

Human-shape classification with artificial backgrounds

Armin Niedermueller¹, Melanie Urban¹ and Ahmet Cihat Bozkurt¹

Abstract— Human Body Dimensions (HBD) Estimation using images becomes increasingly important. Datasets consisting of 2d images of humanoids have been successfully generated and then used for training such networks. However, those images all have a single coloured grey background. The HBD estimation of such images generates very good results, but what if an artificial background is introduced into the dataset and the humans are visually not easily separable? In this work, we introduce random coloured backgrounds as well as textured backgrounds into an existing dataset of humanoids. We create and train a modified network architecture and then calculated the HBDs. Our modified network achieved far better results (~ 50% improvement) as in the original paper, even with the texture and background dataset.

I. INTRODUCTION

Calculating human dimensions from 3D human meshes is possible, but that is another research topic. In this study, we tried to develop a model that predicts body measurements from images of humans and monsters in different poses and related body measurements. Therefore we used the architecture of the *Neural Anthropometer* (NA) as basis and also the provided dataset [7].

Since the model trained with uniform background images will have a narrow usage area, we added artificial backgrounds to the images in order to train a model that can make accurate predictions in different scenarios. This setting is much more realistic and can be used in a variety of application areas. For example an online shop could provide a tool for uploading an image by oneself and afterwards announce the dress size.

The background we add are either a random color or texture images. We tested our model with plain images of humans or monsters and respectively images with RGB background or with texture background. We evaluate the results by calculating the average model accuracy, MAD (Mean Absolute Deviation) for each body dimensions and RPE (Relative Percentage Error) for each body dimensions.

In this paper will first provide an overview of related work and the used dataset. Then we introduce our model and show the differences to the original Neural Anthropometer. Afterwards, we evaluate and discuss the results of testing our model.

II. HUMAN BODY DIMENSION ESTIMATION

There are plenty ways of estimating human body dimensions and just as many different input types. We orientated our model on the descriptions made in the paper *A Neural Anthropometer Learning from Body Dimensions Computed on Human 3D Meshes* [7]. They used images of 3D meshes synthesized with the Skinned Multi-Person Linear Model (SMPL) [5] together with the calculated HBDs as ground truth.

But our NA is modified and outperforms the original architecture by far. Our code is publicly available at https://github.com/nerovalerius/humanoids_cnn and the code of the original NA and also the dataset at: <https://github.com/neoglez/neural-anthropometer>

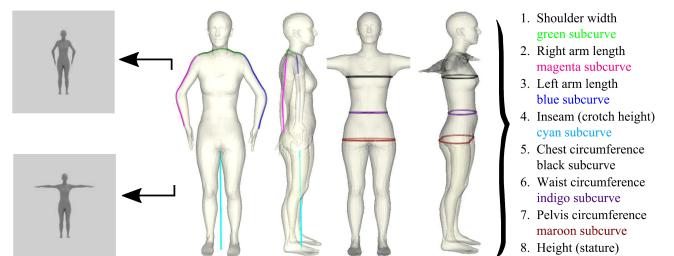


Fig. 1. Method of generating 2D images from 3D meshes which was used in the original paper [7]. On the right side, the HBDs are listed. While the humanoids are stored as 2D grey scale images, the HBD are provided as json data.

A. Dataset

The used dataset consists of 12.000 grey-scale 2D images of humans and 2.000 additionally provided 2D images of monsters, with a size of 200 x 200 each.

The underlying 3D meshes were formed by providing shape parameters to the SMPL. To get bodies with a huge diversity of measurements the parameters were uniformly varied. [7] But there is also the possibility to end up with so-called monster shapes. We also tested these monsters, however, separately from the humans.

Two different poses are used, pose 0 and pose 1 (see Fig. 2). Pose 0 shows the humanoid with the arms stretched out to both sides and pose 1 shows the hands hanging down. Furthermore, the dataset consists of females and males. We have 3500 females in pose 0, 3500 females in pose 1, 3500 males in pose 0, 3500 males in pose 1.

¹Pattern Recognition II, Paris Lodron Universität Salzburg, armin.niedermueller@mailbox.org
melanie.urban@stud.sbg.ac.at
s1079921@stud.sbg.ac.at

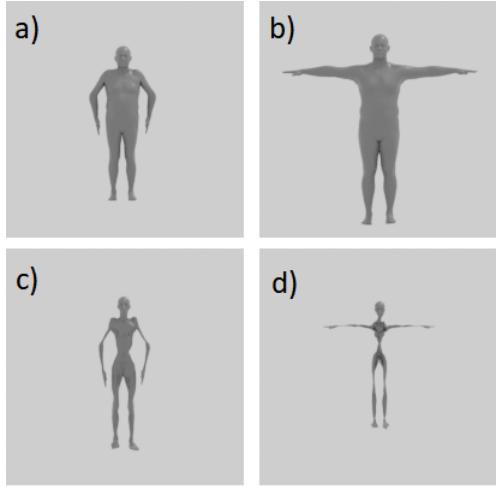


Fig. 2. Different Humanoids from the original dataset [7]. a) shows a human male in pose 0. b) shows a human male in pose 1. c) shows a monster male in pose 0. d) shows a monster female in pose 1.

Furthermore, the dataset provides the correct annotations for the HBDs of each image in order to let our network be able to learn. The eight HBDs are as follows:

Shoulder width	Right arm length
Left arm length	Inseam (crotch height)
Chest circumference	Waist circumference
Pelvis circumference	Height (stature)

B. Artificial Backgrounds

We augment the original gray scale dataset by adding two types of backgrounds and thus generating two new datasets. For the first new dataset, we replaced the single color background by inserting a background where all pixels have the same random RGB value. At the second dataset, we replace the background with textures of the DTD-dataset, which consists of 5640 different images [3]. Since our dataset consists of 14.000 images, the DTD-images will be used multiple times. However, this should not affect the results in any way, when a DTD-image is used twice, it is attached to a different humanoid. A subset of the resulting images are shown in Fig. 3.

At the end of the process, we have 3 different datasets: human/monster images with the original gray background, a second data set with grey texture backgrounds and the third dataset with random grey scale backgrounds.

III. APPROACH

A. Pre-Processing

To keep the overall CNN architecture simple we restrict the input channels to 1. Since we introduced RGB backgrounds to the datasets in Section II-B, we converted them back to grey scale before training using OpenCV [2].



Fig. 3. Artificial Backgrounds which are introduced into the original dataset. Row 1 and 2 show texture backgrounds from the DTD-dataset

B. Network Architecture

Our neural network is shown in Fig. 4 and implemented with Pytorch [6]. We used the original architecture as in [7] which is described as follows:

The input layer processes 200 x 200 x 1 images. A 2D-convolution with a 5-pixels square kernel produces a feature map of size 196 x 196 x 8, which is afterwards passed through a rectified linear unit (ReLU) [1] and then is batch normalized. After a max pooling with stride 2, a second 2D-convolution is applied (again squared kernel of size 5), which results in a tensor of size 94 x 94 x 16. After another max pooling, the output is flattened to a tensor of size 35344. At last, this tensor is passed to a fully connected layer and through a ReLU. The last layer is a regressor that outputs the eight human body dimensions in meters [7, p. 6].

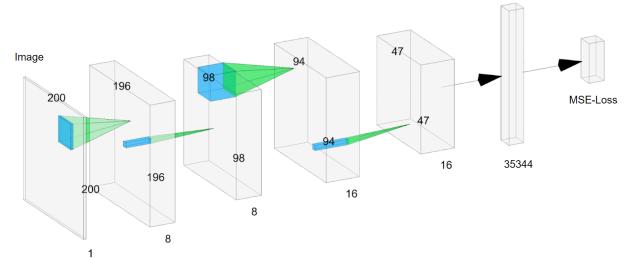


Fig. 4. Our CNN architecture based on the original NA framework [7]

Once we implemented our network and selected hyperparameters providing the best results, we also compared our model to the meanwhile published original Neural Anthropometer. There are several differences which lead to better results:

LayerNorm

The original NA uses a BatchNorm2d-Layer [4] right after the output of the first convolution is passed through the ReLU-Layer. We instead normalize with LayerNorm [9]. As shown in Fig. 5 the BatchNorm normalizes the

tensor across the batch-size and the spatial dimensions for each channel while the LayerNorm uses for normalization the values computed across all channels for each sample. Because at this stage of the NA we already have eight different channels it might be useful to normalize per channel. But these channels were derived from a single input image without channels at all in the first place. If the eight derived channels do not include all the informations needed to compute accurate statistics for the final eight output dimensions a layer normalization might handle this.

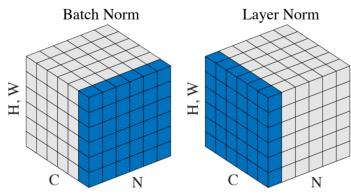


Fig. 5. Different normalization methods for a tensor with batch axis N, channel axis C, the height H and the width W. The blue pixels were normalized by calculating the needing values from these pixels. [8]

Pooling Layer

Both NAs use a MaxPool2D-Layer [6] twice. But instead of a kernel size of 2 we reduced it to size 1. This means the amount of pixels is still shrunk to the half but without max pooling.

Fully Connected Layer

After the second max pooling the output is flattened to size 1x35344 and passed to a fully connected layer. This might now maintain this dimension or already shrink the output. Shrinking the output also significantly decrease the size of the model. After cross-validating different output sizes we fixed it to 512. The original NA uses an output size of 84. It seems that the fast reduction of size before passing the tensor to the last ReLu leads to too much information loss.

IV. EXPERIMENTS AND RESULTS

We take the whole dataset (See II-A) once for humans and once for monsters and use a k-fold cross-validation ($k = 5$) for each separately. We randomly select images from the dataset and train our network for 20 epochs with a mini-batch size of 100 and a learning rate of $\eta = 0.01$. **MSE** between the actual and predicted HBDs is minimized during training. Furthermore we use stochastic-gradient-descent with momentum term ($\alpha = 0.9$) as optimization method. Those hyperparameters are identical to those provided in the paper of the original NA [7].

The experiment is performed on a Computer with an AMD Ryzen 9 5950X CPU and a NVIDIA RTX 3090 with 24 GB of GPU Memory.

A. Quantitative Evaluation

We evaluate our model with a k-fold cross validation ($k = 5$). Each fold j consists of $a = 2400$ images for the humans

TABLE I
COMPARISON ORIGINAL NA AND OUR NETWORK

HBD	MAD			RPE (%)		
	Org.	Ours	Diff. (%)	Org.	Ours	Diff. (%)
Shoulder width	12.54	7.33	- 0.42	4.93	3.41	- 0.31
Right arm length	12.98	7.33	- 0.44	2.22	1.26	- 0.43
Left arm length	13.48	7.64	- 0.43	2.34	1.33	- 0.43
Inseam/crotch h.	22.17	12.20	- 0.45	3.12	1.78	- 0.43
Chest circumf.	25.22	10.95	- 0.57	2.51	1.09	- 0.57
Waist circumf.	27.53	10.58	- 0.62	3.67	1.26	- 0.66
Pelvis circumf.	25.85	8.88	- 0.66	2.40	0.84	- 0.65
Height	27.34	6.70	- 0.75	1.58	0.39	- 0.75
AMAD	20.89	8.95	- 0.57			
ARPE				2.84	1.42	- 0.50

TABLE II
RESULTS HUMANS

HBD	Plain		RGB		Texture	
	MAD	RPE (%)	MAD	RPE (%)	MAD	RPE (%)
Shoulder width	7.33	3.41	9.59	3.89	17.02	5.85
Right arm length	7.33	1.26	9.37	1.62	17.90	3.10
Left arm length	7.64	1.33	9.77	1.71	17.89	3.14
Inseam/crotch height	12.20	1.78	15.27	2.22	29.99	4.21
Chest circumference	10.95	1.09	15.25	1.51	36.30	3.61
Waist circumference	10.58	1.26	15.31	1.80	38.46	4.65
Pelvis circumference	8.88	0.84	11.96	1.14	31.01	3.01
Height	6.70	0.39	10.28	0.60	31.78	1.86
AMAD	8.95		12.10		27.54	
ARPE (%)		1.42		1.81		3.68

dataset and $a = 400$ for the monsters dataset. We estimate the eight HBDs using our altered model and use the same metrics as in the original paper [7]:

Statistics (here for the humans dataset) are based on a results tensor of shape:

$$k \times a \times |\hat{D}_i, D_i| \times 8 = 5 \times 2400 \times 2 \times 8 \quad (1)$$

The estimation error e_{MAD}^j for each HBD i is the Mean Absolute Difference over the j folds between the predicted and actual HBDs \hat{D}_i, D_i :

$$e_{MAD}^j = \frac{1}{a} \sum_{l=1}^a |\hat{D}_l - D_l|, \quad e_{MAD}^j = \frac{1}{k} \sum_{j=1}^k e_{MAD}^j \quad (2)$$

Furthermore, we consider the Relative Percentage Error (RPE) e_{RPE}^j for each HBD i and its average (ARPE).

$$e_{RPE}^j = \frac{1}{a} \sum_{l=1}^a \left| \frac{\hat{D}_l - D_l}{D_l} \right|, \quad e_{RPE}^j = \frac{1}{k} \sum_{j=1}^k e_{RPE}^j \quad (3)$$

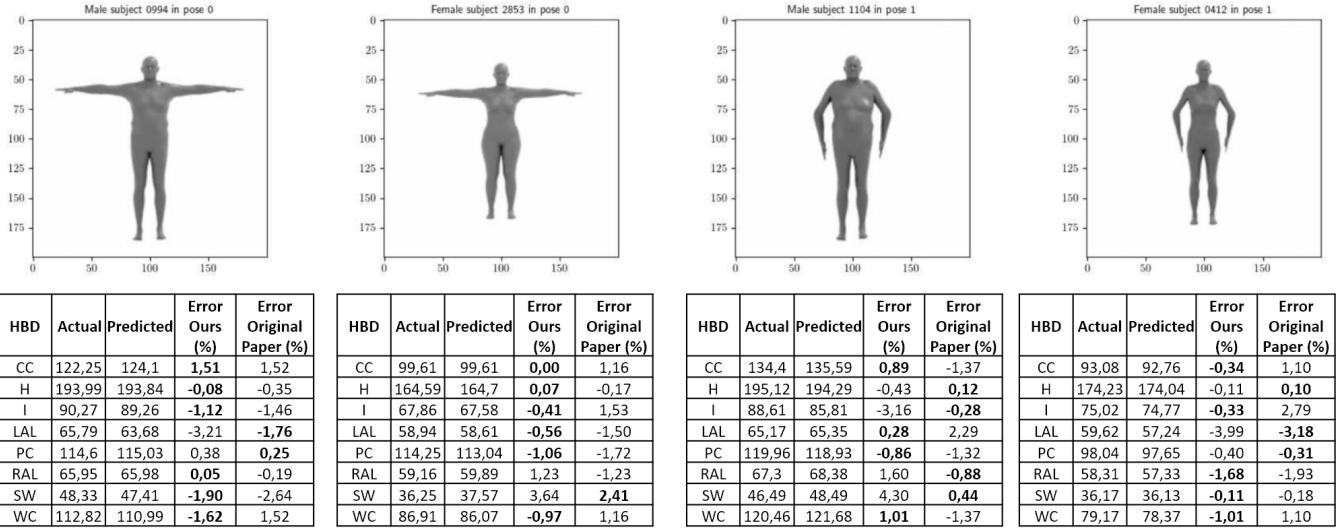


Fig. 6. Estimation results of the same four subjects which are also selected in the original paper. The error of our net and the original net is also given. Bold numbers indicate a better value compared to the Two females and two males. Abbreviations: CC: chest circumference, H: height, I: inseam, LAL: left arm length, PC: pelvis circumference, RAL: right arm length, SW: shoulder width, WC: waist circumference.

TABLE III
RESULTS MONSTERS

HBD	Plain		RGB		Texture	
	MAD	RPE (%)	MAD	RPE (%)	MAD	RPE (%)
Shoulder width	22.50	7.10	26.10	8.23	32.19	9.93
Right arm length	32.87	12.92	37.15	13.43	42.21	13.01
Left arm length	37.75	16.73	42.59	13.95	48.27	21.21
Inseam/crotch height	30.99	inf	37.32	inf	61.72	inf
Chest circumference	36.96	3.72	46.60	4.59	68.96	6.75
Waist circumference	31.67	4.43	41.76	5.94	72.68	10.74
Pelvis circumference	25.80	2.59	35.42	3.59	60.92	6.38
Height	21.31	1.26	29.21	1.74	64.75	3.87
AMAD	29.98		37.02		56.46	
ARPE (%)		inf		inf		inf

We report the Average Mean Absolute Difference (AMAD) over the HBDs:

$$e_{AMAD} = \frac{1}{8} \sum_{j=1}^8 e_{AMAD}^j \quad (4)$$

For having a baseline and compare our modified NA to the original NA we first evaluated the human dataset with plain background. The results are shown in Table I. It shows that the AMAD declines from 20.89 to 8.95 millimeters using our NA, an improvement of 57%. Furthermore, ARPE improves by 50%, showing a fall from 2.82 to 1.42 mm.

The tables II and III show the MAD and RPE for each HBD and in average for the humans dataset and the monsters dataset respectively. They also contain all three different setting: plain, RGB or texture background.

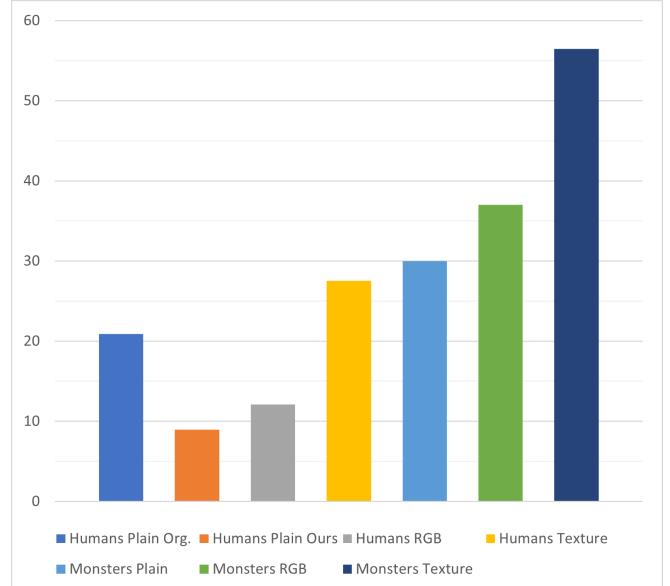


Fig. 7. AMAD on the test set after 20 epochs of training

Also for the humans dataset the AMAD for the images with RGB background is about 25% higher than that of the ones with plain background it is still better than the results of the original NA. The same applies to the ARPE. The ARPE for the texture background is twice at high than the one achieved with RGB background. Shoulder width, left and right arm length can be estimated always quite well while the pelvis circumference and height also was estimated considerable well at the images with plain and RGB background.

The monsters dataset instead was indeed difficult for the

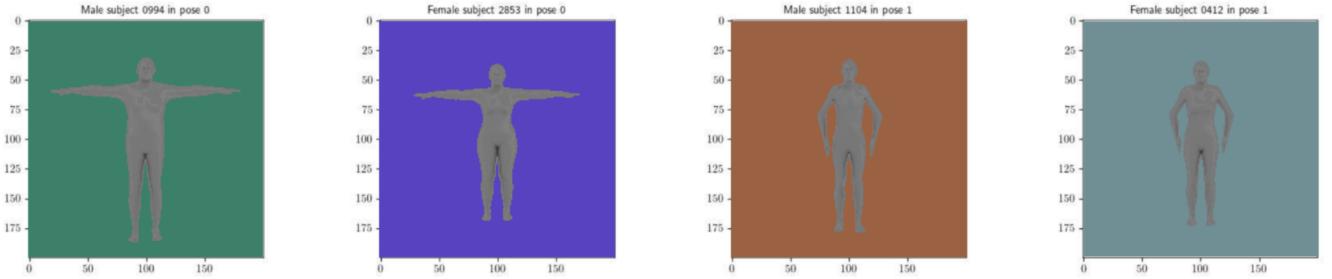


Fig. 8. Our input mini batch now with an added RGB background. Via preprocessing the image is converted from RGB to grey, which results in images with different grey values.

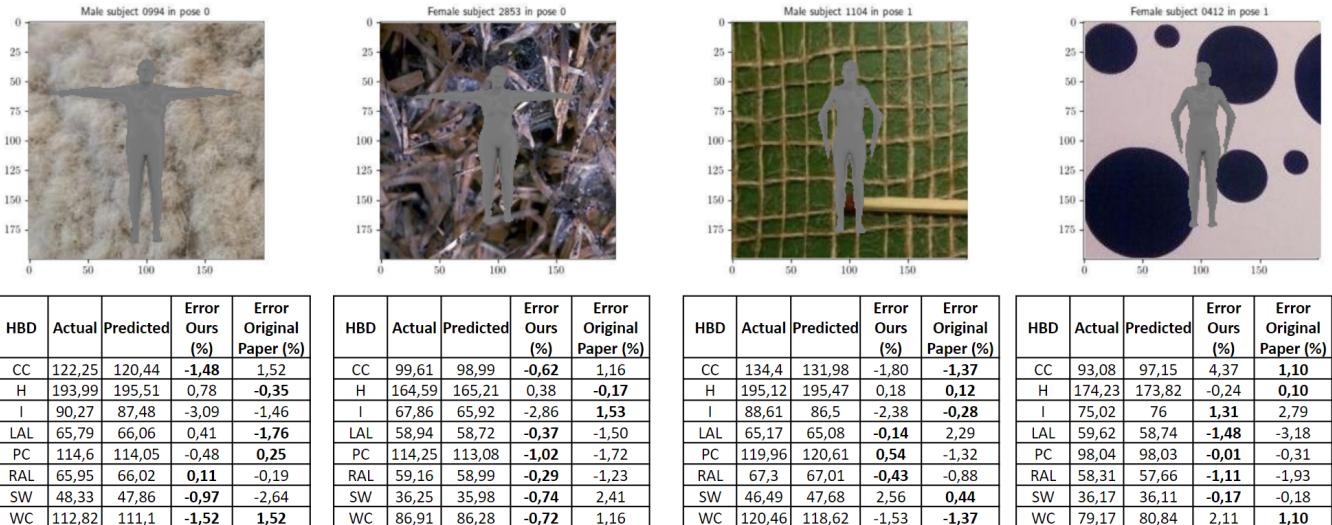


Fig. 9. The mini batch shown with the added DTD backgrounds.

modified NA. Also the RPE for the inseam could not be calculated at all. That might be because of a hugh difference between the real and the estimated value in some cases. This difference grows to infinity while passed on during training.

B. Visualisation and Discussion

Figure 7 shows the AMAD of all six different datasets and also the one achieved by the original NA for the dataset with humans without background. It shows the modified NA can take on the results of the original one even with images of humanoids with textured background.

As already stated out, the monster HBDs are difficult to estimate. But we do recommend to keep a hugh varity of different body shapes due to diversity. Nevertheless one

should only include images showing survivable humans.

Fig 6 shows $k = 4$ instances from a mini batch of the humans dataset with plain background. The models match the very same images discussed in the original paper [7] to give a comparison between the two implementations. The tables beneath the images contain the real and predicted value and the resulting relative percentage error. In the last column there is the achieved relative percentage error reported in the original paper. The shoulder with was also hard to estimate for the modified network but not only in the case of an human in pose 0. The height also achieved one of the lowest errors across all subjects and HBDs

The same humans with an RGB background achieve only slightly worse results as with the plain background dataset

and even the images with added textures show very good results. When looking at the height, the plain background dataset showed an error range from -0.43 % to 0.07 %, the RGB background dataset -1,21 % to 0.69 % and the texture background dataset showed -0,24 to 0,78.

Taking the shoulder width, we achieve with the plain background dataset -1,90 % to 4,30 %, with the RGB background dataset -1,77 % to 3,48 % and with the texture background dataset we reach an error of -0,97 % to 2,56 %.

Also some errors might seem not irrelevant at first glance a look at the actual and predicted values relativizes it. Almost no errors of the four presented examples across all different background types exceeds 3 cm.

V. CONCLUSION

Our goal was to create datasets with artificial backgrounds and evaluate the performance of the Neural Anthropometer [7] when using those augmented datasets. We changed the number of connections in the last linear layer, the normalization method and the pooling strategy. The rest of the CNN Architecture remained the same. However, we achieved far better results as in the original paper, even with the texture and background dataset. AMAD was improved by 57% and ARPE by 50%.

REFERENCES

- [1] A. F. Agarap, “Deep learning using rectified linear units (relu),” 2018.
- [2] G. Bradski, “The OpenCV Library,” *Dr. Dobb’s Journal of Software Tools*, 2000.
- [3] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi, “Describing textures in the wild,” in *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [4] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” 2015.
- [5] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black, “SMPL: A skinned multi-person linear model,” *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, vol. 34, no. 6, pp. 248:1–248:16, Oct. 2015.
- [6] Paszke, et. al., “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds. Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
- [7] Y. G. Tejeda and H. A. Mayer, “A neural anthropometer learning from body dimensions computed on human 3d meshes,” 2021.
- [8] Y. Wu and K. He, “Group normalization,” 2018.
- [9] J. Xu, X. Sun, Z. Zhang, G. Zhao, and J. Lin, “Understanding and improving layer normalization,” 2019.