

# Supplementary Document

## MoSculp: Interactive Visualization of Shape and Time

### HIDDEN MARKOV MODEL IN KEYPOINT DETECTION

As mentioned in the paper [3], we enforce temporal smoothness in the 2D keypoint detections across all frames using a Hidden Markov Model (HMM).

Specifically, given  $T$  video frames  $\{x^1, \dots, x^T\}$ , we compute the maximum marginal likelihood estimate of joint  $i$ 's pixel location at time  $t$ , denoted by  $y_i^t$ , as

$$\arg \max_{y_i^t} p(y_i^t | x^1, \dots, x^t) \propto p(y_i^t, x^1, \dots, x^t) \quad (1)$$

$$= \int_{y_i^{t-1}} p(y_i^{t-1}, y_i^t, x^1, \dots, x^t) dy_i^{t-1} \quad (2)$$

$$= \int_{y_i^{t-1}} p(y_i^t, x^t | y_i^{t-1}, x^1, \dots, x^{t-1}) \quad (3)$$

$$p(y_i^{t-1}, x^1, \dots, x^{t-1}) dy_i^{t-1} \quad (4)$$

$$= \int_{y_i^{t-1}} p(y_i^t, x^t | y_i^{t-1}) \quad (5)$$

$$p(y_i^{t-1}, x^1, \dots, x^{t-1}) dy_i^{t-1} \quad (6)$$

$$= \int_{y_i^{t-1}} p(y_i^t | y_i^{t-1}) p(x^t | y_i^t) \quad (7)$$

$$p(y_i^{t-1}, x^1, \dots, x^{t-1}) dy_i^{t-1}, \quad (8)$$

where the transition probability  $p(y_i^t | y_i^{t-1})$  is a bivariate Gaussian centered at  $y_i^{t-1}$  with a standard deviation of three pixels, and the emission probability  $p(x^t | y_i^t)$  is approximated as the heatmap predicted for joint  $i$ , independently at frame  $t$ . Strictly speaking, this heatmap resembles more the posterior  $p(y_i^t | x^t)$ , *i.e.*, how likely joint  $i$  lies at each pixel location given frame  $t$ . Considering the simple Bayes' rule

$$p(x^t | y_i^t) = \frac{p(y_i^t | x^t) p(x^t)}{p(y_i^t)},$$

by using the heatmaps as  $p(x^t | y_i^t)$ , we essentially assume a uniform  $p(y_i^t)$ , the prior distribution on joint  $i$ 's location.

Notice how  $p(y_i^t, x^1, \dots, x^t)$  in Equation 1 gets expressed in terms of  $p(y_i^{t-1}, x^1, \dots, x^{t-1})$  in Equation 8. This recurrent structure allows us to solve the HMM efficiently by message passing. More specifically, to find the  $y_i^t$  that maximizes  $p(y_i^t, x^1, \dots, x^t)$ , we integrate (spatially) over frame  $t$  the product of the Gaussian distribution for smoothness, the predicted heatmap for this joint in frame  $t$ , and where this joint is believed to be in frame  $t - 1$ , *i.e.*,  $p(y_i^{t-1}, x^1, \dots, x^{t-1})$ .

In cases where some of the person's joints are not detected locally, we linearly interpolate their locations from neighboring frames before running the HMM.

### JOINT OPTIMIZATION OF SHAPE, POSE, AND GLOBAL TRAJECTORY OVER TIME

Here, we provide the mathematical and implementation details of our formulation.

$$\begin{aligned} \mathcal{L}(\{T^t\}, \{\theta^t\}, \beta) &= \sum_{t=1}^N \mathcal{L}_{\text{data}}(T^t, \theta^t, \beta) + \alpha_1 \mathcal{L}_{\text{prior}}(\theta^t, \beta) \\ &+ \alpha_2 \sum_{t=1}^{N-1} \mathcal{L}_{\text{temporal}}(T^t, T^{t+1}, \theta^t, \theta^{t+1}, \beta). \end{aligned} \quad (9)$$

As mentioned in the paper, this formulation can be seen as an extension of SMPLify [1], a single-image 3D human pose and shape estimation algorithm, to videos. Therefore, we will elaborate only on the newly added term—the temporal smoothness prior  $\mathcal{L}_{\text{temporal}}$ .

This term encourages the 3D model to be temporally smooth; it penalizes changes in the human's global translations (Equation 11), local vertex locations (Equation 12), and pose parameters (Equation 13). More specifically,

$$\mathcal{L}_{\text{temporal}} = \lambda_1 \mathcal{L}_{\text{global}} + \lambda_2 \mathcal{L}_{\text{local}} + \lambda_3 \mathcal{L}_{\text{rotation}} \quad (10)$$

$$\mathcal{L}_{\text{global}}(T^t, T^{t+1}) = \|T^t - T^{t+1}\|_2^2 \quad (11)$$

$$\mathcal{L}_{\text{local}}(\theta^t, \theta^{t+1}, \beta) = \|V(\theta^t, \beta) - V(\theta^{t+1}, \beta)\|_F^2 \quad (12)$$

$$\mathcal{L}_{\text{rotation}}(\theta^t, \theta^{t+1}) = \left\| \begin{bmatrix} \cos(\theta^t) - \cos(\theta^{t+1}) \\ \sin(\theta^t) - \sin(\theta^{t+1}) \end{bmatrix} \right\|_2^2, \quad (13)$$

where  $V(\cdot)$  are the vertices' local 3D coordinates given the pose and shape, and the  $\lambda$ 's are constant weights that roughly match the orders of magnitude of the three losses.

Intuitively,  $\mathcal{L}_{\text{global}}$  requires the body's global trajectory to be smooth;  $\mathcal{L}_{\text{local}}$  further requires vertices of the human mesh to translate smoothly; and  $\mathcal{L}_{\text{rotation}}$  imposes additional rotational

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

UIST '18 October 14–17, 2018, Berlin, Germany

© 2018 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-5948-1/18/10.

DOI: <https://doi.org/10.1145/3242587.3242592>

smoothness in the parameter space, which is necessary for producing natural pose evolution.

### Optimization

We adopt a two-stage optimization procedure, in which we first ignore the temporal loss, optimizing only the per-frame loss:  $\mathcal{L}_{\text{data}} + \alpha_1 \mathcal{L}_{\text{spatial}}$ . This runs in a fully parallelized fashion and provides good initializations to, thereby speeding up, the subsequent joint optimization.

More importantly, this two-step procedure allows us to address effectively the “pose-flipping problem”: tendency of the joint optimization getting stuck with flipped facing directions when the person is captured in a side view throughout, often due to the inherent ambiguity (see *Run, Forrest, Run!* for an example). To avoid such local minima, the algorithm tries both directions that the human model could face, when the left-right shoulder or hip joints are less than 100 pixels apart (a heuristic), and then initializes the joint optimization with the pose that gives a lower  $\mathcal{L}_{\text{prior}}$ . We minimize this loss using a non-linear least squares approach [2].

### REFERENCES

1. Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J. Black. 2016. Keep It SMPL: Automatic Estimation of 3D Human Pose and Shape from a Single Image. In *European Conference on Computer Vision*.
2. Matthew M. Loper and Michael J. Black. 2014. OpenDR: An Approximate Differentiable Renderer. In *European Conference on Computer Vision*.
3. Xiuming Zhang, Tali Dekel, Tianfan Xue, Andrew Owens, Qiurui He, Jiajun Wu, Stefanie Mueller, and William T. Freeman. 2018. MoSculp: Interactive Visualization of Shape and Time. In *ACM Symposium on User Interface Software and Technology*.