

Estimating Articulated Human Motion With Covariance Scaled Sampling

Cristian Sminchisescu

Bill Triggs

INRIA Rhône-Alpes, GRAVIR-CNRS, 655 avenue de l'Europe, 38330 Montbonnot, France

{Cristian.Sminchisescu,Bill.Triggs}@inrialpes.fr; www.inrialpes.fr/movi/people/{Sminchisescu,Triggs}

Abstract

We present a method for recovering 3D human body motion from monocular video sequences based on a robust image matching metric, incorporation of joint limits and non-self-intersection constraints, and a new sample-and-refine search strategy guided by rescaled cost-function covariances. Monocular 3D body tracking is challenging: besides the difficulty of matching an imperfect, highly flexible, self-occluding model to cluttered image features, realistic body models have at least 30 joint parameters subject to highly nonlinear physical constraints, and at least a third of these degrees of freedom are nearly unobservable in any given monocular image. For image matching we use a carefully designed robust cost metric combining robust optical flow, edge energy, and motion boundaries. The nonlinearities and matching ambiguities make the parameter-space cost surface multi-modal, ill-conditioned and highly nonlinear, so searching it is difficult. We discuss the limitations of CONDENSATION-like samplers, and describe a novel hybrid search algorithm that combines inflated-covariance-scaled sampling and robust continuous optimization subject to physical constraints and model priors. Our experiments on challenging monocular sequences show that robust cost modeling, joint and self-intersection constraints, and informed sampling are all essential for reliable monocular 3D motion estimation.

Keywords: 3D human body tracking, particle filtering, high-dimensional search, constrained optimization, robust matching.

1 Introduction

Extracting 3D human motion from natural *monocular* video sequences poses difficult modeling and computation problems:

- (i) Even a minimal human model is very complex, with at least 30 joint parameters and many more body shape ones, subject to highly nonlinear joint limits and non-self-intersection constraints.
- (ii) Matching a complex, imperfectly known, self-occluding model to a cluttered scene is inherently difficult. Typical

loose clothing only complicates matters.

(iii) In contrast to simplified 2D approaches (Cham and Rehg, 1999; Ju et al., 1996) and the multi-camera 3D case (Kakadiaris and Metaxas, 1996; Gavrilu and Davis, 1996; Bregler and Malik, 1998; Delamarre and Faugeras, 1999; Plankers and Fua, 2001; Drummond and Cipolla, 2001), the estimation problem is extremely ill-conditioned, with at least 1/3 of the 30+ degrees of freedom remaining very nearly unobservable in any given monocular image. The most important non-observabilities are motions of major body segments in depth (*i.e.* towards or away from the camera — these account for 1/3 of the 3D d.o.f.), but others include rotations of near-cylindrical limbs about their axes, and internal motions of compound joints like the spine or shoulder that are difficult to observe even with 3D data.

(iv) In addition to being ill-conditioned, the monocular estimation problem is highly multi-modal. In particular, for any given set of image projections of the 3D joint centers, there are typically some thousands of possible inverse kinematics solutions for the 3D body configuration¹. Under any reasonable model-image matching cost metric, each kinematic solution produces a corresponding local minimum in configuration space, and correspondence ambiguities only compound this number of minima. In practice, choosing the wrong minimum rapidly leads to mistracking, so reliable tracking requires a powerful multiple hypothesis tracker capable of finding and following a significant number of minima. The development of such a tracker is one of the main contributions of this paper. Some more recent work, not reported

¹For each body segment, for any given depth for its top (innermost) endpoint, the bottom endpoint can be aligned with its image projection either in a 'sloped forwards' configuration, or in a 'sloped backwards' one. A full body model contains at least 10 main body segments, and hence has at least $2^{10} = 1024$ possible inverse kinematics solutions (sets of forwards/backwards segment configurations). See Lee and Chen (1985) and the empirical confirmations in Sminchisescu and Triggs (2002a,b).

here, further enhances tracking reliability by explicitly enumerating the possible kinematic minima (Sminchisescu and Triggs, 2003).

Also note that these four difficulties interact strongly in practice. For example, minor modeling or matching errors tend to lead to large compensatory biases in hard-to-estimate depth parameters, which in turn cause mis-prediction and tracking failure. Hence, we believe that a successful monocular 3D body tracking system must pay attention to all of them.

Organisation: §1.1 discusses several existing approaches to human articular tracking, explaining why we believe that they are not suitable for the difficult 3D-from-monocular case and informally motivating our new tracker. §2 briefly describes our 3D body model, which includes full 3D occlusion prediction, joint angle limits and body non-self-intersection constraints. §3 discusses our robust model-image matching framework, which combines robust optical flow, edge energy, and motion boundaries. §4 details our hybrid search / tracking scheme, which combines a mixture density propagation tracker with carefully shaped cost-sensitive sampling, with robust constraint-respecting local optimization. §5 briefly describes the local optimization schedule we use to find initial 3D body poses and internal body proportions from model/image joint correspondence input. §7 details some experiments on challenging monocular sequences. These illustrate the need for each of robust cost modeling, joint and self-intersection constraints, and well-controlled sampling plus local optimization. We end the paper with discussions of the effect of the sampling regime on search efficiency (§7) and approximation accuracy (§8), and ideas for future work.

1.1 High-Dimensional Tracking Strategies

Locating good poses in a high-dimensional body configuration space is intrinsically difficult. Three main classes of search strategies exist: **local descent** incrementally improves an existing estimate, *e.g.* using local Newton strategies to predict good search directions (Bregler and Malik, 1998; Rehg and Kanade, 1995; Kakadiaris and Metaxas, 1996; Wachter and Nagel, 1999); **regular sampling** evaluates the cost function at a predefined pattern of points in (a slice of) parameter space, *e.g.* a local rectangular grid (Gavrila and Davis, 1996); and **stochastic sampling** generates random

sampling points according to some hypothesis distribution encoding “good places to look”, *e.g.* (Deutscher et al., 2000; Sidenbladh et al., 2000). Densely sampling the entire parameter space would in principle guarantee a good solution, but it is infeasible in more than 2–3 dimensions. In 30 dimensions any feasible sample must be extremely sparse, and hence likely to miss significant cost minima. Local descent does at least find a local minimum, but with multimodality there is no guarantee that the globally most representative ones are found. Whichever method is used, effective focusing is the key to high-dimensional search. This is an active research area (Deutscher et al., 2000; Heap and Hogg, 1998; Cham and Rehg, 1999; Merwe et al., 2000), but no existing method can guarantee global minima.

During tracking the search method is applied time-recursively, the starting point(s) for the current search being obtained from the results at the previous time step, perhaps according to some noisy dynamical model. To the (often limited!) extent that the dynamics and the image matching cost are statistically realistic, Bayes-law propagation of a probability density for the true state is possible. For linearized unimodal dynamics and observation models under least squares / Gaussian noise, this leads to Extended Kalman Filtering. For likelihood-weighted random sampling under general multimodal dynamics and observation models, bootstrap filters (Gordon et al., 1993; Gordon and Salmond, 1995) or CONDENSATION (Isard and Blake, 1998) result. In either case various model parameters must be tuned and it sometimes happens that physically implausible settings are needed for acceptable performance. In particular, to control mistracking caused by correspondence errors, selection of slightly incorrect inverse kinematics solutions, and similar model identification errors, visual trackers often require exaggerated levels of dynamical noise. The problem is that even quite minor errors can pull the state estimate a substantial distance from its true value, especially if they persist over several time steps. Recovering from such an error requires a state space jump greater than any that a realistic random dynamics is likely to provide, whereas using an exaggeratedly noisy dynamics provides an easily controllable degree of local randomization that often allows the mistracked estimate to jump back onto the right track. Boosting the dynamical noise does have the side effect of reducing the information propagated from past observations, and hence increasing the

local uncertainty associated with each mode. But this is a small penalty to pay for reliable tracking lock, and in any case the loss of accuracy is often minor in visual tracking, where weak dynamical models (*i.e.* short integration times: most of the state information comes from current observations and dynamical details are unimportant) are common.

In summary, in multi-modal problems, sample based Bayesian trackers often get trapped into following incorrect local minima, and some form of explicit local (but not *too* local!) search must be included to rescue them. For trackers operating in this “memoryless step and search” regime, the machinery of Bayes-law propagation is superfluous — the dynamical model is not correct in any case — and it is simpler to think in terms of sequential local search rather than tracking and noisy dynamics. It seems that many, if not most, existing Bayesian trackers in vision operate essentially in this regime, and the current paper is no exception. Hence, we will assume only weak zeroth order dynamical models and use the language of search rather than tracking. But this is largely a matter of terminology, and more elaborate dynamical models are trivial to incorporate if desired.

Many existing human trackers silently inflate the dynamical noise as a local search mechanism, *e.g.* (Cham and Rehg, 1999; Heap and Hogg, 1998; Deutscher et al., 2000). But in each of these papers, it is only one component of the overall search strategy. The randomization provided by noise inflation is an effective search strategy only for relatively low-dimensional problems, where the samples can cover the surrounding neighborhood fairly densely. In high dimensions, volume increases very rapidly with radius, so any sample that is spread widely enough to reach nearby minima must necessarily be extremely sparse. Hence, the samples are most unlikely to hit the small core of low cost values surrounding another minimum: if they fall into its basin of attraction at all, they are much more likely do so at a high cost point, simply because high cost points are far more common. This is fatal for CONDENSATION-style weighted resampling: high cost points are very unlikely to be resampled, so the new minimum is almost certain to be missed even though an independent track started at the sample would eventually condense to the minimum. The moral is that in high dimensions, random sampling alone does not suffice: some form of local optimization of the samples, or at least a delayed decision about whether they are viable or not, is essential

to prevent mistracking. Cham and Rehg (1999); Heap and Hogg (1998) and the current work use explicit descent-based local optimization for this, while (Deutscher et al., 2000) use a simulated annealing like process (which is usually less efficient, although better aligned with the point-based sample-and-evaluate philosophy of pure particle tracking).

The 3D-from-monocular problem has characteristic ill-conditioning associated with depth degrees of freedom, whereas transverse degrees of freedom are directly observable and hence relatively well conditioned. It also has large numbers of kinematic local minima related by motions in depth, in addition to the minima in transversal directions produced by correspondence ambiguities. Hence, we would like to ensure a thorough, perhaps even a preferential, search along the hard-to-estimate depth degrees of freedom. The problem is that the two sets of directions have very different properties and scales. Precisely because they have such similar image appearances, related kinematic minima may cause confusion even if they are separated by significant distances in parameter space, whereas false-correspondence minima only cause confusion if they are relatively nearby. In other words, the natural metric for tracker confusion — and hence for the sampling distribution of the randomized local search — is perceptual image distance, not parameter space distance. This holds notwithstanding the fact that large jumps in configuration (depth) are improbable under natural human dynamics: the tracker may have been gradually misled over a period of time, and it is essential that it should be able to jump far enough to recover before tracking fails entirely².

This suggests that we need to inflate the dynamical noise preferentially along the depth directions. But these depend strongly on where the model is viewed from, so no constant (configuration or camera-position independent) noise inflation suffices here. The simplest way to adapt the noise to the configuration/camera-position is to estimate the covariance of the posterior likelihood and use this for noise scaling. (In fact, we advocate inflating the *prior* covariance — the previous posterior after dynamics with *physically realistic* noise levels — *i.e.* there should be both realistic dynamics and some degree of deliberate random search). Evaluating

²Ideally, a subsequent smoothing process would push the corrective jump back in time to where the error first occurred (where the jump presumably becomes small). But whether or not this is done, the likelihood penalty for following an incorrect path arbitrarily far forwards in time is likely to be greater than that for any single corrective jump, bad as this may be.

covariances might be burdensome in a conventional particle tracking framework where we only had point samples of likelihoods, but we have already seen that some form of local refinement of the samples is practically essential in high dimensions, and efficient local optimizers require (and in the case of quasi-Newton style methods, even provide) information equivalent to covariance estimates.

The emphasize how much difference covariance scaling can make, consider the 32 d.o.f. cost spectrum in fig. 5 on page 12, which has a 2000:1 range of principal standard deviations. For inflation large enough to double the sampling radius along the most uncertain direction (*e.g.*, for a modest search for local minima along this cost valley), a scaling based on uniform dynamical noise would produce a search volume 10^{54} times larger than that of our prior-based one, and an overwhelming fraction of these samples would have extremely high cost and images implausibly different from the source image (see also fig. 1 on page 6). Such wastage factors are clearly untenable. In practice, samplers based on inflating non-covariance-based dynamical noises simply can not sample deeply enough along the most uncertain (depth) directions to find the local minima there, and frequent mis-tracking is the result.

Finally, given that we are including a component of covariance-scaled but inflated noise expressly as a local search mechanism, what kinds of noise distributions will give the most efficient search? Basically, we need to keep a reasonably large proportion of the samples focused on the current track, while scattering the others fairly widely in the hope of finding other good tracks. Also, volume increases very rapidly with radius in high dimensions, so (even with local optimization) we can not hope to sample densely enough to provide effective search coverage at large inflation factors. It is preferable to choose a moderate inflation level, even though this only provides access to relatively nearby local minima.

In summary, owing to its high dimensionality and the ill-conditioning and multi-modality associated with unobservable depth degrees of freedom, we believe that reliable 3D-from-monocular human body tracking requires deliberate sampling (or some other form of local search) in a region shaped by, but significantly larger than, the local state covariance, followed by local optimization of the samples before any resampling step.

1.2 Previous Work

Below we will compare our method to several existing ones, which we briefly summarize here without attempting a full literature review. 3D body tracking from monocular sequences is significantly harder than 2D (Cham and Rehg, 1999; Ju et al., 1996) or multi-camera 3D (Kakadiaris and Metaxas, 1996; Gavrilu and Davis, 1996; Bregler and Malik, 1998; Delamarre and Faugeras, 1999; Plankers and Fua, 2001; Drummond and Cipolla, 2001) tracking, and surprisingly few works have addressed it (Deutscher et al., 2000; Sidenbladh et al., 2000; Wachter and Nagel, 1999; Gonglaves et al., 1995; Howe et al., 1999; Brand, 1999).

Deutscher et al. (2000) uses a sophisticated ‘annealed sampling’ strategy and a cross-over operator (Deutscher et al., 2001) to speed up CONDENSATION. He reports very good results for unconstrained full-body motion, but for his main sequence he uses 3 cameras and a black background to limit the impact of the alternative minima produced by clutter and depth ambiguities. Sidenbladh et al. (2000) uses a similar importance sampling technique with a strong learned prior walking model or a database of motion snippets (Sidenbladh et al., 2002) to track a walking person in an outdoor monocular sequence. Subsequent work (Sidenbladh and Black, 2001) integrates flow, edge and ridge cues using Laplace like error distributions learned from training data, and shows improved upper body tracking for a subject performing planar motion in a cluttered scene, acquired with a moving camera. Our current method uses no motion model — we optimize static poses — but it is true that when they hold, prior motion models are very effective tracking stabilizers. It is possible, but expensive, to track using a bank of motion models (Blake et al., 1999). Partitioned sampling (MacCormick and Isard, 2000) is another notable sampling technique for articulated models, under certain labeling assumptions (MacCormick and Isard, 2000; Deutscher et al., 2000).

Several authors address the difficulty that the sampling-based searches of pure particle filtering converge rather slowly to modes (Pitt and Shephard, 1997; Heap and Hogg, 1998; Cham and Rehg, 1999; Merwe et al., 2000; Choo and Fleet, 2001), especially when the observation likelihood peaks deep in the tail of the prior. This is especially problematic in high dimensions, where prohibitively long sampling runs are often required for convergence. Heap and Hogg (1998); Cham and Rehg (1999); Merwe et al. (2000)

all combine CONDENSATION-style sampling with either local optimization or Kalman filtering, while Pitt and Shephard (1997) samples discretely using the current observation likelihood (and not the transition prior). The visual trackers of Heap and Hogg (1998) and Cham and Rehg (1999) combine CONDENSATION-style sampling with least-squares optimization, but they only consider the simpler (and much better conditioned) case of 2D tracking. Cham & Rehg combine their heuristic 2D Scaled Prismatic Model (SPM) body representation with a first order motion model and a piecewise Gaussian resampling method for the CONDENSATION step. The Gaussian covariances are estimated from the Gauss-Newton approximation at the fitted optima³, but the search region widths are controlled by the traditional method of adding a large dynamical noise ((Cham and Rehg, 1999) section 3.2).

Choo and Fleet (2001) use a stick model (without any shape model) for which 3D-2D joint to image correspondences from motion capture data are available and propose a (gradient-based) Hybrid Monte Carlo sampler that is more efficient than (point-based) CONDENSATION. The method provides more efficient local descent towards the minima, but it is still prone to trapping in sub-optimal local minima.

Wachter and Nagel (1999) use articulated kinematics and a shape model built from truncated cones, and estimate motion in a monocular sequence, using edge and intensity (optical flow) information using an extended Kalman filter. Anatomical joint limits are enforced at the level of the filter prediction, but not during the update step, where they could be violated. They show experiments in an unconstrained environment for a subject wearing normal clothing, tracking motion parallel with the image plane using articulated models with 10-15 d.o.f.

Both Brand (1999) and Howe et al. (1999) pose 3D estimation as a learning and inference problem, assuming that some form of 2D tracking (stick 2D model positions or silhouettes) is available over an entire time-series. Howe et al. (1999) learn Gaussian distributions over short “snippets” of observed human motion trajectories, then use these as priors

³The covariance estimates of nonlinear least-squares optimizers as used by (Heap and Hogg, 1998; Cham and Rehg, 1999) are not robust to model/image matching errors and incorrect (*i.e.* biased) for natural image statistics that have highly non-Gaussian shape with high kurtosis and long tails (Zhu and Mumford, 1997; Sidenbladh and Black, 2001). We use an observation likelihood and a robust local continuous optimizer based on heavy tail error distributions (see §3.1 and §4.1) to address these problems.

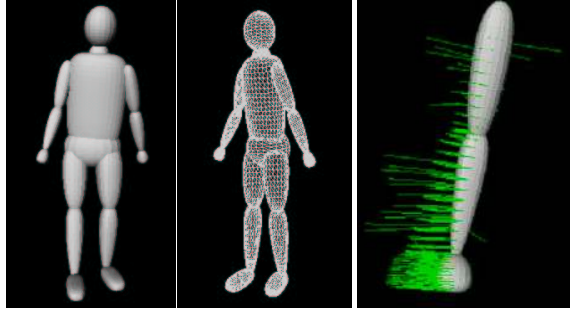


Figure 2: Different body models using for tracking (a,b,c). In (c) the prediction errors for a model configuration are also plotted (per node, for a contour and intensity cost function, see text).

in an EM-based Bayesian MAP framework to estimate new motions. Brand (1999) learns a HMM with piecewise linear states and solves for the MAP estimate using an entropy minimization framework. As presented, these methods are basically monomodal so they can not accommodate multiple trajectory interpretations, and they also rely heavily on their learned-prior temporal models to stabilize the tracking. Nevertheless, they provide a powerful higher-level learning component that is complementary to the framework proposed in this paper.

2 Human Body Model

Our human body model (fig. 2a,b,c) consists of a kinematic ‘skeleton’ of articulated joints controlled by angular **joint parameters** \mathbf{x}_a , covered by ‘flesh’ built from superquadric ellipsoids with additional tapering and bending parameters (Barr, 1984). A typical model has around 30 joint parameters, plus 8 **internal proportion** parameters \mathbf{x}_i encoding the positions of the hip, clavicle and skull tip joints, plus 9 **deformable shape** parameters for each body part, gathered into a vector \mathbf{x}_d . A complete model can be encoded as a single large parameter vector $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_d, \mathbf{x}_i)$. During tracking we usually estimate only joint parameters, but during initialization the most important internal proportions and shape parameters are also optimized, subject to a soft prior based on standard humanoid dimensions obtained from Group (2002) and updated using collected image evidence. This model is far from photorealistic, but it suffices for high-level interpretation and realistic occlusion prediction, offering a good

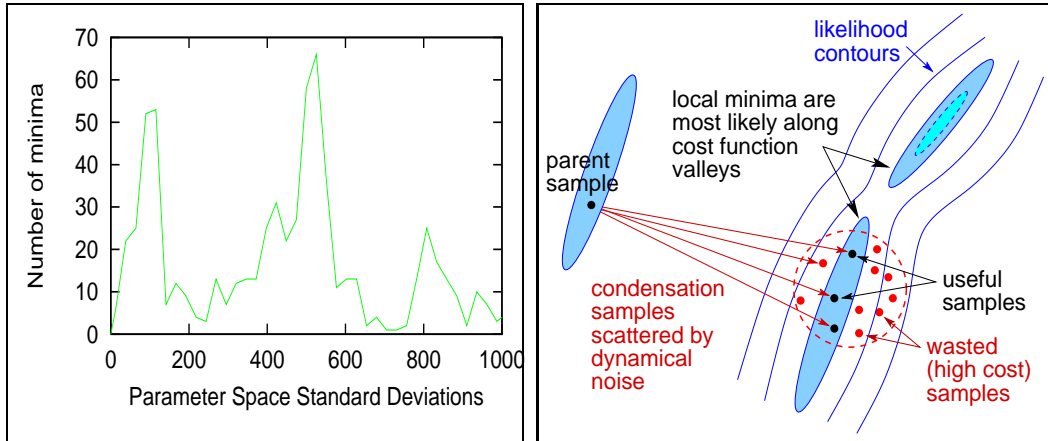


Figure 1: (a) Typical parameter space minima distribution measured with respect to an arbitrary minimum. Notice that the minima are far from each other in parameter space so wide sampling is necessary to find them. However, boosting the dynamics by sampling from the transition prior (as in particle filtering) leads to inefficiencies (b).

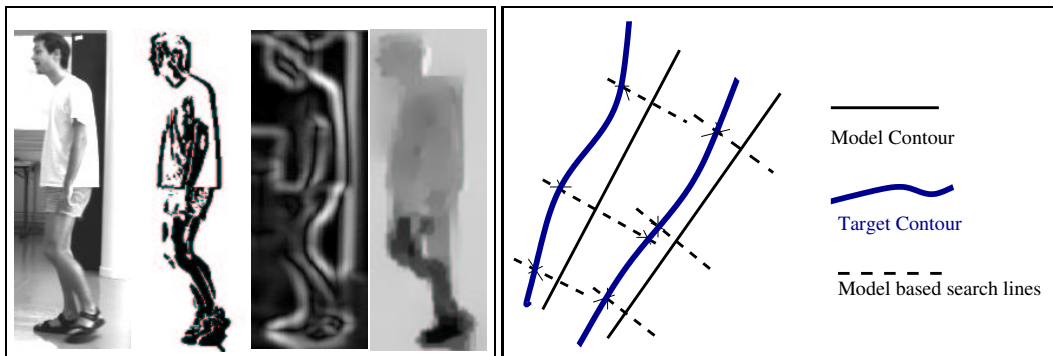


Figure 3: Examples of our robust low-level feature extraction: original image (a), motion boundaries (b), intensity-edge energy (c), robust horizontal flow field (d) and the model-based edge matching process (e). Multiple edge matches found along individual search lines (model projected contour normals) are fused using a probabilistic assignment strategy (see text).

trade-off between computational complexity and coverage.

The model is used as follows. Superquadric surfaces are discretized as meshes parameterized by angular coordinates in a 2D topological domain. Mesh nodes \mathbf{u}_i are transformed into 3D points $\mathbf{p}_i = \mathbf{p}_i(\mathbf{x})$ and then into predicted image points $\mathbf{r}_i = \mathbf{r}_i(\mathbf{x})$ using composite nonlinear transformations:

$$\mathbf{r}_i(\mathbf{x}) = P(\mathbf{p}_i(\mathbf{x})) = P(A(\mathbf{x}_a, \mathbf{x}_i, D(\mathbf{x}_d, \mathbf{u}_i))) \quad (1)$$

where D represents a sequence of parametric deformations that construct the corresponding part in its own reference frame, A represents a chain of rigid transformations that map it through the kinematic chain to its 3D position, and P represents perspective image projection. During model estimation, robust prediction-to-image matching cost metrics are

evaluated for each predicted image feature \mathbf{r}_i , and the results are summed over all features to produce the image contribution to the overall parameter space cost function. We use both direct image-based cost metrics such as robustified normalized edge energy, and extracted feature based ones. The latter associate the predictions \mathbf{r}_i with one or more nearby image features $\bar{\mathbf{r}}_i$ (with additional subscripts if there are several matches). The cost is then a robust function of the prediction errors $\Delta\mathbf{r}_i(\mathbf{x}) = \bar{\mathbf{r}}_i - \mathbf{r}_i(\mathbf{x})$.

3 Problem Formulation

We aim towards a probabilistic interpretation and optimal estimates of the model parameters by maximizing the total

probability according to Bayes rule:

$$p(\mathbf{x}|\bar{\mathbf{r}}) \propto p(\bar{\mathbf{r}}|\mathbf{x})p(\mathbf{x}) = \exp(-\sum_i e(\bar{\mathbf{r}}_i|\mathbf{x})) p(\mathbf{x}) \quad (2)$$

where $e(\bar{\mathbf{r}}_i|\mathbf{x})$ is the cost density associated with the observation of node i and $p(\mathbf{x})$ is a prior on model parameters. In our MAP approach, we discretize the continuous problem and attempt to minimize the negative log-likelihood for the total posterior probability, expressed as the following cost function:

$$f(\mathbf{x}) = -\log(p(\bar{\mathbf{r}}|\mathbf{x})p(\mathbf{x})) \quad (3)$$

$$= -\log p(\bar{\mathbf{r}}|\mathbf{x}) - \log p(\mathbf{x}) = f_o(\mathbf{x}) + f_p(\mathbf{x}) \quad (4)$$

3.1 Observation Likelihood

Whether continuous or discrete, the search process depends critically on the observation likelihood component of the parameter space cost function. Besides smoothness properties, the likelihood should be designed to limit the number of spurious local minima in parameter space. Our method employs a combination of robust edge and intensity information on top of a multiple assignment strategy based on a weighting scheme that focuses attention towards motion boundaries. Our likelihood term is also based on robust (heavy-tailed) error distributions. Note that both robustly extracted image cues and robust parameter space estimation are used: the former provides ‘‘good features to track’’, while the latter directly addresses the model-image association problem.

Robust Error Distributions: Robust parameter estimation can be viewed as the choice of a realistic total likelihood model for the combined inlier and outlier distributions for the observation. We model the total likelihood in terms of robust radial terms ρ_i , where $\rho_i(s)$ can be any increasing function with $\rho_i(0) = 0$ and $\frac{d}{ds}\rho_i(0) = \frac{\nu}{\sigma^2}$. These model error distributions corresponding to a central peak with scale σ , and a widely spread background of outliers ν . Here we used the ‘Lorentzian’ $\rho_i(s, \sigma) = \nu \log(1 + \frac{s^2}{\sigma^2})$ and ‘Leclerc’ $\rho_i(s, \sigma) = \nu(1 - \exp(-\frac{s^2}{\sigma^2}))$ robust error potentials.

The cost for the observation i , expressed in terms of corresponding model prediction is $e(\bar{\mathbf{r}}_i|\mathbf{x}) = \frac{1}{N\nu}e_{ui}(\mathbf{x})$, where N is the total number of model nodes, \mathbf{W}_i is a positive definite weighting matrix associated to the assignment i , and:

$$e_i(\mathbf{x}) = \begin{cases} \frac{1}{2}\rho_i(\Delta\mathbf{r}_i(\mathbf{x})\mathbf{W}_i\Delta\mathbf{r}_i(\mathbf{x})^\top) & \text{if } i \text{ is assigned} \\ \nu_{bf} = \nu & \text{if back-facing} \\ \nu_{occ} = k\nu, \quad k > 1 & \text{if occluded} \end{cases} \quad (5)$$

The robust observation likelihood contribution is thus:

$$f_o(\mathbf{x}) = -\log p(\bar{\mathbf{r}}|\mathbf{x}) \quad (6)$$

$$= f_a(\mathbf{x}) + N_{bf}\nu_{bf} + N_{occ}\nu_{occ} \quad (7)$$

where $f_a(\mathbf{x})$ represents the term associated with the image assigned model nodes, while N_{occ} and N_{bf} are the numbers of occluded and back-facing (self-occluded) model nodes.

Notice that occluded model predictions are not simply ignored. They contribute a constant penalty to the overall observation likelihood. This is necessary in order to build likelihoods that preserve their response properties under occlusion and viewpoint change. For instance, good fits from both frontal and side views should ideally have similar peak responses, but it is clear that the number of occluded model points is in general larger in a side view than in a frontal one. This can lead to down-weighting of peaks for side views if only the visible nodes are taken into account. An additional difficulty arises, for example, in cases where the legs pass each other (in a side-view) and the model ‘locks’ both of its legs onto the same image leg. To avoid such situations, we include all of the model nodes when fusing the likelihood, but we slightly penalize occluded ones in order to make them less attractive. A way to choose the occlusion penalty ν is to fit the model to the data and compute an approximate error per node. By using a slightly higher value for occluded nodes, we make them more attractive than a bad fit but less attractive than other non-occluded states that can exist in the neighborhood of the parameter space. We find this heuristic gives good results in practice⁴, although a more rigorous treatment of occlusion would be desirable in the general case. At present, this is computationally too expensive, but interesting approximations can be found in MacCormick and Blake (1998).

Cue Integration and Assigned Image Descriptors: We use both edge and intensity features in our cost function (see Sminchisescu (2002b) for details). For edges, the images are smoothed with a Gaussian kernel, contrast normalized, and a Sobel edge detector is applied. For intensities, a robust multi-scale optical flow method based on Black and Anandan (1996) implementation gives both a flow field and an associated outlier map (see fig. 3b). The outlier map is processed

⁴This is particularly effective when combined with the Covariance Scaled Sampling (CSS) algorithm presented in §4. Loss of visibility of certain body parts leads to increased uncertainty in related parameters, and CSS automatically ensures broader sampling in those parameter space regions.

similar to edges, to obtain a smooth 2D potential field S_p . It conveys useful information about the motion boundaries and is used to weight the significance of edges (see fig. 3b). We typically use diagonal weighting matrices \mathbf{W}_i , associated with the predicted feature \mathbf{r}_i and corresponding matched observation $\bar{\mathbf{r}}_i$, of the form $\mathbf{W}_i(\mathbf{r}_i) = 1 - kS_p(\bar{\mathbf{r}}_i)$, where k is a constant that controls the emphasis and confidence in the motion boundary estimation. (The smoothed motion boundary image is a real image with values between 0 and 1 as in fig. 3b. For instance, $k = 0$ will weight all the edges uniformly, while $k = 1$ will entirely exclude the edge responses that are not on motion boundaries). In practice, we found that values of k in the range $k = 0.2$ – 0.4 worked well. For visible nodes on model occluding contours (\mathcal{O}), we perform line search along the normal and retain all possible assignments within the search window (see fig. 3e), weighting them by their importance qualified by the motion boundary map \mathbf{W} . For visible model nodes lying inside the object (\mathcal{I}), we use the correspondence field derived from the robust optical flow at their corresponding image prediction. This acts as a residual measurement error at each visible model node (see Sminchisescu (2002b) for details). The assigned data term (6) thus becomes:

$$f_a(\mathbf{x}) = \frac{1}{2} \sum_{i \in \mathcal{O}, e \in \mathcal{E}_i} \rho_{i_e}(\Delta \mathbf{r}_{i_e}(\mathbf{x}) \mathbf{W}_{i_e} \Delta \mathbf{r}_{i_e}(\mathbf{x})^\top) \quad (8)$$

$$+ \frac{1}{2} \sum_{j \in \mathcal{I}} \rho_{j_f}(\Delta \mathbf{r}_{j_f}(\mathbf{x}) \mathbf{W}_{j_f} \Delta \mathbf{r}_{j_f}(\mathbf{x})^\top) \quad (9)$$

where the subscripts “ i_e ” denote multiple edges \mathcal{E}_i assigned to model prediction i , and “ j_f ” denote the flow term assigned to model prediction j .

3.2 Model Priors

The complete prior penalty over model parameters is a sum of negative log likelihoods $f_p = f_{an} + f_s + f_{pa}$ corresponding to the following prior densities p_{an}, p_s, p_{pa} :

Anthropometric data p_{an} : The internal proportions for a standard humanoid (based on statistical measurements) are collected from (Group, 2002) and used effectively as a Gaussian prior, $p_{an} = \mathcal{N}(\boldsymbol{\mu}_{an}, \boldsymbol{\Sigma}_{an})$, to estimate a concrete model for the subject to be tracked. Left-right symmetry of the body is assumed: only “one side” of the internal proportions parameters are estimated while collecting image measurements from the entire body.

Parameter stabilizers p_s : Certain modeling details are far more important than one might think. For example, it is impossible to track common turning and reaching motions unless the clavicle joints in the shoulder are modeled accurately. However, these parameters have fairly well defined equilibrium positions and leaving them unconstrained would often lead to ambiguities that produce nearly singular (flat) cost surfaces. We control these hard-to-estimate parameters with long-tailed “sticky prior” stabilizers scaling their Gaussian equilibria, $p_s = \mathcal{N}(\boldsymbol{\mu}_s, \boldsymbol{\Sigma}_s)$. This ensures that in the absence of strong observations, the parameters are constrained to lie near their default values, whereas stronger observations can “unstuck” them from the defaults and effectively turn off the prior.

Anatomical joint angle limits C_{bl} : 3D consistency requires that the values of joint angles evolve within anatomically consistent intervals. Also, when estimating internal body proportions during initialization, we ensure that they remain within a certain range of deviation from the standard humanoid (typically 10%). We model this with a set of inequalities of the form $C_{bl} \cdot \mathbf{x} < 0$, where C_{bl} is a ‘box-limit’ constraint matrix.

Body part interpenetration avoidance p_{pa} : Physical consistency requires that different body parts do not interpenetrate during estimation. We avoid this by introducing repulsive potentials that decay rapidly outside the surface of each body part, $f_{pa} = \exp(-f(\mathbf{x}) |f(\mathbf{x})|^{p-1})$, where $f(\mathbf{x}) < 0$ defines the interior of the part and p controls the decay rate.

3.3 Distribution Representation

We represent parameter space distributions as sets of separate modes $m_i \in \mathcal{M}$, each having an associated overall probability, mean and covariance matrix $m_i = (\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i, c_i)$. These can be viewed as Gaussian mixtures. Cham & Rehg (Cham and Rehg, 1999) also use multiple Gaussians, but they had to introduce a special piecewise representation as their modes seem to occur in clusters after optimization. We believe that this is an artifact of their cost function design. In our case, as the modes are the result of robust continuous optimization, they are necessarily either separated or confounded. Our 3D-from-monocular application also requires a more effective sampling method than the 2D one of Cham and Rehg (1999), as explained in §4.2.

3.4 Temporal Propagation

Equation 2 reflects the search for the model parameters in a static image, under likelihood terms and model priors but without a temporal or initialization prior. For temporal observations $\mathbf{R}_t = \{\bar{\mathbf{r}}_1, \bar{\mathbf{r}}_2, \dots, \bar{\mathbf{r}}_t\}$, and sequence of states $\mathbf{X}_t = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t\}$, the posterior distribution over model parameters becomes:

$$p(\mathbf{x}_t | \mathbf{R}_t) \propto p(\bar{\mathbf{r}}_t | \mathbf{x}_t) p(\mathbf{x}_t) \int_{\mathbf{x}_{t-1}} p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{x}_{t-1} | \mathbf{R}_{t-1}) \quad (10)$$

where $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ is a dynamical prior and $p(\mathbf{x}_{t-1} | \mathbf{R}_{t-1})$ is the prior distribution from $t - 1$. Together they form the temporal prior $p(\mathbf{x}_t | \mathbf{R}_{t-1})$ for initializing the static image search (2)⁵.

4 Search Algorithm

Our parameter search technique combines robust constraint-consistent local optimization with a more global discrete sampling method.

4.1 Mode Seeking using Robust Constrained Continuous Optimization

The cost function is a negative log likelihood. In order to optimize a sample \mathbf{x} to find the center of its associated likelihood peak, we employ an iterative second order robust constrained local optimization procedure. At each iteration, the log-likelihood gradients and Hessians of the observations and the soft priors⁶ are assembled from (3):

$$\mathbf{g} = \frac{d f}{d \mathbf{x}} = \mathbf{g}_o + \nabla f_{an} + \nabla f_s + \nabla f_{pa} \quad (11)$$

$$\mathbf{H} = \frac{d^2 f}{d \mathbf{x}^2} = \mathbf{H}_o + \nabla^2 f_{an} + \nabla^2 f_s + \nabla^2 f_{pa} \quad (12)$$

For local optimization we use a second order trust region method, where a descent direction is chosen by solving the

⁵In practice, at any given time step we work on a negative log-likelihood ‘energy’ function that is essentially static, being based on both the current observation likelihood and the parameter space priors, as in (3) on page 7. The samples from the temporal prior $p(\mathbf{x}_t | \mathbf{R}_{t-1})$ are used as initialization seeds for local energy minimization. The different minima found will represent the components of the posterior mixture representation.

⁶‘Soft’ means that these terms are part of the cost surface, whereas ‘hard’ constraints such as joint limits restrict the range of variation of their corresponding parameters.

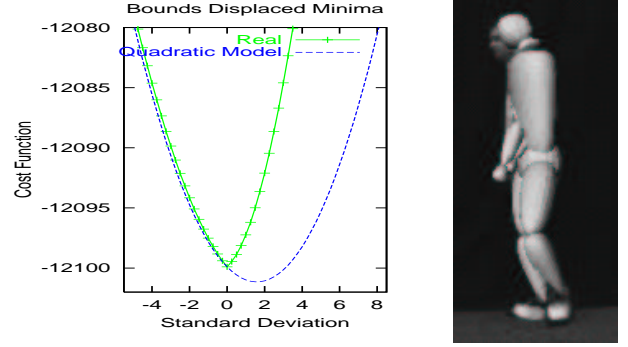


Figure 4: (a) Displaced minimum due to joint limits constraints, (b) Joint limits without body non-self-intersection constraints do not suffice for physical consistency.

regularized subproblem (Fletcher, 1987):

$$(\mathbf{H} + \lambda \mathbf{W}) \delta \mathbf{x} = -\mathbf{g} \quad \text{subject to} \quad C_{bl} \cdot \mathbf{x} < 0 \quad (13)$$

where \mathbf{W} is a symmetric positive definite damping matrix and λ is a dynamically chosen weighting factor. Joint limits C_{bl} are handled as hard bound constraints in the optimizer, by projecting the gradient onto the currently active (*i.e.* currently unlimited) variables. The joint constraints change the character of the cost function and the minima reached very significantly. Fig. 4 plots a 1D slice through the constrained cost function together with a second order Taylor expansion of the unconstrained cost. Owing to the presence of the bounds, the cost gradient is nonzero (orthogonal to the active constraints) at the constrained minimum. The unconstrained cost function is smooth, but the constrained one changes gradient abruptly when a constraint is hit, essentially because the active-set projection method changes the motion direction to maintain the constraint.

4.2 Covariance Scaled Sampling

Although representations based on propagating multiple modes, hypotheses or samples tend to increase the robustness of model estimation, the great difficulty with high-dimensional distributions is finding a sampleable proposal density that often hits their **typical sets** — the areas where most of their probability mass is concentrated. Here we develop a proposal density based on local parameter estimation uncertainties⁷. The local sample optimizations give us

⁷Related variable metric ideas can be found in global optimization, in the context of continuous annealing (Vanderbilt and Louie, 1984) and have been

not only local modes, but also their (robust, constraint consistent) Hessians and hence estimates of the local posterior parameter estimation uncertainty at each mode⁸.

The main insight is that alternative cost minima are most likely to occur along local valleys in the cost surface, *i.e.* along highly uncertain directions of the covariance. It is along these directions that cost-modeling imperfections and noise, and 3D nonlinearities and constraints, have the most likelihood of creating multiple minima, as the cost function is shallowest and the 3D movements are largest there. This is particularly true for monocular 3D estimation, where the covariance is unusually ill-conditioned owing to the many poorly observable motion-in-depth d.o.f. Some examples of such multimodal behavior along high covariance eigen-directions are given in fig. 7. Also, it is seldom enough to sample at the scale of the estimated covariance. Samples at this scale almost always fall back into the same local minimum, and significantly deeper sampling is necessary to capture nearby but non-overlapping modes lying further up the valley⁹. Hence, we sample according to rescaled covariances, typically scaling by a factor of 8 or so. Finally, one can sample either randomly, or according to a regular pattern¹⁰. For the experiments showed here, we use random sampling using CSS with Gaussian tails. Fig. 6 summarizes the resulting covariance-scaled search method.

Given the explanations above, we must implement the following steps:

- (i) Generate fair samples from a prior with known modes.

applied by Black (1992) to low-dimensional (2D) optical flow computation.

⁸A sample is optimized to convergence to obtain the corresponding mode. The Hessian matrix at the convergence mode gives the principal curvature directions and magnitude around the mode and its inverse gives the covariance matrix, reflecting the cost local uncertainty structure. The Hessian is estimated by the algorithm §4.1 during optimization (using (11)), and the covariance is readily obtained from there.

⁹In part this is due to imperfect modeling, which easily creates biases greater than a few standard deviations, particularly in directions where the measurements are weak. Also, one case in which multiple modes are likely to lie so close together in position and cost that they cause confusion is when a single mode fragments due to smooth evolutions of the cost surface. In this case, singularity ('catastrophe') theory predicts that generically, exactly two modes will arise (bifurcation) and that they will initially move apart very rapidly (at a speed proportional to $1/\sqrt{t}$). Hence, it is easy for one mode to get lost if we sample too close to the one we are tracking.

¹⁰For efficiency purposes, an implementation could sample regularly, in fact only along lines corresponding to the lowest few covariance eigen-directions. Although this gives a very sparse sampling indeed, this is an avenue that can be explored in practice.

This is easy. In our case we propagate Gaussian mixtures¹¹ can be used as importance sampling distributions, and correction weighting is readily performed. Mixtures provide a compact, explicitly multi-modal representation and accurate localization, advantages emphasized by Heap and Hogg (1998) and Cham and Rehg (1999) (§5.1). However, both papers use sampling stages based on the unmodified process model (*i.e.* dynamics with fixed, near-spherical noise), which therefore have trapping and sample wastage problems analogous to CONDENSATION.

(ii) Recover new modes of a distribution for which only *some* of the modes are known. This is significantly more difficult. A-priori, the distribution of unknown modes is not available, nor are the boundaries of the basins of attraction of the existing modes (in order to find their neighbors). Also, such likelihood peaks are often well-separated in configuration space (*e.g.* the forwards/backwards flipping ambiguities for human pose, or the cascades of incorrect matches when a model limb is assigned to the incorrect side of an image limb). For typical distributions of minima in parameter space and in cost, see fig. 13 on page 18 and the results in table 1 on page 19. For well-separated peaks, sampling based purely on the known (and potentially incomplete) ones is inadequate, as most of the samples will simply fall back into the peaks they arose from¹². So broader sampling is necessary, but it is also important to focus the samples in relatively low cost regions (see also fig. 1). To achieve this we propose to use the local cost surface to shape a broad sampling distribution. As expected on theoretical grounds, this turns out to give significantly improved results for CSS (for sample cost

¹¹There are at least two ways to obtain a mixture. One is by clustering a set of posterior samples generated, *e.g.*, by CONDENSATION updates. This may produce centers that are not necessarily well-separated, and that may not actually reflect the true modes of the posterior owing to sampling artifacts. Another possibility, followed here, is to optimize the samples locally. In this case the modes found are true local peaks that are, necessarily, either separated or confounded.

¹²Several metrics exist for assessing the efficiency of particle filters (Liu, 1996; MacCormick and Isard, 2000). The 'survival diagnostic' (also called 'effective sample size') measures how many particles will survive a resampling operation. If the weights are unbalanced very few may survive, thus reducing search diversity. But balanced weights do not imply that all peaks have been well explored: samples trapped in a single mode have reasonably well-balanced weights. The same criticism applies to the 'survival rate'. This tries to characterize the ratio of the volume of support of the posterior to that of the prior. Low values suggest that the filter may produce inaccurate density estimates, but again, trapping leaves the survival rate reasonably high.

median, number of minima found, their cost) than competing methods based on either pure prior-based sampling or prior-based sampling plus spherical ‘dynamical’ noise (see table 1 on page 19).

(iii) Sample a prior under dynamic observations but without making restrictive assumptions on the motion of its peaks. In this case the modes from time $t - 1$ are available, and it is critical that the sampling procedure cover the peaks of the observation likelihood in the next time step t . This means that samples should be generated in the basins of attraction of the density peaks *after* applying the dynamical update. In the absence of knowledge about the peaks’ motion (*i.e.* known system dynamics), we exploit the local uncertainty structure in the distribution, and shape the search region based on it. Again, broader sampling is necessary, as the tracked object moves between frames. Also, as explained above, the mode tracking process is not one-to-one. New modes might emerge or split under the effect of increased uncertainty, and it is important that the sampling process does not miss such events by sampling too close to a given mode core, which may both move and split between two temporal observations. Our quantitative results in §6 directly support such findings *e.g.* for mode splitting reflecting bi-modality generated by locally-planar versus in-depth motion explanations (see below).

In this paper, we have not used specific motion models as we want to be able to track general human motions (see for instance, the sequences given in fig. 9–11). For the experiments shown in the next section, we used trivial driftless diffusion dynamics, so CSS has to account for local uncertainty and sample widely enough to cover moving peaks. One could also use constant velocity dynamics, or more sophisticated learned motion models such as walking Rohr (1994); Deutscher et al. (2000); Sidenbladh et al. (2000). When they hold, such models can significantly stabilize tracking, but note that they often turn out to be misleading, *e.g.* when the subject makes unexpected motions like turning or switching activities.

To build up intuition about the shape of our cost surface, we studied it empirically by sampling along uncertain covariance directions (in fact eigenvectors of the covariance matrix), for various model configurations. With our carefully selected image descriptors, the cost surface is smooth apart from the apparent gradient discontinuities caused by active-

set projection at joint constraint activation points. Hence, our local optimizer reliably finds a local minimum. We find that multiple modes do indeed occur for certain configurations, usually separated by cost barriers that a classical (uninflated) sampling strategy would have difficulty crossing. For example, fig. 7 shows the two most uncertain modes of the fig. 9 human tracking sequence at times 0.8 s and 0.9 s. (These are minima only within the sampled slice of parameter space, but they do lie in the attraction zones of full parameter space minima). Secondary minima like those shown here occur rather often, typically for one of two reasons. The first is incorrect registration and partial loss of track when both edges of a limb model are attracted to the same image edge of the limb. This is particularly critical when there is imperfect body modeling and slightly misestimated depth. The second occurs when the character of a motion in depth is misinterpreted. Image registration is maintained until the incorrect 3D interpretation becomes untenable, at which point recovery is difficult. This situation occurs in fig. 7 (see also fig. 12). Identifying and tracking such ambiguous behaviors is critical, as incorrect depth interpretations quickly lead to tracking failure.

Fig. 8a shows some typical slices along cost eigendirections at much larger scales in parameter space. Note that we recover the expected robust shape of the matching distribution, with some but not too many spurious local minima. This is crucial for efficiency and robustness, as the tracker can only follow a limited number of possible minima.

5 Model Initialization

Our tracker starts with a set of initial hypotheses produced by a model initialization process. Correspondences need to be specified between model joint locations and approximate joint positions of the subject in the initial image, and a non-trivial optimization process is run to estimate certain body dimensions and the initial 3D joint angles. Previous approaches to single-view model initialization (Taylor, 2000; Barron and Kakadiaris, 2000) do not fully address the generality and consistency problems, failing to enforce the joint limit constraints, and assuming either restricted camera models or restricted human poses in the image. An algorithm like the one we propose could also probably be bootstrapped using estimates of 2D joint positions derived from learned mod-

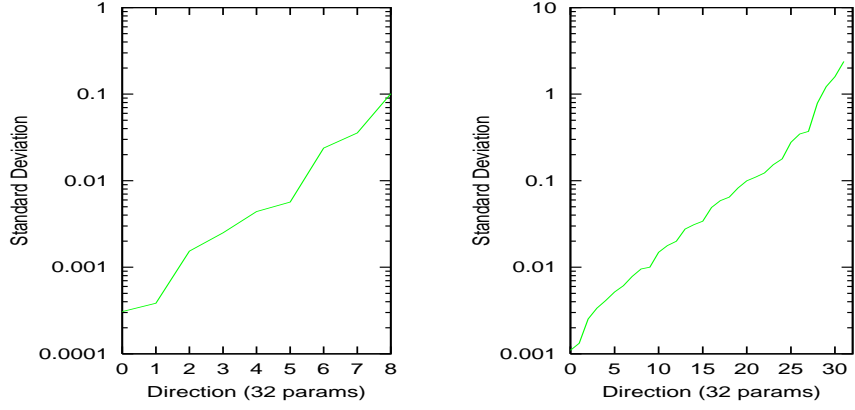


Figure 5: Typical covariance eigenvalue spectra plotted on a logarithmic scale, for a local minimum. $\sigma_{\max}/\sigma_{\min}$ is 350 for the 8 d.o.f. arm model, and 2000 for the 32 d.o.f. body one.

From the ‘old’ mixture prior $p(\mathbf{x}_{t-1}|\mathbf{R}_{t-1}) = \sum_{i=1}^K \pi_i^{t-1} \mathcal{N}(\mu_i^{t-1}, \Sigma_i^{t-1})$, at time $t - 1$, build ‘new’ mixture posterior $p(\mathbf{x}_t|\mathbf{R}_t) = \sum_{i=1}^K \pi_i^t \mathcal{N}(\mu_i^t, \Sigma_i^t)$, at time t , as follows:

1. Build covariance scaled proposal density $p_{t-1}^* = \sum_{i=1}^K \pi_i^{t-1} \mathcal{P}(\mu_i^{t-1}, \Sigma_i^{t-1})$. For the experiments we have used Gaussian tails. The covariance scaled Gaussian component proposals are $\mathcal{P} = \mathcal{N}(\mu_i^{t-1}, s\Sigma_i^{t-1})$ with $s=4-14$ in our experiments.

2. Generate components of the posterior at time t by sampling from p_{t-1}^* as follows. Iterate over $j = 1 \dots N$ until the desired number of samples N are generated:

2.1. Choose component i from p_{t-1}^* with probability π_i^{t-1} .

2.2. Sample from $\mathcal{P}(\mu_i^{t-1}, \Sigma_i^{t-1})$ to obtain \mathbf{s}_j .

2.3. Optimize \mathbf{s}_j over the observation likelihood at time t , $p(\mathbf{x}|\bar{\mathbf{r}}_t)$ defined by (2), using the local continuous optimization algorithm (§4.1). The result is the parameter space configuration at convergence μ_j^t , and the covariance matrix $\Sigma_j^t = \mathbf{H}(\mu_j^t)^{-1}$. If the μ_j^t mode has been previously found by a local descent process, discard it (For notational clarity, without any loss of generality, consider all the modes found are different).

3. Construct an un-pruned posterior for time t as: $p_t^u(\mathbf{x}_t|\mathbf{R}_t) = \sum_{j=1}^N \pi_j^t \mathcal{N}(\mu_j^t, \Sigma_j^t)$ where $\pi_j^t = \frac{p(\mu_j^t|\bar{\mathbf{r}}_t)}{\sum_{j=1}^N p(\mu_j^t|\bar{\mathbf{r}}_t)}$.

4. Prune the posterior p_t^u to keep the best K components with highest probability π_j^t (rename indices $j = 1 \dots N$ into the set $k = 1 \dots K$) and renormalize the distribution as follows: $p_t^p(\mathbf{x}_t|\mathbf{R}_t) = \sum_{k=1}^K \pi_k^t \mathcal{N}(\mu_k^t, \Sigma_k^t)$ where $\pi_k^t = \frac{p(\mu_k^t|\bar{\mathbf{r}}_t)}{\sum_{j=1}^K p(\mu_j^t|\bar{\mathbf{r}}_t)}$.

5. For each mixture component $j = 1 \dots K$ in p_t^p , find the closest prior component i in $p(\mathbf{x}_{t-1}|\mathbf{R}_{t-1})$, according to a Bhattacharyya distance $\mathcal{B}_{ij}(\mu_i^{t-1}, \Sigma_i^{t-1}, \mu_j^t, \Sigma_j^t) = (\mu_i^{t-1} - \mu_j^t)^T \left[\frac{\Sigma_i^{t-1} + \Sigma_j^t}{2} \right]^{-1} (\mu_i^{t-1} - \mu_j^t) + \frac{1}{2} \log \frac{|\Sigma_i^{t-1} + \Sigma_j^t|}{\sqrt{|\Sigma_i^{t-1}| |\Sigma_j^t|}}$. Recompute $\pi_j^t = \pi_j^t * \pi_i^{t-1}$. Discard the component i of $p(\mathbf{x}_{t-1}|\mathbf{R}_{t-1})$ from further consideration.

6. Compute the posterior mixture $p(\mathbf{x}_t|\mathbf{R}_t) = \sum_{k=1}^K \pi_k^t \mathcal{N}(\mu_k^t, \Sigma_k^t)$ where $\pi_k^t = \frac{\pi_k^t}{\sum_{j=1}^K \pi_j^t}$.

Figure 6: The steps of our covariance-scaled sampling algorithm.

els of silhouette appearance (Rosales and Sclaroff, 2000).

For stability, parameters are initialized in three stages,

each based on the formulation described in §4.1. Hard joint

limits are enforced at all stages by the constrained optimiza-

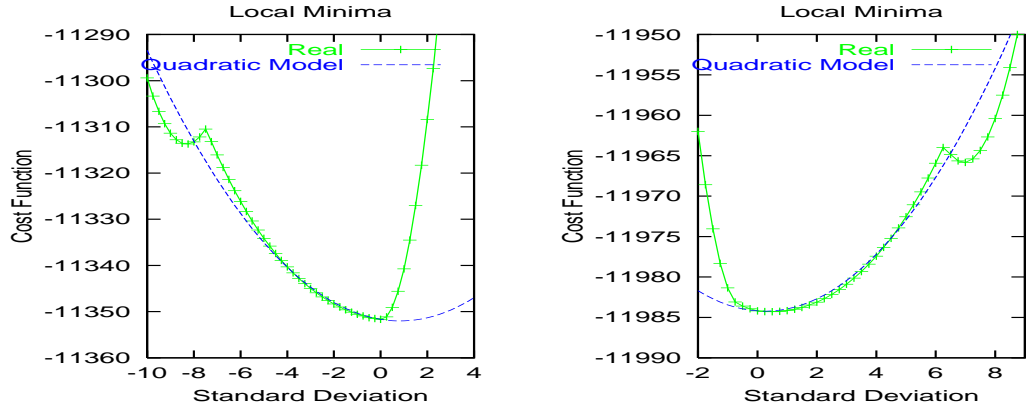


Figure 7: Multimodality along several uncertain eigen-directions (0.8 and 0.9 s in cluttered body tracking sequence).

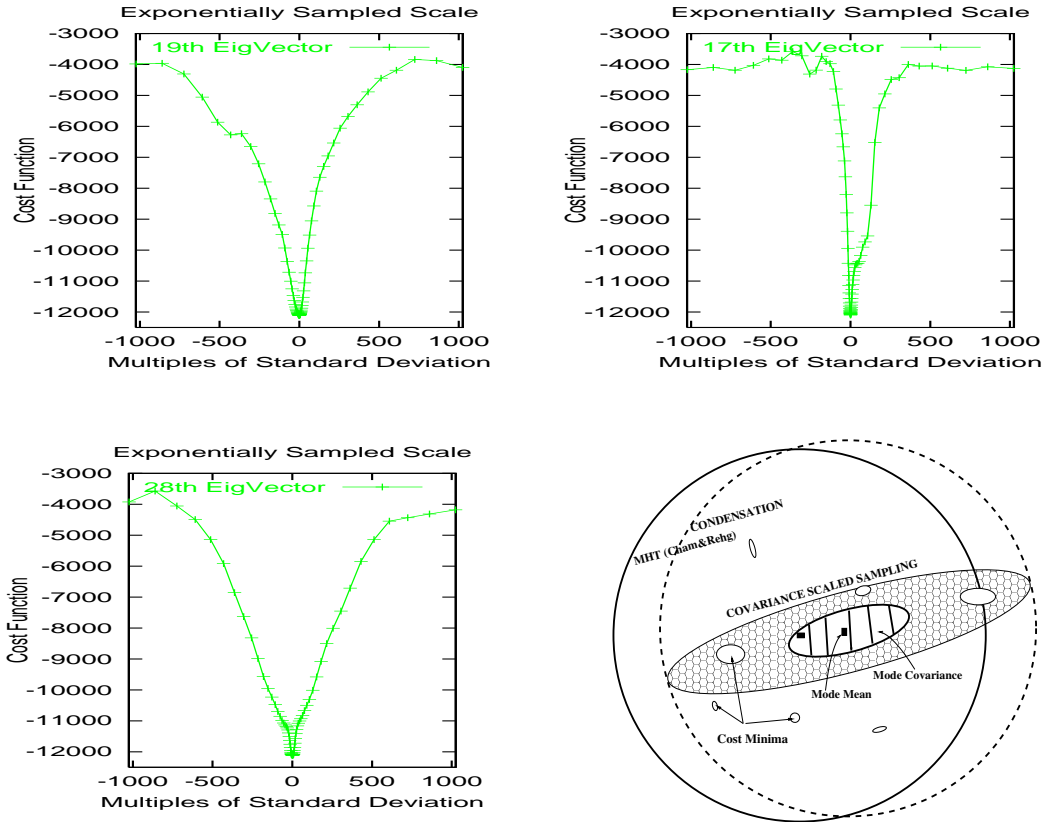


Figure 8: (a,b,c) Cost function slices at large scales, (d) Comparison of sampling methods: (1) CONDENSATION (dashed circle coverage) randomizes each sample by dynamic noise, (2) MHT, solid circle ((Cham and Rehg, 1999), section 3.2 page 3) samples within covariance support (dashed ellipse) and applies the same noise policy as (1), finally, our (3) *Covariance Scaled Sampling* (pattern ellipse) targets good cost minima (flat filled ellipses) by inflating or heavy tail sampling the local robust covariance estimation (dashed ellipse)).

tion procedure, and corresponding parameters on the left and right sides of the body are held equal, whereas measurements are collected from the entire body (see below). The first stage estimates joint angles \mathbf{x}_a , internal proportions \mathbf{x}_i and a few simple shape \mathbf{x}_d parameters, subject to the given 3D to 2D joint correspondences and prior intervals on the internal proportions and body part sizes. The second stage uses both the given joint correspondences and the local contour signal from image edges to optimize the remaining volumetric body parameters (limb cross-sections and their tapering parameters \mathbf{x}_d) while holding the other parameters fixed. Finally, we refine the full model (\mathbf{x}) using similar image information to the second stage. The covariance matrix corresponding to the final estimate is used to generate an initial set of hypotheses, which are propagated in time using the algorithm described in §4. While the process is heuristic, it gives a balance between stability and flexibility. In practice we find that enforcing the joint constraints, mirror information and prior bounds on the variation of body parameters gives far more stable and satisfactory results. However, with monocular images, the initialization always remains ambiguous and highly uncertain in some parameter space directions, especially under 3D-2D joint correspondence data. In our case, we employ a suitable coarse pose initialization and use the above process for fine refinement, but if available, one could fuse pose information from multiple images.

6 Experiments

For the experiments shown here we use an edge and intensity based cost function and a body model incorporating priors and constraints as explained in §3.1 and §3.2. We use Gaussian tails for CSS. A quantitative evaluation of different Gaussian scalings appears in table 1 on page 19.

To illustrate our method we show results for an 8 second arm tracking sequence and two full body ones (3.5 s and 4 s). All three sequences contain both self-occlusion and significant relative motion in depth. The first two (fig. 9) were shot at 25 frames (50 fields) per second against a cluttered, unevenly illuminated background. The third (fig. 11) is at 50 non-interlaced frames per second against a dark background, but involves a more complex model and motions. In our unoptimized implementation, a 270 Mhz SGI O2 required about 5 s per field to process the arm experiment and 180 s

per field for the full body ones, most of the time being spent in cost function evaluation. The figures show the current best candidate model overlaid on the original images. We also explore the characteristic failure modes of various tracker components, as follows. By a *Gaussian single mode tracker* we mean a single hypothesis tracker doing local continuous optimization based on Gaussian error distributions and without enforcing any physical constraints. A *Robust single mode tracker* improves this by using robust matching distributions. A *Robust single mode tracker with joint limits* also enforces physical constraints. For multimodal trackers, the sampling strategy can be either CONDENSATION-based or CSS-based, as introduced in previous sections.

Cluttered background sequences: These sequences explore 3D estimation behavior with respect to image assignment and depth ambiguities, for a bending rotating arm under an 8 d.o.f. model and a pivoting full-body motion under a 30 d.o.f. one. They have cluttered backgrounds, specular lighting and loose fitting clothing. In the arm sequence, the deformations of the arm muscles are significant and other imperfections in our arm model are also apparent.

The *Gaussian single mode tracker* manages to track 2D frontoparallel motions in moderate clutter, although it gradually slips out of registration when the arm passes the strong edges of the white pillar (0.5 s and 2.2 s for the arm sequence and 0.3 s for the human body sequence). Any significant motion in depth is untrackable.

The *robust single mode tracker* tracks frontoparallel motions reasonably well even in clutter, but quickly loses track during in-depth motions, which it tends to misinterpret as frontoparallel ones. In the arm tracking sequence, shoulder motion towards the camera is misinterpreted as frontoparallel elbow motion, and the error persists until the upper bound of the elbow joint is hit at 2.6 s and tracking fails. In the full body sequence, the pivoting of the torso is underestimated, being partly interpreted as quasi-frontoparallel motion of the left shoulder and elbow joints. Despite the presence of anatomical joint constraints, the fist eventually collapses into the body if non-self-intersection constraints are not present.

The *robust joint-limit-consistent CSS multi-mode tracker* tracks the motion of the entire arm and body sequence without failure. We retain just the 3 best modes for the arm sequence and the 7 best modes for the full human body se-

quence. As discussed in §4.2, multimodal behavior occurs mainly during significantly non-frontoparallel motions, between 2.2–4.0 s for the arm sequence, and over nearly the entire full body sequence (0.2–1.2 s). For the latter, the modes mainly reflect the ambiguity between true pivoting motion and its incorrect “frontoparallel explanation”.

We also compared our method with a 3D version of that of (Heap and Hogg, 1998; Cham and Rehg, 1999). These methods were developed for 2D tracking and we were interested in how well they would behave in the far less well controlled monocular 3D case. We used a parametric Gaussian mixture representation, local descent for mode refinement (as in Heap and Hogg (1998); Cham and Rehg (1999)) and a process model based on constant velocity plus dynamical noise sampling as in Cham and Rehg (1999) (section 3.2, page 3), on the cluttered full body tracking sequence. However, note that unlike the original methods, ours uses robust (rather than least squares) image matching and robust optimization by default, and also incorporates physical constraints and model priors. We used 10 modes to represent the distribution over our 30 d.o.f. 3D configurations, whereas Cham and Rehg (1999) used 10 for their 38 d.o.f. 2D SPM model. Our first set of experiments used a non-robust SSD image matching metric and a Levenberg-Marquardt routine for local sample optimization (as in Cham and Rehg (1999), except that we use analytical Jacobians). With this cost function, we find that outliers cause large fluctuations, bias and frequent convergence to physically invalid configurations. Registration is lost early in the turn (0.5 s), as soon as the motion becomes significantly non-frontoparallel. Our second experiments used our robust cost function and optimizer, but still with sampling as in Cham and Rehg (1999). The track survived further into the turn, but was lost at 0.7 s when the depth variation became larger. As expected, we find that a dynamical noise large enough to provide sufficiently deep sampling along uncertain in-depth directions produces much too deep sampling along well-controlled transversal ones, so that most of the samples are lost on uninformative high-cost configurations. Similar arguments apply to standard CONDENSATION, as can be seen in the monocular 3D experiments of Deutscher et al. (2000).

Black background sequence: In this experiment we focus on 3D errors, in particular depth ambiguities and the influence of physical constraints and parameter stabilization pri-

ors. We use an improved body model with 34 d.o.f. The four extra parameters control the left and right clavicle joints in the shoulder complex, which we find to be essential for following many arm motions. Snapshots from the full 4 s sequence are shown in fig. 11, and various failures modes in fig. 12.

The *Gaussian single mode tracker* manages to follow near-frontoparallel motions fairly reliably owing to the absence of clutter, but it eventually loses track after 0.5 s (fig. 12a–d). The *robust single mode tracker* tracks the non-frontoparallel motion somewhat longer (about 1 s), although it significantly misestimates the depth (fig. 12e,f — the right leg and shoulder are pushed much too far forward and the head is pushed forward to match subject contour, *c.f.* the “correct” pose in fig. 11). It eventually loses track during the turn. The *robust multi-mode tracker with joint-limits* is able to track quite well, but, as body non-self-intersection constraints are not enforced, the modes occasionally converge to physically infeasible configurations (fig. 12g) with terminal consequences for tracking. Finally, the *robust fully constrained multi-mode tracker* is able to deal with significantly more complex motions and tracks the full sequence without failure (fig. 11).

7 Sampling Distributions

We also ran some more quantitative experiments aimed at studying the behavior of the different sampling regimes, particularly the efficiency with which they locate minima or low-cost regions of parameter space. We are interested in how the sampling distribution, as characterized by the shape of its core and the width of its tails, impacts the search efficiency. For the study here we used the simple, but still highly multi-modal, 3D joint to image joint likelihood surface that we use for initializing our 34 d.o.f. articulated model. We only estimated joint parameters, not body dimensions. We ran experiments involving Covariance Scaled Sampling (CSS) and Spherical Sampling (SS) for Gaussian distributions with scalings 1,2 8. To allow a fair comparison, at each scale we kept the volume of the sphere (proportional to R^n) equal to the volume of the corresponding rescaled unit covariance CSS ellipsoid (proportional to $\lambda_1 \dots \lambda_n$, the product of eigenvalues). Also note that the final sampling distributions are not exactly Gaussian — in fact they are often noticeably multimodal — because our sampler preserves

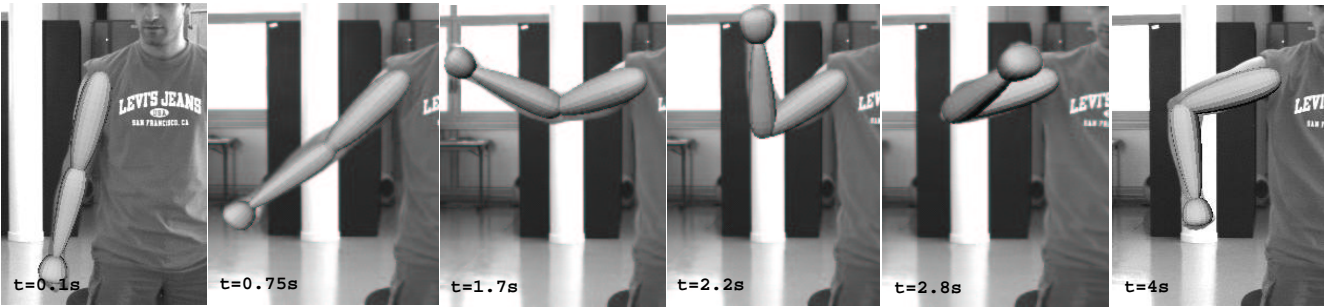


Figure 9: Arm tracking against a cluttered background.

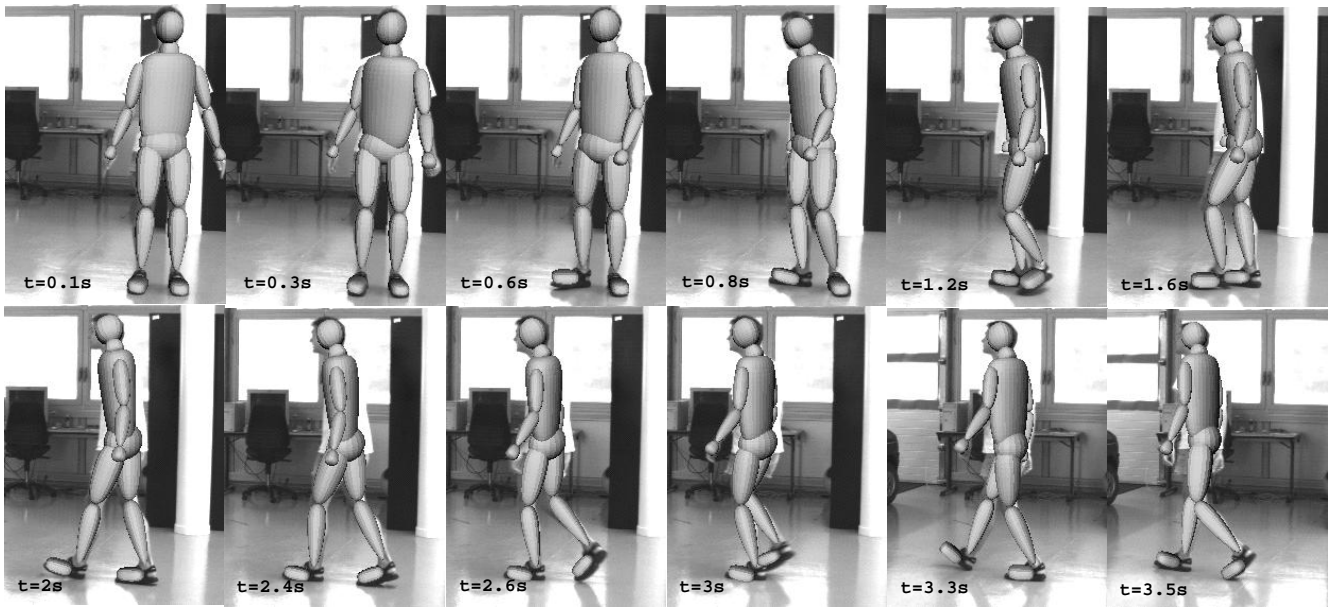


Figure 10: Human tracking against a cluttered background. See plate 14 on page 23 for details.

the physical constraints by projecting inadmissible samples back onto the constraint surface. Once made, the samples are locally optimized subject to the physical constraints using the method of §4.1. We report the number of minima found by each method, and the medians and standard deviations of their parameter space distances and cost differences, in table 1. Fig. 13 shows distributions of numbers of samples and minima versus parameter space distance, standard deviation and cost, for scaling 8. Note that CSS finds significantly more minima, and also places samples at positions of significantly lower cost, than SS. One can also see the large cost difference between optimized and unoptimized samples. SS appears to find minima of slightly lower median cost than CSS, but this is misleading. CSS still finds the few minima found by SS, but it also finds many other more distant ones,

which, being further away, tends to increase its median cost.

8 Approximation Accuracy

The tracking experiments in §6 illustrated the practical behavior and failure modes of some of the components of the CSS algorithm, and §7 presented a more quantitative evaluation. Now we turn to more technical points.

The CSS algorithm involves both local continuous optimization and somewhat more global covariance-scaled sampling. It therefore has a natural mechanism to trade-off speed and robustness. When tracking fails, both the number of modes used to represent the distribution and the number of samples produced in the sampling stage can be increased.

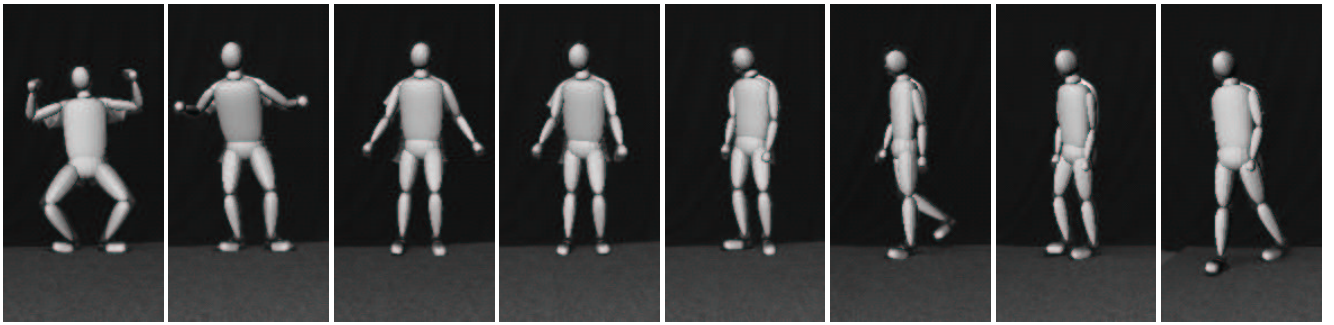


Figure 11: Human tracking under complex motion. See figure 15 on page 24 for details.

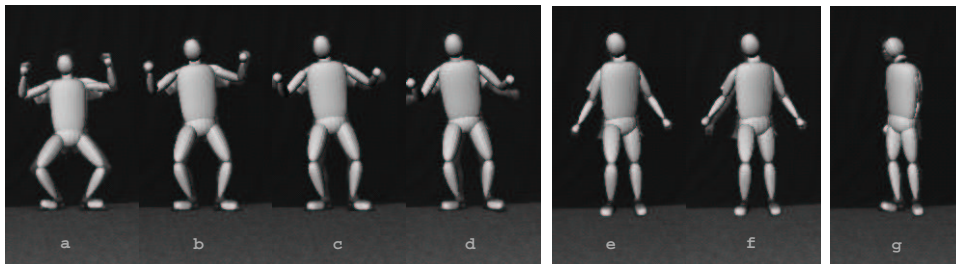


Figure 12: Failure modes of various components of the tracker (see text).

This increases the computational cost, but it may allow the tracker to follow more difficult portions of the image sequence. In principle, a sufficiently long run of any sampling method would visit every region of the parameter space, so that the basin of attraction of each mode was sampled and all minima were found. It has been argued that mixed continuous/discrete trackers (Heap and Hogg, 1998; Cham and Rehg, 1999) will 'diverge' if the visual information is ambiguous and converge to a 'best' mode when the target in the image is easily detectable. However, this kind of divergence is not that important here. We are working with likelihood surfaces that have multiple peaks with individual probabilities. Local optimization methods can converge to any of these peaks and sampling methods will eventually 'condense' near them if they use enough samples. Given the sampling/dynamics stage, both methods have a chance of jumping between peaks (*i.e.* escaping spurious ones), although this may be a very slow process. The method presented here is designed to address the problems of more efficient and systematic multi-modal exploration. Note also that CSS can be viewed as an importance sampling distribution and correction weighting for fair sample generation can be performed with respect to the true prior.

A second issue concerns the algorithm's efficiency versus its bias behavior. In tracking, assuming temporal coherence, one may want to confine the search to the neighborhood of the configurations propagated from the previous tracking step. This can be done implicitly by designing a likelihood surface that emphasizes local responses¹³, or by tuning the search process for locality, using short range dynamics. In either case, a global estimate of the posterior distribution is too expensive, both for sampling and optimization, while restricting attention to nearby states carries the risk of missing distant but significant peaks.

A third issue concerns the approximation accuracy of a Gaussian mixture for arbitrary multi-modal distributions. The mixture model is likely to be inaccurate away from the mode cores, and this may affect the accuracy of statistical calculations based on it. However for tracking and localization applications we are mainly interested in the modes themselves, not the low-probability regions in their remote tails. Sampling methods are non-parametric so in principle

¹³For example, an optical flow correspondence field, like that described in §3.1 but based on least squares brightness matching, can behave as a locality prior, forcing local image velocity explanations and pruning away remote, potentially 'objective' multi-modality.

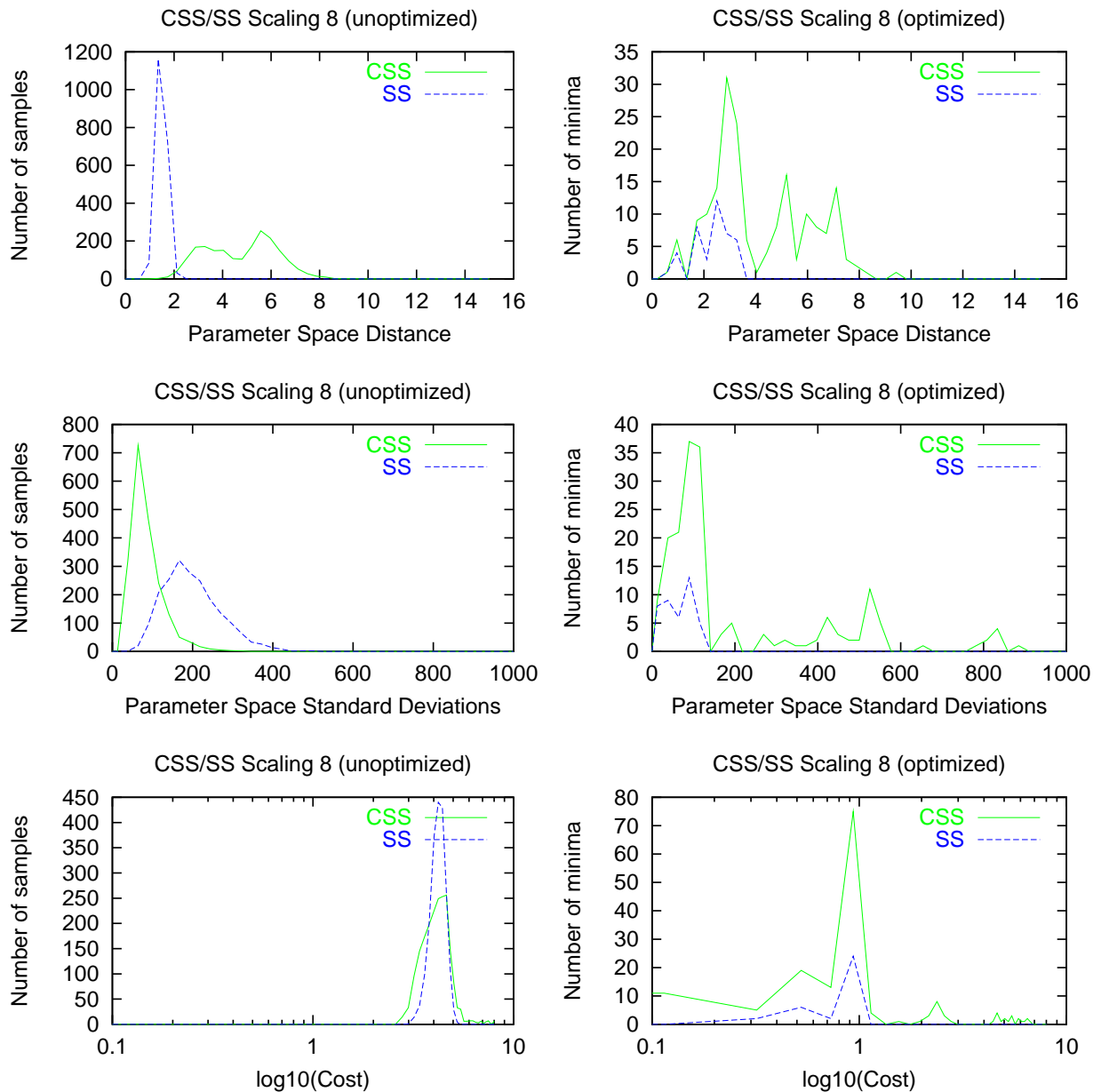


Figure 13: Optimized and unoptimized sample statistics for Spherical (SS) and Covariance Scaled (CSS) Sampling with scaling factor 8 and runs of 2000 samples. Note the significantly larger number of minima found by CSS than by SS, and also that the samples are placed at much lower cost.

they do not have this limitation, but in practice so few samples fall deep in the tails that noisiness of the estimates is a problem. Pure sample-based representations also provide little insight into the structure of the uncertainty and the degree of multi-modality of the likelihood surface. In any case, the issue of approximation accuracy in low-probability re-

gions is not a main concern here. Provided initial seeds are available in the individual mode's basins of attraction, sampling methods can generate fair samples from the modes and optimization methods can precisely identify their means and covariances by local descent. The two techniques can be used inter-changeably, depending on the application. It is the pro-

| METHOD | SCALE | NUMBER OF MINIMA | PARAMETER DISTANCE MEDIAN | | STANDARD DEVIATIONS MEDIAN | | COST MEDIAN | |
|--------|-------|------------------|---------------------------|---------|----------------------------|---------|-------------|---------|
| | | | UNOPT | OPT | UNOPT | OPT | UNOPT | OPT |
| CSS | 1 | 8 | 1.148 | 2.55242 | 10.9351 | 47.6042 | 116.951 | 8.49689 |
| CSS | 4 | 59 | 3.21239 | 2.9474 | 35.2918 | 55.3163 | 1995.12 | 6.98109 |
| CSS | 8 | 180 | 4.969 | 3.34661 | 75.1119 | 109.813 | 16200.8 | 7.09866 |
| SS | 1 | 0 | 0.199367 | - | 24.5274 | - | 273.509 | - |
| SS | 4 | 11 | 0.767306 | 2.04928 | 96.1519 | 39.0745 | 4291.12 | 6.28014 |
| SS | 8 | 42 | 1.47262 | 2.54884 | 188.157 | 56.8268 | 16856.1 | 6.96481 |

Table 1: Quantitative results for the distribution of minima found. Note again that CSS finds more minima and places raw samples at lower cost than SS.

cess of finding the initial seeds for each mode that represents the major difficulty for high-dimensional multi-modal distributions.

A fourth and important practical issue concerns the properties of the likelihood function. For many complex models a good likelihood is difficult to build, and the one used may be a poor reflection of the desired observation density. In these situations, the strength of true and spurious responses is similar and this may affect the performance of the tracking algorithm, irrespective how much computational power is used. In such contexts, it can be very difficult to identify the true tracked trajectory in a temporal flow of spurious responses. This is a particularly complex problem, since many likelihoods commonly used in vision degrade ungracefully under occlusion/disocclusion events and viewpoint change. At present, we do not have good mechanisms for detecting disocclusion events in complex backgrounds. The CSS algorithm has an elegant mechanism that accounts for high-uncertainty if particular degrees of freedom are not observed (like occluded limbs, *etc.*) and it will automatically sample broadly there. However, for sub-sequences with long occlusion events, it is still likely to attach occluded limbs to background clutter, rather than maintaining them as occluded. Global silhouettes, or a human contour detector (Papageorgiou and Poggio, 1999), or higher-order matching consistency (Sminchisescu, 2002a) may help here. As an indication of the potential benefits of this, we currently use foreground/background segmentation and the motion boundaries from the robust optical flow computation to weight the importance of contours, and this significantly improves the re-

sults in the sequences we have analyzed.

9 Conclusions and Future Work

We have presented a new method for monocular 3D human body tracking, based on optimizing a robust model-image matching cost metric combining robustly extracted edges, flow and motion boundaries, subject to 3D joint limits, non-self-intersection constraints, and model priors. Optimization is performed using Covariance Scaled Sampling, a novel high-dimensional search strategy based on sampling a hypothesis distribution followed by robust constraint-consistent local refinement to find a nearby cost minima. The hypothesis distribution is determined by propagating the posterior at the previous time step (represented as a Gaussian mixture defined by the observed cost minima and their Hessians / covariances) through the assumed dynamics (here trivial) to find the prior at the current timestep, then inflating the prior covariances and resampling to scatter samples more broadly. Our experiments on real sequences show that this is significantly more effective than using inflated dynamical noise estimates as in previous approaches, because it concentrates the samples on low-cost points, rather than points that are simply nearby irrespective of cost. In future work, it should also be possible to extend the benefits of CSS to CONDENSATION by using inflated (diluted weight) posteriors and dynamics for sample generation, then re-weighting the results. Our human tracking work will focus on incorporating better pose and motion priors as well as designing likelihoods that are

better adapted for human localization in the image.

Acknowledgments

This work was supported by an EIFFEL doctoral grant and the European Union under FET-Open project VIBES. We would like to thank Alexandru Telea for stimulating discussions and implementation assistance and Frédéric Martin for helping with the video capture and posing as a model.

References

- A. Barr. Global and Local Deformations of Solid Primitives. *Computer Graphics*, 18:21–30, 1984.
- C. Barron and I. Kakadiaris. Estimating Anthropometry and Pose from a Single Image. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 669–676, 2000.
- M. Black. *Robust Incremental Optical Flow*. PhD thesis, Yale University, 1992.
- M. Black and P. Anandan. The Robust Estimation of Multiple Motions: Parametric and Piecewise Smooth Flow Fields. *Computer Vision and Image Understanding*, 6(1):57–92, 1996.
- A. Blake, B. North, and M. Isard. Learning Multi-Class Dynamics. *Advances in Neural Information Processing Systems*, 11: 389–395, 1999.
- M. Brand. Shadow Puppetry. In *IEEE International Conference on Computer Vision*, pages 1237–1244, 1999.
- C. Bregler and J. Malik. Tracking People with Twists and Exponential Maps. In *IEEE International Conference on Computer Vision and Pattern Recognition*, 1998.
- T. Cham and J. Rehg. A Multiple Hypothesis Approach to Figure Tracking. In *IEEE International Conference on Computer Vision and Pattern Recognition*, volume 2, pages 239–245, 1999.
- K. Choo and D. Fleet. People Tracking Using Hybrid Monte Carlo Filtering. In *IEEE International Conference on Computer Vision*, 2001.
- Q. Delamarre and O. Faugeras. 3D Articulated Models and Multi-View Tracking with Silhouettes. In *IEEE International Conference on Computer Vision*, 1999.
- J. Deutscher, A. Blake, and I. Reid. Articulated Body Motion Capture by Annealed Particle Filtering. In *IEEE International Conference on Computer Vision and Pattern Recognition*, 2000.
- J. Deutscher, A. Davidson, and I. Reid. Articulated Partitioning of High Dimensional Search Spacs associated with Articulated Body Motion Capture. In *IEEE International Conference on Computer Vision and Pattern Recognition*, 2001.
- J. Deutscher, B. North, B. Bascle, and A. Blake. Tracking through Singularities and Discontinuities by Random Sampling. In *IEEE International Conference on Computer Vision*, pages 1144–1149, 1999.
- T. Drummond and R. Cipolla. Real-time Tracking of Highly Articulated Structures in the Presence of Noisy Measurements. In *IEEE International Conference on Computer Vision*, 2001.
- R. Fletcher. Practical Methods of Optimization. In *John Wiley*, 1987.
- D. Gavrilu and L. Davis. 3-D Model Based Tracking of Humans in Action: A Multiview Approach. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 73–80, 1996.
- L. Gonglaves, E. Bernardo, E. Ursella, and P. Perona. Monocular Tracking of the human arm in 3D. In *IEEE International Conference on Computer Vision*, pages 764–770, 1995.
- N. Gordon and D. Salmond. Bayesian State Estimation for Tracking and Guidance Using the Bootstrap Filter. *Journal of Guidance, Control and Dynamics*, 1995.
- N. Gordon, D. Salmond, and A. Smith. Novel Approach to Non-linear/Non-Gaussian State Estimation. *IEE Proc. F*, 1993.
- Group. *Specifications for a Standard Humanoid*. Hanim-Humanoid Animation Working Group, <http://www.hanim.org/Specifications/H-Anim1.1/>, 2002.
- T. Heap and D. Hogg. Wormholes in Shape Space: Tracking Through Discontinuities Changes in Shape. In *IEEE International Conference on Computer Vision*, pages 334–349, 1998.
- N. Howe, M. Leventon, and W. Freeman. Bayesian Reconstruction of 3D Human Motion from Single-Camera Video. *Neural Information Processing Systems*, 1999.
- M. Isard and A. Blake. CONDENSATION – Conditional Density Propagation for Visual Tracking. *International Journal of Computer Vision*, 1998.
- S. Ju, M. Black, and Y. Yacoob. Cardboard People: A Parameterized Model of Articulated Motion. In *2nd Int. Conf. on Automatic Face and Gesture Recognition*, pages 38–44, October 1996.

- I. Kakadiaris and D. Metaxas. Model-Based Estimation of 3D Human Motion with Occlusion Prediction Based on Active Multi-Viewpoint Selection. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 81–87, 1996.
- H. J. Lee and Z. Chen. Determination of 3D Human Body Postures from a Single View. *Computer Vision, Graphics and Image Processing*, 30:148–168, 1985.
- J. Liu. Metropolized Independent Sampling with Comparisons to Rejection Sampling and Importance Sampling. *Statistics and Computing*, 6, 1996.
- J. MacCormick and A. Blake. A Probabilistic Contour Discriminant for Object Localisation. In *IEEE International Conference on Computer Vision*, 1998.
- J. MacCormick and M. Isard. Partitioned Sampling, Articulated Objects, and Interface-Quality Hand Tracker. In *European Conference on Computer Vision*, volume 2, pages 3–19, 2000.
- R. Merwe, A. Doucet, N. Freitas, and E. Wan. The Unscented Particle Filter. Technical Report CUED/F-INFENG/TR 380, Cambridge University, Department of Engineering, May 2000.
- C. Papageorgiu and T. Poggio. Trainable Pedestrian Detection. In *International Conference on Image Processing*, 1999.
- M. Pitt and N. Shephard. Filtering via Simulation: Auxiliary Particle Filter. *Journal of the American Statistical Association*, 1997.
- R. Plankers and P. Fua. Articulated Soft Objects for Video-Based Body Modeling. In *IEEE International Conference on Computer Vision*, pages 394–401, 2001.
- J. Rehg and T. Kanade. Model-Based Tracking of Self Occluding Articulated Objects. In *IEEE International Conference on Computer Vision*, pages 612–617, 1995.
- K. Rohr. Towards Model-based Recognition of Human Movements in Image Sequences. *Computer Vision, Graphics and Image Processing*, 59(I):94–115, 1994.
- R. Rosales and S. Sclaroff. Inferring Body Pose without Tracking Body Parts. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 721–727, 2000.
- H. Sidenbladh and M. Black. Learning Image Statistics for Bayesian Tracking. In *IEEE International Conference on Computer Vision*, 2001.
- H. Sidenbladh, M. Black, and D. Fleet. Stochastic Tracking of 3D Human Figures Using 2D Image Motion. In *European Conference on Computer Vision*, 2000.
- H. Sidenbladh, M. Black, and L. Sigal. Implicit Probabilistic Models of Human Motion for Synthesis and Tracking. In *European Conference on Computer Vision*, 2002.
- C. Sminchisescu. Consistency and Coupling in Human Model Likelihoods. In *IEEE International Conference on Automatic Face and Gesture Recognition*, pages 27–32, Washington D.C., 2002a.
- C. Sminchisescu. *Estimation Algorithms for Ambiguous Visual Models—Three-Dimensional Human Modeling and Motion Reconstruction in Monocular Video Sequences*. PhD thesis, Institute National Polytechnique de Grenoble (INRIA), July 2002b.
- C. Sminchisescu, D. Metaxas, and S. Dickinson. Improving the Scope of Deformable Model Shape and Motion Estimation. In *IEEE International Conference on Computer Vision and Pattern Recognition*, volume 1, pages 485–492, Hawaii, 2001.
- C. Sminchisescu and A. Telea. Human Pose Estimation from Silhouettes. A Consistent Approach Using Distance Level Sets. In *WSCG International Conference for Computer Graphics, Visualization and Computer Vision*, Czech Republic, 2002.
- C. Sminchisescu and B. Triggs. A Robust Multiple Hypothesis Approach to Monocular Human Motion Tracking. Technical Report RR-4208, INRIA, 2001a.
- C. Sminchisescu and B. Triggs. Covariance-Scaled Sampling for Monocular 3D Body Tracking. In *IEEE International Conference on Computer Vision and Pattern Recognition*, volume 1, pages 447–454, Hawaii, 2001b.
- C. Sminchisescu and B. Triggs. Building Roadmaps of Local Minima of Visual Models. In *European Conference on Computer Vision*, volume 1, pages 566–582, Copenhagen, 2002a.
- C. Sminchisescu and B. Triggs. Hyperdynamics Importance Sampling. In *European Conference on Computer Vision*, volume 1, pages 769–783, Copenhagen, 2002b.
- C. Sminchisescu and B. Triggs. Kinematic Jump Processes for Monocular 3D Human Tracking. In *IEEE International Conference on Computer Vision and Pattern Recognition*, 2003.
- C. J. Taylor. Reconstruction of Articulated Objects from Point Correspondences in a Single Uncalibrated Image. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 677–684, 2000.
- B. Triggs, P. McLauchlan, R. Hartley, and A. Fitzgibbon. Bundle Adjustment - A Modern Synthesis. In Springer-Verlag, editor, *Vision Algorithms: Theory and Practice*, 2000.

- D. Vanderbilt and S. G. Louie. A Monte Carlo Simulated Annealing Approach over Continuous Variables. *J. Comp. Physics*, 56, 1984.
- S. Wachter and H. Nagel. Tracking Persons in Monocular Image Sequences. *Computer Vision and Image Understanding*, 74(3): 174–192, 1999.
- S. C. Zhu and D. Mumford. Learning Generic Prior Models for Visual Computation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(11), 1997.



Figure 14: Clutter human tracking sequence detailed results for the CSS algorithm in §6, on page 14.

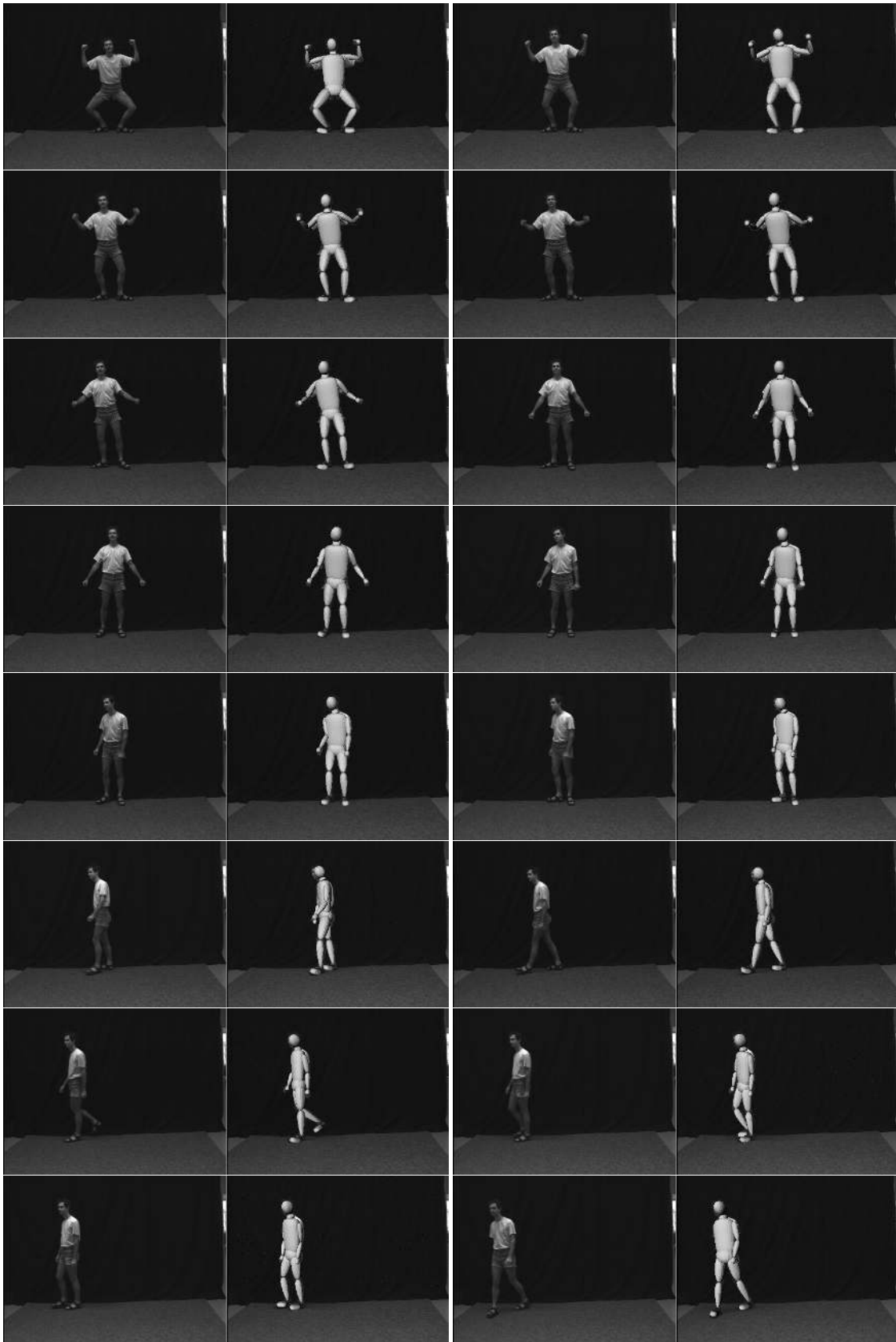


Figure 15: Complex motion tracking sequence detailed results for the CSS algorithm in §6 on page 14.