

Improving 2D Human Pose Estimation in Rare Camera Views with Synthetic Data

Supplementary Material

6. Pose spaces analysis

6.1. Definitions

The spaces associated with pose estimation can be conceptualized through the popular SMPL model [24]. The fixed kinematic chain allows us to describe each pose by the rotations of its joints. Every pose inherent to the SMPL model can be uniquely denoted as a vector positioned in a 63-dimensional space (21 joints, 3 axes of rotations each), each dimension bounded in the range of $(-\pi, \pi)$. We call this space P^{all} .

Definition 6.1 (P^{all}). Space of all poses possible with a given kinematic chain defined by the SMPL model and joints ranges between $(-\pi, \pi)$.

Once the P^{all} is defined, we can concretize other spaces to analyze prior synthetic data generation techniques.

Definition 6.2 ($P^{bounded}$). Space of all poses possible with a given kinematic chain defined by the SMPL model with joints ranges between (α_j, β_j) for each joint j . Angles of joints are independent.

Definition 6.3 ($P^{anatomical}$). Space of all poses that at least one human can achieve.

Definition 6.4 (P^{AMASS}). Space of all poses captured in the AMASS dataset.

All defined spaces along with examples are visualized in the Fig. 6.

In the space P^{all} , we use Euclidean distance to measure the similarity of the two poses. Using the Euclidean distance is not trivial; we justify it in two steps.

(1) With the given kinematic chain of the SMPL model, the two poses differ only by angles of joints. In the rendered image, a camera viewpoint is crucial for human-perceived pose similarity (and also 2D pose similarity). To discard the influence of the camera, we treat each element of the pose vector independently with the same weight.

(2) Euclidean distance would not work for angles on the border between $-\pi$ and π where Euclidean distance is biggest while the angle distance is low. We theorize that most human joints have a range of movement smaller than π , and the mentioned overflow of the distance will not happen between poses close to the human anatomy. It is worth noting that the Euclidean distance would fail for two unbounded poses.

Using Euclidean distance allows us to use off-the-shelf algorithms like k-means clustering or dimensionality reduction techniques for further analysis.

6.2. Synthetic data generation through the pose spaces

Previous works like AMASS [26] focus on generating anatomically plausible data from the space $P^{anatomical}$. They do so by fitting the SMPL model into 3D scans of humans. The approach has two challenges.

First, an error can occur during the SMPL model fitting, and the resulting pose would be outside of the $P^{anatomical}$. An example of such an error is the Fig. 6 B) where the left elbow is rotated by more than 45° along the axis where no rotation is anatomically possible.

The second and bigger challenge is that the approach requires many expensive 3D scans of people in various poses. Due to the nature of the 3D scanning techniques, most poses are not dynamic (jumping, running, etc.), and the researcher in the lab must design rare poses. There is a high probability that a lot of the $P^{anatomical}$ space remains unsampled.

On the other hand, sampling poses from the $P^{bounded}$ space has its challenges. Foremost, how to design the bounds and sample the space to balance rare poses (as is the case of Fig. 6 F) with real-world ones. The bounds should not be too tight to allow sampling of very rare poses but also not too loose to generate a lot of anatomically impossible poses (like the pose D) in Fig. 6). The optimal bounds would cause spaces $P^{anatomical}$ and $P^{bounded}$ to overlap almost completely. There will always be poses impossible to capture by the $P^{bounded}$ - for example, contortionist visualized in Fig. 6 C).

6.3. RePoGen vs. AMASS

When comparing the RePoGen-generated data with the AMASS dataset through the spaces introduced in Sec. 6.1, we can measure their similarity (using the Euclidean distance) and the number of rare poses.

To measure the ratio of rare poses as in the [14], we clustered AMASS poses using the k-means. Poses further than the threshold from all cluster centers are classified as rare. Setting up the threshold such that the AMASS dataset has 5% of rare poses results in over 90% of rare poses in the RePoGen data. RePoGen data are, therefore, different from AMASS data by a big margin. RePoGen data are rare and usually completely new instead of weighting rare poses as in [14]. To justify using the RePoGen data instead of

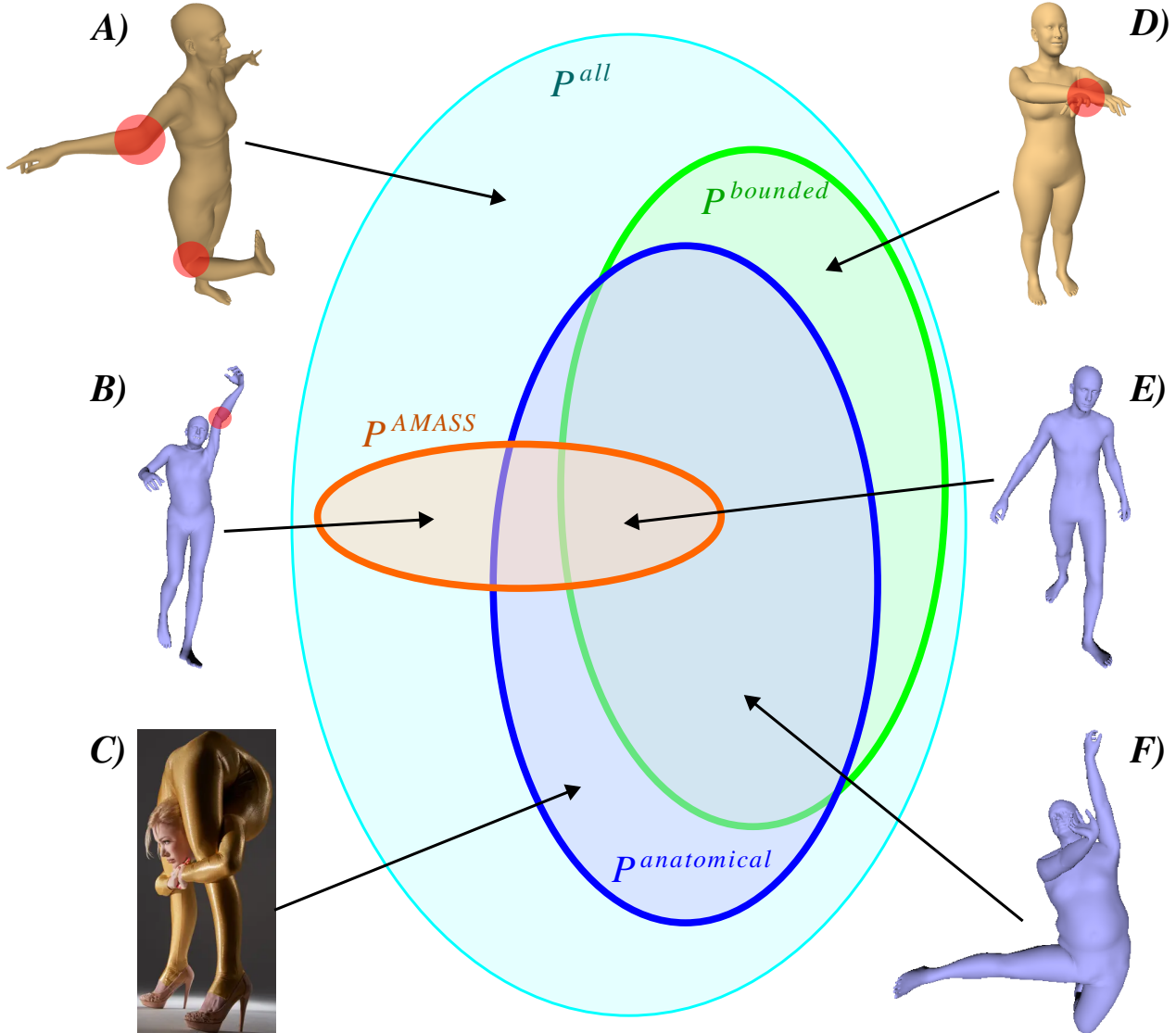


Figure 6. Visualization of described pose spaces and examples of poses from them. Pose spaces are described in detail in Sec. 6.1. The examples are A) impossible pose, B) pose from AMASS dataset with impossible joint rotation, C) anatomical but not bounded pose, D) bounded pose with self-intersection, E) AMASS pose with low pose variance, F) RePoGen pose not present in the AMASS dataset. Impossible rotations and self-intersections emphasized by red circles.

AMASS, the Fig. 7 shows the most similar AMASS pose to selected RePoGen poses. We can see that while RePoGen generates poses similar to standard ones, it also generates ones not in the AMASS dataset.

7. RePo Dataset

Here we describe a new RePo dataset of manually annotated real images. The dataset focuses on extreme poses in top and bottom views typically encountered in sports. Images come primarily from public sports videos on YouTube. The dataset is split into two parts - one focusing on the bottom view with 187 images, the other focusing on the top view

with 91 images. Each part is divided into sets described in Tab. 1 in the paper.

One professional annotator annotated the whole dataset. We created a custom annotation environment allowing for an easier understanding of the scene necessary for annotating extreme views. Since the visibility of the joints is not defined in detail in the COCO dataset [20], we defined it as follows:

Visibility 0. The keypoint is not visible, and we cannot reliably tell its precise location. The words *reliably* and *precise* are crucial in situations where a keypoint is not visible. We can guess its location from the context but cannot be

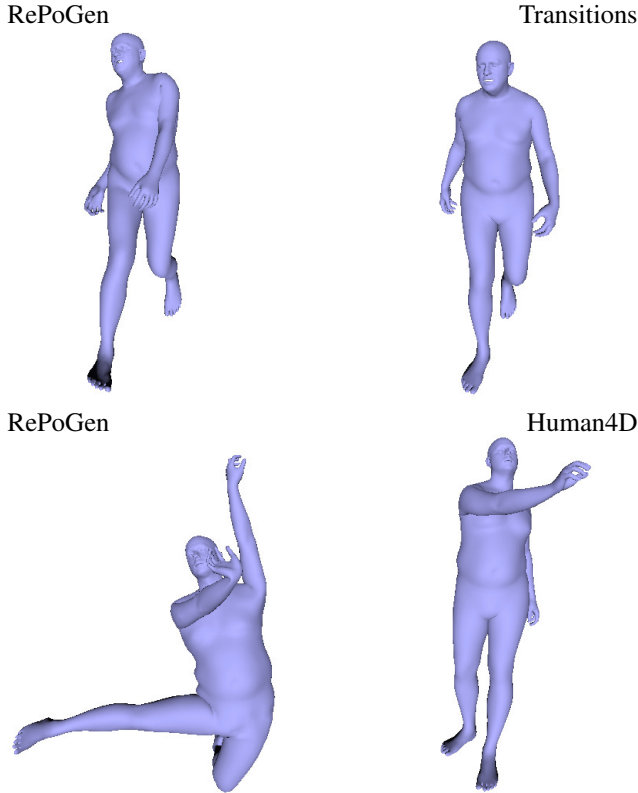


Figure 7. Selected RePoGen poses and their most similar AMASS counterpart. The first row shows two similar poses generated by the RePoGen and AMASS pipelines, while the second row proves the RePoGen pipeline generates poses unseen in AMASS. Each AMASS pose is labeled by its respective subset.

sure the guess was correct.

Visibility 1. The keypoint is not visible, but we can reliably tell its precise location from the context and other keypoints.

Visibility 2. The keypoint is visible in the image.

Further, as the extreme views pose additional challenges, we stick to the annotation of joint projection to the image plane if it is unambiguous. A typical example would be the ankle which is rarely visible from the bottom view (we see the heel instead). Without this relaxation of definition, almost no keypoints would be annotated as visible (visibility 2).

Examples of images from all sets of the dataset are in the Fig. 11, Fig. 12 and Fig. 13.

8. RePoGen Dataset

The RePoGen dataset was created with the proposed RePoGen method and was used to train the best-performing model. There are 3 variants of the RePoGen dataset, all meeting the description in Tab. 6. The parameter distinguishing the 3 variants is the camera viewpoint distribution.

parameter	value
number of poses	1 500
number of views per pose	2
pose variance	1.0
with texture	✓
with background	✓
default SMP pose	✗
distribution	skewed Gaussian

Table 6. Parameters used for RePoGen dataset generation.

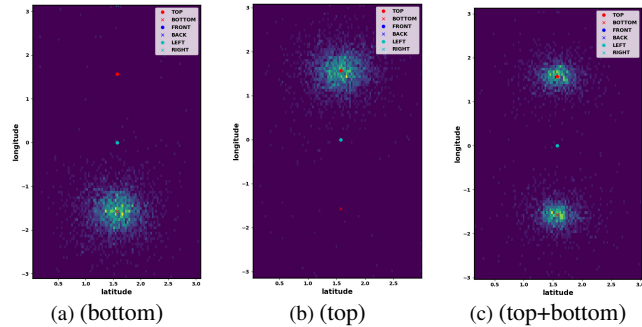


Figure 8. Camera viewpoint distributions for various RePoGen datasets. Vertical axis - latitude, horizontal - longitude.

The RePoGen (bottom) and RePoGen (top) are sampled with a normal distribution centered around the bottom view and top view respectively. The RePoGen (top+bottom) is sampled from a combination of these two distributions. Visualization is in the Fig. 8, where 3D coordinates are projected to the latitude and longitude.

The Figs. 9 and 10 contain images from the RePoGen dataset.

9. Additional results

We also offer additional qualitative results on the proposed RePo dataset. The RePo Bottom Test is in the Fig. 11, and Val set in Fig. 12. The model fine-tuned with RePoGen data struggles the most with head keypoints and strong motion blur. The Fig. 13 shows results on the Bottom Seq set where we show performance on a video. We show every third frame from a video.

10. Code

The code is available in the supplementary material. The repository builds on the SMPL-X project [28] and uses the same dependencies. For more details, see the enclosed README.md file and the code itself.



Figure 9. Images from the RePoGen (bottom) dataset.

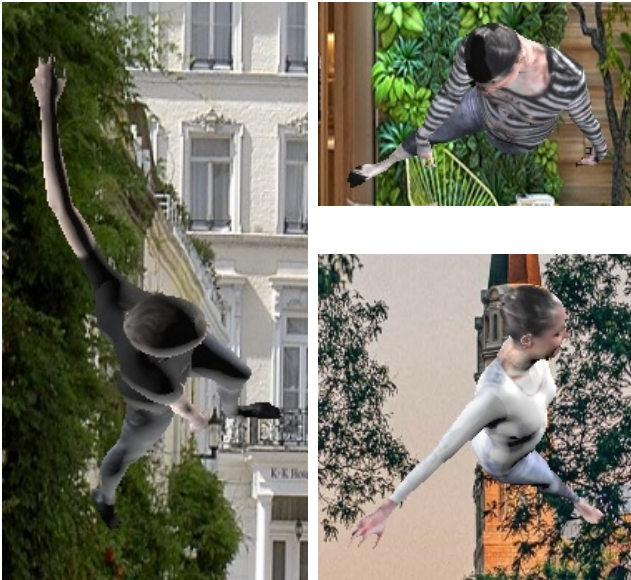


Figure 10. Images from the RePoGen (top) dataset.



Figure 11. Examples from the RePo bottom test set. ViTPose estimates when trained on COCO (left) and on RePoGen data (right). Colors – right hand, right leg, left hand and left leg

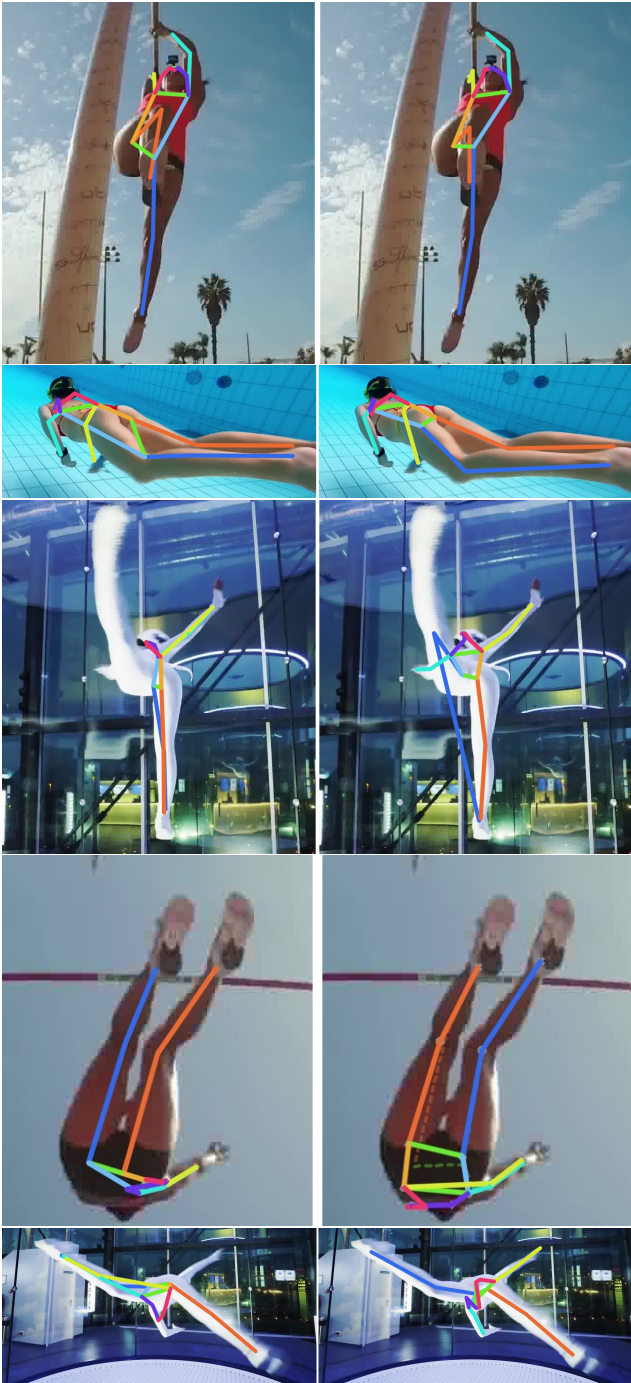


Figure 12. Examples from the RePo bottom val set. ViTPose-s estimates when trained on COCO (left) and on RePoGen data (right). Colors – right hand, right leg, left hand and left leg



Figure 13. Examples from the RePo bottom seq set. ViTPose-s estimates when trained on COCO (left) and on RePoGen data (right). Images from a consecutive sequence, taking every third frame. Colors – right hand, right leg, left hand and left leg