

DSNet: Dynamic skin deformation prediction by Recurrent Neural Network

Hyewon SEO¹, Kaifeng ZOU¹, Frédéric CORDIER²

¹ ICube laboratory, Université de Strasbourg, France

² IRIMAS, Université de Haute Alsace, France



Introduction

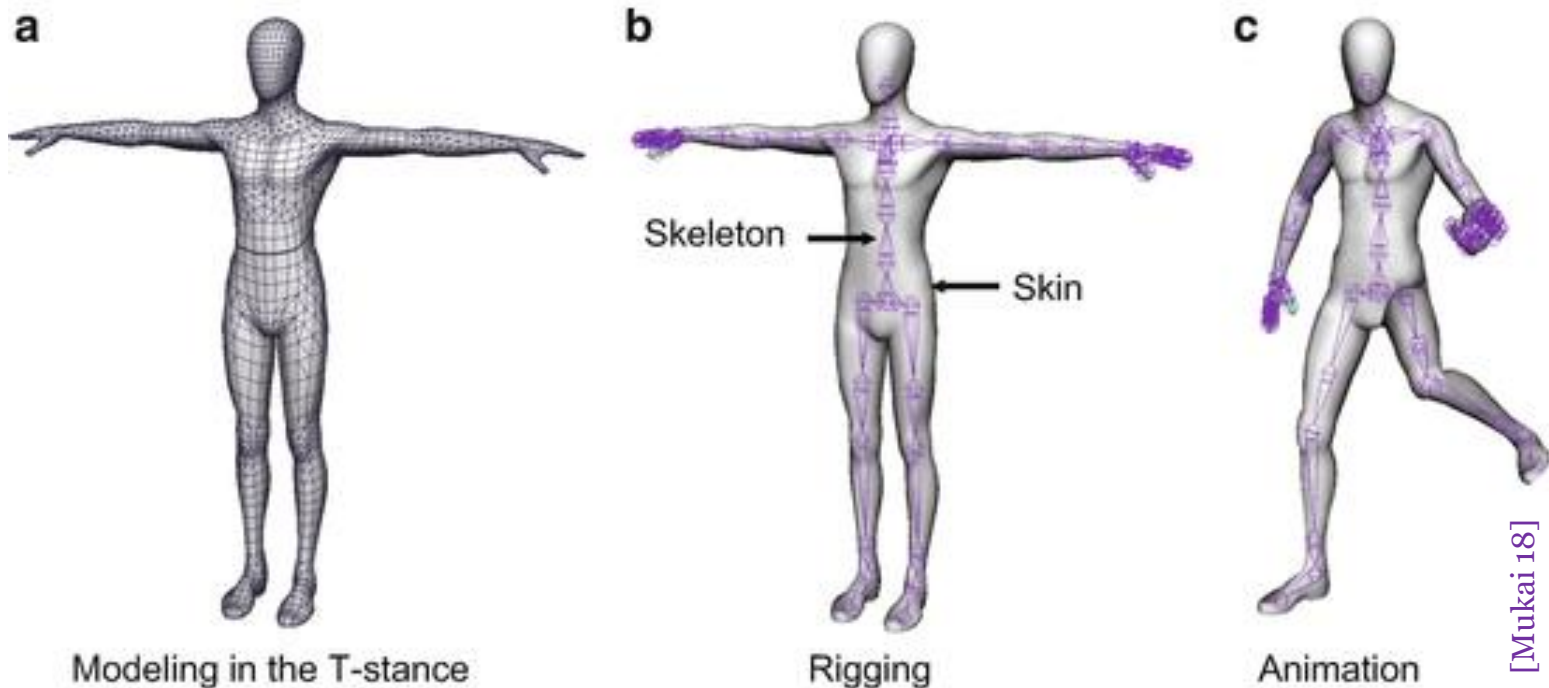
Dynamic skin deformation contributes to the enriched realism of character models in rendered scenes.



It has a long tradition in CG and CA...

Introduction

Linear blend skinning: [MTT91] $\mathcal{W}(\theta; M, J, W)$



Get the skin surface M .

Define the skeleton J .

Map vertices to the skeleton: W

App
skel
Rep

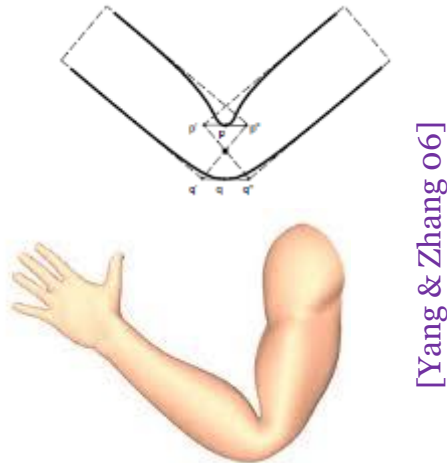
[MTT91] Magnenat-Thalmann N., Thalmann D., “Human Body Deformations Using Joint-dependent Local Element Theory”, Making Them Move, N.Badler, B.A.Barsky, D.Zeltzer, eds, Morgan Kaufmann, San Mateo, 1991.

[Mukai18] Tomohiko Mukai, Example-Based Skinning Animation, pp 2093-2112, Handbook of Human-Computer Interaction, 2018.



Introduction

Limitations of LBS



Unnatural deformations
at certain poses



Impossible to express
nonlinear deformation
i.e. muscle bulging

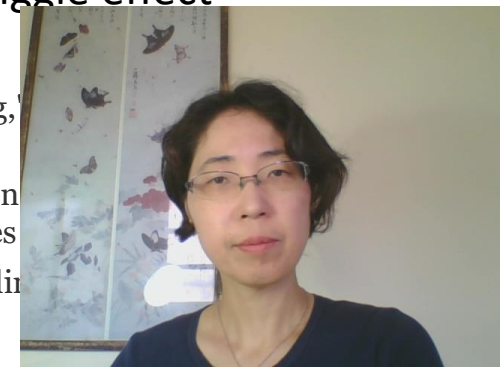


Impossible to simulate
skin dynamics
i.e. jiggle effect

[Yang & Zhang 06] Xiaosong Yang and J. J. Zhang, "Stretch It - Realistic Smooth Skinning," Computer Graphics, Imaging and Visualisation (CGIV'06), Sydney, Qld., 2006, pp. 323-328.

[Lewis et al 06] J. P. Lewis, Matt Cordner, and Nickson Fong. 2000. Pose space deformation interpolation and skeleton-driven deformation. Proc Computer graphics and interactive techniques

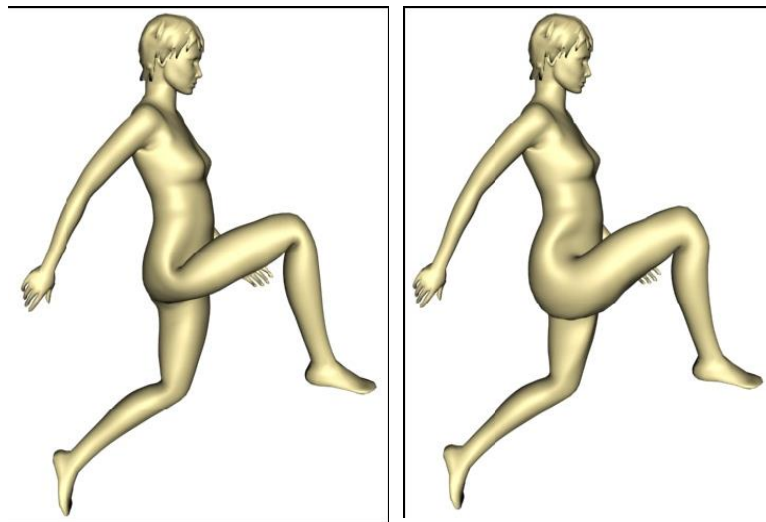
[Romero et al 20] Romero, Cristian & Otaduy, Miguel & Casas, Dan & Perez, Jesus. (2020). Modeling Skin Mechanics for Animated Avatars. Computer Graphics Forum. 39. pp. 77-88.



Previous work

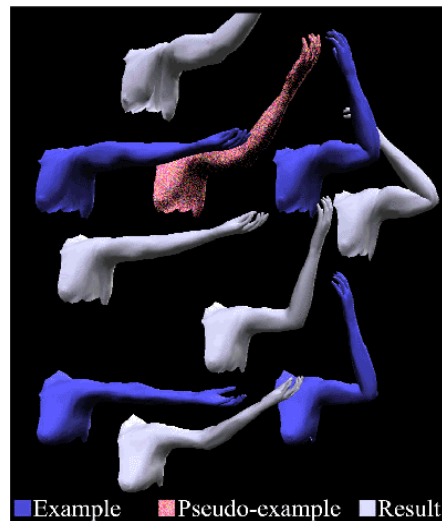
Solutions: Previous work

Geometric



[Magnenat-Thalmann et al. 04]

Example-based



[Sloan et al. 01]

Physics-based



[Ziva Dynamics]

[Magnenat-Thalmann et al. 04] N Magnenat-Thalmann, F Cordier, H Seo, G Papagianakis, Modeling virtual environments, 2004 International Conference on Cyberworlds, 201-208

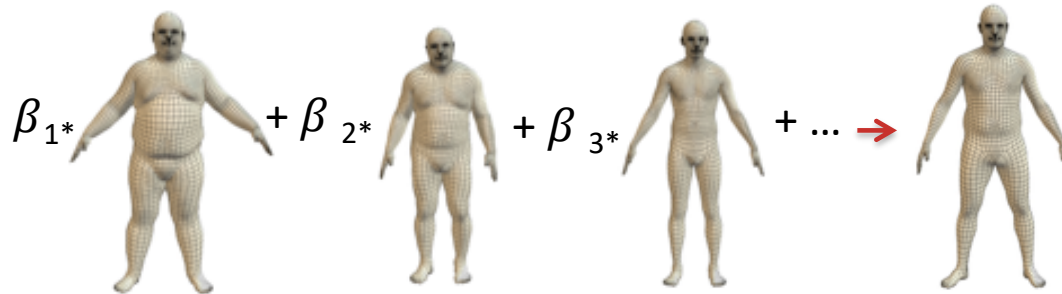
[Sloan et al. 01] P. P. Sloan, C. Rose and M. Cohen, "Shape by Example", ACM SIGGRAPH 2001 Computer Graphics, NC, USA, pp. 135-143, 2001.



Previous work

Data-driven body shape modelers [SMT03, ASK+05]

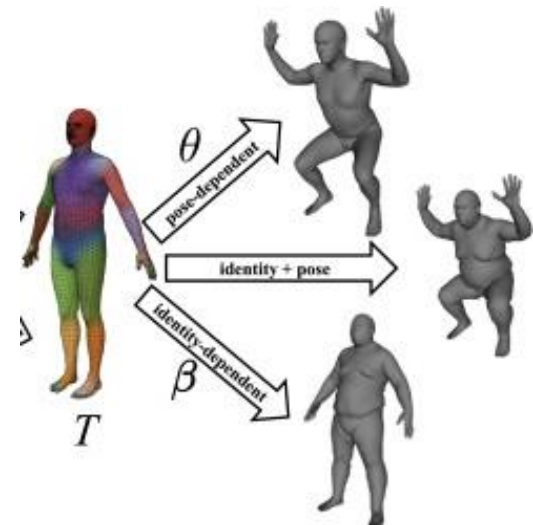
$$M(\vec{\beta}) = \bar{T} + \sum \beta_n S_n$$



Basis shape vectors



A unifying framework for subject- & pose-dependent shapes [HLRB12, LMRP+15]



[SMT03] Seo H., and Magnenat-Thalmann N., “An Automatic Modeling of Human Bodies from Sizing Parameters”, ACM SIGGRAPH 2003 Symposium on Interactive 3D Graphics (April), pp.19-26, Monterey, USA, 2003.

[ASK+05] D. Anguelov, P. Srinivasan, D. Koller, S. Thrun, J. Rodgers, and J. Davis J., SCAPE: Shape Completion and Animation of People. ACM Trans. Graph. (Proc. SIGGRAPH 24, 3, 408–416) 2005.

[HLRB12] D. Hirshberg, M. Loper, E. Rachlin, and M. Black, Coregistration: Simultaneous alignment and modeling of articulated 3D shape. In European Conf. on Computer Vision (ECCV), LNCS 7577, Part IV, 242–255, 2012.

[LMRP+15] 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), 2015.

Previous work

Data-driven dynamic human shape modelers [PMR+15, CO18]



$$M_t = LBS_t + \Delta_t$$

$$LBS_t = f_{linear}(\theta_t; \bar{M}(\beta), J(\beta), W)$$

$$\Delta_t = g_{linear}(v_t, a_t, \dot{\theta}_t, \ddot{\theta}_t, \Delta_{t-1}, \Delta_{t-2}; \bar{M}(\beta))$$

[PMR+15] Pons-Moll G., Romero J., Mahmood N., and Black M. J.: Dyna: a model of dynamic human shape in motion. *ACM Trans. Graph.* 34, 4, Article 120 (July 2015).

[CO18] Casas, D. & Otaduy, M. (2018). Learning Nonlinear Soft-Tissue Dynamics for Interactive Avatars. *Proc. ACM Computer Graphics and Interactive Techniques*. 1. 1-15.

DS-Net: Overview

Our goal is to learn a function $f(\{\underline{\theta}_t\}) = \{\Delta_t\}, t = 1, \dots, T$

c.f. $\varphi_t = \{v_t, a_t, \dot{\theta}_t, \ddot{\theta}_t\}$

 Both input and outputs are **sequences!!**

- We deploy LSTM network to learn our function.
- The results of frame t depend on the results of previous frames $t-1, t-2, \dots$
- We also consider subject specificity i.e. β .

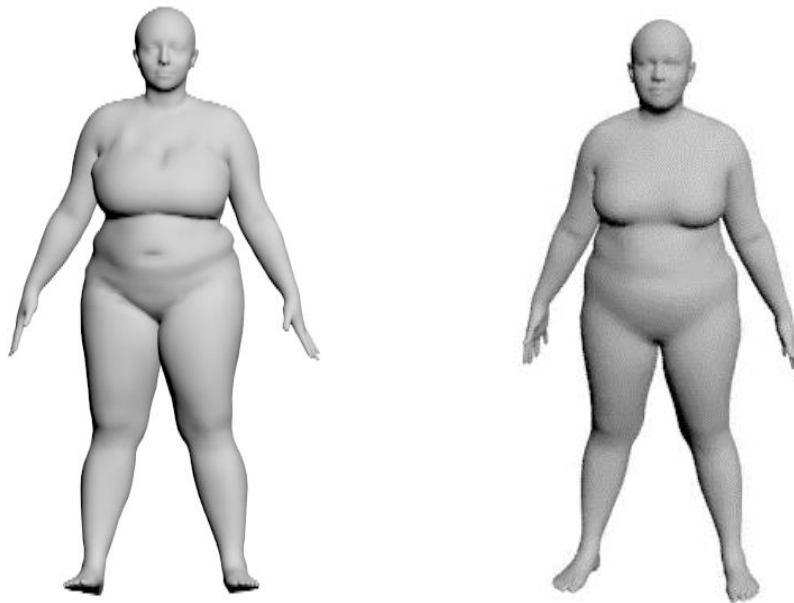
$$\Rightarrow \Delta_t = f(\theta_t, f(\theta_{t-1}), \beta)$$

A common shape space is required: SMPL! (A Skinned Multi-Person Linear Model) [LMRP+15]

DS-Net: dataset

Dyna dataset [PRMB15]

- Captured shapes exhibiting dynamic skin deformation
- 5 (female) subjects, 10~14 motions each
- Inter-, intra-subject correspondence with $N=6890$ vertices, 13776 triangles
- The duration of each sequence varies: 2 ~15 sec.



DS-Net: dataset

Dyna [PRMB15] : *training & validation*

Mosh [LMB14] : *test*

dataset	subjects	motions	fps	No. sequences (men/women)
Dyna	5 men, 5 women	10~14 motions for each subject: one-leg jumping, light hoping, jumping jacks, shake hips, running in place, etc.	60	66 / 67
Mosh	Same subjects as above	Includes some skin-dynamics inducing motions (side-to-side hoping, basketball, kicking) that are not included Dyna.	100	24 / 30

[PRMB15] Pons-Moll G., Romero J., Mahmood N., and Black M. J.: Dyna: a model of dynamic human shape in motion. ACM Trans. Graph. 34, 4, Article 120 (July 2015), 14 pages.

[LMB14] M. Loper, N. Mahmood, and M. J. Black. MoSh: Motion and Shape Capture from Sparse Markers. ACM Trans. Graph., 33(6):220:1–220:13, Nov. 2014.

Generation of training data

Extraction of SMPL parameters + residuals Δ , from each mesh.

For each motion sequence m :

1. Compute the best matching SMPL parameters $(\boldsymbol{\beta}, \boldsymbol{\theta}_1)$ at frame 1.

$$\min_{\boldsymbol{\beta}, \boldsymbol{\theta}_1} \left\| \mathcal{W}(\boldsymbol{\theta}_1, \bar{\mathbf{T}} + M_S(\boldsymbol{\beta}) + M_P(\boldsymbol{\theta}_1)) - S_1 \right\|_2.$$

2. Compute the best matching SMPL parameters $\boldsymbol{\theta}_t$ for each frame $t > 1$.

$$\min_{\boldsymbol{\theta}_t} \left\| \mathcal{W}(\boldsymbol{\theta}_t, \bar{\mathbf{T}} + M_S(\boldsymbol{\beta}^*) + M_P(\boldsymbol{\theta}_t)) - S_t \right\|_2.$$

Fixed throughout all frames $t > 1$.

3. The displacement vector is considered as the dynamic skin component.

$$\Delta_t = \mathcal{W}^{-1}(\boldsymbol{\theta}_t, S_t) - (\bar{\mathbf{T}} + M_S(\boldsymbol{\beta}^*))$$

Unposing operation: transforms a body mesh to its rest pose.

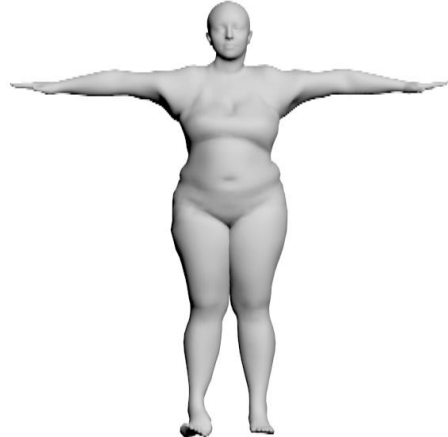
The training data is a set of input and output pairs : $\{(\boldsymbol{\beta}^m, \boldsymbol{\theta}_t^m, \Delta_t^m)\}, m=1\dots 65$.

Generation of training data

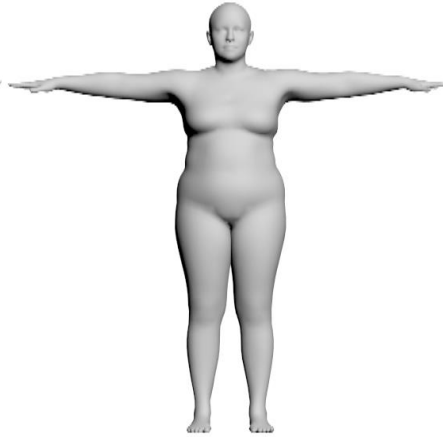
Mesh alignment results



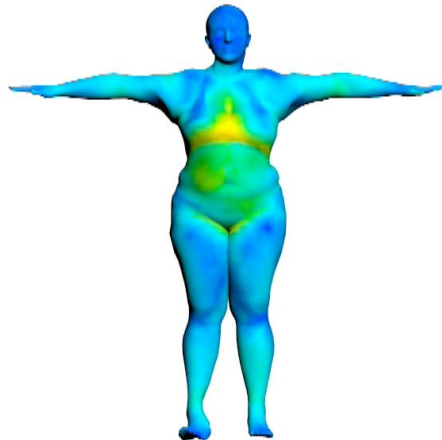
(a) S



(b) $\mathcal{W}^{-1}(\theta_t, S)$



(c) $\bar{T} + M_S(\beta)$



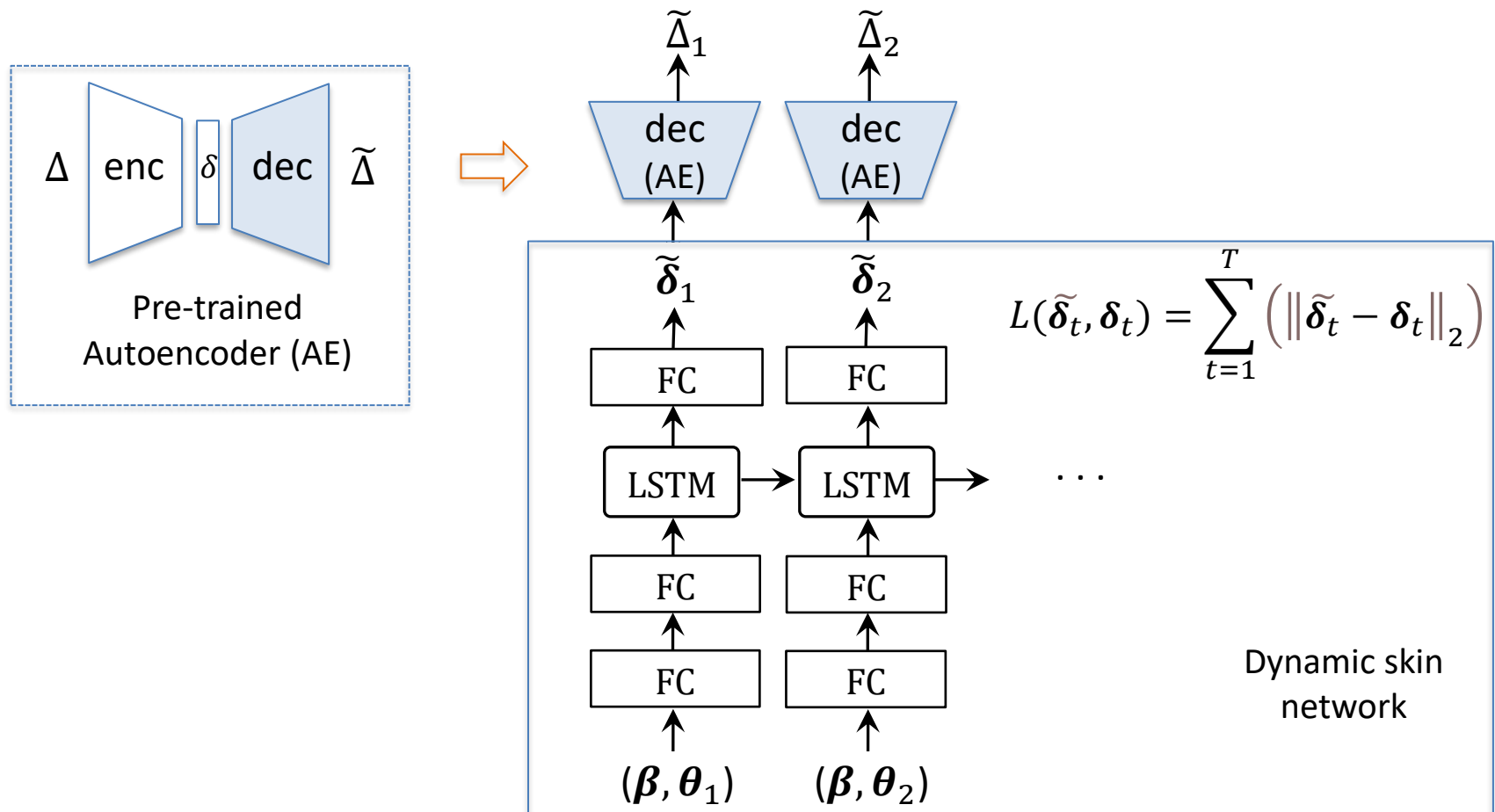
(d) $\bar{T} + M_S(\beta) + \Delta_t$

Skin displacements contributed by the dynamic skin deformation are recorded at a canonical pose θ^0 .

DS-Net: Architecture

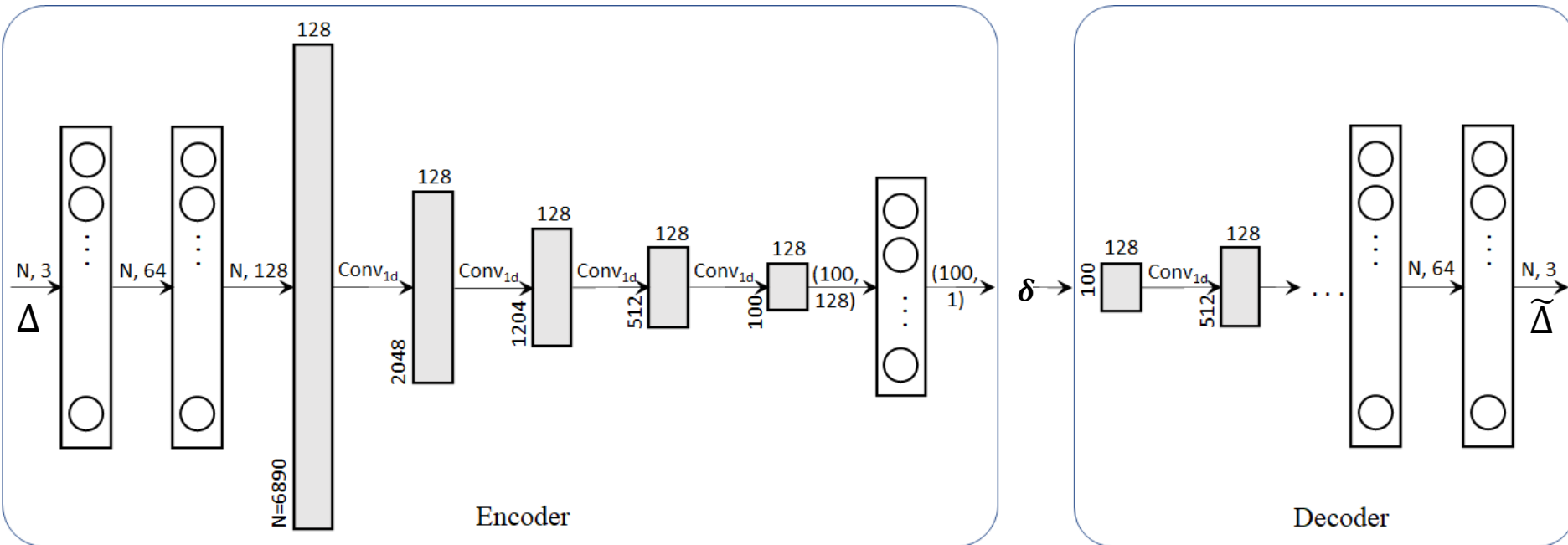
DSNet: Dynamic skin prediction

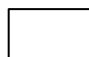

- The original data space resides in a high dimensional space: $\Delta_t \in R^{N \times 3}$ ($> 20\,000$)
- We represent them in a latent space by using an autoencoder: $\delta_t \in R^{100}$
- The DSNet LSTM [HS97] is trained on the latent space



Data dimension reduction

Displacement mesh autoencoder (AE):



 : dense layer
 : data

The dimension of the original mesh $3N$ ($3 \times 6890 = 20,670$) is reduced to **100!!**

DS-Net: AE details

Displacement mesh autoencoder (AE):

- The input data Δ has been normalized to $[-1,1]$.
- Pytorch implementation of Adam optimizer.
- Batch size 64, learning rate 0,0001.
- 11,8% of network parameters, compared to the other AE

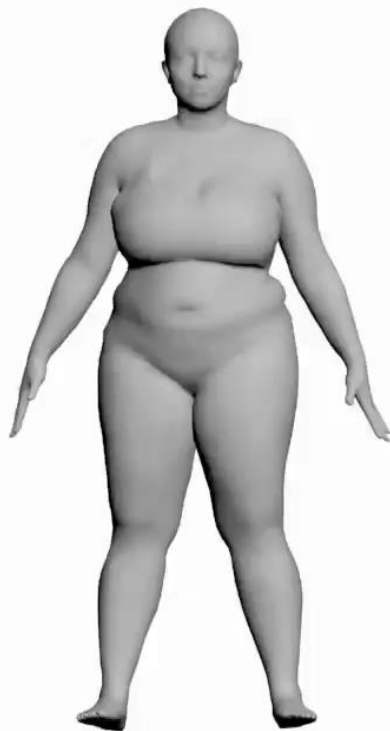
[CO18].

=> much more efficient to train!!

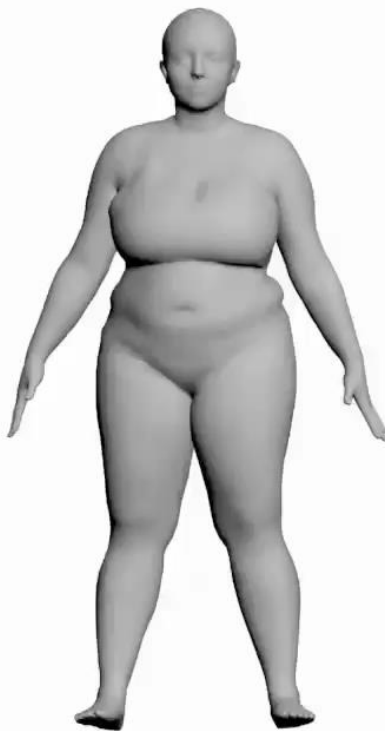
[CO18] Casas, D. & Otaduy, M. (2018). Learning Nonlinear Soft-Tissue Dynamics for Interactive Avatars. Proceedings of the ACM on Computer Graphics and Interactive Techniques

DS-Net: AE results

Reconstruction results: min 0 cm, max 1.033 cm



ground-truth



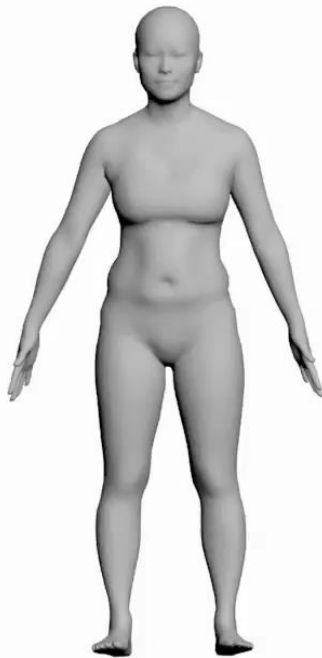
network output



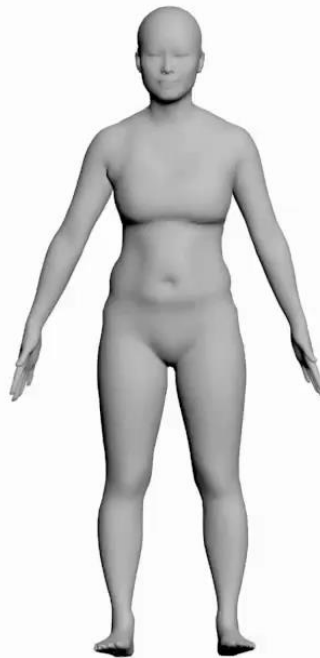
per-vertex error

DS-Net: AE results

Reconstruction results: min 0 cm, max 1.000 cm



ground-truth



network output

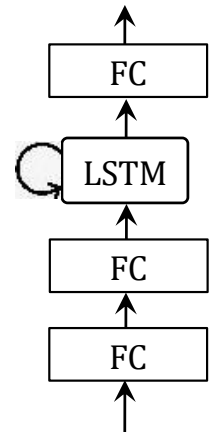


per-vertex error

Implementation details

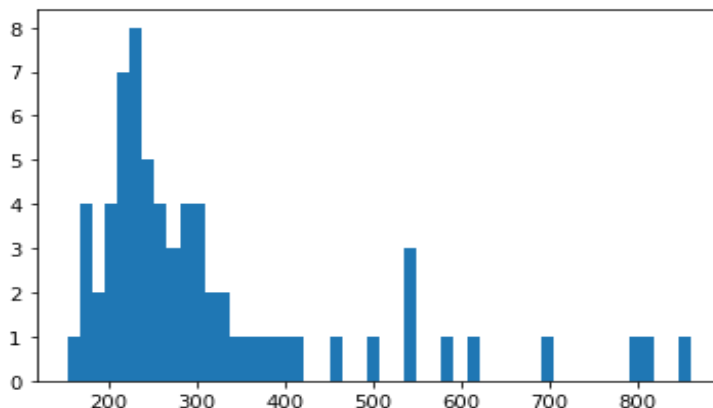
Implementation details

- Tensorflow 2.0 implementation of Adam optimizer
- 3rd dimensions of output vectors: 64, 128, 60, 100
- Activation functions: linear, tanh, (bath normalization), linear
- Batch size=16, lr= 0.0001.
- 0.05 sec/epoch on an Ubuntu machine with Nvidia GeForce RTX 2080 Super



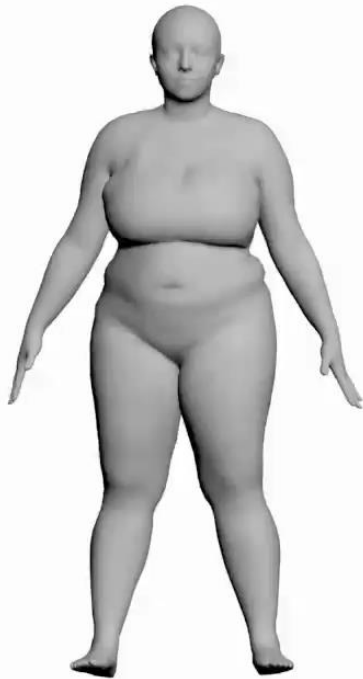
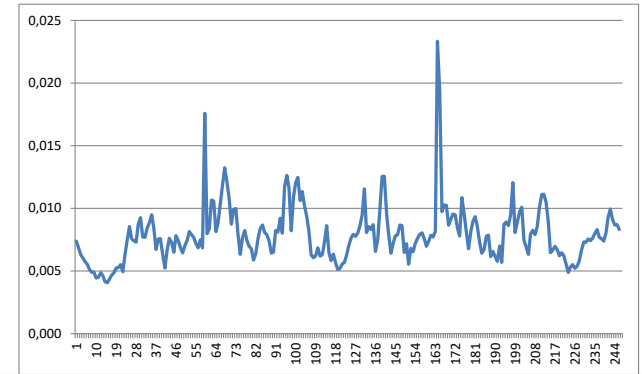
Data preprocessing

- Uniformize the frame lengths (to 300) by zero padding or tail clipping.

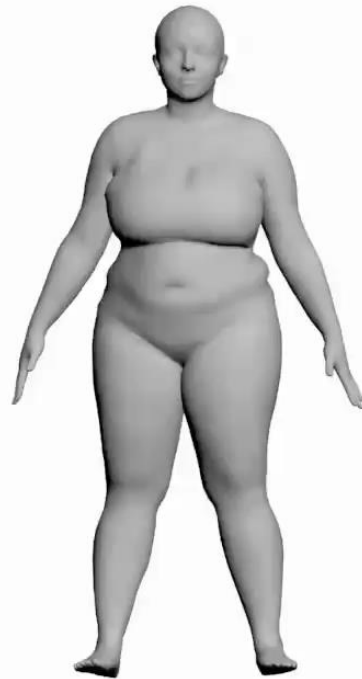


DS-Net: Prediction results

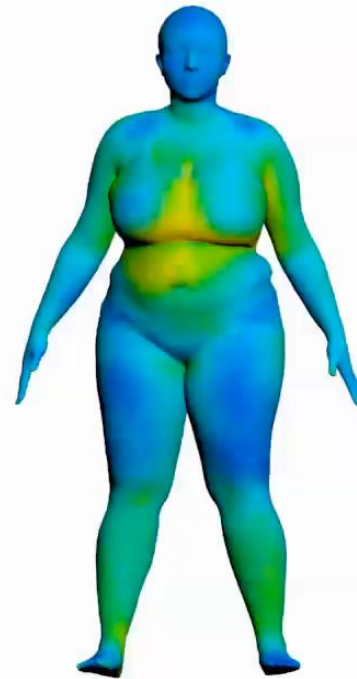
On validation data :



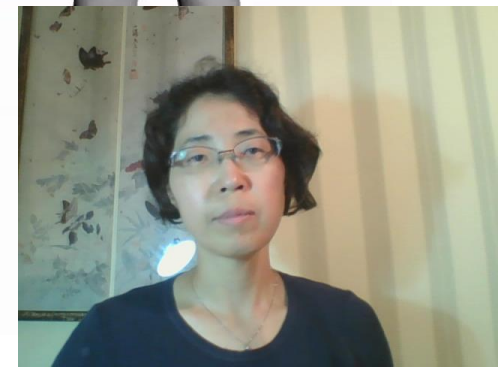
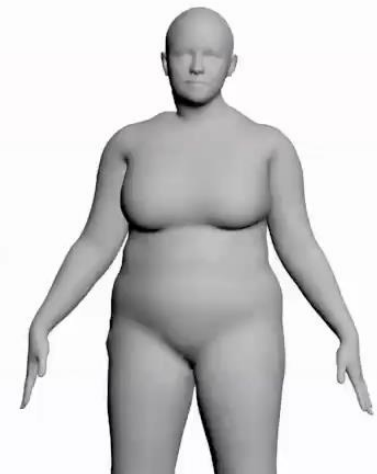
ground-truth



SMPL(β, θ) +
DSNet(β, θ)

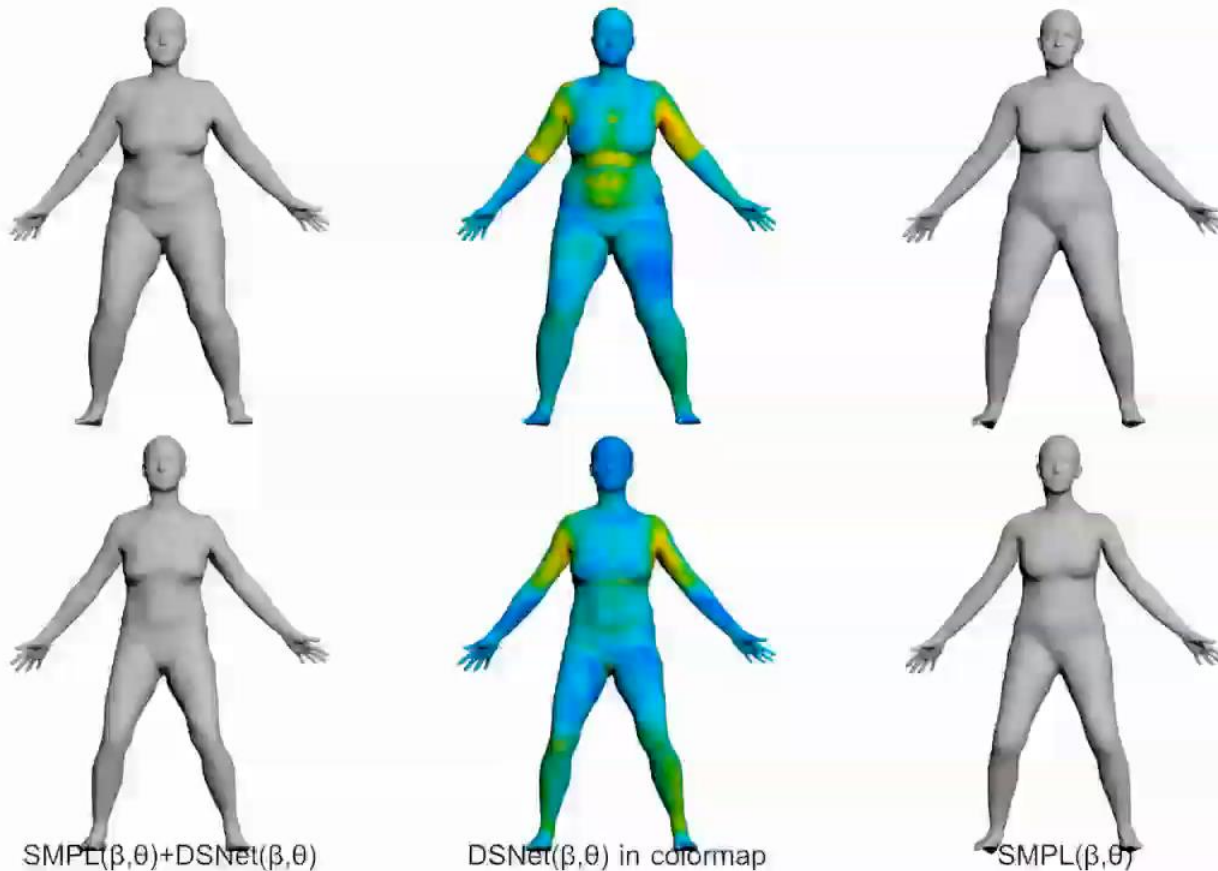


DSNet(β, θ)
in colormap



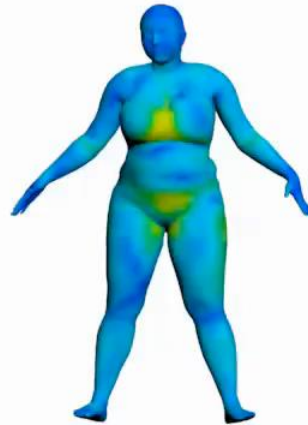
DS-Net: Prediction results

On validation data :



DS-Net: Prediction results

On unseen motions :



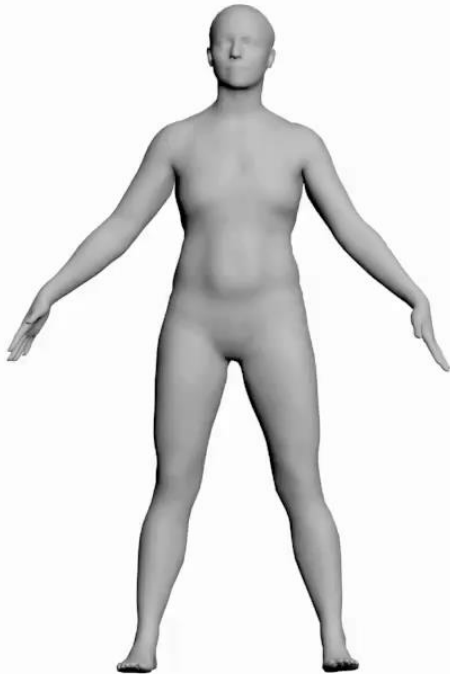
SMPL(β, θ)+DSNet(β, θ)

DSNet(β, θ) in colormap

SMPL(β, θ)

DS-Net: Prediction results

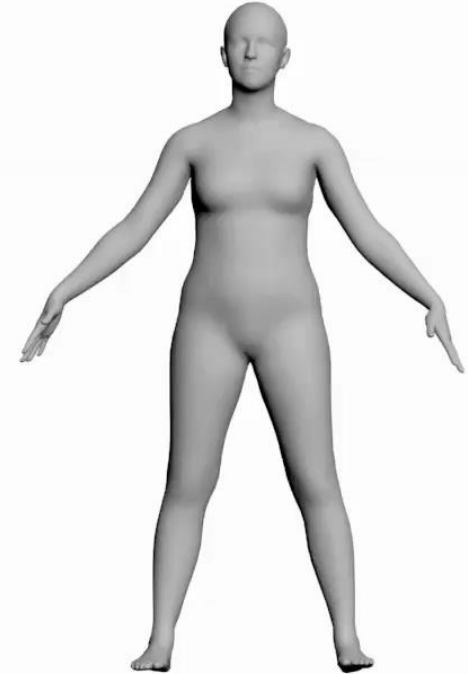
On unseen motions & unseen subjects:



SMPL(β, θ) +
DSNet(β, θ)



DSNet(β, θ)
in colormap



SMPL(β, θ)

Conclusion

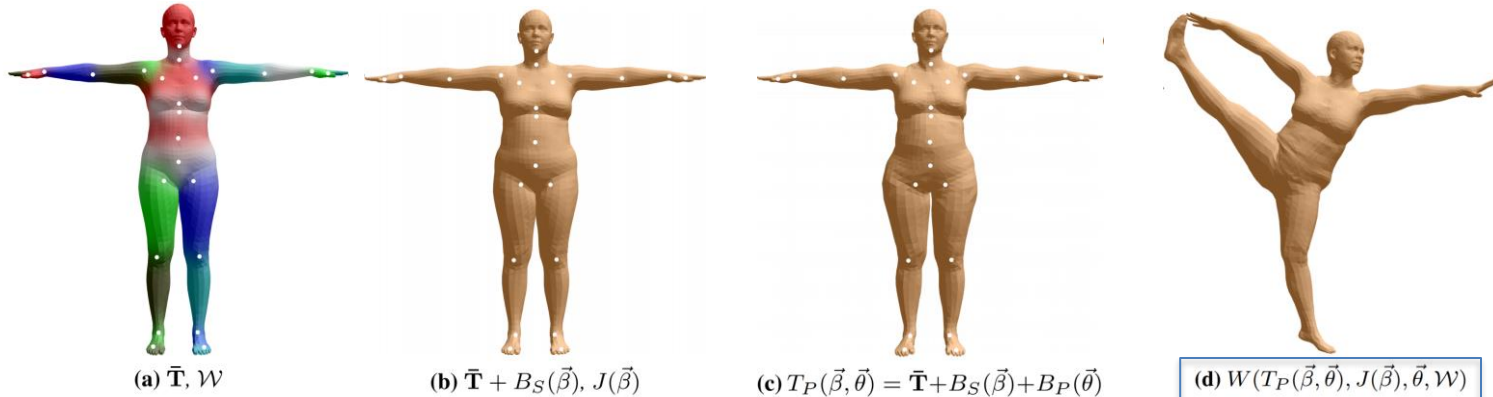
- A learning based method to the estimation of quality dynamic skin deformation.
- The dynamic skin deformation has been modeled as a time series data, as a function of pose, body shape, and the results of previous time steps.
 - => An LSTM based NN has been developed, trained on sequences of triangular meshes captured from real people.
- Also developed has been an AE, which builds a compact space for the intrinsic representation of skin displacement, allowing a very efficient operation of the DSNet.

Thank you!

Acknowledgement: ANR Human4D (ANR-19-CE23-0020) by the French Agence Nationale de la Recherche

DS-Net: Body model

SMPL: A Skinned Multi-Person Linear Model [LMRP+15]



$$M(\boldsymbol{\beta}, \boldsymbol{\theta}) = \mathcal{W}(\boldsymbol{\theta}, \bar{M}(\boldsymbol{\beta}, \boldsymbol{\theta}), J(\boldsymbol{\beta}), \mathbf{W})$$

linear blend skinning

$$\bar{M}(\boldsymbol{\beta}, \boldsymbol{\theta}) = \bar{\mathbf{T}} + M_S(\boldsymbol{\beta}) + M_P(\boldsymbol{\theta})$$

Template model Shape blend shape Pose blend shape

$$M_S(\boldsymbol{\beta}) = \boldsymbol{\mu}_S + \sum_{n=1}^{|\boldsymbol{\beta}|} \beta_n \mathbf{s}_n$$

$$M_P(\boldsymbol{\theta}) = \sum_{n=1}^{9K} (R_n(\boldsymbol{\theta}) - R_n(\boldsymbol{\theta}^0)) \mathbf{P}_n$$

DS-Net : LSTM

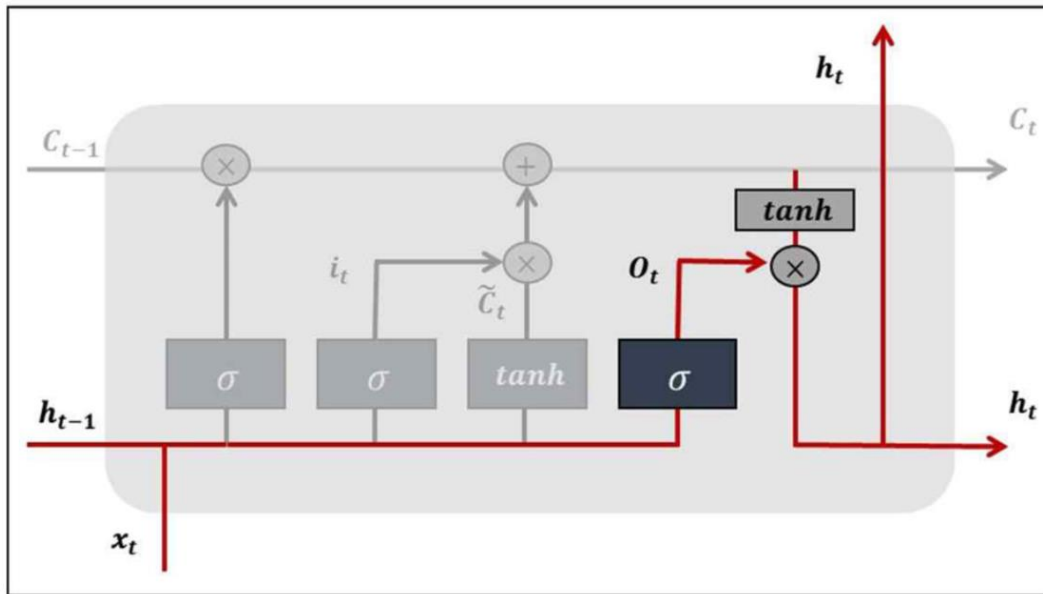
1

Long Short Term Memory network [HS97]

- It's an RNN, network with recurrent edges
- One or more layer is connected to itself
 - Self connections allow the network to build an **internal representation** of past inputs
 - In effect they serve as network **memory**

Our function

$$\Delta_t = f(x_t, f(x_{t-1}))$$



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

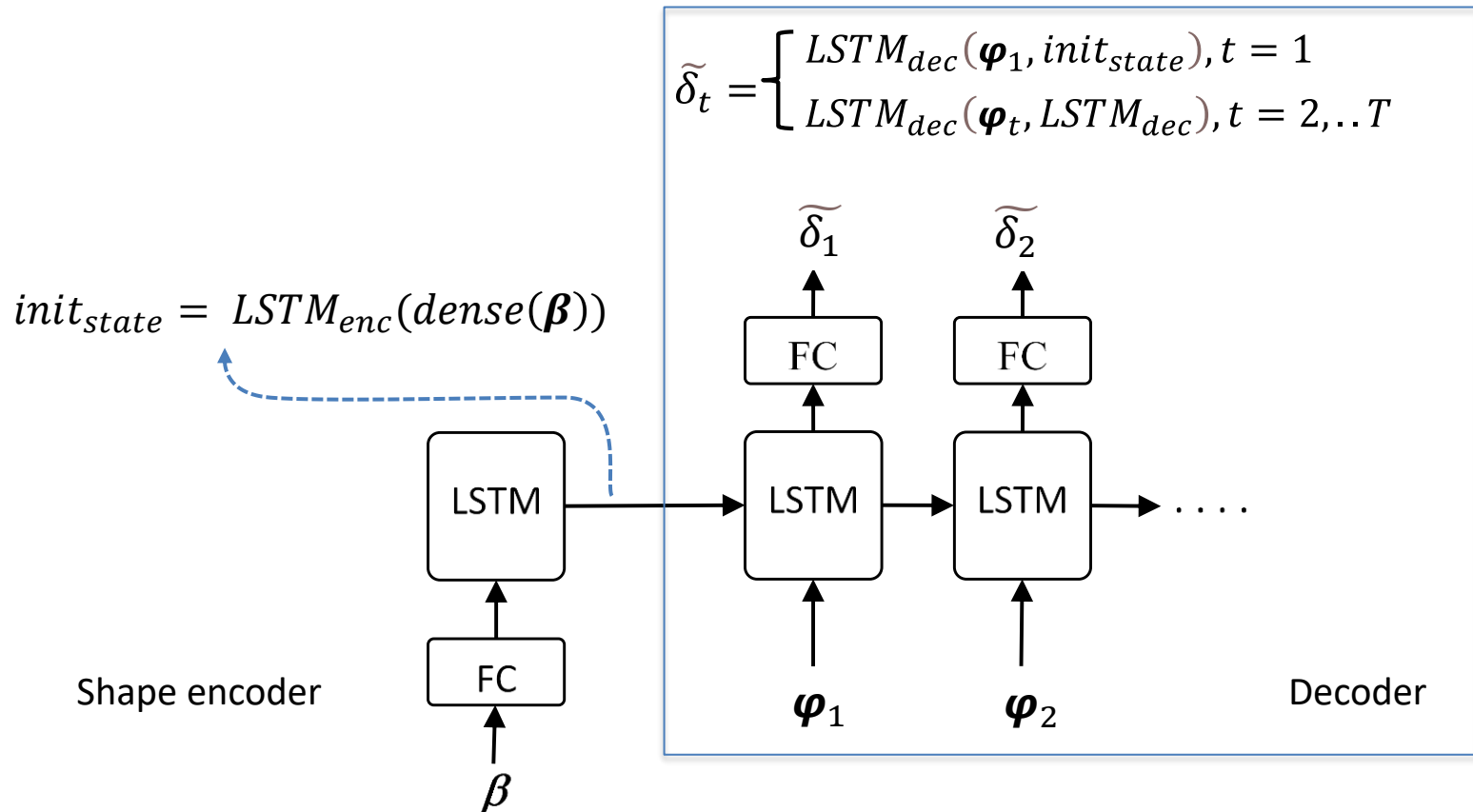
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

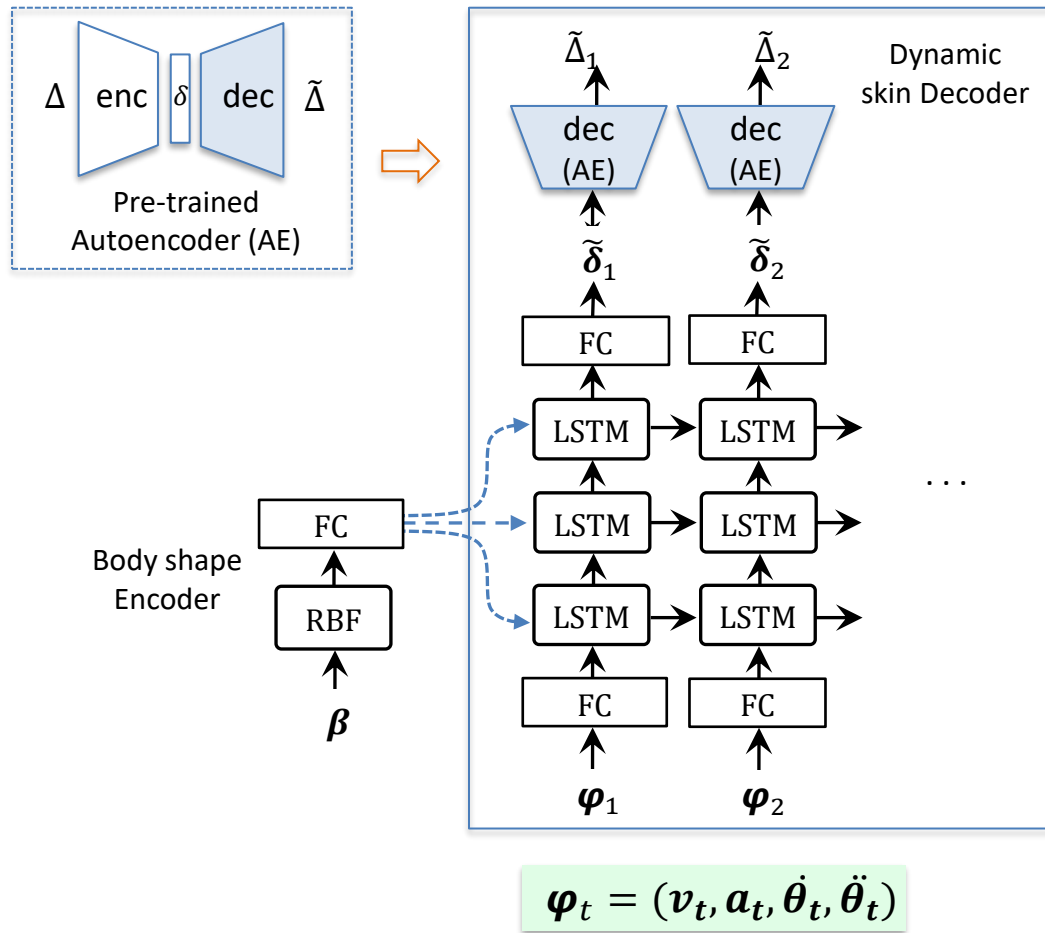
DS-Net: Architecture

DSNet: Earlier versions II



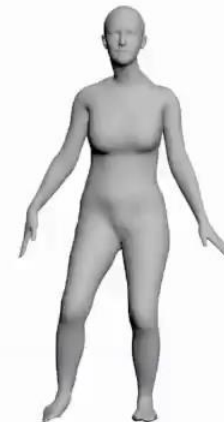
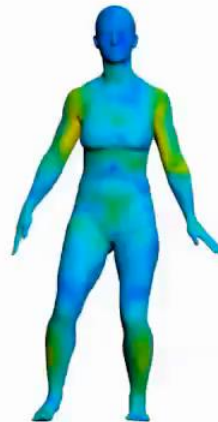
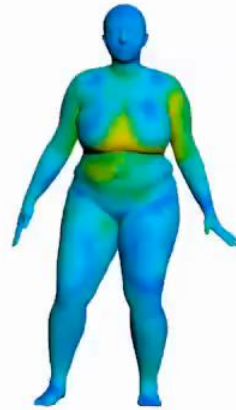
DS-Net: Architecture

DSNet: Earlier versions I



DS-Net: Prediction results

On validation data :



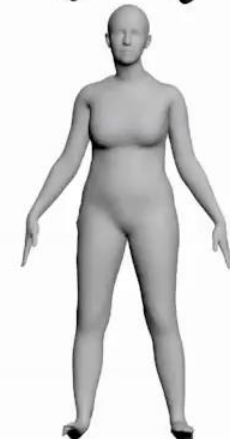
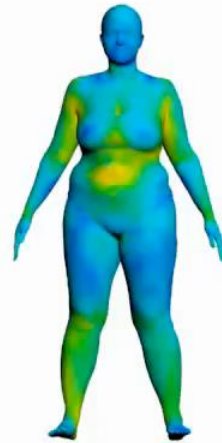
SMPL(β, θ)+DSNet(β, θ)

DSNet(β, θ) in colormap

SMPL(β, θ)

DS-Net: Prediction results

On validation data :



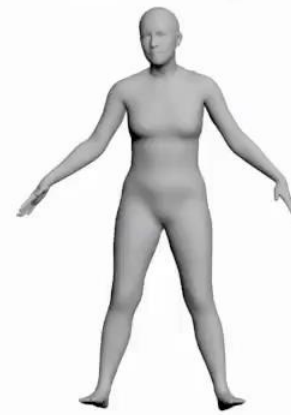
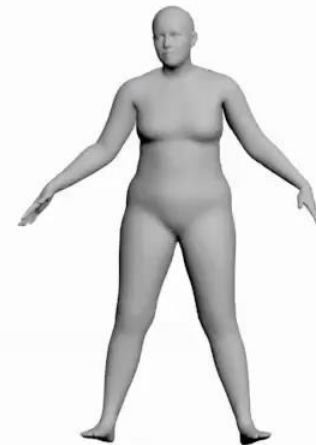
SMPL(β, θ)+DSNet(β, θ)

DSNet(β, θ) in colormap

SMPL(β, θ)

DS-Net: Prediction results

On unseen motions :



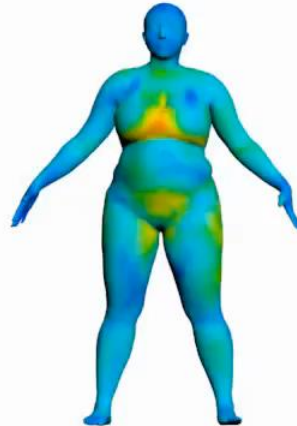
SMPL(β, θ)+DSNet(β, θ)

DSNet(β, θ) in colormap

SMPL(β, θ)

DS-Net: Prediction results

On unseen motions :



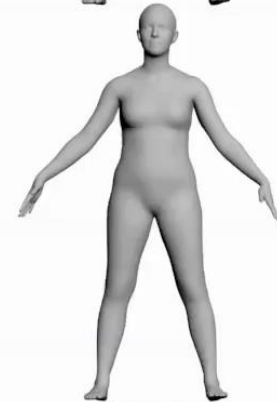
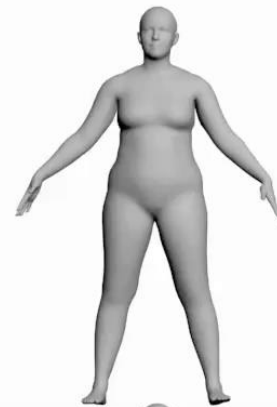
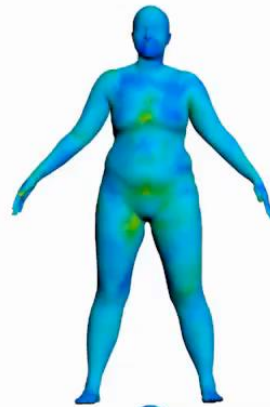
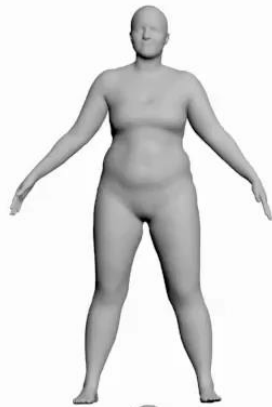
SMPL(β, θ)+DSNet(β, θ)

DSNet(β, θ) in colormap

SMPL(β, θ)

DS-Net: Prediction results

On unseen motions :



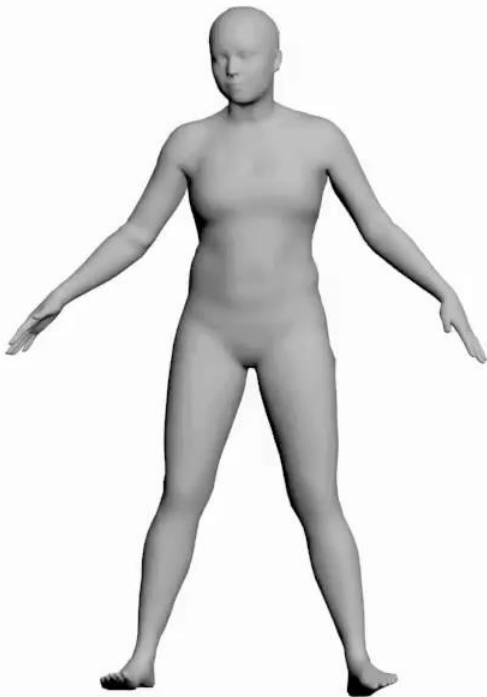
SMPL(β, θ)+DSNet(β, θ)

DSNet(β, θ) in colormap

SMPL(β, θ)

DS-Net: Prediction results

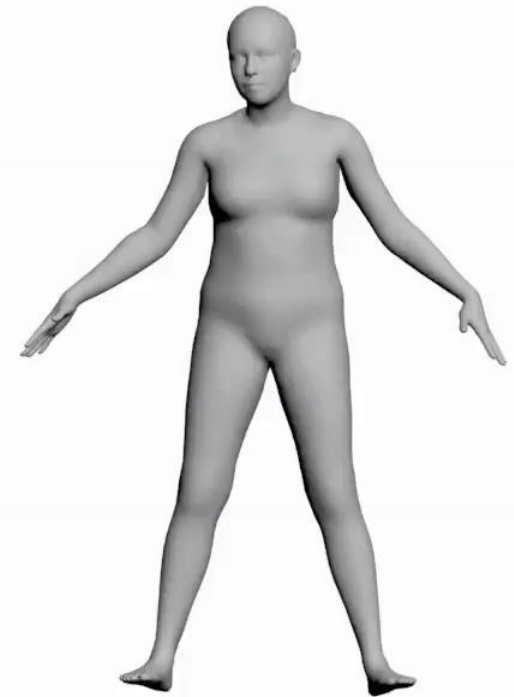
On unseen motions & unseen subjects:



SMPL(β, θ)+
DSNet(β, θ)



DSNet(β, θ)
in colormap

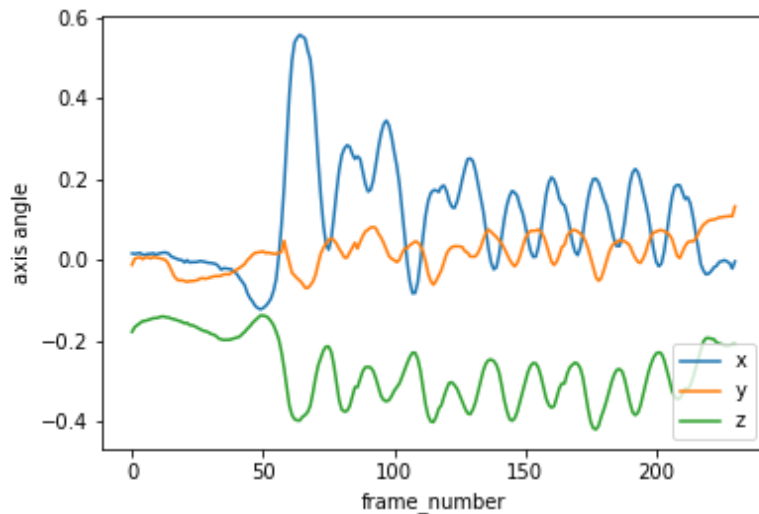


SMPL(β, θ)

Conclusion

A note on the training data

- We observed that the dynamics dependent shapes had been partly absorbed by the pose-dependent shape..!!



- 'spine 2' joint angles during 'Jiggling on toes' motion

- This means that our training data do not fully capture the observed dynamics...