



Tracking Many Objects Using Subordinated CONDENSATION

David Tweed and Andrew Calway
Computer Science Department, Bristol University, UK
{tweed, andrew}@cs.bris.ac.uk

Abstract

We describe a novel extension to the CONDENSATION algorithm for tracking multiple objects of the same type. Previous extensions for multiple object tracking do not scale effectively to large numbers of objects. The new approach – subordinated CONDENSATION – deals effectively with arbitrary numbers of objects in an efficient manner, providing a robust means of tracking individual objects across heavily populated and cluttered scenes. The key innovation is the introduction of bindings (*subordination*) amongst particles which enables multiple occlusions to be handled in a natural way within the standard CONDENSATION framework. The effectiveness of the approach is demonstrated by tracking multiple animals of the same species in cluttered wildlife footage.

1 Introduction

Tracking objects over the course of an image sequence is one of the basic tasks in Computer Vision. The resulting trajectory can be either of interest in its own right or used as the foundation for a higher level analysis. Applications include surveillance and recognition systems [11, 6, 9] and advanced human-computer interaction[8]. However, designing robust tracking algorithms is difficult, requiring mechanisms to deal with issues such as weak distinguishing image features, background clutter, erratic and discontinuous motion, multiple and occluding objects, and many other problems.

Tracking algorithms use two sources of information to tackle these issues: a model of the dynamical behaviour of the object being tracked; and a model of its appearance within an image. The former can be obtained either from simplifying assumptions or by using *a priori* information obtained, e.g., that obtained from exemplars, whilst the latter will involve a combination of assumptions about the imaging geometry and the 3-D structure of the objects being tracked. The tracking process then involves finding a model configuration which balances consistency of the dynamical model against image support within the individual frames, leading to robust estimates of the object position.

The classic formulation of this approach is via the Kalman filter, although its assumption of a unimodal probability distribution and standard Gaussian interpretation means that its performance is limited in difficult tracking scenarios. A more robust approach and one that has found considerable success recently is particle filtering [3] and in particular the CONDENSATION algorithm [2]. These methods are based on creating an approximation to the full probability distribution of the object's configuration over all



number of samples derived from s_i is proportional to an *importance function* f but then using $\pi_i/f(s_i)$ as the weight of the resulting samples. However the greater the difference between π_i and $f(s_i)$ the lower the weight of the final particle even when the image observations give strong support, so this technique can only be used to adjust the sampling behaviour by a limited amount.

Probability decay. We measure the support for a particle s using a fixed observation model \mathcal{O} . This is inevitably a simplified model which, whilst giving high support for all objects, will have unpredictable variations between objects. Consider the idealised case of objects A and B which *consistently* give responses of p and q respectively. Then after t time-steps the ratio of the state density weight for B compared to that of A will be approximately $(q/p)^t$, i.e., it has decayed geometrically. This issue is difficult to deal with in a standard numerical way since by definition the degree to which an object matches the learned model is *a priori* unknown. Rather, we want to deal with it by periodically ‘normalizing’ the confidence of any hypotheses which are sufficiently supported.

We deal with these two issues in an integrated way based upon finding clusters within the state density and using those to construct a Voronoi tessellation [10] based upon these cluster centres. For these purposes only the 2-D positions of the reference point are used. Within each of these cells the distribution is describing primarily one object. Thus, we use the following to avoid sample impoverishment and probability decay using the following steps respectively:

- Every step, build an importance function which results in equal numbers of samples being taken in each Voronoi cell.
- Every N steps rescale the weights in each cell so that the peak weight is 1.

Empirically we have found that $N = 5$ works well. This intuitive scheme needs to be modified for the case of subordinated particles, and we do this in a simplistic way of treating each ‘depth-level’ independently, i.e., by computing a separate Voronoi diagram for each occupied depth level.

3 Models for animal tracking

In this section we briefly describe our representation of animals used within particles and then describe our implementation of observation exclusivity for this representation. Experiments using subordinated CONDENSATION for tracking birds in the sequence shown in Fig. 3 are then presented.

State representation. Our basis is a set of 2-D points linked together to form a ‘skeleton’ as shown for a bird in Fig. 4a–b. (Note the skeleton models image-plane appearance and *not necessarily* the anatomical skeleton.) It was manually chosen so that the configuration of the animal away from the nodes is well approximated by the links between them. For example, between the joint in the middle of the wing and the end of the wing the line between them follows the bird; this would not be the case if the middle of wing joint were removed and the end of the wing linked directly to the node on the side of the body.

The nodes on the skeleton were then divided by hand into suitable levels for partitioned sampling. Thus, Fig. 4b shows that the midpoint of the body is sampled over first,

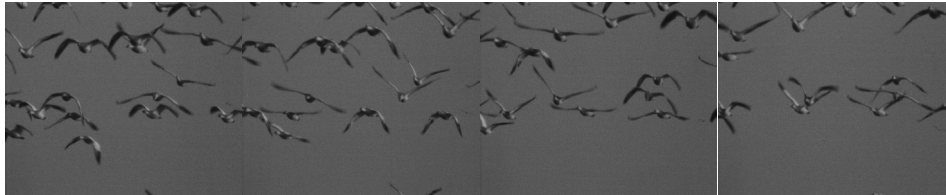


Figure 3: Frames 6, 18, 30 & 42 from a sequence of birds flying against an orange sky

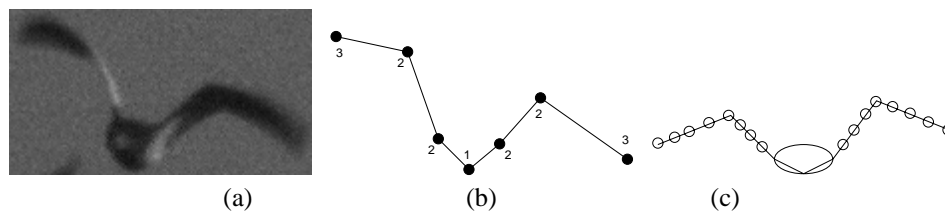


Figure 4: Representing animals: (a) original bird; (b) points representing (c) observation model consisting of ellipse for body and evenly spaced points along wings

then the positions of the upper halves of the wings, and finally the tips of the wings. This strategy of sampling over dependent points after sampling over those they depend upon is effective in reducing the total number of samples required, but does require manual analysis to determine it. This is not a drawback as our application is dealing with *potentially* previously unseen animals and so we require a minimal level of manual analysis to collect data for learning the motion model (described next).

Stochastic motion model. We present only an overview of the motion model for space reasons; see [12] for more details. We assume that (i) all the animals are moving in the same way and (ii) the motion can be split into a global position change and a periodic relative motion of the body parts. As we are tracking many animals over potentially hundreds of frames it is reasonable to use manual markup of points to be tracked on *one* animal over approximately two cycles of its motion to learn a global motion model which is then used for all the animals. We assume ‘constant velocity plus noise’ for the position of a reference point on the animal – the midpoint of the underside in the case of the bird in fig 4b – since the camera is often jerkily panning to keep the animals framed in the shot and a separate model to predict the positions of the points on the skeleton relative to the reference point. The limb model essentially treats the position in motion cycle as a *latent* variable which increments cyclically each frame and constructs a ‘body-configuration to body-configuration’ predictor at each point in the cycle.

Instance and model coordinates. There are two kinds of inaccuracy in the 2-D limb model described above, namely the inherent noise in the motion and the much larger variations due to differences in 3-D orientation and physical size between the animal the model was learned from and the animal being tracked. Rather than use a 3-D model (which would be difficult in practice given the relative scarcity of data and the fact the birds are small objects far from the camera) we attempt to take each animal’s coordinates as defined in a separate *instance coordinate system* and the single global motion model

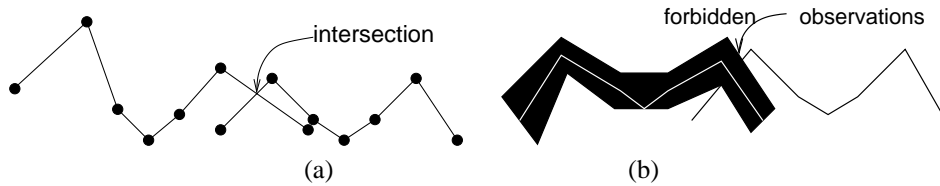


Figure 5: Schematic of methods used for (a) deciding if occlusion must be occurring and (b) deciding which image features cannot be used by the occluded object

as defined in a *model coordinate system*. Then if we attach to each animal a warping f from instance coordinates to model coordinates, we can perform prediction by mapping the state into model coordinates, applying the model and mapping the prediction back into instance coordinates. We use the simple warping

$$(x_i, y_i) \xleftarrow{f} (\bar{x}_i, \bar{y}_i) = ((x_i \cos \theta - y_i \sin \theta) / x_{scale}, (x_i \sin \theta + y_i \cos \theta) / y_{scale}) \quad (1)$$

with two scale factors as different sized animals tend to scale differently in vertical and non-vertical directions, and θ permitting physical rotation. These easily interpretable parameters can thus be restricted to lie in ‘plausible’ regions (e.g., birds can be allowed to fly at an angle of up to 45° whilst prohibiting unrealistic configurations such as flying upside down), and a least squares solution from $\{(x_i, y_i)\}$ and $\{(\bar{x}_i, \bar{y}_i)\}$ can be found analytically.

Measuring image support. For the initial experiments described in Section 4 we used one crudely hand-segmented frame to learn a bird/background likelihood model based upon pixel RGB values, modelling each class with a 6-component Gaussian mixture model [1]. The first level of the sampling should strongly localise the overall location of the bird using the large body. We do this by finding (via hill-climbing) the ellipse through the main point of the bird which maximises the sum of the log likelihoods of pixels inside the ellipse belonging to the bird and outside the ellipse belonging to the background – where the initial values are given by static ‘per-bird’ parameters – and using this as the observation weight. The two lower levels progressively localise the wings, so their support is measured by summing the log likelihood of a small number of points equally spaced along the wings being in the bird class. This is summarised in Fig. 4c.

3.1 Implementing observation exclusivity

The algorithm in Section 2.1 requires a method for determining if one particle could be occluding another and a method for preventing measurements of the support for an occluded object using features which belong to its occluder. Although these are related tasks, deciding possible occlusion is *much* more common and should therefore use a relatively inexpensive technique.

As we are assuming that the animals can be represented by a skeleton of points, we use whether the two skeletons intersect as a test for whether one of the animals must be occluding the other (as shown in Fig. 5a) since determining intersection can be performed quickly. This is a reasonable approximation for ‘thin-limbed’ animals but would be less appropriate for ‘bulky’ animals such as bears.

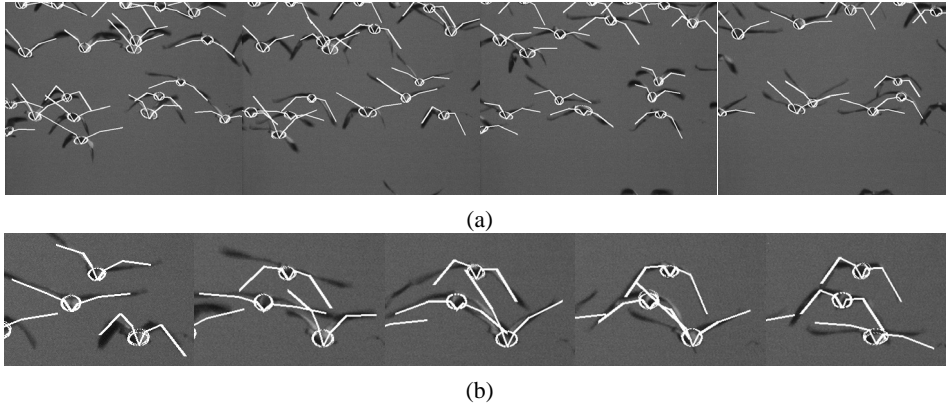


Figure 6: Tracking results: (a) frames 6, 18, 30 & 42 from the birds sequence; (b) subregion of frames 19, 21, 23, 25 & 27 showing an example of occlusion

To exclude foreground object features being used in background objects we simply prevent image pixels corresponding to marked ones in the bitmap being used when getting support for the background object (as shown in Fig. 5b). To approximate the set of pixels belonging to the foreground object we use a very simple technique: each point on the bird is linked to another, so that there is a well-defined normal vector to this edge at each point. The offset (in positive and negative directions) along the normal which maximises where the probability of belonging to the bird (measured using the same likelihood model as the basic observations) first drops below $1/2$ can be found. These points define a (non-convex) polygon which can be filled using standard graphics techniques [4].

4 Experiments

We demonstrate the performance of the algorithm on 50 frames from the sequence in figure 3 consisting of over 15 birds flying against an orange sky. As our elementary image feature is a bird/background likelihood model on pixel RGB values, the background is ‘cluttered’ in as much as many background pixels support belonging to the bird class. In addition several of the birds pass behind others during the sequence, testing the occlusion-resilience of the algorithm. The dynamical model was learned from two cycles of *one* manually marked up bird. The initial state density was generated by marking the approximate reference point of each bird in the initial frame (along with the body ellipse static parameters as in Section 3) and placing an equal number of samples in the vicinity random choices for both the position in the motion cycle and the parameters for the instance-to-model coordinate-mapping. The tracker was then run using 1200, 6000 and 6000 particles respectively on each of the three levels of the model (i.e., very roughly 400 particles per individual bird on the lower levels). Fig. 6a shows the mean configurations within the cluster at equally spaced frames throughout the sequence. Although not shown, the tracker quickly converges onto the correct position in the cycle and mapping parameters for each bird over the first five frames leading the result for frame 6. The results show the robust tracking of the birds in the sequence despite heavy occlusion between birds, an



detailed example of which is shown in Fig. 6b. (Note that the unmarked bird on the right of the final image in Fig. 6a moves into shot midway through the sequence, a situation the algorithm currently does not deal with.)

Note that whilst the wings sometimes drift from the correct position, this is essentially due to the somewhat simplistic measurement model being used, and in all cases the model locks onto the wings again later in the sequence.

5 Conclusions and future work

We have presented an extension to the CONDENSATION algorithm which can be used for efficiently tracking both individual and multiple objects, based upon forming *subordination* links between particle and enforcing an observation exclusion principle on linked particles. Tracking results on the bird sequence suggest its potential as a robust and efficient method for tracking multiple objects.

One unsatisfactory aspect is the use of a relatively strong observation model in the form of an object/background likelihood model learned from the sequence. In future work we intend to find ‘image features’ better suited to the issues in wildlife footage and are amenable to being learned automatically. We will also be investigating bootstrapping of the models required for tracking, possibly by using generic models (e.g., bird, biped, quadruped, etc) from a library, deforming them to improve tracking performance. Another improvement would be to use some feedback about how well the tracker is performing to determine the number of particles to be used in the next frame.

References

- [1] C M Bishop. *Neural Networks for Pattern Recognition*. OUP, 1995.
- [2] A Blake and M Isard. CONDENSATION – conditional density propagation for visual tracking. *Int J Computer Vision*, 29(1):5–28, 1998.
- [3] A Doucet, N de Freitas, and N Gordon, editors. *Sequential Monte Carlo Methods in Practice*. Springer-Verlag, 2001.
- [4] Foley, van Dam, Feiner, and Hughes. *Computer Graphics: Principles and Practice*. Addison-Wesley, 1990.
- [5] C Hue, J Le Cadre, and P Perez. Tracking multiple objects with particle filtering. Technical Report 1361, IRISA, October 2000.
- [6] M Isard and J MacCormick. BraMBLe: A Bayesian multiple-blob tracker. In *ICCV*, pages 34–41, 2001.
- [7] J MacCormick and A Blake. A probabilistic exclusion principle for tracking multiple objects. In *ICCV*, pages 572–578, 1999.
- [8] J MacCormick and M Isard. Partitioned sampling, articulated objects and interface quality hand tracking. In *ECCV*, pages 3–19, 2000.
- [9] D Ormoneit, H Sidenbladh, M J Black, and T Hastie. Learning and tracking cyclic human motion. In *NIPS*, pages 894–900, 2000.
- [10] R Sedgewick. *Algorithms*. Addison-Wesley, 1992.
- [11] H Tao, H S Sawhney, and R Kumar. A sampling algorithm for detecting and tracking multiple objects. In *Proc. Vision Algorithms (associated ICCV)*, 1999.
- [12] D Tweed and A Calway. Tracking multiple animals in wildlife footage. In *ICPR*, 2002.