# iMapper: Interaction-guided Joint Scene and Human Motion Mapping from Monocular Videos

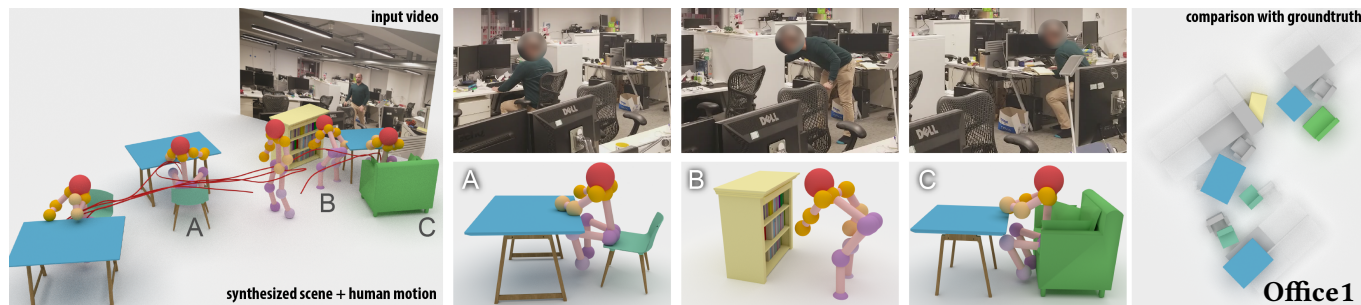ARON MONSZPART, PAUL GUERRERO, DUYGU CEYLAN, ERSIN YUMER, NILOY J. MITRA

Fig. 1. We present iMapper, a method that reasons about the interactions of humans with objects, to recover both a plausible scene arrangement and human motions, that best explain an input monocular video (see inset). We fit characteristic interactions called *scenelets* (*e.g.,* A, B, C) to the video and use them to reconstruct a plausible object arrangement and human motion path (left). The key challenge is that reliable fitting requires information about occlusions, which are unknown (*i.e.,* latent). (Right) We show an overlay (from top-view) of our result over manually annotated groundtruth object placements. Note that object meshes are placed based on estimated object category, location, and size information.

A long-standing challenge in scene analysis is the recovery of scene arrangements under moderate to heavy occlusion, directly from monocular video. While the problem remains a subject of active research, concurrent advances have been made in the context of human pose reconstruction from monocular video, including image-space feature point detection and 3D pose recovery. These methods, however, start to fail under moderate to heavy occlusion as the problem becomes severely under-constrained. We approach the problems differently. We observe that people *interact similarly in similar scenes*. Hence, we exploit the correlation between scene object arrangement and motions performed in that scene in both directions: first, typical motions performed when interacting with objects inform us about possible object arrangements; and second, object arrangements, in turn, constrain the possible motions.

We present iMapper, a data-driven method that focuses on identifying human-object interactions, and *jointly* reasons about objects and human movement over space-time to recover both a plausible scene arrangement and consistent human interactions. We first introduce the notion of *characteristic interactions* as regions in space-time when an informative human-object interaction happens. This is followed by a novel *occlusion-aware matching* procedure that searches and aligns such characteristic snapshots from an interaction database to best explain the input monocular video. Through extensive evaluations, both quantitative and qualitative, we demonstrate that iMapper significantly improves performance over both dedicated state-of-the-art scene analysis and 3D human pose recovery approaches, especially under medium to heavy occlusion.

Additional Key Words and Phrases: shape analysis, interaction, scene layout, 3D pose estimation, monocular video, occlusion

## 1 INTRODUCTION

Digitizing the physical world is critical for many emerging fields such as virtual and augmented reality, smart home systems, or robotics. Such applications require access to not only reconstructions of the physical spaces, but also an *understanding* of the common human actions performed in such spaces. For example, our future personal robot assistants should not only know the object layout of our lounges, but also our working habits in these spaces.

Traditionally, researchers have tackled scene estimation and human performance capture as two separate problems. On the one hand, scene reconstruction methods such as Kinect Fusion [Newcombe et al. 2011] and Bundle Fusion [Dai et al. 2017b] can produce high-quality static indoor reconstructions, and the likes of DynamicFusion [Newcombe et al. 2015] can capture non-rigidly deforming scenes. These methods require the sensor to be manually moved to see around occlusions making the capture process cumbersome. Moreover, the process needs to be repeated each time scene objects are moved. On the other hand, there are various options for reconstructing 3D human performances, either using multiple sensors [von Marcard et al. 2017] or monocular video input, based on a CNN pose regressor along with kinematic skeleton fitting [Mehta et al. 2017b]. These methods assume performances to be free from object-induced occlusions (see Figure 2 and Section 8).

While indoor scene configurations can be extremely rich and diverse, we observe that a large fraction of them are linked by a common thread — *they are regularly inhabited by humans*. Moreover, in similar scene configurations, social behavioral rules lead to humans typically performing similar actions (*cf.,* [Krasner 2013]). Examples of such actions include sitting on sofas, picking up books from shelves, or walking around obstacles. Hence, instead of tackling scene modeling and human reconstruction as separate problems, we propose to *jointly* estimate both plausible scene layouts and consistent human performances from monocular video.

A fundamental challenge in 3D reconstruction that can benefit from such an approach is working with scenes under *occlusion*. A successful solution needs to tackle two problems: (i) How to hallucinate information about partially or fully hidden scene objects
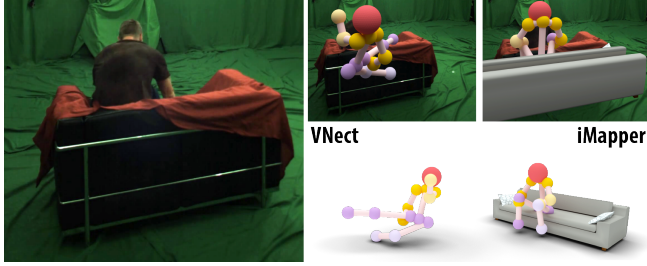
Fig. 2. Comparison of state-of-the-art 3D human pose detection from monocular video VNect [Mehta et al. 2017b] with iMapper. Note how VNect breaks down in regions of occlusion (VNect was not designed to handle occlusions), while our approach continues to produce plausible results because of explicit occlusion detection and handling. Bottom row shows the recovered human pose from another camera view for better visibility.

in monocular input? (ii) How to determine human performance that is strongly occluded by various (unknown) scene objects? It is believed that, as humans, we focus on the interactions of actors with the objects in the scene (referred to as *anticipation* in [Neisser 1976]), instead of separately identifying objects and human performances. Our experience allows us to compensate for missing information in *both* objects and performances, using subtle hints arising from their interaction. For example, in the video for the scene shown in Figure 1, we can 'see' the person walk behind the desk and sit down – from that, we can imagine both the person's sitting pose over time *and* the location of the unseen chair; similarly, for the person picking up an object from the shelf (see also supplementary video).

We propose iMapper, a data driven method that accomplishes a similar feat by jointly reasoning about interaction detection and object placement. We rely on humans behaving consistently near similar object configurations at different times. Hence, as a data-prior, we use a database of 3D interactions (extracted from the PiGraph dataset [Savva et al. 2016]) between humans over time and local object configurations that we call *scenelets*. Human actions are easier to robustly detect than objects using state-of-the-art methods, and the scenelets in our dataset give us typical object configurations associated with actions. Additionally, a scenelet also tells us which parts of human actors are likely to be occluded when viewed from different angles.

Given a monocular video, we fit such scenelets to each actor's visible 2D joints, and ensure that the joints missing in the video match the occlusions we expect from the candidate scene arrangement. However, we have a cyclic dependency: the recovered scene objects depend on the quality of the human pose estimates; while, the estimation of human poses, in turn, needs to know which 2D feature points remain unoccluded by the discovered scene objects. Hence, we first generate candidate scenelets matching potentially informative segments of the video, and then solve a global selection problem to extract a consistent set of scenelets fitted to the observed interactions. This simultaneously produces an unoccluded 3D human performance and a set of static objects that plausibly fit the observed interactions (see Figure 1).

We extensively evaluated iMapper on a range of scenes of varying complexity and occlusion. We provide both qualitative and quantitative evaluation on real data demonstrating that our simultaneous estimation (*i.e.,* scene layout and human pose) improves upon dedicated state-of-the-art methods that treat the two problems independently. In summary, our main contributions are: (i) proposing the first method that jointly reasons about 3D scene modeling and plausible human object interactions given monocular video input of scenes with occlusion; (ii) detecting which human motions in a sequence are informative, when such *characteristic* pose sequences occur, and matching them to the scenelet database while accounting for (unknown) occlusion; and (iii) combining detected scenelets into a consistent 3D scene and human performance by a novel optimization that reasons about several cues including number of objects and object categories, position and orientation of objects, realistic room layouts, and compatible human movements.

## 2 RELATED WORK

Our system is related to prior work on scene analysis and synthesis, human-centric shape analysis, as well as human pose estimation. We now discuss selected papers in these categories to better position our approach.

**Scene analysis and synthesis.** With the advances in acquisition technologies, several large scale indoor reconstruction datasets have been created [Chang et al. 2017; Dai et al. 2017a]. Such datasets, together with the availability of 3D scene collections, are attracting more attention to 3D scene analysis. Several previous works focus on analyzing inter-object relationships [Fisher et al. 2011; Huang et al. 2016; Zhao et al. 2016] and hierarchical grammars [Liu et al. 2014] given a 3D scene collection. Xu et al. [2014] introduce the concept of *focal points* that represent frequently occurring sub-arrangements of objects across heterogeneous collections of scenes. For synthesis purposes, one line of related work focuses on modeling scenes from single images [Izadinia et al. 2017; Poirson et al. 2016; Satkin and Hebert 2013] or RGBD scans [Nan et al. 2012; Shao et al. 2012] by matching individual 3D objects. Other work [Chen et al. 2014; Schwing et al. 2013] further explores spatial relationships between objects. Given a room layout and object instances, Yeh et al. [2012] provide a Markov Chain Monte Carlo based algorithm to synthesize object arrangements that take spatial constraints into consideration. Fisher et al. [2012] present an example-based synthesis approach that uses a probabilistic graphical model to encode object relationships. Del Pero et al. [2013] represent 3D objects as a collection of primitives to reconstruct more accurate room geometry from images.

**Human-centric scene synthesis.** Recent work in scene synthesis incorporates human actions into scene analysis (i.e., object arrangements, scene layout) to infer where specific actions can take place [Savva et al. 2014] or where a new object may be placed [Jiang et al. 2016] in a scene. Frank et al. [2015] recognize certain human actions to insert 3D objects into a SLAM-based reconstruction output. Ma et al. [2016] model both object-object and object-human actions to refine an input 3D scene. The recent work of Fu et al. [2017b] analyzes a collection of floor plans annotated with possible human actions and 3D objects to guide the synthesis of a scene given an empty room layout and a few object categories. In contrast to these
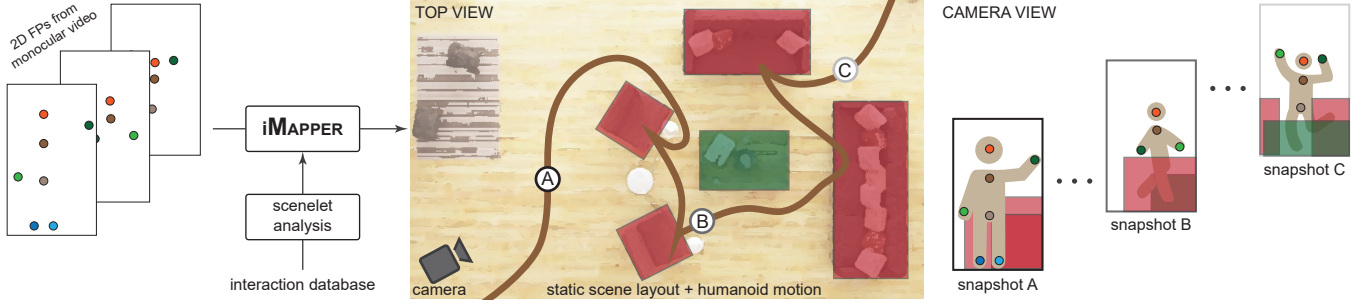
Fig. 3. Starting from a monocular video and an interaction database, iMapper produces a *plausible* and *consistent* static scene layout and humanoid 3D motion path (shown as brown curve) over time. In a plausible solution, the projection of the discovered human motion under the input camera agrees with per frame 2D feature points in the unoccluded parts. In a consistent solution, we expect the scene and human motions to have matching interactions in space-time. By jointly reasoning over occlusions and interactions, we are able to improve both scene layout (*e.g.,* depth estimation) and pose estimation (*e.g.,* robustness under latent occlusion). For example, in this illustration, the inferred woman's pose remains unoccluded in 'snapshot A', while she is partly occluded in 'snapshot B' and 'snapshot C' by inferred scene objects.

approaches, our work does not assume any initial input about the scene geometry or layout.

One of the earlier works that uses motion cues to reason about the scene geometry is the work of Brostow and Essa [1999]. Given an input video, this work represents motion in a general sense by segmenting each video frame into blobs that are classified as static or active. This classification is then used to extract depth layers for the scene yielding a 2.5D representation. Fouhey et al. [2012] combine human pose estimates with appearance and other geometric cues to estimate the room layout and free space in a single image or a time-lapse video. Compared to these approaches, we focus on a fundamentally different problem: instead of recovering a faithful geometry reconstruction, our goal is to generate a plausible scene as well as a human motion that that can explain the input video.

One of the previous systems most closely related to our approach is the work of Fisher et al. [2015] that utilizes a scene template computed from a given 3D scan, together with an activity classifier, to model plausible layouts. There are two main differences between our approach and this work. First, they require an initial 3D scan of the scene to be provided in order to predict possible actions that can take place in the scene. In contrast, the input to our system is the human motion sequence with no prior knowledge of the scene geometry. Second, instead of operating at the level of high-level activity labels, which rely on a pre-defined set of activity classes, we use *scenelets, i.e.,* short human pose-object interaction sequences, where any pose example provides a cue for the scene geometry, *i.e.,* presence of a specific object or void space. While we bootstrap the scene population process by identifying *characteristic poses* that imply the presence of typical objects, any input pose provides constraints for identifying the occupied and void regions of the scene.

The recent work of Savva et al. [2016] analyzes interaction snapshots, *i.e.,* action and pose labeled RGBD sequences to learn *prototypical interaction graphs (PiGraphs)* which link the attributes of the human pose to the surrounding object geometries. They show how PiGraphs can be utilized to generate scenes that correspond to static interaction snapshots (*e.g.,* lie on bed). In contrast, we focus

on observing a dynamic motion sequence that consists of different actions. We combine *scenelets* which depict short sequences of human actions to synthesize a scene that is in agreement for the whole duration of the motion. Finally, Kang et al. [2017] focus on a similar goal of scene synthesis that explores motion cues. However, their input is a 3D human motion sequence free of occlusions. In contrast, we use monocular videos with moderate to severe occlusions as input. Thus, our goal is not only to synthesize plausible scene objects but also recover the 3D human motion where occlusions happen.

**Human-centric shape analysis.** Earlier work that uses observations of how humans interact with objects focuses on tasks such as object and event recognition [Delaitre et al. 2012; Gupta et al. 2009; Wei et al. 2013] and action detection [Yao et al. 2011]. Kim et al. [2014] propose a shape analysis tool based on a human-object affordance model that can be used for many applications including correspondence estimation, shape retrieval, and view selection. In a followup project, Fu et al. [2017a] utilize a similar model to generate new objects by combining functional parts of existing objects. The recent work of Pirk et al. [2017b], introduces the concept of *interaction landscapes* which provides a descriptor of an object based on the type of interactions it can be involved in. In a follow-up paper [Pirk et al. 2017a], they extend this notion to compute descriptors for individual interactions. More recently, Gkioxari et al. [2018] predict human-verb-object instances from a single image to characterize human-object interactions. In our work, on the other hand, instead of focusing on individual object-based inference tasks, we use observed human motion to generate plausible scenes from a diverse and rich set of possibilities.

**Human pose estimation.** With the recent success of deep learning, we have seen advances both in 2D [Cao et al. 2017; Insafutdinov et al. 2016; Newell et al. 2016; Toshev and Szegedy 2014; Wei et al. 2016] and 3D pose estimation [Huang et al. 2014; Rogez et al. 2018; Tekin et al. 2016; Tome et al. 2017; Zhou et al. 2016]. In particular, the recent VNect system [Mehta et al. 2017b] demonstrates state-of-the-art results for global pose estimation from monocular video. Many of these approaches, however, do not specifically tackle the occlusion problem. The few existing works that focus on predicting pose in
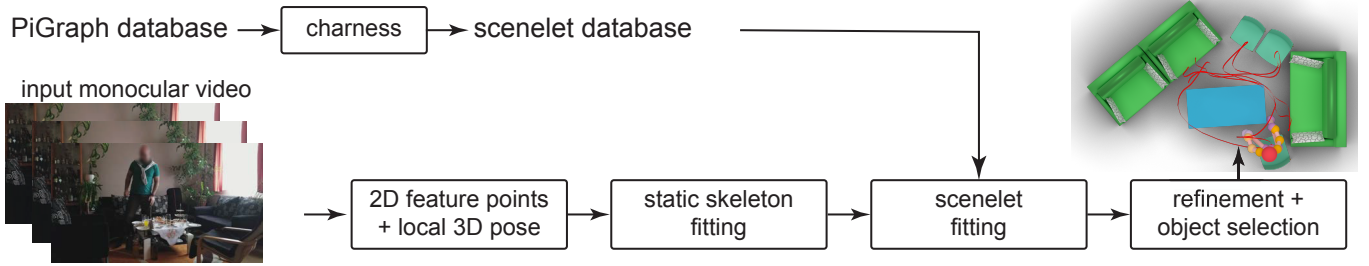
Fig. 4. Overview of our approach. Please refer to the text for details of the individual steps.

the presence of occlusions either consider only 2D pose [Fu et al. 2015], or represent 3D human pose as sparse linear combinations of known 3D poses and recover these blend weights from partially observed data [Huang and Yangc 2009]. Such approaches, however, are limited to image input, do not consider the temporal dynamics of the human motion, and do not reason explicitly about occlusion arising due to human-object interactions. In contrast, given an initial pose estimate (acquired by off-the-shelf 3D pose estimators that do not consider global positioning of the detected skeleton), iMAPPER jointly reasons about scene geometry and human pose to synthesize both plausible scenes and human motion even in case of moderate to heavy occlusions. In Section 8, we show comparisons of the 3D human motion recovered by our method to the recent 3D pose estimation methods.

## 3 OVERVIEW

The input to our method[1] is a monocular video showing a person interacting with objects. Our goal is to synthesize human performance that fits this input video and a static object arrangement that is plausible and consistent with the inferred human performance.

When watching performance of a human actor in a scene, there are several cues that one can use to recover plausible explanations for objects and performance sequences, even under partial occlusion. Interactions with objects, for example, give cues about both the objects and the motions that are part of the interaction. Walking gives cues about occluded empty space in a room, while occlusion of joints give cues about the location of objects relative to the actor. Such cues, when combined with prior knowledge of typical human-object interactions, help recover plausible human performance and object arrangements (see Figure 3) as described below.

To represent prior information (see Section 4) about typical interactions with objects, we build a dataset using PiGraphs of human-object interactions called *scenelets*. Each scenelet contains a short motion clip and a set of static objects in close proximity to the motion. These scenelets capture relationships between the motions of an actor and object arrangements. Additionally, we get prior information about typical skeleton poses from a selection of pre-trained models of human poses [Rogez et al. 2018; Tome et al. 2017].

We then synthesize an output by fitting selection of these models to each part of the video, optimizing an energy function (see Section 6) that quantifies the consistency of the fitted models with the video and with each other. By consistency, we measure how well

the fitted models explain the presence or absence of skeleton joint detections in each video frame, and on other plausibility criteria, such as path smoothness and intersection avoidance. We formulate the problem as a global optimization (see Section 5).

We decompose the optimization into several simpler sub-problems to obtain a robust initialization before the final optimization. We start by fitting scenelets in our dataset to a set of regularly spaced time intervals in the video. Fitting is done independently for each scenelet to avoid a combinatorial explosion of possible scenelet combinations. We continue to fit skeletons to all frames not covered by scenelets. This initialization provides sufficient information to commit to a subset of scenelets that constitute our synthesized scene. Finally, we optimize the placement of all chosen scenelets and skeletons with the full fitting energy, including interactions between all fitted models (see Section 7).

## 4 HUMAN POSES AND SCENELET PRIORS

In occluded parts of a video sequence, iMAPPER recovers information about objects and actor motions by fitting interaction models to the video sequence. In unoccluded sequences, we fit static skeleton poses instead. The models we fit to the video sequence represent our prior information about valid human poses and interactions. Next, we will describe these models.

### 4.1 Human Pose Prior

The space of valid static poses is defined with the 3D skeleton pose model described by Tome et al. [2017] or Rogez et al. [2018]. Such models are trained on a large dataset of poses, carefully removing factors such as rotation in the ground plane and left-right symmetry. The authors also take care to capture less frequently occurring poses, giving this model a good coverage of the valid human pose space. Please refer to the original papers for details.

### 4.2 Scenelet Prior

We use the PiGraphs dataset [Savva et al. 2016] to model interactions between humans and objects. The original dataset contains a set of scenes with a commodity depth sensor, each containing a human performance and a set of labeled objects captured. Labels group objects into a small set of categories, such as tables, sofas, chairs, and bookshelves. From this dataset, we extract short sequences $\mathcal{L}_1, \ldots, \mathcal{L}_m$ showing interactions between the actor and objects. We call such short sequences *scenelets*.

Each scenelet consists of a short motion clip with known 3D joint positions and a set of static objects. We denote with $s_k^{lt}$ the location

---

[1]Project code/data will be released for research use.

of skeleton joint $k$ in frame $t$ of scenelet $l$. Objects $O^l = \{o_1^l \dots o_n^l\}$, of scenelet $l$ are defined by a placement $p$, a rough approximation of their geometry $\kappa$, and a label $b$; *i.e.,* each object is encoded as triplet $o = \{p, \kappa, b\}$. We assume objects can only rotate around the up direction, giving us four degrees of freedom for the placement of an object: $p = (x, y, z, \theta)$, where $x$, $y$, and $z$ are the location, and $\theta$ the orientation of the object. Similar to the original dataset, we approximate geometry of objects by unions of cuboids; and the label $b_i^l$ describes the object type from the predefined set of categories. Both motion clip and objects are stored in the local coordinate frame of the scenelet, defined by the pelvis location and the forward-facing direction of the skeleton in the center frame of the motion clip.

**Scenelet parameterization.** When constructing scenelets, we make a design choice regarding the length of the motion clip used for each scenelet based on the following considerations.

First, the speed at which an interaction is performed should not affect the contents of a scenelet. For example, if a scenelet captures a fast 'sitting-down' performance, then a slower version of the same interaction should also be captured by a single scenelet. This property is necessary to ensure that interactions captured by scenelets are comparable.

Second, scenelets should represent interactions that are local in space. Keeping spatial locality simplifies the subsequent optimization, since it reduces the number of potential interactions between scenelets. We have found that using scenelets where the skeleton (*e.g.,* the pelvis joint) traverses constant arc length satisfies both time invariance and locality. To reduce the effect of noise, we compute the arc length on a smoothed pelvis trajectory, using 10 iterations of a moving average with an arc length radius of 1 cm.

**Scenelet construction.** The PiGraphs dataset describes human performance with a set of 16 joint locations per frame. We start by sampling the center of each scenelet's motion clip at regularly spaced intervals on the arc length of the pelvis joint's trajectory. The start and end of the motion clip in each scenelet is then defined as this center point plus/minus half the scenelet length. The objects of the scenelet are chosen as the subset of objects roughly within arm's reach of the actor at any point in the motion clip, *i.e.,* the objects the actor can potentially interact with. In our test, we pick objects within 1 m radius when projected to the ground plane.

**Scenelet descriptors.** In order to compare two scenelets and compute their distance, we define two descriptors for each scenelet: a *motion descriptor* $\Psi$ and an *object descriptor* $\Phi$.

The *motion descriptor* $\Psi$ compactly describes the motion clip of a scenelet with a fixed-length vector. It is the concatenation of a fixed number of static pose descriptors $\Psi := (\psi_1, \dots, \psi_k)$, sampled evenly over the motion clip. In our experiments, we use $k = 15$ samples. Static pose descriptors are based on 14 robust joint-line distances (see supplemental for details) as suggested in Zhang et al. [2016]. Distances between these descriptors are defined using a weighted $L_2$ distance, assigning more weight to center frames. The descriptors $\psi_i$ should evenly cover the motion, and similar to the length of a scenelet, they should be invariant to the speed of the motion. We evenly distribute these samples along the trajectory of the motion clip in a 17D space of the combined pose descriptor and global skeleton location (taken to be the 3D location of the pelvis).
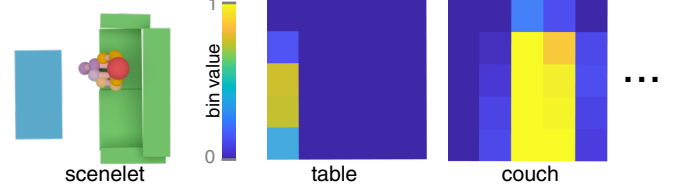


Fig. 5. Object descriptors compactly represent the object arrangement of a scenelet. One $5 \times 5$ histogram per category stores the layout of objects of this category relative to the scenelet center.

The *object descriptor* $\Phi$ captures the object arrangement as a set of histograms. We store one histogram per object category. The histograms capture the 2D placement of objects, projected to the ground plane. Our histograms are $5 \times 5$ square grids (with center poses facing the forward direction) and each bin $\Phi_j$ describes to what extent any object of the same category in the scenelet is located in this bin. We define the value in a bin as the maximum coverage of the bin by any object. To handle both objects smaller and larger than the bin, we normalize by the smaller of either the bin area or the object area:

$$\Phi_j := \max_i \left\{ A\left(\Lambda(o_i) \cap \phi_j\right) / \min\left(A(\Lambda(o_i)), \ A(\phi_j)\right) \right\},$$

where $\Lambda(o_i)$ is the projection of object $o_i$ to the ground plane, $\phi_i$ is the part of the ground plane covered by bin $i$ and $A(x)$ is the area of $x$. Figure 5 shows an example of an object descriptor.

**Charness.** The characteristicness, or *charness* for short, of a bin in an object descriptor, describes how typical an activation of this bin is for similar motion clips. For example, for a sitting-down motion, a couch or a chair at the center bin of the histogram will have high charness. This charness score helps to distinguish between objects that are related to a given interaction, and objects that are near the motion, but unrelated to the interaction. The charness of an object descriptor bin is computed as a weighted average of that bin's activation over similar motion clips, where similarity is defined with a Gaussian kernel in the space of motion descriptors:

$$h_j^l := \frac{\sum_{k=1}^m \Phi_j^k \ \mathcal{G}(d(\Psi_k, \Psi_l)|0, \sigma)p_k^{-1}}{\sum_{k=1}^m \mathcal{G}(d(\Psi_k, \Psi_l)|0, \sigma)p_k^{-1}},$$

where $h_j^l$ is the characteristicness of bin $j$ in scenelet $l$, $\Phi_j^l$ is the bin value of scenelet $l$, $\mathcal{G}$ is a Gaussian kernel taken over the distance $d$ between the motion clip descriptors defined earlier (we empirically set $\sigma = 13$), and $p_l$ is the *density* of scenelets at the origin of scenelet $l$. We measure the density of scenelets in the space of the original PiGraph scene that scenelet $l$ was obtained from. It is defined as the spatial density of the origins of all scenelets that were obtained from this scene. The intuition behind dividing by this density is to remove the bias we would otherwise introduce due to multiple scenelets taken from nearby parts of the scene. Note that multiple scenelets showing the same part of the same scene are correlated and would bias the charness towards this particular scene arrangement. Finally, we define the charness of a scenelet as the maximum bin characteristicness:
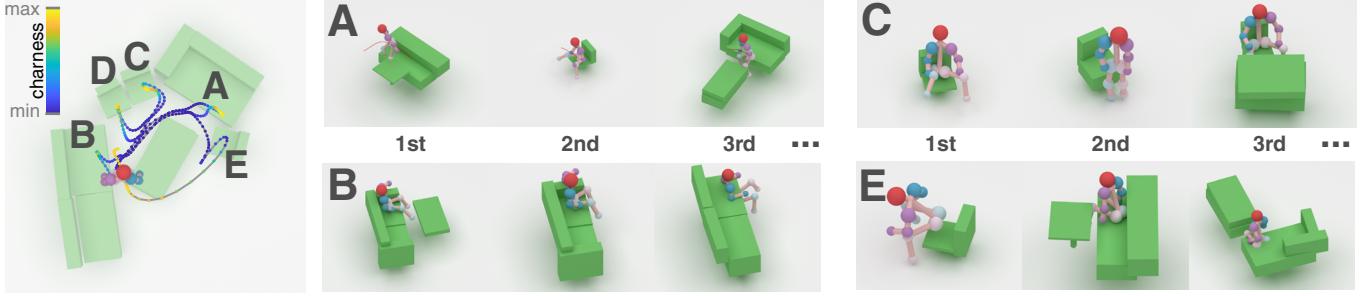
$$H^l := \max_j (h_j^l).$$

Fig. 6. [LvRoom] Fitting scenelets to characteristic sequences. Scenelets are fit to charness maxima in the video sequence. On the left we show the charness of the human's pelvis trajectory in one of our synthesized scenes. The scenelets on the right are the best 3 candidates for an independent fit to each of the gaps, without interaction between scenelets. Note that most of these candidates plausibly fit their video sequence.

Scenelets with high charness are likely to contain interactions with objects, since they have objects within interaction range that are typical for the scenelet's motion. We will use the charness of a scenelet to determine which sequences of a video are likely to contain interactions between the actor and objects (see Figure 7).
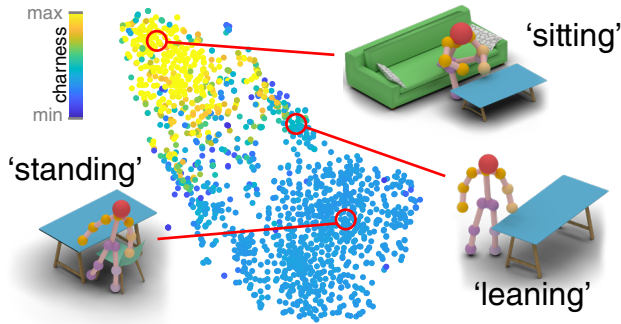


Fig. 7. Characteristicness of our scenelet dataset over pose descriptor space. We show a t-SNE embedding with our descriptor distance. Warmer colors denote more characteristic scenelets. The close-ups show scenelets with a standing sequence, a sitting sequence, and a leaning sequence. Standing sequences have low charness because they are not specific to the neighboring objects.

## 5 FITTING SCENELETS AND SKELETONS

The pose space and the scenelets discussed in the previous section give us strong priors for human poses, and interactions with objects, respectively. In occluded parts of a video sequence, we can fit these models to the sparse set of observations, giving us likely explanations for the parts that could not directly be observed. In this section we describe this fitting problem as a search in a high-dimensional parameter space, and in the next section we define an energy in this space that we minimize to get an optimal fit.

In each frame of the input video, we detect 2D image-space skeleton of the actor, consisting of $n_j$ joint locations using an off-the-shelf pose detector. In our experiments we use either (i) 2D keypoints detected by CPM [Wei et al. 2016] and grouped based on the grouping heuristic of Tome et al. [2017] or (ii) LCR-Net++ [Rogez et al. 2018].

Given a video with $n_v$ frames, for each joint $k$ detected in frame $t$, $u_k^t \in \mathbb{R}^2$ denotes its location and with confidence $c_k^t \in [0, 1]$. While confidences are directly provided by the method of CPM, when using LCR-Net++ we compute them based on the per-joint pose proposal variance (see Appendix B for detail).

Our goal is to synthesize a scene consisting of 3D joint locations $q_k^t \in \mathbb{R}^3$ for each video frame, describing the human performance, and a set of objects $O = \{o_1, \ldots o_{n_o}\}$. The 3D joint locations at each frame are obtained by fitting either a scenelet or a static skeleton to the 2D joint detections in the video, and objects are taken from the fitted scenelets. Which of the two models we fit to a given part of the video depends on two factors: the estimated amount of joint occlusion observed in the video (i.e., the confidence of the joint detection signal) and the estimated probability of object interactions.

**(i) Fitting scenelets.** Video sequences that contain interactions with objects or occluded human performance, are modeled with scenelets allowing allow us to both populate the scene with objects involved in interactions and explain occlusions of joints due to these objects. Thus, joint occlusions can help in both choosing and placing the scenelets: for a given video sequence, we would like to choose and place a scenelet such that the scenelet objects explain the joint occlusions observed in the video sequence.

We start by modeling the assignment of scenelets in our dataset to time intervals in the video. Given a video with $n_v$ frames and a dataset with $m$ scenelets, only a single scenelet can start at any frame of the video. The scenelet assignment can therefore be expressed with a binary matrix $\mathcal{X} \in \{0, 1\}^{m \times n_v}$

$$\mathcal{X}_{lt} = \begin{cases} 1 & \text{if scenelet } l \text{ starts at video frame } t \\ 0 & \text{otherwise.} \end{cases}$$

The constraint that scenelets should not overlap in time can then be formulated as

$$\eta_t = \sum_l \sum_{i=1}^{\max(t, n_l)} \mathcal{X}_{l(1+t-i)} \leq 1, \quad t = 1 \ldots n_v,$$

where $\eta_t$ is the number of scenelets assigned to frame $t$, and $n_l$ is the number of frames in scenelet $l$. Since only a single scenelet can start at any frame $t$, we model scenelet placement with one set of parameters per frame $P = \{P_1 \ldots P_{n_v}\}$, where $P_t = (x, y, z, \theta)$ is the placement of the scenelet starting at $t$, with $x$, $y$, and $z$ the

location and $\theta$ the orientation of the scenelet. The 3D joint locations $\hat{q}_k^t$ in video sequences covered by scenelets can then be defined as a function of the placement $P$ and the scenelet assignment $\mathcal{X}$:

$$\hat{q}_k^t(P, \mathcal{X}) := \sum_l \sum_{i=1}^{\max(t, n_l)} \mathcal{X}_{l(1+t-i)} \, T(P_{(1+t-i)}) s_k^{li}, \qquad (1)$$

where $s_k^{li}$ is the 3D position of joint $k$ in frame $i$ of scenelet $l$, and $T(P_t)$ is the transformation to placement $P_t$.

Finally, the objects in the scene are obtained from all scenelets that have been assigned to the scene:

$$O(P, \mathcal{X}) := \bigcup_{\{(l,t) \mid \mathcal{X}_{lt}=1\}} T(P_t, O^l),$$

where we denote with $T(P, O)$ the transformation of objects in $O$ to the placement $P$, i.e., $T(P, O^l) = \{(T(p), \kappa, b) \mid (p, \kappa, b) \in O^l\}$.

**(ii) Fitting static skeletons.** For parts of the video that contain an unoccluded human performance without object interactions, we can fit static skeletons to each frame. Since the number of degrees of freedom for human poses is smaller than for human-object interactions, the space of possible human poses can be covered more accurately than the space of possible human-object interactions. Thus, fitting static skeletons to the video gives us better performance in unoccluded sequences that do not contain interactions.

The aforementioned 3D pose reconstruction methods [Rogez et al. 2018; Tome et al. 2017] retrieve the best matching 3D skeleton pose for a given frame. This pose is defined in the *local* space of the skeleton and does not give us the placement of the skeleton in the scene. We fit the retrieved 3D skeleton to our video by optimizing the 3D placement of the skeleton, using the fitting energy described later in Section 6.

Skeletons are only fitted to frames that do not have any scenelet assignment. The joint locations $\check{q}_k^t$ for video sequences that are not covered by scenelets are then defined as:

$$\check{q}_k^t(P, \mathcal{X}) = (1 - \eta_t) \, T(P_t) a_k^t, \qquad (2)$$

where the first term is only non-zero if no scenelet is assigned to frame $t$, and $a_k^t$ is the local skeleton pose computed by using Tome et al. or LCR-Net++, $P_t = (x, y, z, \theta)$ is the placement of the skeleton in frame $t$, and $T(P_t)$ is the transformation to placement $P_t$. Combining Equations 1 and 2, we define the location of any joint $q_k^t$ in the video as:

$$q_k^t(P, \mathcal{X}) = \hat{q}_k^t(P, \mathcal{X}) + \check{q}_k^t(P, \mathcal{X}). \qquad (3)$$

In the following, we will omit the explicit dependence of $q_k^t(P, \mathcal{X})$ and $o_i(P, \mathcal{X})$ on $P$ and $\mathcal{X}$ for a less cluttered notation.

## 6 FITTING ENERGY

We have now set up our search space over possible configurations of objects and actor motions, parameterized through the scenelet and pose placements $P$ and the assignment matrix $\mathcal{X}$. Next, we define an energy in this space that can be minimized to obtain a plausible configuration of objects and actor motions given the observations in the video. We quantify the quality of a given fit as consistency of the fitted models with the video and consistency between the

fitted models. Thus, we define an energy function that penalizes inconsistency:

$$\underset{P, \mathcal{X}}{\arg\min} \, L := w_r L_r + w_o L_o + w_s L_s + w_c L_c + w_m L_m, \qquad (4)$$

where $L_r$ is the *reprojection* error measuring the difference between the 2D joints and the projection of 3D joint locations, $L_o$ penalizes the presence of *occlusions* of skeleton joints in the video that are not explained by occlusions in the synthesized scene, $L_s$ encourages *smoothness* among human performance, $L_c$ penalizes *intersections* between objects, and $L_m$ penalizes *intersections* between the motion clip and objects. Our goal is to synthesize joint locations $q_k^t$ and objects $o_i$ by minimizing this energy over placements $P$ and assignments $\mathcal{X}$ of all fitted models, while ensuring a valid human performance. We next describe each energy term in detail.

**Reprojection term ($L_r$)** penalizes the distance from the 2D joint locations detected in the video to the corresponding output joints $q$ projected to screen space as commonly defined residue:

$$L_r = \sum_t \sum_k \left\| \Pi q_k^t - u_k^t \right\|_2^2 c_k^t, \qquad (5)$$

where $\Pi$ is the camera projection matrix, $u_k^t$ are our input 2D joint locations, and $c_k^t$ is the confidence of each detection.

**Occlusion term ($L_o$)** enforces consistency between joint occlusions observed in the video and occlusions of joints induced by synthesized scene objects. Thus, we require the synthesized objects to explain observed occlusions. We assume the person stays in the camera frame for the entire duration of the video. Hence, missing joint detections occur either due to false negatives in the detection method, or due to occluding objects. The reverse is, however, *not* true: the joint detector may, in some cases, also predict the position of occluded joints with high confidence. Therefore, we define an asymmetric occlusion error:

$$L_o = \sum_t \sum_k F(v(q_k^t, O), c_k^t), \qquad (6)$$

where $v(q_k^t, O)$ denotes the visibility of joint $q_k^t$ given the scene objects $O$. We obtain non-zero gradients that are necessary for our gradient-based solver by defining $v$ as the signed distance of joint $q_k^t$ to the *occlusion volume* induced by $O$, which is the volume that is not visible from the camera. The asymmetric occlusion error, $F$, for a joint and a set of objects is then defined as:

$$F(v, c) = \begin{cases} (c - 0.5)^2 v^2 & \text{if } c - 0.5 < 0 \text{ and } v > 0 \\ 0 & \text{otherwise,} \end{cases}$$

where $c \in [0, 1]$. Note that this function is non-zero *only* when low-confidence joint detections are explained by visible joints.

**Smoothness term ($L_s$)** ensures continuity of the synthesized motion by measuring the time-derivative of the synthesized joint locations. We approximate this derivative with finite differences:

$$L_s = \sum_t \left\| q_\lambda^t - q_\lambda^{t-1} \right\|_2^2, \qquad (7)$$

where $\lambda$ is the index of the pelvis joint at video time of frame $t$.

**Object intersection term ($L_c$)** discourages object-object penetration. In our flat scene assumption, all objects are placed on the
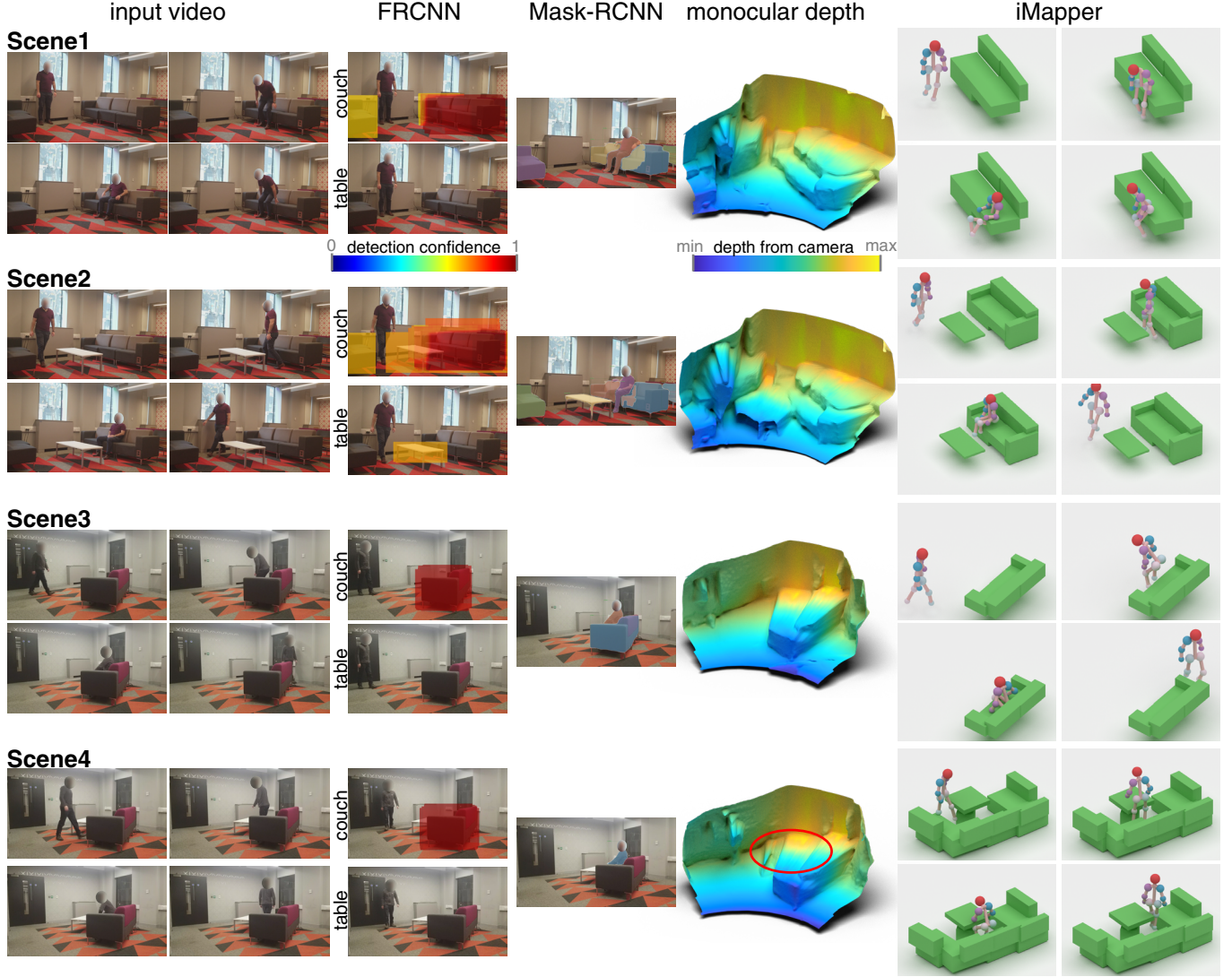
Fig. 8. Qualitative comparisons to state-of-the-art object detection methods. Note that FRCNN and Mask-RCNN both produce only 2D image-space segments. iMapper synthesizes scenes that plausibly match the video, even in sequences where objects and the human performance are occluded. Note, the table is not detected by Mask-RCNN in Scene4. Encircled in red in the bottom depth map are the typical errors near depth discontinuities obtained from image depth estimators.

ground plane. We approximate intersections in 2D, using the projections of objects to the ground plane. To obtain non-zero gradients, which are necessary to resolve intersections in a gradient-based solver, we quantify the amount of penetration using signed distance functions:

$$L_c = - \sum_{b_i \neq b_j \wedge \theta_i \neq \theta_j} \left( \int_{\Lambda(o_i)} \delta^-_{o_j}(x)\, dx + \int_{\Lambda(o_j)} \delta^-_{o_i}(x)\, dx \right), \quad (8)$$

where $\delta^-_{o_i}$ is the negative part of the signed distance function of object $o_i$, $\Lambda(o_i)$ is the projection of object $o_i$ to the ground plane, $x$ is a point on the ground plane, $b_i$ is the label of object $o_i$, and $\theta_i$ its orientation. We do not penalize intersecting objects that have the same label and orientation, since we assume these to be representations of the same object placed by different scenelets. For example, a scene with two couches facing one table may be constructed by two scenelets, each placing a couch and the same table. We identify objects to be compatible if they have the same label and orientation.

**Motion intersection term ($L_m$)** discourages humans going through objects. The trajectory of the human motion provides information about empty regions in the scene, since objects may not intersect the motion trajectory. For efficiency, we compute the intersection in 2D on the ground plane and focus on three joints only: the pelvis joint and the two knee joints. In practice, we have found that taking the maximum 2D distance of these three joints to objects allows a

Fig. 9. Qualitative comparisons to Tome et al. [2017], which produces image-space and local 3D poses. We compare the reconstructed 3D skeleton on four input videos, shown in the top row. Note that Tome and colleagues do not compute world space positions of the skeletons. For better comparison, we position them in world coordinates using the hip locations as estimated by iMapper. Relevant differences between the methods are marked with red arrows. Note that our method gives plausible skeleton poses in many cases where the method of Tome et al. fails, especially in occluded areas.

reasonable estimation of full 3D intersections in our scenes, since sitting motions can be handled correctly. We use,

$$L_m = \sum_t \max_{q \in \{q_\lambda^t,\ q_{\Gamma l}^t,\ q_{\Gamma r}^t\}} \min_i \delta_{o_i}(q),$$

where $\delta_o$ is the signed distance function of object $o$, and $q_{\Gamma l}^t$, $q_{\Gamma r}^t$ are the left- and right-knee joints, respectively.

Note that directly optimizing Equation 4 over placements of objects $o_i$ and joint positions $q_k^t$ would neither guarantee the realism of the synthesized human performance, nor its compatibility with the synthesized objects. However, our approach to minimizing this energy by fitting scenelets and static skeletons, as described in the previous section, introduces a good initialization of valid human performances and object interactions that favours realistic scenes and compatibility between objects and human performance. We obtain joint locations $q_k^t$ and the placements of objects $o_i$ from fitted scenelets and skeletons, and our goal is to optimize the fitting energy over the placements $P$ of skeletons and scenelets (see Equations 1 and 2) as well as the assignments $X$.

## 7  SCENE SYNTHESIS

Minimizing the energy $L$ gives use a plausible set of objects and actor motions. However, due to a difficult parameter domain including

Fig. 10. (Top-row) Plausible object layout and human movement as predicted by iMapper on various monocular videos. (Bottom-row) For qualitative evaluation, we overlay, shown from top-view, estimated scene layout versus annotated groundtruth. For quantitative evaluation, please refer to Tables 1 and 2. Please refer to supplmental for videos and results.

integer parameters for scenelet choice $\mathcal{X}$ and a highly non-linear energy function, it is not feasible to minimize the full energy over all parameters in a single optimization. Instead, we approximate the solution by decomposing the optimization. We start with a large set of candidates and evaluate a selection of computationally-efficient energy terms. More complex terms are added in stages, where each stage allows us to filter out low-scoring scenelet candidates. In each stage, we perform an optimization over the placement parameters $P$. The scenelet assignment $\mathcal{X}$ is optimized indirectly by filtering out candidates in each stage of the decomposed optimization instead of directly invoking an integer program.

**Static skeleton fitting.** Initially, we do not know if any given sequence of the video contains interactions (since object selection and placements are unknown). Therefore, we start by fitting static skeletons to all frames of the video, *i.e.,* we initialize $\mathcal{X}$ to zeros. We optimize for the skeleton placements $P_t$ in each frame $t$. Note that this initial scene does not yet contain objects, so only the reprojection and smoothness terms need to be minimized.

Depending on the method used to generate input pose detections, we may or may not get valid 3D skeleton guesses for occluded frames. For example, the local 3D pose detector of Tome et al. [2017] returns highly unlikely poses in occluded regions, such as poses having their knees above their head. We conservatively discard such estimates and label the corresponding frames as occluded. In addition, we perform outlier detection on the 2D keypoint estimates and define frames as occluded based on a 5% Winsorization of the joint velocities between frames. For such occluded frames, we initially interpolate the joints linearly from unoccluded frames. LCR-Net++ provides 3D joint guesses even for such frames, thus we keep these estimates even though they are unrealistic (e.g., ankles detected below the floor level). Later, these interpolated or unreliable pose estimates will be replaced by scenelets.

**Scenelet fitting.** To optimize the scenelet assignment $\mathcal{X}$ and optimize placement parameters $P$ for assigned scenelets, we start by identifying frames of the video that contain interactions.

We fit the scenelets in our dataset (1500) to each video frame in the regularly spaced subset $t \in T'$ and use the characteristicness of the scenelets weighted by their matching quality to determine the probability of an interaction at time $t$. Since we only want to fit scenelets to parts of the video that are occluded or contain interactions, we perform non-maximum suppression of the charness over the video frames and only keep frames that are at charness maxima and above a minimum charness, in addition to frames without static skeletons.

Scenelets are fitted independently, that is, for each scenelet, we evaluate the energy of a scene containing the previously fitted static skeletons plus the single fitted scenelet. Since we are only interested in evaluating the characteristicness of the motion in the video, we do not include the occlusion term in this optimization.

In addition to the charness, the partial fitting energy obtained in this step provides a lower bound for the full fitting energy. We can therefore discard high-energy scenelets. In our experiments, we keep the top 200 scenelets for each charness maximum. Figure 6 shows some example fits on different parts of a test sequence.

After this first stage of scenelet fitting, we add the motion intersection term to our energy and re-optimize the reduced set of scenelets candidates. The occlusion term is expensive to compute, we evaluate it once for the fitted scenelet candidates and add it to the energy. Based on this energy, we pick the top 3 scenelets of each charness maximum.

**Refinement.** Having committed to a set a smaller set of candidates, we can optimize placements $P$ of all fitted models in the scene, both static skeletons and scenelets, using the full energy term. For simpler scenes, we perform one optimization for all combinations of the remaining candidates. In our experiments, the maximum number of characteristicness maxima was 5, giving us a maximum of $3^5$ combinations to evaluate. We keep the combination that results in the scene with the lowest fitting energy. For more complex scenes, *e.g.,* multi-person scenes with more ambiguous motion (*e.g.,* sitting all the time) we list the top 5 diverse candidates. By diversity we prefer scenes not containing same scenelets in nearby times.

**Object selection.** Finally, we resolve intersections between objects. Recall that the object intersection term does not penalize intersections between compatible objects, *i.e.,* objects with the same orientation and category. For each such intersection, between objects $o_A$ and $o_B$, we remove the object that results in a scene with higher energy. The result of this step is our final scene.

## 8 RESULTS AND DISCUSSION

We tested iMapper on a range of input monocular videos of varying complexity, including both in-house and out-of-house sequences, such as old movies. Table 1 shows statistics of these videos. For in-house sequences, we first describe how groundtruth annotations were created, both for object placements and for actor movement. Based on this benchmark dataset, we then qualitatively and quantitatively evaluate the performance of iMapper separately for object placement and actor pose accuracy. We further compare our method against dedicated object detection and pose estimation methods. For out-of-house sequences, such as old movies, suitable groundtruth for objects and actor poses is not always available, in these cases we only provide qualitative evaluation. Please refer to the supplemental for the full input videos and ground truth annotations.

Table 1. Statistics of scenes presented in the paper showing frame count, fraction of frames with occlusion (frac. frames), number of objects in the scene with interactions (objs.). For comparison, we list number of objects detected by Mask-RCNN (MR) while iMapper detected *all* the objects with interactions. Also, we show quality of depth estimation error in cm by MonoDepth (MonoD) and iMapper as mean (s.d.) ($\mu(\sigma)$) compared against groundtruth annotations.

| scene | frame # | frac. frames | objs. | MR | MonoD $\mu(\sigma)$ | iMapper $\mu(\sigma)$ |
|---|---|---|---|---|---|---|
| Office1 | 348 | 1.00 | 7 | 5 | 263 (117) | **81** (43) |
| LvRoom | 380 | 0.84 | 5 | 4 | 98 (36) | **57** (32) |
| Library1 | 539 | 0.49 | 6 | 5 | 174 (73) | **83** (29) |
| Garden | 430 | 0.86 | 5 | 5 | 242 (68) | **58** (28) |
| Lobby1 | 254 | 0.65 | 3 | 1 | 243 (98) | **70** (55) |
| Lobby2 | 224 | 0.97 | 4 | 1 | 259 (101) | **51** (21) |
| Scene1 | 80 | 0.00 | 1 | 2 | 120 (-) | **45** (-) |
| Scene2 | 130 | 0.57 | 2 | 5 | 120 (10) | **72** (6) |
| Scene3 | 120 | 0.82 | 1 | 1 | 191 (-) | **25** (-) |
| Scene4 | 115 | 0.77 | 2 | 2 | 203 (12) | **70** (57) |
| Scene5a | 49 | 0.84 | 2 | 2 | 197 (24) | **135** (13) |
| Scene5b | 148 | 0.93 | 3 | 2 | 230 (52) | **72** (41) |
| Scene5c | 182 | 1.00 | 4 | 2 | 229 (48) | **53** (49) |
| Scene5d | 189 | 0.98 | 4 | 2 | 216 (50) | **69** (38) |

**Groundtruth dataset.** In order to create a groundtruth benchmark dataset, we annotated both object locations and 3D world-space poses for the in-house video sequences.

For object locations, we physically measured the scene objects' dimensions and positioned objects in the scene to minimize video reprojection error, using known camera intrinsics. We also added labels (*e.g.,* . 'chair') for each individual object.

For human poses, manually annotating the 3D pose in each frame is not feasible, so we used an assisted approach to generate the groundtruth. We started with estimated 2D joint locations [Wei

et al. 2016] and then manually corrected them. These corrected 2D locations were then lifted to 3D using the reprojection and smoothness energy terms described in Section 6. Finally, we inspected the output, manually corrected 3D poses, and added these corrections as additional constraints to the optimization. This process was repeated until we found no more significant errors. The number of corrections depended on the amount of occlusion in the scene.

**Multi-actor scenes.** We described our method in the context of single actors. However, for scenes with multiple actors, our method generalizes easily when input 2D tracks come with correspondence information over time. We developed a semi-automatic labeling method that optimizes a Markov Random Field (MRF) energy term to find correspondences between skeleton detections, given some manual guidance (see Appendix A for details). This labeling tool was employed for the following multi-person scenes: Office2, Library2, Angrymen, and Grease.

### 8.1 Qualitative Evaluation

First, we show qualitative results of our method on several example scenes. We demonstrate that iMapper finds plausible object arrangements and actor poses for occluded parts of the scene by rendering our 3D result from a different viewpoint than the source video. Note that in these views, the recreated human motions may appear noisy as our scenelets are based on raw Kinect-based captures included the PiGraph dataset.

Figures 1, 10, and 11, we show results for 9 different scenes (more examples/videos in the supplemental). For each of these scenes, we show a reference video frame in the background on a plane facing the camera. The trajectory of the actor's pelvis is shown as a colored line and the colored skeleton shows an occluded actor pose. Please refer to supplemental for full videos of the reconstructed actor motions. To better visualize the detected orientation of objects, we use object bounding boxes to scale and place category specific proxy object geometries in our scenelets – please note the meshes are *not* output of our method. In the future we plan to update our dataset to use detailed models in the scenelets, making this step unnecessary.

In Figure 10, we show scenes from two old movies: '12 Angry Men' and 'Grease'. To our knowledge there is currently no method available to give a plausible reconstruction of these scenes, due to significant occlusions. Our method finds a plausible object layout and actor motion in these scenes based on the observed interactions. The remaining examples in the figure show results on in-house scenes with varying amounts of occlusion.

Note how much information is encoded in the interactions, enabling us to generate scenes close to the original scenes. We can, however, not hallucinate objects that are not interacted with, or obtain information about exact dimensions of objects from interactions.

**Scene exploration over time.** One advantage of iMapper is that the reconstruction quality of the scene improves over time as more interactions 'reveal' the true underlying scene. As shown in Figure 12, as the same environment is explored over time, our system recovers larger parts of the object arrangement. Further, perturbations to the input (*e.g.,* in the form of the same action being

Fig. 11. Scene layout and human motion estimated by iMapper on monocular videos containing multiple actors (see supplementary video). Note that in case of actors not moving sufficiently during the sequence (e.g., Angrymen, Grease), the estimated human 'motion' can have drifts.

Table 2. Comparison of pose estimates by iMapper against Tome et al. [2017], LCR-Net++3D [Rogez et al. 2018], and VNect [Mehta et al. 2017b]. Note that Tome et al. and LCR-Net++3D return only image-space and local coordinates (LC); while, VNect is the only other method returning world coordinates (WC). All units are in cm.

| | Tome3D | | | LCR-Net++3D | | | Vnect | | | iMapper | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | WC | LC | 2D | WC | LC | 2D | WC | LC | 2D | WC | LC | 2D |
| LvRoom | x | 22.1 | 346.0 | x | **12.0** | **77.0** | x | x | x | **71.5** | 20.9 | 197.6 |
| Lobby1 | x | 28.1 | 122.5 | x | **21.9** | **63.6** | x | x | x | **60.8** | 29.2 | 148.9 |
| Lobby2 | x | 19.6 | 81.0 | x | 20.1 | 73.6 | x | x | x | **32.0** | **19.2** | **68.6** |
| Scene3 | x | 22.4 | 70.6 | x | **19.2** | 83.5 | x | x | x | **56.7** | 19.7 | **65.4** |
| Scene4 | x | 25.1 | 91.6 | x | 20.7 | 104.1 | x | x | x | **61.8** | **13.6** | **46.3** |
| GrRoom | x | 35.5 | 237.9 | x | **20.0** | 150.4 | 68.3 | 28.9 | 171.1 | **40.3** | 23.5 | **138.5** |

performed by different people, or at different times) lead to slightly different, but still plausible and consistent reconstructions of the scene and the interactions.

## 8.2 Quantitative and Qualitative Comparisons

In the following, we compare both the scene layout and the actor motion recovered by iMapper to dedicated object detection and pose estimation methods.

**Scene Layout.** We compare against two types of methods: per-frame region detection using FRCNN [Ren et al. 2017] and Mask-RCNN [He et al. 2017], and per-pixel monocular depth estimation [Chakrabarti et al. 2016].

*Region detection methods:* For the region detection methods, qualitative comparisons are shown in Figure 8, and quantitative results are given in Table 1, column 'MR'. In Table 1, we count the number of objects where at least 50% of the object's region was detected and correctly labeled on average, over all frames. We provide the count of objects that are participating in at least one interaction. Since iMapper discovers objects through interactions, this is the set of objects we can potentially detect. FRCNN and Mask-RCNN (MR) are designed to detect visible objects, so they naturally fail to detect any objects 'hidden' behind visible objects. For MR this is reflected in the table by a low number of detected objects. Note that other systems that rely on FRCNN or MR as their primary building blocks will have similar problems in occluded regions. Our method can more reliably recover these occluded objects if they participate in interactions.

*Depth estimation methods:* For per-pixel depth estimation [Chakrabarti et al. 2016] method, as seen in Figure 8, fourth column, even for visible regions the estimated depths are smoothed out and fail to capture the object specific layout. Table 1, column 'MonoD', shows the mean and standard deviation of the distances between the predicted and the ground truth object centroids in the scene. For the monocular depth map, we approximate the object centroid as the mean world position of all samples that are inside the 2D region of an object. Objects without a single visible pixel are ignored. Again, the depth map contains only limited information about partially or

fully occluded objects, resulting in large errors. In contrast, iMapper produces plausible objects along with their spatial locations. These locations, in turn, provide better occlusion information for 3D human movement, as evaluated next.

**Actor Poses.** Most monocular 3D pose detection methods compute only *local* 3D poses, i.e., joint locations relative to pelvis, limiting our choice of baselines. We compare to Vnect [Mehta et al. 2017b], Tome et al. [2017], and LCRNet++3D [Rogez et al. 2018]. Both Tome et al. and LCRNet++3D output local 3D joint locations, and do not provide world-space coordinates, only VNect also provides world coordinates. For all quantitative comparisons, we report the root mean square error (RMSE) over 14 joint locations. This error is computed in three spaces: 2D image space, (with a resolution of 1920 × 1080), local 3D (pelvis-relative) space, and world space, where available. *Vnect:* A qualitative comparison to Vnect is given in Figure 2. Since we do not have direct access to their source code, we compare to this method on the *MPI-INF-3DHP* dataset [Mehta et al. 2017a], where a relatively high-quality ground truth is available. VNect was designed for unoccluded human motion, so results in occluded areas are not robust. Quantitative results are given in Table 2. As expected, iMapper can improve upon Vnect in occluded scenes, resulting in better local and world-space 3D predictions.
*Tome et al. and LCRNet++3D:* Figure 9 shows qualitative comparisons to Tome et al. While non-occluded poses are very close to our method (up to the smoothness term), in occluded frames, Tome3D returns unrealistic poses. iMapper, however, by fitting scenelets to these occluded sequences returns more realistic (partial) poses that are visually closer to the pose observed in the video. Note that in the last column of **D**, the reprojection is slightly worse for our method, since the scenelet database did not contain a sufficiently close match to describe this kind of motion.

Table 2 shows quantitative comparisons to both Tome et al. and LCRNet++3D. LCR-Net performs well in the less heavily occluded scenes shown near the top of the table, but performance drops in the presence of occlusions. We observe a similar trend for Tome3D, but with a larger error on average. iMapper can improve upon both methods in more heavily occluded scenes, where we can rely on interaction priors to give us information about hidden joints that the two other methods cannot correctly predict.

### 8.3 Ablation study

In Figure 13 we report the effect of removing some of the terms from our optimization. The reprojection and smoothness terms have the highest effect on the result, since they are used throughout the entire pipeline. Omitting the smoothness term allows for strong path deformations, while omitting the reprojection term removes the anchoring of the motions to the video and permits the smoothness term to strongly contract the path. The other terms have a more situational effect on the results: how much they influence the result depends on the scene configuration. The occlusion term, for example, has a heavy influence if multiple joint locations are occluded in the video. Recall that our occlusion term is assymmetric, so that more occlusion always has less cost. The couch is thus optimized to be closer to the camera to occlude more of the scene (the camera is at the bottom center of the image).

**Hold-one-out validation** We evaluated our method on a hold out-set taken from the PiGraphs dataset. We pick a single scene and remove the scenelets that were generated from this scene from our database, accounting for ∼ 10% of our scenelets. We compare to the groundtruth objects given in the PiGraphs dataset through manually established correspondence. Our mean reconstruction error was 110 cm (std: 56 cm) with all relevant objects detected.

### 8.4 Limitations

We currently assume the input camera to be fixed during the entire duration of the capture. Such a fixed viewpoint results in heavy occlusion in crowded scenes and makes the job of iMapper unnecessarily challenging. In large-scale environments, one can imagine having a network of fixed cameras to address this issue and jointly process the information from the individual cameras. An obvious limitation of iMapper is failing to react when it 'sees' an interaction that is missing in its interaction database. This is an unavoidable
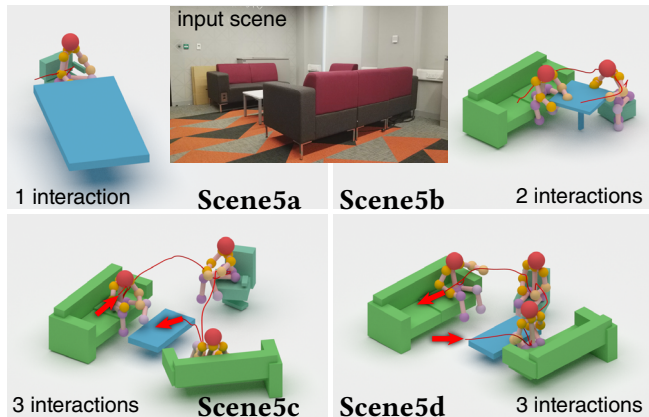


Fig. 12. Scenes can be explored over time. We show results of four different videos taken from the same scene. As interactions with more objects are made available, we can recompute the results to synthesize additional objects. Variations of scene explorations, for example performing the interactions in reverse order, as shown in the bottom row, give slightly different, but comparable and plausible results.
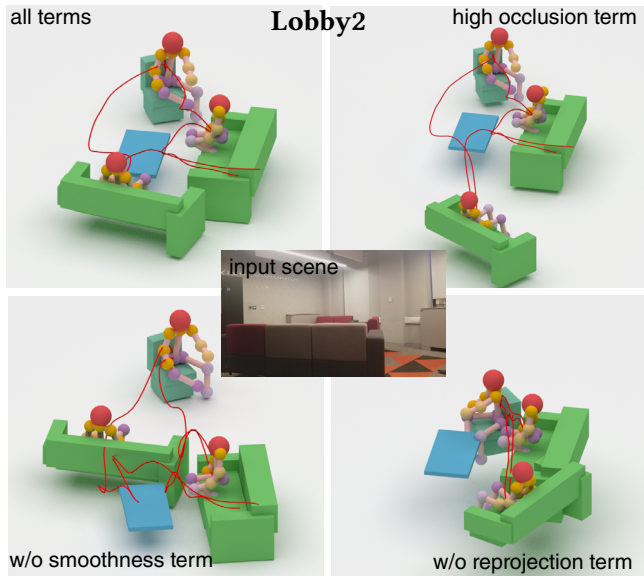


Fig. 13. We test the effect of various terms on the solution for the scene shown in the center. Reprojection and smoothness generally have the largest influence, since they are used throughout the entire pipeline.

problem in any data-driven approach unless we are able to additionally synthesize new interactions, which is a significantly more challenging problem. Finally, since our method builds on the expectation that people react similarly in similar settings, it will naturally get confused when this assumption is broken. For example, if a person decides to hand-walk, or use a sofa as a bed, *etc.* We expect this to be partially addressed by richer interaction databases with associated probability priors on actions, and spatio-temporal reasoning that detects and ignores 'outlier' interactions by observing a scene over longer time intervals.

## 9 CONCLUSION AND FUTURE WORK

Based on the observation that humans (both the same person or different people) interact similarly in similar scenes, we presented iMapper to *jointly* reason about static object arrangements and human movements. At the heart of iMapper lies a novel data-driven method to define *characteristic* poses that help identify space-time moments when matched human poses provide reliable cues about the surrounding scene arrangement. Our method links such partial scenelets, fitted to monocular video under unknown occlusion, and assembles them to form a global scene layout and 3D human pose estimates. By extensive evaluation we demonstrate that iMapper improves both the quality of scene layout estimation as well as 3D pose estimates, especially in scenes with low to medium levels of occlusions.

Exciting research directions lay ahead as we are only starting to capture, analyze, and understand the space of (human) interactions, or *interaction landscapes* (*cf.,* [Pirk et al. 2017b]). Below we discuss some of the immediate issues.

*Capturing richer interaction databases:* Current datasets only capture limited variety of interactions. Limited refers to both different types of interactions, and variance for each interaction type. For example,

we miss examples of interactions with small objects (*e.g.*, picking up a cup/glass, using pots and pans in kitchens, lifting a bag or suitcase), or examples of the different ways that people sit in sofas, couches, chairs, etc. While significant progress has been made in capturing static environments at high geometric detail, capturing interactions remains *fundamentally* difficult due to heavy occlusion arising due to the interactions. One possibility is to separate the capture of static geometry (*e.g.*, with mobile 3D scanners) from the capture of interactions using a mix of sensors like IMU sensors, RGBD scanners, markers, *etc.* We expect such data gathering efforts to happen in the near future.

*Utilizing scene priors:* We used signals only from interactions. However, in scenes with heavy occlusion, scenelet matching with partial (occluded) information is not sufficient to accurately ground object positions. Also, we cannot directly distinguish between objects that afford similar interactions, such as couches and chairs. One direction would be to additionally use scene statistics and local context, as has been heavily utilized in scene synthesis research, to regularize the interaction-based reconstruction problem.

*Recovering interactions over large timescales:* As shown in Figure 12, iMapper has only a chance of recovering scene arrangements once people interact with various parts of the environment. This suggests that the approach gets better as we 'observe' the scene over larger timescales, ideally days or weeks. However, then our static scene assumption breaks down as objects are going to be shifted and moved around. Hence, we would like to extend our approach to also capture space-time object movements, starting with rigid movement of objects, such as a moving chair.

## REFERENCES

Gabriel J Brostow and Irfan A Essa. 1999. Motion based decompositing of video. *IEEE ICCV* (1999).

Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. 2017. Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields. In *IEEE CVPR*.

Ayan Chakrabarti, Jingyu Shao, and Gregory Shakhnarovich. 2016. Depth from a Single Image by Harmonizing Overcomplete Local Network Predictions. *CoRR* abs/1605.07081 (2016).

Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. 2017. Matterport3D: Learning from RGB-D Data in Indoor Environments. *3DV*.

Kang Chen, Yu-Kun Lai, Yu-Xin Wu, Ralph Martin, and Shi-Min Hu. 2014. Automatic Semantic Modeling of Indoor Scenes from Low-quality RGB-D Data Using Contextual Information. *ACM SIGGRAPH Asia* 33, 6, Article 208 (Nov. 2014), 12 pages.

Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. 2017a. ScanNet: Richly-annotated 3D Reconstructions of Indoor Scenes. In *IEEE CVPR*.

Angela Dai, Matthias Nießner, Michael Zollöfer, Shahram Izadi, and Christian Theobalt. 2017b. BundleFusion: Real-time Globally Consistent 3D Reconstruction using On-the-fly Surface Re-integration. *ACM TOG* (2017).

Luca Del Pero, Joshua Bowdish, Bonnie Kermgard, Emily Hartley, and Kobus Barnard. 2013. Understanding Bayesian Rooms Using Composite 3D Object Models. In *IEEE CVPR*.

V. Delaitre, D. Fouhey, I. Laptev, J. Sivic, A. Gupta, and A. Efros. 2012. Scene semantics from long-term observation of people. *ECCV* (2012).

Matthew Fisher, Daniel Ritchie, Manolis Savva, Thomas Funkhouser, and Pat Hanrahan. 2012. Example-based Synthesis of 3D Object Arrangements. In *ACM SIGGRAPH Asia*.

Matthew Fisher, Manolis Savva, and Pat Hanrahan. 2011. Characterizing structural relationships in scenes using graph kernels. In *ACM SIGGRAPH*, Vol. 30. 34.

Matthew Fisher, Manolis Savva, Yangyan Li, Pat Hanrahan, and Matthias Nießner. 2015. Activity-centric Scene Synthesis for Functional 3D Scene Modeling. *ACM SIGGRAPH* 34, 6 (2015).

David F. Fouhey, Vincent Delaitre, Abhinav Gupta, Alexei A. Efros, Ivan Laptev, and Josef Sivic. 2012. People Watching: Human Actions as a Cue for Single-View Geometry. *ECCV* (2012).

B. Frank, M. Ruhnke, M. Tatarchenko, and W. Burgard. 2015. 3D-reconstruction of indoor environments from human activity. In *IEEE ICRA*. 4644–4649.

Lianrui Fu, Junge Zhang, and Kaiqi Huang. 2015. Beyond Tree Structure Models: A New Occlusion Aware Graphical Model for Human Pose Estimation. In *IEEE ICCV*.

Q. Fu, X. Chen, X. Su, and H. Fu. 2017a. Pose-Inspired Shape Synthesis and Functional Hybrid. *IEEE TVCG* 23, 12 (2017), 2574–2585.

Qiang Fu, Xiaowu Chen, Xiaotian Wang, Sijia Wen, Bin Zhou, and Hongbo Fu. 2017b. Adaptive Synthesis of Indoor Scenes via Activity-Associated Object Relation Graphs. *ACM SIGGRAPH Asia* 36, 6 (2017), Article No. 201.

Georgia Gkioxari, Ross Girshick, Piotr Dollár, and Kaiming He. 2018. Detecting and Recognizing Human-Object Interactions. *CVPR* (2018).

Abhinav Gupta, Aniruddha Kembhavi, and Larry S. Davis. 2009. Observing Human-Object Interactions: Using Spatial and Functional Compatibility for Recognition. *IEEE PAMI* 31, 10 (Oct. 2009), 1775–1789.

K. He, G. Gkioxari, P. DollÃąr, and R. Girshick. 2017. Mask R-CNN. In *2017 IEEE International Conference on Computer Vision (ICCV)*. 2980–2988. https://doi.org/10.1109/ICCV.2017.322

C. H. Huang, E. Boyer, N. Navab, and S. Ilic. 2014. Human Shape and Pose Tracking Using Keyframes. In *IEEE CVPR*. 3446–3453. https://doi.org/10.1109/CVPR.2014.440

Jia-Bin Huang and Ming-Hsuan Yangc. 2009. Estimating Human Pose from Occluded Images. In *ACCV*.

Shi-Sheng Huang, Hongbo Fu, and Shi-Min Hu. 2016. Structure guided interior scene synthesis via graph matching. *Graphical Models* 85 (2016), 46 – 55.

Eldar Insafutdinov, Leonid Pishchulin, Bjoern Andres, Mykhaylo Andriluka, and Bernt Schiele. 2016. DeeperCut: A Deeper, Stronger, and Faster Multi-Person Pose Estimation Model. *ECCV* (Oct 2016).

Hamid Izadinia, Qi Shan, and Steven M Seitz. 2017. IM2CAD. In *CVPR*.

Yun Jiang, Hema S. Koppula, and Ashutosh Saxena. 2016. Modeling 3D Environments Through Hidden Human Context. *IEEE PAMI* 38, 10 (Oct. 2016), 2040–2053.

Changgu Kang and Sung-Hee Lee. 2017. Scene reconstruction and analysis from motion. *Graphical Models* 94 (2017), 25 – 37.

Vladimir G. Kim, Siddhartha Chaudhuri, Leonidas Guibas, and Thomas Funkhouser. 2014. Shape2Pose: Human-Centric Shape Analysis. *ACM SIGGRAPH* (Aug. 2014).

Leonard Krasner. 2013. *Environmental Design and Human Behavior*. Elsevier.

Tianqiang Liu, Siddhartha Chaudhuri, Vladimir G. Kim, Qixing Huang, Niloy J. Mitra, and Thomas Funkhouser. 2014. Creating Consistent Scene Graphs Using a Probabilistic Grammar. *ACM SIGGRAPH Asia* 33, 6, Article 211 (Nov. 2014), 12 pages.

Rui Ma, Honghua Li, Changqing Zou, Zicheng Liao, Xin Tong, and Hao Zhang. 2016. Action-driven 3D Indoor Scene Evolution. *ACM SIGGRAPH Asia* 35, 6, Article 173 (Nov. 2016), 13 pages.

Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. 2017a. Monocular 3D Human Pose Estimation In The Wild Using Improved CNN Supervision. In *3D Vision (3DV), 2017 Fifth International Conference on*. IEEE. https://doi.org/10.1109/3dv.2017.00064

Dushyant Mehta, Srinath Sridhar, Oleksandr Sotnychenko, Helge Rhodin, Mohammad Shafiei, Hans-Peter Seidel, Weipeng Xu, Dan Casas, and Christian Theobalt. 2017b. VNect: Real-time 3D Human Pose Estimation with a Single RGB Camera. *ACM SIGGRAPH* 36, 4, 14.

Liangliang Nan, Ke Xie, and Andrei Sharf. 2012. A Search-classify Approach for Cluttered Indoor Scene Understanding. *ACM SIGGRAPH Asia* 31, 6, Article 137 (Nov. 2012), 10 pages.

Ulric Neisser. 1976. *Environmental Design and Human Behavior*. W. H. Freeman, SF.

Richard A. Newcombe et al. 2011. KinectFusion: Real-time dense surface mapping and tracking. In *IEEE ISMAR*. 127–136.

Richard A. Newcombe, Dieter Fox, and Steven M. Seitz. 2015. DynamicFusion: Reconstruction and tracking of non-rigid scenes in real-time. *IEEE CVPR* (2015).

Alejandro Newell, Kaiyu Yang, and Jia Deng. 2016. Stacked Hourglass Networks for Human Pose Estimation. *ECCV* (2016).

Sören Pirk, Olga Diamanti, Boris Thibert, Danfei Xu, and Leonidas J. Guibas. 2017a. SHAPE-AWARE SPATIO-TEMPORAL DESCRIPTORS FOR INTERACTION CLASSIFICATION. *ICIP* (2017).

Sören Pirk, Vojtech Krs, Kaimo Hu, Suren Deepak Rajasekaran, Hao Kang, Yusuke Yoshiyasu, Bedrich Benes, and Leonidas J. Guibas. 2017b. Understanding and Exploiting Object Interaction Landscapes. *ACM SIGGRAPH Asia* 36, 3, Article 31 (June 2017), 14 pages.

P. Poirson, P. Ammirato, C. Y. Fu, W. Liu, J. Kosecka, and A. C. Berg. 2016. Fast Single Shot Detection and Pose Estimation. In *3DV*. 676–684.

Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2017. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *IEEE PAMI* 39, 6 (2017), 1137–1149.

Grégory Rogez, Philippe Weinzaepfel, and Cordelia Schmid. 2018. LCR-Net++: Multi-person 2D and 3D Pose Detection in Natural Images. *CoRR* abs/1803.00455v1 (2018).

S. Satkin and M. Hebert. 2013. 3DNN: Viewpoint Invariant 3D Geometry Matching for Scene Understanding. In *IEEE CVPR*. 1873–1880. https://doi.org/10.1109/ICCV.2013.235

Manolis Savva, Angel X. Chang, Pat Hanrahan, Matthew Fisher, and Matthias Nießner. 2014. SceneGrok: Inferring Action Maps in 3D Environments. *ACM SIGGRAPH Asia* 33, 6 (2014).

Manolis Savva, Angel X. Chang, Pat Hanrahan, Matthew Fisher, and Matthias Nießner. 2016. PiGraphs: Learning Interaction Snapshots from Observations. *ACM SIGGRAPH* 35, 4 (2016).

A. G. Schwing, S. Fidler, M. Pollefeys, and R. Urtasun. 2013. Box in the Box: Joint 3D Layout and Object Reasoning from Single Images. In *IEEE ICCV*. 353–360. https://doi.org/10.1109/ICCV.2013.51

Tianjia Shao, Weiwei Xu, Kun Zhou, Jingdong Wang, Dongping Li, and Baining Guo. 2012. An Interactive Approach to Semantic Modeling of Indoor Scenes with an RGBD Camera. *ACM SIGGRAPH Asia* 31, 6, Article 136 (Nov. 2012), 11 pages.

Bugra Tekin, Artem Rozantsev, Vincent Lepetit, and Pascal Fua. 2016. Direct Prediction of 3D Body Poses from Motion Compensated Sequences. In *IEEE CVPR*.

Denis Tome, Chris Russell, and Lourdes Agapito. 2017. Lifting from the Deep: Convolutional 3D Pose Estimation from a Single Image. *IEEE CVPR* (2017).

Alexander Toshev and Christian Szegedy. 2014. DeepPose: Human Pose Estimation via Deep Neural Networks. In *IEEE CVPR*.

T. von Marcard, B. Rosenhahn, M. J. Black, and G. Pons-Moll. 2017. Sparse Inertial Poser: Automatic 3D Human Pose Estimation from Sparse IMUs. *CGF Eurographics* 36, 2 (May 2017), 349–360.

Ping Wei, Yibiao Zhao, Nanning Zheng, and Song-Chun Zhu. 2013. Modeling 4D Human-Object Interactions for Event and Object Recognition. In *IEEE ICCV*.

Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. 2016. Convolutional pose machines. In *IEEE CVPR*.

Kai Xu, Rui Ma, Hao Zhang, Chenyang Zhu, Ariel Shamir, Daniel Cohen-Or, and Hui Huang. 2014. Organizing Heterogeneous Scene Collections Through Contextual Focal Points. *ACM SIGGRAPH* 33, 4, Article 35 (July 2014), 12 pages.

Bangpeng Yao, Aditya Khosla, and Li Fei-Fei. 2011. Classifying Actions and Measuring Action Similarity by Modeling the Mutual Context of Objects and Human Poses. In *ICML*.

Yi-Ting Yeh, Lingfeng Yang, Matthew Watson, Noah D. Goodman, and Pat Hanrahan. 2012. Synthesizing Open Worlds with Constraints Using Locally Annealed Reversible Jump MCMC. *ACM SIGGRAPH* 31, 4, Article 56 (July 2012), 11 pages.

Hong-Bo Zhang, Qing Lei, Bi-Neng Zhong, Ji-Xiang Du, and JiaLin Peng. 2016. A Survey on Human Pose Estimation. *Intelligent Automation and Soft Computing* 22, 3 (2016), 483–489.

Xi Zhao, Ruizhen Hu, Paul Guerrero, Niloy Mitra, and Taku Komura. 2016. Relationship Templates for Creating Scene Variations. *ACM SIGGRAPH Asia* 35, 6, Article 207 (Nov. 2016), 13 pages.

Xiaowei Zhou, Menglong Zhu, Spyridon Leonardos, Kosta Derpanis, and Kostas Daniilidis. 2016. Sparseness Meets Deepness: 3D Human Pose Estimation from Monocular Video. *IEEE CVPR* (2016).

## A  MULTI-PERSON TRACKER

In case of multi-person videos, we require the user to inspect the detected 2D skeletons over time and mark the ones that need tracking in the first frame they appear. The input pose detections are not temporally consistent, so we perform a naive actor-reidentification in the form of an MRF optimization. Specifically, given the number of actors $n_a$, the set of skeletons $S(t)$ detected in frame $t$ of the video, our goal is to find the values of the binary variables $\xi^t_{a \to s} \in \{0, 1\}$ where $\xi^t_{a \to s} = 1$ if actor $a$ is associated with skeleton $s \in S(t)$. To allow some actors to be identified as *invisible* in certain portions of the video, we define a dummy skeleton $\bar{s}$ such that if $\xi^t_{a \to \bar{s}} = 1$, $a$ is marked as invisible. We represent the user annotations for a new actor in the first frame they appear as hard constraints in our optimization. In addition, we define unary and binary terms. The unary term, $E_u(a, s, t)$, measures the cost of assigning an actor $a$ to a skeleton $s$ at a frame $t$ and is based on per-joint confidence measures, $c^t_{k,s}$:

$$E_u(a, s, t) = -\xi^t_{a \to s} \begin{cases} \frac{1}{n_{joints}} \sum_k^{n_{joints}} c^t_{k,s} & \text{if } s \neq \bar{s} \\ 1 & \text{else.} \end{cases} \quad (9)$$

Assigning the dummy skeleton, $\bar{s}$, to any actor results in a fixed cost of 1. The binary term, $E_b(a, s_0, s_1, t)$, measures the cost of assigning an actor $a$ to the skeleton $s_0$ and $s_1$ in frames $t$ and $t + 1$

respectively:

$$E_b(a, s_0, s_1, t) = \xi^t_{a \to s_0} \xi^{t+1}_{a \to s_1} C(a, s_0, s_1, t) \quad (10)$$

$$C(a, s_0, s_1, t) =$$

$$\begin{cases} 1 & \text{if } \bar{s} \in \{s_0, s_1\} \\ \frac{1}{n_{joints}} \sum_k^{n_{joints}} \left\| \frac{u^t_{k,s_0} - u^{t+1}_{k,s_1}}{diag} \right\|_2^2 c^t_{k,s_0} c^{t+1}_{k,s_1} & \text{else,} \end{cases}$$

where *diag* refers the the half of the diagonal of the image and transitioning from/to the dummy skeleton results in a fixed cost of 1. Finally, we optimize for the following energy function:

$$\operatorname*{argmin}_\xi E = - \sum_t \sum_{a \leq n_a, s \in S(t)} E_u(a, s, t)$$
$$+ w_{pw} \sum_t \sum_{a \leq n_a} \sum_{s_0 \in S(t)} \sum_{s_1 \in S(t+1)} E_b(a, s_0, s_1, t) \quad (11)$$

subject to the constraints:

$$\forall_t \forall_s \sum_{a \leq n_a} \xi^t_{a \to s} = 1 \,, \; \forall_t \forall_a \sum_{s \neq \bar{s}} \xi^t_{a \to s} \leq 1. \quad (12)$$

We set the relative weight of the binary term in Equation 11 as $w_{pw} = 10^3$ and optimize using a binary discrete optimizer (Gurobi).

For scenes with lots of occlusion and crossings, we require additional manual constraints. Once, multi-actor/2D skeleton associations are established, we optimize for 3D poses by enforcing smoothness between poses that belong to the same actor only.

## B  CONFIDENCE OF KEYPOINT DETECTION USING LCR-NET++

When using LCR-Net++, we confidence estimate 2D detection keypoints as

$$v_k = \frac{\underset{q^i_k \in pose\ proposals}{var} \left( q^i_k \right)}{1 + \exp\left(-0.2s' + 3.5\right)} \quad (13)$$

$$c_k(v_k) = \frac{1}{1 + \exp\left(-10 \exp\left(\frac{log(v_k)}{P_{99}(log(v_k))}\right) + 24\right)}, \quad (14)$$

where $var_k$ denotes the variance of the 3D joint position among the grouped pose proposals, and $P_{99}$ denotes the 99th percentile of *log* joint variances over the whole recording, assigning high confidence to low variance joint estimates, and $s'$ is a per-pose score defined in Equation 6 in [Rogez et al. 2018].