

Learning Multi-Human Optical Flow

Anurag Ranjan · David T. Hoffmann · Dimitrios Tzionas · Siyu Tang · Javier Romero · Michael J. Black

Received: date / Accepted: date

Abstract The optical flow of humans is well known to be useful for the analysis of human action. Recent optical flow methods focus on training deep networks to approach the problem. However, the training data used by them does not cover the domain of human motion. Therefore, we develop a dataset of multi-human optical flow and train optical flow networks on this dataset. We use a 3D model of the human body and motion capture data to synthesize realistic flow fields in both single- and multi-person images. We then train optical flow networks to estimate human flow fields from pairs of images. We demonstrate that our trained networks are more accurate than a wide range of top methods on held-out test data and that they can generalize well to real image sequences. The code, trained models and the dataset are available for research.

*Anurag Ranjan and David Hoffmann contributed equally.

Anurag Ranjan^{*}¹
 anurag.ranjan@tue.mpg.de

David T. Hoffmann^{*}¹
 david.hoffmann@tue.mpg.de

Dimitrios Tzionas¹
 dimitris.tzionas@tue.mpg.de

Siyu Tang¹
 siyu.tang@tue.mpg.de

Javier Romero²
 javier@amazon.com

Michael J. Black¹
 black@tue.mpg.de

¹ Max Planck Institute for Intelligent Systems, Germany

² Amazon Inc.

This work was done when JR was at MPI-IS. MJB has received research gift funds from Intel, Nvidia, Adobe, Facebook, and Amazon. While MJB is a part-time employee of Amazon, his research was performed solely at, and funded solely by, MPI. MJB has financial interests in Amazon and Meshcapade GmbH.

1 Introduction

A significant fraction of videos on the Internet contain people moving [4] and the literature suggests that optical flow plays an important role in understanding human action [5, 6]. Several action recognition datasets [6, 7] contain human motion as a major component. The 2D motion of humans in video, or *human optical flow*, is an important feature that provides a building block for systems that can understand and interact with humans. Human optical flow is useful for various applications including analyzing pedestrians in road sequences, motion-controlled gaming, activity recognition, human pose estimation system, etc.

Despite this, optical flow has previously been treated as a generic, low-level, vision problem. Given the importance of people, and the value of optical flow in understanding them, we develop a dataset and trained models that are specifically tailored to humans and their motion. Such motions are non-trivial since humans are complex, articulated objects that vary in shape, size and appearance. They move quickly, adopt a wide range of poses, and self-occlude or occlude in multi-person scenarios.

Our goal is to obtain more accurate 2D motion estimates for human bodies by training a flow algorithm specifically for human movement. To do so, we create a large and realistic dataset of humans moving in virtual worlds with ground truth optical flow (Fig. 1(a)), called the *Human Optical Flow* dataset. This is comprised of two parts; the *Single-Human Optical Flow* dataset (SHOF), where the image sequences contain only one person in motion and the *Multi-Human Optical Flow* dataset (MHOF) where images contain multiple people involving significant occlusion between them. We analyse the performance of SPyNet [2] and PWC-Net [3] by training (fine-tuning) them on both the SHOF and MHOF dataset. We observe that the optical flow performance of the networks improves on sequences containing human scenes,



Fig. 1 (a) We simulate human motion in virtual worlds creating an extensive dataset with images (top row) and flow fields (bottom row); color coding from [1]. (b) We train SPyNet [2] and PWC-Net [3] for human motion estimation and show that they performs better when trained on our dataset and (c) can generalize to human motions in real world scenes. Columns show single-person and multi-person cases alternately.

both qualitatively and quantitatively. Furthermore we show that the trained networks generalize to real video sequences (Fig. 1(c)). Several datasets and benchmarks [1, 8, 9] have been established to drive the progress in optical flow. We argue that these datasets are insufficient for the task of human motion estimation and, despite its importance, no attention has been paid to datasets and models for human optical flow. One of the main reasons is that dense human motion is extremely difficult to capture accurately in real scenes. Without ground truth, there has been little work focused specifically on estimating human optical flow. To advance research on this problem, the community needs a dataset tailored to human optical flow.

A key observation is that recent work has shown that optical flow methods trained on synthetic data [2, 10, 11] generalize relatively well to real data. Additionally, these methods obtain state-of-the-art results with increased realism of the training data [12, 13]. This motivates our effort to create a dataset designed for human motion.

To that end, we use the SMPL [14] and SMPL+H [15] models, that capture the human body alone and the body together with articulated hands respectively, to generate different human shapes including hand and finger motion. We then place humans on random indoor backgrounds and simulate human activities like running, walking, dancing etc. using motion capture data [16, 17]. Thus, we create a large virtual dataset that captures the statistics of natural human motion in multi-person scenarios. We then train on this deep neural networks and evaluate their performance for estimating human motion. While the dataset can be used to train any flow method, we focus specifically on networks based on spatial pyramids, namely SpyNet [2] and PWC-Net [3], because they are compact and computationally efficient.

A preliminary version of this work appeared in [18] that presented a dataset and model for human optical flow for

the *single-person* case with a *body-only* model. The present work extends [18] for the *multi-person* case, as images with multiple occluding people have different statistics. It further employs a holistic model of the *body together with hands* for more realistic motion variation. This work also extends training SPyNet [2] and PWC-Net [3] using the new dataset in contrast to training only SPyNet in the earlier work [18]. Our experiments show both qualitative and quantitative improvements.

In summary, our major contributions in this extended work are: 1) We provide the *Single-Human Optical Flow* dataset (SHOF) of human bodies in motion with realistic textures and backgrounds, having 146,020 frame pairs for single-person scenarios. 2) We provide the *Multi-Human Optical Flow* dataset (MHOF), with 111,312 frame pairs of multiple human bodies in motion, with improved textures and realistic visual occlusions, but without (self-)collisions or intersection of body meshes. These two datasets together comprise the *Human Optical Flow* dataset. 3) We fine-tune SPyNet [18] on SHOF and show that its performance improves by about 43% (over the initial SPyNet), while it also outperforms existing state of the art by about 30%. Furthermore, we fine-tune SPyNet and PWC-Net on MHOF and observe improvements of 10 – 20% (over the initial SPyNet and PWC-Net). Compared to existing state of the art, improvements are particularly high for human regions. After masking out the background, we observe improvements of up to 13% for human pixels. 4) We provide the dataset files, dataset rendering code, training code and trained models¹ for research purposes.

¹ <https://humanflow.is.tue.mpg.de>

2 Related Work

Human Motion. Human motion can be understood from 2D motion. Early work focused on the movement of 2D joint locations [19] or simple motion history images [20]. Optical flow is also a useful cue. Black et al. [21] use principal component analysis (PCA) to parametrize human motion but use noisy flow computed from image sequences for training data. More similar to us, Fablet and Black [22] use a 3D articulated body model and motion capture data to project 3D body motion into 2D optical flow. They then learn a view-based PCA model of the flow fields. We use a more realistic body model to generate a large dataset and use this to train a CNN to directly estimate dense human flow from images.

Only a few works in pose estimation have exploited human motion and, in particular, several methods [23, 24] use optical flow constraints to improve 2D human pose estimation in videos. Similar work [25, 26] propagates pose results temporally using optical flow to encourage time consistency of the estimated bodies. Apart from its application in warping between frames, the structural information existing in optical flow alone has been used for pose estimation [27] or in conjunction with an image stream [28, 29].

Learning Optical Flow. There is a long history of optical flow estimation, which we do not review here. Instead, we focus on the relatively recent literature on learning flow. Early work looked at learning flow using Markov Random Fields [30], PCA [31], or shallow convolutional models [32]. Other methods also combine learning with traditional approaches, formulating flow as a discrete [33] or continuous [34] optimization problem.

The most recent methods employ large datasets to estimate optical flow using deep neural networks. Voxel2Voxel [35] is based on volumetric convolutions to predict optical flow using 16 frames simultaneously but does not perform well on benchmarks. Other methods [2, 10, 11] compute two frame optical flow using an end-to-end deep learning approach. FlowNet [10] uses the Flying Chairs dataset [10] to compute optical flow in an end-to-end deep network. FlowNet 2.0 [11] uses stacks of networks from FlowNet and performs significantly better, particularly for small motions. Ranjan and Black [2] propose a Spatial Pyramid Network that employs a small neural network on each level of an image pyramid to compute optical flow. Their method uses a much smaller number of parameters and achieves similar performance as FlowNet [10] using the same training data. Sun et al. [3] use image features in a similar spatial pyramid network achieving state-of-the-art results on optical flow benchmarks. Since the above methods are not trained with human motions, they do not perform well on our Human Optical Flow dataset.

Optical Flow Datasets. Several datasets have been developed to facilitate training and benchmarking of optical

flow methods. Middlebury is limited to small motions [1], KITTI is focused on rigid scenes and automotive motions [8], while Sintel has a limited number of synthetic scenes [9]. These datasets are mainly used for evaluation of optical flow methods and are generally too small to support training neural networks.

To learn optical flow using neural networks, more datasets have emerged that contain examples on the order of tens of thousands of frames. The Flying Chairs [10] dataset contains about 22,000 samples of chairs moving against random backgrounds. Although it is not very realistic or diverse, it provides training data for neural networks [2, 10] that achieve reasonable results on optical flow benchmarks. Even more recent datasets [12, 13] for optical flow are especially designed for training deep neural networks. Flying Things [12] contains tens of thousands of samples of random 3D objects in motion. The Monkaa and Driving scene datasets [12] contain frames from animated scenes and virtual driving respectively. Virtual KITTI [13] uses graphics to generate scenes like those in KITTI and is two orders of magnitude larger. Recent synthetic datasets [36] show that synthetic data can train networks that generalize to real scenes.

For human bodies, the SURREAL dataset [37] uses 3D human meshes rendered on top of color images to train networks for depth estimation, and body part segmentation. While not fully realistic, they show that this data is sufficient to train methods that generalize to real data. In a similar fashion, [38, 39] and [40] render synthetic color images for 3D hand pose estimation and 3D hand-object reconstruction, accordingly. We go beyond these works to address the problem of optical flow.

3 The Human Optical Flow Dataset

Our approach generates a realistic dataset of synthetic human motions by simulating them against different realistic backgrounds. We use parametric models [15, 41] to generate synthetic humans with a wide variety of different human shapes. We employ Blender² and its Cycles rendering engine to generate realistic synthetic image frames and optical flow. In this way we create the *Human Optical Flow* dataset, that is comprised of two parts. We first create the *Single-Human Optical Flow* (SHOF) dataset [18] using the body-only SMPL model [41] in images containing a single synthetic human. However, image statistics are different for the single- and multi-person case, as multiple people tend to occlude each other in complicated ways. For this reason we then create the *Multi-Human Optical Flow* (MHOF) dataset to better capture this realistic interaction. To make images even more realistic for MHOF, we replace SMPL [41] with the SMPL+H [15] model that models the body together with articulated fingers,

² <https://www.blender.org>

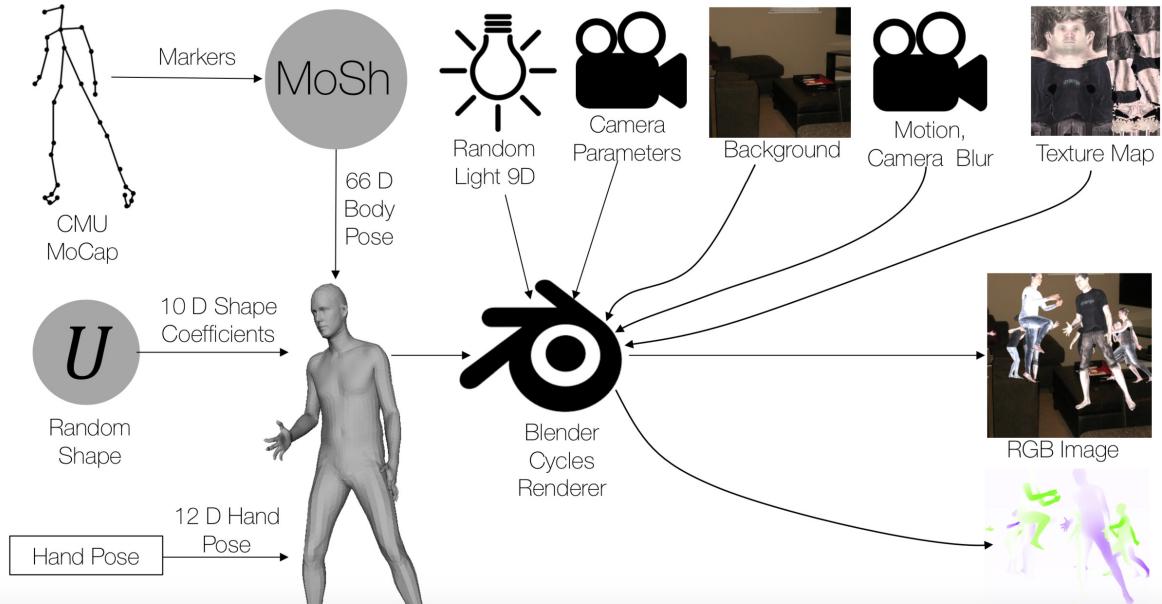


Fig. 2 Pipeline for generating the RGB frames and ground truth optical flow for the *Multi-Human Optical Flow dataset*. The datasets used in this pipeline are listed in Table 1, while the various rendering component are summarized in Table 2.

to have richer motion variation. In the rest of this section, we describe the components of our rendering pipeline, shown in Figure 2. For easy reference, in Table 1 we summarize the data used to generate the SHOF and MHOF datasets, while in Table 2 we summarize the various tools, Blender passes and parameters used for rendering. In the rest of the section, we describe the modules used for generating the data.

3.1 Human Body Generation

Body Model. A parametrized body model is necessary to generate human bodies in a scene. In the SHOF dataset, we use SMPL [41] for generating human body shapes. For the MHOF dataset we, use SMPL+H [15] that parametrizes the human body together with articulated fingers, for increased realism. The models are parameterized by pose and shape parameters to change the body posture and identity, as shown in Figure 2. They also contain a UV appearance map that allows us to change the skin tone, face features and clothing texture of the resulting virtual humans.

Body Poses. The next step is articulating the human body with different poses, to create moving sequences. To find such poses, we use 3D MoCap datasets [42, 43, 44] that capture 3D MoCap marker positions, glued onto the skin surface of real human subjects. We then employ MoSh [16, 17] that fits our body model to these 3D markers by optimizing over parameters of the body model for articulated pose, translation and shape. The pose specifically is a vector of axis-angle parameters, that describes how to rotate each body part around its corresponding skeleton joint.

For the SHOF dataset, we use the Human3.6M dataset [42], that contains five subjects for training (S1, S5, S6, S7, S8) and two for testing (S9, S11). Each subject performs 15 actions twice, resulting in 1,559,985 frames for training and 550,727 for testing. These sequences are subsampled at a rate of 16 \times , resulting in 97,499 training and 34,420 testing poses from Human3.6M.

For the MHOF dataset, we use the CMU [43] and HumanEva [44] MoCap datasets to increase motion variation. From CMU MoCap dataset, we use 2,605 sequences of 23 high-level action categories. From the HumanEva dataset, we use more than 10 sequences performing actions from 6 different action categories. To reduce redundant poses and allow for larger motions between frames, sequences are subsampled to 12 fps resulting in 321,873 poses. As a result the final MHOF dataset has 254,211 poses for training, 32,670 for validation and 34,992 for testing.

Hand Poses. Traditionally MoCap systems and datasets [42, 43, 44] record the motion of body joints, and avoid the tedious capture of detailed hand and finger motion. However, in natural settings, people use their body, hands and fingers to communicate social cues and to interact with the physical world. To enable our methods to learn such subtle motions, it should be represented in our training data. Therefore, we use the SMPL+H model [15] and augment the body-only MoCap datasets, described above, with finger motion. Instead of using random finger poses that would generate unrealistic optical flow, we employ the *Embodied Hands* dataset [15] and sample continuous finger motion to generate realistic optical flow. We use 43 sequences of hand motion with 37,232 frames recorded at 60 Hz by [15]. Similarly to body MoCap, we subsample hand MoCap to 12 fps to reduce overlapping poses without sacrificing variability.

Body Shapes. Human bodies vary a lot in their proportions, since each person has a unique body shape. To represent this in our dataset, we first learn a gender specific Gaussian distribution of shape parameters, by fitting SMPL to 3D CAESAR scans [45] of both genders. We then sample random body shapes from this distribution to generate a large number of realistic body shapes for rendering. However, naive sampling can result in extreme and unrealistic shape parameters, therefore we bound the shape distribution to avoid unlikely shapes.

For the SHOF dataset we bound the shape parameters to the range of $[-3, 3]$ standard deviations for each shape coefficient and draw a new shape for every subsequence of 20 frames to increase variance.

For the MHOF dataset, we account explicitly for collisions and intersections, since intersecting virtual humans would result in generation of inaccurate optical flow. To minimize such cases, we use similar sampling as above with only small differences. We first use shorter subsequences of 10 frames for less frequent inter-human intersections. Furthermore, we bound the shape distribution to the narrower range of $[-2.7, 2.7]$ standard deviations, since re-targeting motion to unlikely body shapes is more prone to mesh self-intersections.

Body Texture. We use the CAESAR dataset [45] to generate a variety of human skin textures. Given SMPL registrations to CAESAR scans, the original per-vertex color in the CAESAR dataset is transferred into the SMPL texture map. Since fiducial markers were placed on the bodies of CAESAR subjects, we remove them from the textures and inpaint them to produce a natural texture. In total, we use 166 CAESAR textures that are of good quality. The main drawback of CAESAR scans is their homogeneity in terms of outfit, since all of the subjects wore grey shorts and the women wore sports bras. In order to increase the clothing variety, we also use textures extracted from our 3D scans (referred as non-CAESAR in the following), to which we register SMPL with 4Cap [51]. A total of 772 textures from 7 different subjects with different clothes were captured. We anonymized the textures by replacing the face by the average face in CAESAR, after correcting it to match the skin tone of the texture. Textures are grouped according to the gender, which is randomly selected for each virtual human.

For the SHOF dataset the textures were split in training and testing sets with a 70/30 ratio, while each texture dataset is sampled with a 50% chance. For the MHOF dataset, we introduce more refined splitting with a 80/10/10 ratio for the train, validation and test sets. Moreover, since we introduce also finger motion, we want to favour sampling non-CAESAR textures, due to the bad quality of CAESAR texture maps for the finger region. Thus each texture is sampled with equal probability.

Hand Texture. Hands and fingers are hard to be scanned due to occlusions and measurement limitations. As a result, texture maps are particularly noisy or might even have holes. Since texture is important for optical flow, we augment the body texture maps to improve hand regions. For this we follow a divide and conquer approach. First, we capture hand-only scans with a 3dMD scanner [15]. Then, we create hand-only textures using the MANO model [15], getting 176 high resolution textures from 20 subjects. Finally, we use the hand-only textures to replace the problematic hand regions in the full-body texture maps.

We also need to find the best matching hand-only texture for every body texture. Therefore, we convert all texture maps in HSV space, and compute the mean HSV value for each texture map from standard sampling regions. For full body textures, we sample face regions without facial hair; while for hand-only textures, we sample the center of the outer palm. Then, for each body texture map we find the closest hand-only texture map in HSV space, and shift the values of the latter by the HSV difference, so that the hand skin tone becomes more similar to the facial skin tone. Finally, this improved hand-only texture map is used to replace the pixels in the hand-region of the full body texture map.

(Self-) Collision. The MHOF dataset contains multiple virtual humans moving differently, so there are high chances of collisions and penetrations. This is undesirable because penetrations are physically implausible and unrealistic. Moreover, the generated ground truth optical flow might have artifacts. Therefore, we employ a collision detection method to avoid intersections and penetrations.

Instead of using simple bounding boxes for rough collision detection, we draw inspiration from [52] and perform accurate and efficient collision detection on the triangle level using bounding volume hierarchies (BVH) [50]. This level of detailed detection allows for challenging occlusions with small distances between virtual humans, that can commonly be observed for realistic interactions between real humans. This method is useful not only for inter-person collision detection, but also for self-intersections. This is especially useful for our scenarios, as re-targeting body and hand motion to people of different shapes might result in unrealistic self-penetrations. The method is applicable out of the box, with the only exception that we exclude checks of neighboring body parts that are always or frequently in contact, e.g. upper and lower arm, or the two thighs.

3.2 Scene Generation

Background texture. For the scene background in the SHOF dataset, we use random indoor images from the LSUN dataset [47]. This provides a good compromise between simplicity and the complex task of generating varied full 3D environments. We use 417,597 images from the LSUN cate-

	SHOF	MHOF	Purpose
MoCap data	Human3.6M [42]	CMU [46], HumanEva [44]	Natural body poses
MoCap → SMPL	MoSh [16,17]	MoSh [16,17]	SMPL parameters from MoCap
Training poses	97,499	254,211	Articulate virtual humans
Validation poses	–	32,670	Articulate virtual humans
Test poses	34,420	34,992	Articulate virtual humans
Hand pose dataset	–	Embodied Hands [15]	Natural finger poses
Body shapes	Sample Gaussian distr. (CAESAR) bounded within $[-3, 3]$ st.dev.	Sample Gaussian distr. (CAESAR) bounded within $[-2.7, 2.7]$ st.dev.	Body proportions of virtual humans
Textures	CAESAR, non-CAESAR	CAESAR (hands improved), non-CAESAR (hands improved)	Appearance of virtual humans
Background	LSUN [47] (indoor) 417,597 images	SUN397 [48] (indoor and outdoor) 30,022 images	Scene background

Table 1 Comparison of datasets and most important data preprocessing steps used to generate the SHOF and MHOF datasets. A short description of the respective part is provided in the last column.

	SHOF	MHOF	Purpose
Rendering	Cycles	Cycles	Synthetic RGB image rendering
Optical flow	Vector pass (Blender)	Vector pass (Blender)	Optical flow ground truth
Segment. masks	Material pass (Blender)	Material pass (Blender)	Body part segment. masks (Fig. 3)
Motion blur	Vector pass (Blender)	Vector pass (Blender)	Realistic motion blur artifacts
Imaging noise	Gaussian blur (pixel space) 1px std.dev. for 30% of images	Gaussian blur (pixel space) 1px std.dev. for 30% of images	Realistic image imperfections
Camera translation	Sampled for 30% of frames from Gaussian with 1 cm std.dev.	Sampled for 30% of subsequences from Gaussian with 1 cm std.dev.	Realistic perturbations of the camera (and resulting optical flow)
Camera rotation	Sampled per frame from Gaussian with 0.2 degrees std.dev.	–	Realistic perturbations of the camera (and resulting optical flow)
Illumination	Spherical harmonics [49]	Spherical harmonics [49]	Realistic lighting model
Subsequence length	20 frames	10 frames	Number of successive frames with consistent rendering parameters
Mesh collision	–	BVH [50]	Detect (self-)collisions on the triangle level to avoid defect Optical Flow

Table 2 Comparison of tools, Blender passes and parameters used to generate the SHOF and MHOF datasets. The last column provides a short description of the respective method.

gories kitchen, living room, bedroom and dining room. These images are placed as billboards, 9 meters from the camera, and are not affected by the spherical harmonics lighting.

In the MHOF dataset, we increase the variability in background appearance. We employ the Sun397 dataset [48] that contains images for 397 highly variable scenes that are both indoor and outdoor, in contrast to LSUN. For quality reasons, we reject all images with resolution smaller than 512×512 px, and also reject images that contain humans using mask-RCNN [53,54]. As a result, we use 30,222 images, split in 24,178 for the training set and 3,022 for each of the validation and test sets. Further, we increase the distance between the camera and background to 12 meters, to increase the space in which the multiple virtual humans can move without colliding frequently to each other, while still being close enough for visual occlusions.

Scene Illumination. We illuminate the bodies with Spherical Harmonics lighting [49] that define basis vectors for light directions. This parameterization is useful for randomizing

the scene light by randomly sampling the coefficients with a bias towards natural illumination. The coefficients are uniformly sampled between -0.7 and 0.7 , apart from the ambient illumination, which has a minimum value of 0.3 to avoid extremely dark images, and illumination direction, which is strictly negative to favour illumination coming from above.

Increasing Image Realism. In order to increase realism, we introduced three types of image imperfections. First, for 30% of the generated images we introduced camera motion between frames. This motion perturbs the location of the camera with Gaussian noise of 1 cm standard deviation between frames and rotation noise of 0.2 degrees standard deviation per dimension in an Euler angle representation. Second, we add motion blur to the scene using the Vector Blur Node in Blender, and integrated over 2 frames sampled with 64 steps between the beginning and end point of the motion. Finally, we add a Gaussian blur to 30% of the images with a standard deviation of 1 pixel.

Scene Compositing. For animating virtual humans, each MoCap sequence is selected at least once. To increase variability, each sequence is split into subsequences. For the first frame of each subsequence, we sample a body and background texture, lights, blurring and camera motion parameters, and re-position virtual humans on the horizontal plane. We then introduce a random rotation around the z -axis for variability in the motion direction.

For the SHOF dataset, we use subsequences of 20 frames, and at the beginning of each one the single virtual human is re-positioned in the scene such that the pelvis is projected onto the image center.

For the MHOF dataset, we increase the variability with smaller subsequences of 10 frames and introduce more challenging visual occlusions by uniformly sampling the number of virtual humans in the range [4, 8]. We sample MoCap sequences S_j with a probability of $p_j = \frac{|S_j|}{\sum_{i=1}^{|S|} |S_i|}$, where $|S_j|$ denotes the number of frames of sequence S_j and $|S|$ the number of sequences. In contrast to the SHOF dataset, for the MHOF dataset the virtual humans are not re-positioned at the center, as they would all collide. Instead, they are placed at random locations on the horizontal plane within camera visibility, making sure there are no collisions with other virtual humans or the background plane during the whole subsequence.

3.3 Ground Truth Generation

Segmentation Masks. Using the material pass of Blender, we store for each frame the ground truth body part segmentation for our models. Although the body part segmentation for both models is similar, SMPL models the palm and fingers as one part, while SMPL+H has a different part segment for each finger bone. Figure 3 shows an example body part segmentation for SMPL+H. These segmentation masks allow us to perform a per body-part evaluation of our optical flow estimation.

Rendering & Ground Truth Optical Flow. For generating images, we use the open source suite Blender and its *vector pass*. The render pass is typically used for producing motion blur, and it produces the motion in image space of every pixel; i.e. the ground truth optical flow. We are mainly interested in the result of this pass, together with the color rendering of the textured bodies.

4 Learning

We train two different network architectures to estimate optical flow on both the SHOF and MHOF dataset. We choose compact models that are based on spatial pyramids, namely SPyNet [2] and PWC-Net [3], shown in Figure 4. We denote the models trained on the SHOF dataset by SPyNet+SHOF

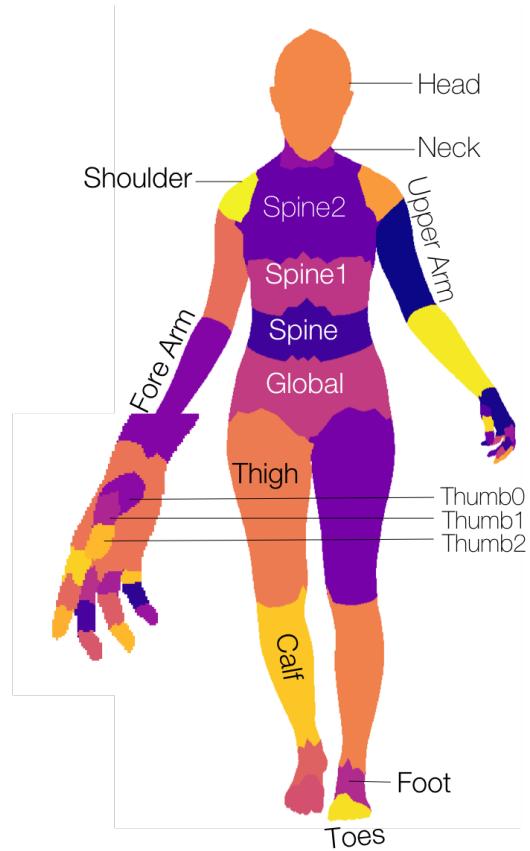


Fig. 3 Body part segmentation for the SMPL+H model. Symmetrical body parts are labeled only once. Finger joints follow the same naming convention as shown for the thumb. (Best viewed in color)

and PWC+SHOF. Similarly, we denote models trained on the MHOF dataset by SPyNet+MHOF and PWC+MHOF.

The spatial pyramid structure employs a convnet at each level of an image pyramid. A pyramid level works on a particular resolution of the image. The top level works on the full resolution and the image features are downsampled as we move to the bottom of the pyramid. Each level learns a convolutional layer d , to perform downsampling of image features. Similarly, a convolution layer u , is learned for decoding optical flow. At each level, the convnet G_k predicts optical flow residuals v_k at that level. These flow residuals get added at each level to produce the full flow, V_K at the finest level of the pyramid.

In SPyNet, each convnet G_k takes a pair of images as inputs along with flow V_{k-1} obtained by upsampling the output of the previous level. The second frame is however warped using V_{k-1} and the triplet $\{I_k^1, w(I_k^2, V_{k-1}), V_{k-1}\}$ is fed as input to the convnet G_k .

In PWC-Net, a pair of image features, $\{I_k^1, I_k^2\}$ is input at a pyramid level, and the second feature map is warped using the flow V_{k-1} from the previous level of the pyramid. We then compute the cost-volume $c(I_k^1, w(I_k^2, V_{k-1}))$ over

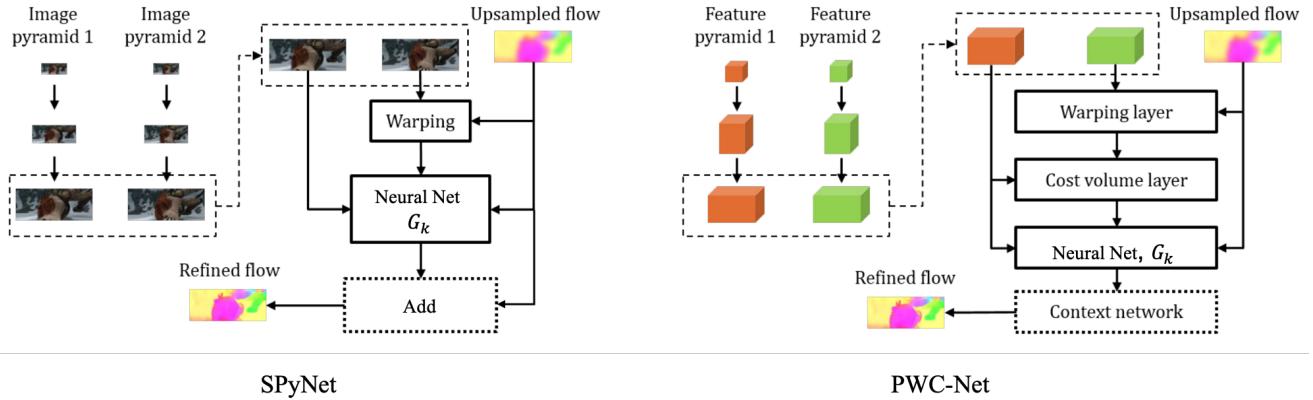


Fig. 4 Spatial Pyramid Network [2] (left) and PWC-Net [3] (right) for optical flow estimation. At each pyramid level, network G_k predicts flow at that level which is used to condition the optical flow at the higher resolution level in the pyramid. Adapted from [3].

feature maps and pass it to network G_k to compute optical flow V_k at that pyramid level.

We use the pretrained weights as initializations for training both SPyNet and PWC-Net. We train both models end-to-end to minimize the average End Point Error (EPE).

Hyperparameters. We follow the same training procedure for SPyNet and PWC-Net. The only exception to this is the learning rate, which is determined empirically for each dataset and network from $\{10^{-6}, 10^{-5}, 10^{-4}\}$. For the SHOF we found 10^{-6} to yield best results for SpyNet. Predictions of PWC on the SHOF dataset do not improve for any of these learning rates. For training on MHOF a learning rate of 10^{-6} and 10^{-4} yield best results for SpyNet and PWC-Net, respectively. We use Adam [55] to optimize our loss with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. We use a batch size of 8 and run 400,000 training iterations. All networks are implemented in the Pytorch framework. Fine-tuning the networks from pretrained weights takes approximately 1 day on SHOF and 2 days on MHOF.

Data Augmentations. We also augment our data by applying several transformations and adding noise. Although our dataset is quite large, augmentation improves the quality of results on real scenes. In particular, we apply scaling in the range of $[0.3, 3]$, and rotations in $[-17^\circ, 17^\circ]$. The dataset is normalized to have zero mean and unit standard deviation using [56].

5 Experiments

In this section, we first compare the SHOF, MHOF and other common optical flow datasets. Next, we show that fine-tuning SPyNet on SHOF improves the model, while we observe that fine-tuning PWC-Net on SHOF does not improve the model further. We then fine-tune the same methods on MHOF and evaluate them. We show that both, SPyNet and PWC-Net improve when fine-tuned on MHOF. We show that the meth-

Dataset	# Train Frames	# Test Frames	Resolution
MPI Sintel [9]	1,064	564	1024×436
KITTI 2012 [8]	194	195	1226×370
KITTI 2015 [57]	200	200	1242×375
Virtual Kitti [13]	21,260	—	1242×375
Flying Chairs [10]	22,232	640	512×384
Flying Things [12]	21,818	4,248	960×540
Monkaa [12]	8,591	—	960×540
Driving [12]	4,392	—	960×540
SHOF (ours)	135,153	10,867	256×256
MHOF (ours)	86,259	13,236	640×640

Table 3 Comparison of the *Human Optical Flow* datasets, namely the *Single-Human Optical Flow* (SHOF) and the *Multi-Human Optical Flow* (MHOF) dataset, with previous optical flow datasets.

ods trained on the MHOF dataset outperform generic flow estimation methods for the pixels corresponding to humans. Finally, we show on qualitative results that both, the models trained on SHOF and models trained on MHOF seem to generalize to real word scenes.

Dataset Details. In comparison with other optical flow datasets, our dataset is larger by an order of magnitude (see Table 3); the SHOF dataset contains 135,153 training frames and 10,867 test frames with optical flow ground truth, while the MHOF dataset has 86,259 training, 13,236 test and 11,817 validation frames. For the single-person dataset we keep the resolution small at 256×256 px to facilitate easy deployment for training neural networks. This also speeds up the rendering process in Blender for generating large amounts of data. We show the comparisons of processing time of different models on the SHOF dataset in Table 4(a). For the MHOF dataset we increase the resolution to 640×640 px to be able to reason about optical flow even in small body parts like fingers, using SMPL+H. Our data is extensive, containing a wide variety of human shapes, poses, actions and virtual backgrounds to support deep learning systems.

Comparison on SHOF. We compare the average End Point Errors (EPEs) of optical flow methods on the SHOF dataset in Table 4, along with the time for evaluation. We show visual comparisons in Figure 5. Human motion is complex and general optical flow methods fail to capture it. Our trained network SPyNet+SHOF outperforms previous methods, and SPyNet [2] in particular.

We observe that FlowNet [10] shows poor generalization on our dataset. Since the results of FlowNet [10] in Table 4 and 6 are very close to the zero flow (no motion) baseline, we cross-verify by evaluating FlowNet on a mixture of Flying Chairs [10] and *Human Optical Flow* and observe that the flow outputs on SHOF is quite random (see Figure 5). The main reason is that SHOF contains a significant amount of small motions and it is known that FlowNet does not perform very well on small motions. SPyNet [2] however performs quite well and is able to generalize to body motions. The results however look noisy in many cases.

Our dataset employs a layered structure where a human is placed against a background. As such layered methods like PCA-layers [31] perform very well on a few images (row 8 in Figure 5) where they are able to segment a person from the background. However, in most cases, they do not obtain good segmentation into layers.

Previous state-of-the-art methods like LDOF [58] and Epic-Flow[34] perform much better than others. They get a good overall shape, and smooth backgrounds. However, their estimation is quite blurred. They tend to miss the sharp edges that are typical of human hands and legs. They are also significantly slower.

In contrast, by fine-tuning on our dataset, the performance of SPyNet+SHOF improves by 40% over SPyNet on the SHOF dataset. We also find that fine-tuning PWC-Net on the SHOF does not improve the model. Empirically, we have seen that PWC-Net already performs well for small motions, that are a significant portion of SHOF. This partially motivates the generation of the MHOF dataset, which includes larger motions and more complex scenes with occlusions.

A qualitative comparison to popular optical flow methods can be seen in Figure 5. Flow estimations of SPyNet+MHOF can be observed to be sharper than those of generic methods. This can especially be seen for edges. Furthermore, it can be seen that fine details like motion of hands are estimated more precisely.

Comparison on MHOF. Training (fine-tuning) on the MHOF dataset improves SPyNet and PWC-Net on average, as can be seen in Table 5. In particular PWC+MHOF outperforms SPyNet+MHOF and also improves over generic state-of-the-art optical flow methods. Large parts of the image are background, whose movements are relatively easy to estimate. However, we are particularly interested in human motions. Therefore, we mask out all errors of background pixels and compute the average EPE only on body pixels (see

Method	AEPE	Time(s)	Learned	Fine-tuned on SHOF
Zero	0.6611	-	-	
FlowNet [10]	0.5846	0.080	✓	✗
PCA Layers [31]	0.3652	10.357	✗	✗
PWC-Net [3]	0.2158	0.024	✓	✗
PWC+SHOF	0.2158	0.024	✓	✓
SPyNet [2]	0.2066	0.022	✓	✗
Epic Flow [34]	0.1940	1.863	✗	✗
LDOF [58]	0.1881	8.620	✗	✗
FlowNet2 [11]	0.1895	0.127	✓	✗
Flow Fields [59]	0.1709	4.204	✗	✗
SPyNet+SHOF	0.1164	0.022	✓	✓

Table 4 EPE comparisons and evaluation times of different optical flow methods on the SHOF dataset. Zero refers to the EPE when zero flow (no motion) is always used for evaluation. Evaluation times are based on the SHOF dataset with 256×256 image resolution. We time all GPU based methods using a Tesla V100-16GB GPU.

Method	Average EPE	Average EPE on <i>body pixels</i>	Fine-tuned on MHOF
FlowNet	0.808	2.574	✗
PCA Layers	0.556	2.691	✗
Epic Flow	0.488	1.982	✗
SPyNet	0.429	1.977	✗
SPyNet+MHOF	0.391	1.803	✓
PWC-Net	0.369	2.056	✗
LDOF	0.360	1.719	✗
FlowNet2	0.310	1.863	✗
PWC+MHOF	0.306	1.620	✓

Table 5 Comparison using End Point Error (EPE) on the *Multi-Human Optical Flow* (MHOF) dataset. We show the average EPE and body-only EPE. The EPE is computed only over segments of the image depicting a human body. Best results are shown in boldface. A comparison of body-part specific EPE can be found in Table 6 .

Table 5). For these pixels, light-weight networks like SpyNet and PWC-Net improve over almost all generic optical flow estimation methods using our dataset (SpyNet+MHOF and PWC+MHOF), including the much larger network FlowNet2. PWC+MHOF is the best performing method.

A more fine grained analysis of EPE across body parts is shown in Table 6. We obtain EPE of these body parts using the segmentation shown in Figure 3. It can be seen that improvements of PWC+MHOF over FlowNet2 are larger for body parts that are at the end of the kinematic tree (i.e. feet, calves, arms and in particular fingers). Differences are less strong for body parts close to the torso. One interpretation of these findings is that movements of the torso are easier to predict, while movements of body parts at the end of the kinematic tree are more complex and thus harder to estimate. In contrast, SPyNet+MHOF outperforms FlowNet2 on body parts close to the torso and does not learn to capture the more complex motions of limbs better than FlowNet2.

Visual comparisons are shown in Figure 6. In particular, PWC+MHOF predicts flow fields with sharper edges than generic methods or SPyNet+ MHOF. Furthermore, the

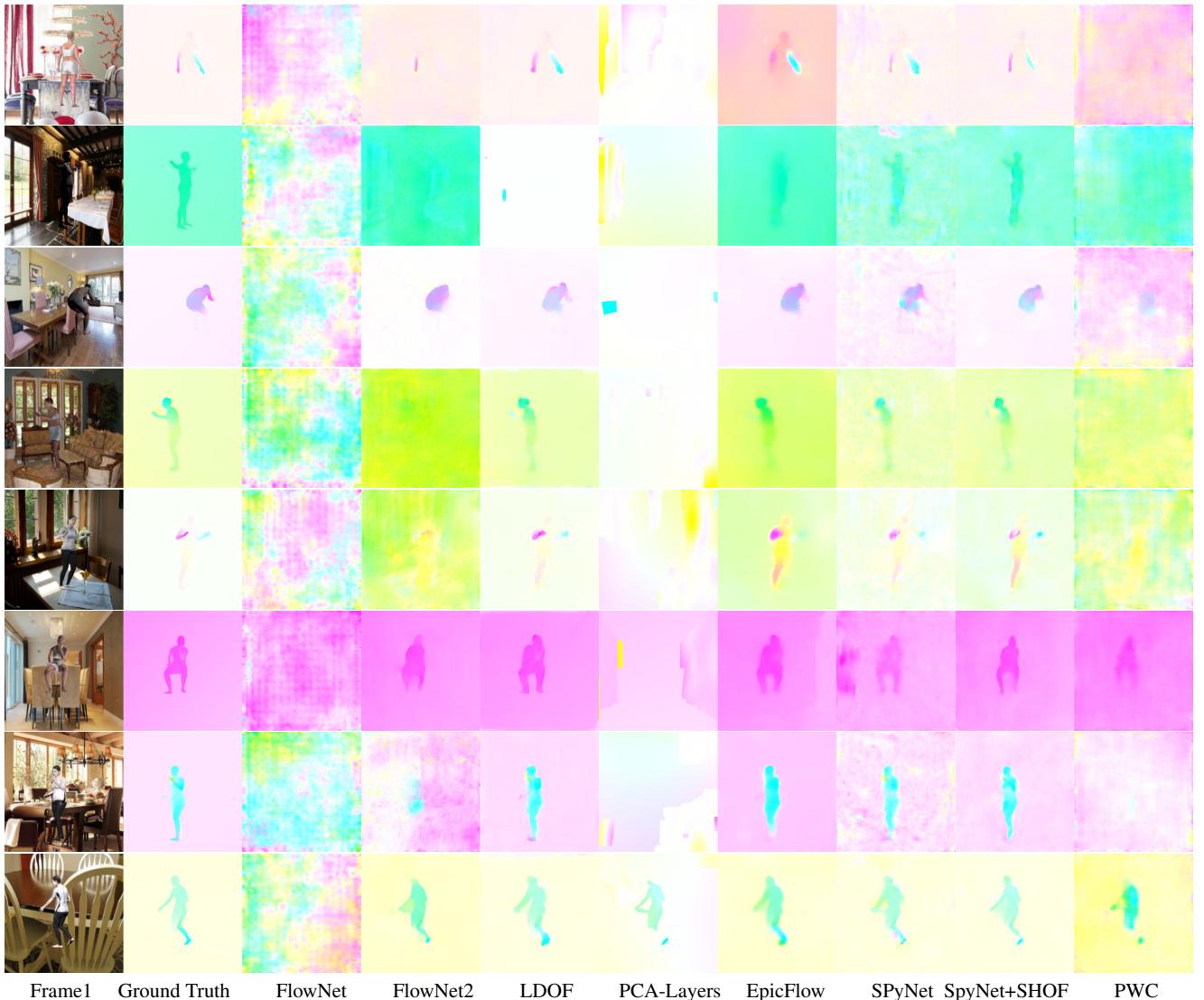


Fig. 5 Visual comparison of optical flow estimates using different methods on the *Single-Human Optical Flow* (SHOF) test set. From left to right, we show Frame 1, Ground Truth flow, results on FlowNet [10], FlowNet2 [11], LDOF [58], PCA-Layers [31], SPyNet [2], EpicFlow [34] , LDOF [58], SPyNet [2], SPyNet+SHOF (ours) and PWC-Net [3]

qualitative results suggest that PWC+MHOFR is better at distinguishing the motion of people, as people can be better separated on the flow visualizations of PWC+MHOFR (Figure 6, row 3). Last, it can be seen that fine details, like the motion of distant humans or small body parts, are better estimated by PWC+MHOFR.

The above observations are strong indications that our *Human Optical Flow* datasets (SHOF and MHOFR) can be beneficial for the performance on human motion for other optical flow networks as well.

Real Scenes. We show a visual comparison of results on real-world scenes of people in motion. For visual comparisons of models trained on the SHOF dataset we collect these scenes by cropping people from real world videos as shown in Figure 7. We use DPM [60] for detecting people and compute

bounding box regions in two frames using the ground truth of the MOT16 dataset [61]. The results for the SHOF dataset are shown in Figure 8. A comparison of methods on real images with multiple people can be seen in Figure 9.

The performance of PCA-Layers [31] is highly dependent on its ability to segment. Hence, we see only a few cases where it looks visually correct. SPyNet [2] gets the overall shape but the results look noisy in certain image parts. While LDOF [58], EpicFlow [34] and FlowFields [59] generally perform well, they often find it difficult to resolve the legs, hands and head of the person. The results from models trained on our *Human Optical Flow* dataset look appealing especially while resolving the overall human shape, and various parts like legs, hands and the human head. Models trained on the *Human Optical Flow* dataset perform well under occlusion

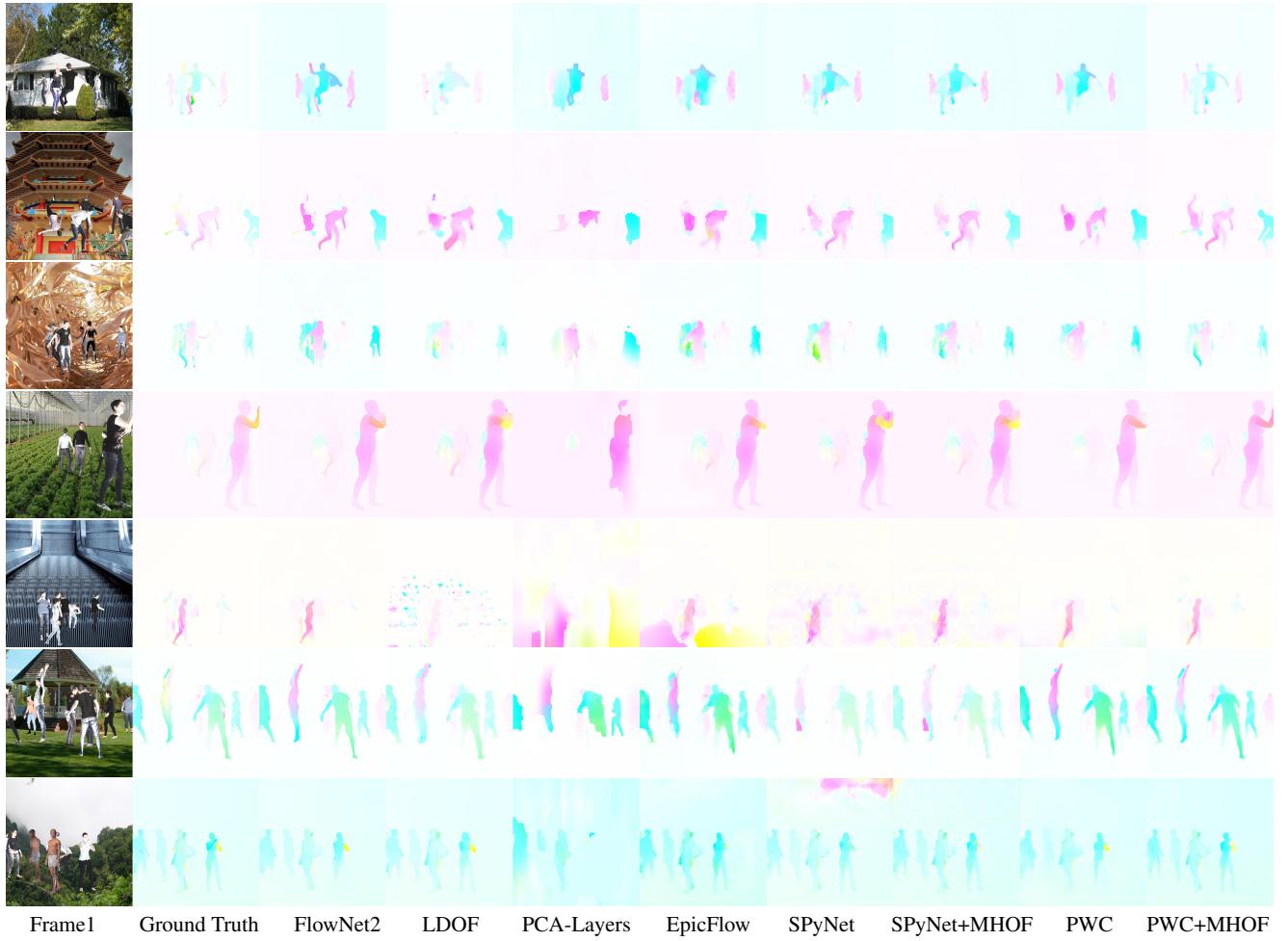


Fig. 6 Visual comparison of optical flow estimates using different methods on the *Multi-Human Optical Flow* (MHOF) test set. From left to right, we show Frame 1, Ground Truth flow, results on FlowNet2 [11], LDOF [58], PCA-Layers [31], EpicFlow [34], SPyNet [2], SPyNet+MHOF (ours), PWC-Net [3] and PWC+MHOF (ours).



Fig. 7 We use the DPM [60] person detector to crop out people from real-world scenes (left) and use SPyNet+SHOF to compute optical flow on the cropped section (right).

(Figure 8, Figure 9). Many examples including severe occlusion can be seen in Figure 9. Besides that, Figure 9 shows that the models trained on MHOF are able to distinguish motions of multiple people and predict sharp edges of humans.

6 Conclusion and Future Work

In summary, we created an extensive *Human Optical Flow* dataset containing images of realistic human shapes in motion together with ground truth optical flow. The dataset is comprised of two parts, the *Single-Human Optical Flow* (SHOF) and the *Multi-Human Optical Flow* (MHOF) dataset. We then train two compact network architectures based on spatial pyramids, namely SpyNet and PWC-Net. The realism and extent of our dataset, together with an end-to-end training scheme, allows these networks to outperform previous state-of-the-art optical flow methods on our new human-specific dataset. This indicates that our dataset can be beneficial for other optical flow network architectures as well. Furthermore, our qualitative results suggest that the networks trained on the *Human Optical Flow* generalize well to real world scenes with humans. The trained models are compact and run in real time making them highly suitable for phones and embedded devices.

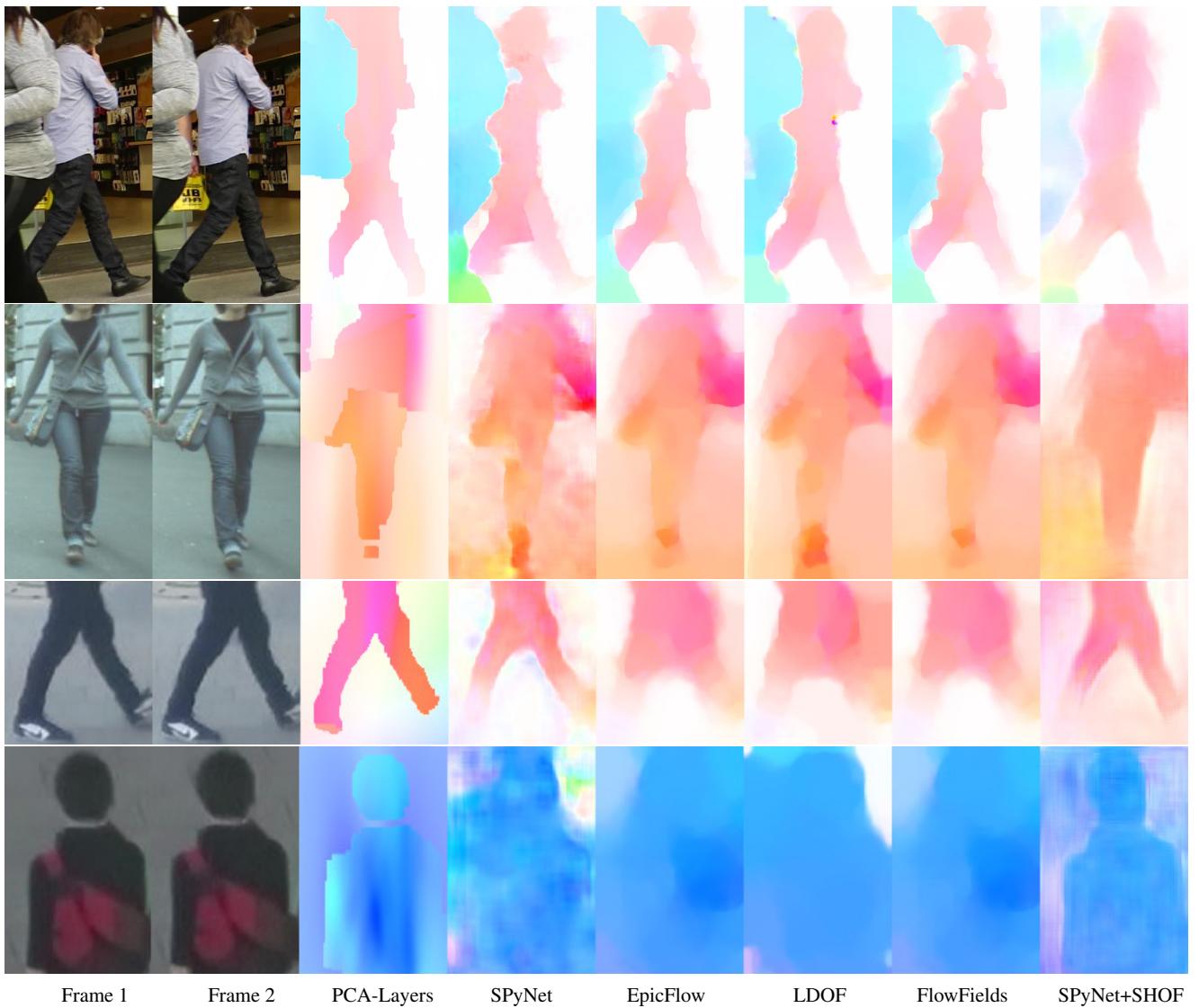


Fig. 8 Single-Human Optical Flow visuals on real images using different methods. From left to right, we show Frame 1, Frame 2, results on PCA-Layers [31], and SPyNet [2], EpicFlow [34], LDOF [58], FlowFields [59] and SPyNet+SHOF (ours).

In future work, we plan to add 3D clothing and accessories to humans in the scene. The dataset and our focus on human optical flow opens up a number of research directions in human motion understanding and optical flow computation. We would like to extend our dataset by modeling more diverse clothing and outdoor scenarios. A direction of potentially high impact for this work is to integrate it in end-to-end systems for action recognition, which typically take precomputed optical flow as input. The real-time nature of the method could support motion-based interfaces, potentially even on devices like cell phones with limited computing power. The dataset, dataset generation code, pretrained models, and training code are available, enabling researchers to use them for problems involving human motion.

Acknowledgements

We thank Yiyi Liao for helping us with optical flow evaluation. We thank Cristian Sminchisescu for the Human3.6M MoCap marker data.

References

1. Simon Baker, Daniel Scharstein, JP Lewis, Stefan Roth, Michael J Black, and Richard Szeliski. A database and evaluation methodology for optical flow. *International Journal of Computer Vision*, 92(1):1–31, 2011.
2. Anurag Ranjan and Michael J Black. Optical flow estimation using a spatial pyramid network. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
3. Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. In *CVPR*, 2018.

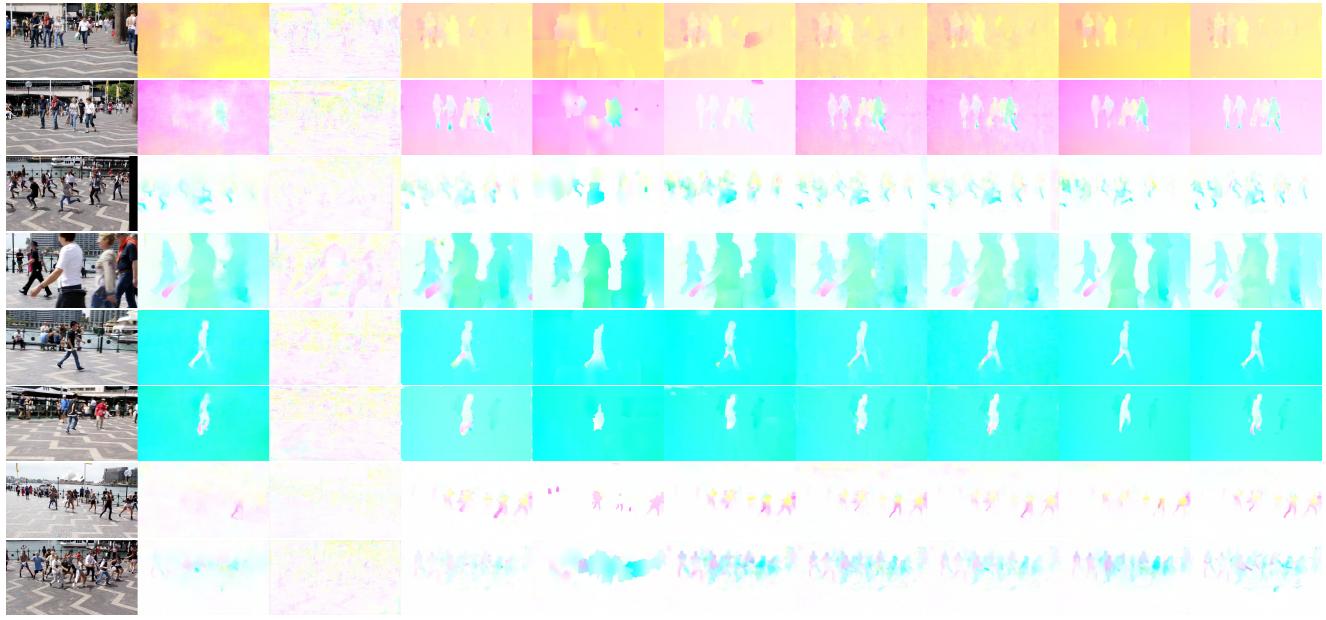


Fig. 9 *Multi-Human Optical Flow* visuals on real images. From left to right, we show Frame 1, results on FlowNet2 [11], FlowNet [10], LDOF [58], PCA-Layers [31], EpicFlow [34], SPyNet [2], SPyNet+MHOF (ours), PWC-Net [3] and PWC+MHOF (ours).

4. Donald Geman and Stuart Geman. Opinion: Science in the age of selfies. *Proceedings of the National Academy of Sciences*, 113(34):9384–9387, 2016.
5. Hueihan Jhuang, Juergen Gall, Silvia Zuffi, Cordelia Schmid, and Michael J. Black. Towards understanding action recognition. In *IEEE International Conference on Computer Vision (ICCV)*, pages 3192–3199, Sydney, Australia, December 2013. IEEE.
6. Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
7. Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large video database for human motion recognition. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 2556–2563. IEEE, 2011.
8. Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the KITTI vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
9. D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black. A naturalistic open source movie for optical flow evaluation. In A. Fitzgibbon et al. (Eds.), editor, *European Conf. on Computer Vision (ECCV)*, Part IV, LNCS 7577, pages 611–625. Springer-Verlag, October 2012.
10. Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Caner Hazirbas, Vladimir Golkov, Patrick van der Smagt, Daniel Cremers, Thomas Brox, et al. Flownet: Learning optical flow with convolutional networks. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 2758–2766. IEEE, 2015.
11. Eddy Ilg, Nikolaus Mayer, Tommoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. *arXiv preprint arXiv:1612.01925*, 2016.
12. N.Mayer, E.Ilg, P.Häusser, P.Fischer, D.Cremers, A.Dosovitskiy, and T.Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. arXiv:1512.02134.
13. A Gaidon, Q Wang, Y Cabon, and E Vig. Virtual worlds as proxy for multi-object tracking analysis. In *CVPR*, 2016.
14. Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J. Black. Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. In *Computer Vision – ECCV 2016*, Lecture Notes in Computer Science. Springer International Publishing, October 2016.
15. Javier Romero, Dimitrios Tzionas, and Michael J. Black. Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6), November 2017.
16. Matthew Loper, Naureen Mahmood, and Michael J Black. Mosh: Motion and shape capture from sparse markers. *ACM Transactions on Graphics (TOG)*, 33(6):220, 2014.
17. Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. AMASS: archive of motion capture as surface shapes. *CoRR*, abs/1904.03278, 2019.
18. Anurag Ranjan, Javier Romero, and Michael J. Black. Learning human optical flow. In *29th British Machine Vision Conference*, September 2018.
19. Gunnar Johansson. Visual perception of biological motion and a model for its analysis. *Perception & Psychophysics*, 14(2):201–211, 1973.
20. J. W. Davis. Hierarchical motion history images for recognizing human motion. In *Detection and Recognition of Events in Video*, pages 39–46, 2001.
21. M. J. Black, Y. Yacoob, A. D. Jepson, and D. J. Fleet. Learning parameterized models of image motion. In *IEEE Conf. on Computer Vision and Pattern Recognition, CVPR-97*, pages 561–567, Puerto Rico, June 1997.
22. R. Fablet and M. J. Black. Automatic detection and tracking of human motion with a view-based representation. In *European Conf. on Computer Vision, ECCV 2002*, volume 1 of *LNCS 2353*, pages 476–491. Springer-Verlag, 2002.
23. K. Fragiadaki, H. Hu, and J. Shi. Pose from flow and flow from pose. In *2013 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2059–2066, June 2013.
24. Silvia Zuffi, Javier Romero, Cordelia Schmid, and Michael J Black. Estimating human pose with flowing puppets. In *IEEE International Conference on Computer Vision (ICCV)*, pages 3312–3319, 2013.

25. Tomas Pfister, James Charles, and Andrew Zisserman. Flowing convnets for human pose estimation in videos. In *ICCV*, pages 1913–1921. IEEE Computer Society, 2015.
26. James Charles, Tomas Pfister, Derek R. Magee, David C. Hogg, and Andrew Zisserman. Personalizing human video pose estimation. In *CVPR*, pages 3063–3072. IEEE Computer Society, 2016.
27. Javier Romero, Matthew Loper, and Michael J. Black. FlowCap: 2D human pose from optical flow. In *Pattern Recognition, Proc. 37th German Conference on Pattern Recognition (GCPR)*, volume LNCS 9358, pages 412–423. Springer, 2015.
28. Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. Convolutional two-stream network fusion for video action recognition. In *CVPR*, pages 1933–1941. IEEE Computer Society, 2016.
29. Xuanyi Dong, Shouo-I Yu, Xinshuo Weng, Shih-En Wei, Yi Yang, and Yaser Sheikh. Supervision-by-registration: An unsupervised approach to improve the precision of facial landmark detectors. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
30. William T. Freeman, Egon C. Pasztor, and Owen T. Carmichael. Learning low-level vision. *International Journal of Computer Vision*, 40(1):25–47, 2000.
31. Jonas Wulff and Michael J Black. Efficient sparse-to-dense optical flow estimation using a learned basis and layers. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 120–130. IEEE, 2015.
32. Sun D., S Roth, JP Lewis, and MJ Black. Learning optical flow. In *ECCV*, pages 83–97, 2008.
33. Fatma Güney and Andreas Geiger. Deep discrete flow. In *Asian Conference on Computer Vision (ACCV)*, 2016.
34. Jerome Revaud, Philippe Weinzaepfel, Zaid Harchaoui, and Cordelia Schmid. EpicFlow: Edge-Preserving Interpolation of Correspondences for Optical Flow. In *Computer Vision and Pattern Recognition*, 2015.
35. Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Deep End2End Voxel2Voxel prediction. In *The 3rd Workshop on Deep Learning in Computer Vision*, 2016.
36. Adrien Gaidon, Zaid Harchaoui, and Cordelia Schmid. Activity representation with motion hierarchies. *International Journal of Computer Vision*, 107(3):219–238, 2014.
37. GÜl Varol, Javier Romero, Xavier Martin, Naureen Mahmood, Michael Black, Ivan Laptev, and Cordelia Schmid. Learning from synthetic humans. In *CVPR*, 2017.
38. Christian Zimmermann and Thomas Brox. Learning to estimate 3D hand pose from single rgb images. In *ICCV*, 2017.
39. Franziska Mueller, Florian Bernard, Oleksandr Sotnychenko, Dushyant Mehta, Srinath Sridhar, Dan Casas, and Christian Theobalt. GANerated hands for real-time 3D hand tracking from monocular RGB. In *CVPR*, 2018.
40. Yana Hasson, Gul Varol, Dimitrios Tzionas, Igor Kalevatykh, Michael J. Black, Ivan Laptev, and Cordelia Schmid. Learning joint reconstruction of hands and manipulated objects. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
41. Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, October 2015.
42. Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7):1325–1339, jul 2014.
43. Carnegie-mellon mocap database.
44. L. Sigal, A. Balan, and M. J. Black. HumanEva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion. *International Journal of Computer Vision*, 87(1):4–27, March 2010.
45. Kathleen M Robinette, Sherri Blackwell, Hein Daanen, Mark Boehmer, and Scott Fleming. Civilian american and european surface anthropometry resource (caesar), final report. volume 1. summary. Technical report, DTIC Document, 2002.
46. Ralph Gross and Jianbo Shi. The cmu motion of body (mobo) database. 2001.
47. Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv:1506.03365*, 2015.
48. Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *Computer vision and pattern recognition (CVPR), 2010 IEEE conference on*, pages 3485–3492. IEEE, 2010.
49. Robin Green. Spherical Harmonic Lighting: The Gritty Details. *Archives of the Game Developers Conference*, March 2003.
50. Matthias Teschner, Stefan Kimmerle, Bruno Heidelberger, Gabriel Zachmann, Laks Raghupathi, Arnulph Fuhrmann, Marie-Paule Cani, François Faure, Nadia Magnenat-Thalmann, Wolfgang Strasser, and Pascal Volino. Collision detection for deformable objects. In *Eurographics*, pages 119–139, 2004.
51. Gerard Pons-Moll, Javier Romero, Naureen Mahmood, and Michael J. Black. Dyna: A model of dynamic human shape in motion. *ACM Transactions on Graphics, (Proc. SIGGRAPH)*, 34(4):120:1–120:14, August 2015.
52. Dimitrios Tzionas, Luca Ballan, Abhilash Srikantha, Pablo Aponte, Marc Pollefeys, and Juergen Gall. Capturing hands in action using discriminative salient points and physics simulation. *International Journal of Computer Vision (IJCV)*, 118(2):172–193, June 2016.
53. Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 2980–2988. IEEE, 2017.
54. Waleed Abdulla. Mask r-cnn for object detection and instance segmentation on keras and tensorflow. https://github.com/matterport/Mask_RCNN, 2017.
55. Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
56. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.
57. Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3061–3070, 2015.
58. Thomas Brox, Christoph Bregler, and Jitendra Malik. Large displacement optical flow. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 41–48. IEEE, 2009.
59. Christian Bailer, Bertram Taetz, and Didier Stricker. Flow fields: Dense correspondence fields for highly accurate large displacement optical flow estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4015–4023, 2015.
60. P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *TPAMI*, 32(9):1627–1645, 2010.
61. Anton Milan, Laura Leal-Taixé, Ian D. Reid, Stefan Roth, and Konrad Schindler. MOT16: A benchmark for multi-object tracking. *arXiv:1603.00831*, 2016.

Parts	Epic Flow	LDOF	FlowNet2	FlowNet	PCA Layers	PWC-Net	PWC+MHOF	SPyNet	SPyNet+MHOF
Average (whole image)	0.488	0.360	0.310	0.808	0.556	0.369	0.306	0.429	0.391
Average (body pixels)	1.982	1.719	1.863	2.574	2.691	2.056	1.620	1.977	1.803
global	1.269	1.257	1.337	2.005	1.920	1.389	1.172	1.356	1.236
head	1.806	1.328	1.626	2.681	2.808	1.881	1.441	1.708	1.519
leftCalf	2.116	1.802	1.787	2.420	2.711	2.109	1.469	1.991	1.796
leftFoot	3.089	2.346	2.476	2.987	3.393	3.002	2.102	2.701	2.566
leftForeArm	3.972	3.231	3.536	4.380	4.778	3.926	3.117	3.945	3.605
leftHand	5.777	4.422	4.823	5.928	6.531	5.634	4.337	5.547	5.040
leftShoulder	1.513	1.429	1.646	2.331	2.336	1.732	1.468	1.560	1.462
leftThigh	1.424	1.338	1.466	2.102	2.150	1.565	1.241	1.517	1.362
leftToes	3.147	2.573	2.755	3.065	3.307	3.100	2.487	2.830	2.784
leftUpperArm	2.215	1.947	2.288	3.005	3.139	2.376	1.951	2.307	2.076
lIndex0	6.199	4.900	5.334	6.254	6.785	6.124	4.796	5.925	5.472
lIndex1	6.367	5.159	5.672	6.340	6.829	6.303	5.142	6.087	5.727
lIndex2	6.315	5.253	5.878	6.203	6.670	6.270	5.367	6.028	5.784
lMiddle0	6.338	4.983	5.331	6.364	6.910	6.211	4.786	6.012	5.544
lMiddle1	6.498	5.239	5.632	6.435	6.927	6.383	5.121	6.143	5.767
lMiddle2	6.266	5.212	5.756	6.130	6.592	6.182	5.245	5.934	5.679
lPinky0	6.048	4.792	5.302	6.035	6.603	5.940	4.833	5.738	5.307
lPinky1	6.106	4.922	5.489	6.038	6.574	6.014	5.012	5.765	5.418
lPinky2	5.780	4.856	5.419	5.655	6.170	5.702	4.905	5.474	5.231
lRing0	6.388	4.973	5.281	6.413	7.010	6.218	4.794	6.064	5.552
lRing1	6.313	5.083	5.391	6.256	6.801	6.168	4.899	5.966	5.558
lRing2	6.047	5.035	5.515	5.924	6.409	5.942	5.007	5.710	5.441
lThumb0	5.415	4.318	4.673	5.473	6.072	5.316	4.272	5.212	4.809
lThumb1	5.636	4.527	5.065	5.698	6.232	5.612	4.616	5.449	5.065
lThumb2	5.825	4.749	5.388	5.820	6.323	5.802	4.934	5.629	5.314
neck	1.336	1.195	1.371	2.151	2.245	1.440	1.227	1.399	1.250
rightCalf	2.243	1.892	1.864	2.530	2.851	2.223	1.539	2.081	1.907
rightFoot	3.270	2.454	2.610	3.149	3.599	3.171	2.237	2.894	2.732
rightForeArm	3.990	3.242	3.554	4.381	4.759	3.928	3.181	4.029	3.641
rightHand	5.735	4.348	4.787	5.837	6.447	5.550	4.307	5.582	4.978
rightShoulder	1.547	1.431	1.670	2.390	2.340	1.735	1.472	1.573	1.462
rightThigh	1.477	1.374	1.512	2.158	2.226	1.624	1.275	1.556	1.407
rightToes	3.395	2.707	2.918	3.293	3.566	3.346	2.679	3.064	2.999
rightUpperArm	2.267	1.974	2.294	3.033	3.148	2.400	2.007	2.346	2.113
rIndex0	6.264	4.875	5.324	6.255	6.800	6.150	4.833	6.003	5.486
rIndex1	6.541	5.210	5.755	6.449	6.951	6.457	5.269	6.237	5.835
rIndex2	6.465	5.320	5.968	6.294	6.776	6.404	5.473	6.149	5.879
rMiddle0	6.509	5.056	5.454	6.470	7.014	6.354	4.931	6.211	5.662
rMiddle1	6.680	5.341	5.777	6.562	7.058	6.537	5.277	6.325	5.895
rMiddle2	6.394	5.261	5.838	6.209	6.713	6.274	5.313	6.038	5.739
rPinky0	5.983	4.750	5.372	5.952	6.504	5.855	4.812	5.741	5.262
rPinky1	6.076	4.905	5.566	5.979	6.533	5.943	4.984	5.809	5.402
rPinky2	5.789	4.813	5.403	5.645	6.220	5.662	4.853	5.532	5.232
rRing0	6.397	4.948	5.350	6.383	6.938	6.215	4.828	6.126	5.565
rRing1	6.395	5.108	5.465	6.290	6.841	6.212	4.973	6.066	5.615
rRing2	6.222	5.129	5.644	6.052	6.610	6.063	5.110	5.889	5.571
rThumb0	5.417	4.304	4.748	5.470	6.057	5.301	4.292	5.247	4.819
rThumb1	5.605	4.465	4.945	5.643	6.210	5.514	4.536	5.434	5.032
rThumb2	5.835	4.748	5.262	5.789	6.328	5.749	4.874	5.639	5.306
spine	1.233	1.271	1.325	1.941	1.856	1.360	1.173	1.322	1.221
spine1	1.330	1.369	1.421	2.028	1.957	1.460	1.279	1.417	1.322
spine2	1.329	1.308	1.439	2.089	2.049	1.480	1.280	1.387	1.309

Table 6 Comparison using End Point Error (EPE) on the *Multi-Human Optical Flow* (MHOF) dataset. We show the average EPE and body part specific EPE, where part labels follow Figure 3. The first two rows are repeated from Tab 5.