

Transferring Dense Pose to Proximal Animal Classes

Artsiom Sanakoyeu*
Heidelberg University

Vasil Khalidov
Facebook AI Research

Maureen S. McCarthy
MPI for Evolutionary Anthropology

Andrea Vedaldi
Facebook AI Research

Natalia Neverova
Facebook AI Research

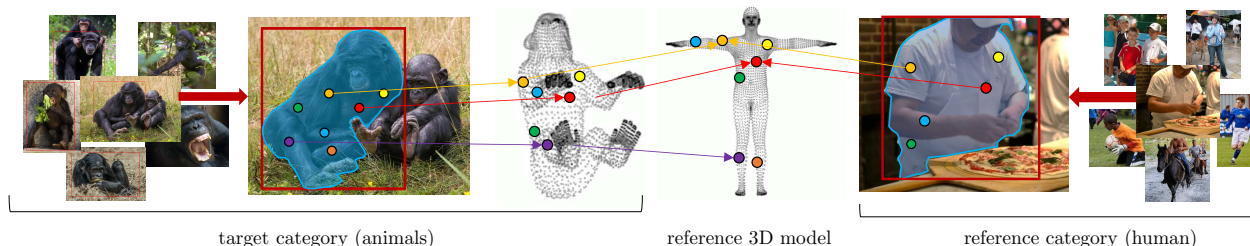


Figure 1: We consider the problem of dense pose labelling in animal classes. We show that, for proximal humans classes such as chimpanzees (left), we can obtain excellent performance by learning an integrated recognition architecture from existing data sources, including DensePose for humans as well as detection and segmentation information from other COCO classes (right). The key is to establish a common reference (middle), which we obtain via alignment of the reference models of the animals. This enables training a model for the target class without having to label a single example image for it.

Abstract

Recent contributions have demonstrated that it is possible to recognize the pose of humans densely and accurately given a large dataset of poses annotated in detail. In principle, the same approach could be extended to any animal class, but the effort required for collecting new annotations for each case makes this strategy impractical, despite important applications in natural conservation, science and business. We show that, at least for proximal animal classes such as chimpanzees, it is possible to transfer the knowledge existing in dense pose recognition for humans, as well as in more general object detectors and segmenters, to the problem of dense pose recognition in other classes. We do this by (1) establishing a DensePose model for the new animal which is also geometrically aligned to humans (2) introducing a multi-head R-CNN architecture that facilitates transfer of multiple recognition tasks between classes, (3) finding which combination of known classes can be transferred most effectively to the new animal and (4) using self-calibrated uncertainty heads to generate pseudo-labels graded by quality for training a model for this class. We also introduce two benchmark datasets labelled in the manner of DensePose for the class chimpanzee and use them to evaluate our approach, showing excellent transfer learning performance.

1. Introduction

In the past few years, computer vision has made significant progress in human pose recognition. Deep networks can effectively detect and segment humans [14], localize their sparse 2D keypoints [36], lift these 2D keypoints to 3D [37], and even fit complex 3D models such as SMPL [20, 21], all from a single picture or video. DensePose [11] has shown that it is even possible to estimate a dense parameterization of pose by mapping individual image pixels to a canonical embedding space for the human body.

Such advances have been made possible by the introduction of large human pose datasets manually annotated with sparse or dense 2D keypoints, or even in 3D by means of capture systems such as domes. For example, the DensePose-COCO dataset [11] contains 50K COCO images manually annotated with more than 5 millions human body points. Clearly, collecting such data is very tedious, but is amply justified by the importance of human understanding in applications. However, the natural world contains much more than just people. For example, as of today scientists have identified 6,495 species of mammals, 60k vertebrates and 1.2M invertebrates [1]. The methods that have been developed for human understanding could likely be applied to most of these animals as well, provided that one is willing to incur the data annotation burden. Unfortunately, while the

*Work done during an internship at Facebook AI Research

Project page: <https://asanakoy.github.io/densepose-evolution>

applications of animal pose recognition in conservation, natural sciences, and business are numerous, just learning about one more animal may be difficult to justify economically, let alone learning about *all* animals.

Yet, there is little reason to believe that these challenges are intrinsic. Humans can understand the pose of most animals almost immediately, with good accuracy, and without requiring any data annotations at all. Furthermore, images and videos of animals are abundant, so the bottleneck is the inability of machines to learn without external supervision.

In this paper, we thus consider the problem of learning to recognize the pose of animals with as little supervision as possible. However, rather than starting from scratch, we want to make use of the rich annotations that are *already* available for several animals, and humans in particular. Thus, we focus on the problem of taking the existing annotated data as well as additional unlabelled images and videos of a target animal species and learn to recognize the pose of the latter. Furthermore, for this study we restrict our attention to an animal species that is reasonably close to the available annotations, and elect to focus on the particular example of chimpanzees due to their evolutionary closeness to humans.* However, the findings in this paper are likely to generalize to many other classes as well.

We make several contributions in this work. First, we introduce a dataset for chimpanzees, *DensePose-Chimps*, labelled in the DensePose fashion, which we mostly use to assess quantitatively the performance of our methods. We carefully design the canonical mapping for chimpanzees to be compatible with the one for humans in the original DensePose-COCO, in the sense that points in the two animal models are in as close a correspondence as possible. This is essential to be able to transfer dense pose recognition results from humans to chimpanzees while being able to assess the quality of the obtained results.

Second, we study in detail several strategies to transfer existing animal detectors, segmenters, and dense pose extractors from the available annotated data to chimpanzees. In particular, while dense pose annotations exist only for humans, bounding box and mask annotations have been collected for several other object categories as well. As a representative source dataset we thus consider COCO and we investigate how the different COCO classes can be combined to train an object detector and segmenter that transfers optimally to chimpanzees. Surprisingly, we find that transfer from humans alone is not optimal, nor human is the best class for training a model for chimpanzees. In addition to the DensePose-Chimps data, we collect human annotations for instance masks on the *Chimp&See*[†] videos of chimpanzees

*The idea is to eventually extend pose recognition to more and more animal species, in an incremental fashion.

[†]Some of these videos are available at <http://www.zooniverse.org/projects/sassydumbledore/chimp-and-see>.

captured with camera traps in the wild to evaluate the detection performance in the most challenging conditions (with severe occlusions, low visibility and motion blur).

Finally, we propose a framework for augmenting and adapting the human DensePose datasets to new species by self-supervision and pseudo-labeling with zero ground truth annotations on the target class.

2. Related work

Human pose recognition. There is abundant work on the recognition of human body pose, both in 2D and in 3D. Given that our focus is 2D pose recognition, we discuss primarily the first class of methods. 2D human pose recognition has flourished by the introduction of deep neural networks [46, 36, 8] trained on large manually-annotated datasets of images and videos such as COCO [27], MPII [3], Leeds Sports Pose Dataset (LSP) [18, 19], PennAction [50] and Posetrack [2]. Furthermore, Dense Pose [11] has introduced a dataset with dense surface point annotations, mapping images to a UV representation of a parametric 3D human model (SMPL) [29].

While all such approaches are strongly-supervised, there are also methods that attempt to learn pose in a completely unsupervised manner [5, 43, 44, 41, 42, 30, 51]. Unfortunately, this technology is not sufficiently mature to compete with strong supervision in the wild.

Animal pose recognition. Also related to our work, several authors have learned visual models of animals for the purpose of detection, segmentation, and pose recognition. Some animals are included in almost all general-purpose 2D visual recognition datasets, and in COCO in particular. Hence, all recent detectors and segmenters have been tested on at least a few animal classes.

For pose recognition, however, the existing body of research is more restricted. Some recent papers have focused on designing pose estimation systems and benchmarks for particular animal species such as Amur tigers [26], cheetahs [33] or drosophila melanogaster flies [12]. There have been a number of large efforts on designing annotation tools for animals, such as DeepLabCut [31] and Anipose [22]. These tools also provide functionality for lifting 2D keypoints to 3D by using multiple views and triangulation. A more detailed overview on applying computer vision and machine learning methodology in neuroscience and zoology is given in [32]. One of the main challenges in this field remains the narrow focus of existing research on specific kinds of animals and particular environments.

There have been few works focusing on the problem of animal understanding from visual data alone and in a more systematic way. This includes the estimation of facial landmarks through domain adaptation [48, 39], and very recently full body pose estimation [7] of four-legged animals by com-

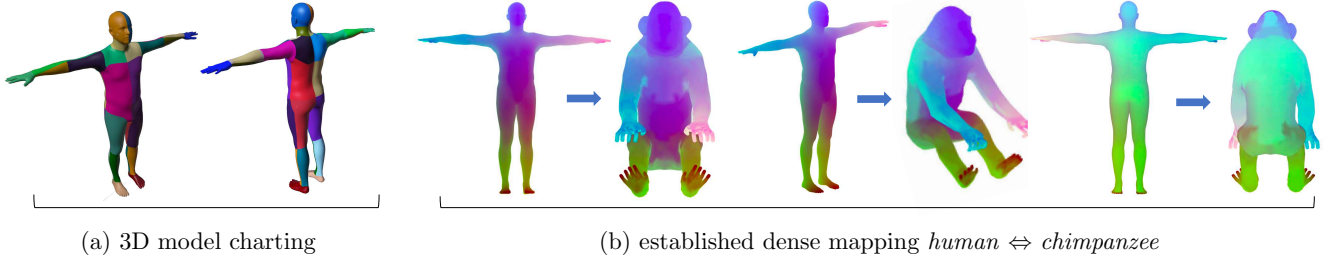


Figure 2: 3D shape re-mapping from the SMPL model for humans to new object categories (chimps). Manually defined semantic charting (a) on both models is used to establish dense correspondences (b) based on continuous semantic descriptors

binning large-scale human datasets with a smaller number of animal annotations in a cross-domain adaptation framework. Finally, a line of work from Zuffi et al. [54, 53, 52] is exploring the problem of model-based 3D pose and shape estimation for animal classes. Their research is based on parametric linear model, Skinned Multi-Animal Linear (SMAL), obtained from 3D scans of toy animals and having the capacity to represent multiple classes of mammals. SMAL is the animal analogous of the popular SMLP [28] model for humans. It has since been used in other publications [6] for 3D animal reconstruction, but these methods may still be insufficiently robust for deployment in the wild.

Unsupervised and less supervised pose recognition. Recent methods such as [43, 44, 42, 17, 51, 30] learn sparse and dense object landmarks for simple classes without making use of any annotation, but are too fragile to be used in our application. Also relevant to our work, Slim DensePose [35] looked at reducing the number of annotations required to learn a good DensePose model for humans.

Self-training for dense prediction. A recent study [47] has demonstrated effectiveness of self-training on the task of image classification when scaled to large amounts of unlabeled data. Pseudo-labeling by averaging predictions from multiple transformed versions of unlabeled samples has been shown effective for keypoint estimation [38]. However, there has been very little research on self-training in the context of dense prediction tasks. A recent work [4] explored the idea of self-training for segmentation of seismic images and showed promising results on this task for the first time.

3. Method

We wish to develop a methodology to learn Dense Pose models for new classes with minimal annotation effort. Existing labelled datasets for object detection, segmentation and pose estimation, provide a significant source of supervision that can be harnessed for this task. For detection and segmentation, COCO provide extensive annotations for a variety of object classes, including several animals. For pose recognition, however, the available supervision is generally

limited to humans, with a few exceptions. Furthermore, for *dense* pose recognition only human datasets are available — the best example of which is DensePose-COCO [11].

In this work, we raise a number of questions most critical for this setup, namely:

- defining learning and evaluation protocols on new animal categories allowing for training class-specific or class-agnostic DensePose models on a variety of species in a unified way (described in Sect. 3.1);
- improving quality of DensePose models and their robustness to unseen data distributions at test time (discussed in Sect. 3.2 and 3.3);
- optimally combining the existing variety of data sources in order to initialize a detection model for a new animal species (discussed in Sect. 3.4);
- defining strategies for mining dense pseudo-labels for gradual domain adaptation from humans to chimpanzees in a teacher-student setting (discussed in Sect. 3.5).

3.1. Annotation through 3D shape re-mapping

While our aim is to learn to reconstruct the dense pose of chimpanzees with zero supervision, a manually-annotated dataset for this class is required for evaluation. Here, we explain how to collect DensePose annotations for a new category, such as chimpanzees.

Dense Pose model. Recall that DensePose-COCO contains images of people collected ‘in the wild’ and annotated with dense correspondences. These dense keypoints are identified as the point $p \in S$ of a reference 3D model $S \subset \mathbb{R}^3$ of the object.[‡] Furthermore, the keypoints $p \in S$ are indexed by triplets $(c, u, v) \in \{1, \dots, C\} \times [0, 1]^2$ where c is the *chart index*, corresponding to one of C model parts, and (u, v) are the coordinates within a chart. The DensePose-COCO dataset [11] contains bounding boxes, pixel-perfect foreground-background and part segmentations, and (c, u, v) annotations for a large number of foreground pixels.

Dense Pose for chimps. We wish to extend the DensePose annotations to the chimpanzee class. In order to do so, we

[‡]Dense Pose uses SMPL [29] to define S due to its popularity

rely on a separate artist-created 3D model[§] of a chimpanzee as a reference for annotators to collect labels for the chimpanzee images (instead of the human model used by the original DensePose).

For each object, we use Amazon Mechanical Turk to collect the object bounding boxes, followed by pixel-perfect foreground/background segmentation masks, and finally the (c, u, v) chart coordinates for a certain number of pixels randomly sampled from the foreground regions. Differently from the original DensePose, we *do not* also collect dense annotations for the body parts as the latter was found to be very challenging for the annotators. Note however that the chart index c reveals the part identity for each of the annotated image pixels.

Semantic alignment. Finally, we wish to align the human and chimpanzee DensePose models by mapping the collected annotations back on the surface of the SMPL model using the mesh re-mapping strategy described below. The latter step unifies the evaluation protocols across different object categories and allows to transfer knowledge and annotations between different species.

In spite of the fact that humans and most mammals share topology and the skeletal structure, establishing precise semantic dense correspondences between the 3D models of humans and different animal species is challenging due to differences in body proportions and local geometry.

As preprocessing, we manually charted the SMPL and the chimp meshes into $L = 32$ semantically-corresponding parts to guide the mapping. Then, for each vertex p of each mesh S , we extracted an adapted version of the continuous semantic descriptor $\mathbf{d}(p)$ proposed by Léon et al. [25]:

$$\mathbf{d}(p) = (d_\ell(p))_{\ell=1}^L, \quad d_\ell(p) = \frac{1}{|S_\ell|} \sum_{s \in S_\ell} g(p, s; S_\ell) \quad (1)$$

where $S_\ell \subset S$ is the set of all vertices in part ℓ of the mesh and $g(p, s)$ is the geodesic distance between two points on S .[¶] With this, the mapping from the human mesh S to the chimp mesh S' is obtained by matching nearest descriptors: $S \rightarrow S', p \mapsto \operatorname{argmin}_{q \in S'} \|\mathbf{d}_S(p) - \mathbf{d}_{S'}(q)\|^2$.

This simple approach yields satisfactory results both in terms of alignment and smoothness, as shown in Fig. 2. It does not require any optimization in 3D space based on model fitting or mesh deformation and works on meshes of arbitrary resolutions. Interestingly, exploiting information about mesh geometry (such as high dimensional SHOT [40] descriptors or their learned variants [13]) instead or in addition to semantic features results in noisy mappings. This can likely be attributed to prominent inconsistencies in local geometry of some body regions between the object categories.

[§]Purchased from <http://hum3d.com/>

[¶]To partially compensate for differences in proportions across different categories, we further normalized the descriptors by their part average: $d_\ell(p) \leftarrow d_\ell(p) / \langle d_\ell(q) \rangle_{q \in S_\ell}$.

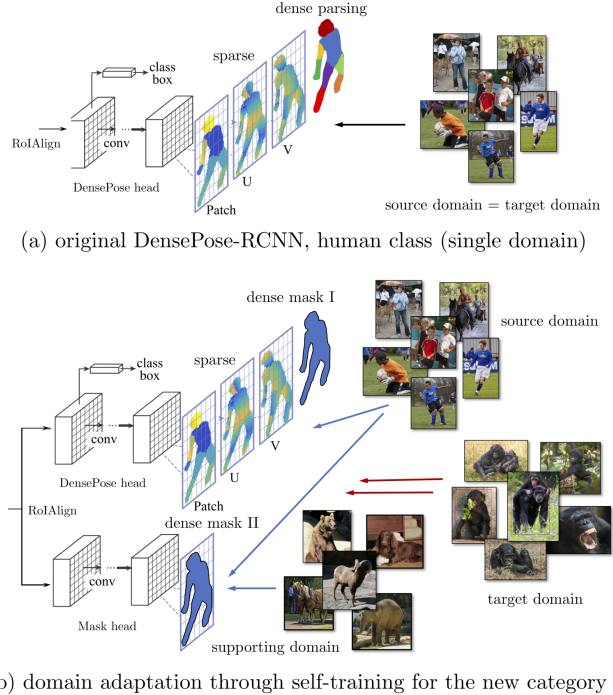


Figure 3: Comparison of the original (a) and our (b) DensePose learning architecture. See Sect. 3.2 for detailed description of the architecture.

3.2. Multi-head R-CNN

Our goal is to develop a DensePose predictor for a new class. Such a predictor must detect the object via a bounding box, segment it from the background, and obtain the DensePose chart and uv -map coordinates for each foreground pixel. We implement this with a *single model* with multiple heads, performing the various tasks on top of the same trunk and shared image features (Fig.3.b).

The base model is R-CNN [14] modified to include the following heads. The first head refines the coordinates of the bounding box. The second head computes a foreground-background segmentation mask in the same way as Mask R-CNN. The third and the final head computes a part segmentation mask I , assigning each pixel to one of the 24 Dense Pose charts, and the uv map values for each foreground pixel.

Class-agnostic model. Compared to the standard Mask R-CNN, our model is *class agnostic*, i.e. trained for only one class type. This is true also when we make use of a Mask R-CNN pre-trained on multiple source classes as the goal is always to only build a model for the final target chimpanzee class — we found that merging classes is an effective way of integrating information.

Heterogeneous training. Our training data can be heterogeneous. In particular, COCO provides segmentation masks for 80 categories, but DensePose-COCO provides DensePose annotations only for humans. While we train a single class-agnostic model, the Dense Pose head is trained only for the class human for which the necessary ground-truth data is available.

Note in particular that both the Mask R-CNN head and the DensePose head contain a foreground-background segmentation component — these are not equivalent, as the DensePose one is only valid (and trainable) for humans, while the Mask R-CNN one is generic (and trainable from all COCO classes). We will see in the experiments that their combination improves performance.

Fine-tuning. As shown later, for fine-tuning the model we generate pseudo-label on chimpanzees imagery. The pseudo-labels are generated for all components of the model (segmentations, uv maps), including in particular both foreground-background segmentation heads.

Other architectural improvements. Our model (Fig. 3.a) has a few mode differences compared to the original Dense Pose (Fig. 3.b) which we found useful to improved accuracy and/or data collection efficiency.

First, both the original and our implementations use dense (pixel-perfect) supervision for the foreground-background masks. However, in our version we *do not* use the pixel-perfect part segmentations in the original DensePose annotations — the part prediction head is trained only from the chart labels for the pixels that are annotated in the data. This is another reason why we do not collect pixel-perfect segmentations for the chimpanzee images.

We further improve the DensePose head by implementing it using Panoptic Feature Pyramid Networks [24], and use a configuration similar to DeepLab [10] that benefits from higher resolution.

3.3. Auto-calibrated R-CNN

As suggested above, pseudo-labelling can be used to fine-tune a pre-trained model on imagery containing the target class, chimpanzees in our case. The idea is to use a model pre-trained on a different class or set of classes to generate labels in the new domain, and then to retrain the model to fit those labels. Due to the domain gap, however, the pseudo-labels are somewhat unreliable. In this section, following [23] we develop a principled manner to let the neural network itself produce a *calibrated measure of uncertainty* which we can use to rank pseudo-labels by reliability.

Classification uncertainty. Our model performs categorical classification for two purposes: to associate a class label to a bounding box, and to classify individual pixels as background, foreground, or as one of the body parts. In order to estimate the uncertainty for these categorical predictions, we

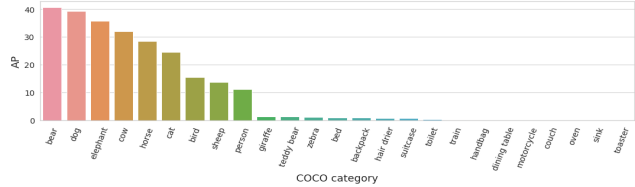


Figure 4: Instance Segmentation score (AP) on DensePose-Chimps for Mask R-CNN models trained using different COCO categories, ranked by decreasing performance.

adopt the *temperature scaling* technique of [16].

Thus let z_y be the score that the neural network associates to hypothesis $y \in \{1, \dots, K\}$ for a given input sample. We extend the network to compute an additional per-sample scalar $\alpha \geq 0$. With this scalar, the posterior probability of hypothesis y is given by the *scaled softmax*

$$\hat{\sigma}(y; z, \alpha) = \frac{\exp(\alpha z_y)}{\sum_{k=1}^K \exp(\alpha z_k)} \quad (2)$$

We can interpret the coefficient $\alpha = 1/T$ as an inverse temperature. A small α means that the model is fairly certain about the prediction, whereas a large α that it is not.

Note that, since α is also estimated by the neural network, we require a mechanism to learn it. This is in fact obtained automatically [16, 34] by simply minimizing the negative log-likelihood of the model, also known in this case as cross-entropy loss: $\ell(y, z, \alpha) = -\log \hat{\sigma}(y; z, \alpha)$.

Regression with uncertainty. Our model performs regression to refine the bounding box proposals (for four scalar outputs, two for each of the two corners of the box) and to obtain the DensePose uv -coordinates (for two scalar outputs for each image pixel in a proposal).

Thus let $y \in \mathbb{R}^D$ be the vector emitted by one of the regression heads (where D depends on the head). Similarly to the classification case, we use the network to also predict an *uncertainty score* $\sigma \in \mathbb{R}^D$. This time, however, we have a different scalar for each element in y (hence, for the uv -maps, we have two uncertainty scores for each pixel, which we can visualize as an image). The vector σ is interpreted as the diagonal variance of the regressed vector y , assuming the latter to have a Gaussian distribution. The uncertainty scores σ can thus be trained jointly with the predictor \hat{y} by minimizing the negative log-likelihood of the model:

$$\ell(y, \hat{y}, \sigma) = \frac{D}{2} \log 2\pi + \frac{1}{2} \sum_{i=1}^D \left(\log \sigma_i^2 + \frac{(\hat{y}_i - y_i)^2}{\sigma_i^2} \right) \quad (3)$$

For a fixed error $|\hat{y}_i - y_i|$, the quantity above is minimized by setting $\sigma_i = |\hat{y}_i - y_i|$ — hence the model is encouraged to guess the magnitude of its own prediction error. However, if $|\hat{y}_i - y_i| = 0$, the quantity above diverges to $-\infty$ for $\sigma_i \rightarrow 0$. Hence, we clamp σ_i from below to a minimum value $\sigma_{\min} > 0$.

model	AP	AP ₅₀	AP ₇₅
DensePose-RCNN	50.88	80.40	54.80
DensePose-RCNN*	51.44	81.44	55.12
DensePose-RCNN* (σ)	54.13	82.32	58.06

model	AP	AP ₅₀	AP ₇₅
DensePose-RCNN	43.84	76.88	45.84
DensePose-RCNN*	43.84	77.52	45.60
DensePose-RCNN* (σ)	45.58	78.79	47.93

Table 1: Detection (left) and instance segmentation (right) performance on DensePose-COCO *minival*.

model	AP	AP ₅₀	AP ₇₅	AP _M	AP _L	AR	AR ₅₀	AR ₇₅	AR _M	AR _L
DensePose-RCNN	46.8	84.5	47.7	41.8	48.0	54.7	89.5	58.9	43.3	55.5
DensePose-RCNN*	47.2	85.8	47.3	42.5	48.4	55.2	91.0	59.1	44.0	55.9
DensePose-RCNN* (σ)	53.2	88.3	57.0	48.6	54.6	61.2	92.4	67.2	50.0	61.9

Table 2: DensePose performance on DensePose-COCO *minival*. * denotes our improved architecture; (σ) denotes the proposed Auto-calibrated version of the network.

Details. For both classification and regression models, the uncertainties α and σ must be positive — in the network, they are obtained via a `softplus` activation.

3.4. Optimal transfer support

In this section, we investigate which object categories in the COCO dataset provide the best support for recognizing a new animal species, chimpanzees in our case. Among the animals in COCO, chimpanzees are most obviously related to humans, and we may thus expect that people may be the most transferable class. However, despite their overall structural similarity, people’s appearance is fairly different, also due to the lack of fur and the presence of clothing. Furthermore, context is also often quite different. It is thus unclear if a deep network trained to recognise humans can transfer well at all on chimpanzees, or whether other object categories might do better.

Class selection. We test what is more important: biological proximity of the species (as a proxy to morphological similarity) or appearance similarity (as a combination of typical poses and textures). We also search for a brute force solution for this particular dataset to back up or disprove our intuition for class selection. In our experiments, we have tested the following selections:

- *person* class only (due to morphological similarity).
- *animal* classes only (due to higher pose and texture similarity): *bear, dog, elephant, cat, horse, cow, bird, sheep, zebra, giraffe, mouse*.
- *top-N* scoring classes on the new category (brute force solution). In this setting, we first train a set of C single-class models for each of the $C = 90$ object classes in the COCO dataset and rank them according to their instance segmentation performance on the DensePose-Chimps dataset (see Fig. 4). Then for each combination of $S \in \{1, \dots, C\}$ top scoring classes we train the same network from scratch. The solution that we found optimal corresponds to $C_{\text{opt}} = 9$, where the top- C scoring classes are: *bear, dog, elephant, cat, horse, cow, bird, person, sheep*.

As shown in Tab. 5, the top- N solution produces similar results compared to combination *person+animals*. *Person* class only is ineffective for training in this setting.

Class fusion. We have also explored the question of class-agnostic vs multi-class training as a trade-off between the number of training samples per class vs granularity of prediction modes. For the task of adapting the new model to a single category (on the given dataset) class-agnostic training showed convincingly stronger results (see Tab. 5).

3.5. Dense label distillation

Finally, we aim at finding an effective strategy for exploiting unlabeled data for the target domain in the teacher-student training setting and performing *distillation* in dense prediction tasks. In our setting, the *teacher* network trained on the selected classes of the COCO dataset with DensePose is used to generate *pseudo-labels* for fine-tuning the *student* network on the augmented data. The *student* network is initialized with *teacher*’s weights.

Once teacher predictions on unlabeled data are obtained, we start by filtering out low confidence detections using calibrated detection scores. After that, the bounding boxes and segmentation masks on remaining samples are used for augmented training. For mining DensePose supervision, we consider three different dense sampling strategies driven by each of the tasks solved by the teacher network, in addition to uniform sampling:

- **uniform sampling** – all points from the selected detections are sampled with equal probability;
- **coarse classification uncertainty [mask-based]** – sampling top k from ranked calibrated posteriors produced by the mask branch for the task of binary classification;
- **fine classification uncertainty [I-based]** – selection of top k from ranked calibrated posteriors from the 24-way segmentation outputs of the DensePose head;
- **regression uncertainty sampling [uv-based]** – sampling of top k points based on ranked confidences in the *uv*-outputs of the DensePose head.

		DensePose-Chimps			Chimp&See	
sampling	k	AP_{DPose}	AP_D	AP_S	AP_D	AP_S
-	-	33.4	62.1	56.4	50.5	43.5
uniform	5	34.5 ± .4	63.3 ± .3	58.0 ± .3	58.9 ± .5	49.0 ± .5
mask-based	5	34.7 ± .4	63.3 ± .3	58.0 ± .2	58.8 ± .6	49.0 ± .5
<i>I</i> -based	5	34.9 ± .6	63.4 ± .3	58.0 ± .2	59.2 ± .4	49.2 ± .5
<i>uv</i> -based	5	34.6 ± .3	63.3 ± .3	58.2 ± .3	59.0 ± .1	49.6 ± .1

Table 3: AP of the *student* network trained with different sampling strategies. Optimal number of sampled points k per detection is reported for each sampling. The first row corresponds to the *teacher* network. $Mean \pm std$ for 20 runs.

In Sect. 4 we provide experimental evidence that sampling based on confidence estimates from fine-grained tasks (*I*-estimation, *uv*-maps) results in the best *student* performance.

4. Experiments

We now describe the results of empirical evaluation and provide detailed descriptions of ablation studies.

4.1. Datasets

We use a combination of human and animal datasets with different kinds of annotations or no annotations at all. A brief description of each of them is provided below.

DensePose-COCO dataset [11]. This is the dataset for human dense pose estimation, that we use for training the teacher model. It contains 50k annotated instances totalling to more than 5 million ground truth correspondences. We also augment the teacher training with other object categories from the original COCO dataset [27].

Chimp&See dataset. For training our models in a self-supervised setting, we used unlabeled videos containing chimpanzees from the *Chimp&See* project^{||}. This data is being collected under the umbrella of The Pan African Programme^{**}: The Cultured Chimpanzee (PanAf) by installing camera traps in more than 40 natural habitats of chimpanzees on different sites in Africa. In this work, we used a subset of the collected data consisting of 18556 video clips, from 10 sec to 1 min long each, captured with cameras in either standard or night vision mode depending on lighting conditions. These recordings were motion triggered automatically by passing animals. As a result, some clips may not contain any chimps beyond first several frames.

For evaluation, we chose videos from one site, sampled frames at 1 fps, removed the near duplicates and collected human annotations for instance masks. This resulted into 1054 images containing 1528 annotated instances, that we use to

^{||}A subset of the videos from the Chimp&See dataset is publicly available at <http://www.zooniverse.org/projects/sassydumbledore/chimp-and-see>.

^{**}<http://panafrican.eva.mpg.de>

		DensePose-Chimps			Chimp&See	
k		$AP_{DensePose}$	AP_D	AP_S	AP_D	AP_S
0		33.8 ± .2	63.1 ± .2	57.9 ± .2	59.0 ± .3	49.2 ± .4
1		34.7 ± .5	63.0 ± .2	57.9 ± .3	59.3 ± .3	49.3 ± .6
2		34.6 ± .6	63.4 ± .3	57.9 ± .3	59.2 ± .4	49.3 ± .4
5		34.9 ± .5	63.4 ± .3	58.0 ± .2	59.2 ± .4	49.2 ± .5
10		34.6 ± .6	63.3 ± .3	58.0 ± .3	59.2 ± .4	49.4 ± .4
1000		33.1 ± .6	63.2 ± .2	57.8 ± .3	59.2 ± .5	49.4 ± .5
10000		27.6 ± 4.6	60.2 ± .4	55.7 ± .5	58.0 ± .7	49.1 ± .6

Table 4: DensePose, detection and instance segmentation AP of the *student* network trained with *I*-sampling for different number of sampled points k . $Mean \pm std$ for 20 runs.

selected COCO object classes	AP	AP ₅₀	AP ₇₅
top-9 classes	57.29	85.63	63.45
bear-only	40.69	70.88	44.23
person-only	9.39	19.32	8.21
animals-only	52.28	80.62	58.60
person + animals	57.34	85.76	63.59
person + animals: class agnostic	57.34	85.76	63.59
person + animals: class specific	50.47	72.85	54.30

Table 5: Instance segmentation AP on DensePose-Chimps for Mask R-CNN trained on different subsets of classes.

benchmark detection performance in our models. However, due to in-the-wild nature of this data and presence of motion blur, severe occlusions, and low resolution in some cases, we found it infeasible to collect precise human annotations at the level of dense correspondences.

DensePose-Chimps test set. For the task of evaluating DensePose performance on this new category, we collected a set of 662 higher quality images from Flickr that contain 933 instances of chimpanzees. We annotated this data with bounding boxes, binary masks, body part segmentation and dense pose correspondences as explained in Sect. 3.1.

4.2. Results

Ablations on architectural choices. First, we compare our model to the original DensePose-RCNN [11] (detection2 implementation). We also ablate our improvements in the architecture and provide results with and without auto-calibration. Tab. 1, 2 show consistent improvements on all tasks for both modifications.

Optimal transfer support. We (a) benchmarked every strategy for class selection described in Sect. 3.4 and (b) experimented with multi-class and class-agnostic models. From Tab. 5 we can see that class agnostic training on the *animals+person* subset shows the best transferability for DensePose-Chimps dataset. Therefore, it was used for training all our DensePose models.

Dense label distillation. We conducted experiments with different sampling strategies and different numbers of sampled points k per detection. In Tab. 3 we show performance

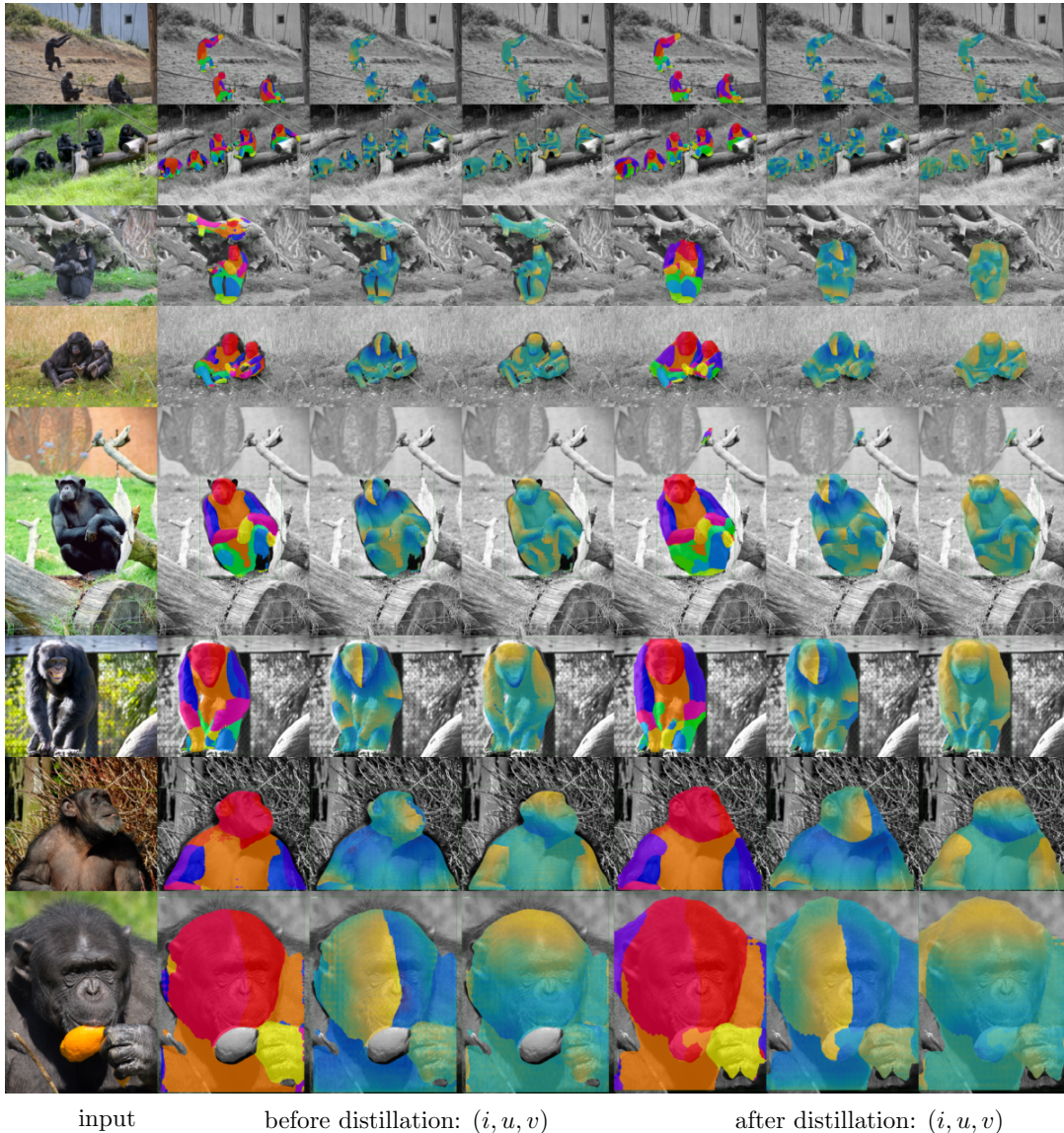


Figure 5: Visual results: (left) *teacher* network predictions vs (right) predictions of *student* network trained using *I*-sampling. The *student* produces more accurate boundaries and *uv*-maps. Zoom-in for details.

of the *teacher* (first row) and the *student* networks trained using different sampling strategies along with the corresponding optimal k . *I*-based sampling showed most impressive gains, followed by *uv*-based sampling. Uniform selection produces poor results. In Tab. 4 we report performance for different number of sampled points in every detection for *I*-based sampling. Qualitative results are shown in Fig. 5.

5. Conclusions

We have studied the problem of extending dense body pose recognition to animal species and suggested that doing this at scale requires learning from unlabelled data. Encouragingly, we have demonstrated that existing detection,

segmentation, and dense pose labelling models can transfer very well to a proximal animal class such as chimpanzee despite significant inter-class differences. We have shown that substantial improvements can be obtained by carefully selecting which categories to use to pre-train the model, by using a class-agnostic architecture to integrate different sources of information, and by modelling labelling uncertainty to grade pseudo-label for self-training. In this manner, we have been able to achieve excellent performance without using a single labelled image of the target class for training.

In the future, we would like to investigate how a limited amount of target supervision can be best used to improve the results, and how other techniques from domain adaptation could also be used for this purpose.

6. Acknowledgements

We thank all parties performing or supporting collection of the Chimp&See dataset, including:

- (a) individual contributors: Theophile Desarmeaux, Kathryn J. Jeffery, Emily Neil, Emmanuel Ayuk Ayimisin, Vincent Lapeyre, Anthony Agbor, Gregory Brazzola, Floris Aubert, Sebastien Regnaut, Laura Kehoe, Lucy DAuvergne, Nuria Maldonado, Anthony Agbor, Emmanuelle Normand, Virginie Vergnes, Juan Lapuente, Amelia Meier, Juan Lapuente, Alexander Tickle, Heather Cohen, Jodie Preece, Amelia Meier, Juan Lapuente, Roman M. Wittig, Dervla Dowd, Sorrel Jones, Sergio Marrocoli, Vera Leinert, Charlotte Coupland, Villard Ebot Egbe, Anthony Agbor, Volker Sommer, Emma Bailey, Andrew Dunn, Inaoyom Imong, Emmanuel Dilambaka, Mattia Bessone, Amelia Meier, Crickette Sanz, David Morgan, Aaron Rundus, Rebecca Chancellor, Felix Mulindahabi, Protais Niyigaba, Chloe Cipoletta, Michael Kaiser, Kyle Yurkiw, Bradley Larson, Alhaji Malikie Siaka, Liliana Pacheco, Manuel Llana, Henk Eshuis, Erin G. Wessling, Mohamed Kambi, Parag Kadam, Alex Piel, Fiona Stewart, Katherine Corogenes, Klaus Zuberbuehler, Kevin Lee, Samuel Angedakin, Kevin E. Langergraber, Christophe Boesch, Hjalmar Kuehl, Mimi Arandjelovic, Paula Dieguez, Mizuki Murai, Yasmin Moebius, Joana Pereira, Silke Atmaca, Kristin Havercamp, Nuria Maldonado, Colleen Stephens;
- (b) funding agencies: Max Planck Society, Max Planck Society Innovation Fund, Heinz L. Krekler Foundation;
- (c) ministries and governmental organizations: Agence Nationale des Parcs Nationaux (Gabon), Centre National de la Recherche Scientifique (CENAREST) (Gabon), Conservation Society of Mbe Mountains (CAMM) (Nigeria), Department of Wildlife and Range Management (Ghana), Direction des Eaux, Forêts et Chasses (Senegal), Eaux et Forêts (Mali), Forestry Commission (Ghana), Forestry Development Authority (Liberia), Institut Congolais pour la Conservation de la Nature (DR-Congo), Instituto da Biodiversidade e das reas Protegidas (IBAP), Makerere University Biological Field Station (MUBFS) (Uganda), Ministère de l'Economie Forestière (R-Congo), Ministère de la Recherche Scientifique et de l'Innovation (Cameroon), Ministère de la Recherche Scientifique (DR-Congo), Ministère de l'Agriculture de l'Élevage et des Eaux et Forêts (Guinea), Ministère de la Recherche Scientifique et Technologique (R-Congo), Ministère des Eaux et Forêts (Cote d'Ivoire), Ministère des Forêts et de la Faune (Cameroon), Ministère de l'Environnement et de l'Assainissement et du Développement Durable du Mali, Ministro da Agricultura e Desenvolvimento Rural (Guinea-Bissau), Ministry of Agriculture, Forestry and Food Security (Sierra Leone), Ministry of Education (Rwanda), National Forestry Au-

- thority (Uganda), National Park Service (Nigeria), National Protected area Authority (Sierra Leone), Rwanda Development Board (Rwanda), Socit Equatoriale d'Exploitation Forestière (SEEF) (Gabon), Tanzania Commission for Science and Technology (Tanzania), Tanzania Wildlife Research Institute (Tanzania), Uganda National Council for Science and Technology (UNCST), (Uganda), Uganda Wildlife Authority (Uganda);
- (d) non-governmental organizations: Budongo Conservation Field Station (Uganda), Ebo Forest Research Station (Cameroon), Fongoli Savanna Chimpanzee Project (Senegal), Foundation Chimbo (Boe), Gashaka Primate Project (Nigeria), Gishwati Chimpanzee Project (Rwanda), Goulougo Triangle Ape Project, Jane Goodall Institute Spain (Dindefelo) (Senegal), Korup Rainforest Conservation Society (Cameroon), Kwame Nkrumah University of Science and Technology (KNUST) (Ghana), Loango Ape Project (Gabon), Lukuru Wildlife Research Foundation (DRC), Ngogo Chimpanzee Project (Uganda), Nyungwe-Kibira Landscape, Rwanda-Burundi (WCS), Projet Grands Singes, La Belgique, Cameroon (KMDA), Station d'Etudes des Gorilles et Chimpanzés (Gabon), Tai Chimpanzee Project (Cote d'Ivoire), The Aspinall Foundation, (Gabon), Ugalla Primate Project (Tanzania), WCS (Conkouati-Douli NP) (R-Congo), WCS Albertine Rift Programme (DRC), Wild Chimpanzee Foundation (Cote d'Ivoire), Wild Chimpanzee Foundation (Guinea), Wild Chimpanzee Foundation (Liberia), Wildlife Conservation Society (WCS) Nigeria (Nigeria), WWF (Campo Maan NP) (Cameroon), WWF Congo Basin (DRC).

References

- [1] IUCN red list of threatened species. <https://www.iucn.org/resources/conservation-tools/iucn-red-list-threatened-species>.
- [2] M. Andriluka, U. Iqbal, E. Ensafutdinov, L. Pishchulin, A. Milan, J. Gall, and Schiele B. PoseTrack: A benchmark for human pose estimation and tracking. *CVPR*, 2018.
- [3] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. *CVPR*, 2014.
- [4] Yauhen Babakhin, Artsiom Sanakoyeu, and Hirotohi Kitamura. Semi-supervised segmentation of salt bodies in seismic images using an ensemble of convolutional neural networks. *German Conference on Pattern Recognition (GCPR)*, 2019.
- [5] Miguel A Bautista, Artsiom Sanakoyeu, Ekaterina Tikhoncheva, and Björn Ommer. Cliques: Deep unsupervised exemplar learning. In *Advances in Neural Information Processing Systems*, pages 3846–3854, 2016.
- [6] Benjamin Biggs, Thomas Roddick, Andrew Fitzgibbon, and Roberto Cipolla. Creatures great and small: Recovering the shape and motion of animals from video. *ACCV*, 2018.
- [7] Jinkun Cao, Hongyu Tang, Fang Hao-Shu, Xiaoyong Shen, Cewu Lu, and Yu-Wing Tai. Cross-domain adaptation for animal pose estimation. *ICCV*, 2019.

- [8] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. *CVPR*, 2017.
- [9] L. Chen, G. Papandreou, F. Schroff, and H. Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. *ECCV*, 2018.
- [10] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- [11] Rıza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. *CVPR*, 2018.
- [12] Semih Günel, Helge Rhodin, Daniel Morales, João H. Campagnolo, Pavan Ramdya, and Pascal Fua. Deepfly3d, a deep learning-based approach for 3d limb and appendage tracking in tethered, adult drosophila. *eLife*, 2019.
- [13] Oshri Halimi, Or Litany, Emanuele Rodola, Alex Bronstein, and Ron Kimmel. Self-supervised learning of dense shape correspondence. *CVPR*, 2019.
- [14] K. He, G. Gkioxari, and P. Dollár and. R. Girshick. Mask R-CNN. *ICCV*, 2017.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [16] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [17] Tomas Jakab, Ankush Gupta, Hakan Bilen, and Andrea Vedaldi. Unsupervised learning of object landmarks through conditional image generation. *NIPS*, 2018.
- [18] Sam Johnson and Mark Everingham. Clustered pose and nonlinear appearance models for human pose estimation. In *BMVC*, 2010.
- [19] Sam Johnson and Mark Everingham. Learning effective human pose estimation from inaccurate annotation. *CVPR*, 2011.
- [20] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. *CVPR*, 2018.
- [21] Angjoo Kanazawa, Jason Y. Zhang, Panna Felsen, and Jitendra Malik. Learning 3d human dynamics from video. *CVPR*, 2019.
- [22] Pierre Karashchuk. lambdaloop/anipose: v0.5.0. *eLife*, 2019.
- [23] A. Kendall and Y. Gal. What uncertainties do we need in bayesian deep learning for computer vision? *NIPS*, 2017.
- [24] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. *CVPR*, pages 6399–6408, 2019.
- [25] V. Leon, N. Bonneel, G. Lavoue, and J.-P. Vandeborre. Continuous semantic description of 3d meshes. *Computer & Graphics*, 2016.
- [26] Shuyuan Li, Jianguo Li, Weiyao Lin, and Hanlin Tang. Amur tiger re-identification in the wild. *arXiv preprint arXiv:1906.05586*, 2019.
- [27] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. *ECCV*, 2014.
- [28] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black and. SMPL: A skinned multi- person linear model. *ACM Trans. on Graphics*, 2015.
- [29] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *TOG*, 2015.
- [30] Dominik Lorenz, Leonard Bereska, Timo Milbich, and Björn Ommer. Unsupervised part-based disentangling of object shape and appearance. *CVPR*, 2019.
- [31] Alexander Mathis, Pranav Mamidanna, Kevin M. Cury, Abe Taiga, Venkatesh N. Murthy, Mackenzie Weygandt Mathis, and Matthias Bethge. Deeplabcut: markerless pose estimation of user-defined body parts with deep learning. *Nature neuroscience*, 2018.
- [32] Mackenzie Weygandt Mathis and Alexander Mathis. Deep learning tools for the measurement of animal behavior in neuroscience. *arXiv preprint arXiv:1909.13868v2*, 2019.
- [33] Tanmay Nath, Alexander Mathis, An Chi Chen, Amir Patel, and Mackenzie W. Bethge, Matthias andd Mathis. Using deeplab-cut for 3d markerless pose estimation across species and behaviors. *Nature protocols*, 2019.
- [34] Lukas Neumann, Andrew Zisserman, and Andrea Vedaldi. Efficient confidence auto-calibration for safe pedestrian detection. *NIPS Workshop on Machine Learning for Intelligent Transportation Systems*, 2018.
- [35] Natalia Neverova, James Thewlis, Rıza Alp Güler, Iasonas Kokkinos, and Andrea Vedaldi. Slim DensePose: Thrifty learning from sparse annotations and motion cues. *CVPR*, 2019.
- [36] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hour-glass networks for human pose estimation. *ECCV*, 2016.
- [37] David Novotny, Nikhila Ravi, Benjamin Graham, Natalia Neverova, and Andrea Vedaldi. C3DPO: Canonical 3d pose networks for non-rigid structure from motion. *ICCV*, 2019.
- [38] Ilija Radosavovic, Piotr Dollar, Ross Girshick, Georgia Gkioxari, and Kaiming He. Data distillation: Towards omniscient supervised learning. *CVPR*, 2018.
- [39] Maheen Rashid, Xiuye Gu, and Yong Jae Lee. Interspecies knowledge transfer for facial keypoint detection. *CVPR*, 2017.
- [40] S. Salti, F. Tombari, and L. Di Stefano. Shot: Unique signatures of histograms for surface and texture description. *Computer Vision and Image Understanding*, 2014.
- [41] Artsiom Sanakoyeu, Miguel A Bautista, and Björn Ommer. Deep unsupervised learning of visual similarities. *Pattern Recognition*, 78:331–343, 2018.
- [42] James Thewlis, Samuel Albanie, Hakan Bilen, and Andrea Vedaldi. Unsupervised learning of landmarks by descriptor vector exchange. *ICCV*, 2019.
- [43] J. Thewlis, H. Bilen, and A. Vedaldi. Unsupervised learning of object landmarks by factorized spatial embeddings. *ICCV*, 2017.
- [44] J. Thewlis, H. Bilen, and A. Vedaldi. Unsupervised object learning from dense invariant image labelling. *NIPS*, 2017.
- [45] X. Wang, R. Girshick, A. Gupta, and K. He. Non-local neural networks. *CVPR*, 2018.
- [46] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. *CVPR*, 2016.
- [47] I. Zeki Yalniz, Hervé Jegou, Kan Chen, Manohar Paluri, and Dhruv Mahajan. Billion-scale semi-supervised learning for image classification. *arXiv preprint arXiv:1905.00546v1*,

2019.

- [48] Heng Yang, Renqiao Zhang, and Peter Robinson. Human and sheep facial landmarks localisation by triplet interpolated features. *WACV*, 2015.
- [49] Lu Yang, Qing Song, Zhihui Wang, and Ming Jiang. Parsing r-cnn for instance-level human analysis. *CVPR*, 2018.
- [50] Weiyu Zhang, Menglong Zhu, and Konstantinos G Derpanis. From actemes to action: A strongly-supervised representation for detailed action understanding. *ICCV*, 2013.
- [51] Yuting Zhang, Yijie Guo, Yixin Jin, Yijun Luo, Zhiyuan He, and Honglak Lee. Unsupervised discovery of object landmarks as structural representations. *CVPR*, 2018.
- [52] Silvia Zuffi, Angjoo Kanazawa, Tanya Berger-Wolf, and Michael J. Black. Three-d safari: Learning to estimate zebra pose, shape, and texture from images "in the wild". *ICCV*, 2019.
- [53] Silvia Zuffi, Angjoo Kanazawa, and Michael J. Black. Lions and tigers and bears: Capturing non-rigid, 3d, articulated shape from images. *ICCV*, 2018.
- [54] Silvia Zuffi, Angjoo Kanazawa, David W. Jacobs, and Michael J. Black. 3d menagerie: Modeling the 3d shape and pose of animals. *CVPR*, 2017.

Appendix

In Section A we provide more details on our implementation of the Multi-head R-CNN network. Then, in Section C we describe additional ablation studies on the advantages of the auto-calibrated training, as well as other architectural choices. Finally, Section D refers the reader to the qualitative results obtained on videos from the Chimp&See dataset.

A. Architecture

We introduced a number of changes and improvements in the DensePose head of the standard DensePose R-CNN architecture of [11] with ResNet-50 [15] backbone. These changes are listed below for the affected branches; other branches remained unchanged and correspond exactly to the Mask R-CNN architecture of [14].

- We have increased the RoI resolution from 14×14 to 28×28 in the DensePose head, as proposed in [49].
- We have replaced the 8-layer DensePose head with the geometric and context encoding (GCE) module [49], combining a non-local convolutional layer [45] with the atrous spatial pyramid pooling (ASPP) [9].
- We have replaced the original FPN of DensePose R-CNN with a Panoptic FPN [24].

Each of these modifications led to increase in network performance due to improved multi-scale context aggregation. We refer the reader to the work of [49] for ablation studies whose results are aligned well with our own observations.

To predict or we simply extend the output layer of the corresponding head by doubling the number of its neurons.

Our codebase, network configuration files for each experiment and pretrained models will be publicly released.

B. Computational cost

Our auto-calibrated model has a negligible computational overhead ($< 1\%$) compared to the baseline model. Before training the *student*, sampling of the pseudo-labels requires one forward pass of the *teacher* network over the unlabeled dataset. The *teacher* and the *student* networks share the same architecture.

C. Ablation studies

First, we report performance of the original Mask R-CNN [14] framework, as well as our auto-calibrated version of the same architecture, on detection and segmentation tasks (see Tab. 6). Training in the auto-calibration setting resulted in minor gains on the COCO dataset that the model was trained on, but, as expected, led to major improvements in performance on the out-of-distribution data (DensePose-Chimps and Chimp&See).

Second, Tab. 7 shows results of replacing the proposed binary foreground-background segmentation in the DensePose head (a) with 15-way coarse body part segmentation as in the original DensePose-RCNN framework [11] (b). We can see that binary segmentation generalizes better than the 15-way. We have also experimented with using the binary mask from the Mask R-CNN head instead of mask produced by the DensePose head (Tab. 7 (c)) *during inference step*. Moreover, even though exploiting the mask from the separate mask head at test time results in better performance, complete removal of the mask from the DensePose head leads to under-training and decreased accuracy of estimation of uv -coordinates (since in this case the DensePose head receives only sparse supervisory signals at the annotated locations).

D. Qualitative results

In addition, we also point the readers to the attached video samples from the Chimp&See dataset showing frame-by-frame predictions produced by our model before (*teacher*) and after self-training (*student*). The results produced by the *student* network are generally significantly more stable.

model	COCO minival		DensePose-Chimps		Chimp&See	
	AP_D	AP_S	AP_D	AP_S	AP_D	AP_S
Mask RCNN	40.98	37.17	48.3	44.92	40.56	33.91
σ -Mask RCNN	41.12 (+0.14)	37.09 (-0.08)	52.05 (+3.75)	47.94 (+3.02)	42.9 (+2.34)	34.74 (+0.82)

Table 6: Auto-calibrated Mask R-CNN [14]: detection, instance segmentation on COCO minival (all classes).

model	Mask in DensePose head	AP	AP_{50}	AP_{75}
a) DensePose-RCNN* (σ)	binary	53.20	88.27	56.98
b) DensePose-RCNN* (σ)	15-way	50.87	86.91	54.49
c) DensePose-RCNN* (σ) + mask from the mask head	binary	54.35	88.58	60.28

Table 7: Ablation study of the mask in the DensePose head. Reports the DensePose performance on DensePose-COCO minival. a) our proposed architecture; b) replace the binary segmentation of the DensePose head with 15-way coarse body part segmentation as in the original DensePose-RCNN framework [11]; c) use the binary mask from the DensePose head during training, but substitute it with the mask from the separate mask head during inference.