# MONET: Multiview Semi-supervised Keypoint Detection via Epipolar Divergence

Yuan Yao
University of Minnesota
yaoxx340@umn.edu

Yasamin Jafarian
University of Minnesota
yasamin@umn.edu
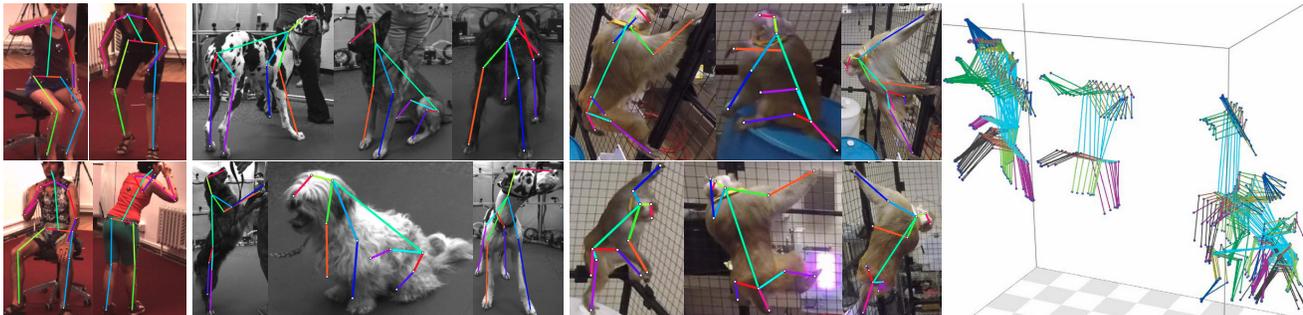
Hyun Soo Park
University of Minnesota
hspark@umn.edu

Figure 1: This paper presents MONET-an semi-supervised learning for keypoint detection, which is able to localize customized keypoints of diverse species, e.g., humans, dogs, and monkeys with very limited number of labeled data without a pre-trained model. The right most figure illustrates 3D reconstruction of monkey movement using our pose detection.

## Abstract

*This paper presents MONET—an end-to-end semi-supervised learning framework for a keypoint detector using multiview image streams. In particular, we consider general subjects such as non-human species where attaining a large scale annotated dataset is challenging. While multiview geometry can be used to self-supervise the unlabeled data, integrating the geometry into learning a keypoint detector is challenging due to representation mismatch. We address this mismatch by formulating a new differentiable representation of the epipolar constraint called epipolar divergence—a generalized distance from the epipolar lines to the corresponding keypoint distribution. Epipolar divergence characterizes when two view keypoint distributions produce zero reprojection error. We design a twin network that minimizes the epipolar divergence through stereo rectification that can significantly alleviate computational complexity and sampling aliasing in training. We demonstrate that our framework can localize customized keypoints of diverse species, e.g., humans, dogs, and monkeys.*

## 1. Introduction

Human pose detection has advanced significantly over the last few years [8, 43, 64, 69], driven in large part to new approaches based on deep learning. But these techniques require large amounts of labeled training data. For this reason, pose detection is almost always demonstrated on humans, for which large-scale datasets are available (e.g.,

MS COCO [37] and MPII [2]). What about pose detectors for other animals, such as monkeys, mice, and dogs? Such algorithms could have enormous scientific impact [41], but obtaining large-scale labeled training data would be a substantial challenge: each individual species may need its own dataset, some species have large intra-class variations, and domain experts may be needed to perform the labeling accurately. Moreover, while there is significant commercial interest in human pose recognition, there may be little incentive for companies and research labs to invest in collecting large-scale datasets for other species.

This paper addresses this annotation challenge by leveraging *multiview image streams*. Our insight is that the manual effort of annotation can be significantly reduced by using the redundant visual information embedded in the multiview imagery, allowing cross-view self-supervision: one image can provide a supervisionary signal to another image through epipolar geometry without 3D reconstruction. To this end, we design a novel end-to-end semi-supervised framework to utilize a large set of unlabeled multiview images using cross-view supervision.

The key challenge of integrating the epipolar geometry for building a strong keypoint (pose) detector lies in a representational mismatch: the geometric quantities such as points, lines, and planes are represented as a *vectors* [18] (Figure 2(a) left) while the *raster* representation via pixel response (heatmap [8, 43, 69]) has been shown strong performance on keypoint detection. For instance, applying the epipolar constraint [40]—a point $\mathbf{x} \in \mathbb{R}^2$ must lie in the
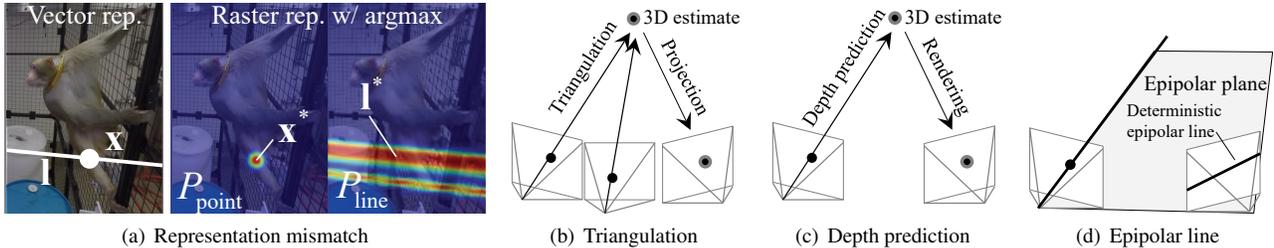
Figure 2: (a) Vector and raster representations describe the epipolar geometry. Note that the raster representation requires a non-differentiable argmax operation to compute $\mathbf{x}^*$ and $\mathbf{l}^*$. (b-d) Various multiview supervision approaches. (b) Keypoint prediction from at least two images can be triangulated and projected to supervise another image. This involves a non-differentiable argmax and RANSAC process [58]. (c) A 3D point [54], mesh [?], and voxel [70] can be predicted from a single view and projected to supervise another image. This requires an additional 3D prediction that fundamentally bounds the supervision accuracy. (d) Our approach precisely transfers a keypoint detection in one image to another image through the epipolar plane for the cross-supervision, and does not require 3D reconstruction.

corresponding epipolar line $\mathbf{l} \in \mathbb{P}^2$—can be expressed as:

$$(\widetilde{\mathbf{x}}^*)^\mathsf{T}\mathbf{l}^* = 0 \quad \text{s.t.} \ \ \mathbf{x}^* = \operatorname*{argmax}_{\mathbf{x}} P_\mathrm{p}(\mathbf{x}), \ \ \mathbf{l}^* = \operatorname*{argmax}_{\mathbf{l}} P_\mathrm{e}(\mathbf{l}),$$

where $\widetilde{\mathbf{x}}$ is the homogeneous representation of $\mathbf{x}$, and $P_\mathrm{p}$ and $P_\mathrm{e}$ are the distributions of keypoints and epipolar lines[1]. Note that the raster representation involves non-differentiable *argmax* operations, which are not trainable. This challenge leads to offline reconstruction [7,58,67], data driven depth prediction [31,53,54,65,74], or the usage of the soft-argmax operation [13], which shows inferior performance (see Figure 6).

In this paper, we formulate a new raster representation of the epipolar geometry that eliminates the argmax operations. We prove that the minimization of geometric error (i.e., $|\widetilde{\mathbf{x}}^\mathsf{T}\mathbf{l}|$) is equivalent to minimizing *epipolar divergence*— a generalized distance from the epipolar lines to the corresponding keypoint distribution. With this measure, we design a new end-to-end semi-supervised network called MONET (Multiview Optical Supervision Network). The network efficiently leverages the unlabeled multiview image streams with limited numbers of manual annotations ($< 1\%$). We integrate this raster formulation into the network by incorporating it with stereo rectification, which reduces the computational complexity and sampling artifacts while training the network.

The key features of MONET include that (1) it does not require offline triangulation that involves non-differentiable argmax and RANSAC operations [58] (Figure 2(b)); (2) it does not require 3D prediction [53,54,70] (Figure 2(c)), i.e., it deterministically transfers keypoint detections in one image to the other via epipolar geometry (Figure 2(d))[2]; (3) it is compatible with any keypoint detector design including CPM [69] and Hourglass [43] which localizes keypoints through a raster representation; and (4) it can apply to general multi-camera systems (e.g., different multi-camera rigs, number of cameras, and intrinsic parameters).

The main contributions of this paper include: (1) introducing a novel measure called the epipolar divergence, which measures the geometric consistency between two view keypoint distributions; (2) a network called MONET that efficiently minimizes the epipolar divergence via stereo rectification of keypoint distributions; (3) a technique for large-scale spatiotemporal data augmentation using 3D reconstruction of keypoint trajectories; (4) experimental results that demonstrate that MONET is flexible enough to detect keypoints in various subjects (humans, dogs, and monkeys) in different camera rigs and to outperform existing baselines in terms of localization accuracy and precision (re-projection error).

## 2. Related Work

The physical and social behaviors of non-human species such as rhesus macaque monkeys have been widely used as a window to study human activities in neuroscience and psychology. While measuring their subtle behaviors in the form of 3D anatomic landmarks is key, implementing marker-based 3D tracking systems is challenging due to the animal's sensitivity to reflective markers and occlusion by fur, which limits its applications to restricted body motions (e.g., body tied to a chair) [1]. Vision-based marker-less motion capture is a viable solution to measure their free ranging behaviors [16,42,55].

In general, the number of 3D pose configurations of a deformable articulated body is exponential with respect to the number of joints. The 2D projections of the 3D body introduces substantial variability in illumination, appearance, and occlusion, which makes pose estimation challenging. But the space of possible pose configurations has structure that can be captured by efficient spatial representations such as pictorial structures [3,4,14,25,50,51,71], hierarchical and non-tree models [12,32,35,57,60,62,68] and convolutional architectures [9,10,33,39,44,48,49,63,64], and inference on these structures can be performed efficiently using clever algorithms, e.g., dynamic programming, convex relaxation, and approximate algorithms. Albeit efficient and accurate on canonical images, they exhibit inferior performance on images in the long-tail distribution, e.g., a pigeon pose of

---

[1]See Section 3.1, respectively, as shown in Figure 2(a) for computation of $P_\mathrm{e}$.

[2]This is analogous to the fundamental matrix computation without 3D estimation [18,40].

yoga. Fully supervised learning frameworks using millions of perceptrons in convolutional neural networks (CNNs) [8, 43, 64, 69] can address this long-tail distribution issue by leveraging a sheer amount of training data annotated by crowd workers [2, 37, 56]. However, due to the number of parameters in a CNN, the trained model can be highly biased when the number of data samples is not sufficient ($<$1M).

Semi-supervised and weakly-supervised learning frameworks train CNN models with limited number of training data [5,23,36,38,45,46,59,61,66,75]. For instance, temporal consistency derived by tracking during training can provide a supervisionary signal for body joint detection [36]. Geometric (such as 3DPS model [5]) and spatial [59] relationship are another way to supervise body keypoint estimation. Active learning that finds the most informative images to be annotated can alleviate the amount of labeling effort [38], and geometric [46] and temporal [23] in 2D [30] and 3D [27, 72] consistency can also be used to augment annotation data.

These approaches embed underlying spatial structures such as 3D skeletons and meshes that regularize the network weights. For instance, motion capture data can be used to jointly learn 2D and 3D keypoints [75], and scanned human body models are used to validate 2D pose estimation via reprojection [17,29,31,73,76], e.g., by using a DoubleFusion system that can simultaneously reconstruct the inner body shape and pose. The outer surface geometry and motion in real-time by using a single depth camera [73] and recovery human meshes that can reconstruct a full 3D mesh of human bodies from a single RGB camera by having 2D ground truth annotations [31]. Graphical models can also be applied for animal shape reconstruction by learning a 3D model based on a small set of 3D scans of toy figurines in arbitrary poses and refining the model and initial registration of scans together, and then generalizing it by fitting the model to real images of animal species out of the training set [76]. Notably, a multi-camera system can be used to cross-view supervise multiview synchronized images using iterative process of 3D reconstruction and network training [54, 58].

Unlike existing methods, MONET does not rely on a spatial model. To our knowledge, this is the first paper that jointly reconstructs and trains a keypoint detector without iterative processes using epipolar geometry. We integrate reconstruction and learning through a new measure of keypoint distributions called epipolar divergence, which can apply to general subjects including non-human species where minimal manual annotations are available.

## 3. MONET

We present a semi-supervised learning framework for training a keypoint detector by leveraging multiview image streams for which $|\mathcal{D}_U| \gg |\mathcal{D}_L|$, where $\mathcal{D}_L$ and $\mathcal{D}_U$ are labeled and unlabeled data, respectively. We learn a network model that takes an input image $\mathcal{I}$ and outputs a keypoint distribution, i.e., $\phi(\mathcal{I}; \mathbf{w}) \in [0,1]^{W \times H \times C}$ where $\mathcal{I}$ is an input image, $\mathbf{w}$ is the learned network weights, and $W$, $H$, and $C$ are the width, height, and the number of keypoints.

To enable end-to-end cross-view supervision without 3D reconstruction, we formulate a novel raster representation of epipolar geometry in Section 3.1, and show how to implement it in practice using stereo rectification in Section 3.2. The full learning framework is described in Section 3.3 by incorporating a bootstrapping prior.

### 3.1. Epipolar Divergence

A point in the $i^{\text{th}}$ image $\mathbf{x}_i \in \mathbb{R}^2$ is *transferred* to form a corresponding epipolar line in the $j^{\text{th}}$ image via the fundamental matrix $\mathbf{F}$ between two relative camera poses, which measures geometric consistency, i.e., the corresponding point $\mathbf{x}_j$ must lie in the epipolar line [18]:

$$D(\mathbf{x}_i, \mathbf{x}_j) = \left| \widetilde{\mathbf{x}}_j^\mathsf{T} \left( \mathbf{F} \widetilde{\mathbf{x}}_i \right) \right| \propto \inf_{\mathbf{x} \in \mathbf{F}\widetilde{\mathbf{x}}_i} \|\mathbf{x} - \mathbf{x}_j\|. \quad (1)$$

The infimum operation measures the distance between the closest point in the epipolar line ($\mathbf{F}\widetilde{\mathbf{x}}_i$) and $\mathbf{x}_j$ in the $j^{\text{th}}$ image.

We generalize the epipolar line transfer to define the distance between keypoint distributions. Let $P_i : \mathbb{R}^2 \to [0,1]$ be the keypoint distribution given the $i^{\text{th}}$ image computed by a keypoint detector, i.e., $P_i(\mathbf{x}) = \phi(\mathcal{I}_i; \mathbf{w})|_{\mathbf{x}}$, and $P_{j \to i} : \mathbb{R}^2 \to [0,1]$ be the keypoint distribution in the $i^{\text{th}}$ image *transferred* from the $j^{\text{th}}$ image as shown in Figure 3(a). Note that we abuse notation by omitting the keypoint index, as each keypoint is considered independently.

Consider a max-pooling operation along a line, $g$:

$$g(\mathbf{l}; P) = \sup_{\mathbf{x} \in \mathbf{l}} P(\mathbf{x}), \quad (2)$$

where $P : \mathbb{R}^2 \to [0,1]$ is a distribution and $\mathbf{l} \in \mathbb{P}^2$ is a 2D line parameter. $g$ takes the maximum value along the line in $P$. Given the keypoint distribution in the $j^{\text{th}}$ image $P_j$, the transferred keypoint distribution can be obtained:

$$P_{j \to i}(\mathbf{x}_i) = g(\mathbf{F}\widetilde{\mathbf{x}}_i; P_j). \quad (3)$$

The supremum operation is equivalent to the infimum operation in Equation (1), where it finds the most likely (closest) correspondences along the epipolar line. The first two images in Figure 3(a) illustrate the keypoint distribution transfer via Equation (3). The keypoint distribution in the $i^{\text{th}}$ image is deterministically transformed to the rasterized epipolar line distribution in the $j^{\text{th}}$ image, i.e., no explicit 3D reconstruction (triangulation or depth prediction) is needed. In fact, the transferred distribution is a posterior distribution of a 3D keypoint given a uniform depth prior.

$P_i$ and $P_{j \to i}$ cannot be directly matched because $P_i$ is a point distribution while $P_{j \to i}$ is a line distribution. A key observation is that points that lie on the same epipolar line in $P_{j \to i}$ have the same probability, i.e., $P_{i \to j}(\mathbf{x}_j) = P_{i \to j}(\mathbf{y}_j)$ if $\mathbf{F}^\mathsf{T} \widetilde{\mathbf{x}}_j \propto \mathbf{F}^\mathsf{T} \widetilde{\mathbf{y}}_j$ as shown in the second image of Figure 3(a). This indicates that the transferred distribution can be parametrized by the slope of an epipolar line, $\theta \in \mathbb{S}$, i.e.,

$$Q_{j \to i}(\theta) = g(\mathbf{l}_i(\theta); P_{j \to i}), \quad (4)$$

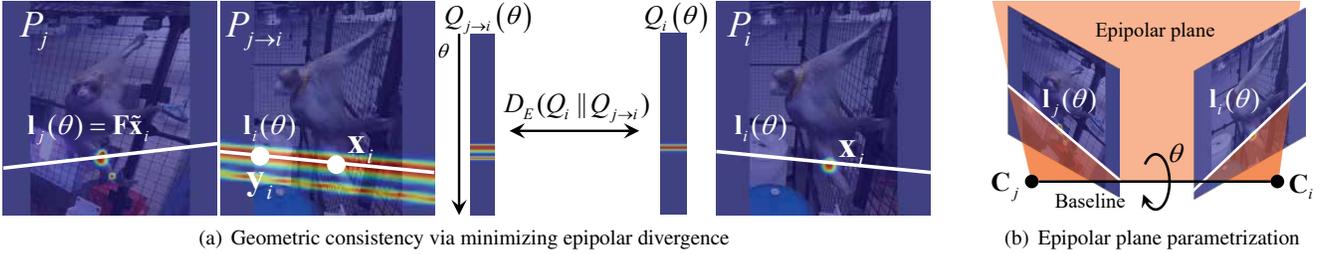(a) Geometric consistency via minimizing epipolar divergence

(b) Epipolar plane parametrization

Figure 3: (a) The keypoint distribution of the knee joint for the $j^{\text{th}}$ image, $P_j$, is transferred to the $i^{\text{th}}$ image to form the epipolar line distribution $P_{j \to i}(\mathbf{x}_i)$. Note that the points that lie in the same epipolar line have the equal transferred distribution, $P_{j \to i}(\mathbf{x}_i) = P_{j \to i}(\mathbf{y}_i)$, and therefore (b) the distribution can be reparametrized by the 1D rotation $\theta \in \mathbb{S}$ about the baseline where $\mathbf{C}_i$ and $\mathbf{C}_j$ are the camera optical centers. We match two distributions: the distribution transferred from the $i^{\text{th}}$ image $Q_{j \to i}(\theta)$ and the distribution of keypoint in the $j^{\text{th}}$ image $Q_i(\theta)$. The minimization of the epipolar divergence $D_E(Q_i || Q_{j \to i})$ is provably equivalent to reprojection error minimization.

where $\mathbf{l}_i(\theta)$ is the line passing through the epipole parametrized by $\theta$ in the $i^{\text{th}}$ image, and $Q_{j \to i} : \mathbb{S} \to [0, 1]$ is a flattened 1D distribution across the line. Similarly, the flattened keypoint distribution of $P_i$ can be defined as $Q_i(\theta) = g(\mathbf{l}_i(\theta); P_i)$.

**Theorem 1.** *Two keypoint distributions $P_i$ and $P_j$ are geometrically consistent, i.e., zero reprojection error, if $Q_i(\theta) = Q_{j \to i}(\theta)$.*

See the proof in Appendix. Theorem 1 states the necessary condition of zero reprojection: the detected keypoints across views must lie in the same epipolar plane in 3D. Figure 11 illustrates the epipolar plane that is constructed by the baseline and the 3D ray (inverse projection) of the detected keypoint. Matching $Q_i$ and $Q_{j \to i}$ is equivalent to matching the probabilities of epipolar 3D planes, which can be parametrized by their surface normal ($\theta$).

To match their distributions, we define an *epipolar divergence* that measures the difference between two keypoint distributions using relative entropy inspired by Kullback–Leibler (KL) divergence [34]:

$$D_{\text{E}}(Q_i || Q_{j \to i}) = \int_{\mathbb{S}} Q_i(\theta) \log \frac{Q_i(\theta)}{Q_{j \to i}(\theta)} d\theta. \quad (5)$$

This epipolar divergence measure how two keypoint distributions are geometrically consistent.

## 3.2. Cross-view Supervision via Rectification

In practice, embedding Equation (5) into an end-to-end neural network is non-trivial because (a) a new max-pooling operation over oblique epipolar lines in Equation (3) needs to be defined; (b) the sampling interval for max-pooling along the line is arbitrary, i.e., uniform sampling does not encode geometric meaning such as depth; and (c) the sampling interval across $\theta$ is also arbitrary. These factors increase computational complexity and sampling artifacts in the process of training.

We introduce a new operation inspired by stereo rectification, which warps a keypoint distribution such that the epipolar lines become parallel (horizontal) as shown the

bottom right image in Figure 4. This rectification allows converting the max-pooling operation over an oblique epipolar line into regular row-wise max-pooling, i.e., epipolar line can be parametrized by its height $\mathbf{l}(v)$. Equation (2) can be re-written with the rectified keypoint distribution:

$$\overline{g}(v; \overline{P}) = g\left(\mathbf{l}(v); \overline{P}\right) = \max_u \overline{P}\left(\begin{bmatrix} u \\ v \end{bmatrix}\right) \quad (6)$$

where $(u, v)$ is the $x, y$-coordinate of a point in the rectified keypoint distribution $\overline{P}$ warped from $P$, i.e., $\overline{P}(\mathbf{x}) = P(\mathbf{H}_r^{-1}\mathbf{x})$ where $\mathbf{H}_r$ is the homography of stereo-rectification. $\overline{P}$ is computed by inverse homography warping with bilinear interpolation [19, 24]. This rectification simplifies the flattening operation in Equation (4):

$$\begin{aligned} \overline{Q}_{j \to i}(v) &= \overline{g}(v; \overline{P}_{j \to i}) = \overline{g}\left(av + b; \overline{P}_j\right), \\ \overline{Q}_i(v) &= \overline{g}(v; \overline{P}_i), \end{aligned} \quad (7)$$

where $a$ and $b$ are re-scaling factors between the $i^{\text{th}}$ and $j^{\text{th}}$ cameras, accounting different camera intrinsic and cropping parameters. See Appendix for more details.

The key innovation of Equation (7) is that $\overline{Q}_{j \to i}(v)$ is no longer parametrized by $\theta$ where an additional sampling over $\theta$ is not necessary. It directly accesses $\overline{P}_j$ to max-pool over each row, which significantly alleviates computational complexity and sampling artifacts. Moreover, sampling over the $x$-coordinate is geometrically meaningful, i.e., uniform sampling is equivalent to disparity, or inverse depth.

With rectification, we model the loss for multiview cross-view supervision:

$$\mathcal{L}_E = \sum_{c=1}^{C} \sum_{i=1}^{S} \sum_{j \in \mathcal{V}_i} \sum_{v=1}^{H} \overline{Q}_i^c(v) \log \frac{\overline{Q}_i^c(v)}{\overline{Q}_{j \to i}^c(v)} \quad (8)$$

where $H$ is the height of the distribution, $P$ is the number of keypoints, $S$ is the number of cameras, and $\mathcal{V}_i$ is the set of paired camera indices of the $i^{\text{th}}$ camera. We use the superscript in $\overline{Q}_i^c$ to indicate the keypoint index. Figure 4 illustrates our twin network that minimizes the epipolar divergence by applying stereo rectification, epipolar transfer,
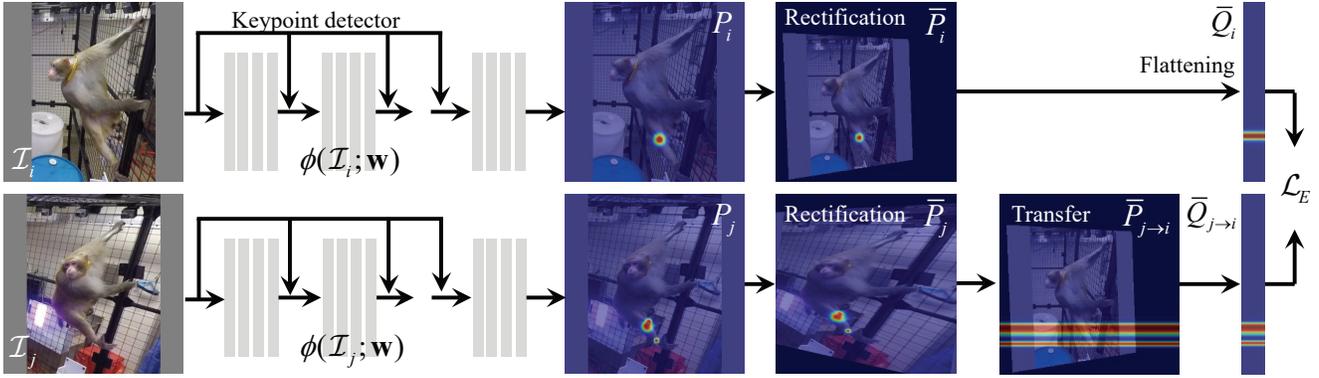
Figure 4: We design a twin network to minimize the epipolar divergence between $\overline{Q}_i$ and $\overline{Q}_{j\to i}$. Stereo rectification is used to simplify the max-pooling operation along the epipolar line, and reduce computational complexity and sampling aliasing.
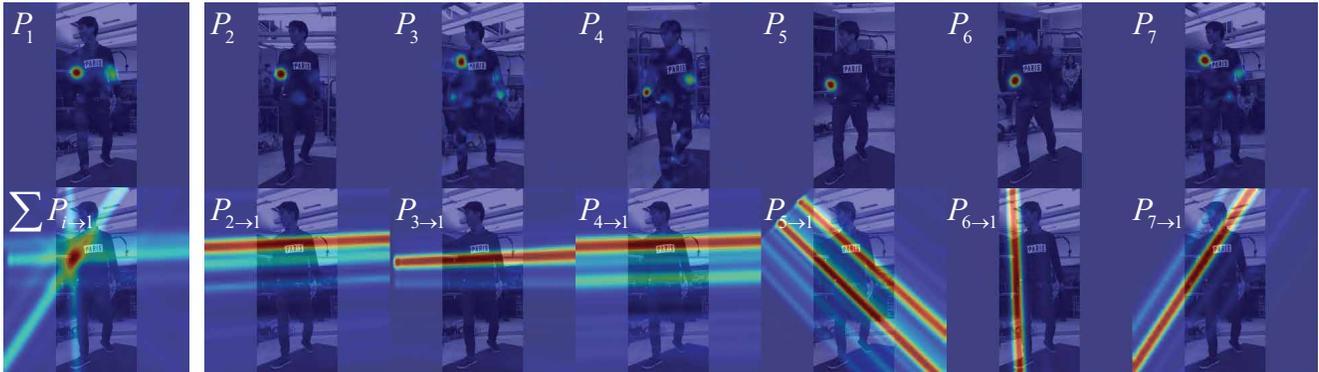


Figure 5: Epipolar cross-view supervision on right elbow on view 1. Top right row shows elbow detections across views, i.e., $P_2, \cdots, P_7$. The transferred distribution to view 1 is shown on the bottom right row, i.e., $P_{2\to1}, \cdots, P_{7\to1}$. These transferred probabilities are used to supervise view 1 where the bottom left image is the summation of cross-view supervisions.

and flattening operations, which can perform cross-view supervision from unlabeled data.

Since the epipolar divergence flattens the keypoint distribution, cross-supervision from one image can constrain in one direction. In practice, we find a set of images given the $i^{\text{th}}$ image such that the expected epipolar lines are not parallel. When camera centers lie on a co-planar surface, a 3D point on the surface produces all same epipolar lines, which is a degenerate case[3]. Figure 5 illustrates cross-view supervision on a right elbow on view 1. Elbow detections from view 2 to 7 (top right row) are transferred to view 1 (bottom right row). These transferred probabilities are used to supervise view 1 where the bottom left image is the summation of cross-view supervisions.

### 3.3. Multiview Semi-supervised Learning

We integrate the raster formulation of the epipolar geometry in Section 3.2 into a semi-supervised learning framework. The keypoint detector is trained by minimizing the following loss:

$$\underset{\mathbf{w}}{\text{minimize}} \ \mathcal{L}_L + \lambda_e \mathcal{L}_E + \lambda_p \mathcal{L}_B, \qquad (9)$$

---

[3]This degenerate case does not apply for 3D point triangulation where the correspondence is known.

where $\mathcal{L}_L$, $\mathcal{L}_E$, and $\mathcal{L}_B$ are the losses for labeled supervision, multiview cross-view supervision, and bootstrapping prior, and $\lambda_e$ and $\lambda_p$ are the weights that control their importance.

Given a set of labeled data ($<1\%$), we compute the labeled loss as follows:

$$\mathcal{L}_L = \sum_{i \in \mathcal{D}_L} \|\phi\left(\mathcal{I}_i; \mathbf{w}\right) - \mathbf{z}_i\|^2 \qquad (10)$$

where $\mathbf{z} \in [0,1]^{W \times H \times C}$ is the labeled likelihood of keypoints approximated by convolving the keypoint location with a Gaussian kernel.

To improve performance, we incorporate with offline spatiotemporal label augmentation by reconstructing 3D keypoint trajectories using the multiview labeled data inspired by the multiview bootstrapping [58]. Given synchronized labeled images, we triangulate each 3D keypoint $\mathbf{X}$ using the camera projection matrices and the 2D labeled keypoints. The 3D reconstructed keypoint is projected onto the rest synchronized unlabeled images, which automatically produces their labels. 3D tracking [27,72] further increases the labeled data. For each keypoint $\mathbf{X}_t$ at the $t$ time instant, we project the point onto the visible set of cameras. The projected point is tracked in 2D using optical flow and triangulated with RANSAC [15] to form $\mathbf{X}_{t+1}$. We compute the visibility

of the point to reduce tracking drift using motion and appearance cues: (1) optical flow from its consecutive image is compared to the projected 3D motion vector to measure motion consistency; and (2) visual appearance is matched by learning a linear correlation filter [6] on PCA HOG [11], which can reliably track longer than 100 frames forward and backward. We use this spatiotemporal data augmentation to define the bootstrapping loss:

$$\mathcal{L}_B = \sum_{i \in \mathcal{D}_U} \|\phi(\mathcal{I}_i; \mathbf{w}) - \widehat{\mathbf{z}}_i\|^2. \qquad (11)$$

where $\widehat{\mathbf{z}} \in [0,1]^{W \times H \times C}$ is the augmented labeled likelihood using bootstrapping approximated by convolving the keypoint location with a Gaussian kernel.

## 4. Result

We build a keypoint detector for each species without a pre-trained model, using the CPM network (5 stages). The code can be found in https://github.com/MONET2018/MONET. To highlight the model flexibility, we include implementations with two state-of-the-art pose detectors (CPM [8] and Hourglass [43]). $\lambda_e = 5$ and $\lambda_p = 1$ are used. Our detection network takes an input image ($368 \times 368$), and outputs a distribution ($46 \times 46 \times C$). In training, we use batch size 30, learning rate $10^{-4}$, and learning decay rate 0.9 with 500 steps. We use the ADAM optimizer of TensorFlow with single nVidia GTX 1080.

**Datasets** We validate our MONET framework on multiple sequences of diverse subjects including humans, dogs, and monkeys. (1) **Monkey subject** 35 GoPro HD cameras running at 60 fps are installed in a large cage ($9' \times 12' \times 9'$) that allows the free-ranging behaviors of monkeys. There are diverse monkey activities include grooming, hanging, and walking. The camera produces $1280 \times 960$ images. 12 keypoints of monkey's pose in 85 images out of 63,000 images are manually annotated. (2) **Dog subjects** Multi-camera system composed of 69 synchronized HD cameras ($1024 \times 1280$ at 30 fps) are used to capture the behaviors of multiple breeds of dogs including Dalmatian and Golden Retrievers. Less than 1% of data are manually labeled. (3) **Human subject I** A multiview behavioral imaging system composed of 69 synchronized HD cameras capture human activities at 30 fps with $1024 \times 1280$ resolution. 30 images out of 20,700 images are manually annotated. This dataset includes a diverse human activities such as dancing, jumping, and sitting. We use a pre-trained CPM model [8] to generate the ground truth data. (4) **Human subject II** We test our approach on two publicly available datasets for human subjects: Panoptic Studio dataset [26] and Human3.6M [22]. For the Panoptic Studio dataset, we use 31 HD videos ($1920 \times 1080$ at 30 Hz). The scenes includes diverse subjects with social interactions that introduce severe social occlusion. The Human3.6M dataset is captured by 4 HD cameras that includes variety of single actor activities, e.g., sitting, running, and eating/drinking.
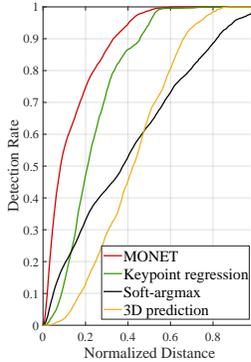


Figure 6: PCK for hypothesis validation

**Hypothesis Validation** We hypothesize that our raster formulation is superior to existing multiview cross-view supervision approaches used for semi-supervised learning because it is an end-to-end system without requiring 3D prediction. We empirically validate our hypothesis by comparing to three approaches on multiview monkey data from 35 views (300 labeled and 600 unlabeled time instances). No pretrained model is used for the evaluation. (1) *Keypoint regression*: a vector representation of keypoint locations is directly regressed from an image. We use DeepPose [64] to detect keypoints and use the fundamental matrix to measure the distance (loss) between the epipolar line and the detected points, $|\widetilde{\mathbf{x}}_2^\top \mathbf{F} \widetilde{\mathbf{x}}_1|$, for the unlabeled data. (2) *Soft-argmax*: a vector representation can be approximated by the raster keypoint distribution using a soft-argmax operation: $\mathbf{x}_{\text{softmax}} = \sum_{\mathbf{x}} P(\mathbf{x})\mathbf{x} / \sum_{\mathbf{x}} P(\mathbf{x})$, which is reasonable when the predicted probability is nearly unimodal. This is differentiable, and therefore end-to-end training is possible. However, its approximation holds when the predicted distribution is unimodal. We use CPM [69] to build a semi-supervised network with epipolar distance as a loss. (3) *3D prediction*: each 3D coordinate is predicted from a single view image where the projection of the 3D prediction is used as cross-view supervison [**?**, 54, 70]. We augment 3D prediction layers on CPM to regress the depth of keypoints [47]. The reprojection error is used for the loss. Figure 6 illustrates the probability of correct keypoint (PCK) curve, showing that our approach using raster epipolar geometry significantly outperforms other approaches.

**Baselines** We compare our approach with 5 different baseline algorithms. For all algorithms, we evaluate the performance on the unlabeled data. (1) *Supervised learning*: we use the manually annotated images to train the network in a fully supervised manner. Due to the limited number of labeled images ($<100$), the existing distillation methods [21, 52] perform similarly. (2) *Spatial augmentation*: the 3D keypoints are triangulated and projected onto the synchronized unlabeled images. This models visual appearance and spatial configuration from multiple perspectives, which can greatly improve the generalization power of keypoint detection. (3) *Spatiotemporal augmentation*: we track the 3D keypoints over time using multiview optical flow [27, 72]. This augmentation can model different geometric configurations of 3D keypoints. (4) *Bootstrapping I*: Given the spatiotemporal data augmentation, we apply the multiview bootstrapping approach [58] to obtain pseudo-labels computed by RANSAC-based 3D triangulation for the unlabeled data. (5) *Bootstrapping II*: the Bootstrapping I model is refined by re-triangulation and re-training. This can reduce the reprojection errors. We evaluate our approach based on

(a) Human subject PCK    (b) Monkey subject PCK    (c) Dog subject PCK    (d) Panoptic PCK    (e) Reprojection error
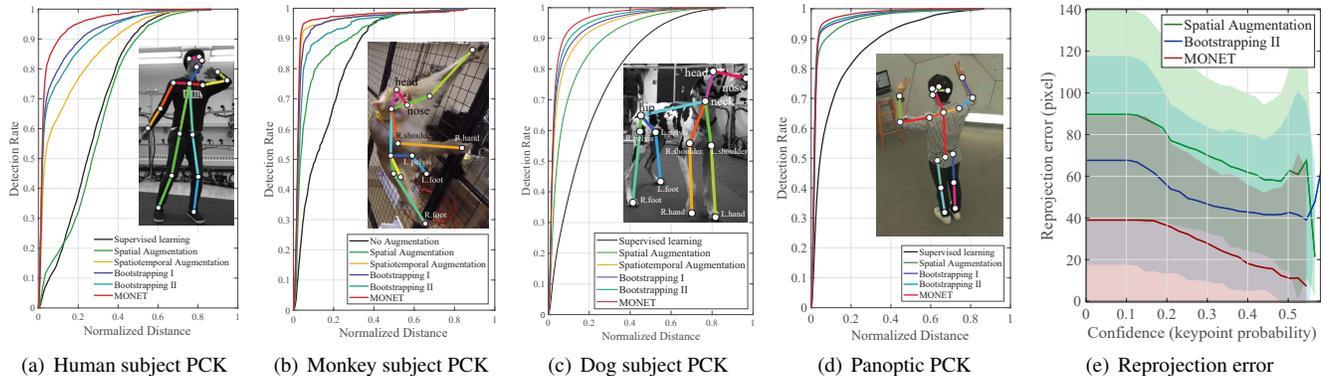
Figure 7: PCK curves for (a) humans, (b) monkeys, (c) dogs and (d) the CMU Panoptic dataset [28]. MONET (red) outperforms 5 baseline algorithms. (e) MONET is designed to minimize the reprojection error, and we achieve far stronger performance as the confidence increases.
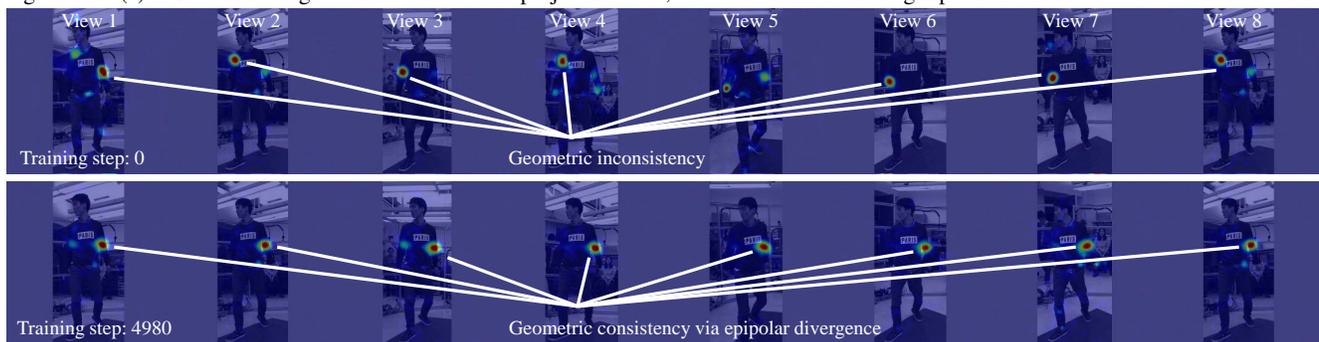


Figure 8: Erroneous elbow detections from multiview images converge to the geometrically consistent location through training.

accuracy and precision: accuracy measures distance from the ground truth keypoint and precision measures the coherence of keypoint detections across views. (6) *Rhodin et al. [54]*: The unlabeled multi-view image pairs are used to generate 3D point cloud of body first during unsupervised training, and then the model is trained with images with 3D ground truth to learn to transfer point cloud to joint positions.

**Accuracy** We use PCK curves to measure the accuracy. The distance between the ground truth keypoint and the detected keypoint is normalized by the size of the width of the detection window (46). Figure 7 shows PCK performance on human, monkey, and dog subjects where no pre-trained model is used. Our MONET (red) model exhibits accurate detection for all keypoints, and outperforms 5 baselines. For the monkey data, higher frame-rate image streams (60 fps) greatly boost the performance of multiview tracking due to smaller displacements, resulting in accurate keypoint detection by spatiotemporal augmentation. We also conducted an experiment on the CMU Panoptic dataset [28] to validate the generalization power of our approach. This dataset differs from ours in terms of camera parameters, placements, and scene (e.g., pose, illumination, background, and subject). MONET outperforms on both accuracy (PCK) and precision (reprojection error) as shown in Figure 7(d).

**Precision** We use reprojection error to evaluate the precision of detection. Given a set of keypoint detections in a synchronized frame and 3D camera poses, we triangulate

| | Human | Monkey | Dog | Panoptic |
|---|---|---|---|---|
| Supervised learning | 77.8±73.3 | 31.1±872 | 88.9±69.9 | 53.2±271.4 |
| Spatial aug. | 69.0±66.2 | 12.9±26.6 | 37.5±47.1 | 22.2±40.4 |
| Spatiotemporal aug. | 50.3±65.4 | 8.10±17.8 | 24.0±36.2 | N/A |
| Bootstrapping I [58] | 28.5±44.7 | 8.68±18.9 | 18.9±31.0 | 15.6±31.7 |
| Bootstrapping II [58] | 35.4±62.4 | 9.97±22.1 | 17.1±29.3 | 13.7±24.6 |
| MONET | **15.0±24.1** | **5.45±11.4** | **10.3±18.7** | **12.8±18.0** |

Table 1: Reprojection error (Mean±Std).

| Labeled / Unlabeled | Hips | R.Leg | R.Arm | Head | L.Hand | L.Foot | R.UpLeg | Neck | Total |
|---|---|---|---|---|---|---|---|---|---|
| S1 / S5,6,7,8 | 13.0 | 3.1 | 3.4 | 1.0 | 6.6 | 6.2 | 10.9 | 1.6 | 5.5 |
| S1,5 / S6,7,8 | 12.7 | 2.2 | 2.9 | 1.0 | 5.2 | **3.3** | 10.9 | 1.6 | 5.2 |
| S1,5,6 / S7,8 | **7.1** | **2.0** | **2.7** | **0.9** | **5.0** | 4.7 | **5.6** | **1.5** | **4.3** |

Table 2: Mean pixel error vs. labeled data size on Human3.6M dataset

the 3D point without RANSAC. The 3D point is projected back to each camera to compute the reprojection error, which measures geometric consistency across all views. MONET is designed to minimize the reprojection error, and it outperforms baselines significantly in Figure 7(e). Our MONET performs better at higher keypoint distribution, which is key for 3D reconstruction because it indicates which points to triangulate. Figure 8 shows how erroneous detections of the left elbow from multiview images converge to geometrically consistent elbow locations as the training progresses. The performance for each subject is summarized in Table 1.

**Robustness** We evaluate the robustness of our approach by varying the amount of labeled data on Human3.6M dataset (four cameras), which provides motion capture ground truth
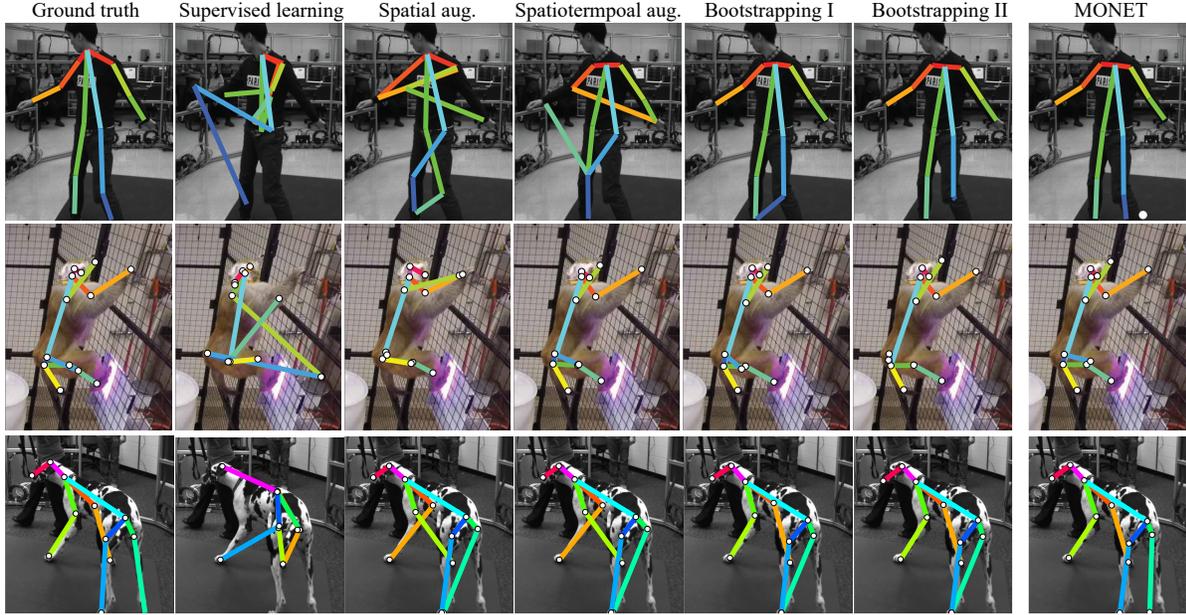
Figure 9: We qualitatively compare our MONET with 5 baseline algorithms on humans, monkeys, and dogs.



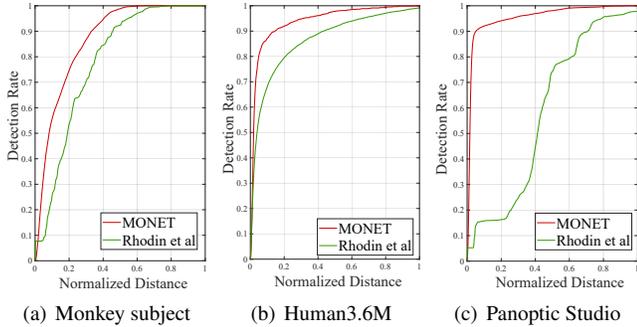(a) Monkey subject    (b) Human3.6M    (c) Panoptic Studio

Figure 10: Comparison with Rhodin et al. [54] that predict 3D points for cross-view supervision on monkey, Human3.6M, and Panoptic Studio datasets.

data. Table 2 summarizes the mean pixel error as varying the labeled and unlabeled subjects. As expected, as the labeled data increases, the error decreases while the minimally labeled S1 (subject 1) still produces less than 15 max pixel error. We also compare to a 3D prediction approach [54], which showed strong performance on Human3.6M dataset. Similar to their experimental setup, we use S1, S5, and S6 as the labeled data, and S7 and S8 as the unlabeled data for training. In addition to Human3.6M dataset, we also conduct the comparison on the Monkey and CMU Panoptic dataset [28]. Figure 10 illustrates the PCK measure on the unlabeled data. Our approach outperforms the baseline on all the datasets. The advantage of our approach is especially reflected on the CMU Panoptic dataset. Full body is not often visible due to the narrow FOV cameras, which makes the explicit 3D reconstruction in [54] of body less efficient.

**Qualitative Comparison** A qualitative comparison can be found in Figure 9. MONET can precisely localize keypoints

by leveraging multiview images jointly. This becomes more evident when disambiguating symmetric keypoints, e.g., left and right hands, as epipolar divergence penalizes geometric inconsistency (reprojection error). It also shows stronger performance under occlusion (the bottom figure) as the occluded keypoints can be visible to other views that can enforce to the correct location.

## 5. Discussion

We present a new semi-supervised framework, MONET, to train keypoint detection networks by leveraging multi-view image streams. The key innovation is a measure of geometric consistency between keypoint distributions called epipolar divergence. Similar to epipolar distance between corresponding points, it allows us to directly compute reprojection error while training a network. We introduce a stereo rectification of the keypoint distribution that simplifies the computational complexity and imposes geometric meaning on constructing 1D distributions. A twin network is used to embed computation of epipolar divergence. We also use multiview image streams to augment the data in space and time, which bootstraps unlabeled data. We demonstrate that our framework outperforms existing approaches, e.g., multiview bootstrapping, in terms of accuracy (PCK) and precision (reprojection error), and apply it to non-human species such as dogs and monkeys. We anticipate that this framework will provide a fundamental basis for enabling *flexible* marker-less motion capture that requires exploiting a large (potentially unbounded) number of unlabeled data.

## 6. Acknowledgments

# References

[1] David Anderson. The nonhuman primate as a model for biomedical research. *Sourcebook of Models for Biomedical Research*, 2008. 2

[2] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *CVPR*, 2014. 1, 3

[3] Mykhaylo Andriluka, Stefan Roth, and Bernt Schiele. Pictorial structures revisited: People detection and articulated pose estimation. In *CVPR*, 2009. 2

[4] Mykhaylo Andriluka, Stefan Roth, and Bernt Schiele. Monocular 3d pose estimation and tracking by detection. In *CVPR*, 2010. 2

[5] Vasileios Belagiannis, Sikandar Amin, Mykhaylo Andriluka, Bernt Schiele, Nassir Navab, and Slobodan Ilic. 3d pictorial structures revisited: Multiple human pose estimation. *TPAMI*, 2016. 3

[6] Vishnu Naresh Boddeti and B.V.K Vijaya Kumar. A framework for binding and retrieving class-specific information to and from image patterns using correlation filters. *TPAMI*, 2013. 6

[7] Arunkumar Byravan and Dieter Fox. SE3-nets: Learning rigid body motion using deep neural networks. In *ICRA*, 2016. 2

[8] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017. 1, 3, 6, 12

[9] Joao Carreira, Pulkit Agrawal, Katerina Fragkiadaki, and Jitendra Malik. Human pose estimation with iterative error feedback. In *CVPR*, 2016. 2

[10] Xianjie Chen and Alan Yuille. Articulated pose estimation by a graphical model with image dependent pairwise relations. In *NIPS*, 2014. 2

[11] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *CVPR*, 2005. 6

[12] Matthias Dantone, Juergen Gall, Christian Leistner, and Luc Van Gool. Human pose estimation using body parts dependent joint regressors. In *CVPR*, 2013. 2

[13] Xuanyi Dong, Shoou-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 *CVPR*, 2018. 2

[14] Pedro Felzenszwalb and Daniel Huttenlocher. Pictorial structures for object recognition. *IJCV*, 2005. 2

[15] Martin Fischler and Robert Bolles. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *ACM Comm.*, 1981. 5

[16] Justin Foster, Paul Nuyujukian, Oren Freifeld, Hua Gao, Ross Walker, Stephen I Ryu, Teresa H Meng, Boris Murmann, Michael J Black, and Krishna V Shenoy. A freely-moving monkey treadmill model. *Journal of Neural Engineering*, 2014. 2

[17] Rıza Alp Guler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. In *CVPR*, 2018. 3

[18] Richard Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, second edition, 2004. 1, 2, 3, 12

[19] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *CVPR*, 2017. 4

[20] João Henriques, Rui Caseiro, Pedro Martins, and Jorge Batista. High-speed tracking with kernelized correlation filters. *TPAMI*, 2015. 12

[21] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In *arXiv:1503.02531*, 2015. 6

[22] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *TPAMI*, 2014. 6

[23] Umar Iqbal, Anton Milan, and Juergen Gall. Posetrack: Joint multi-person pose estimation and tracking. In *CVPR*, 2017. 3

[24] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer networks. In *NIPS*, 2015. 4

[25] Sam Johnson and Mark Everingham. Clustered pose and nonlinear appearance models for human pose estimation. In *BMVC*, 2010. 2

[26] Hanbyul Joo, Hao Liu, Lei Tan, Lin Gui, Bart Nabbe, Iain Matthews, Takeo Kanade, Shohei Nobuhara, and Yaser Sheikh. Panoptic studio: A massively multiview system for social motion capture. In *ICCV*, 2015. 6

[27] Hanbyul Joo, Hyun Soo Park, and Yaser Sheikh. Map visibility estimation for large-scale dynamic 3d reconstruction. In *CVPR*, 2014. 3, 5, 6

[28] Hanbyul Joo, Tomas Simon, Xulong Li, Hao Liu, Lei Tan, Lin Gui, Sean Banerjee, Timothy Godisart, Bart Nabbe, Iain Matthews, Takeo Kanade, Shohei Nobuhara, and Yaser Sheikh In. Panoptic studio: A massively multiview system for social interaction capture. *TPAMI*, 2017. 7, 8

[29] Hanbyul Joo, Tomas Simon, and Yaser Sheikh. Total capture: A 3d deformation model for tracking faces, hands, and bodies. In *CVPR*, 2018. 3

[30] Zdenek Kalal, Krystian Mikolajczyk, and Jiri Matas. Tracking-learning-detection. *TPAMI*, 2012. 3

[31] Angjoo Kanazawa, Michael Black, David Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018. 2, 3

[32] Leonid Karlinsky and Shimon Ullman. Using linking features in learning non-parametric part models. In *ECCV*, 2012. 2

[33] Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 2

[34] Solomon Kullback and Richard Leibler. On information and sufficiency. *Annals of Mathematical Statistics*, 1951. 4

[35] Xiangyang Lan and Daniel Huttenlocher. Beyond trees: Common-factor models for 2d human pose recovery. In *ICCV*, 2005. 2

[36] Mude Lin, Liang Lin, and Xiaodan Liang. Recurrent 3d pose sequence machines. In *CVPR*, 2017. 3

[37] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollàr, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 1, 3

[38] Buyu Liu and Vittorio Ferrari. Active learning for human pose estimation. In *CVPR*, 2017. 3

[39] Jonathan Long, Evan Shelhamer, and Trevor Darrel. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 2

[40] Hugh Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections. *Nature*, 1981. 1, 2

[41] Alexander Mathis, Pranav Mamidanna, Kevin Cury, Taiga Abe, Venkatesh Murthy, Mackenzie Mathis, and Matthias Bethge. Deeplabcut: markerless pose estimation of user-defined body parts with deep learning. *Nature Neuroscience*, 2018. 1

[42] Tomoya Nakamura, Jumpei Matsumoto, Hiroshi Nishimaru, Rafael Vieira Bretas, Yusaku Takamura, Etsuro Hori, Taketoshi Ono, and Hisao Nishijo. A markerless 3d computerized motion capture system incorporating a skeleton model for monkeys. *Plos ONE*, 2016. 2

[43] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, 2016. 1, 2, 3, 6

[44] Wanli Ouyang, Xiao Chu, and Xiaogang Wang. Multi-source deep learning for human pose estimation. In *CVPR*, 2014. 2

[45] Seyoung Park, Bruce Xiaohan Nie, and Song-Chun Zhu. Attribute and-or grammar for joint parsing of human pose, parts and attributes. *TPAMI*, 2017. 3

[46] Georgios Pavlakos, Xiaowei Zhou, Konstantinos G. Derpanis, and Kostas Daniilidis. Harvesting multiple views for markerless 3d human pose annotations. In *CVPR*, 2017. 3

[47] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3d human pose and shape from a single color image. In *CVPR*, 2018. 6

[48] Tomas Pfister, James Charles, and Andrew Zisserman. Flowing convnets for human pose estimation in videos. In *ICCV*, 2015. 2

[49] Pedro Pinheiro and Ronan Collober. Recurrent convolutional neural networks for scene labeling. In *ICML*, 2014. 2

[50] Leonid Pishchulin, Mykhaylo Andriluka, Peter Gehler, and Bernt Schiele. Poselet conditioned pictorial structures. In *CVPR*, 2013. 2

[51] Leonid Pishchulin, Mykhaylo Andriluka, Peter Gehler, and Bernt Schiele. Strong appearance and expressive spatial models for human pose estimation. In *ICCV*, 2013. 2

[52] Ilija Radosavovic, Piotr Dollár, Ross Girshick, Georgia Gkioxari, and Kaiming He. Data distillation: Towards omni-supervised learning. In *arXiv:1712.04440*, 2017. 6

[53] Helge Rhodin, Mathieu Salzmann, and Pascal Fua. Unsupervised geometry-aware representation for 3d human pose estimation. In *ECCV*, 2018. 2

[54] Helge Rhodin, Jörg Spörri, Isinsu Katircioglu, Victor Constantin, Frédéric Meyer, Erich Müller, Mathieu Salzmann, and Pascal Fua. Learning monocular 3d human pose estimation from multi-view images. In *CVPR*, 2018. 2, 3, 6, 7, 8

[55] William Sellers and Eishi Hirasaki. Markerless 3d motion capture for animal locomotion studies. *Biology Open*, 2014. 2

[56] Jamie Shotton, Andrew Fitzgibbon, Mat Cook, Toby Sharp, Mark Finocchio, Richard Moore, Alex Kipman, and Andrew Blake. Real-time human pose recognition in parts from single depth images. In *CVPR*, 2011. 3

[57] Leonid Sigal and Michael Black. Measure locally, reason globally: Occlusion-sensitive articulated pose estimation. In *CVPR*, 2006. 2

[58] Tomas Simon, Hanbyul Joo, Iain Matthews, and Yaser Sheikh. Hand keypoint detection in single images using multiview bootstrapping. In *CVPR*, 2017. 2, 3, 5, 6, 7

[59] Jie Song, Limin Wang, Luc Van Gool, and Otmar Hilliges. Thin-slicing network: A deep structured model for pose estimation in videos. In *CVPR*, 2017. 3

[60] Min Sun and Silvio Savarese. Articulated part-based model for joint object detection and pose estimation. In *ICCV*, 2011. 2

[61] Bugra Tekin, Artem Rozantsev, Vincent Lepetit, and Pascal Fua. Direct prediction of 3d body poses from motion compensated sequences. In *CVPR*, 2016. 3

[62] Yuandong Tian, Lawrence Zitnick, and Srinivasa Narasimhan. Exploring the spatial hierarchy of mixture models for human pose estimation. In *ECCV*, 2012. 2

[63] Jonathan Tompson, Arjun Jain, Yann LeCun, and Christoph Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. In *NIPS*, 2014. 2

[64] Alexander Toshev and Christian Szegedy. Deeppose: Human pose estimation via deep neural networks. In *CVPR*, 2014. 1, 2, 3, 6

[65] Shubham Tulsiani, Tinghui Zhou, Alexei A. Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *CVPR*, 2017. 2

[66] Norimichi Ukita and Yusuke Uematsu. Semi- and weakly-supervised human pose estimation. In *CVIU*, 2018. 3

[67] Sudheendra Vijayanarasimhan, Susanna Ricco, Cordelia Schmid, Rahul Sukthankar, and Katerina Fragkiadaki. Sfmnet: Learning of structure and motion from video. In *arXiv:1704.07804*, 2017. 2

[68] Yang Wang and Greg Mori. Multiple tree models for occlusion and spatial constraints in human pose estimation. In *ECCV*, 2008. 2

[69] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *CVPR*, 2016. 1, 2, 3, 6

[70] Xinchen Yan, Jimei Yang, Ersin Yumer, Yijie Guo, and Honglak Lee. Perspective transformer nets: Learning single-view 3d object reconstruction without 3d supervision. In *NIPS*. 2016. 2, 6

[71] Yi Yang and Deva Ramanan. Articulated pose estimation with flexible mixtures-of-parts. In *CVPR*, 2011. 2

[72] Jae Shin Yoon, Ziwei Li, and Hyun Soo Park. 3d semantic trajectory reconstruction from 3d pixel continuum. In *CVPR*, 2018. 3, 5, 6

[73] Tao Yu, Zerong Zheng, Kaiwen Guo, Jianhui Zhao, Qionghai Dai, Hao Li, Gerard Pons-Moll, and Yebin Liu. Doublefusion: Real-time capture of human performances with inner body shapes from a single depth sensor. In *CVPR*, 2018. 3

[74] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G. Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, 2017. 2

[75] Xingyi Zhou, Qixing Huang, Xiao Sun, Xiangyang Xue, and Yichen Wei. Towards 3d human pose estimation in the wild: a weakly-supervised approach. In *ICCV*, 2017. 3

[76] Silvia Zuffi, Angjoo Kanazawa, David Jacobs, and Michael Black. 3d menagerie: Modeling the 3d shape and pose of animals. In *CVPR*, 2017. 3
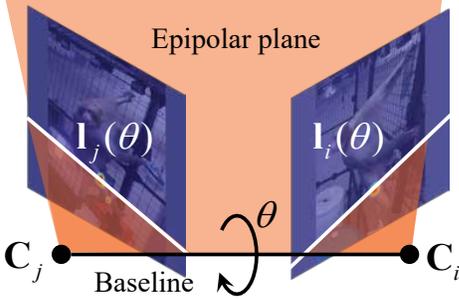
# Supplementary Material



Figure 11: Two epipolar lines are induced by an epipolar plane, which can be parametrized by the rotation $\theta$ about the baseline where $\mathbf{C}_i$ and $\mathbf{C}_j$ are the camera optical centers.

## A. Proof of Theorem 1

*Proof.* A point in an image corresponds to a 3D ray $\mathbf{L}$ emitted from the camera optical center $\mathbf{C}$ (i.e., inverse projection), and $\lambda$ corresponds to the depth. $\mathbf{K}$ is the intrinsic parameter. The geometric consistency, or zero reprojection error, is equivalent to proving $\mathbf{L}_i^*, \mathbf{L}_j^* \in \mathbf{\Pi}$ where $\mathbf{\Pi}$ is an epipolar plane rotating about the camera baseline $\overline{\mathbf{C}_i \mathbf{C}_j}$ as shown in Figure 11, and $\mathbf{L}_i^*$ and $\mathbf{L}_j^*$ are the 3D rays produced by the inverse projection of correspondences $\mathbf{x}_i^* \leftrightarrow \mathbf{x}_j^*$, respectively, i.e., $\mathbf{L}_i^* = \mathbf{C}_i + \lambda \mathbf{R}_i^\mathsf{T} \mathbf{K}^{-1} \widetilde{\mathbf{x}}_i^*$. The correspondence from the keypoint distributions are:

$$\mathbf{x}_i^* = \underset{\mathbf{x}}{\operatorname{argmax}} \, P_i(\mathbf{x}) \tag{12}$$

$$\mathbf{x}_j^* = \underset{\mathbf{x}}{\operatorname{argmax}} \, P_j(\mathbf{x}), \tag{13}$$

$Q_i(\theta) = Q_{j\to i}(\theta)$ implies:

$$
\begin{aligned}
\theta^* &= \underset{\theta}{\operatorname{argmax}} \, \underset{\mathbf{x} \in \mathbf{l}_i(\theta)}{\sup} \, P_i(\mathbf{x}) \\
&= \underset{\theta}{\operatorname{argmax}} \, \underset{\mathbf{x} \in \mathbf{l}_i(\theta)}{\sup} \, P_{j\to i}(\mathbf{x}) \\
&= \underset{\theta}{\operatorname{argmax}} \, \underset{\mathbf{x} \in \mathbf{l}_j(\theta)}{\sup} \, P_j(\mathbf{x}).
\end{aligned}
\tag{14}
$$

This indicates the correspondence lies in epipolar lines induced by the same $\theta^*$, i.e,. $\mathbf{x}_i^* \in \mathbf{l}_i(\theta^*)$ and $\mathbf{x}_j^* \in \mathbf{l}_j(\theta^*)$. Since $\mathbf{l}_j(\theta^*) = \mathbf{F}\widetilde{\mathbf{x}}_i^*$, $\mathbf{l}_i(\theta^*)$ and $\mathbf{l}_j(\theta^*)$ are the corresponding epipolar lines. Therefore, they are in the same epipolar plane, and the reprojection error is zero. $\square$

## B. Cropped Image Correction and Stereo Rectification

We warp the keypoint distribution using stereo rectification. This requires a composite of transformations because the rectification is defined in the full original image. The transformation can be written as:

$$\overline{h}\mathbf{H}_h = \left(\overline{h}\mathbf{H}_{\overline{c}}\right)\left(\overline{c}\mathbf{H}_{\overline{b}}\right)\mathbf{H}_r\left(\overline{c}\mathbf{H}_b\right)^{-1}\left(\overline{h}\mathbf{H}_c\right)^{-1}. \tag{15}$$

The sequence of transformations takes a keypoint distribution of the network output $P$ to the rectified keypoint distribution $\overline{P}$: heatmap→cropped image→original image→rectified image→rectified cropped image→rectified heatmap.

Given an image $\mathcal{I}$, we crop the image based on the bounding box as shown in Figure 12: the left-top corner is $(u_x, u_y)$ and the height is $h_b$. The transformation from the image to the bounding box is:

$$
{}^c\mathbf{H}_b = \begin{bmatrix} s & 0 & w_x - su_x \\ 0 & s & w_y - su_y \\ 0 & 0 & 1 \end{bmatrix}
\tag{16}
$$

where $s = h_c/h_b$, and $(w_x, w_y)$ is the offset of the cropped image. It corrects the aspect ratio factor. $h_c = 364$ is the height of the cropped image, which is the input to the network. The output resolution (heatmap) is often different from the input, $s_h = h_h/h_c \neq 1$, where $h_h$ is the height of the heatmap. The transformation from the cropped image to the heatmap is:

$$
{}^h\mathbf{H}_c = \begin{bmatrix} s_h & 0 & 0 \\ 0 & s_h & 0 \\ 0 & 0 & 1 \end{bmatrix}
\tag{17}
$$

The rectified transformations $\left(\overline{h}\mathbf{H}_{\overline{c}}\right)$ and $\left(\overline{c}\mathbf{H}_{\overline{b}}\right)$ can be defined in a similar way.

The rectification homography can be computed as $\mathbf{H}_r = \mathbf{K}\mathbf{R}_n\mathbf{R}^\mathsf{T}\mathbf{K}^{-1}$ where $\mathbf{K}$ and $\mathbf{R} \in SO(3)$ are the intrinsic parameter and 3D rotation matrix and $\mathbf{R}_n$ is the rectified rotation of which x-axis is aligned with the epipole, i.e., $\mathbf{r}_x = \dfrac{\mathbf{C}_j - \mathbf{C}_i}{\|\mathbf{C}_j - \mathbf{C}_i\|}$ where $\mathbf{R}_n = \begin{bmatrix} \mathbf{r}_x^\mathsf{T} \\ \mathbf{r}_y^\mathsf{T} \\ \mathbf{r}_z^\mathsf{T} \end{bmatrix}$ and other axes can be computed by the Gram-Schmidt process.

The fundamental matrix between two rectified keypoint distributions $\overline{P}_i$ and $\overline{P}_j$ can be written as:

$$
\begin{aligned}
\mathbf{F} &= \mathbf{K}_j^{-\mathsf{T}} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}_\times `\mathbf{K}_i^{-1} \\
&= \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1/f_y^j \\ 0 & 1/f_y^i & p_y^j/f_y^j - p_y^i/f_y^i \end{bmatrix}
\end{aligned}
\tag{18}
$$

where $[\cdot]_\times$ is the skew symmetric representation of cross product, and

$$
\mathbf{K}_i = \begin{bmatrix} f_x^i & 0 & p_x^i \\ 0 & f_y^i & p_y^i \\ 0 & 0 & 1 \end{bmatrix}.
\tag{19}
$$

$(w_x, w_y)$

$(u_x, u_y)$

$h_b$

${}^{c}\mathbf{H}_b$

$h_c$

${}^{h}\mathbf{H}_c$

$h_h$

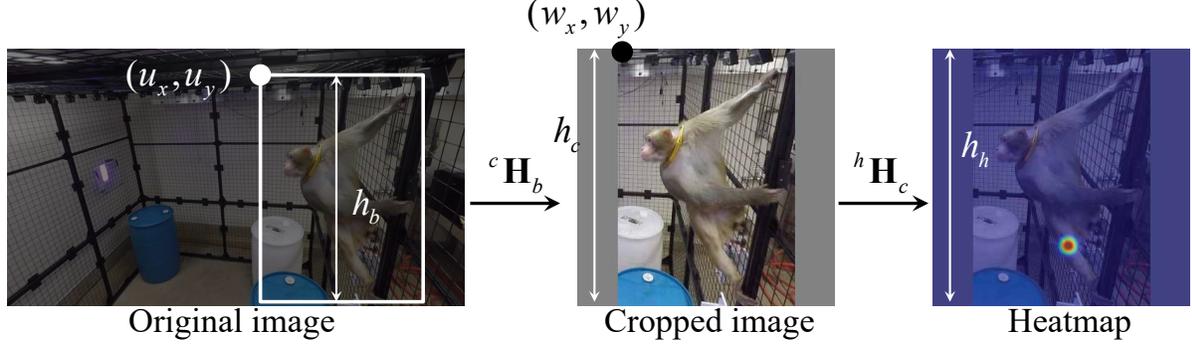Original image       Cropped image       Heatmap

Figure 12: A cropped image is an input to the network where the output is the keypoint distribution. To rectify the keypoint distribution (heatmap), a series of image transformations need to be applied.

| Subjects | $P$ | $|\mathcal{D}_L|$ | $|\mathcal{D}_U|$ | $|\mathcal{D}_L|/|\mathcal{D}_U|$ | $C$ | FPS | Camera type |
|----------|-----|-------------------|-------------------|-----------------------------------|-----|-----|-------------|
| Monkey   | 13  | 85   | 63,000  | 0.13% | 35 | 60 | GoPro 5 |
| Humans   | 14  | 30   | 20,700  | 0.14% | 69 | 30 | FLIR BlackFly S |
| Dog I    | 12  | 100  | 138,000 | 0.07% | 69 | 30 | FLIR BlackFly S |
| Dog II   | 12  | 75   | 103,500 | 0.07% | 69 | 30 | FLIR BlackFly S |
| Dog III  | 12  | 80   | 110,400 | 0.07% | 69 | 30 | FLIR BlackFly S |
| Dog IV   | 12  | 75   | 103,500 | 0.07% | 69 | 30 | FLIR BlackFly S |

Table 3: Summary of multi-camera dataset where $P$ is the number of keypoints, $C$ is the number of cameras, $|\mathcal{D}_L|$ is the number of labeled data, and $|\mathcal{D}_U|$ is the number of unlabeled data.

This allows us to derive the re-scaling factor of $a$ and $b$ in Equation (7):

$$a = \frac{s^i f_y^i}{s^j f_y^j} \tag{20}$$

$$b = s_h s^i \left( \left( \overline{u}_y^j - p_y^j \right) \frac{f_y^i}{f_y^j} + p_y^i - \overline{u}_y^i \right) \tag{21}$$

where $\overline{u}_y^i$ is the bounding box offset of the rectified coordinate.

## C. Evaluation Dataset

All cameras are synchronized and calibrated using structure from motion [18]. The input of most pose detector models except for [8] is a cropped image containing a subject, which requires specifying a bounding box. We use a kernelized correlation filter [20] to reliably track a bounding box using multiview image streams given initialized 3D bounding box from the labeled data.