# Towards Accurate 3D Human Body Reconstruction from Silhouettes

## Supplementary Material

Brandon M. Smith, Visesh Chari , Amit Agrawal, James Rehg, Ram Sever
Amazon.com Inc.
`{smithugh,viseshc,aaagrawa,jamerehg,severram}@amazon.com`

## 1. Synthetic Data Generation Details

Large segmentation errors are created by selecting regions of the silhouette at random and filling them with random splotches. The process is illustrated in Figure 1. This mimics the behavior of the segmentation network in uncertain image regions, and is generated as follows: (1) select a random silhouette boundary pixel, (2) construct a randomly sized region of interest (ROI) around the selected pixel (from 10 to 200 pixels on a side), (3) select all boundary pixels withing the ROI, (4) dilate these pixels using either a square or round kernel with random diameter between 10 and 30 pixels, (5) fill the dilated region with white noise, (4) smooth the dilated region with a random-bandwidth Gaussian filter, and (5) threshold, which creates splotches of various sizes. These splotchy regions closely mimic the kinds of segmentation errors observed in low confidence regions in real images, where foreground and background labels are predicted randomly. Once created, noise patches are added to, subtracted from, or replace the corresponding region of the silhouette image at random. Empirically we modeled low confidence with a score of $0.2$, and high confidence with a score of $0.9$. We use these synthetic images for pre-training our network, and then subsequently fine-tune using real segmentation results on photorealistic CAESAR [1] scans rendered in front of random background images from the LSUN dataset [2].

## 2. Repeatability Examples

We investigated the repeatability of our system on more realistic CAESAR[1] renderings. We rendered 1k color scan in front of ten different backgrounds, each from a different virtual camera viewpoint. This produced a dataset of 10k instances. To better handle the subtleties of real segmentation noise and confidence, we fine-tuned our full network on 9k instances. On the remaining 1k instances (100 unique scans outside the training set) we measured the standard deviation of each vertex across the 10 different renderings of each scan. The average vertex standard deviation was 3.09mm. For reference, 3mm corresponds to approximately 1 pixel for an average height person occupying 600 vertical pixels of a VGA ($640 \times 480$) image.

Three example repeated sessions are shown in Figure 2. The estimated 3D model are consistent (low standard deviation) across a wide range of inputs for the same person. The plot in Figure 2 shows that our method produces results with low *intra*-subject measurement variance compared to the amount of variance between different subjects.



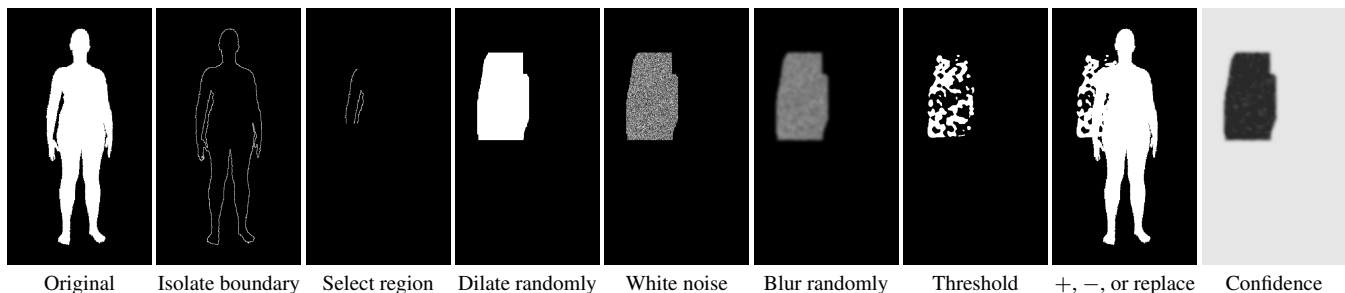| Original | Isolate boundary | Select region | Dilate randomly | White noise | Blur randomly | Threshold | $+, -$, or replace | Confidence |

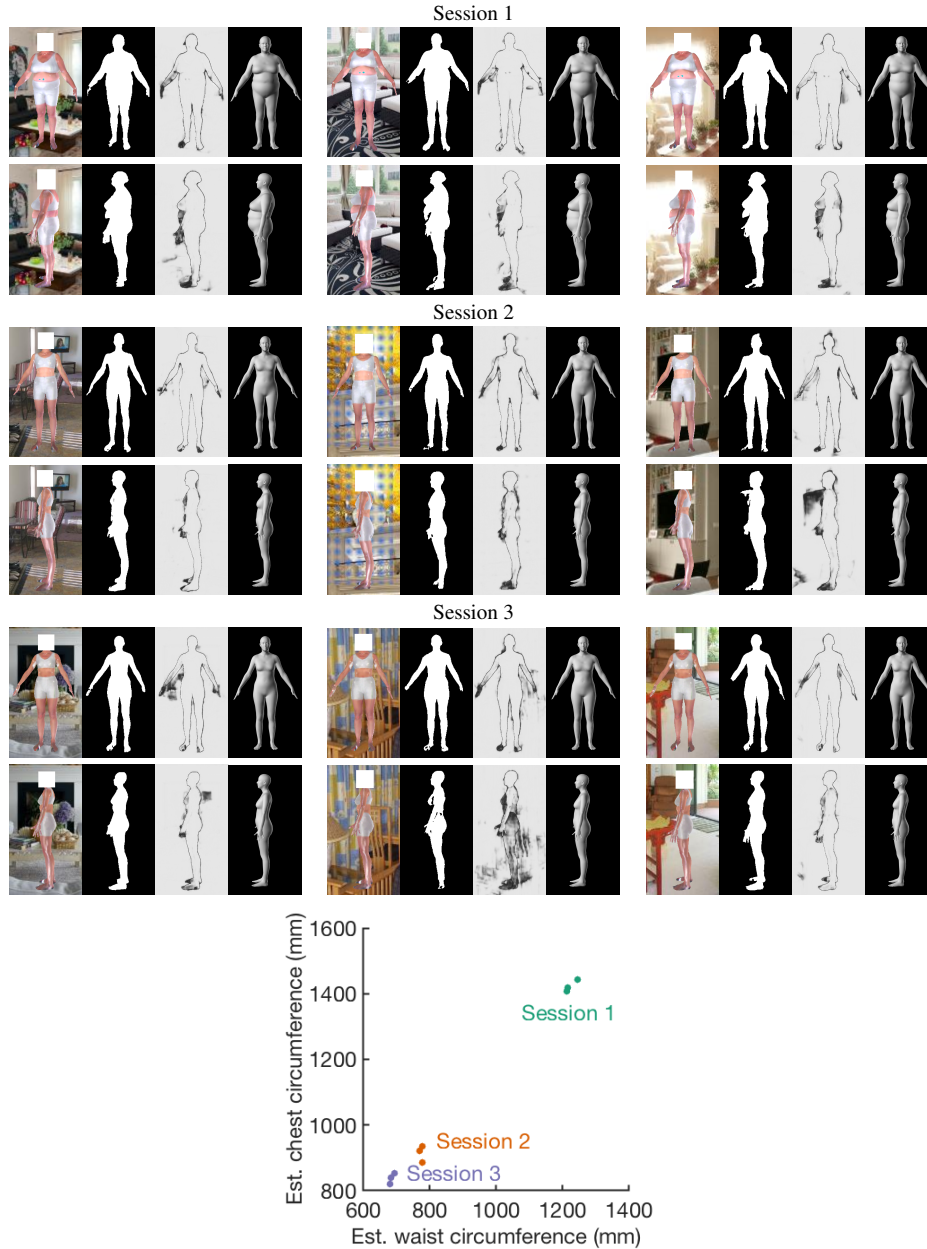Figure 1. An example that illustrates synthetic segmentation noise and confidence mask generation.

Figure 2. Repeatability examples. **Top:** Three example repeated sessions (input RGB images, binary segmentation masks, segmentation confidence, and estimated 3D body model). Despite significant camera variation and segmentation noise, the estimated 3D models are consistent within each session. **Bottom:** For repeated captures of the same person with different camera positions and backgrounds, our method produces results with low intra-subject measurement variance compared to the amount of variance between sessions/subjects, *i.e.*, clusters are tight and well separated.

## 3. Qualitative Results on CAESAR Renderings

Here we present additional qualitative results from the experiment described in Section 4.4 of the main paper. Each input example is an RGB scan from the CAESAR dataset [1] rendered in front of a random background sampled from the LSUN image dataset [2]. Figures 3, 4, 5, and 6 highlight good results, and Figure 7 highlights several problem cases.

## 4. Architecture Details

Our proposed system was implemented in Keras with a TensorFlow backend. In Figure 8 we provide our network architecture. Please see Figure 1 in the main paper for a high-level overview.

## References

[1] K. M. Robinette, S. Blackwell, H. Daanen, M. Boehmer, S. Fleming, T. Brill, D. Hoeferlin, and D. Burnsides. Civilian American and European Surface Anthropometry Resource (CAESAR) final report. *Tech. Rep. AFRL-HEWP-TR-2002-0169, US Air Force Research Laboratory*, 2002. 1, 2

[2] F. Yu, Y. Zhang, S. Song, A. Seff, and J. Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. 1, 2
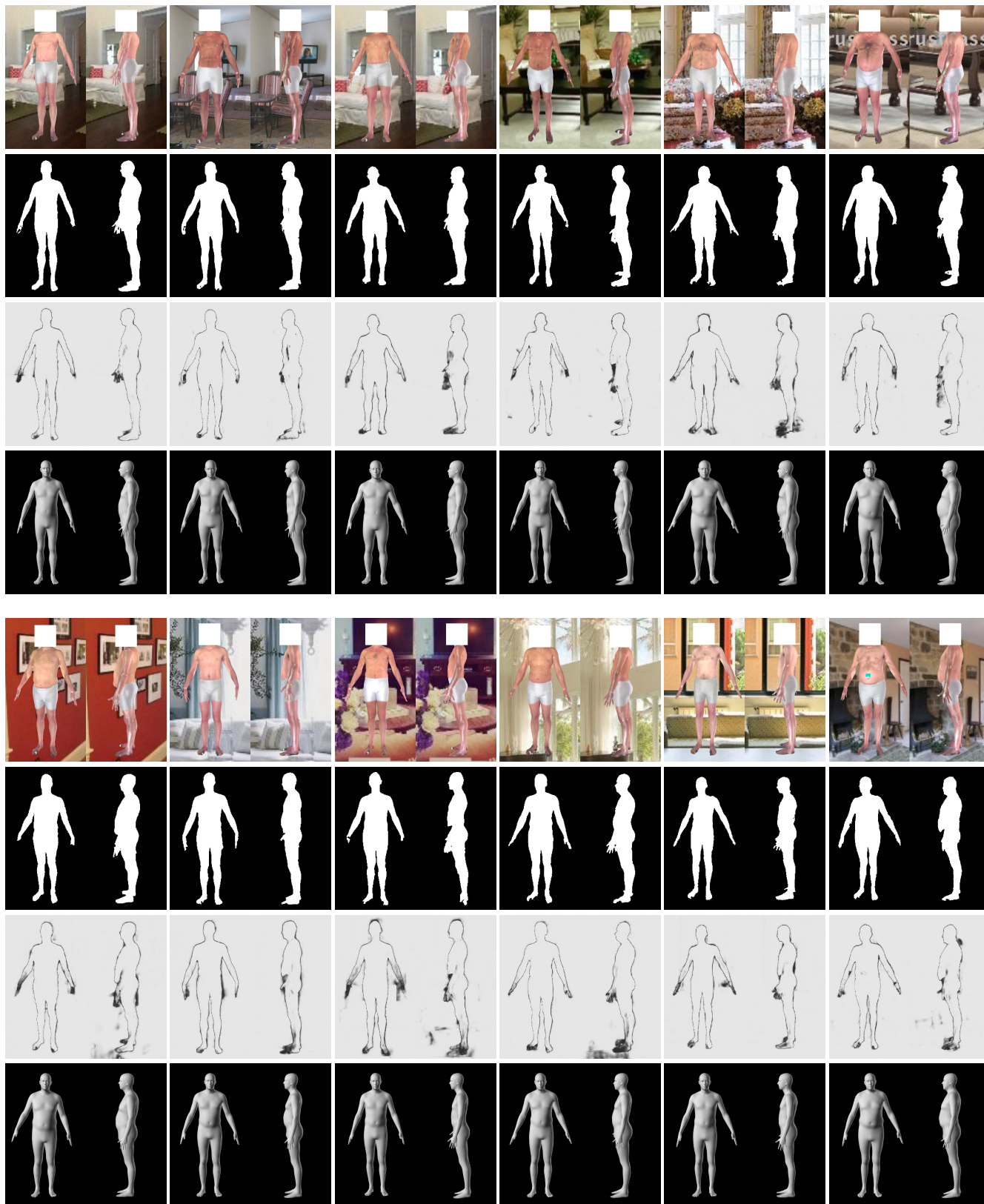
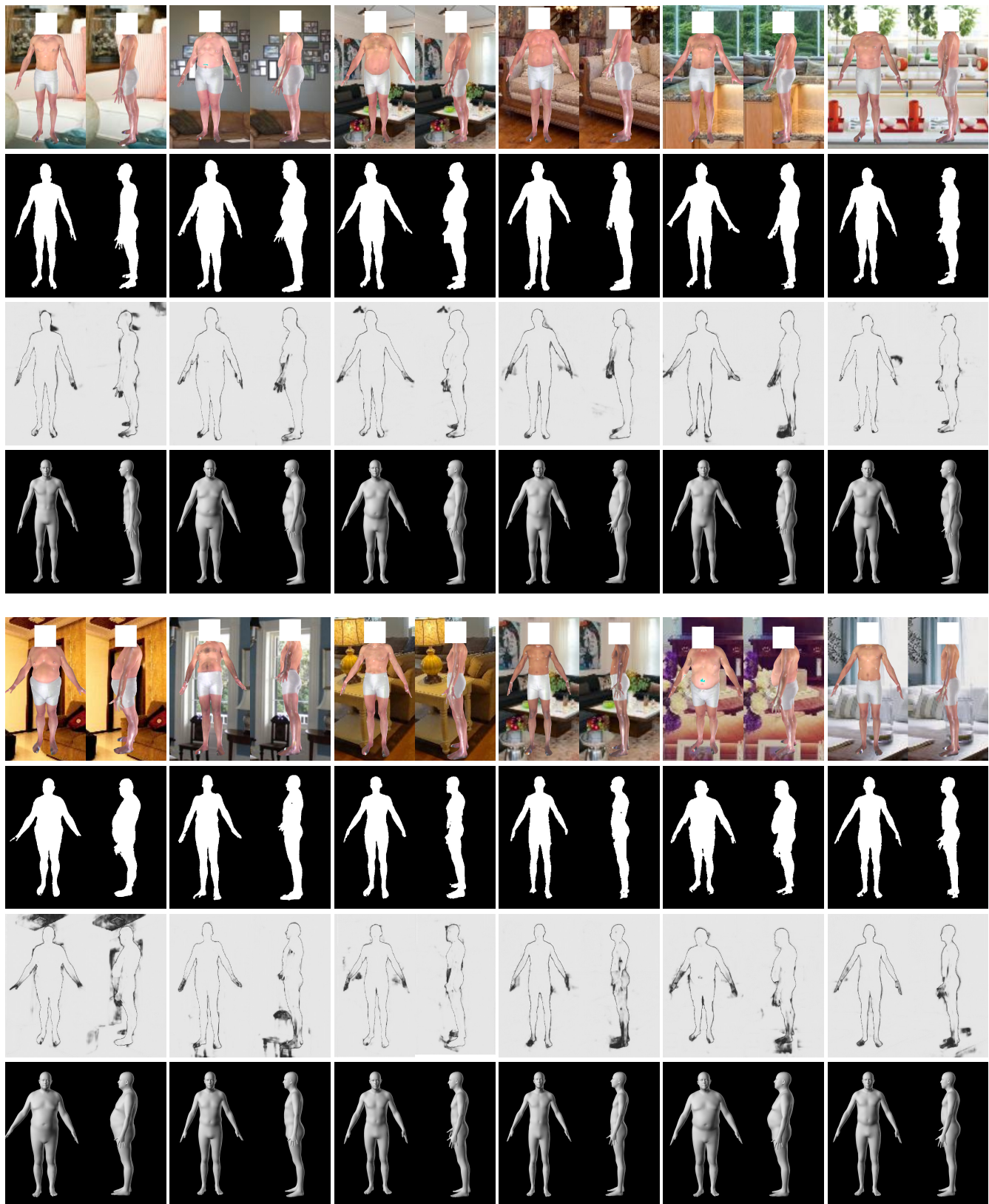Figure 3. Qualitative results on male CAESAR images, part 1.

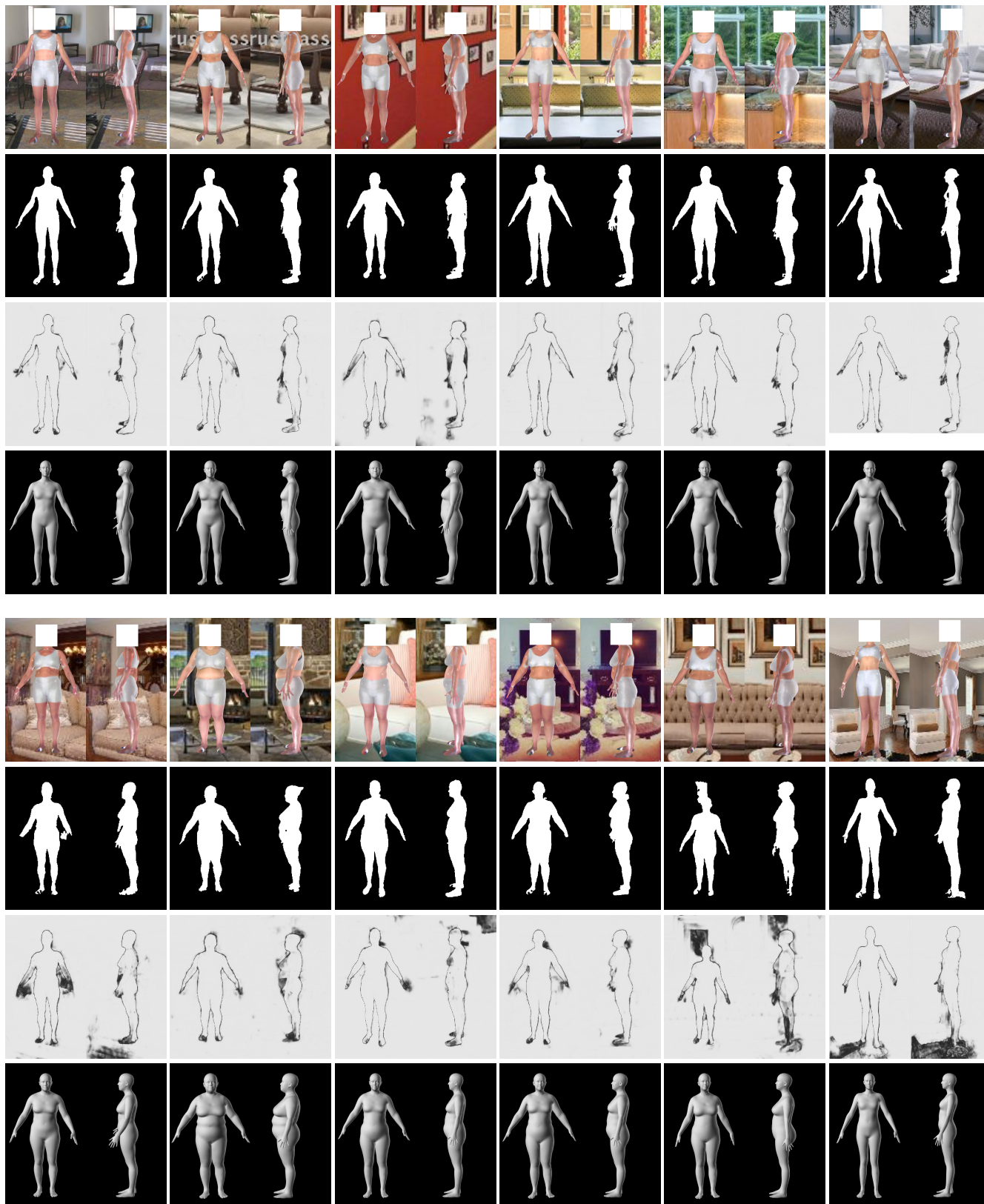Figure 4. Qualitative results on male CAESAR images, part 2.

Figure 5. Qualitative results on female CAESAR images, part 1.

Figure 6. Qualitative results on female CAESAR images, part 2.

Figure 7. Problem cases. Large segmentation errors, although rare, can lead BaReNet to estimate 3D body shapes that don't accurately reflect the image.
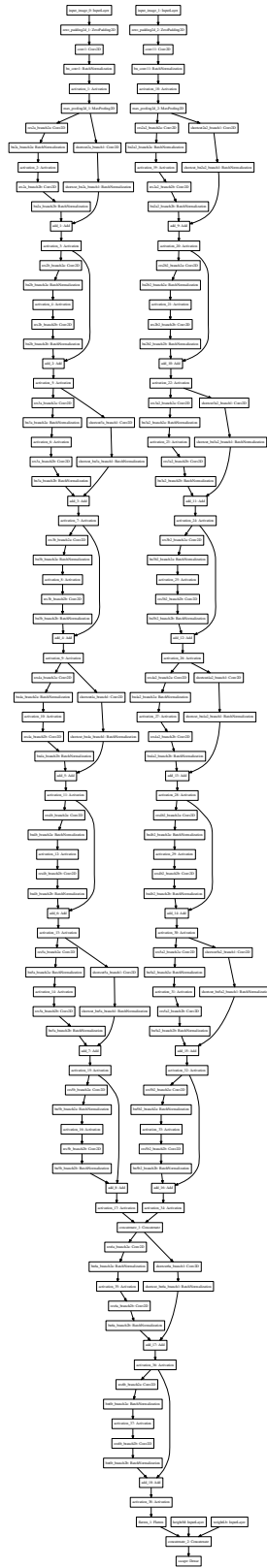
Figure 8. SfSNet Keras architecture details, as generated by the Keras model visualization tool.