

# U4D: Unsupervised 4D Dynamic Scene Understanding

Armin Mustafa

Chris Russell

Adrian Hilton

CVSSP, University of Surrey, United Kingdom

{a.mustafa, c.russell, a.hilton}@surrey.ac.uk

## Abstract

We introduce the first approach to solve the challenging problem of unsupervised 4D visual scene understanding for complex dynamic scenes with multiple interacting people from multi-view video. Our approach simultaneously estimates a detailed model that includes a per-pixel semantically and temporally coherent reconstruction, together with instance-level segmentation exploiting photo-consistency, semantic and motion information. We further leverage recent advances in 3D pose estimation to constrain the joint semantic instance segmentation and 4D temporally coherent reconstruction. This enables per person semantic instance segmentation of multiple interacting people in complex dynamic scenes. Extensive evaluation of the joint visual scene understanding framework against state-of-the-art methods on challenging indoor and outdoor sequences demonstrates a significant ( $\approx 40\%$ ) improvement in semantic segmentation, reconstruction and scene flow accuracy.

## 1. Introduction

With the advent of autonomous vehicles and rising demand for immersive content in augmented and virtual reality, understanding dynamic scenes has become increasingly important. In this paper we propose an unsupervised framework for 4D dynamic scene understanding to address this demand. By “4D Scene understanding” we refer to a unified framework that describes: 3D modelling; motion/flow estimation; and semantic instance segmentation on a per frame basis for an entire sequence. Recent advances in pose estimation [8, 46] and recognition [21, 56, 10] using deep learning have achieved excellent performance for complex images. We exploit these advances to obtain 3D human-pose and an initial semantic instance segmentation from multiple view videos to bootstrap the detailed 4D understanding and modelling of complex dynamic scenes captured with multiple static or moving cameras (see Figure 1). Joint 4D reconstruction allows us to understand how people move and interact, giving contextual information in general scenes.

Existing multi-task methods for scene understanding perform per frame joint reconstruction and semantic in-

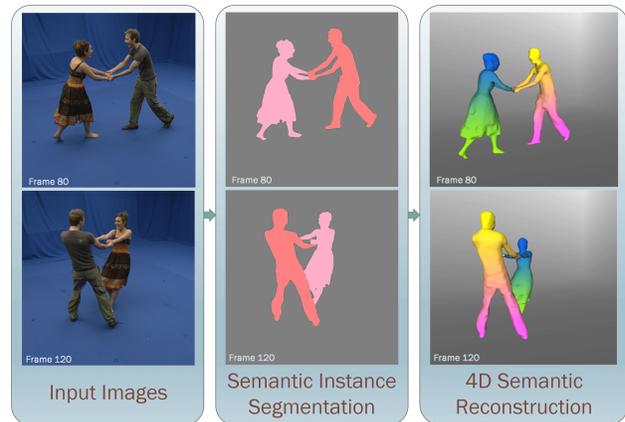


Figure 1. Joint 4D semantic instance segmentation and reconstruction exploiting 3D human-pose of interacting people in dynamic scenes. Shades of pink in segmentation represents instances of people. Colour assigned to reconstruction of frame 80 is reliably propagated to frame 120 using proposed temporal coherence.

stance segmentation from a single image [25], showing that joint estimation can improve each task. Other methods have fused semantic segmentation with reconstruction [36] or flow estimation [42] demonstrating significant improvement in both semantic segmentation and reconstruction/scene flow. We exploit the joint estimation to understand dynamic scenes by simultaneous reconstruction, flow and segmentation estimation from multiple view video.

The first category of methods in joint estimation for dynamic scenes generate segmentation and reconstruction from multi-view [37] and monocular video [16, 30] without any output scene flow estimate. The second category of methods segment and estimates motion in 2D [42], or give spatio-temporal aligned segmentation [11, 34, 12] from multiple views without retrieving the shape of the objects. The third category of methods in 4D temporally coherent reconstruction either align meshes using correspondence information between consecutive frames [58] or extract the scene flow by estimating the pairwise surface correspondence between reconstructions at successive frames [53, 5]. However methods in these three categories do not exploit semantic information of the scene. The fourth category of joint estimation methods exploit semantic information

by introducing joint semantic segmentation and reconstruction for general dynamic scenes [19, 56, 27, 49, 36] and street scenes [13, 50]. However these methods give per-frame semantic segmentation and reconstruction with no motion estimate leading to unaligned geometry and pixel level incoherence in both segmentation and reconstruction for dynamic sequences. Other methods for semantic video segmentation classify objects exploiting spatio-temporal semantic information [48, 34, 11] but do not perform reconstruction. We address this gap in the literature by proposing a novel unsupervised framework for joint multi-view 4D temporally coherent reconstruction, semantic instance segmentation and flow estimation for general dynamic scenes.

Methods in the literature have exploited human-pose information to improve results in semantic segmentation [55] and reconstruction [22]. However existing joint methods for dynamic scenes (with multiple people) do not exploit human-pose information often detecting interacting people as a single object [36]. Table 1 shows a comparison between the tasks performed by state-of-the-art methods. We exploit advances in 3D human-pose estimation to propose the first approach for 4D (3D in time) human-pose based scene understanding of general dynamic scenes with multiple interacting dynamic objects (people) with complex non-rigid motion. 3D human-pose estimation makes full use of multi-view information and is used as a prior to constrain the shape, segmentation and motion in space and time in the joint scene understanding estimation to improve the results. Our contributions are:

- High-level 4D scene understanding for general dynamic scenes from multi-view video.
- Joint instance-level segmentation, temporally coherent reconstruction and scene flow with human-pose priors.
- Robust 4D temporal coherence and per-pixel semantic coherence for dynamic scenes containing interactions.
- An extensive performance evaluation against 15 state-of-the-art methods demonstrating improved semantic segmentation, reconstruction and motion estimation.

## 2. Joint 4D dynamic scene understanding

This section describes our approach to joint 4D scene understanding, with different stages shown in Figure 2. The input to the joint optimisation is multi-view video, per-view initial semantic instance segmentation [21] and 3D human-pose estimation [47]. To achieve stable long-term 4D understanding a set of unique key-frames are detected exploiting multi-view information. Sparse temporal feature tracks are obtained per view between key-frames to initialise the joint estimation. This allows robust 4D understanding in the presence of large non-rigid motion between frames. An initial reconstruction is obtained for each object in the scene combining the initial semantic instance segmentation with the sparse reconstruction [36]. The ini-

	Semantic	Segment	Instance	3D	Motion	Pose
[25, 49, 13]	✓	✓	✓	✓	×	×
[42]	✓	✓	✓	×	✓	×
[36, 19, 27]	✓	✓	×	✓	×	×
[55]	✓	✓	✓	×	×	✓
[22]	×	×	×	✓	✓	✓
[16]	✓	✓	×	✓	✓	×
[30, 41]	×	×	✓	✓	✓	×
[37]	×	✓	×	✓	✓	×
[48, 34, 11]	✓	✓	×	×	✓	×
<b>Proposed</b>	✓	✓	✓	✓	✓	✓

Table 1. Comparison of tasks state-of-the-art methods are solving against the proposed method.

tial reconstruction and semantic instance segmentation is refined for each object instance through novel joint optimisation of segmentation, shape, and motion constrained by 3D human-pose (Section 2.1). Key-frames are used to introduce robust temporal coherence in the joint estimation across long-sequences with large non-rigid deformation. Depth, motion and semantic instance segmentation is combined across views between frames for 4D temporally coherent reconstruction and dense per-pixel semantic coherence for final 4D understanding of scenes (Section 3).

### 2.1. Joint per-view optimisation

Existing methods for semantic segmentation do not give instance level segmentation of the scene. Previous approach either segment the image followed by a per-segment object category classification [35, 18], give deep per-pixel CNN features followed by per-pixel classification in the image [15, 20] or predict semantic segmentation from raw pixels [32] followed by conditional random fields [28, 60]. A recent state-of-the-art method gives a good estimate of initial semantic instance segmentation masks from an image of complex sequence [21]. We employ this approach to predict initial semantic instance segmentation pre-trained parameters on MS-COCO[31] and PASCAL VOC12 [14] for each view. Per-view semantic instance segmentation is combined across views with sparse reconstruction to obtain an initial reconstruction for each frame [36], this is refined through a joint scene understanding optimisation.

The goal of the joint estimation is to refine initial semantic instance segmentation and reconstruction by assigning a label from a set of classes obtained from initial semantic instance segmentation  $\mathcal{L} = \{l_1, \dots, l_{|\mathcal{L}|}\}$  ( $|\mathcal{L}|$  is the total number of classes), a depth value from a set of depth values  $\mathcal{D} = \{d_1, \dots, d_{|\mathcal{D}|-1}, \mathcal{U}\}$  (each depth value is sampled on the ray from camera and  $\mathcal{U}$  is an unknown depth value to handle occlusions), and a motion flow field  $\mathcal{M} = \{m_1, \dots, m_{|\mathcal{M}|}\}$  simultaneously for the region  $\mathcal{R}$  of each object per view.  $|\mathcal{M}|$  is the pre-defined discrete flow-fields for pixel  $p = (x, y)$  in image  $I$  by  $m = (\delta x, \delta y)$  in time. Joint semantic instance segmentation, reconstruction and motion estimation is achieved by global optimisation of a cost function over unary  $E_{unary}$  and pairwise  $E_{pair}$

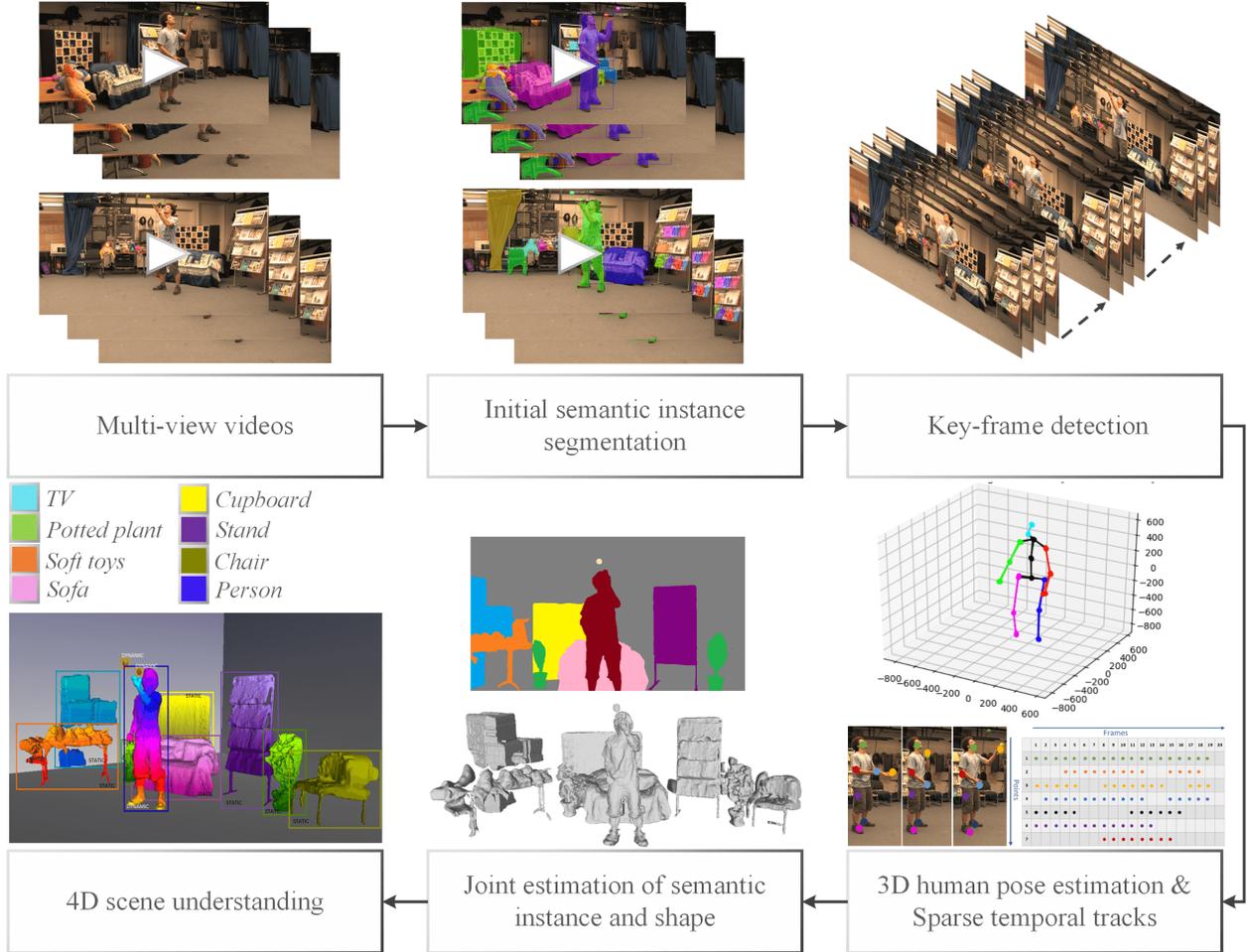


Figure 2. Unsupervised 4D scene understanding framework for dynamic scenes from multi-view video.

terms, defined as:

$$E(l, d, m) = E_{unary}(l, d, m) + E_{pair}(l, d, m) \quad (1)$$

$$E_{unary} = \lambda_d E_d(d) + \lambda_a E_a(l) + \lambda_{sem} E_{sem}(l) + \lambda_f E_f(m)$$

$$E_{pair} = \lambda_s E_s(l, d) + \lambda_c E_c(l) + \lambda_r E_r(l, m) + \lambda_p E_p(l, d, m)$$

where,  $d$  is the depth,  $l$  is the class label, and  $m$  is the motion at pixel  $p$ . Novel terms are introduced for flow  $E_f$ , motion regularisation  $E_r$  and human-pose  $E_p$  costs, explained in Section 2.1.3 and 2.1.2 respectively. Results of the joint optimisation with and without pose ( $E_p$ ) and motion ( $E_f$ ,  $E_r$ ) information are presented in Figure 3, showing the improvement in results. Ablative analysis on individual costs in Section 4 show the improvement in performance with the novel introduction of motion and pose constraints in the joint optimisation. Standard unary terms for depth ( $E_d$ ), semantic ( $E_{sem}$ ), and appearance ( $E_a$ ) costs are used [36], explained in Section 2.1.5. Standard pairwise terms colour contrast ( $E_c$ ) is used to assist segmentation and smoothness ( $E_s$ ) cost ensures that depth vary smoothly in a neighbourhood, are explained in **Appendix A** of the supplementary material.

Global optimisation of Equation 1 is performed over all terms simultaneously, using the  $\alpha$ -expansion algorithm by iterating through the set of labels in  $\mathcal{L} \times \mathcal{D} \times \mathcal{M}$  [7]. Each iteration is solved by graph-cut using the min-cut/max-flow algorithm [6]. Convergence is achieved in 7-8 iterations.

### 2.1.1 Spatio-temporal coherence in the optimisation

Constraints are applied on the spatial and temporal neighborhood to enforce consistency in the appearance, semantic label, 3D human pose and motion across views and time.

**Spatial coherence:** Multi-view spatial coherence is enforced in the optimisation such that the motion, shape, appearance, 3D pose and class labels are consistent across views using an 8-connected spatial neighbourhood  $\psi_S$  for each camera view such that the set of pixel pairs  $(p; q)$  belong to the same frame.

**Temporal coherence:** Temporal coherence is enforced in the joint optimisation by enforcing coherence across key-frames to handle large non-rigid motion and to reduce errors in sequential alignment for long sequences in the 4D scene understanding. Sparse temporal feature correspondences are used for key-frame detection and robust initiali-

sation of the joint optimisation. They measure the similarity between frames and unlike optical flow are robust to large motions and visual ambiguity. To achieve robust temporal coherence in the 4D scene understanding framework for large non-rigid motion, sparse temporal feature correspondences in 3D are obtained across the sequence.

The temporal neighbourhood is defined for each frame between its respective key-frames. Sparse temporal correspondence tracks define the temporal neighbourhood  $\psi_T = \{(p, q) \mid q = p + e_{i,j}\}$ ; where  $j = \{t-1, t+1\}$  and  $e_{i,j}$  is the displacement vector from image  $i$  to  $j$ .

### 2.1.2 Human-pose constraints $E_p(l, d, m)$

We use 3D human-pose to constrain joint optimisation and improve the flow, reconstruction and instance segmentation, in both 2D and 3D for dynamic scenes with multiple interacting people (see Figure 1). 3D human-pose is used as it is consistent across multiple views unlike 2D human-pose. A state-of-the-art method for 3D human-pose estimation from multiple cameras [47] is used in the paper. Previous work on 3D pose estimation [46] iteratively builds a 3D model of human-pose consistent with 2D estimates of joint locations and prior knowledge of natural body pose. In [47], multiple cameras are used when estimating the 3D model; this then feeds back into new estimates of the 2D joint locations in each image. This approach allows us to take full advantage of 3D estimates of pose, consistent across all cameras when finding fine grained 2D correspondences between images, and leading to more lifelike, vivid human reconstructions.

Initial semantic reconstruction is updated if the 3D pose of the person lies outside the region  $\mathcal{R}$  by dilating the boundary to include the missing joints. This allows for more robust and complete reconstruction and segmentation. We use a standard set of 17 joints [47] defined as  $\mathcal{B}$ . A circle  $\mathcal{C}_i$  is placed around the joint position in 2D and a sphere  $\mathcal{S}_i$  is placed around the joint position in 3D based on the confidence map to identify the nearest neighbour vertices for every joint  $b_i$ .

$$E_p(l, d, m) = \sum_{b_i \in \mathcal{B}} \lambda_{2d} e_{2d}(l, m) + \lambda_{3d} e_{3d}(d) \quad (2)$$

$$e_{2d}(l, m) = e_{2d}^L(l) + e_{2d}^S(l) + e_{2d}^M(m)$$

$$e_{3d}(d) = e_{3d}^M(d) + e_{3d}^S(d), \text{ if } d_p \neq \mathcal{U} \text{ else } 0$$

**3D shape term:** This term constrains the reconstruction in 3D such that the neighbourhood points around the joints do not move far from the respective joints, and is defined as:

$$e_{3d}^S(d) = \exp\left(-\frac{1}{|\sigma_{S_D}|} \sum_{\Phi(p) \in \mathcal{S}_i} \|O\|_F^2\right)$$

where  $\Phi(p)$  is the 3D projection of pixel  $p$ . The Frobenius norm  $\|O\|_F = \|\begin{bmatrix} \Phi(p) & b_i \end{bmatrix}\|_F$  is applied on the 3D points in all directions to obtain the net motion at each pixel within

$$\mathcal{S}_i \text{ and } \sigma_{S_D} = \left\langle \frac{\|O\|_F^2}{\vartheta_{\Phi(p), b_i}} \right\rangle.$$

**3D motion term:** This enforces as rigid as possible [43] constraint on 3D points in the neighbourhood of each joint  $b_i$  in space and time. An optimal rotation matrix  $R_i$  is estimated for each  $b_i$  by minimising the energy defined as:

$$e_{3d}^M(d) = \sum_{\Phi(p) \in \mathcal{S}_i} \left\| (b_i^{t+1} - \Phi(p)^{t+1}) - R_i (b_i^t - \Phi(p)^t) \right\|_2^2 + \lambda_{3d}^p \|p - e_{3d}^M\|_2^2$$

**2D term:** 3D poses are back-projected in each view to constrain per view appearance ( $e_{2d}^L$ ), semantic segmentation ( $e_{2d}^S$ ) and motion estimation ( $e_{2d}^M$ ) in 2D. If  $p \in \mathcal{C}_i$ ,

$$e_{2d}^L(l) = \exp\left(-\sum_{p \in \psi_S} \sum_{p \in \psi_T} \frac{\|I(\Pi(b_i)) - I(p)\|^2}{|\sigma_{S_L}|}\right)$$

$$e_{2d}^S(l) = \exp\left(-\sum_{p \in \psi_S} \sum_{p \in \psi_T} \frac{\|\Pi(b_i) - p\|^2}{|\sigma_{S_S}|}\right)$$

$$e_{2d}^M(m) = \exp\left(-\sum_{p \in \psi_S} \sum_{k \in \psi_T} \frac{\left\| \vartheta_{p, \Pi(b_i^k)} - \vartheta_{p+m_p, \Pi(b_i^{k+1})} \right\|^2}{|\sigma_{S_M}|}\right)$$

where,  $\Pi$  is the back-projection of 3D poses to 2D,  $N_{pose}$  is the number of nearest neighbours,  $\sigma_{S_L} = \left\langle \frac{\|\Pi(b_i) - q\|^2}{\vartheta_{\Pi(b_i), q}} \right\rangle$  and,  $\sigma_{S_S}$  and  $\sigma_{S_M}$  is defined similarly.  $e_{2d}^L(l)$  and  $e_{2d}^S(l)$  ensures that the pixels around projected 3D pose  $\Pi(b_i)$  have the same semantic label and appearance across views ( $\psi_S$ ) and time ( $\psi_T$ ) thereby ensuring spatio-temporal appearance and semantic consistency respectively.

### 2.1.3 Motion constraints- $E_f(m)$ and $E_r(l, m)$

**Flow term:** This term is obtained by integrating the sum of three penalisers over the reference image domain inspired from [45], defined as:

$$E_f(p, m_p) = e_F^T(p, m_p) + e_F^V(p, m_p) + e_F^S(p, m_p)$$

where,  $e_F^T(p, m_p) = \sum_{i=1}^{N_v} \|(I_i(p, t) - I_i(p + m_p, t + 1))\|^2$  penalises deviation from the brightness constancy assumption in a temporal neighbourhood for the same view;  $e_F^V(p, m_p) = \sum_{t \in \psi_T} \sum_{i=2}^{N_v} \|(I_1(p, t) - I_i(p + m_p, t))\|^2$  penalises deviation in appearance from the brightness constancy assumption between the reference view and other views at other time instants; and  $e_F^S(p, m_p) = 0$  if  $p \in N$  otherwise  $\infty$  which forces the flow to be close to nearby sparse temporal correspondences.  $I_i(p, t)$  is the intensity at point  $p$  at time  $t$  in camera  $i$ . The flow vector  $m$  is located within a window from a sparse constraint at  $p$  and it forces the flow to approximate the sparse 2D temporal correspondences.

**Motion regularisation term:** This penalises the absolute difference of the flow field to enforce motion smoothness

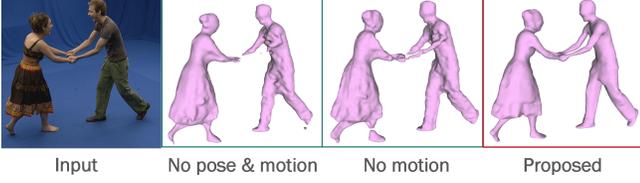


Figure 3. Comparison of reconstruction without pose and motion in the optimisation framework, proposed result is best.

and handle occlusions in areas with low confidence [45].

$$E_r(l, m) = \sum_{p, q \in N_p} \|\Delta m\|^2 \lambda_r^L e_r^L(p, q, m_p, m_q, l_p, l_q) + \lambda_r^A e_r^A(p, q, m_p, m_q, l_p, l_q)$$

where  $\Delta m = m_p - m_q$  and;

$$e_r^X = \begin{cases} \forall_{l_p=l_q} \text{mean}_{q \in N_p} E_X(q, m_q) - \min_{q \in N_p} E_X(q, m_q) & \text{else } 0. \end{cases}$$

We compute  $e_r^L$  (semantic regularisation) and  $e_r^A$  (appearance regularisation) as the minimum subtracted from the mean energy within the search window  $N_p$  for each pixel  $p$ .

### 2.1.4 Long-term temporal coherence

**Sparse temporal correspondences:** The sparse 3D points projected in all views are matched between frames  $N_f^i$  and key-frames across the sequence using nearest neighbour matching [33] followed by a symmetry test which employs forward and backward match consistency by performing two-way matching to remove the inconsistent correspondences. This gives sparse temporal feature correspondence tracks per frame for each object:  $F_i^c = \{f_1^c, f_2^c, \dots, f_{R_i^c}^c\}$ , where  $c = 1$  to  $N_v$ .  $R_i^c$  are the 3D points visible at each frame  $i$ . Exhaustive matching is done, such that each frame is matched to every other frame to handle appearance, reappearance and disappearance of points between frames.

**Key-frame detection:** Previous work [40, 39] showed that sparse key-frames allow robust long-term correspondence for 4D reconstruction. In this work we introduce the additional use of pose in the detection and sparse temporal feature correspondence across key-frames to prevent the accumulation of errors in long sequences. 4D scene alignment between key-frames is explained in Section 3.

**Key-frame similarity metric** is defined as:

$$KS_{i,j} = 1 - \frac{1}{5N_v} \sum_{c=1}^{N_v} (M_{i,j}^c + L_{i,j}^c + D_{i,j}^c + P_{i,j}^c + I_{i,j}^c) \quad (3)$$

Key-frame detection exploits sparse correspondence ( $M_{i,j}^c$ ), pose ( $P_{i,j}^c$ ), shape ( $I_{i,j}^c$ ), semantic ( $L_{i,j}^c$ ) and distance ( $D_{i,j}^c$ ) information across views  $N_v$  between frame  $i$  and  $j$  for each object in view  $c$ , to improve the long-term temporal coherence of the proposed method, using similar frames across the sequence, illustrated in Figure 4. All frames with similarity  $> 0.75$  in a sequence are selected as key-frames defined as  $K = \{k^1, k^2, \dots, k^{N_k}\}$  where  $N_k$  is the number of key-frames and  $N_f^i$  is the number of frames between  $K_i$  and  $K_{i+1}$ . All the metrics used in 3 and an ablation study for key-frame detection is given in detail in **Appendix B** of supplementary material.

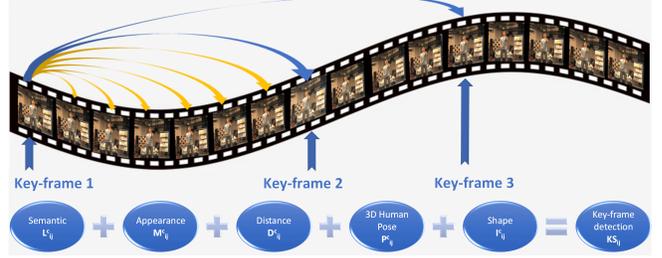


Figure 4. An illustration of key-frame detection and matching across a short sequence for stable long-term temporal coherence.

Features at view  $c$  frame  $i$ ,  $F_i^c$  are matched to features at view  $c$  to frames  $j = \{i+1, \dots, N_f^i\}$  to give correspondences for all the frames  $N_f^i$  with key-frame  $K_i$ . The corresponding joint locations from the 3D pose are back-projected in each view and added to sparse temporal tracks in between key-frames. Any new point-tracks are added to the list of point tracks for key-frame  $K_i$ .

### 2.1.5 Unary terms - $E_{unary}(l, d, m)$

**Depth term:** This gives a measure of photo-consistency between views  $E_d(d) = \sum_{p \in \psi_S} e_d(p, d_p)$ , defined as:

$$e_d(p, d_p) = \begin{cases} M(p, q) = \sum_{i \in \mathcal{O}_k} m(p, q), & \text{if } d_p \neq \mathcal{U} \\ M_{\mathcal{U}}, & \text{if } d_p = \mathcal{U} \end{cases}$$

where  $M_{\mathcal{U}}$  is the fixed cost of labelling pixel unknown and  $q$  denotes the projection of the hypothesised point  $P$  (3D point along the optical ray passing through pixel  $p$  located at a distance  $d_p$  from the camera) in an auxiliary camera.  $\mathcal{O}_k$  is the set of the  $k$  most photo-consistent pairs with reference camera and  $m(p, q)$  is inspired from [37].

**Appearance term:** This term is computed using the negative log likelihood [6] of the colour models (GMMs with 10 components) learned from the initial semantic mask in the temporal neighbourhood  $\psi_T$  and the foreground markers obtained from the sparse 3D features for the dynamic objects. It is defined as:

$$E_a(l) = \sum_{p \in \psi_T} \sum_{p \in \psi_S} -\log P(I_p | l_p)$$

where  $P(I_p | l_p = l_i)$  denotes the probability of pixel  $p$  belonging to layer  $l_i$ .

**Semantic term:** This term is based on the probability of the class labels at each pixel based on [10], defined as:

$$E_{sem}(l) = \sum_{p \in \psi_T} \sum_{p \in \psi_S} -\log P_{sem}(I_p | l_p)$$

where  $P_{sem}(I_p | l_p = l_i)$  denotes the probability of pixel  $p$  being in layer  $l_i$  in the reference image obtained from initial semantic instance segmentation [21].

## 3. 4D scene understanding

The final 4D scene model fuses the semantic instance segmentation, depth information and dense flow across views and in time between frames ( $N_f^i$ ) and key-frames ( $K_i$ ). The initial instance segmentation, human pose and motion information for each object is combined to obtain final instance segmentation of the scene. The depth informa-

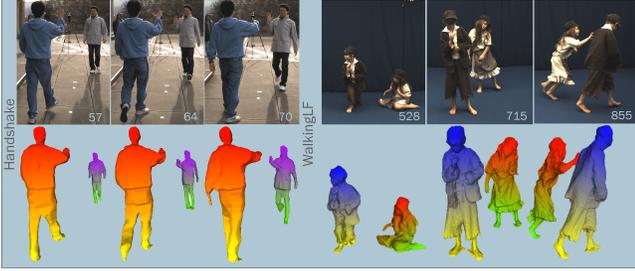


Figure 5. Example of 4D scene reconstruction for two datasets

Datasets	Resolution	$N_v$	Baseline	L	KF	Tracks
Handshake[26]	$1920 \times 1080$	8(all S)	$15^\circ$ - $30^\circ$	125	15	1945
Meetup[17]	$1920 \times 1080$	16(all S)	$25^\circ$ - $35^\circ$	100	9	1341
Juggler2[4]	$960 \times 544$	6(all M)	$15^\circ$ - $45^\circ$	300	16	1278
Handstand[51]	$1600 \times 1200$	8(all S)	$25^\circ$ - $45^\circ$	174	12	1056
Rachel[2]	$3840 \times 2160$	16(all S)	$20^\circ$ - $30^\circ$	270	15	1978
Juggler1[2]	$1920 \times 1080$	8(2 M)	$15^\circ$ - $30^\circ$	253	17	2083
Dance[1]	$780 \times 582$	8(all S)	$35^\circ$ - $45^\circ$	60	7	732
Magician[4]	$960 \times 544$	6(all M)	$15^\circ$ - $45^\circ$	300	10	1312
Human3.6[23]	$1000 \times 1000$	4(all S)	$25^\circ$ - $30^\circ$	250	14	994
MagicianLF[39]	$2048 \times 2048$	25(all S)	$5^\circ$ - $8^\circ$	350	5	1312
WalkLF[39]	$2048 \times 2048$	20(all S)	$5^\circ$ - $8^\circ$	221	7	1934

Table 2. Properties of all datasets:  $N_v$  is the number of views, L is the sequence length, KF gives number of key-frames, and Tracks gives the number of sparse temporal correspondence tracks averaged over the entire sequence for each object (S stands for static cameras and M for moving cameras).

tion is combined across views using Poisson surface reconstruction [24] to obtain a mesh for each object in the scene. 4D temporally coherent meshes are obtained by combining the most consistent motion information from all views for each 3D point. This is combined with spatial semantic instance information to give per-pixel semantic and temporal coherence. Appearing, disappearing, and reappearing regions are handled by using the sparse temporal tracks and their respective motion estimate. The dense flow and semantic instance segmentation together with 3D models of each object in the scene gives the final 4D understanding of the scenes. Examples are shown in Figure 1 and 5 on two datasets, where objects are coloured in one key-frame and colours are propagated reliably between frames and key-frames across the sequence for robust 4D scene modelling.

## 4. Results and evaluation

Joint semantic instance segmentation, reconstruction and flow estimation (section 2) is evaluated quantitatively and qualitatively against 15 state-of-the-art methods on a variety of publically available multi-view indoor and outdoor dynamic scene datasets, detailed in Table 2. More results are provided in supplementary material *Appendix C*.

Algorithm parameters listed in Table 3 are the same for all outdoor datasets, and for indoor datasets parameters depend on the number of cameras ( $N_v$ ). Pairwise costs are constant  $\lambda_p = 0.9$ ,  $\lambda_c = \lambda_s = \lambda_r = 0.5$  for all datasets.

	$\lambda_d$	$\lambda_a$	$\lambda_{sem}$	$\lambda_f$	$\lambda_s^t/\lambda_s^s$	$\lambda_{ca}/\lambda_{ct}$	$\lambda_r^L/\lambda_r^C$	$\lambda_{2d}/\lambda_{3d}$
Outdoor	1.2	0.5	0.5	0.4	1.0	5.0	0.6	7.5
I, $N_v < 6$	1.0	0.7	0.5	0.6	0.4	5.0	0.4	7.5
I, $6 \leq N_v < 20$	1.0	0.7	0.2	0.4	0.4	5.0	0.4	5.0
I, $N_v \geq 20$	1.0	1.0	0.5	0.5	0.2	5.0	0.4	5.0

Table 3. Parameters for all datasets. I is Indoor

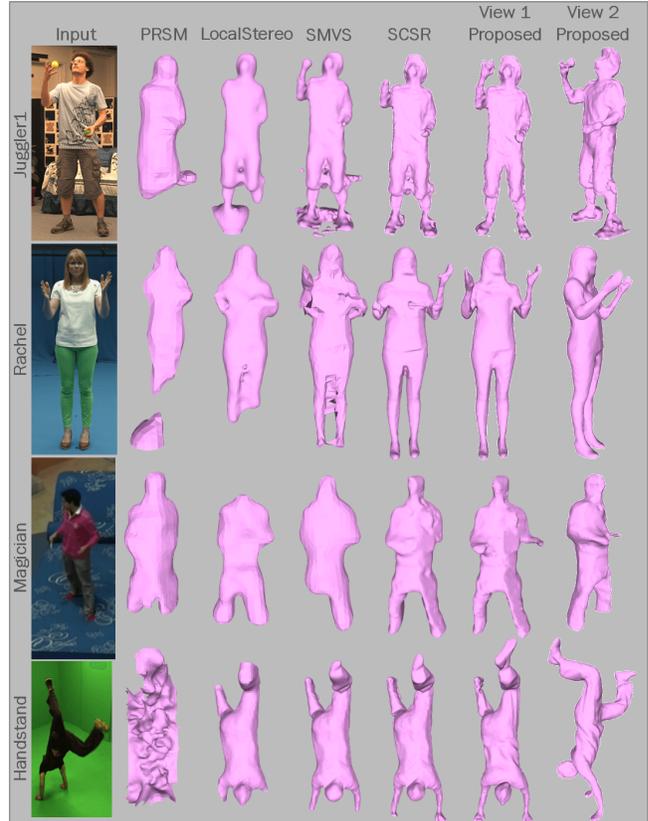


Figure 6. Reconstruction evaluation against existing methods. Two different views of 3D model are shown for proposed method.

### 4.1. Reconstruction evaluation

The proposed approach is compared against state-of-the-art approaches for semantic co-segmentation and reconstruction (SCSR) [36], piecewise scene flow (PRSM) [52], multi-view stereo (SMVS) [29], and deep learning based stereo approaches (LocalStereo) [44]. Qualitative comparison with 2 views of proposed method are shown in Figure 6. Pre-trained parameters were used for LocalStereo and per-view depth maps were fused using Poisson reconstruction. The quality of surface obtained using proposed method is improved compared to state-of-the-art methods. In contrast to previous approaches, limbs of people are reliably reconstructed because of the exploitation of human-pose and temporal information (motion) in the joint optimisation.

For quantitative comparison to state-of-the-art methods, we project the reconstruction onto different views and compute the projection errors shown in Table 4. A significant improvement is obtained in projected surface completeness with the proposed approach.

Methods	Handshake	Handstand	Rachel	Juggler1	Juggler2	Magician	Dance	Meetup	Human3.6	MagicianLF	WalkLF
PRSM [52]	1.56	1.79	1.51	1.57	1.68	1.72	1.79	1.98	2.01	1.59	1.41
LS [44]	1.24	1.38	1.15	1.21	1.18	1.33	1.46	1.47	1.64	1.20	1.23
SMVS [29]	0.84	0.97	0.73	0.75	0.85	0.92	0.85	0.96	1.19	0.94	0.88
SCSR [36]	0.70	0.84	0.67	0.69	0.73	0.78	0.77	0.87	0.92	0.77	0.71
$P_{PS}$	0.73	0.87	0.65	0.70	0.71	0.75	0.74	0.88	0.90	0.78	0.70
$P_{PM}$	0.71	0.85	0.64	0.68	0.69	0.73	0.72	0.85	0.87	0.75	0.68
$P_P$	0.57	0.71	0.56	0.59	0.61	0.64	0.62	0.75	0.77	0.67	0.63
$P_S$	0.59	0.69	0.59	0.57	0.63	0.66	0.60	0.73	0.76	0.65	0.60
$P_M$	0.55	0.68	0.55	0.54	0.59	0.61	0.59	0.74	0.73	0.62	0.59
Proposed	<b>0.46</b>	<b>0.55</b>	<b>0.47</b>	<b>0.49</b>	<b>0.51</b>	<b>0.53</b>	<b>0.55</b>	<b>0.57</b>	<b>0.60</b>	<b>0.49</b>	<b>0.44</b>

Table 4. Reconstruction evaluation: Projection error across views against state-of-the-art methods, LS is LocalStereo.  $P_P = E - E_p$ ,  $P_M = E - E_f - E_r$ ,  $P_{PM} = E - E_f - E_r - E_p$ ,  $P_S = E - E_{sem}$  and  $P_{PS} = E - E_{sem} - E_p$ , where  $E$  is defined in Equation 1.

Methods	Handshake	Handstand	Rachel	Juggler1	Juggler2	Magician	Dance	Meetup	Human3.6	MagicianLF	WalkLF
CRFRNN [60]	62.7	55.8	61.6	40.5	68.7	52.4	49.3	41.1	42.9	60.8	63.6
Segnet [3]	47.9	51.1	55.2	45.1	61.9	55.3	53.9	43.9	49.4	59.3	65.9
JSR [17]	67.8	58.7	58.4	56.2	66.0	61.3	57.9	50.2	53.4	62.3	68.9
SCV [48]	56.4	52.6	48.8	49.5	59.1	59.2	56.7	42.0	49.1	58.2	65.7
Dv3+ [9]	63.8	58.9	64.0	48.8	69.7	58.9	57.6	48.4	54.8	69.6	69.1
MRCNN [21]	65.2	59.6	67.4	50.3	70.5	60.5	58.7	47.2	53.4	69.5	70.2
PSP [59]	74.7	64.5	75.5	67.9	81.2	73.4	71.5	62.6	65.3	74.6	82.5
SCSR [36]	81.8	75.2	78.4	81.4	89.3	88.2	85.1	78.9	70.4	82.2	86.7
$P_{PM}$	85.7	75.9	78.6	81.8	89.6	88.5	85.5	79.2	70.6	82.9	87.5
$P_P$	86.3	77.4	80.7	82.6	90.1	89.1	87.6	80.8	76.3	86.1	89.3
$P_M$	87.6	79.1	81.7	83.5	90.5	89.6	86.4	81.9	75.4	85.2	88.1
Proposed	<b>89.6</b>	<b>83.3</b>	<b>85.8</b>	<b>88.2</b>	<b>91.1</b>	<b>90.9</b>	<b>88.5</b>	<b>84.7</b>	<b>81.1</b>	<b>89.4</b>	<b>91.8</b>

Table 5. Segmentation comparison against state-of-the-art methods using the *Intersection-over-Union* metric.

## 4.2. Segmentation evaluation

Our approach is evaluated against a variety of state-of-the-art multi-view (SCV [48], SCSR [36], and JSR [17]) and single-view (Dv3+ [9], MRCNN [21], PSP [59], CRF RNN [60], and Segnet [3]) segmentation methods, shown in Figure 7. For fair evaluation against single-view semantic segmentation methods, multi-view consistency is applied for segmentation estimated from each view to obtain multi-view consistent semantic segmentation using dense multi-view correspondence. Colour in the results is kept from the original papers. Only MRCNN and the proposed approach gives instance segmentation.

Quantitative evaluation against state-of-the-art methods is measured by *Intersection-over-Union* with ground-truth, shown in Table 5. Ground-truth is available on-line for most of the datasets and obtained by manual labelling for other datasets. Pre-trained parameters were used for semantic segmentation methods. The semantic instance segmentation results from the joint optimisation are significantly better compared to the state-of-the-art methods ( $\approx 20 - 40\%$ ).

## 4.3. Motion evaluation

Flow from the joint estimation is evaluated against state-of-the-art methods: (a) Dense flow algorithms DCflow [57] and Deepflow [54]; (b) Scene flow methods PRSM [52]; and (c) Non-sequential alignment of partial surfaces 4DMatch [38] (requires a prior 3D mesh of the object as input for 4D reconstruction). The key-frames of sequence are coloured and the colour is propagated using dense flow from the joint optimisation throughout the sequence. The

red regions in 2D dense flow in Figure 8 are the regions for which reliable correspondences are not found. This demonstrates improved performance using the proposed method. The colours in the 4D alignment in Figure 9 are not reliably propagated by DCflow for limbs.

We also compare the silhouette overlap error ( $S_e$ ) across frames, key-frames and views to evaluate long-term temporal coherence in Table 6 for all datasets. This is defined as 
$$S_e = \frac{1}{N_v N_k N_f^i} \sum_{i=1}^{N_k} \sum_{j=1}^{N_f^i} \sum_{c=1}^{N_v} \frac{\text{Area of intersection}}{\text{Area of semantic segmentation}}$$
. Dense flow in time is used to obtain the propagated mask for each image. The propagated mask is overlapped with semantic segmentation at each time instant to evaluate the accuracy of the propagated mask. The lower the  $S_e$  the better. Our approach gives the lowest error demonstrating higher accuracy compared to the state-of-the-art methods.

## 4.4. Ablation study on Equation 1

We perform an ablation study on Equation 1, such that we remove motion  $E_f, E_r$ , pose  $E_p$  and semantic  $E_{sem}$  constraints from the equation, defining  $P_M = E - E_f - E_r$ ,  $P_P = E - E_p$ ,  $P_{PM} = E - E_f - E_r - E_p$ ,  $P_S = E - E_{sem}$  and  $P_{PS} = E - E_{sem} - E_p$ . Reconstruction, flow and semantic segmentation is obtained with removed constraints, and the results are shown in Tables 4, 6 and 5 respectively. The proposed approach gives best performance with joint pose, motion and semantic constraints.

## 4.5. Limitations

Gross errors in initial semantic instance segmentation and 3D pose estimation lead to degradation in the quality of

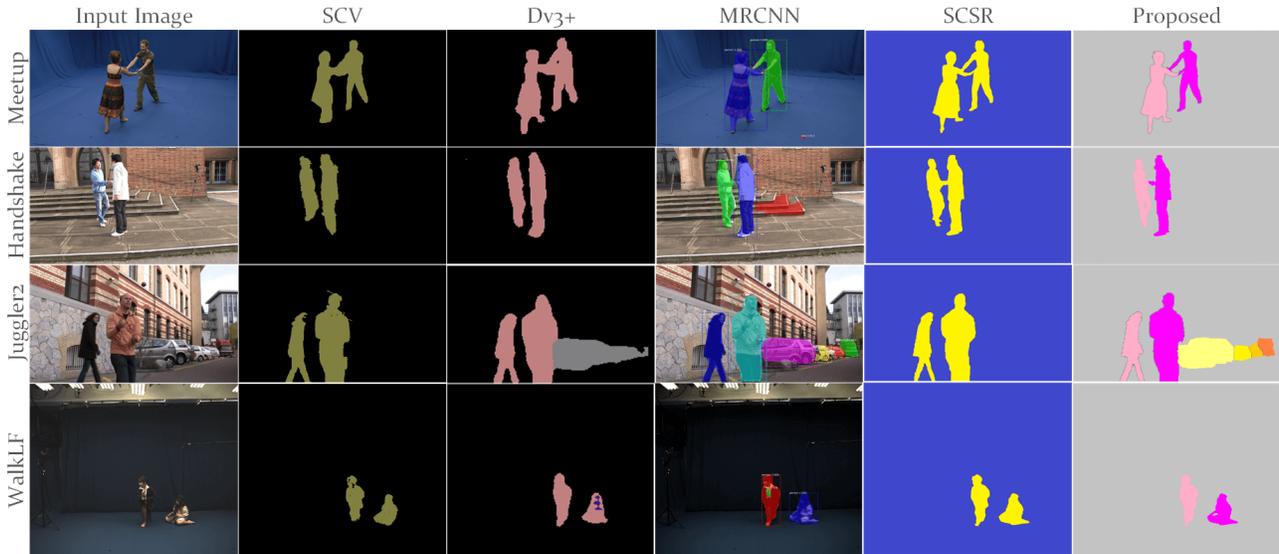


Figure 7. Semantic segmentation comparison against state-of-the-art methods. In the proposed method shades of pink depicts instances of humans and shades of yellow depict instances of cars.

Methods	Handshake	Handstand	Rachel	Juggler1	Juggler2	Magician	Dance	Meetup	Human3.6	MagicianLF	WalkLF
PRSM [57]	1.80	2.15	1.54	1.65	1.79	1.96	1.87	2.11	2.34	1.87	1.52
Deepflow [54]	1.15	1.48	1.01	1.08	1.16	1.27	1.21	1.37	1.52	1.05	0.81
DCFlow [52]	0.90	1.17	0.97	0.87	0.93	1.03	0.96	1.12	1.21	0.83	0.79
4DMatch [38]	0.79	0.98	0.75	0.69	0.87	0.81	0.77	0.87	0.94	0.80	0.77
$P_{PS}$	0.75	1.01	0.85	0.78	0.91	0.93	0.86	0.99	1.07	0.81	0.78
$P_P$	0.71	0.93	0.80	0.73	0.84	0.87	0.78	0.92	0.99	0.76	0.73
$P_S$	0.64	0.77	0.63	0.61	0.65	0.72	0.65	0.76	0.81	0.64	0.61
Proposed	<b>0.51</b>	<b>0.61</b>	<b>0.48</b>	<b>0.49</b>	<b>0.52</b>	<b>0.58</b>	<b>0.55</b>	<b>0.63</b>	<b>0.68</b>	<b>0.53</b>	<b>0.44</b>

Table 6. Silhouette overlap error for multi-view datasets for evaluation of long-term temporal coherence, where .

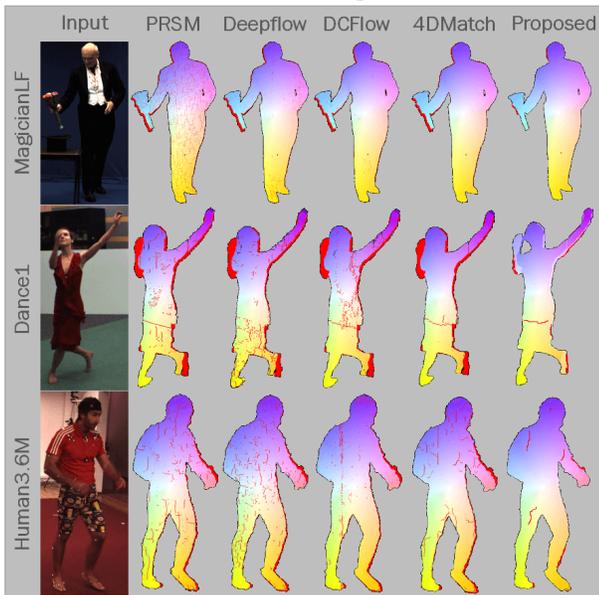


Figure 8. Temporal coherence evaluation against existing methods. results (e.g. the cars in Juggler2 - Figure 7). Although 3D human pose helps in robust 4D reconstruction of interacting people in dynamic scenes, current 3D pose estimation is unreliable for highly crowded environments resulting in degradation of the proposed approach.

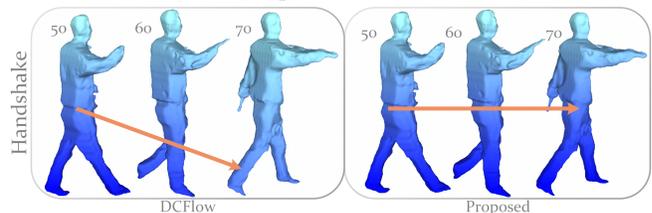


Figure 9. 4D alignment evaluation against DCFlow [57].

## 5. Conclusions

This paper introduced the first method for unsupervised 4D dynamic scene understanding from multi-view video. A novel joint flow, reconstruction and semantic instance segmentation estimation framework is introduced exploiting 2D/3D human-pose, motion, semantic, shape and appearance information in space and time. Ablation study on the joint optimisation demonstrates the effectiveness of the proposed scene understanding framework for general scenes with multiple interacting people. The semantic, motion and depth information per view is fused spatially across views for 4D semantically and temporally coherent scene understanding. Extensive evaluation against state-of-the-art methods on a variety of complex indoor and outdoor datasets with large non-rigid deformations demonstrates a significant improvement in the accuracy in semantic segmentation, reconstruction, motion estimation and 4D alignment.

## References

- [1] 4d repository, <http://4drepository.inrialpes.fr/>. In *Institut national de recherche en informatique et en automatique (INRIA) Rhone Alpes*. 6
- [2] Multiview video repository, <http://cvssp.org/data/cvssp3d/>. In *Centre for Vision Speech and Signal Processing, University of Surrey, UK*. 6
- [3] V. Badrinarayanan, A. Kendall, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *TPAMI*, 2017. 7
- [4] L. Ballan, G. J. Brostow, J. Puwein, and M. Pollefeys. Unstructured video-based rendering: Interactive exploration of casually captured videos. *ACM Trans. Graph.*, 29(4):1–11, 2010. 6
- [5] T. Basha, Y. Moses, and N. Kiryati. Multi-view scene flow estimation: A view centered variational approach. In *CVPR*, pages 1506–1513, 2010. 1
- [6] Y. Boykov and V. Kolmogorov. An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision. *TPAMI*, 26(11):1124–1137, 2004. 3, 5
- [7] Y. Boykov, O. Veksler, and R. Zabih. Fast approximate energy minimization via graph cuts. *TPAMI*, 23(11):1222–1239, 2001. 3
- [8] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017. 1
- [9] L. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. *CoRR*, abs/1802.02611, 2018. 7
- [10] L. C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *CoRR*, abs/1606.00915, 2016. 1, 5
- [11] W.-C. Chiu and M. Fritz. Multi-class video co-segmentation with a generative multi-video model. In *CVPR*, 2013. 1, 2
- [12] A. Djelouah, J.-S. Franco, E. Boyer, P. Pérez, and G. Dretakis. Cotemporal Multi-View Video Segmentation. In *3DV*, 2016. 1
- [13] F. Engelmann, J. Stückler, and B. Leibe. Joint object pose estimation and shape reconstruction in urban street scenes using 3D shape priors. In *GCPR*, 2016. 2
- [14] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>. 2
- [15] C. Farabet, C. Couprie, L. Najman, and Y. LeCun. Learning hierarchical features for scene labeling. *TPAMI*, 35(8):1915–1929, 2013. 2
- [16] G. Floros and B. Leibe. Joint 2d-3d temporally consistent semantic segmentation of street scenes. In *CVPR*, pages 2823–2830, 2012. 1, 2
- [17] J. Y. Guillemaut and A. Hilton. Joint Multi-Layer Segmentation and Reconstruction for Free-Viewpoint Video Applications. *IJCV*, 93:73–100, 2010. 6, 7
- [18] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik. *Learning Rich Features from RGB-D Images for Object Detection and Segmentation*, pages 345–360. 2014. 2
- [19] C. Hane, C. Zach, A. Cohen, and M. Pollefeys. Dense semantic 3d reconstruction. *TPAMI*, page 1, 2016. 2
- [20] B. Hariharan, P. A. Arbelaz, R. B. Girshick, and J. Malik. Hypercolumns for object segmentation and fine-grained localization. In *CVPR*, pages 447–456, 2015. 2
- [21] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. In *ICCV*, 2017. 1, 2, 5, 7
- [22] Y. Huang, F. Bogo, C. Lassner, A. Kanazawa, P. V. Gehler, J. Romero, I. Akhter, and M. J. Black. Towards accurate marker-less human shape and pose estimation over time. In *3DV*, 2017. 2
- [23] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *TPAMI*, 36(7):1325–1339, jul 2014. 6
- [24] M. Kazhdan, M. Bolitho, and H. Hoppe. Poisson surface reconstruction. In *Eurographics Symposium on Geometry Processing*, pages 61–70, 2006. 6
- [25] A. Kendall, Y. Gal, and R. Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. *CoRR*, abs/1705.07115, 2017. 1, 2
- [26] H. Kim, J. Guillemaut, T. Takai, M. Sarim, and A. Hilton. Outdoor Dynamic 3-D Scene Reconstruction. *T-CSVT*, 22(11):1611–1622, 2012. 6
- [27] A. Kundu, Y. Li, F. Dellaert, F. Li, and J. M. Rehg. Joint semantic segmentation and 3d reconstruction from monocular video. In *ECCV*, volume 8694, pages 703–718, 2014. 2
- [28] A. Kundu, V. Vineet, and V. Koltun. Feature space optimization for semantic video segmentation. In *CVPR*, pages 3168–3175, 2016. 2
- [29] F. Langguth, K. Sunkavalli, S. Hadap, and M. Goesele. Shading-aware multi-view stereo. In *ECCV*, 2016. 6, 7
- [30] E. Larsen, P. Mordohai, M. Pollefeys, and H. Fuchs. Temporally consistent reconstruction from multiple video streams using enhanced belief propagation. In *ICCV*, pages 1–8, 2007. 1, 2
- [31] T.-Y. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014. 2
- [32] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 2
- [33] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 60:91–110, 2004. 5
- [34] B. Luo, H. Li, T. Song, and C. Huang. Object segmentation from long video sequences. In *ACM Multimedia*, pages 1187–1190, 2015. 1, 2
- [35] M. Mostajabi, P. Yadollahpour, and G. Shakhnarovich. Feed-forward semantic segmentation with zoom-out features. In *CVPR*, pages 3376–3385, 2015. 2

- [36] A. Mustafa and A. Hilton. Semantically coherent co-segmentation and reconstruction of dynamic scenes. In *CVPR*, 2017. 1, 2, 3, 6, 7
- [37] A. Mustafa, H. Kim, J.-Y. Guillemaut, and A. Hilton. Temporally coherent 4d reconstruction of complex dynamic scenes. In *CVPR*, 2016. 1, 2, 5
- [38] A. Mustafa, H. Kim, and A. Hilton. 4d match trees for non-rigid surface alignment. In *ECCV*, 2016. 7, 8
- [39] A. Mustafa, M. Volino, J.-Y. Guillemaut, and A. Hilton. 4d temporally coherent light-field video. In *3DV*, 2017. 5, 6
- [40] R. A. Newcombe, D. Fox, and S. M. Seitz. Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time. *CVPR*, pages 343–352, 2015. 5
- [41] A. Roussos, C. Russell, R. Garg, and L. Agapito. Dense multibody motion estimation and reconstruction from a handheld camera. In *ISMAR*, 2012. 2
- [42] L. Sevilla-Lara, D. Sun, V. Jampani, and M. J. Black. Optical flow with semantic segmentation and localized layers. In *CVPR*, pages 3889–3898, 2016. 1, 2
- [43] O. Sorkine and M. Alexa. As-rigid-as-possible surface modeling. In *SGP*, pages 109–116, 2007. 4
- [44] T. Tani, Y. Matsushita, Y. Sato, and T. Naemura. Continuous 3D Label Stereo Matching using Local Expansion Moves. *TPAMI*, 40(11):2725–2739, 2018. 6, 7
- [45] M. W. Tao, J. Bai, P. Kohli, and S. Paris. Simpleflow: A non-iterative, sublinear optical flow algorithm. *Computer Graphics Forum (Eurographics 2012)*, 31(2), May 2012. 4, 5
- [46] D. Tome, C. Russell, and L. Agapito. Lifting from the deep: Convolutional 3d pose estimation from a single image. In *CVPR*, July 2017. 1, 4
- [47] D. Tomè, M. Toso, L. Agapito, and C. Russell. Rethinking pose in 3d: Multi-stage refinement and recovery for markerless motion capture. In *3DV*, 2018. 2, 4
- [48] Y.-H. Tsai, G. Zhong, and M.-H. Yang. Semantic co-segmentation in videos. In *ECCV*, pages 760–775, 2016. 2, 7
- [49] A. O. Ulusoy, M. J. Black, and A. Geiger. Semantic multi-view stereo: Jointly estimating objects and voxels. In *CVPR*, 2017. 2
- [50] V. Vineet, O. Miksik, M. Lidegaard, M. Nießner, S. Golodetz, V. A. Prisacariu, O. Kähler, D. W. Murray, S. Izadi, P. Perez, and P. H. S. Torr. Incremental dense semantic stereo fusion for large-scale semantic scene reconstruction. In *ICRA*, 2015. 2
- [51] D. Vlasic, I. Baran, W. Matusik, and J. Popović. Articulated mesh animation from multi-view silhouettes. *ACM Trans. Graph.*, 27(3), Aug. 2008. 6
- [52] C. Vogel, K. Schindler, and S. Roth. 3d scene flow estimation with a piecewise rigid scene model. pages 1–28, 2015. 6, 7, 8
- [53] A. Wedel, T. Brox, T. Vaudrey, C. Rabe, U. Franke, and D. Cremers. Stereoscopic scene flow computation for 3d motion understanding. *IJCV*, 95(1):29–51, 2011. 1
- [54] P. Weinzaepfel, J. Revaud, Z. Harchaoui, and C. Schmid. Deepflow: Large displacement optical flow with deep matching. In *ICCV*, pages 1385–1392, 2013. 7, 8
- [55] F. Xia, P. Wang, X. Chen, and A. L. Yuille. Joint multi-person pose estimation and semantic part segmentation. In *CVPR*, 2017. 2
- [56] J. Xie, M. Kiefel, M.-T. Sun, and A. Geiger. Semantic instance annotation of street scenes by 3d to 2d label transfer. In *CVPR*, 2016. 1, 2
- [57] J. Xu, R. Ranftl, and V. Koltun. Accurate Optical Flow via Direct Cost Volume Processing. In *CVPR*, 2017. 7, 8
- [58] A. Zanfir and C. Sminchisescu. Large displacement 3d scene flow with occlusion reasoning. In *ICCV*, 2015. 1
- [59] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *CVPR*, 2017. 7
- [60] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. Torr. Conditional random fields as recurrent neural networks. In *ICCV*, 2015. 2, 7