of motion capture, for tasks such as gait analysis and computer animation, involves attaching a number of retro-reflective markers to a subject's body and viewing the motion of the markers over time using a set of calibrated cameras. The easily-recovered image positions of the markers are transformed into 3D trajectories via triangulation of the measurements, and a parameterised representation of the subject's movements can be calculated.

The use of markers is intrusive and restrictive, and necessitates the use of potentially expensive, specialised capture hardware. The goal of markerless motion capture is to reproduce the performance of marker-based methods in a system using conventional cameras and without the use of special apparel or other equipment. For this reason recent years have seen a huge growth in research in the computer vision community with the aim of recovering motion data directly from

images, without markers. However, full-body tracking from standard images is a challenging problem, and markerless system presented to date rarely achieve the following combination of capabilities of current marker-based systems: full 3D motion recovery; robust tracking of rapid, arbitrary movement; high accuracy; easy application to new scenarios; on-line model acquisition; real-time, or near real-time processing.

A major problem which confronts all attempts to satisfy these criteria is the high dimensionality of the configuration space, and the exponentially increasing computational cost that results. A realistic articulated model (see Fig. 4) of the human body usually has at least 25 DOF. The model used in this paper for example has between 29 and 34 DOF, and models employed for commercial character animation usually have over 40.

In this paper we describe a multi-camera system for markerless human motion capture which goes some way to achieving the goals above. The work described combines and extends our previous efforts published

56	in short form in Deutscher et al. (2000, 2001). Our ap-	101
57	proach is characterised by the following: 1. articulated	102
58	body model, 2. weak dynamical modelling, 3. edge	103
59	and background subtraction image measurements, and	104
60	4. a particle-based stochastic search algorithm. The key	105
61	contributions comprise:	
62	• The development of a novel, particle-based stochas-	106
63	tic search algorithm, called annealed particle filter-	107
64	ing. The method uses a continuation principle, based	108
65	on annealing, to introduce the influence of narrow	109
66	peaks in the fitness function gradually. This is intro-	110
67	duced in Section 3.	111
68	• The annealed particle filter is applied to the problem	112
69	of markerless human motion capture, and shown to	113
70	be more effective and efficient than, for example,	114
71	Condensation (Blake and Isard, 1998), at localising	115
72	the pose. Section 4 discusses implementation issues	116
73	and Section 5 shows results.	117
74	• We demonstrate how adaptive selection of the vari-	118
75	ances/covariances which control the diffusion dur-	119
76	ing annealing can lead to what can be thought of as a	120
77	“soft” hierarchical partitioning of the configuration	121
78	space, and hence to further gains in efficiency.	122
79	• We introduce a crossover operator, analogous that	123
80	that found in Genetic Algorithms, into the particle	124
81	filtering framework. We demonstrate that this opera-	125
82	tor improves the ability of the algorithm to search the	126
83	configuration spaces of objects whose articulations	127
84	can be modelled as a kinematic tree. In particular	128
85	we show results for reliable and efficient tracking	129
86	for walking, running and jumping with no special	130
87	training of the dynamics of any activity.	131
88	These latter two developments are discussed in	132
89	Section 6.	133
90	We begin with a review of relevant literature in	134
91	Section 2, including a detailed discussion of the two	135
92	most closely associated technologies: particle filtering	136
93	and simulated annealing.	137
94	<b>2. Background</b>	138
95	<i>2.1. Visual Tracking and Particle Filters</i>	139
96	Full-body motion capture is an example of model-	140
97	-based tracking; i.e. the process of sequentially esti-	141
98	imating the parameters of a model of a target over	142
99	time from visual data. Typically a priori knowledge	143
100	about the target’s observable properties (such as its	144
	geometry) are compared with the visual data from an	145
	image stream, to estimate a best fit for each frame in the	146
	scene. Thus the principal components usually comprise	147
	a target model, an image search method, and a dynam-	148
	ical model.	149
	The system of Plänkners and Fua (2003) represents	150
	one of the best examples of this paradigm in its simplest	
	form. Their system does not have a (strong) model of	
	the person’s dynamics (in contrast to Sidenbladh et al.	
	(2000), e.g., see below) or have a sophisticated multi-	
	modal search algorithm such as we describe. Rather,	
	the key to the success of their system is in much more	
	careful modelling of the shape and appearance than in	
	most other work, and in the use of binocular disparity	
	maps as well as silhouette data. Unlike most other work	
	(including our own) their system estimates the size of	
	the body as well as the pose parameters.	
	In considering the role of the other two system com-	
	ponents, search and dynamics it is useful to discuss the	
	influential work of Harris (1992). He showed how rapid	
	motions can be tracked by constraining the search area	
	via a predicted motion of the object. Harris used rigid	
	polyhedral models (and simple surfaces of revolution)	
	and sought the 6 DOF pose of object. Given a predicted	
	location, the system searches from predicted edge lo-	
	cations along 1D profiles to find “actual” edges. These	
	1D measurements are then combined to obtain a pose	
	update.	
	Drummond and Cipolla (2001) showed how many	
	of the ideas in Harris’ system can be extended to artic-	
	ulated objects by effectively tracking body parts using	
	Harris’ method but enforcing global consistency via	
	kinematic constraints.	
	A second and arguably more important innovation in	
	Harris (1992) was to place the tracking system within	
	the framework of a Kalman Filter, a provably optimal	
	recursive estimator for linear systems which can be	
	thought of as an algorithm for sequential propagation	
	of Gaussian probability densities.	
	A natural step would be to consider the use of a	
	Kalman Filter, (or its extension to non-linear measure-	
	ments and/or dynamics, the Extended Kalman Filter	
	or EKF) for articulated body tracking. Wachter and	
	Nagel (1999) demonstrated this for single view track-	
	ing using image motion and edges (though the results	
	show only motion parallel to the image plane). More	
	recently Mikic et al. (2001) demonstrated the extrac-	
	tion and filtering of pose parameters from a volumetric	
	model obtained by “carving” space using silhouettes	
	from multiple cameras.	

151 While Gaussian uncertainty is sufficient for mod-  
 152 elling many motion and measurement noise sources,  
 153 the Kalman Filter has been shown to fail catastrophi-  
 154 cally in cases where the true probability function has  
 155 a very different shape. In particular attempts to track  
 156 objects moving against a very cluttered background,  
 157 where measurement densities include the chance of  
 158 detecting erroneous image features and are therefore  
 159 multi-modal, lead to tracking failure for Harris' algo-  
 160 rithm and many of its ilk.

161 An alternative, more general approach is particle fil-  
 162 tering, in which arbitrary densities are approximated.  
 163 This was introduced in the context of visual tracking  
 164 in the form of the Condensation algorithm (Isard and  
 165 Blake, 1996). A posterior density  $p(\mathbf{X} | \mathbf{Z}_t)$  represent-  
 166 ing current knowledge about the model state  $\mathbf{X}$  after  
 167 incorporation of all measurements up to the current  
 168 time-step  $t$ ,  $\mathbf{Z}_t$ , is represented by a finite set of *nor-*  
 169 *malised weighted particles*, or samples,

$$\{(\mathbf{s}_t^{(0)}, \pi_t^{(0)}) \cdots (\mathbf{s}_t^{(N)}, \pi_t^{(N)})\}.$$

170 An estimate of the state  $\mathcal{X}_t$  at each time-step  $t$  can  
 171 easily be estimated by the sample mean of the posterior  
 172 density,  $p(\mathbf{X} | \mathbf{Z}_t)$ ,

$$\mathcal{X}_t = \mathcal{E}[\mathbf{X}] = \sum_{n=1}^N \pi_t^{(n)} \mathbf{s}_t^{(n)} \quad (1)$$

or the mode

$$\mathcal{X}_t = \mathcal{M}[\mathbf{X}] = \mathbf{s}_t^{(j)}, \quad \pi_t^{(j)} = \max(\pi_t^{(n)}). \quad (2)$$

173 Essentially, a smooth probability density function is  
 174 approximated by a finite collection of weighted sam-  
 175 ple points, and it can be shown that as the number of  
 176 points tends to infinity the behaviour of the particle set  
 177 is indistinguishable from that of the smooth function.  
 178 Tracking with a particle filter works by:

- 179 1. Resampling, in which a weighted particle set is  
 180 transformed into a set of evenly weighted particles  
 181 distributed with concentration dependent on proba-  
 182 bility density;
- 183 2. Stochastic movement and dispersion of the particle  
 184 set in accordance with a motion model to represent  
 185 the growth of uncertainty during movement of the  
 186 tracked object;
- 187 3. Measurement, in which the likelihood function is  
 188 evaluated at each particle site, producing a new

weight for each particle proportional to how well it  
 fits image data. The weighted particle set produced  
 represents the new (posterior) probability density  
 after movement and measurement.

Particle filtering works well for tracking in clutter  
 because it can represent arbitrary functional shapes  
 and propagate multiple hypotheses. Less likely model  
 configurations will not be thrown away immediately  
 but given a chance to prove themselves later on, re-  
 sulting in more robust tracking. In its original imple-  
 mentation, Condensation demonstrated robust tracking  
 in low-dimensional configuration spaces (up to about  
 10 DOF) in the presence of significant clutter. Even  
 in the absence of a cluttered background, the compli-  
 cated nature of the observation process during human  
 motion capture causes the posterior density to be non-  
 Gaussian and multi-modal as shown experimentally  
 by Deutscher et al (1999). Condensation has indeed  
 been implemented successfully for short human mo-  
 tion capture sequences (see Deutscher et al. (2000)  
 and Sidenbladh et al. (2000)), however, in the high-  
 dimensional configuration spaces occurring in human  
 motion capture and other domains, there are serious  
 problems with Condensation arising from the inabil-  
 ity of a manageable size particle set (of, say, a few  
 thousand particles), adequately to populate the space  
 and represent an arbitrary density. In fact it has been  
 shown by MacCormick and Isard (2000) that  $N \geq \frac{\mathcal{D}_{\min}}{\alpha^d}$   
 where  $N$  is the number of particles required,  $d$  is the  
 number of dimensions. The survival diagnostic  $\mathcal{D}_{\min}$   
 and the particle survival rate  $\alpha$  are both constants with  
 $\alpha \ll 1$ . Clearly when  $d$  is large normal particle filtering  
 becomes infeasible.

Cham and Rehg (1999) proposed the use of a mul-  
 tiple hypothesis tracker which represented the poste-  
 rior distribution as a piecewise Gaussian. As only local  
 modes are propagated between frames, the solution is  
 computationally much cheaper than Condensation, but  
 they avoid the pitfalls of a single hypothesis tracker.  
 Unlike our work, in which we explicitly model the joint  
 angles and overall pose degrees of freedom, they use a  
 so-called scaled prismatic model which explicitly mod-  
 els 2D in-plane translation and rotation, but models out  
 of plane rotation via a per-link independent scaling.

Partitioned sampling was developed by  
 MacCormick and Blake (1999) as a variation on  
 Condensation to reduce the number of particles  
 needed to track more than one object, and applied  
 by MacCormick and Isard (2000) to the problem

238 of tracking articulated objects. Using partitioned  
 239 sampling reduces the number of particles required  
 240 to  $N \geq \frac{D_{\min}}{\alpha}$  making the problem tractable. How-  
 241 ever, this assumes that the configuration space can  
 242 be sliced so that one can construct an observation  
 243 density  $p(\mathbf{Z}_t | x_i)$  for each dimension  $x_i$  of the model  
 244 configuration vector  $\mathbf{X} = \{x_0 \dots x_d\}$ . This assumption,  
 245 that it is possible independently to localise separate  
 246 parts of an articulated model, is similar to that made  
 247 by Gavrilu and Davis (1996) to enable a hierarchical  
 248 search.

249 Another variation on the standard particle filter used  
 250 to reduce the number of particles needed to effectively  
 251 represent a posterior density has been developed by  
 252 Sullivan et al. (1999). Called *layered sampling* it is  
 253 centred around the concept of importance resampling.  
 254 The technique we present in this paper bears some sim-  
 255 ilarity to layered sampling, but experimental evidence  
 256 suggests our technique is more effective at reducing the  
 257 number of particles required when the dimensionality  
 258 of the search space approaches 30.

259 Two successful recent approaches which use parti-  
 260 cle filtering are due to Sminchisescu and Triggs (2001)  
 261 and Sidenbladh et al. (2000). Both are concerned with  
 262 monocular tracking (in some important ways more dif-  
 263 ficult than the multi-camera case) but in other respects  
 264 the problem is essentially the same: how can a high dimen-  
 265 sional space be adequately populated with a particle set  
 266 of manageable size? Their approaches to this problem  
 267 are quite different and in some ways complementary.

268 The former introduces the idea of *covariance sam-*  
 269 *pling*, spreading particles in areas where there is least  
 270 confidence in the localisation. This idea is very closely  
 271 related to our soft partitioning approach developed in  
 272 Section 6.1. More recently they have extended this  
 273 work explicitly to take into account the particular ambi-  
 274 guities that arise from human kinematics, “scattering”  
 275 particles into areas of potential ambiguity and therefore  
 276 making better use of the particle set Sminchisescu and  
 277 Triggs (2003).

278 The latter work (Sidenbladh et al., 2000, 2002) on  
 279 the other hand, takes the approach that dynamical mod-  
 280 elling can be used to obtain strong, predictive priors,  
 281 reducing the search space to manageable proportions.  
 282 Indeed in Sidenbladh et al. (2000) tracking was res-  
 283 tricted, via the learnt dynamics, to the case of walk-  
 284 ing people. More recently however (Sidenbladh et al.,  
 285 2002) showed how a database of motions could be con-  
 286 structed and efficiently indexed in order to obtain pre-  
 287 dictions over a wide class of motions.

288 In addition to the problems of representing PDFs  
 289 via particle sets in high dimensional spaces, a second  
 290 difficulty is associated with constructing a valid obser-  
 291 vation model  $p(\mathbf{Z}_t | \mathbf{X})$  as a normalised probability den-  
 292 sity distribution. Even if such a likelihood model can be  
 293 constructed the cost of evaluating it can be prohibitive.<sup>1</sup>  
 294 Often an intuitive weighting function  $w(\mathbf{Z}_t, \mathbf{X})$  can be  
 295 constructed that approximates the probabilistic likeli-  
 296 hood  $p(\mathbf{Z}_t | \mathbf{X})$  but which requires much less computa-  
 297 tional effort to evaluate.

298 In this paper we reduce the problem from propagat-  
 299 ing the conditional density  $p(\mathbf{X} | \mathbf{Z}_t)$  using  $p(\mathbf{Z} | \mathbf{X})$ ,  
 300 simply to finding the configuration  $\mathcal{X}_t$  which returns the  
 301 maximum value from a simple and efficient weighting  
 302 function  $w(\mathbf{Z}_t, \mathbf{X})$  at each time  $t$ , given  $\mathcal{X}_{t-1}$ . By doing  
 303 this gains are made on two fronts: (i) it is possible to  
 304 make do with fewer likelihood (or weighting function)  
 305 evaluations because the function  $p(\mathbf{X} | \mathbf{Z}_t)$  no longer  
 306 has to be fully represented; and (ii) an evaluation of a  
 307 simple weighting function  $w(\mathbf{Z}_t, \mathbf{X})$  requires less com-  
 308 putational effort when compared to an evaluation of  
 309 the observation model  $p(\mathbf{Z}_t | \mathbf{X})$ . The main disadvan-  
 310 tage is that we no longer work within a truly Bayesian  
 311 framework.

312 We retain the use of a particle based stochastic frame-  
 313 work because of its ability to handle multi-modal like-  
 314 lihoods, or in the case of a weighting function, one  
 315 with many local maxima. In order most effectively to  
 316 optimise the non-convex weighting function we use an  
 317 approach similar to that of simulated annealing.

## 2.2. Simulated Annealing 318

319 The Markov chain based method of simulated anneal-  
 320 ing was proposed by Kirkpatrick et al. (1983) as a  
 321 means to optimise a multi-modal objective function  
 322  $U(\mathbf{x})$ . It proceeds by defining a distribution over the  
 323 function values

$$P(\mathbf{x}) = \text{const } e^{-\lambda U(\mathbf{x})}$$

324 The aim is then to generate samples  $\mathbf{x}_i$  from this dis-  
 325 tribution, in the knowledge that as  $\lambda \rightarrow \infty$ , the prob-  
 326 ability mass concentrates on the minimum of  $U$ , and  
 327 hence the samples  $\mathbf{x}_i$  will cluster around the minimum  
 328 value state.

329 Samples from the distribution are generated in a  
 330 straightforward fashion using the Metropolis-Hastings  
 331 algorithm (Metropolis et al., 1953) which generates  
 332 a Markov sequence of points whose distribution will

333 converge to  $P$ : a new candidate point  $\mathbf{x}'$  in a sequence is  
 334 generated “at random”, and accepted with probability:

$$\min \left( 1, \frac{P(\mathbf{x}')}{P(\mathbf{x})} \right)$$

335 i.e. the candidate point is accepted if it improves  $U$  or  
 336 with probability  $e^{-\lambda[U(\mathbf{x})-U(\mathbf{x}')]}$ .

337 Simply using a large value of  $\lambda$  and generating a  
 338 sequence starting at a random  $\mathbf{x}_0$  yields poor results if  
 339  $U$  has isolated minima since the sequence can easily  
 340 become trapped in a local mode of  $P$  (e.g. the closest  
 341 to  $\mathbf{x}_0$ ).

342 The annealing process is a heuristic for avoiding this.  
 343 The initial value of  $\lambda$  is set to be small (or in more phys-  
 344 ical language, the temperature, which is inversely pro-  
 345 portional to  $\lambda$ , is initially high). This results in a broad  
 346 distribution  $P$  and hence allows free exploration of the  
 347 search space. Samples are generated from this distribu-  
 348 tion, and then the value of  $\lambda$  increased. Samples are then  
 349 generated from the new distribution starting from the  
 350 final state of the previous sequence, and so on. Each  
 351 increase of  $\lambda$  successively excludes (in a probabilistic  
 352 sense) regions that contain little of the probability mass  
 353 of the distribution.

354 The set of values for  $\lambda = \lambda_M \dots \lambda_0$  is known as  
 355 the annealing schedule. This schedule needs to be de-  
 356 signed as a compromise between speed and efficacy:  
 357 slow annealing is more likely to find a globally optimal  
 358 solution, but is also prohibitively expensive.

359 The similarity with particle-based methods arises  
 360 when we view this process one of generating samples  
 361 from a sequence of distributions,  $P_{\lambda_M} \dots P_{\lambda_0}$ , where  
 362  $P_{\lambda_m}(\mathbf{x}) \propto P_{\lambda_0}(\mathbf{x})^{\beta_m}$ , for  $1 = \beta_0 > \beta_1 > \dots > \beta_M$ ,  
 363 and where  $\beta_m = \lambda_m/\lambda_0$  (as described by Neal (2001)  
 364 whose algorithm ours resembles). Moreover the algo-  
 365 rithm exhibits exactly the kind of behaviour needed  
 366 for the our purposes: one wants to move towards the  
 367 global maximum of the weighting function  $w(\mathbf{Z}_t, \mathbf{X})$ ,  
 368 using the overall trend of the matching function as a  
 369 guide, without becoming misguided by local maxima  
 370 as seen in Fig. 1. The idea of annealing for optimisation  
 371 is now adapted to perform a particle based stochastic  
 372 search within the framework of an annealed particle  
 373 filter.

### 374 3. Annealed Particle Filter

375 A series of weighting functions  $w_0(\mathbf{Z}, \mathbf{X})$  to  $w_M(\mathbf{Z}, \mathbf{X})$   
 376 is employed in which each  $w_m$  differs only slightly

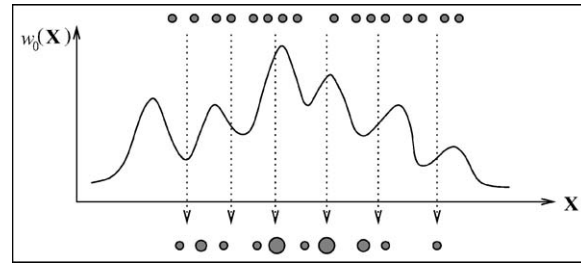


Figure 1. Illustration of the annealed particle filter with  $M = 1$ . Even though a large number of particles are used (so that an equivalent number of weighting function evaluations are made as in Fig. 2), the search is misdirected by local maxima. From the resulting weighted set it is very hard to tell where the global maximum of  $w_0$  lies.

from  $w_{m-1}$  (see Fig. 2, where  $M = 3$ ). The function 377  
 $w_M$  is designed to be very broad, representing the over- 378  
 all trend of the search space while  $w_0$  should be very 379  
 peaked, emphasising local features. This is achieved by 380  
 setting 381

$$w_m(\mathbf{Z}, \mathbf{X}) = w(\mathbf{Z}, \mathbf{X})^{\beta_m}, \quad (3)$$

for  $\beta_0 > \beta_1 > \dots > \beta_M$ , where  $w(\mathbf{Z}, \mathbf{X})$  is the original 382  
 weighting function, as suggested by the discussion in 383  
 Section 2.2. Because it is not the aim to sample from 384  
 $w(\mathbf{Z}, \mathbf{X})$ , but only to find its maximum it is not required 385  
 that  $\beta_0 = 1$ . 386

A large  $\beta_m$  produces a peaked weighting function  $w_m$  387  
 resulting in a high rate of annealing. Small values of  $\beta_m$  388  
 will have the opposite effect. If the rate of annealing is 389  
 too high the influence of local maxima will distort the 390  
 estimate of  $\mathcal{X}_t$  as seen in Fig. 1. If the rate is too low  $\mathcal{X}_t$  391  
 will not be determined with enough resolution (unless 392  
 more layers are used wasting computational resources). 393  
 The manner in which the rate of annealing is influenced 394  
 by the sequence  $\beta_M, \dots, \beta_0$  is discussed in Section 3.1. 395

One annealing run is performed at each time  $t$  using 396  
 image observations  $\mathbf{Z}_t$ . The state of the tracker after 397  
 each layer  $m$  of an annealing run is represented by a 398  
 set of  $N$  weighted particles 399

$$\mathcal{S}_{t,m}^\pi = \{(\mathbf{s}_{t,m}^{(0)}, \pi_{t,m}^{(0)}) \dots (\mathbf{s}_{t,m}^{(N)}, \pi_{t,m}^{(N)})\}. \quad (4)$$

An unweighted set of particles will be denoted 400

$$\mathcal{S}_{t,m} = \{(\mathbf{s}_{t,m}^{(0)}) \dots (\mathbf{s}_{t,m}^{(N)})\}. \quad (5)$$

Each particle in the set  $\mathcal{S}_{t,m}^\pi$  is considered as an 401  
 $(\mathbf{s}_{t,m}^{(i)}, \pi_{t,m}^{(i)})$  pair in which  $\mathbf{s}_{t,m}^{(i)}$  is an instance of the 402

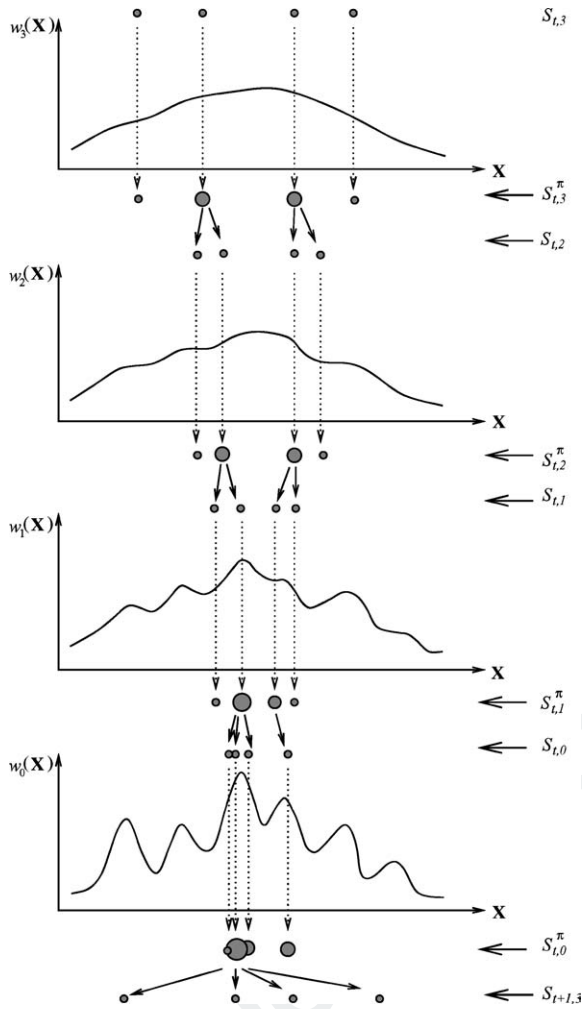


Figure 2. Illustration of the annealed particle filter with  $M = 3$ . With a multi-layered search the sparse particle set is able to migrate gradually towards the global maximum without being distracted by local maxima. The final set  $S_{t,0}^\pi$  provides a good indication of the weighting function's global maximum.

403 multi-variate model configuration  $\mathbf{X}$ , and  $\pi_{t,m}^{(i)}$  is the  
 404 corresponding particle weighting. Each annealing run  
 405 can be broken down as follows (the process is illus-  
 406 trated in Fig. 2).

- 407 1. For every time step  $t$  an annealing run is started at
- 408 layer  $M$ , with  $m = M$ .
- 409 2. Each layer of an annealing run is initialised by a set
- 410 of un-weighted particles  $S_{t,m}$ .
- 411 3. Each of these particles is then assigned a weight

$$\pi_{t,m}^{(i)} \propto w_m(\mathbf{Z}_t, \mathbf{s}_{t,m}^{(i)}) \quad (6)$$

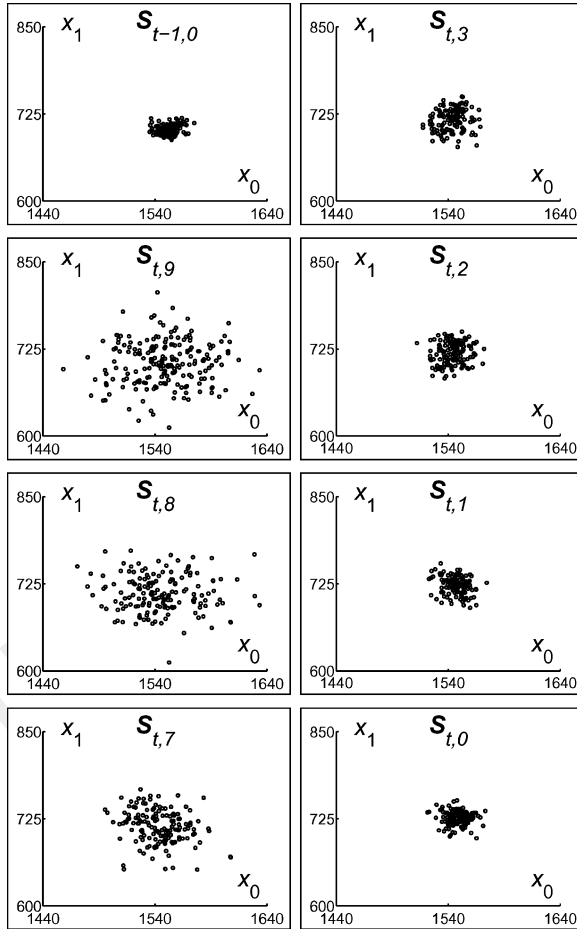


Figure 3. Annealed particle filter in progress. The sets  $S_{t,m}$  are plotted here, taken while tracking the walking person as seen in Fig. 9. Only the horizontal translation components  $x_0$  and  $x_1$  of the model configuration vector  $\mathbf{X}$  are shown. Starting with  $S_{t-1,0}$  from the previous time step the particles are diffused to form  $S_{t,9}$  which easily covers the expected range of translational movement of the subject. The particles are then slowly annealed over 10 layers (the sets  $S_{t,6}$  to  $S_{t,4}$  are omitted for brevity) to produce  $S_{t,0}$  which is clustered around the maximum of the weighting function.

- which are normalised so that  $\sum_N \pi_{t,m}^{(i)} = 1$ . The set of weighted particles  $S_{t,m}^\pi$  has now been formed.
4.  $N$  particles are drawn randomly from  $S_{t,m}^\pi$  with replacement and with a probability equal to their weighting  $\pi_{t,m}^{(i)}$ . As the  $n$ th particle  $\mathbf{s}_{t,m}^{(n)}$  is chosen it is used to produce the particle  $\mathbf{s}_{t,m-1}^{(n)}$  using

$$\mathbf{s}_{t,m-1}^{(n)} = \mathbf{s}_{t,m}^{(n)} + \mathbf{B}_m \quad (7)$$

where  $\mathbf{B}_m$  is a multi-variate Gaussian random variable with covariance  $\mathbf{P}_m$  and mean  $\mathbf{0}$ .

Au: Pls. cite Fig. 3 in text.

420 5. The set  $\mathcal{S}_{t,m-1}$  has now been produced which can be  
 421 used to initialise layer  $m-1$ . The process is repeated  
 422 until we arrive at the set  $\mathcal{S}_{t,0}^\pi$ .  
 423 6.  $\mathcal{S}_{t,0}^\pi$  is used to estimate the optimal model configu-  
 424 ration  $\mathcal{X}_t$  using

$$\mathcal{X}_t = \sum_{i=1}^N \mathbf{s}_{t,0}^{(i)} \pi_{t,0}^{(i)}. \quad (8)$$

425 7. The set  $\mathcal{S}_{t+1,M}$  is then produced from  $\mathcal{S}_{t,0}^\pi$  using

$$\mathbf{s}_{t+1,M}^{(n)} = \mathbf{s}_{t,0}^{(n)} + \mathbf{B}_0. \quad (9)$$

426 This set is then used to initialise layer  $M$  of the next  
 427 annealing run at  $t_{+1}$ .

Note that Step 7, where the particle set for the next time-step is generated, incorporates no dynamic model. There is nothing in the algorithm that precludes the use of dynamics: simply replace Eq. (9) with the more general

$$\mathbf{s}_{t+1,M}^{(n)} = f(\mathbf{s}_{t,0}^{(n)}) + \mathbf{B}_0 \quad (10)$$

428 where the function  $f$  represents the dynamical model.  
 429 We have not done so since our focus is on tracking  
 430 previously unseen/unmodelled agile motions. While a  
 431 dynamical model is certainly beneficial during “steady  
 432 state” tracking, it can be a hindrance if the model is  
 433 poor (as it often is for agile motions). The price we pay  
 434 for this is a less economical use of particles than would  
 435 be ideal, and the potential for jittery trajectories. The  
 436 latter could be addressed by smoothing the recovered  
 437 pose/joint trajectories.

### 438 3.1. Setting the Tracking Parameters

439 It remains to consider how best to set the free paramete-  
 440 rs of the algorithm, and in particular, to consider how  
 441 to influence the annealing schedule. In Eq. (3) it is the  
 442 value of  $\beta_m^k$  that determines the rate of annealing at  
 443 each layer.

444 To see how and why this is so, first note that the  
 445 equivalent of temperature in our particle-based frame-  
 446 work is the particle survival rate: the ratio of effective  
 447 particles to total number of particles. If the probabil-  
 448 ity mass is all concentrated in a few particles then the  
 449 number of effective particles is low, and conversely, an  
 450 even distribution of probability mass amongst particles  
 451 signals a large number of effective particles. A sensible

measure of the effective number of particles that will be  
 chosen for propagation to the next layer is the survival  
 diagnostic  $\mathcal{D}$  (taken from MacCormick (2000)) where

$$\mathcal{D} = \left( \sum_{n=1}^N (\pi^{(n)})^2 \right)^{-1} \quad (11)$$

and from this the particle survival rate  $\alpha$  can be esti-  
 mated MacCormick (2000)

$$\alpha = \frac{\mathcal{D}}{N}. \quad (12)$$

In the case of traditional annealing, the temperature  
 acts like a barrier, restricting the movement of sam-  
 ples: the cooler the temperature, the fewer the number  
 of samples with a low function value  $U(\mathbf{x})$  (energy) that  
 will be generated. In the context of a particle set, a high  
 survival rate corresponds to an even spread probability  
 mass, while a low one suggests the mass is concen-  
 trated in a few particles. Hence decreasing the survival  
 rate has the same effect as cooling the temperature in  
 traditional annealing.

Now  $\mathcal{D}$  is clearly a monotonic decreasing function  
 of  $\beta_m^k$ . At a given layer, we therefore adjust the value  
 of  $\beta_m^k$  to change the value of  $\mathcal{D}(\beta_m^k)$  so that  $\alpha = \mathcal{D}/N$   
 approaches a desired value. This is trivially done by  
 searching over  $\beta_m^k$  (using the value from the previous  
 time step  $\beta_m^{k-1}$  as the starting point) to find the value  
 that solves the equation

$$\alpha_{\text{desired}} = \mathcal{D}(\beta_m^k)/N$$

i.e. produces the desired rate of annealing.

Note that this does not mean the weights have to  
 be completely re-evaluated each time  $\beta_m^k$  is adjusted  
 during the search. Since  $w_m(\mathbf{Z}, \mathbf{X}) = w(\mathbf{Z}, \mathbf{X})^{\beta_m}$  the  
 values  $w(\mathbf{Z}, \mathbf{X} = \mathbf{s}_{t,m}^{(i)})$ ,  $i : 1 \dots N$  can be stored for  
 each set  $\mathcal{S}^{k,m}$  and  $\beta_m^k$  applied to each individual weight  
 as appropriate to produce  $\mathcal{S}_{t,m}^\pi$ .

How then are the appropriate values for  $\alpha_0 \dots \alpha_M$   
 determined? There are also a number of other track-  
 ing parameters that need to be set before tracking can  
 begin, including the number of particles  $N$ , the num-  
 ber of annealing layers  $M$  and the diffusion covariance  
 matrices  $\mathbf{P}_M \dots \mathbf{P}_0$ . A tentative framework has been de-  
 veloped to allocate values to these parameters although  
 it is acknowledged that more work needs to be done in  
 this area.

- 490 1. The first step is to decide on how many anneal-  
 491 ing layers are needed. It was found that doubling  
 492 the number of annealing layers reduces the number  
 493 of particles needed for successful tracking by more  
 494 than half. This will only work up to a point how-  
 495 ever as there seems to be a minimum number ( $N$ )  
 496 of particles needed for tracking no matter how many  
 497 layers are used. Using a 30 DOF model it was found  
 498 that setting  $M = 10$  with  $N \geq 200$  worked well.  
 499 2. Each diagonal element in  $\mathbf{P}_0$  is allocated a value  
 500 equal to half the maximum expected movement of  
 501 the corresponding model configuration parameter  
 502 over one time step. In this way the set  $\mathcal{S}_{t+1,M}$  should  
 503 cover all possible movements of the subject between  
 504 time  $t$  and  $t + 1$ . The amount of diffusion added  
 505 to each successive annealing layer should decrease  
 506 at the same rate as the resolution of the set  $\mathcal{S}_{t,m}$   
 507 increases. Our early experiments used

$$\mathbf{P}_m = \alpha_M \times \dots \times \alpha_{m-1} \times \mathbf{P}_0 \quad (13)$$

- 508 and produced decent results, but a better, adaptive  
 509 method for setting the  $\mathbf{P}$  is described in Section 6.1.  
 510 3. The appropriate rates of annealing  $\alpha_M \dots \alpha_1$  are in-  
 511 fluenced by the number of annealing layers used.  
 512 With a higher number of annealing layers a lower  
 513 rate of annealing can be used to obtain the desired  
 514 resolution. It was found that while using 10 anneal-  
 515 ing layers setting  $\alpha_M = \dots = \alpha_1 = 0.5$  provided  
 516 sufficient resolution of  $\mathcal{X}_t$ .

## 517 4. Implementation

### 518 4.1. The Model

519 The articulated model of the human body used in this  
 520 paper is built around the framework of a kinematic tree,  
 521 as seen in Fig. 4. Each limb is fleshed out using cones  
 522 with elliptical cross-sections. Such a model has a num-  
 523 ber of advantages including computational simplicity,  
 524 high-level interpretation of output and compact repre-  
 525 sentation.

### 526 4.2. The Weighting Function

527 The basic annealed particle filter is a general optimi-  
 528 sation tool and can be used for a variety of purposes  
 529 (for another application see Deutscher et al. (2002))  
 530 with different weighting functions. In the present work

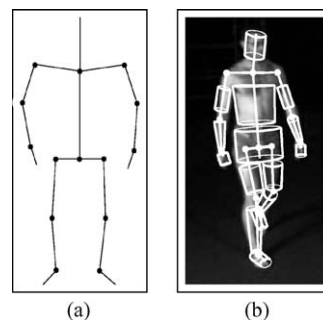


Figure 4. The model is based on a kinematic tree consisting of 17 segments (a). Six degrees of freedom are given to base translation and rotation. The shoulder and hip joints are treated as sockets with 3 degrees of freedom, the clavicle joints are given 2 degrees of freedom (they are not allowed to rotate about their own axis and are assumed to be coupled) and the remaining joints are modelled as hinges requiring only one. This results in a model with 29 degrees of freedom and a configuration vector  $\mathbf{X} = \{x_1 \dots x_{29}\}$ . The model is fleshed out by cones with elliptical cross-sections (b).

we have constructed the weighting function on the 531  
 basis of two image features—edges and foreground 532  
 silhouette—chosen for their joint virtues of simplic- 533  
 ity (i.e. easy and efficient to extract), and a degree 534  
 of invariance to imaging conditions. While these fea- 535  
 tures are not fully general (in particular the silhou- 536  
 ette relies on a knowledge of the background which 537  
 may not be available in more general environments) 538  
 they suffice for our purposes. Even without a large 539  
 degree of background clutter distracting edge mea- 540  
 surements, there remains a challenging, multi-modal 541  
 search problem because of self occlusions and fore- 542  
 ground clutter (i.e. unmodelled markings on the tar- 543  
 get). Other features such as optic flow could equally be 544  
 used. 545

4.2.1. Edges. The strongest continuous edges pro- 546  
 duced by a human subject in an image usually provide 547  
 a good outline of visible arms and legs and are mostly 548  
 invariant to colour, clothing texture, lighting and pose. 549  
 In severely cluttered environments or when the subject 550  
 is wearing very baggy clothes edges may lose some 551  
 of their usefulness, however in most situations they 552  
 provide a good basis for a weighting function. A gradi- 553  
 ent based edge detection mask is used to detect edges. 554  
 The result is thresholded to eliminate spurious edges, 555  
 smoothed with a Gaussian mask and remapped between 556  
 0 and 1. This produces a pixel map (Fig. 5(b)) in which 557  
 each pixel is assigned a value related to its proximity 558  
 to an edge. 559



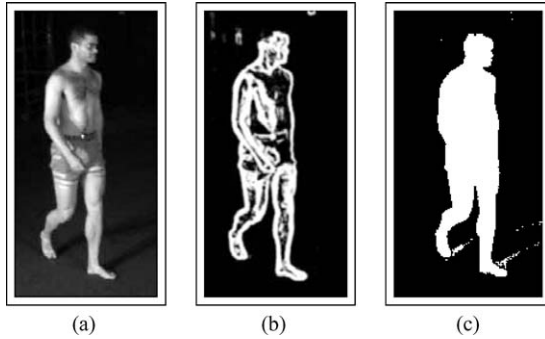


Figure 5. Feature extraction. A gradient based edge detection mask is used to find edges. The result is thresholded to eliminate spurious edges and smoothed using a Gaussian mask to produce a pixel map (b) in which the value of each pixel is related to its proximity to an edge. The foreground is segmented using thresholded background subtraction to produce the pixel map (c) used in the weighting function.

560 A sum-squared difference (SSD) function  $\Sigma^e(\mathbf{X}, \mathbf{Z})$   
 561 is then computed using

$$\Sigma^e(\mathbf{X}, \mathbf{Z}) = \frac{1}{N} \sum_{i=1}^N (1 - p_i^e(\mathbf{X}, \mathbf{Z}))^2 \quad (14)$$

562 where  $\mathbf{X}$  is the model's configuration vector and  $\mathbf{Z}$  is the  
 563 image from which the pixel map is derived.  $p_i(\mathbf{X}, \mathbf{Z})$   
 564 are the values of the edge pixel map at the  $N$  sampling  
 565 points taken along the model's silhouette as seen in  
 566 Fig. 6(a).

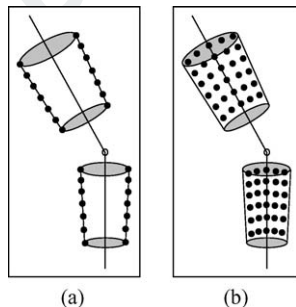


Figure 6. Configurations of the pixel map sampling points  $p_i(\mathbf{X}, \mathbf{Z})$  for the edge based measurements (a) and the foreground segmentation measurements (b). The sampling points for the edge measurements are located along the occluding contours of the model's truncated cones that have been projected into the image. The sampling points for the foreground segmentation measurements are taken from a grid within these occluding contours.

4.2.2. *Silhouette.* The second feature extraction performed on the image is foreground-background segmentation. Thresholded background subtraction was used here to separate the subject from the background and typical results can be seen in Fig. 5(c). This may be inappropriate in some environments with a lot of background movement where more sophisticated methods may have to be employed. Most foreground segmentation techniques are largely invariant to clothing, lighting, pose motion and environment and as such provide an excellent image feature for a general human motion capture system. Once again a pixel map is constructed, this time with foreground pixels set to 1 and background to 0 (Fig. 5(b)), and an SSD is computed

$$\Sigma^r(\mathbf{X}, \mathbf{Z}) = \frac{1}{N} \sum_{i=1}^N (1 - p_i^r(\mathbf{X}, \mathbf{Z}))^2 \quad (15)$$

where  $p_i(\mathbf{X}, \mathbf{Z})$  are the values of the foreground pixel map at the  $N$  sampling points taken from the interior of the truncated cones as seen in Fig. 6(b).

To combine the edge and region measurements the two SSD's are added together and the result exponentiated to give

$$w(\mathbf{X}, \mathbf{Z}) = \exp -(\Sigma^e(\mathbf{X}, \mathbf{Z}) + \Sigma^r(\mathbf{X}, \mathbf{Z})). \quad (16)$$

An equal weighting to each component was determined empirically.

When there is more than one camera the measurements are combined in a similar way, giving

$$w(\mathbf{X}, \mathbf{Z}) = \exp -\left( \sum_{i=1}^C (\Sigma_i^e(\mathbf{X}, \mathbf{Z}) + \Sigma_i^r(\mathbf{X}, \mathbf{Z})) \right) \quad (17)$$

where  $C$  is the number of cameras and  $\Sigma_i^*(\mathbf{X})$  is from camera  $i$ . An example of the output of this weighting function can be seen in Fig. 7.

## 5. Results

Two examples illustrate the system: in the first a subject walks in a circle as seen in Fig. 9; in the second the subject steps over a box, turns 180° on the spot before stepping over it again as seen in Fig. 10.

Three cameras were used to capture the motion and all three views can be seen in the corresponding figures. The same tracking parameters were used in all three

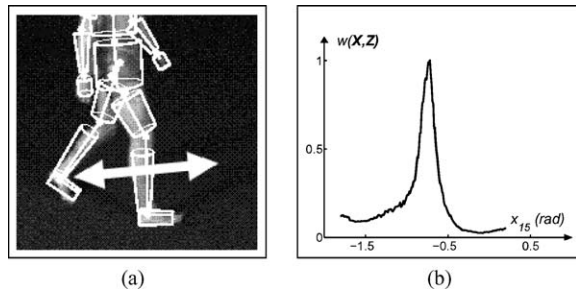


Figure 7. Example output of the weighting function obtained by varying only component  $x_{15}$  of  $\mathbf{X}$  (the right knee angle) using the image and model configuration seen in (a). The function is highly peaked around the correct angle of  $-0.7$  radians (b).

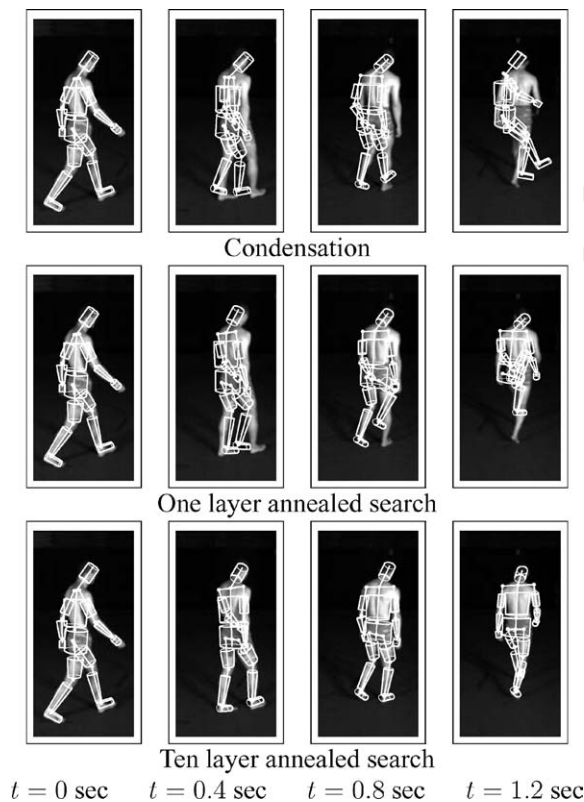


Figure 8. A comparison of condensation with the annealed particle filter. At top the results of tracking with 4000 particles using standard condensation can be seen. Tracking gradually deteriorates until terminal failure after 1.2 seconds. Experiments with 40000 particles were carried out taking over 30 hours to process just 4 seconds of video, still with negative results. An annealed search using 4000 particles with one layer fails little better (middle), also suffering terminal failure after 1.2 seconds. An annealed search using 400 particles and 10 layers (i.e., 4000 weighting function evaluations per frame) tracks very well.

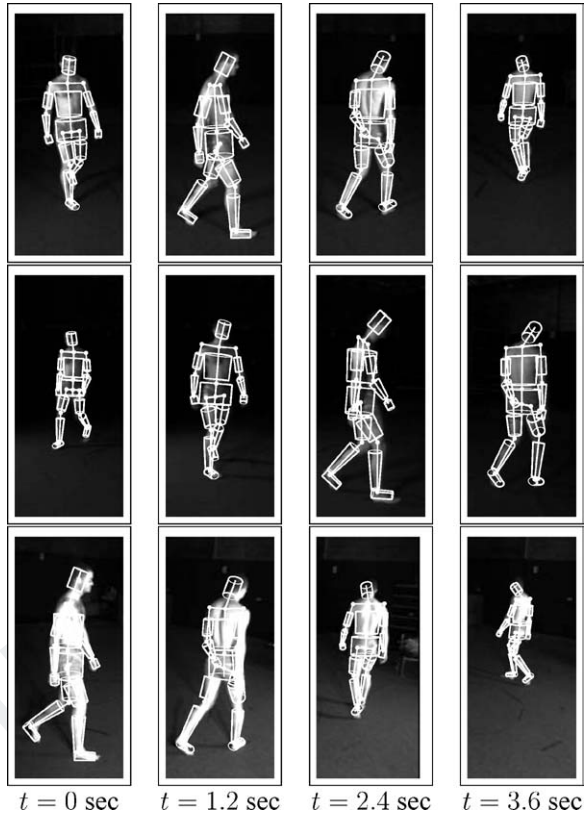


Figure 9. Walking in a circle. Using three cameras (arrayed here from top to bottom) a person is tracked over 4 seconds while walking in a circle. The tracker maintains an accurate lock throughout. 10 annealing layers were used with 200 particles for this sequence.

sequences, which demonstrate the tracker's ability to follow a wide range of human movement.

A comparison of the annealed particle filter with standard Condensation can be seen in Fig. 8. Direct comparison is complicated by the fact that in Annealed Particle Filtering we use a simplified weighting function rather than a "correct" likelihood taking expected clutter into account (such as is derived in Blake and Isard (1998)). For this experiment the likelihood for Condensation comprised the edge based likelihood of Blake and Isard (1998), fused with a silhouette observation. The pose shown in each frame is the sample mean of the particle set. The one layer annealed search represents a similar experiment. It differs from Condensation in using the simplified weighting function (exactly the same as for the full Annealed Particle Filter experiment), and in propagating only the mode of the distribution between frames. The former difference accounts, remarkably, for a four-fold increase in speed

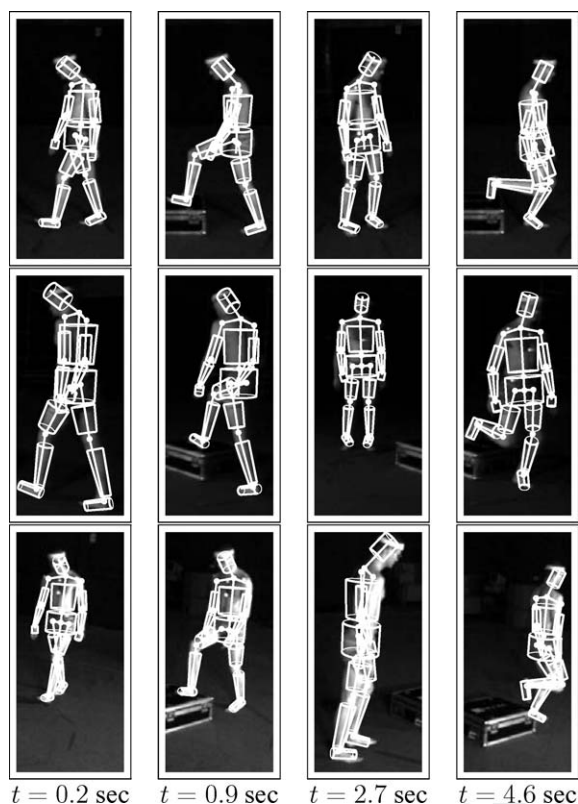


Figure 10. Stepping over a box. Using three cameras (arrayed here from top to bottom) a person is tracked over 5 seconds while stepping over a box, turning around and stepping over the box again. The tracker maintains an accurate lock throughout. 10 annealing layers were used with 200 particles for this sequence.

621 of execution. The final part of the experiment shows  
622 tracking using the full Annealed Particle Filter.

623 Each algorithm used a total of 4000 likelihood evalu-  
624 ations. In the final case this was divided as 400 particles  
625 over 10 layers. It was found in practice that good re-  
626 sults on this sequence could be achieved with as few as  
627 100 particles. While not being a strictly fair comparison  
628 between Condensation and the Annealed Particle Fil-  
629 ter, the experiment gives an indication of the improved  
630 tracking performance of the APF given equivalent  
631 computational resources.

## 632 6. Algorithm Extensions: Hierarchical Search

633 The Annealed particle filter (APF), introduced in the  
634 previous sections, directly addresses the problem of  
635 searching high-dimensional configuration spaces, and  
636 has been demonstrated to be an effective and robust

tracking framework for Human Motion Capture. How- 637  
ever it remains a computationally intensive technique. 638  
The promise of further improvements is held out by the 639  
fact that the model is structured as a kinematic tree. 640

One way to reduce the effective volume of the config- 641  
uration space is to perform a hierarchical search. If one 642  
part of an articulated model can be localised independ- 643  
ently then it can be used as a constraint for restricting 644  
the rest of the search. This straightforward idea has 645  
been applied in a heuristic fashion by (among others) 646  
Gavrila and Davis (1996), who localised the torso us- 647  
ing colour cues and used this information to constrain 648  
the search for the limbs, and more recently by Mikic 649  
et al., who first locate the head in order to limit their 650  
subsequent search. Although this approach is usually 651  
sound, without the assistance of colour cues (or other 652  
labelling cues) it is often very hard independently and 653  
reliably to localise specific body parts in realistic sce- 654  
narios. Furthermore, failure of the first heuristic search 655  
can easily lead to catastrophic, unrecoverable failure. 656

A more formal approach to hierarchical search was 657  
proposed by MacCormick and Isard (2000). That work 658  
applied partitioned sampling to tracking articulated ob- 659  
jects, but crucially assumed that the configuration space 660  
can be sliced so that one can construct an observation 661  
density for each dimension of the model configuration 662  
vector—effectively that it is possible independently to 663  
localise separate parts of an articulated model. 664

When using all but the simplest kinematic models, 665  
the optimal partitioning may not be obvious and it may 666  
indeed change over time as the degree of interaction 667  
between different segments of a model changes—such 668  
as when the legs cross during walking. Rather than 669  
impose a hierarchy on the search, we seek instead 670  
to develop a method for soft or fuzzy partitioning 671  
in which there is no need to commit to a particular 672  
hierarchy. Cham and Rehg (1999) capture this spirit in 673  
describing a search which is sequential in the degrees 674  
of freedom of the body. Their crucial innovation is 675  
to allow the order to be flexible: the search for body 676  
parts is conducted on a “best”-first basis, where best 677  
is defined as the component which can be found with 678  
minimum effort, usually meaning minimum variance. 679

While motivated by similar desires, our solution is 680  
rather different from theirs. Our approach improves 681  
upon and extend the APF in two ways. First we in- 682  
troduce a means to make the diffusion step in the APF 683  
adaptive, so that effort is not wasted in those places 684  
where the algorithm is already confident of doing well, 685  
and is concentrated on localising parts whose location 686

687 is uncertain. The effect of this can be interpreted as a  
 688 hierarchical search strategy which *automatically* par-  
 689 titions the search space in a soft way, without any ex-  
 690 plicit representation of the partitions (Section 6.1). Sec-  
 691 ond, we introduce a crossover operator (similar to that  
 692 found in Genetic Algorithms) which improves the abil-  
 693 ity of the tracker to search different partitions in parallel  
 694 (Section 6.2).

695 We present results for simple examples to demon-  
 696 strate the new algorithm's implementation and effec-  
 697 tiveness, and show that these measures together  
 698 increase the tracker efficiency by a factor of 4 and in-  
 699 crease agility of the motion that can be tracked.

700 We apply the tracker to the complex problem of Hu-  
 701 man Motion Capture with 34 degrees of freedom. Extra  
 702 degrees of freedom have been added to the model in  
 703 Fig. 4 in the back (2) to allow arching that would not  
 704 normally be encountered in everyday walking (and was  
 705 not necessary in our earlier experiments), in the neck (1)  
 706 to account for head nodding, and the clavicles are given  
 707 independent motion (2 each).

708 6.1. Adaptive Diffusion  
 709 and Hierarchical Partitioning

710 Consider the simple task of tracking an articulated arm  
 711 as seen in Fig. 11. The arm consists of four segments,  
 712 each joined by a swivelling joint with one end rooted

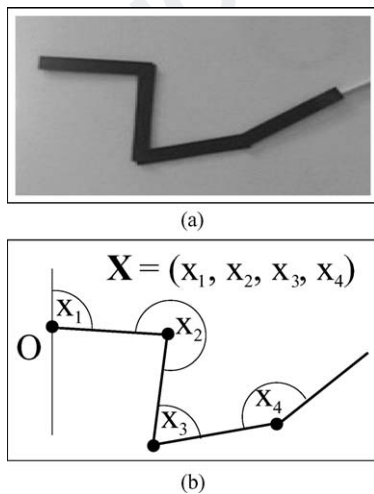


Figure 11. A planar articulated arm with 4 DOF is shown (a). It consists of four links connected by swivelling joints and rooted at  $O$ . The configuration of the arm is described by  $\mathbf{x} = (x_1, x_2, x_3, x_4)$  as seen in (b).

on the spot. A configuration of the arm is described by  
 an instance of the state variable  $\mathbf{x} = (x_1, x_2, x_3, x_4)$ .  
 The weighting function  $w(\mathbf{z}, \mathbf{x})$  required for the APF  
 is computed by a Sum of Squared Differences (SSD)  
 measure between a model template and a silhouette  
 image (the detail to the regional correlation portion of  
 the observation model in Eq. (15)).

The set  $\mathcal{S}_{t,m}$  is initialised with particles uniformly  
 distributed over a range of  $\mathbf{x}$  that we know to con-  
 tain the actual position of the arm. This results in a  
 large and similar variance for each parameter of  $\mathbf{x}$  over  
 all the particles in  $\mathcal{S}_{t,m}$  as can be seen in Fig. 12(a).

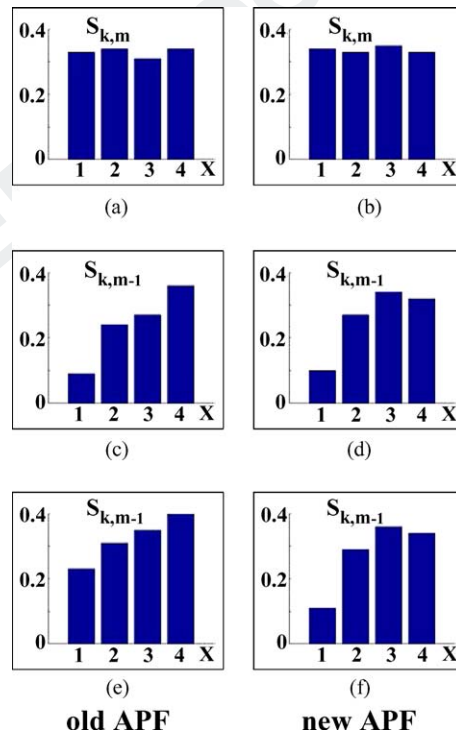


Figure 12. Parameter variance over one annealing layer: new APF vs. old APF. On the left graphs a, b and c plot the variance of each parameter of  $\mathbf{x} = (x_1, x_2, x_3, x_4)$  through the first annealing run of the APF when tracking the articulated arm seen in Fig. 11. Graphs d, e and f show the same information for the improved APF as described in Section 6.1. Graphs a and d show the variances of the initial set  $\mathcal{S}_{t,m}$ , displaying equal variances for each parameter. Graphs b and e show the variances of the set  $\mathcal{S}_{t,m-1}$  before the addition of diffusion noise. Note that in both b and e,  $x_1$  has a very small variance indicating advanced localisation, however the variance of  $x_2, x_3$  and  $x_4$  has been reduced only a little. Up to this point the algorithms are the same and any differences between b and e are random. After the addition of noise in the original APF the localisation of  $x_1$  has been greatly degraded as seen in graph c, however when noise is added in proportion to each parameter's variance the localisation of  $x_1$  is preserved as seen in graph f.

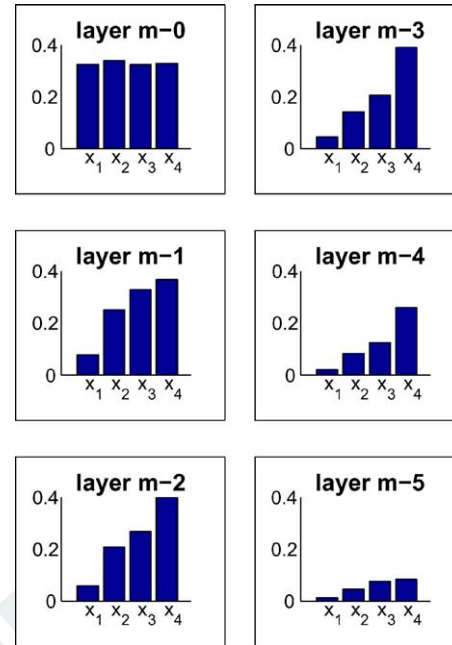
725 After calculating a weight  $\pi_{t,m}^{(i)}$  for each particle using  
 726  $w_m(\mathbf{z}_t, \mathbf{s}_{t,m}^{(i)})$  we then proceed to Step 4 of the APF  
 727 and draw  $N$  particles from  $\mathcal{S}_{t,m}^\pi$  with replacement and  
 728 probability proportional to each particle's weight.

729 Consider the set  $\mathcal{S}_{t,m}$  so produced before the addition  
 730 of any noise. In a typical annealing run the individual  
 731 parameters of each particle were found to have variance  
 732 as detailed in Fig. 12(b). Note here that the variance of  
 733  $x_1$  has been greatly reduced while the other parameters  
 734  $x_2, x_3$  and  $x_4$  have been hardly reduced at all. The  
 735 variance of any parameter can be considered (with a  
 736 number of acceptable caveats) to be directly related to  
 737 the degree to which the optimal value for that parameter  
 738 has been localised. Figure 12(b) shows that  $x_1$  has  
 739 been localised down to a very small area of its range  
 740 simply because it dominates the topology of the search  
 741 space whereas each particle's values for  $x_2, x_3$  and  $x_4$   
 742 had very little influence on whether it was selected or  
 743 not. In effect we see here an automatic partitioning of  
 744 the state space into soft partitions according each parameter's  
 745 topological dominance.

746 The weakness of the original APF (indeed any particle  
 747 filter) arises with the addition of diffusion noise to each  
 748 particle upon selection. According to Eqs. (7)  
 749 and (13) an equal amount of noise should be added to  
 750 each parameter. This results in a parameter variance  
 751 profile like that seen in Fig. 12(c) with the localisation  
 752 of  $x_1$  seen in Fig. 12(b) all but wiped out by the  
 753 excessive addition of noise.

754 If instead the amount of randomness added to the  
 755 parameters of each selected particle is proportional to  
 756 the variance of that parameter over the entire set of  
 757 particles, these gains will be protected from disruption.  
 758 Instead we will arrive at the situation seen in Fig. 12(f)  
 759 where enough noise has been added to each parameter  
 760 to allow the thorough diffusion of the particles into the  
 761 spaces between repeatedly selected particles, but not  
 762 enough to increase the variance of any given parameter  
 763 which would erase any localisation gains made up to  
 764 that point.

765 If this new method for determining the elements of  
 766  $\mathbf{P}_t$  (the covariance matrix of  $\mathbf{B}$  from Eq. (7)), is continued  
 767 through all the annealing layers we can see that  
 768 each parameter is localised in turn, with some degree  
 769 of overlap as seen in Fig. 13. This can be compared  
 770 to the pattern of variance reduction for the original  
 771 APF algorithm seen in Fig. 14. This is exactly the  
 772 kind of hierarchical soft partitioning that was desired  
 773 and no explicit partition boundaries or functions were  
 774 required.



775 *Figure 13.* Variance reduction with the improved APF. Here we  
 776 see the orderly reduction of each of the four parameter's variances  
 777 from most dominant ( $x_1$ ) to least dominant ( $x_4$ ) over 6 layers of the  
 778 annealing process while tracking the simple articulated arm. Using  
 779 the improved APF results in a 2-fold increase in efficiency over the  
 780 classical APF. Tracker efficiency was measured by the minimum  
 781 number of particles needed to successfully track the articulated arm  
 782 over 40 frames.

783 Sminchisescu and Triggs (2001) independently arrived at a very similar  
 784 idea, although in that work they were concerned with most effective use  
 785 of particles between frames in order to recover from "ambiguous" poses.

786 The changes to the APF are almost trivial, and can be formalised as follows.  
 787 Step 4 of the APF algorithm described in Section 3 is amended so that at  
 788 layer  $m$ ,  $\mathbf{P}_m$  is set to be proportional to the covariance of the particles  
 789 in  $\mathcal{S}_{t,m}$  as it exists before the addition of noise, i.e.,

$$\mathbf{P}_m \propto \frac{1}{N} \sum_{i=1}^N (s_{t,m}^{(i)} - s_{t,m}^{av}) \cdot (s_{t,m}^{(i)} - s_{t,m}^{av})^T. \quad (18)$$

786 where  $s_{t,m}^{av}$  is the sample mean of the particle set.

787 Using this modification enabled successful tracking with the APF with fewer  
 788 than half the number of particles; i.e. a 2-fold increase in efficiency.

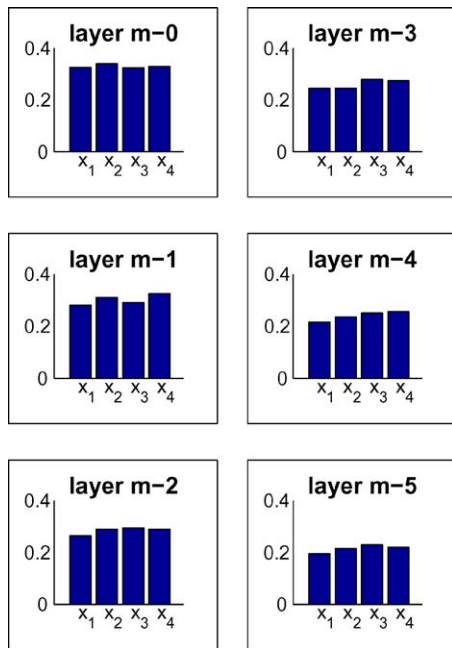


Figure 14. Variance reduction with the conventional APF. The even reduction in variance over 6 layers of the annealing process is shown in contrast to Fig. 13. There is little evidence of hierarchical partitioning and more annealing layers will be required to find the optimal configuration.

790 6.2. A Crossover Operator and Parallel Partitions

791 Now consider the articulated object found in Fig. 15  
 792 which consists of two articulated arms joined at a sta-  
 793 tionary hinge. This configuration is a much simplified  
 794 version of that found in Human Motion Capture when  
 795 using a model with arms and legs.

796 The soft hierarchical partitioning described in  
 797 Section 6.1 provides some increase in efficiency over  
 798 conventional APF when applied to tracking this assem-  
 799 bly, localising  $x_1$  and  $x_4$  together, then  $x_2$  and  $x_5$   
 800 and finally  $x_3$  and  $x_6$ . However if we were to decouple the  
 801 search space and localise each arm independently the  
 802 computational effort required for tracking would be re-  
 803 duced considerably.

804 One possibility, of course, would be to introduce a  
 805 hard partition between the two arms and conduct two  
 806 separate searches. However, in keeping with our phi-  
 807 losophy of adaptive partitioning, we seek to avoid com-  
 808 mitment to specific partitions.

809 Many people comment on the similarity between  
 810 particle filters and Genetic Algorithms. Both employ  
 811 a set (population) of particles (individuals) coded by  
 812 a state vector (genetic sequence) from which the best  
 813 particles (individuals) are chosen to be propagated to  
 814 the next time-step (generation) in the hope of finding  
 815 the maximum of some function (fittest possible indi-  
 816 vidual).

817 One glaring difference between GA's and a typical  
 818 particle filter is the lack of a crossover operator in the  
 819 particle filter which in a conventional GA is meant to  
 820 simulate the breeding of individuals and the sharing of  
 821 genetic information. The use of the crossover operator  
 822 encourages the survival of short, highly fit sections of  
 823 the parameter space known in some GA literature as  
 824 building blocks. This is done in the hope that when  
 825 highly fit building blocks are brought together they  
 826 will have a good chance of forming a very fit com-  
 827 plete individual. These building blocks are effectively  
 828 optimised in parallel without any specification of their  
 829 boundaries or appropriate building block (partition)

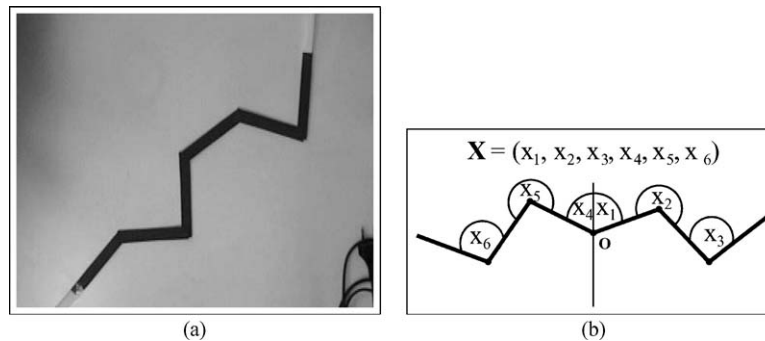


Figure 15. A pair of planar articulated arms consisting of 3 segments each and each rooted to point  $O$  (as seen in b) are used to demonstrate the effectiveness of the crossover operator. The configuration of the arms is described by  $\mathbf{x} = (x_1, \dots, x_6)$  as seen in (b).

830 weighting functions, exactly the kind of behaviour we  
 831 are looking for.

832 We now describe how to incorporate the crossover  
 833 operator into the framework of the APF and examine  
 834 the effect via a simple example.

835 **6.2.1. Inclusion of the Crossover Operator in the APF.**

836 The inclusion of the crossover operator can be for-  
 837 malised as follows. In Step 4 of the APF (as described  
 838 in Section 3) at annealing layer  $m$ , the  $i$ th particle  
 839 of  $\mathcal{S}_{t,m-1}$  is created by drawing two particles from  
 840  $\mathcal{S}_{t,m}^\pi$  with probability proportional to their respective  
 841 weights. Two parameter indices  $\gamma$  and  $\epsilon$  are chosen ran-  
 842 domly and the two selected particles  $\mathbf{s}_{t,m}^{(a)} = (x_1^a \dots x_L^a)$   
 843 and  $\mathbf{s}_{t,m}^{(b)} = (x_1^b \dots x_L^b)$  are combined to form the new  
 844 particle  $\mathbf{s}_{t,m-1}^{(i)}$  where

$$\mathbf{s}_{t,m-1}^{(i)} = (x_1^a, \dots, x_\gamma^a, x_{\gamma+1}^b, \dots, x_\epsilon^b, x_{\epsilon+1}^a, \dots, x_L^a). \quad (19)$$

845 Noise is then added to each particle as detailed in  
 846 Section 6.1.

847 **6.2.2. Testing the Crossover Operator.**

To assess the benefit to the crossover operator two articulated objects were tracked: the first (Fig. 11), was used in the experiment from Section 6.1, an un-branched articulated arm; the second as seen in Fig. 15 is two articulated arms rooted to the same position.

853 As seen in Fig. 16, the object consisting of branched  
 854 arms was more effectively localised by the APF that  
 855 employed the crossover operator whereas there was no  
 856 difference when it was applied to the non-branched ob-  
 857 ject. A good graphical illustration of what the crossover  
 858 operator is actually doing—i.e. partitioning sections of  
 859 the search space which can be tracked in parallel—is  
 860 evident in Figs. 17 and 18 where the parameters lo-  
 861 calised best first are those closest to the root of the  
 862 tree.

863 A good indication of the increased speed provided  
 864 by the crossover operator when tracking branched  
 865 objects is again the number of particles needed for  
 866 successful tracking. This number was reduced by a  
 867 factor of 2 with the introduction of the crossover  
 868 operator.

869 **6.3. Results for Full-Body Tracking**

870 Although less clear-cut than the results for the “toy”  
 871 example in the previous section, Figs. 19 and 20

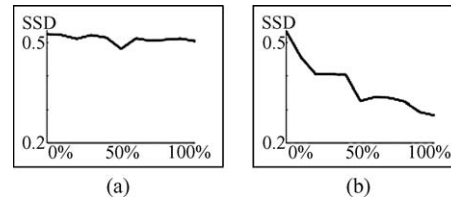


Figure 16. The crossover operator in action. The Sum of Squared Differences (SSD) match between model and image obtained after a set number of annealing layers is plotted against the percentage of particles generated using the crossover operator at each annealing layer. Graph (a) shows the result for the articulated arm seen in Fig. 11 where no benefit to using the crossover operator is seen either (i.e. the SSD does not increase). Graph (b) shows the result for the articulated arms seen in Fig. 15 where a steady improvement in tracking performance is seen when increasing the percentage of particles produced using the crossover operator. This shows that the crossover operator is able to decouple sections of the search space effectively and enables the APF to search them in parallel, improving tracker performance.

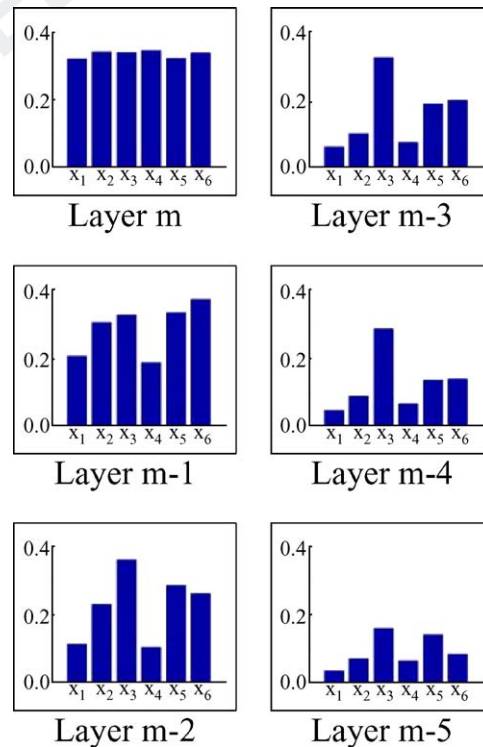


Figure 17. Variance reduction for the parallel arms. When the APF with crossover operator is applied to the articulated arms seen in Fig. 15 we get the pattern of variance reduction seen above. The graphs show the parameters describing each arm ( $x_1, x_2, x_3$  and  $x_4, x_5, x_6$ ) being localised in order of decreasing topological dominance, from the fixed point of the articulated arms, progressing outward.

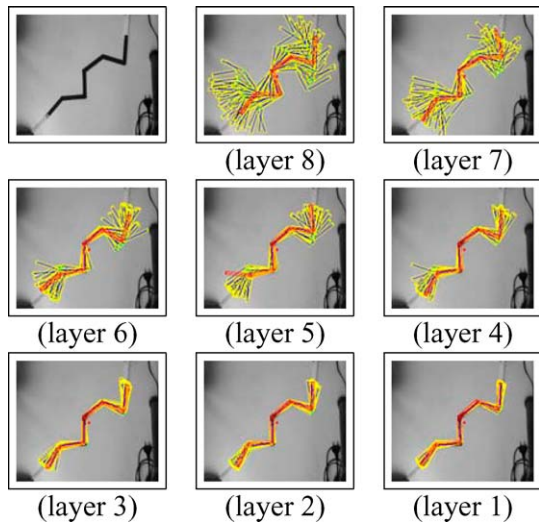


Figure 18. Particle distribution for the branched articulated arm over 8 annealing layers. The entire set of particles is drawn at each annealing layer. The hierarchical localisation of each model segment from the hinge joint outwards is clearly seen.

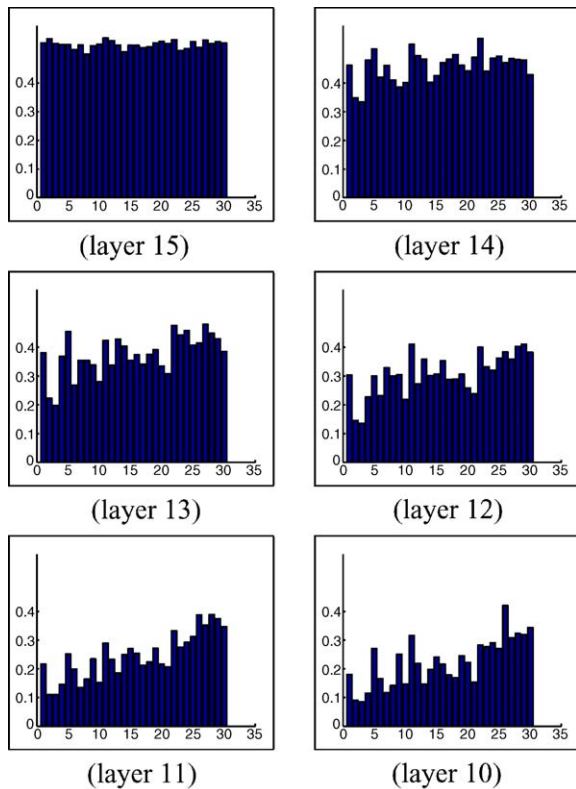


Figure 19. Annealed particle filter particle variance for a fully-body model. The difference in rates of variance reduction for each parameter can clearly be seen. As expected a more complicated pattern of reduction than that seen for the simple articulated arm is evident.

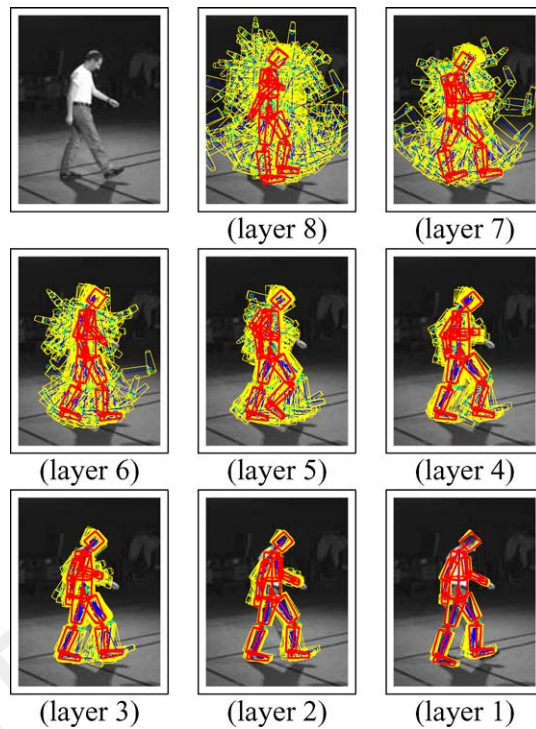


Figure 20. Particle distribution for a full-body model. The entire set of particles is drawn at each annealing layer for one frame. The hierarchical reduction of each parameter from torso rotation and translation out to the limb joint angles is evident.

show a similar process of variance reduction when 872  
 the PAPF with crossover is applied to full-body 873  
 tracking. 874

The algorithm was applied to a variety of challeng- 875  
 ing sequences of human movement including walk- 876  
 ing with turns (Fig. 21), running around in a random 877  
 fashion (Fig. 22) and handstands (Fig. 23). The se- 878  
 quences for these experiments were generated using 879  
 three evenly spaced cameras, calibrated and hardware 880  
 synchronised. 881

We define successful tracking qualitatively as occur- 882  
 ring when the algorithm locks onto the body 883  
 and limbs for the duration of the sequence, return- 884  
 ing sensible values (i.e. ones that can be used for 885  
 re-animation, for example) for the pose and articula- 886  
 tion parameters. Our tests measured the number 887  
 of particles needed to achieve such successful 888  
 tracking. This number represents a sensible mea- 889  
 sure of algorithm speed since the number of like- 890  
 likelihood evaluations dominates the processing time. 891



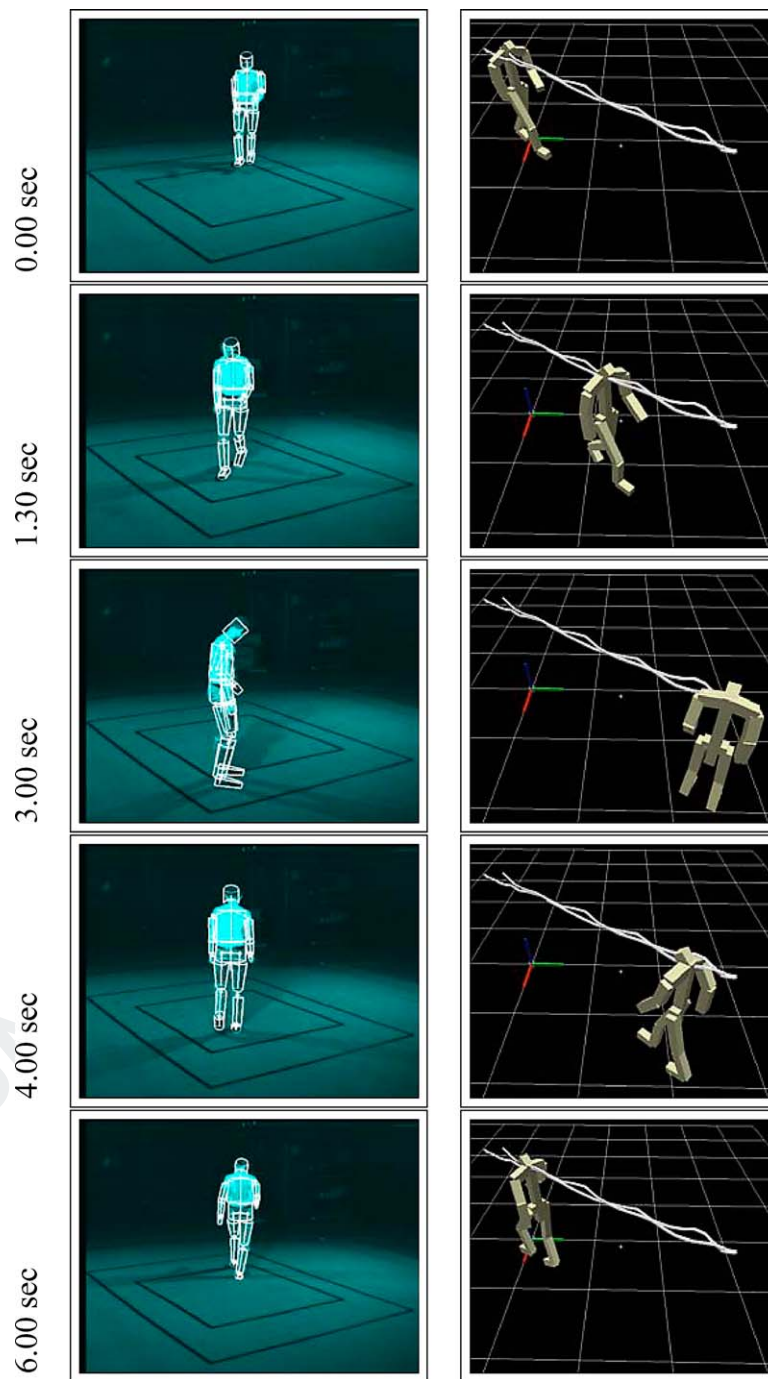


Figure 21. Tracking a walking person.

892 We observed an improved by a factor of 4 when  
 893 comparing the new PAPP to the original APF  
 894 (i.e. successful tracking achieved with one quar-  
 895 ter the number of particles). As a result the PAPP

required on average 15 seconds to process one 896  
 frame whereas the APF required around 60 seconds 897  
 when run on a single processor 1 GHz pIII Linux 898  
 box. 899

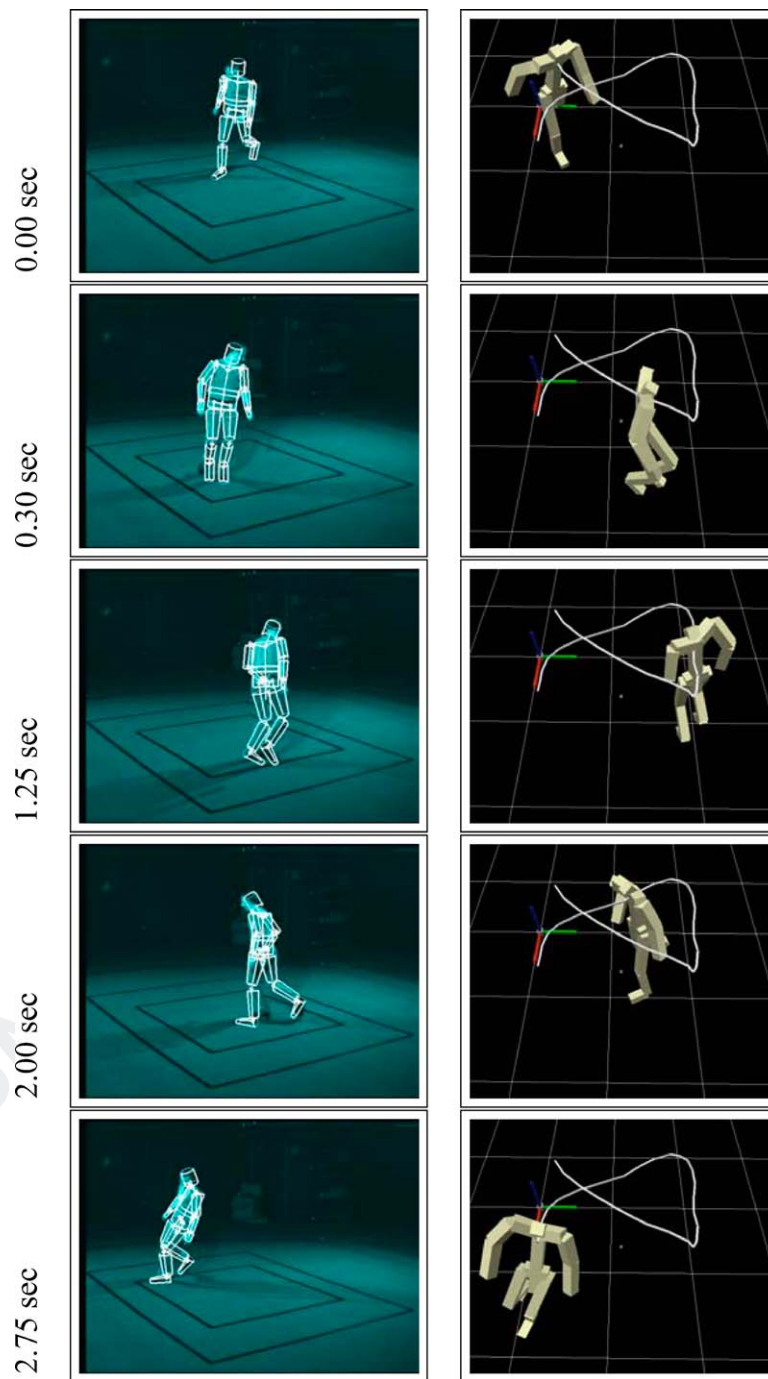


Figure 22. Tracking a running person.

900 We have also built a parallel implementation,  
 901 in which particles are farmed out to independ-  
 902 ent processors which compute the weight/likelihood  
 903 function. This achieves the sort of speed-ups

that are to be expected, with processing time 904  
 decreasing linearly in the number of proces- 905  
 sors (with a constant of proportionality around 906  
 0.8). 907

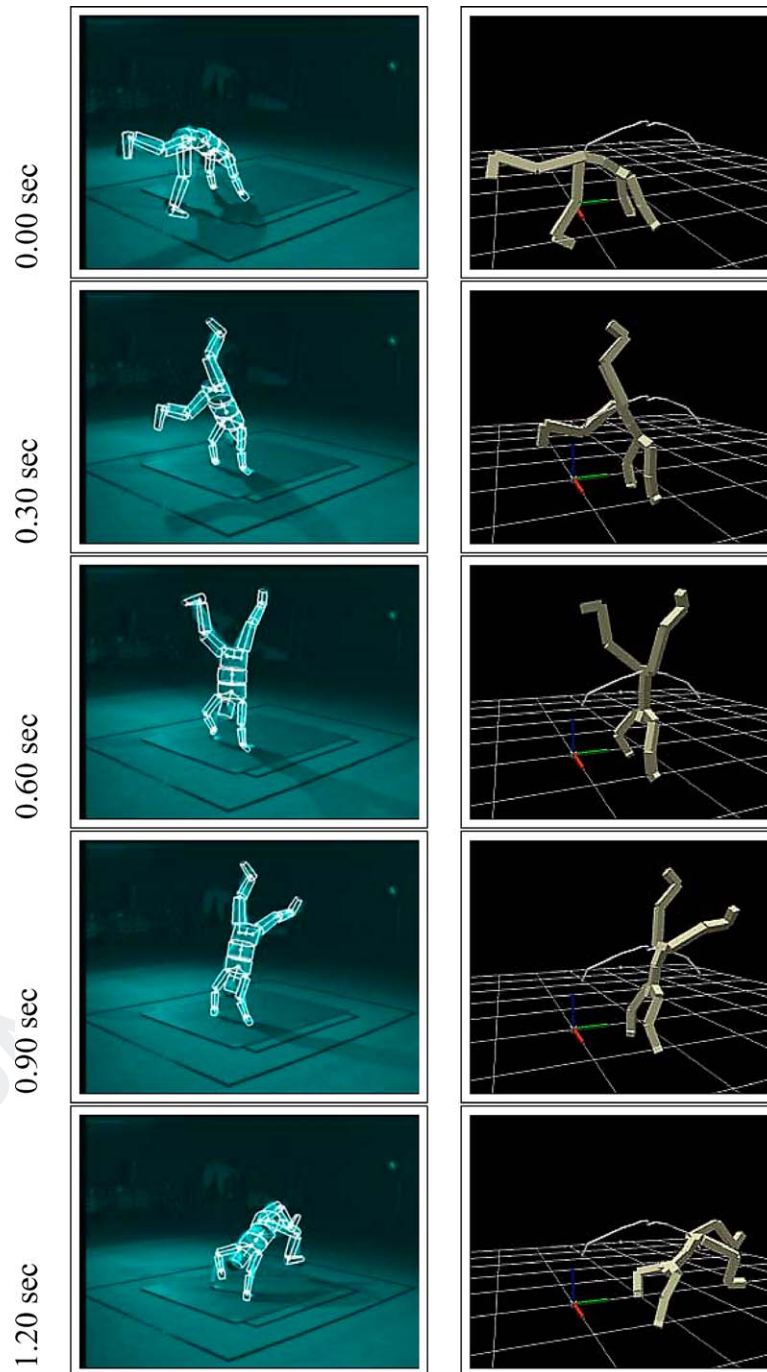


Figure 23. Tracking a person performing a hand-stand.

908 **7. Discussion and Conclusion**

909 We have developed a general algorithm for searching  
 910 large configuration spaces which is more efficient than

traditional particle filters, but which retains a number 911  
 of their significant advantages. The algorithm has been 912  
 applied to the problem of visually tracking a person in 913  
 multiple cameras. In this context we have demonstrated 914

915 reliable tracking of complex human motion using simple  
916 image features, and without the need for a strong  
917 dynamic model of the motion.  
918 We have also introduced two novel improvements  
919 to the algorithm, soft hierarchical partitioning, and a  
920 crossover operator, which have the combined effect of  
921 improving performance and increasing efficiency.  
922 The results, especially Figs. 21–23, show a robust-  
923 ness of tracking human motion achieved by very few  
924 other algorithms. Of particular note in the sequences  
925 shown are the points where the subject turns rapidly  
926 on the spot (shown in both Figs. 21 and 22), and the  
927 unusual and rapid motion of a handstand.  
928 Our primary effort has been concentrated on the  
929 search technique. It seems clear that improvements in  
930 the modelling process, such as published in Plänklers  
931 and Fua (2003), and in dynamic modelling Sidenbladh  
932 et al. (2000), would improve tracking reliability and  
933 applicability further.  
934 Though in the experiments shown the background  
935 lacks a large degree of clutter (but is not entirely  
936 clean either), tracking agile motions, even with mul-  
937 tiple cameras, remains a difficult problem. We have  
938 performed experiments with other sequences with a  
939 greater amount of clutter with similar results, but the  
940 exact degree of clutter that can be tolerated is an open  
941 question. No doubt the use of background subtraction  
942 to obtain silhouette information assists in this signifi-  
943 cantly. The algorithm exhibits some robustness to er-  
944 rors in this data, but in cases where poor contrast results  
945 in poor silhouettes and a lack of edges we have observed  
946 tracker failure.  
947 Our results to date have made use of 3 cameras, and  
948 tracking using a single camera raises issues with re-  
949 gard to ambiguity. Experiments with using the APF  
950 monocularly (Lyons, 2002) suggest that in the monoc-  
951 ular case further sophistication in the placement of par-  
952 ticles is required to overcome the inherent ambiguities  
953 and avoid all associated local minima. Some progress  
954 in this respect has been made recently by Sminchisescu  
955 and Triggs (2002, 2003).

## 956 Acknowledgments

957 This work was supported by Vicon Systems Ltd. and  
958 EPSRC grant GR/M15262. We would also like to thank  
959 Andrew Blake, Ben North, Andrew Davison and David  
960 Murray, for many useful discussions and the anony-  
961 mous referees for insightful comments.

## Note

1. Note, for example, that although (Blake and Isard, 1998) derives the full multi-modal likelihood model for edge-normal observations in the presence of clutter, the implementation makes a much simplified assumption of a unimodal likelihood for each individual observation.

## References

- Blake, A. and Isard, M. 1998. *Active Contours*. Springer. 968  
Cham, T.-J. and Rehg, J.M. 1999. Dynamic feature ordering for efficient registration. In *Proc. 7th Int'l Conf. on Computer Vision*, Corfu, vol. 2, pp. 1084–1091. 969  
Cham, T.-J. and Rehg, J. M. 1999. A multiple hypothesis approach to figure tracking. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 239–245. 970  
Deutscher, J., Blake, A., North, B., and Bascle, B. 1999. Tracking through singularities and discontinuities by random sampling. In *Proc. 7th Int. Conf. on Computer Vision*, vol. 2, pp. 1144–1149. 971  
Deutscher, J., Blake, A., and Reid, I.D. 2000. Articulated body motion capture by annealed particle filtering. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, vol. 2, pp. 126–133. 972  
Deutscher, J., Davison, A.J., and Reid, I.D. 2001. Automatic partitioning of high dimensional search spaces associated with articulated body motion capture. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, vol. 2, pp. 669–676. 973  
Deutscher, J., Isard, M., and MacCormick, J. 2002. Automatic camera calibration from a single manhattan image. In *Proc. 7th European Conf. on Computer Vision*, Copenhagen, vol. 4, pp. 175–188. 974  
Drummond, T. and Cipolla, R. 2001. Real-time tracking of highly articulated structures in the presence of noisy measurements. In *Proc. 8th Int'l Conf. on Computer Vision*, Vancouver, pp. 315–320. 975  
Gavrila, D. and Davis, L.S. 1996. 3d model-based tracking of humans in action: A multi-view approach. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 73–80. 976  
Harris, C.G. 1992. Tracking with rigid models. In *Active Vision*. A. Blake and A. Yuille (Eds.), MIT Press: Cambridge, MA. 977  
Isard, M.A. and Blake, A. 1996. Visual tracking by stochastic propagation of conditional density. In *Proc. 4th European Conf. on Computer Vision*, Cambridge, England, p. 343–356. 978  
Kirkpatrick, S., Gellatt, C.D., and Vecchi, M.P. 1983. Optimisation by simulated annealing. *Science*, 220(4598):671–680. 979  
Lyons, D. 2002. A qualitative approach to computer sign language recognition. Master's thesis, University of Oxford. 980  
MacCormick, J. 2000. Probabilistic models and stochastic algorithms for visual tracking. PhD thesis, University of Oxford. 981  
MacCormick, J. and Blake, A. 1999. A probabilistic exclusion principle for tracking multiple objects. In *Proc. 7th Int. Conf. on Computer Vision*, vol. 1, pp. 572–578. 982  
MacCormick, J. and Isard, M. 2000. Partitioned sampling, articulated objects and interface-quality hand tracking. In *Proc. 6th European Conf. on Computer Vision*, Dublin, vol. 2, pp. 3–19. 983  
Metropolis, N., Rosenbluth, A.W., Rosenbluth, M.N., Teller, A.H., and Teller, E. 1953. Equations of state calculations by fast computing machine. *J. Chem. Phys.*, 21:1087–1091. 984  
Mikic, I., Trivedi, M., Hunter, E., and Cosman, P. 2001. Articulated body posture estimation from multi-camera voxel data. In *Proc.* 1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016

## Articulated Body Motion Capture by Stochastic Search 205

- 1017 of the *IEEE Conf. on Computer Vision and Pattern Recognition*,  
 1018 vol. 1, pp. 455–462.
- 1019 Neal, R.M. 2001. Annealed importance sampling. *Statistics and  
 1020 Computing*, (11):125–139.
- 1021 Plänkers, R. and Fua, P. 2003. Articulated soft objects for multi-view  
 provide 1022 shape and motion capture. *IEEE Transactions on Pattern Analysis  
 page 1023 and Machine Intelligence*, 25(10).
- range 1024 Sidenbladh, H., Black, M.J., and Fleet, D.J. 2000. Stochastic track-  
 1025 ing of 3D human figures using 2D image motion. In *Proc. 6th  
 1026 European Conf. on Computer Vision*, Dublin, vol. 2, pp. 702–  
 1027 718.
- 1028 Sidenbladh, H., Black, M.J., and Sigal, L. 2002. Implicit probabilistic  
 1029 models of human motion for synthesis and tracking. In *Proc. 7th  
 1030 European Conf. on Computer Vision*, Copenhagen, vol. 1, pp. 784–  
 1031 800.
- Sminchisescu, C. and Triggs, B. 2001. Covariance scaled sampling 1032  
 for monocular 3d body tracking. In *Proc. of the IEEE Conf. on 1033  
 Computer Vision and Pattern Recognition*, vol. 1, pp. 447–454. 1034
- Sminchisescu, C. and Triggs, B. 2002. Hyperdynamics importance 1035  
 sampling. In *Proc. 7th European Conf. on Computer Vision*, 1036  
 Copenhagen, vol. 1, pp. 769–783. 1037
- Sminchisescu, C. and Triggs, B. 2003. Kinematic jump processes 1038  
 for monocular 3d human tracking. In *Proc. of the IEEE Conf. on 1039  
 Computer Vision and Pattern Recognition*, vol. 1, pp. 69–76. 1040
- Sullivan, J., Blake, A., Isard, M., and MacCormick, J. 1999. Object 1041  
 localization by bayesian correlation. In *Proc. 7th Int. Conf. on 1042  
 Computer Vision*, vol. 2, pp. 1068–1075. 1043
- Wachter, S. and Nagel, H. 1999. Tracking persons in monocular 1044  
 image sequences. *Computer Vision and Image Understanding*, 1045  
 74(3):174–192. 1046

UNCORRECTED PROOF