

# Human shape estimation with artificial backgrounds

Armin Niedermueller<sup>1</sup>, Melanie Urban<sup>1</sup> and Ahmet Cihat Bozkurt<sup>1</sup>

**Abstract—Human Body Dimension (HBD) Estimation using images becomes increasingly important.** Datasets consisting of images of synthetically created human body meshes have been successfully used for training convolutional neural networks (CNN) to estimate HBDs. However, those images all have a uniform grey background which is rather improbable in practical application. But adding an artificial background to the image might make it harder to clearly identify the human shape. In this work, we introduce uniformly colored and textured background into an existing dataset of images of human body meshes. We develop and train an existing network architecture and then let it predict the HBDs. Our modified network achieved better results ( $\sim 50\%$  less relative percentage error) than those reported in the original work we focused on. And it even accomplish remarkable results with different backgrounds.

## I. INTRODUCTION

In this study, we developed a model that predicts human body measurements. Therefore, we refined the architecture of the convolutional neural network *Neural Anthropometer* (NA) [8] and used the provided dataset [8] to be able to compare the estimations. This dataset contains 2D images of synthetically created 3D human body meshes in two different poses and related body measurements. The dark grey coloured body meshes are in front of a light grey coloured background.

Since the model trained with uniform background images will have a limited area of application, 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 store might provide a tool to upload a picture of the customer and then provide the appropriate dress size.

We tested our model with images of different body shapes with *Red-Green-Blue colour space* (RGB) background or with texture background. We evaluate the results by calculating the average model accuracy, *Mean Absolute Deviation* (MAD) for each body dimensions and *Relative Percentage Error* (RPE) for each body dimension.

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. Then, we evaluate and discuss the results of testing our model.

## II. HUMAN BODY DIMENSION ESTIMATION

There are many ways to estimate human body dimensions and just as many different input types. We based our model on the descriptions made in the paper *A Neural Anthropometer Learning from Body Dimensions Computed on Human 3D Meshes* [8]. They used images of 3D meshes synthesized with the Skinned Multi-Person Linear Model (SMPL) [5] 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/hbd\\_estimation\\_cnn](https://github.com/nerovalerius/hbd_estimation_cnn) and the code of the original NA and the dataset at: <https://github.com/neoglez/neural-anthropometer>

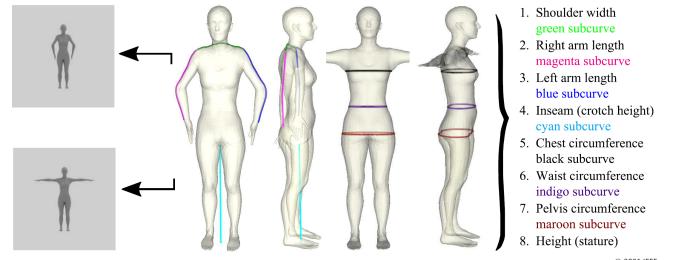


Fig. 1. Method of generating 2D images from 3D meshes which was used in the original paper [8]. The generated human body dimensions are: shoulder width, right and left arm length, inseam and calvis. The meshes present females and males in two poses, which then are transformed to grey scale 200x200 images.

### A. Images of Human Body Meshes

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. See Fig. 1 for more details.

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. [8] But there is the possibility to end up with so-called *monsters*. Those are body shapes with unrealistic body measurements. We also used these monster images, but separately from the humans.

Two different poses are used, pose 0 and pose 1 (Fig. 2). Pose 0 shows the body with arms stretched out on both sides and pose 1 shows the same body with arms slightly bent. Thus, the dataset consists 3500 females respectively males each in pose 0 and in pose 1.

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

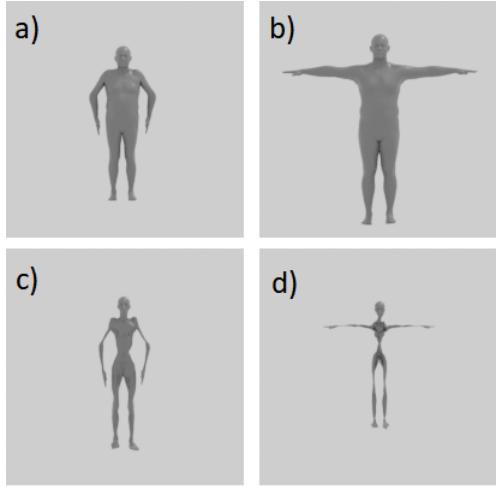


Fig. 2. Different images of human body meshes from the original dataset [8]. a) shows a human male in pose 0. b) shows a human male in pose 1. c) shows a monster male in pose 0. d) shows a monster female in pose 1.

Furthermore, the dataset provides the correct annotations for the HBDs of each image in order to let our network be able to train. 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 images by adding two types of backgrounds thus, generating two new datasets. For the first new dataset, we replaced the uniform grey background by inserting a background where all pixels have the same random RGB value. For the second dataset, we replace the background with textures of the *Describable Textures Dataset* (DTD), 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. If a DTD image is used twice, it is attached to a different body shape. Examples 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. METHODS

### 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. Textured backgrounds added to the human body mesh images. Row 1 and 2 show texture backgrounds from the Describable Textures Dataset.

### B. Network Architecture

Our neural network implemented with Pytorch [6] is shown in Fig. 4. We used the architecture from [8] which is described as follows:

The input layer processes  $200 \times 200 \times 1$  images. A convolution with eight  $5 \times 5$  filters produces a feature map of size  $196 \times 196 \times 8$ , which is passed through a rectified linear unit (ReLU) layer [1] and is then batch normalized. After max pooling with stride 2, a second convolution is applied (again squared kernel of size 5), which results in a tensor of size  $94 \times 94 \times 16$ . After another max pooling, the output is flattened to a vector of size 35344. This vector is passed to a fully connected and through a ReLU layer. Finally the last linear layer outputs the eight human body dimensions in meters [8, p. 6].

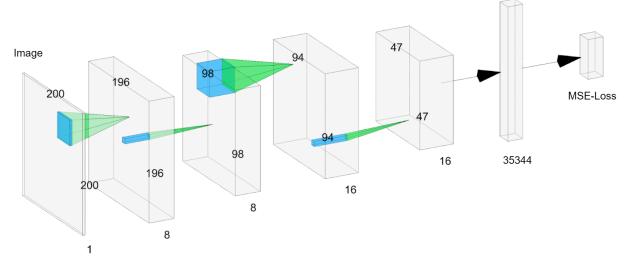


Fig. 4. CNN architecture based on the original framework [8].

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

#### Layer Normalization

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 [10]. As shown in Fig. 5 the Batch Normalization normalizes the tensor across the batch-size and the spatial dimensions

for each channel while the Layer Normalization uses for normalization the values computed across all channels for each sample [9]. 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 with only one channel - because the pictures are only in shades of gray. 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 [10].

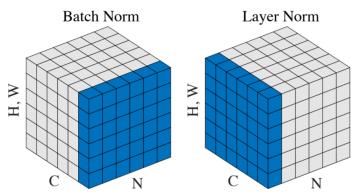


Fig. 5. Different normalization methods for a tensor with batch axis N, channel axis C, the height H and the width W. The blue pixels will normalized by calculating the values required for this from them. [9]

#### Pooling Layer

Both NAs use a MaxPool2D-Layer [6] twice. But instead of stride 2 and kernel size 2 we used also a stride 2 but a kernel size of 1. This means the amount of pixels is still downsampled to the half by but without max pooling. This strategy of down sampling by striding and not by max-pooling increases the amount of trainable parameters on the one hand. But on the other hand replacing pooling layers at all by strided convolution layers improves and stabilizes the model's performance. [7] To stay closer to the original NA we used the MaxPool2D-Layer to achieve the same.

#### Fully Connected Layer

After the second max pooling layer the output is flattened to size 1x35344 and passed to a fully connected layer. After cross-validating different output sizes we fixed it to 512. The original NA uses an output size of 84. It seems that rapidly reducing the vector's size before propagating to the final ReLu results in too much information loss.

## IV. EXPERIMENTS AND RESULTS

We take once the human and once the monster dataset (section II-A) and use k-fold cross-validation ( $k = 5$ ). We train our network for 20 epochs with a mini-batch size of 100. *Mean Squared Error* between the actual and predicted HBDs is minimized during training. Furthermore we use stochastic-gradient-descent with momentum term ( $\alpha = 0.9$ ) and learning rate ( $\eta = 0.01$ ) as optimization method. Those hyperparameters are identical to those provided in the paper of the original NA [8].

The experiment is performed on a computer with an AMD Ryzen 9 5950X CPU and a NVIDIA RTX 3090 GPU. Computing 20 epochs each over  $k = 5$  folds finished in about 5 minutes.

TABLE I  
HBD ESTIMATIONS COMPARED TO [8]

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  
HBD ESTIMATIONS WITH DIFFERENT BACKGROUNDS (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

#### 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 human 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 [8]:

Statistics (here for the humans dataset) are based on a result 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).

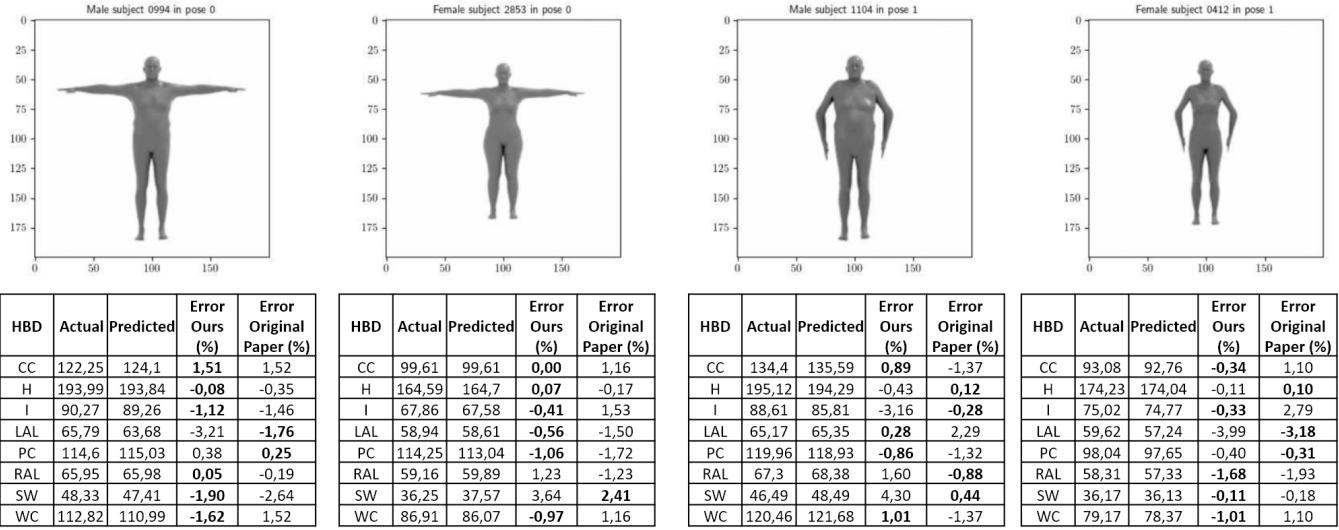


Fig. 6. Estimation results and relative percentage errors of the four example images also selected in the paper [8]. Bold numbers indicate the better value of our approach, when compared to the reference results. (CC: chest circumference, H: height, I: inseam, LAL: left arm length, PC: pelvis circumference, RAL: right arm length, SW: shoulder width, WC: waist circumference)

TABLE III  
HBD ESTIMATIONS WITH DIFFERENT BACKGROUNDS (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

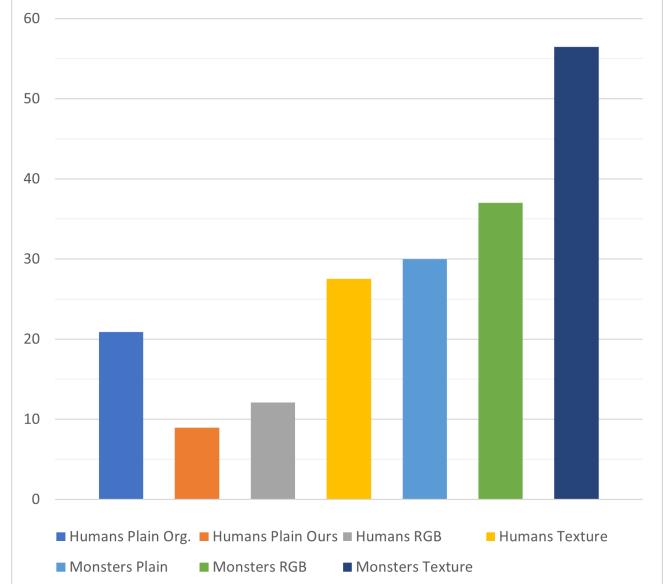


Fig. 7. AMAD on the test set after 20 epochs of training. Clearly the best results are achieved with the plain and single color RGB background. When looking at the human-dataset with texture background, the AMAD almost doubles.

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

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

For having a baseline and compare our modified NA to the [8] 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 decrease from 2.82 to 1.42 mm.

The tables II and III show the MAD and RPE for each HBD for the human dataset and the monster dataset respectively. They also contain all three different backgrounds: texture, RGB or plain.

Also for the human 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 as high than the one achieved

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

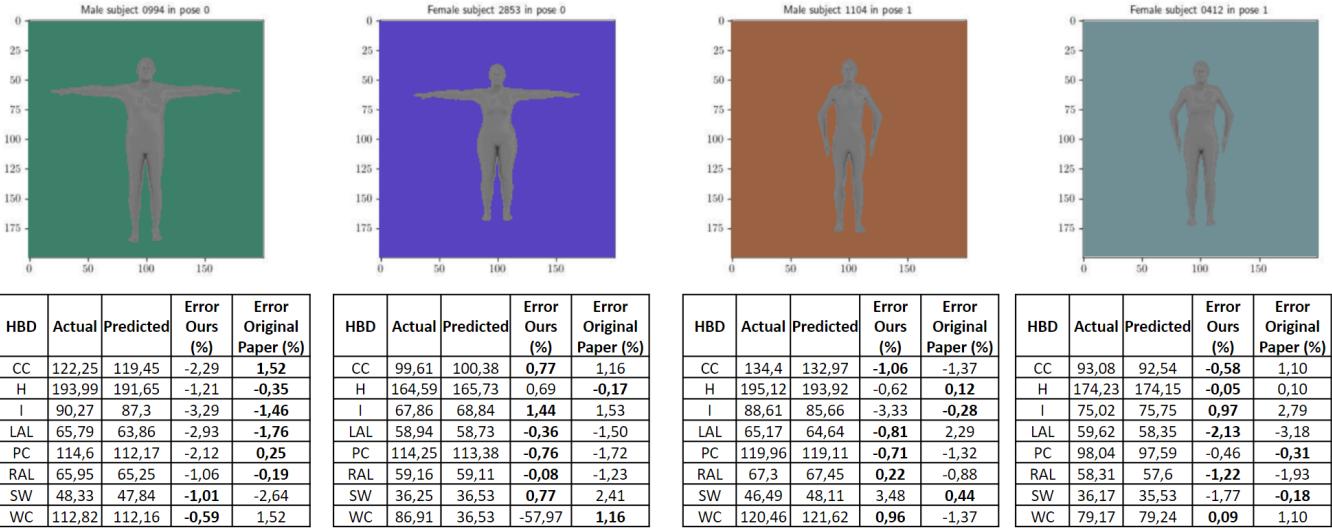


Fig. 8. Estimation results of the four example images from 6 with added RGB background. After the image is converted from colour to grayscale, the background results in a uniform shade of gray.

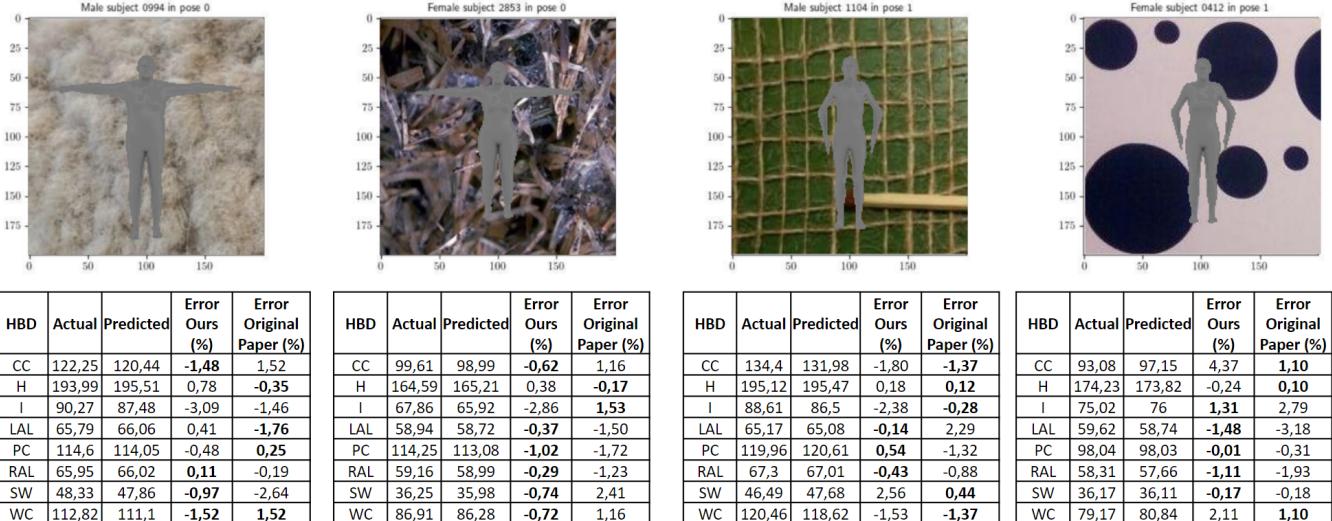


Fig. 9. Estimation results of the four example images from 2 with added textured backgrounds.

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 monster dataset was indeed difficult for the modified NA. 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 that our modified

architecture yields better results than the original one. It even shows lower AMAD values on the human-dataset with RGB and texture backgrounds compare to the results of using plain backgrounds on the original architecture.

As already stated, the monster HBDs are difficult to estimate and (expectedly) exhibit worse results.

Fig 6 shows four example images with plain background. The example images match the very same images discussed in the original paper [8] to give a comparison between the two implementations. The tables beneath the images contain the real and predicted values and the resulting relative percentage error reported in the original paper is presented. The shoulder width 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 example images with a RGB background achieve only slightly worse results than with a plain background and even the example images with added textures show very good results. When looking at the height, the plain background showed an error range from -0.43 % to 0.07 %, the RGB background -1.21 % to 0.69 % and the texture background showed -0.24 % to 0.78 %.

Inspecting the results for shoulder width, we achieve with the plain background -1.90 % to 4.30 %, with the RGB background -1.77 % to 3.48 % and with the texture background -0.97 % to 2.56 %.

Some errors might seem not irrelevant at first glance but a look at the actual and predicted values relativizes it. Almost no error of the four presented examples across all different backgrounds exceeds 3 cm.

## V. CONCLUSION

Our goal was to enhance a existing dataset of images of human body meshes with artificial backgrounds in order to make the prediction of human body dimensions more difficult and realistic. For data augmentation we added uniformly colored and textured backgrounds, thus creating two new datasets. Afterwards, we built the original network [8] and then modified its architecture. 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. The next step was to train the

modified using the three datasets – plain background, color background, texture background – splitted in  $k = 5$  folds and trained over 20 epochs per fold. We then evaluated the results of our architecture and found that we achieved better results ( $\approx 50\%$  less relative percentage error) than those reported in the original work. Specifically, AMAD was improved by 57% and ARPE by 50% in the case of the original dataset (i.e. plain backgrounds).

## 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.
- [7] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, “Striving for simplicity: The all convolutional net,” 2014.
- [8] Y. G. Tejeda and H. A. Mayer, “A neural anthropometer learning from body dimensions computed on human 3d meshes,” in *2021 IEEE Symposium Series on Computational Intelligence (SSCI)*, 2021, pp. 1–8.
- [9] Y. Wu and K. He, “Group Normalization,” 2018.
- [10] J. Xu, X. Sun, Z. Zhang, G. Zhao, and J. Lin, “Understanding and Improving Layer Normalization,” 2019.