# BodyNet: Volumetric Inference of 3D Human Body Shapes

Gül Varol[1], Duygu Ceylan[2], Bryan Russell[2], Jimei Yang[2], Ersin Yumer[3], Ivan Laptev[1] and Cordelia Schmid[1]

[1]Inria    [2]Adobe Research    [3]Argo AI

**ECCV 2018** — European Conference on Computer Vision

## INTRODUCTION

▶ Goal & Contributions
 ▷ Predicting 3D human body **pose and shape** given a **single RGB** image as input.
 ▷ Demonstrating advantages of auxiliary body-related tasks in an end-to-end multi-task setting.

▶ Motivation
 ▷ Volumetric representation of human bodies in the context of neural networks is not studied.
 ▷ Volumetric representation is flexible, e.g. can capture clothing.

input   output voxels   output parts

## BODYNET APPROACH

▶ The architecture benefits from the **multi-task** training of:
 ▷ a **volumetric** 3D loss,
 ▷ a **multi-view re-projection** loss,
 ▷ intermediate supervision of **2D pose, 2D part segmentation, and 3D pose**.

2D pose loss $\mathcal{L}_j^{2D}$   3D pose loss $\mathcal{L}_j^{3D}$   Volumetric loss $\mathcal{L}_v$

volumetric shape   SMPL fit

2D segmentation loss $\mathcal{L}_s$

Re-projection loss $\mathcal{L}_p^{FV}$   Re-projection loss $\mathcal{L}_p^{SV}$

end-to-end   $\mathcal{L}_s + \mathcal{L}_j^{2D} + \mathcal{L}_j^{3D} + \mathcal{L}_v + \mathcal{L}_p^{FV} + \mathcal{L}_p^{SV}$   optimization

▶ We gradually increase the difficulty of the task to go from 2D to 3D:

RGB   2D pose&segm   3D pose   voxels   SMPL

## EXTENDING TO 3D BODY PART SEGMENTATION

▶ Last layer weights are duplicated as many times as the number of parts to initialize training for **part voxels**.

## ARCHITECTURE STUDY

▶ Effect of **additional inputs**
 ▷ 3D shape estimation on SURREAL

Voxel IOU (%) — *higher is better*

| Input | Voxel IOU |
|---|---|
| 2D pose | 47.7 |
| RGB | 51.8 |
| Segm | 54.6 |
| 3D Pose | 56.3 |
| Segm + 3D Pose | 56.4 |
| RGB + 2D Pose + Segm + 3D Pose | 58.1 |

 ▷ 3D pose estimation (mm)

| Input | SURREAL | Human3.6M |
|---|---|---|
| RGB | 49.1 | 51.6 |
| 2D pose | 55.9 | 57.0 |
| Segm | 48.1 | 58.9 |
| 2D pose + Segm | 47.7 | 56.3 |
| RGB + 2D pose + Segm | **46.1** | **49.0** |

input image   2D predictions   3D pose prediction   3D voxels prediction   SMPL fit   Ground truth

## RESULTS: SURREAL dataset [Varol et al. CVPR 2017]

▶ Effect of multi-view **re-projection loss**
▶ Effect of **end-to-end** training
▶ Comparison with **alternative methods**

3D surface error (mm) — *lower is better*

| Alternative methods | | | BodyNet variants | | | | | |
|---|---|---|---|---|---|---|---|---|
| SMPLify++ | Tung 2017 | Shape parameter regression | no re-projection | FV | FV+SV | FV | FV+SV | FV | AS+SV |
| 75.3 | 74.5 | 74.3 | 73.6 | 69.9 | 68.2 | 72.7 | 70.5 | 67.7 | 65.8 |

no end-to-end   end-to-end without intermediate tasks   end-to-end with intermediate tasks

input image   original view   other view

Input   Shape parameter regression   SMPLify++   BodyNet   Ground truth

## RESULTS: Unite the People dataset [Lassner et al. CVPR 2017]

▶ Effect of the **re-projection type**

| | | 2D metrics | | | 3D metrics (mm) | |
|---|---|---|---|---|---|---|
| | | Acc. (%) | IOU | F1 | Landmarks | Surface |
| 3D ground truth | (Lassner et al.) | 92.17 | - | 0.88 | 0 | 0 |
| Decision forests | (Lassner et al.) | 86.60 | - | 0.80 | - | - |
| HMR | (Kanazawa et al.) | 91.30 | - | 0.86 | - | - |
| SMPLify, UP-P91 | (Lassner et al.) | 90.99 | - | 0.86 | - | - |
| SMPLify on DeepCut | (Bogo et al.) | 91.89 | - | 0.88 | - | - |
| BodyNet *(SMPL projections)* | | 92.75 | 0.73 | 0.84 | 83.3 | 102.5 |
| BodyNet *(manual segmentations)* | | **94.67** | **0.80** | **0.89** | | |
| 3D ground truth | (Lassner et al.) | 95.00 | 0.82 | - | 0 | 0 |
| Indirect learning | (Tan et al.) | 95.00 | **0.83** | - | 190.0 | - |
| Direct learning | (Tan et al.) | 91.00 | 0.71 | - | 105.0 | - |
| BodyNet *(SMPL projections)* | | 92.97 | 0.75 | 0.86 | **69.6** | **80.1** |
| BodyNet *(manual segmentations)* | | **95.11** | 0.82 | **0.90** | | |

RGB   GT silhouette   predicted silhouette   predicted voxels (front view) (other view)   predicted silhouette   predicted voxels (front view) (other view)

trained with manual annotation   trained with SMPL projection

## INTERMEDIATE TASKS

▶ **All tasks improve** with end-to-end training.

| | Segmentation mean parts IOU (%) | 2D pose PCKh@0.5 | 3D pose mean joint distance (mm) |
|---|---|---|---|
| Independent single-task | 59.2 | 82.7 | 46.1 |
| Joint multi-task | **69.2** | **90.8** | **40.8** |

▶ Weight **balancing** is important.

loss   2D segm (IOU)   2D pose (PCK)   3D pose (mm)   Voxels (IOU)

balanced / not balanced

## LIMITATIONS

3D ambiguity   multi-person

## CONCLUSIONS

▶ Volumetric representation is **flexible** and effective.
▶ Re-projection loss is critical to obtain **confident body surface**.
▶ Multi-task training of **relevant tasks** helps.

## REFERENCES

▶ *SMPLify:* Bogo et al. ECCV 2016
▶ *Unite the People:* Lassner et al. CVPR 2017
▶ *Indirect learning:* Tan et al. BMVC 2017
▶ *Self-supervised learning:* Tung et al. NIPS 2017
▶ *HMR:* Kanazawa et al. CVPR 2018

Code is available!