

Video Pop-up: Monocular 3D Reconstruction of Dynamic Scenes*

Chris Russell**, Rui Yu**, and Lourdes Agapito

University College London

<http://www0.cs.ucl.ac.uk/staff/lagapito/research/youtube3d>

Abstract. Consider a video sequence captured by a single camera observing a complex dynamic scene containing an unknown mixture of multiple moving and possibly deforming objects. In this paper we propose an unsupervised approach to the challenging problem of simultaneously segmenting the scene into its constituent objects and reconstructing a 3D model of the scene. The strength of our approach comes from the ability to deal with real-world dynamic scenes and to handle seamlessly different types of motion: rigid, articulated and non-rigid. We formulate the problem as hierarchical graph-cut based segmentation where we decompose the whole scene into background and foreground objects and model the complex motion of non-rigid or articulated objects as a set of overlapping rigid parts. We evaluate the motion segmentation functionality of our approach on the Berkeley Motion Segmentation Dataset. In addition, to validate the capability of our approach to deal with real-world scenes we provide 3D reconstructions of some challenging videos from the *YouTube-Objects* dataset.

1 Introduction

With the emergence of video cameras on phones and laptops and the rise of video libraries (e. g. *YouTube Action*, *YouTube Objects*) the use of 3D information for recognition tasks has experienced a resurgence. While structure from motion (sfM) techniques exist that can reliably reconstruct a static scene, most scenes of interest contain multiple moving objects or even articulated or non-rigid objects. Motion segmentation and non-rigid scene reconstruction from monocular video have become more important than ever. This paper proposes a refocusing of 3D reconstruction towards reconstructing videos of dynamic scenes.

Multibody sfM and non-rigid structure from motion (NRSfM) have addressed some of the limitations of sfM and have seen sustained progress in dealing with dynamic scenes [21,25] or creating vivid life-like reconstructions of deformable objects [12]. However, they remain far behind their rigid counterparts. Multibody sfM approaches can segment the scene into multiple rigidly moving objects,

* This research was funded by the European Research Council under the ERC Starting Grant agreement 204871-HUMANIS.

** The first two authors assert equal contribution and joint first authorship.

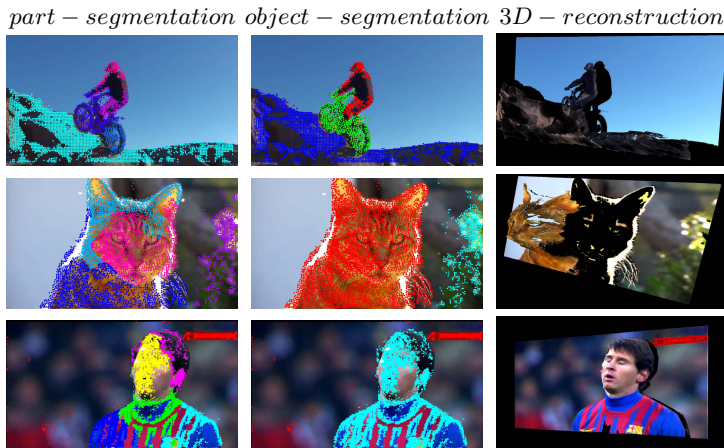


Fig. 1. Segmentation and 3D reconstruction results of two dynamic sequences of the *Youtube-Objects* Dataset [23] and a football sequence downloaded from YouTube. **Left:** segmentation into *parts* (rigid models). **Centre:** segmentation into *objects*. **Right:** densified 3D video pop-up from a novel viewpoint. The *motor-bike* sequence, acquired with a moving camera, shows articulated motion. The *cat* sequence is a non-rigid sequence occluding a static background. Bottom row shows a reconstruction of football footage. For videos see the project website <http://www0.cs.ucl.ac.uk/staff/lagapito/research/youtube3d>.

however they cannot deal simultaneously with the presence of deformable or articulated objects in the scene. Although NRSfM algorithms can reconstruct a single pre-segmented deformable surface moving in front of a camera [12,33], they require manual segmentation of the scene into background and foreground.

Piecewise approaches to non-rigid and articulated reconstruction have been successfully applied to explain the complex motion of 2D tracks on a single non-rigid surface or an articulated object as a network of *overlapping* parts [9,26,33]. However, if applied to an entire scene with foreground/background objects occluding one another, depth boundaries between objects would not be respected and neighbouring models in the image would be forced to overlap irrespective of whether or not they belong to the same physical object.

Contributions: The main contribution of this paper is to offer a solution to the problem of *scene reconstruction* for real-world dynamic monocular videos that deals seamlessly with the presence of non-rigid, articulated or pure rigid motion. In an entirely unsupervised approach, we reorganise/segment the scene into a constellation of object parts, recognise which parts are likely to constitute objects, join them together, and reconstruct the scene. We offer solutions to some of the problems of previous approaches to dynamic scene reconstruction: (i) Our approach is able to adapt the topology of the neighbourhood graph by breaking edges where necessary to preserve boundaries between objects. In this way our approach can deal with an entire scene where objects might occlude one another

and not just pre-segmented objects; *(ii)* Our work results in a hierarchical approach to dynamic scene analysis. At the higher level of the hierarchy the scene is explained as a set of *objects* that are detached from the background and from each other. At the lower level of the hierarchy, each *object* can be explained as a set of overlapping *parts* that can model more complex motion.

2 Related Work

Most works in dynamic scene reconstruction [7,10,27,21] follow a pipeline approach where first feature point tracks or dense optical flow is estimated, followed by a motion segmentation step to cluster trajectories into different independent motions before 3D reconstruction is applied independently to each of the objects. The first approach to **multibody reconstruction** [7] extended the classic affine factorisation algorithm for static scenes [30] to the case of multiple independently moving rigid objects. While the original approach [7] was unable to deal with dependencies in the motions it was later extended to deal with degenerate [38] and articulated motions [32,36]. More recent approaches to multibody s f M such as Ozden *et al.* [21] are able to perform simultaneous tracking, segmentation and reconstruction of a few feature points on realistic sequences. Roussos *et al.* [25] proposed a dense approach to multibody s f M in which depth values are estimated for every pixel in the image. However, none of these approaches can deal with non-rigidity or articulation in each of the objects which are assumed to be rigid.

Providing robust solutions to **video and motion segmentation** is a fundamental problem in computer vision and often a preliminary step towards 3D reconstruction. A wealth of motion segmentation algorithms for multi-rigid scenes have been proposed including algebraic frameworks such as GPCA [34] and methods that can deal with noise and outliers [24]. Motion segmentation has also been cast as a motion subspace clustering problem, first applied to the affine camera case [8,15] and later extended to the case of perspective scenes [18]. Approaches such as Brox and Malik's [6] exploit the consistency of point trajectories over time and can deal with non-rigid motion. On the other hand, superpixel [11] and supervoxel [35] methods for video segmentation can produce high quality video over-segmentations that respect object boundaries, are temporally consistent and are aligned with objects. However, since their aim is to segment non-rigid and articulated objects as a single segment, they are not appropriate for piecewise 3D reconstruction.

Non-rigid structure from motion (NRS f M) approaches reconstruct 3D models of non-rigid objects from monocular video, typically by fitting a global low-rank shape model [31,22] to 2D tracks. Piecewise reconstruction has also been successfully applied to NRS f M [26,33] and articulated reconstruction [9] by fitting local models. However all existing methods can only reconstruct a single presegmented object and can not resolve entire scenes.

Our approach is most closely related to the paradigm of *multiple model fitting* where tracks, that might contain outliers, belong to an unknown number of models. The assignment of tracks to models and the estimation of model parameters is then optimised jointly [14,26] to minimise a geometric cost subject to the

constraint that neighbouring tracks must belong to the same model. The cost also incorporates a minimum description length (MDL) cost that prefers sparse solutions. The cost function is optimised by alternating between a discrete graph-cuts algorithm to solve the labelling problem and continuous optimisation to update the model parameters. This approach has previously been applied to computer vision problems such as stereo [2]; motion segmentation [14]; 3D reconstruction of non-rigid [26] and articulated objects [9]; and multi-body reconstruction [25].

Our approach departs from previous work in geometric multiple model fitting in multiple ways: *(i)* Our model offers segmentation at two granularities: object-level and part-level. At the object-level, we segment the scene into a small number of disjoint *objects*. At the part-level, objects are further divided into a set of overlapping parts; *(ii)* Our model uses a combination of appearance and geometry cues for segmentation which encourages salient foreground objects to be separated accurately from the background even when the motion is not distinctive enough; *(iii)* Our geometric cost uses a perspective camera model and is able to deal with perspective effects and incomplete tracks.

3 Simultaneous Segmentation and Reconstruction

We consider a monocular video sequence, possibly downloaded from the web, captured by a single camera observing a complex dynamic scene that contains an unknown mixture of multiple moving and possibly deforming objects. First, we extract a set $\mathcal{T} = [1, \dots, T]$ of feature point tracks using Sundaram *et al.*'s publicly available code [29]. Although the tracker aims to provide long-term video correspondences, the length of tracks is variable and not all points tracked are visible in all the frames. We make no assumptions about the number of objects in the scene or their motions which could be rigid, articulated or non-rigid. Our goal is to estimate the 3D coordinates for all feature points in every frame.

3.1 Piecewise Reconstruction with Overlapping Models

The works [26,9] proposed a novel piecewise approach to the problem of 3D reconstruction of non-rigid objects. Rather than attempting to reconstruct objects by fitting a global low-rank shape model [31,22] that is sufficiently expressive to capture deformations, but sufficiently low-rank to discourage overfitting, they automatically segmented the object to be reconstructed into a set of parts, each of which could be expressed by a simple model – either local rigid reconstructions [9] or local quadratic deformations [26]. By forcing these parts to overlap, and to agree about the reconstruction of the region of overlap, per part depth/scale and sign-flip ambiguities can be resolved. Figure 2 shows an illustration of the segmentation of an articulated object into overlapping rigid parts.

The problem was formulated as a labelling one where the assignment of tracks to models and the fitting of models to tracks were jointly optimised to minimise a geometric fitting cost subject to the spatial constraint that neighbouring tracks should also belong to the same model.

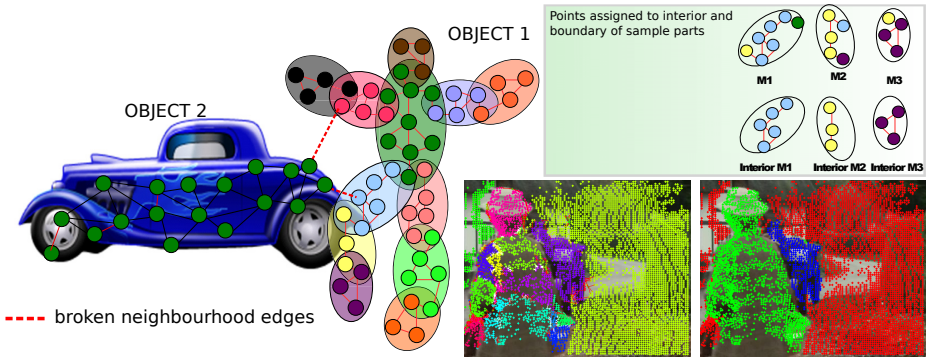


Fig. 2. **Left:** Conceptual illustration of our approach to 3D reconstruction of complex dynamic scenes. The image shows a person occluding a car. In the original neighbourhood graph, some point-tracks on the car are path connected with tracks on the person. Our approach reasons about object boundaries by adapting the neighbourhood, and breaking edges where necessary to detach parts from other occluding objects. **Top Right:** Illustration of the concept of overlapping models and *interior points* [26]. A tracked point belongs to the interior of a model (points with the same colour) if all its neighbours also belong to that model (though not necessarily as interior points). **Bottom right:** Real world example of segmentation into parts (left) and two objects and background (right).

Assignment of Point Tracks to Models. Let \mathcal{T} refer to a set of point tracks and \mathcal{M} a set of models. We use the notation $\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$ to refer to a labelling, where \mathbf{x}_i is the *set* of models assigned to track i . Assuming a known topology, or graph, which connects tracks together in a neighbourhood structure \mathcal{N} , the following objective was proposed by [9,26]

$$C(\mathbf{x}) = \sum_{i \in \mathcal{T}} \sum_{m \in \mathbf{x}_i} U_i(m) + \text{MDL}(\mathbf{x}) \tag{1}$$

where tracks are allowed to belong to multiple models in \mathcal{M} . The unary term $U_i(m)$ is the cost of assigning track i to model m and the term $\text{MDL}(\mathbf{x})$ is a label cost that encourages sparse solutions. In [9] local rigid models were used where each model was parameterised with the rotation and translation associated with a rigid motion and the unary cost $U_i(m)$ was defined as the image reprojection error under orthographic projection for that point given the model parameters. The optimisation of (1) was subject to the constraint that each track must be an interior point of some model i.e. that for every track there is a model such that that track and all its neighbours belong to that model (Figure 2 illustrates the concept of *interior point*), or more formally $\forall i \exists \alpha : \alpha = I_i$ and

$$I_i = \alpha \rightarrow \forall j \in N_i, \alpha \in \mathbf{x}_j \tag{2}$$

where $\mathbf{I} = \{I_1, I_2 \dots I_T\}$ refers to the assignment of each track i to the interior of one model I_i and N_i is the neighbourhood of track i . Russell *et al.* [26] showed how this problem could be formulated as a labelling problem over the assignment

of tracks to the interior of models and efficiently solved using a novel variant of α -expansion. Starting from an excess of models the optimisation followed a hill climbing approach that alternates between assigning tracks to models, and refitting the models to minimise the geometric error (image reprojection error).

3.2 Obstacles to Reconstruction in the Wild

Although these multiple model fitting approaches based on overlapping models do provide a robust approach to non-rigid [26] and articulated [9] reconstruction, they have shortcomings. First, they cannot deal with whole scenes in which the neighbourhood graph maintains connections between tracks of different objects (see Figure 2) – the constraints (2) combined with a bad neighbourhood structure can force parts to straddle multiple objects, leading to an error that can not be recovered from. Secondly, the unary terms of [9,26] minimise a geometric cost based on multiview affine factorisation. Therefore, they have difficulty dealing with incomplete tracks. In real-world videos, tracks are likely not to persist for a large number of frames. Finally, a further limitation of the above approaches comes from the fact that only motion cues are used for the segmentation. Combining motion and appearance cues is useful to encourage object boundaries to be respected. Besides, these cues complement each other particularly if there are frames in the sequence with small motion.

The main contribution of our work is to offer solutions to these three limitations: (i) Our approach adapts the topology of the neighbourhood graph by breaking edges where necessary to preserve boundaries between objects. This allows our approach to deal with complete video footage where objects might occlude one another and not just singular pre-segmented objects. (ii) Our geometric unary cost is based on frame-to-frame fundamental matrices, an approach that leads itself naturally to handling incomplete tracks. (iii) Our data term combines geometric and appearance costs. We use the saliency score provided by [28] to encourage parts of similar saliency to belong to the same object.

4 Scene Reconstruction with an Adaptive Neighbourhood

We propose a novel cost that allows us to modify the topology of the original neighbourhood by deleting edges between point tracks that belong to different physical objects, and should not overlap. Our new cost has four terms

$$C(\mathbf{x}) = E_{data} + E_{edge_break} + E_{sparse} + E_{mdl} \quad (3)$$

$$= \sum_{i \in \mathcal{T}} \sum_{m \in \mathbf{x}_i} U_i(m) + \sum_{i \in \mathcal{T}} \sum_{j \in N_i} d_{i,j} \Delta(j \notin N'_i) \quad (4)$$

$$+ \sum_{m \neq n \in \mathcal{M}} \Delta(\exists i : I_i = m, n \in \mathbf{x}_i) + \text{MDL}(\mathbf{x}) \quad (5)$$

where as before \mathbf{x}_i is the *set* of models that point i belongs to; $\Delta(\cdot)$ is the indicator function, taking value 1 if the statement is true and 0 otherwise; and N'_i the modified

neighbourhood of track i . This optimisation is subject to the constraints that every track i belongs to the interior of one model I_i , or more formally

$$\forall i, I_i = \alpha \rightarrow \forall j \in N'_i, \alpha \in \mathbf{x}_j \quad (6)$$

We now describe in detail each term of our cost function.

4.1 Unary Costs (E_{data})

Our unary term is the sum of two costs i.e. $U_i(m) = G_i(m) + P_i(m)$, that encourage tracks that both move consistently as a rigid object and have similar saliency scores, to belong to the same model. The geometric term $G_i(m)$ evaluates the cost of assigning track i to a rigid model m as the deviation from the epipolar geometry across all pairs of consecutive frames. The second term $P_i(m)$ computes a saliency score for each pixel in every frame and encourages tracks with similar saliency scores, to belong to the same model.

Rigidity Term G_i . Given a set of point tracks assigned to the same rigid part, we parameterise the rigid model m associated with them as a set of $F - 1$ fundamental matrices $\mathbf{F}_m = \{\mathbf{F}_m^{1,2}, \dots, \mathbf{F}_m^{f,f+1}, \dots, \mathbf{F}_m^{F-1,F}\}$ for every pair of consecutive frames in the sequence $f = \{1, \dots, F - 1\}$. The cost of associating track i to a specific rigid model m is the Sampson error [13] added over all pairs of fundamental matrices

$$G_i(m) = \sum_{f < F} \gamma^{-1} (u_i^{f+1T} \mathbf{F}_m^{f,f+1} u_i^f)^2 \quad (7)$$

where u_i^f encodes the homogeneous image coordinates of track i in frame f and u_i^{f+1} its corresponding position in frame $f + 1$ and γ is the Sampson weight [13]. This cost is summed over all frames in which the track is visible. To estimate the fundamental matrices, we use the eight-point algorithm embedded in a Ransac scheme followed by non-linear refinement of (7). This fitting cost has several clear advantages over the affine factorisation cost used by [9]. First, it allows to model perspective effects which are often present in unconstrained videos and to perform perspective reconstruction given an estimate of the camera calibration matrix. Second, it behaves better in the presence of missing data or short tracks, as it computes frame-to-frame geometric costs only for the frames where the track is visible rather than the multiframe factorisation cost of [9].

Saliency Term. The work [28] provides a fully unsupervised method for object detection in an image I , using a novel saliency map S_I . While [28] made use of both the statistics taken from a large corpus of unlabelled images, and from the image itself, we only make use of the statistics of the single image (this measure is termed *within image saliency* in [28]). We compute saliency maps S_{I_f} for each

frame f in the video sequence and define the saliency cost $P_i(m)$ of point i belonging to model m as the distance from the mean saliency of model m

$$P_i(m) = \lambda_s \sum_{f \leq F} (S_{I_f}(i) - \bar{S}_m)^2 \quad (8)$$

where \bar{S}_m is the mean saliency of all tracks that currently belong to model m , $S_{I_f}(i)$ is the saliency score of point i in frame f and λ_s a weight on the importance of this term.

4.2 Topologically Adaptive Neighbourhood (E_{edge_break})

The cost (1) proposed in [26] was internally represented as a local MDL prior defined over the set of interior labels present in a local neighbourhood, and took the cost

$$\sum_{i \in \mathcal{T}} \sum_{m: \exists j \in N_i \cap m = I_j} U_i(m) \quad (9)$$

As discussed, in order to separate connected objects from one another, we wish to discard edges from the neighbourhood N_i with a per edge cost of $d_{i,j}$. As such, the new cost will be of the form

$$\sum_{i \in \mathcal{T}} \sum_{m: \exists j \in N_i \cap m = I_j} \min \left(\sum_{j: I_j = m} d_{i,j}, U_i(m) \right) \quad (10)$$

Here, the weights $d_{i,j}$ are found by passing the distance between points i and j in the image and velocity spaces through a sigmoid function.

4.3 Overlap Sparsity Term (E_{sparse})

By itself, discarding edges from the neighbourhood graph improves the quality of the parts found, and allows more objects to be found. However, it does not correctly separate objects from the background. In almost all sequences, we find that one or two ambiguous tracks exist that could be easily explained as either object or background parts. These ambiguous tracks act as junctions, or regions of overlap between foreground and background objects, connecting the two and making it impossible to distinguish between foreground and background.

To eliminate this leaking, we introduce a novel sparsity term that penalises the total number of models that overlap and encourages regions with limited overlap to disconnect. We formulate this penalty as a count of the number of pairs of models (m, n) such that there exists a track belonging to the interior of model m and also to model n , i.e.

$$\sum_{m \neq n \in \mathcal{M}} \Delta(\exists i : I_i = m, n \in \mathbf{x}_i) \quad (11)$$

As this cost does not depend on the number of tracks in the region of overlap, it dominates in small regions of overlap or where the cost of discarding edges is small, and is ignored elsewhere.

5 Efficient Optimisation

As with other multiple model fitting approaches [9,14,26], we initialise with an excess of models which are generated by sampling randomly groups of ten feature tracks and computing the frame-to-frame fundamental matrices using the eight-point algorithm [13]. We then optimise the cost (3) using a hill-climbing approach alternating between: (i) fixing the parameters \mathbf{F}_m and optimising the labelling that assigns tracks to a set of parts (models) $\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$ and (ii) fixing the labelling and optimising \mathbf{F}_m for all models.

Alpha expansion [5] finds a local optimum of a difficult to optimise cost function by iteratively moving from a current labelling to the lowest-cost solution obtained by relabelling some of the variables as α . Finding an optimal move is formulated as a pseudo Boolean optimisation [3] and solved using graph-cuts [4]. We follow [26] in considering expansion moves over the interior of labels. We use $A \in 2^T$ to refer the found expansion move, with A_i taking value 1 if variable I_i transitions to label α in the move, and 0 otherwise. Unlike [26] we will need to explicitly keep track of whether or not tracks belong to models at all (either as interior or boundary tracks) and for a particular expansion move on label α this will be done by means of binary variables $M_i^\alpha = 1$ if $\alpha \in \mathbf{x}_i$ and a complementary set of variables M_i^β , such that $\beta \neq \alpha$ and $M_i^\beta = 0$ if $\beta \in \mathbf{x}_i$.

Optimisation of the costs E_{data} and E_{mdl} can be done using the techniques of [26]. We now deal with the modifications to the optimisation required by the terms E_{break} and E_{sparse} . Although exact optimisation of either of these costs is relatively straightforward, optimising both together is challenging, and we make use of the convex-concave procedure (CCP) [37,20], and find an optimisable cost that is tight at the current location, but an over-estimate elsewhere.

5.1 Exactly Optimising E_{break}

We can rewrite cost (10) in terms of the auxiliary variables

$$\sum_{i \in \mathcal{T}} \sum_{\substack{\beta \in \mathcal{M} \\ \beta \neq \alpha}} \sum_{j: I_j = \beta} \min_{M_i^\beta} (A_j d_{i,j} (1 - M_i^\beta) + U_i(\beta) (1 - M_i^\beta)) \quad (12)$$

$$+ \sum_{i \in \mathcal{T}} \sum_{j: I_j = \alpha} \min_{M_i^\alpha} ((1 - A_j) d_{i,j} M_i^\alpha + U_i(\alpha) M_i^\alpha) \quad (13)$$

This change can be seen as a robustification of the local co-occurrence potentials of [26] analogous to the robust Pn model [16]. As with the Pn potentials, it can be formulated as a graph-cut problem simply by adjusting the edge weights used as shown in Figure 3, left and centre left.

5.2 Approximately Minimising E_{sparse}

For the following section it is more convenient to use sets to describe which points belong to which models. We use \mathbf{M}^β to refer to the set of points belonging to

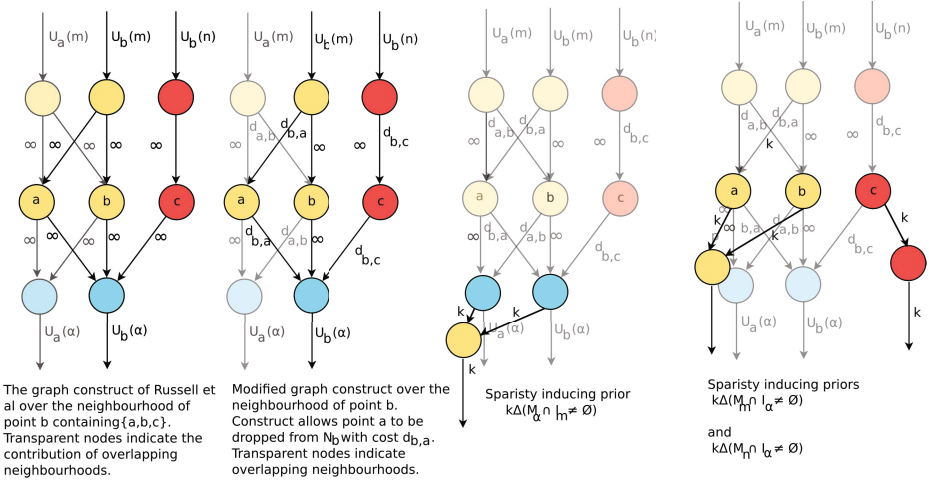


Fig. 3. Graph constructs. Of the main graph construct repeated in all subfigures, the top row contains all auxiliary variables indicating if tracks, a, b or c belong to model m or n . The middle row contains the standard expansion variables which govern whether or not a variable transitions to the interior of model α , while the bottom row indicates if a variable belongs to model α .

model β , \mathbf{I}^β for the interior of model β , $\mathbf{M}_{\text{last}}^\beta$ for the region (fixed throughout the move) that was assigned to model β by the previous move, and $\mathbf{I}_{\text{last}}^\beta$ for points previously belonging to the interior of model β . Performing an expansion move on label α , we have three cases to consider:

1. The cost is a direct function of the interior labels \mathbf{I}^α
2. The cost depends on tracks in the boundary of points belonging α : $\mathbf{M}^\alpha \setminus \mathbf{I}^\alpha$
3. The cost is not a function of α and depends on: $\mathbf{I}^\beta \cap \mathbf{M}^\alpha$, where $\beta, \gamma \neq \alpha$.

For an expansion move on label α , \mathbf{I}^α is monotone increasing while \mathbf{I}^β is monotone decreasing. If one of either the sparsity costs, or the edge breaking of the previous subsection was not used, the labelling of \mathbf{M}^α and \mathbf{M}^β would also be guaranteed to be monotone increasing/decreasing, but together the situation is more complex. In the following discussion, we artificially constrain the set of possible moves of \mathbf{M}^β to be monotone decreasing, and allow \mathbf{M}^α to change arbitrarily. Let us deal with these costs by turn:

Interior of α cost: We consider the localised MDL costs

$$\Delta(\mathbf{M}_{\text{last}}^\beta \cap \mathbf{I}^\alpha \neq \emptyset) + \Delta(\mathbf{I}^\beta \neq \emptyset) - 1. \tag{14}$$

This cost is 1 if \mathbf{I}^α expands into $\mathbf{M}_{\text{last}}^\beta$ without completely removing model β (which can only be done by making sure no tracks belong to the interior of model β) and 0 otherwise. Clearly this is an over-estimate as the true \mathbf{M}^β in

the set of all moves considered is always smaller than $\mathbf{M}_{\text{last}}^\beta$, and tight at the current location. As this cost is simply two MDL costs defined over subregions of the graph, it can be optimised using the techniques of [17]. As these move costs satisfy the CCP criteria, they reduce the original cost function.

Boundary of α cost: A similar argument can be made for the above cost. Instead of directly optimising it, we solve the over-approximation

$$\Delta(\mathbf{M}^\alpha \cap \mathbf{I}_{\text{last}}^\beta \neq \emptyset) + \Delta(\mathbf{I}^\beta \neq \emptyset) - 1. \quad (15)$$

This can be formulated as a local MDL prior over the auxiliary variable of the previous section and an MDL cost over label β . What is more interesting is the quality of the approximation of these terms. If we assume that all tracks belonging to the same model at every iteration of graph-cuts are path-connected¹, and if no edges are discarded² then the over-estimate is tight. There are 3 straightforward cases to consider:

1. No variables in $\mathbf{M}_{\text{last}}^\beta$ are in \mathbf{I}^α . Here the cost is trivially correct.
2. No variable belongs to \mathbf{M}^β . Again, the cost is trivially correct.
3. At least one variable i belongs to \mathbf{I}^α and $\mathbf{M}_{\text{last}}^\beta$, and one variable j belongs to \mathbf{I}^β . As i and j are path-connected in $\mathbf{M}_{\text{last}}^\beta$, and for any possible move variables in $\mathbf{M}_{\text{last}}^\beta$ must either stay as variable β or move to label α , we have a chain of variables $\{i, k, l, \dots, j\}$ belonging to either \mathbf{I}^α or \mathbf{I}^β such that $i \in N_k, k \in N_l, l \in \dots, N_j$. As $i \in \mathbf{I}^\alpha, j \in \mathbf{I}^\beta$, there must be at least one pair where $k \in N_l, k \in \mathbf{I}^\alpha$ and $l \in \mathbf{I}^\beta$ – and the cost is tight. \square

Although the proof does not hold where edges are discarded it does provide intuition as to how alpha-expansion minimises the cost. In the first iterations, regions are swept out without breaking all edges, and finding solutions with excessive overlap between models. In subsequent iterations, the region belonging to the boundaries of models contracts cleanly separating parts.

Costs not dependent on α : The local co-occurrence potentials considered here, fall into the class of potentials that can not be exactly optimised by an expansion move over label α . Instead we follow the strategy of [17] and optimise the cost

$$0.5\Delta(\mathbf{I}^\gamma \cap \mathbf{M}_{\text{last}}^\beta \neq \emptyset) + 0.5\Delta(\mathbf{I}_{\text{last}}^\gamma \cap \mathbf{M}^\beta \neq \emptyset) \quad (16)$$

5.3 Merging Parts into Objects: Object-Level Segmentation

The final result of our scene segmentation algorithm is the labelling \mathbf{x} which assigns each feature track to a set of rigid parts. Figure 6 shows some results of the part segmentation (second row) for five videos of the Berkeley Motion Segmentation Dataset [6]. To segment the scene into objects we label connected components of overlapping parts as object detections.

¹ In practice this is almost always true, due to the regularisation caused by overlapping models or pairwise terms.

² Again, most edges are not discarded. The majority of models that would overlap without edge discarding, continue overlapping in the solution found.

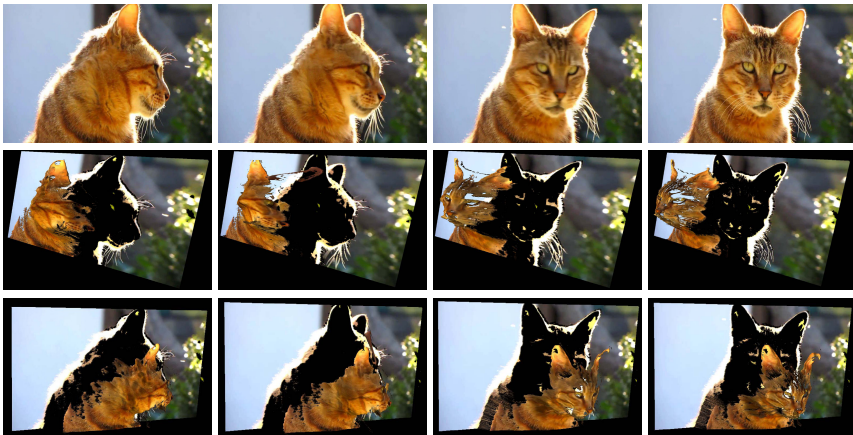


Fig. 4. Reconstruction results for a cat sequence of the *Youtube-Objects* Dataset [23]

Table 1. Evaluation results on the Berkeley Motion Segmentation Dataset using the metrics of [6] *Fayad et al.* shows performance without discarding edges, this the same optimisation as in [9,26].

	Density	overall error	average error	over-segmentation	extracted objects
First 10 frames(26 sequences)					
Brox Malik	3.34%	7.75%	25.01%	0.54	24
Fayad <i>et al.</i>	3.28%	15.23%	51.89%	0.23	7
Our method	3.28%	8.00%	25.46%	1.00	22
First 50 frames(15 sequences)					
Brox Malik	3.27%	7.13%	34.76%	0.53	9
Fayad <i>et al.</i>	3.25%	24.95%	63.67%	0.20	0
Our method	3.25%	5.93%	27.84%	3.70	13
First 200 frames(7 sequences)					
Brox Malik	3.43%	7.64%	31.14%	3.14	7
Fayad <i>et al.</i>	3.42%	28.81%	66.78%	0.29	0
Our method	3.42%	13.28%	39.86%	8.60	4

6 3D Reconstruction

The optimisation of our cost function results in the labelling of rigid models or parts. Using the information about the regions of overlap and the saliency scores, we also have a decomposition of the scene into different objects. In addition, our optimisation estimates model parameters for each rigid model m . Each rigid model is parameterised as the set of fundamental matrices \mathbf{F}_m that describe the epipolar geometry between every pair of consecutive frames.

The 3D reconstruction of each object is then carried out using a piecewise rigid reconstruction approach. For each object we have a list of its constituent parts and a rigid model (set of fundamental matrices) for each of part. First each part is reconstructed independently using the estimated fundamental matrices \mathbf{F}_m . If the calibration of the camera is known, each fundamental matrix can be decomposed into the relative rotation and translation between frames and

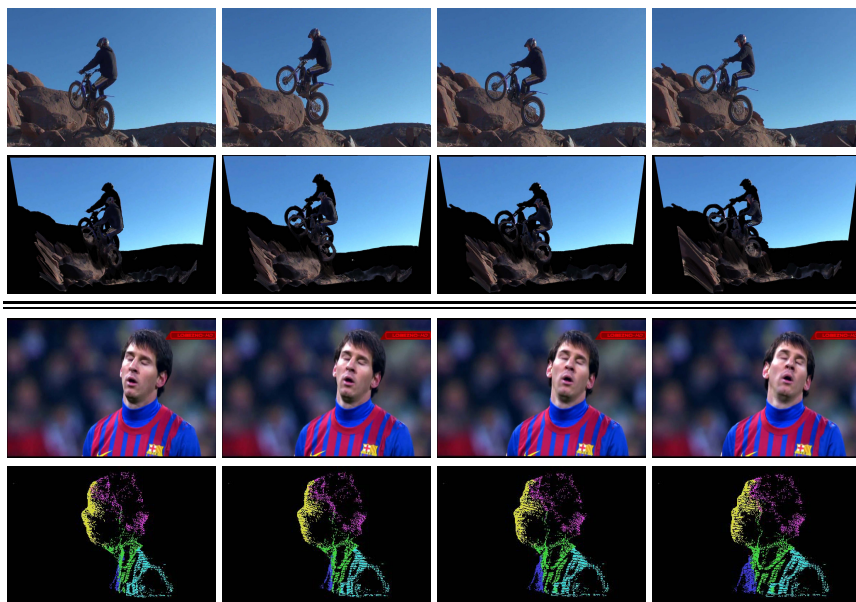


Fig. 5. **Top:** Reconstruction results for a motorbike sequence *Youtube-Objects* Dataset [23]. **Bottom:** Sparse reconstruction of football footage, showing both the assignment of tracks to parts and the quality of reconstruction before densification.

an initial estimate of the shape is obtained using the DLT algorithm [13]. The shape and motion parameters are then refined using the sparse implementation of bundle adjustment [19]. If the camera calibration is unknown the shape is initialized using a factorization algorithm followed by per-frame motion estimation using the PnP algorithm [13]. A final refinement of the shape, motion and focal length parameters is then carried out via bundle adjustment.

Aligning Overlapping Segments: Objects are segmented as a set of overlapping parts that require a final *stitching* step using the areas of overlap to enforce global consistency on the 3D surface. As we use a perspective camera model, the only existing ambiguity between parts is a depth/scale ambiguity which can be resolved by enforcing that tracks belonging to two or more parts should be reconstructed at the same depth by each part model.

Depth-Map Densification: Our reconstruction algorithm is based on sparse feature tracks. To densify the 3D reconstruction, we apply Gaussian filtering on the sparse 3D tracks in xy-RGB image space using the fast implementation of [1] that performs filtering using the permutohedral lattice. Regions of the video far from any tracks in the xy-RGB space are assigned to a flat background billboard.

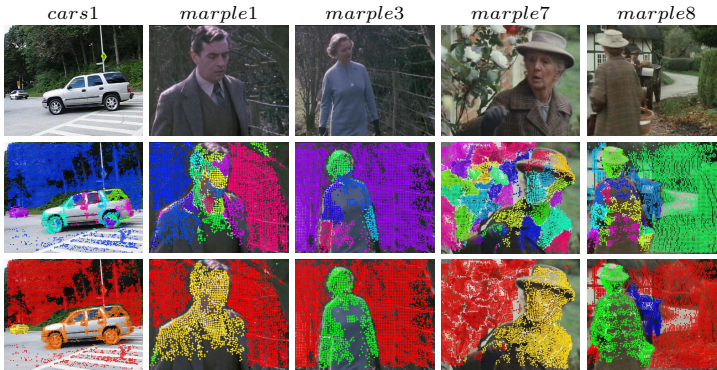


Fig. 6. Motion segmentation results on five sample sequences of the Berkeley Motion Segmentation Dataset [6]. **Second row:** Part segmentation. **Third row:** Object segmentation.

7 Experimental Results

Since we recover both a segmentation of the scene into multiple moving objects and a 3D model for each object, we evaluate both of these steps independently.

Evaluation of the Motion Segmentation Step: We evaluate the results of our *object-level* segmentation on the Berkeley Motion Segmentation Dataset using the tracks and evaluation tool proposed in [6]. Table 1 shows a comparison between the scores of our approach and the results from Brox and Malik’s motion segmentation algorithm [6]. The results show that our method exhibits comparable performance to [6]. While our *over-segmentation* error is higher than [6], the *overall error* and *average error* are very close, and in some cases lower. Although our algorithm can be used for motion segmentation exclusively, it is geared towards 3D reconstruction of complex dynamic scenes. Providing object boundaries are respected, our 3D reconstruction method is unharmed by a slight over-segmentation given that we perform piecewise reconstruction. The same set of parameters was used for all the experiments. The results of Fayad *et al.* [9] show how our algorithm would perform without the novel edge breaking and sparsity terms. Objects are never discovered in sequences longer than 10 frames, and in the majority of the 10 frame long sequences no objects are discovered.

Evaluation of the 3D Reconstruction: We demonstrate our approach on videos from the *Youtube-Objects* Dataset [23]. These are unconstrained real-world videos downloaded from YouTube, with the purpose of object detection in video [23]. Figure 1 shows reconstructions of a *cat*, a *motorbike*, and a *footballer*. We show the decomposition into *parts*, *objects* and a 3D model of the objects from a novel viewpoint for one frame. Figures 4 and 5 show 3D reconstructions for further frames of the three sequences. Our algorithm shows a good segmentation of the scenes and a convincing 3D reconstruction of these challenging videos.

8 Conclusion

In this paper we propose a fully unsupervised approach to the challenging problem of simultaneously segmenting a dynamic scene into its constituent objects and reconstructing a 3D model of the scene. We focus on the reconstruction of real-world videos downloaded from the web or acquired with a single camera observing a complex dynamic scene containing an unknown mixture of multiple moving and possibly deforming objects. Our results show examples of segmentation and 3D reconstruction on videos from the *Youtube Objects* dataset.

References

1. Adams, A., Baek, J., Davis, A.: Fast high-dimensional filtering using the permutohedral lattice. In: Eurographics (2010) 595
2. Bleyer, M., Rother, C., Kohli, P.: Surface stereo with soft segmentation. In: CVPR (2010) 586
3. Boros, E., Hammer, P.L.: Pseudo-boolean optimization. *Discrete Applied Mathematics*, 155–225 (2002) 591
4. Boykov, Y., Kolmogorov, V.: An Experimental Comparison of Min-Cut/Max-Flow Algorithms for Energy Minimization in Vision. *PAMI* 26(9), 1124–1137 (2004) 591
5. Boykov, Y., Veksler, O., Zabih, R.: Fast approximate energy minimization via graph cuts. *PAMI* 23 (2001) 591
6. Brox, T., Malik, J.: Object segmentation by long term analysis of point trajectories. In: Daniilidis, K., Maragos, P., Paragios, N. (eds.) ECCV 2010, Part V. LNCS, vol. 6315, pp. 282–295. Springer, Heidelberg (2010) 585, 593, 594, 596
7. Costeira, J., Kanade, T.: A multi-body factorization method for motion analysis. In: ICCV, pp. 1071–1076 (1995) 585
8. Elhamifar, E., Vidal, R.: Sparse subspace clustering. In: CVPR (2009) 585
9. Fayad, J., Russell, C., Agapito, L.: Automated articulated structure and 3D shape recovery from point correspondences. In: IEEE International Conference on Computer Vision (ICCV), Barcelona, Spain (November 2011) 584, 585, 586, 587, 588, 589, 591, 594, 596
10. Fitzgibbon, A.W., Zisserman, A.: Multibody structure and motion: 3-D reconstruction of independently moving objects. In: Vernon, D. (ed.) ECCV 2000. LNCS, vol. 1842, pp. 891–906. Springer, Heidelberg (2000) 585
11. Galasso, F., Cipolla, R., Schiele, B.: Video segmentation with superpixels. In: Lee, K.M., Matsushita, Y., Rehg, J.M., Hu, Z. (eds.) ACCV 2012, Part I. LNCS, vol. 7724, pp. 760–774. Springer, Heidelberg (2013) 585
12. Garg, R., Roussos, A., Agapito, L.: Dense variational reconstruction of non-rigid surfaces from monocular video. In: CVPR (2013) 583, 584
13. Hartley, R., Zisserman, A.: *Multiple View Geometry in Computer Vision*. Cambridge University Press (2000) 589, 591, 595
14. Isack, H., Boykov, Y.: Energy-based geometric multi-model fitting. *International Journal of Computer Vision (IJCV)* 97(2) (2012) 585, 586, 591
15. Kanatani, K.: Motion segmentation by subspace separation and model selection. In: ICCV, Vancouver, Canada, vol. 2, pp. 301–306 (July 2001) 585
16. Kohli, P., Ladicky, L., Torr, P.: Robust higher order potentials for enforcing label consistency. In: CVPR (2008) 591
17. Ladický, L., Russell, C., Kohli, P., Torr, P.H.: Inference methods for crfs with co-occurrence statistics. *International Journal of Computer Vision* 103(2), 213–225 (2013) 593

18. Li, Z., Guo, J., Cheong, L.-F., Zhou, Z.: Perspective motion segmentation via collaborative clustering. In: ICCV (2013) 585
19. Lourakis, M.A., Argyros, A.: SBA: A Software Package for Generic Sparse Bundle Adjustment. *ACM Trans. Math. Software* (2009) 595
20. Narasimhan, M., Bilmes, J.A.: A submodular-supermodular procedure with applications to discriminative structure learning. arXiv preprint arXiv:1207.1404 (2012) 591
21. Ozden, K., Schindler, K., van Gool, L.: Multibody structure-from-motion in practice. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)* (2010) 583, 585
22. Paladini, M., Del Bue, A., Xavier, J., Agapito, L., Stosic, M., Dodig, M.: Factorization for Non-Rigid and Articulated Structure using Metric Projections. *IJCV* (2012) 585, 586
23. Prest, A., Leistner, C., Civera, J., Schmid, C., Ferrari, V.: Learning object class detectors from weakly annotated video. In: CVPR (2012) 584, 594, 595, 596
24. Rao, S., Tron, R., Vidal, R., Ma, Y.: Motion segmentation in the presence of outlying, incomplete or corrupted trajectories. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)* 32(10), 1832–1845 (2010) 585
25. Roussos, A., Russell, C., Garg, R., Agapito, L.: Dense multibody motion estimation and reconstruction from a handheld camera. In: ISMAR (2012) 583, 585, 586
26. Russell, C., Fayad, J., Agapito, L.: Energy based multiple model fitting for non-rigid structure from motion. In: CVPR (2011) 584, 585, 586, 587, 588, 590, 591, 594
27. Schindler, K., Suter, D., Wang, H.: A model selection framework for multibody structure-and-motion of image sequences. *International Journal of Computer Vision (IJCV)* 79(2), 159–177 (2008) 585
28. Siva, P., Russell, C., Xiang, T., Agapito, L.: Looking beyond the image: Unsupervised learning for object saliency and detection. In: CVPR (2013) 588, 589
29. Sundaram, N., Brox, T., Keutzer, K.: Dense point trajectories by GPU-accelerated large displacement optical flow. In: Daniilidis, K., Maragos, P., Paragios, N. (eds.) ECCV 2010, Part I. LNCS, vol. 6311, pp. 438–451. Springer, Heidelberg (2010) 586
30. Tomasi, C., Kanade, T.: Shape and motion from image streams: a factorization method - part 3 detection and tracking of point features. Technical Report CMU-CS-91-132, Computer Science Department, Carnegie Mellon University, Pittsburgh, PA (April 1991) 585
31. Torresani, L., Hertzmann, A., Bregler, C.: Non-rigid structure-from-motion: Estimating shape and motion with hierarchical priors. *PAMI*, 878–892 (2008) 585, 586
32. Tresadern, P., Reid, I.: Articulated structure from motion by factorization. In: CVPR, vol. 2, pp. 1110–1115 (June 2005) 585
33. Varol, A., Salzmann, M., Tola, E., Fua, P.: Template-free monocular reconstruction of deformable surfaces. In: ICCV (2009) 584, 585
34. Vidal, R., Ma, Y., Sastry, S.: Generalized principal component analysis (gpca). In: CVPR, pp. 621–628 (2003) 585
35. Xu, C., Corso, J.J.: Evaluation of super-voxel methods for early video processing. In: CVPR (2012) 585
36. Yan, J., Pollefeys, M.: A factorization-based approach for articulated non-rigid shape, motion and kinematic chain recovery from video. *PAMI* 30(5) (May 2008) 585
37. Yuille, A.L., Rangarajan, A.: The concave-convex procedure (cccp). In: NIPS (2002) 591
38. Zelnik-Manor, L., Irani, M.: Degeneracies, dependencies and their implications in multi-body and multi-sequence factorizations. In: CVPR, vol. 2, pp. 287–293 (June 2003) 585