

Robust Human Body Shape and Pose Tracking

Chun-Hao Huang¹ Edmond Boyer² Slobodan Ilic¹

¹ Technische Universität München

² INRIA Grenoble Rhône-Alpes

Marker-based motion capture (mocap.)

- Advantages:

- precision, reliability
- little data (couple of kB/frame for 50 cameras, 5 people)
- real-time processing, visualization & retargeting.



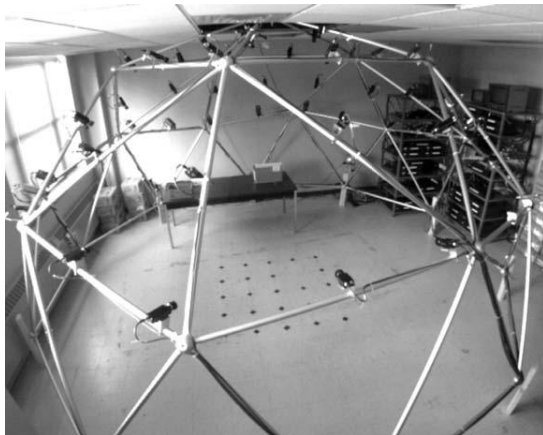
Giant Studios (L.A. Noire set)

- Disadvantages:

- Attaching, removing, re-attaching markers is tedious.
- Markers can interfere with the movement.
- Markers prevent the **simultaneous acquisition of shape and motion.**

Marker-less mocap.

- Multiple camera setup is usually required.
- Acquisition of both motion and shape.



3D dome (CMU)



Grimage (INRIA)

Motivation

- Methods that assume a skeleton usually produce **skinning artifacts**, require **2nd stage shape refinement**.

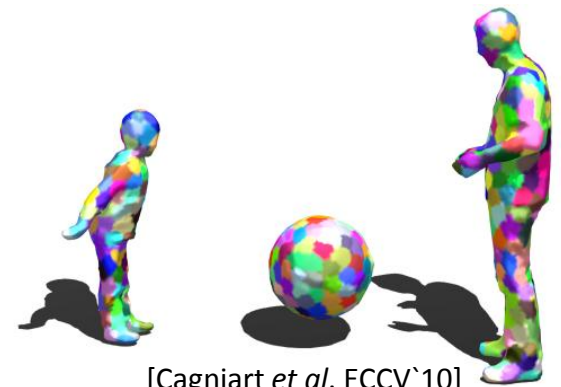
Degree of freedom (DoF): $N_J \times 6$ ($< 10^2$)

- Purely-surface-based methods handle non-rigid surface deformation better, but **do not provide the pose**.

DoF: $N_P \times 6$ ($\cong 10^3$)



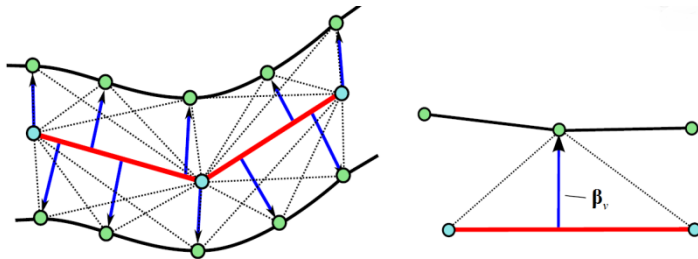
[Vlasic *et al.* ToG'08]
Skeleton-based



[Cagniard *et al.* ECCV'10]
Surface-based

Contribution

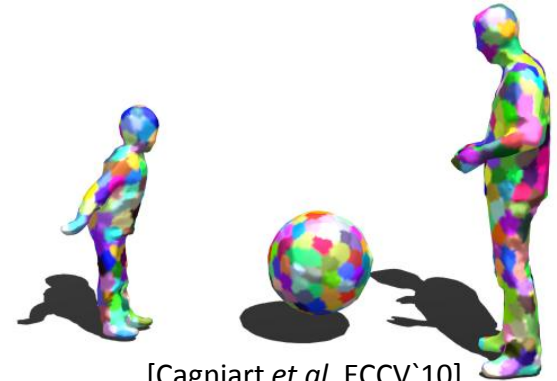
- bone differential coordinate



[Straka *et al.* ECCV'12]



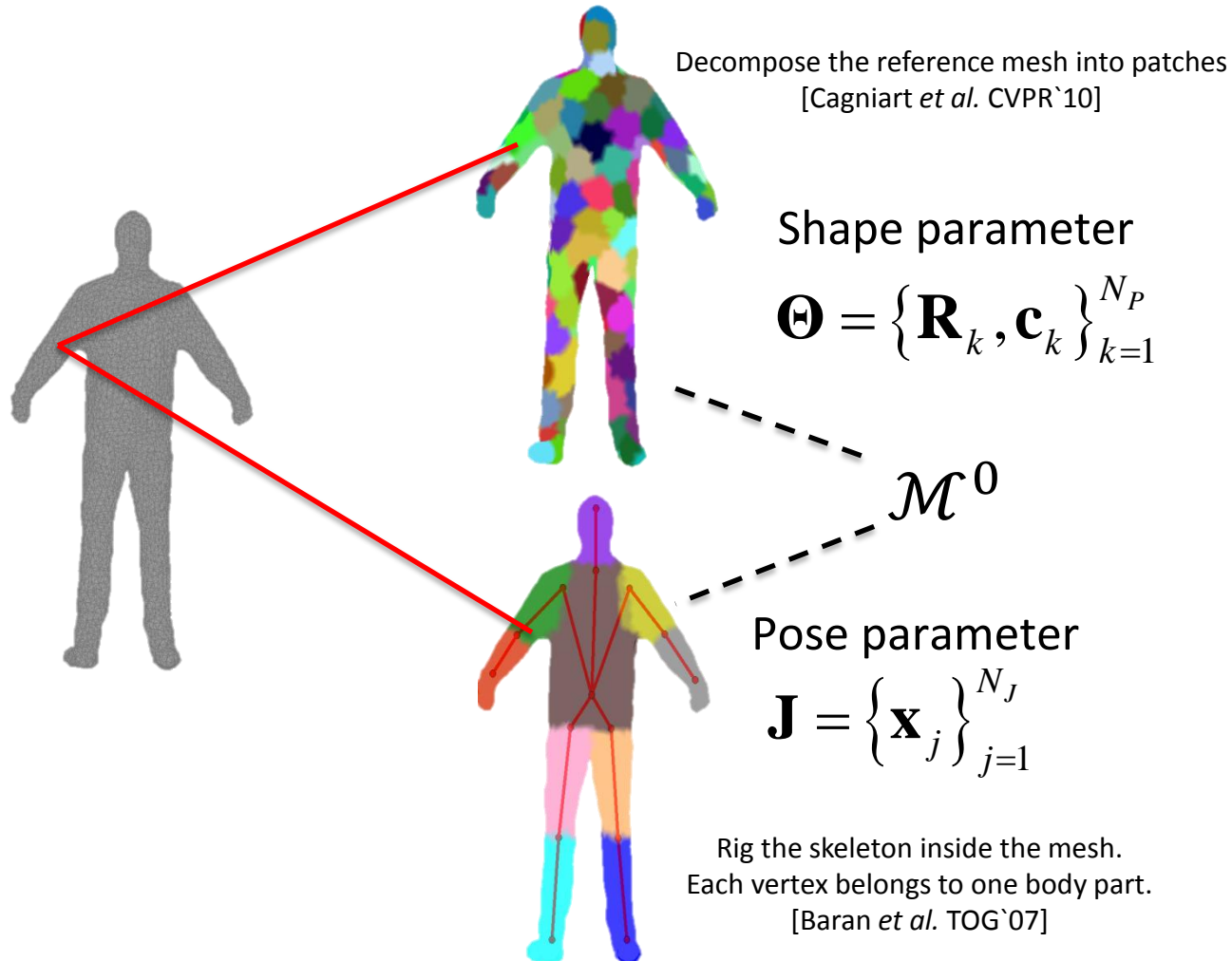
probabilistic surface deformation



[Cagniard *et al.* ECCV'10]

- A learning-based outlier rejection scheme.

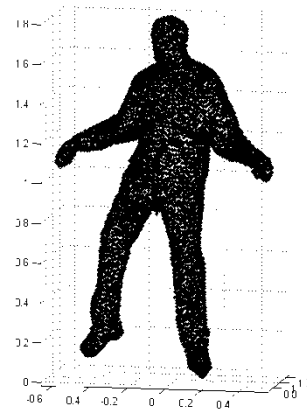
Preprocessing step: model



Preprocessing step: input data

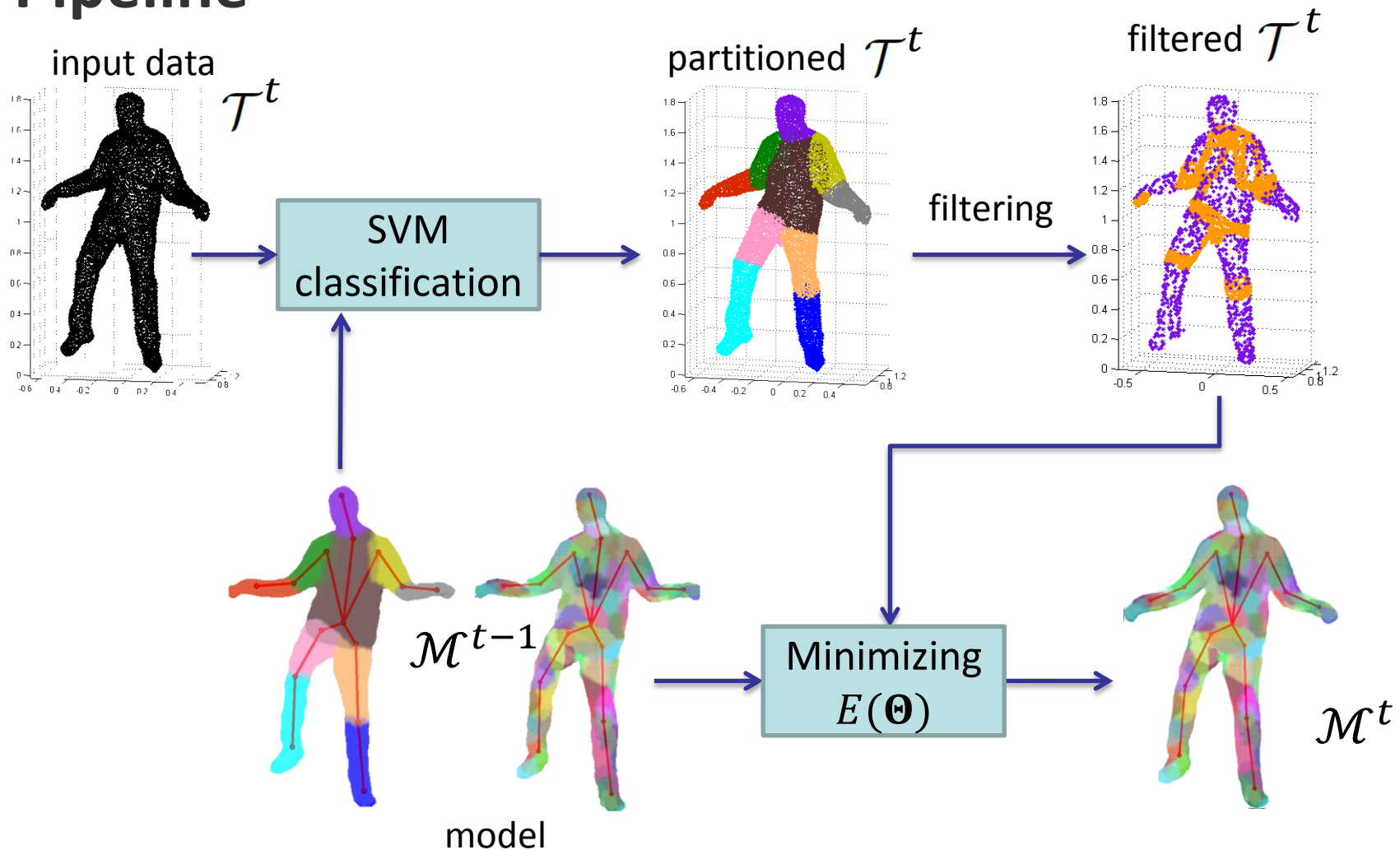


$$\mathcal{T} = \{y_i\}_{i=1}^{N_T}$$

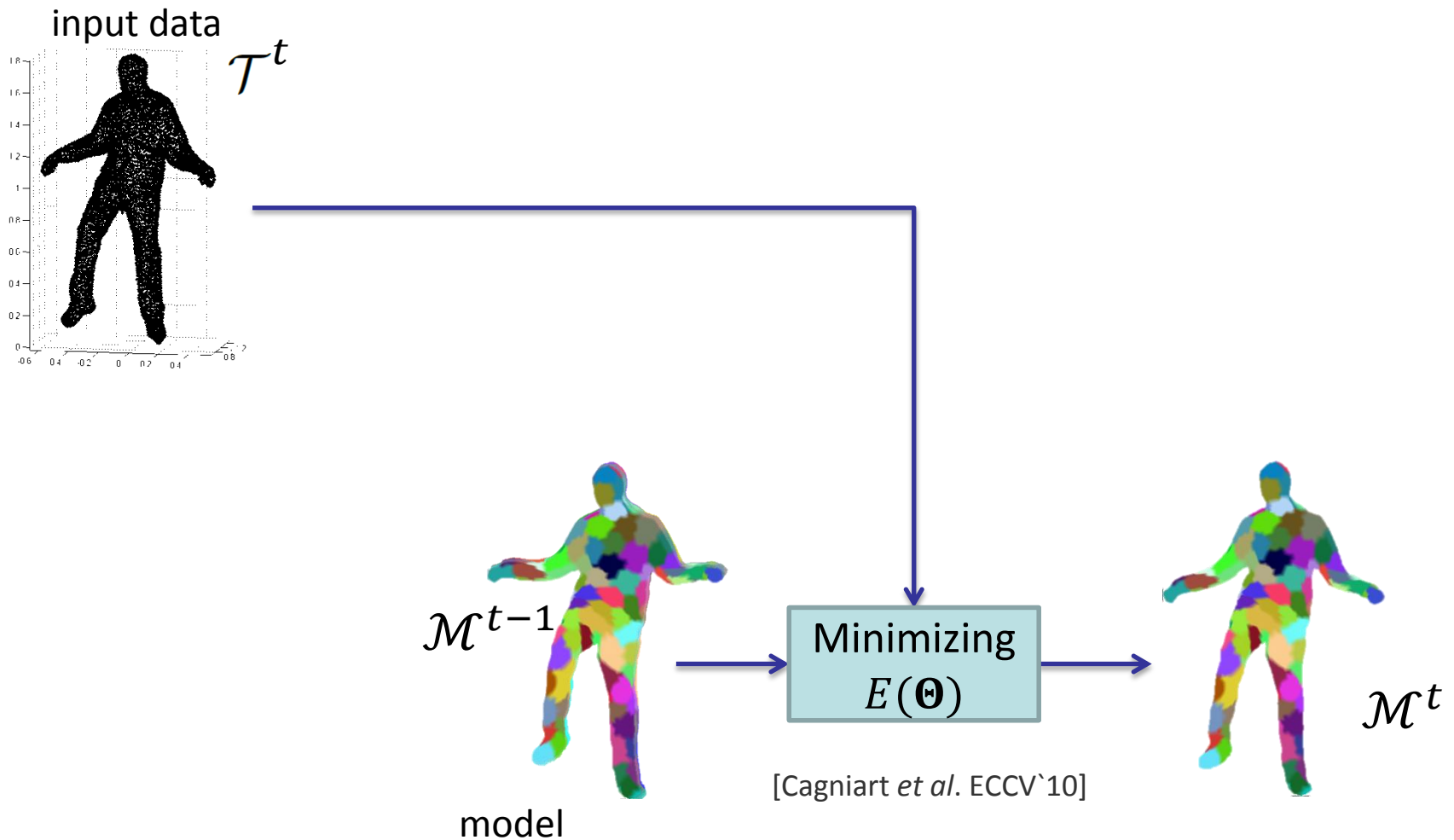


For each time stamp t , visual hull is reconstructed from silhouettes, which serves as our observations

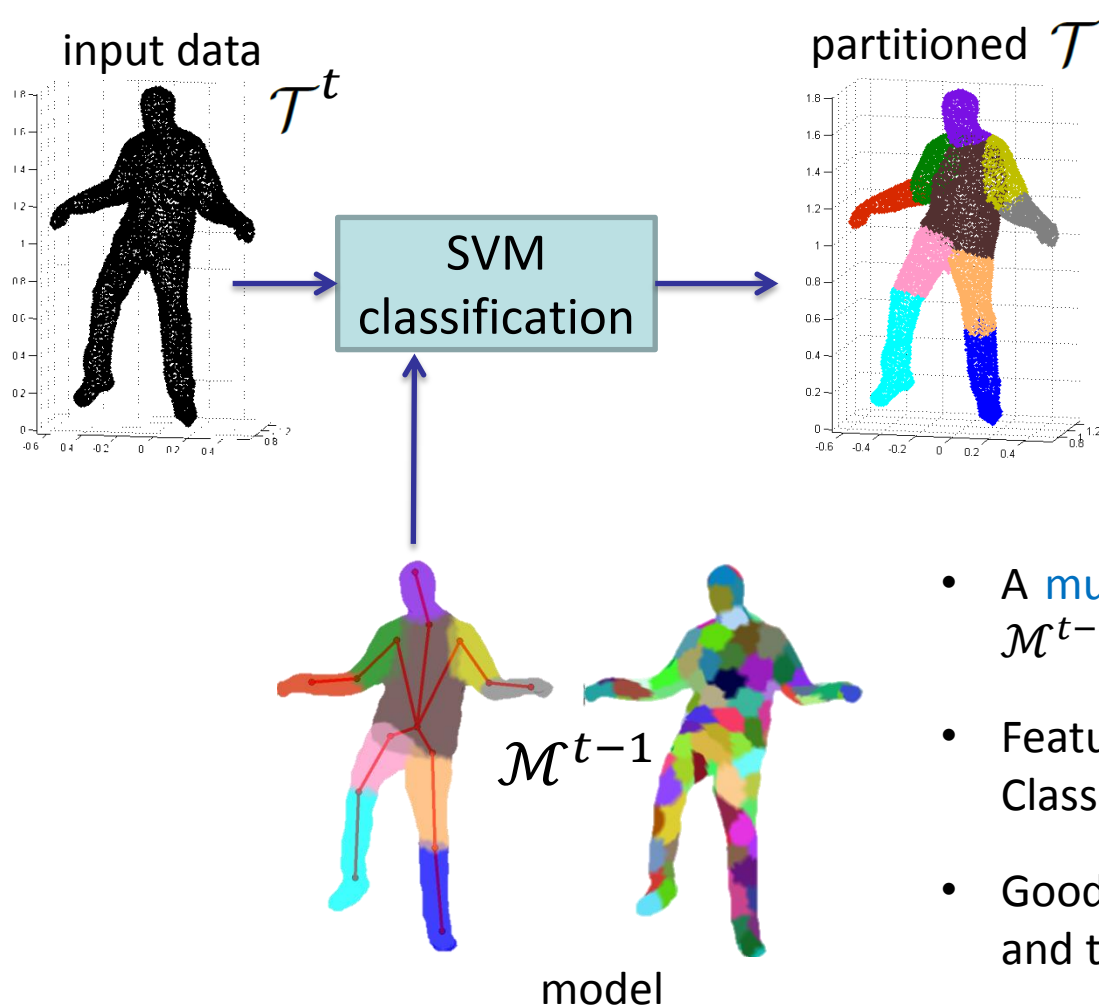
Pipeline



Pipeline



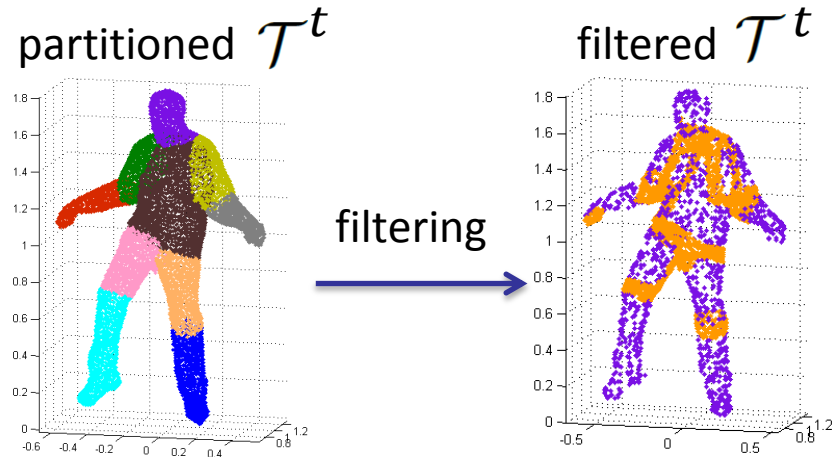
SVM-based body part classification



- A **multi-class linear SVM** is trained on \mathcal{M}^{t-1} and tested on \mathcal{T}^t
- Feature: 3D coordinate of vertices. Class label: **rigid body part label**.
- Good compromise between accuracy and training time.

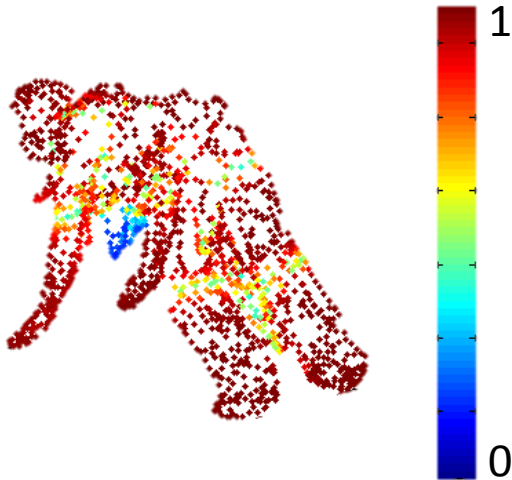
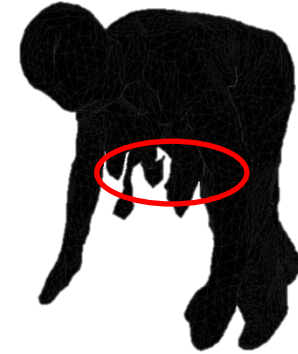
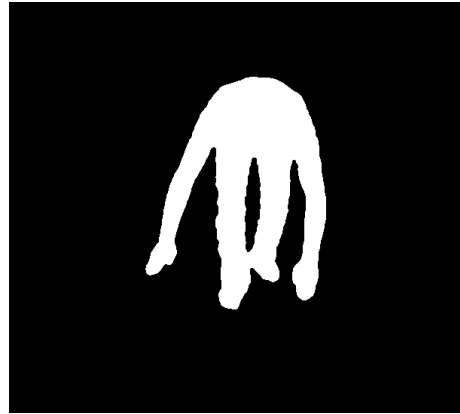
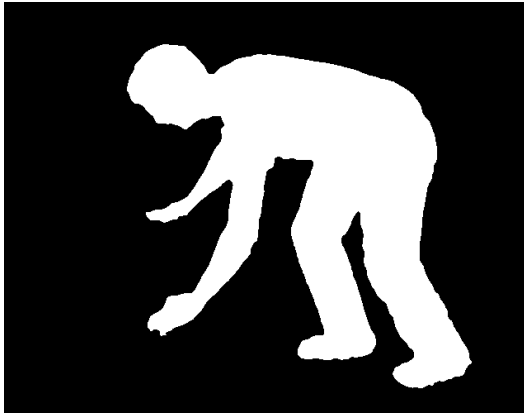
Filtering point cloud

- **Bone \mathcal{T}_b :**
patches on the bone often move rigidly together.
→ sub-sample the observations.
- **Joint \mathcal{T}_g :**
patches on the joint have non-rigid deformation.
→ keep all the observations.
- **Outlier \mathcal{T}_o :**
abandon all the observations.

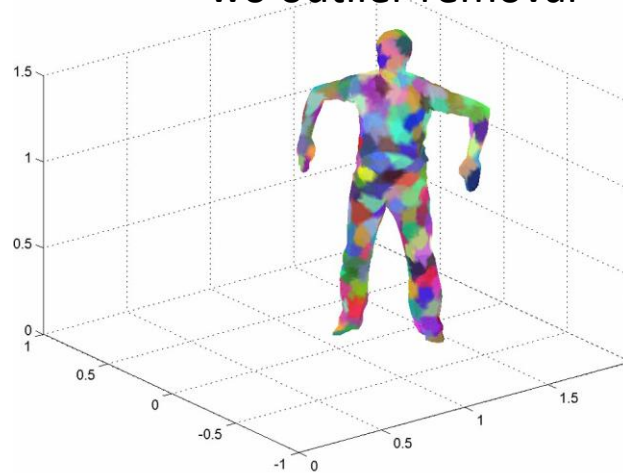


$$y_i \in \begin{cases} \mathcal{T}_b & \text{if } 0.9 < p_i \\ \mathcal{T}_g & \text{if } 0.5 < p_i \leq 0.9 \\ \mathcal{T}_o & \text{if } p_i \leq 0.5, \end{cases}$$

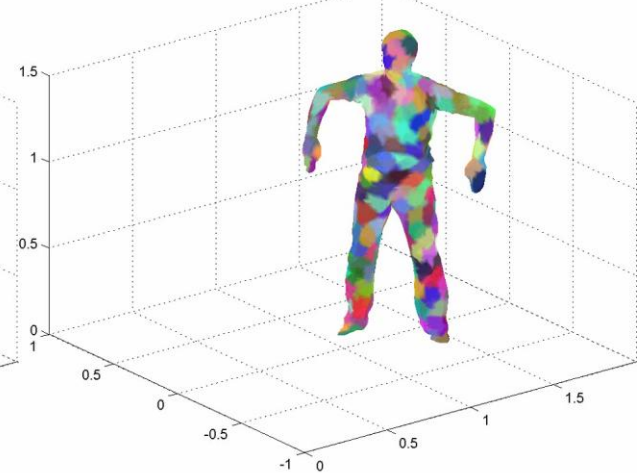
Benefit of SVM classification – outlier removal



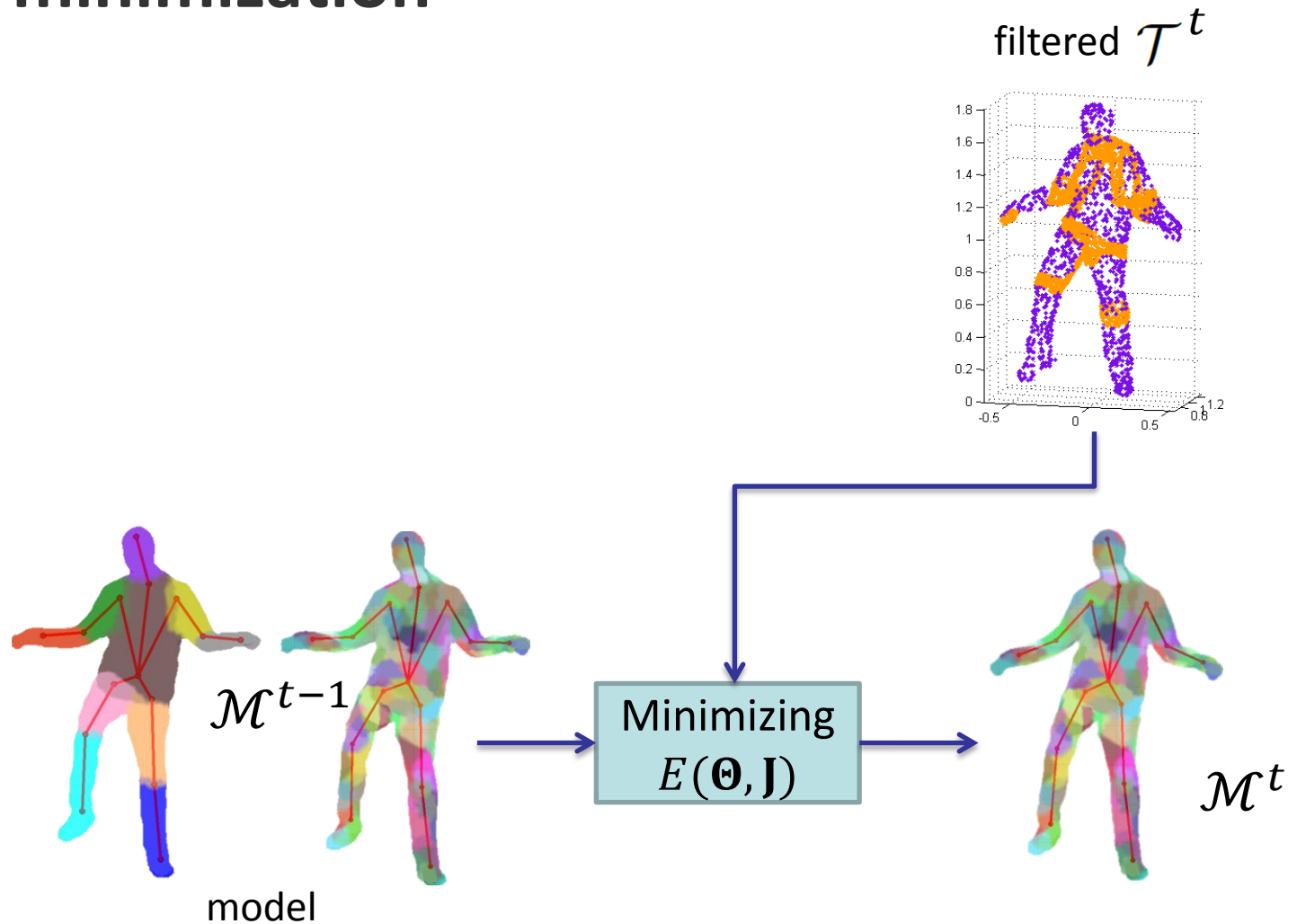
wo outlier removal



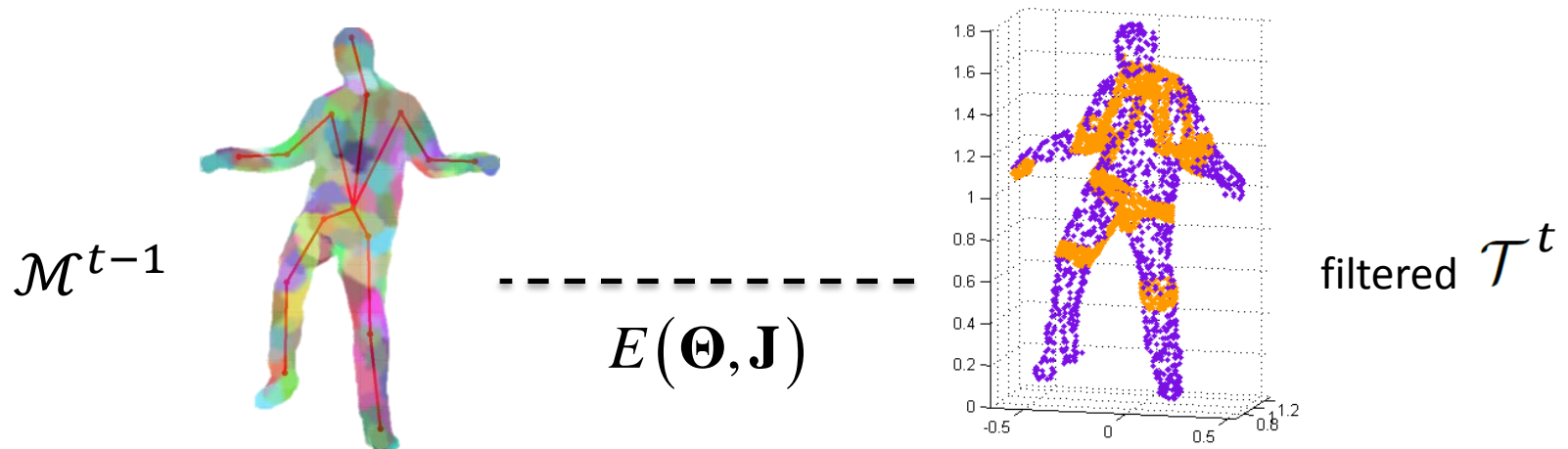
with outlier removal



Energy minimization



Energy function



$$\underline{E(\Theta, \mathbf{J})} = E_{\text{data}}(\Theta) + E_{\text{rigid}}(\Theta) + E_{\text{bone}}(\Theta, \mathbf{J})$$

- $E_{\text{data}}(\Theta)$: how well the surface explains the observations.
- $E_{\text{rigid}}(\Theta)$: smooth the motion of neighboring patches.
- $E_{\text{bone}}(\Theta, \mathbf{J})$: keep the relationship between the mesh and the skeleton.

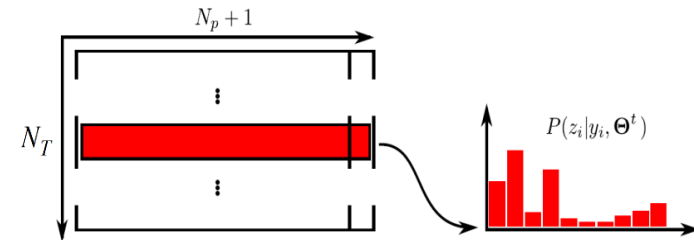
Data term

- $E_{\text{data}}(\Theta)$: how well the surface explains the observations.
- A probabilistic Iterative Closest Point (ICP) approach.
- Each observation has a **soft assignment** to every patch, **updated in each iteration**.
- Let observation i correspond to vertex v_i^k in P_k with a soft assignment w_i^k .

$$E_{\text{data}}(\Theta) = \sum_{i=1}^{N_T} \sum_{k=1}^{N_P+1} w_i^k \left\| \mathbf{y}_i - \mathbf{x}_{v_i^k} \right\|^2$$

A probabilistic point of view [Cagniard et al. ECCV'10]

- Can be interpreted as **EM algorithm**.
- The likelihood: Gaussian mixture model



$$P(\mathbf{y}_i | \Theta) = \sum_{k=1}^{N_p+1} \Pi_k \underline{P(\mathbf{y}_i | z_i = k, \Theta)}$$

- E-step: update the soft assignment.

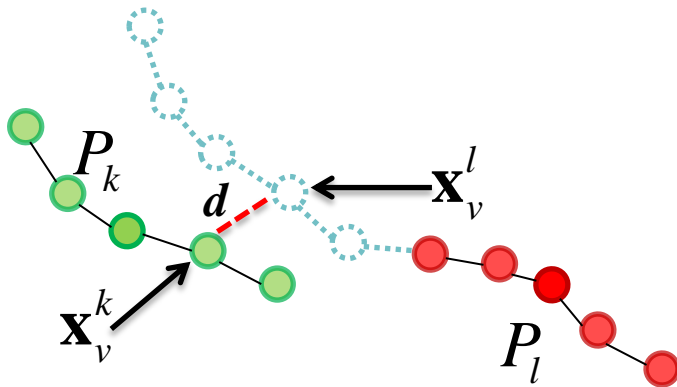
$$\underline{P(z_i = k | \mathbf{y}_i, \Theta)} = \frac{\Pi_k P(\mathbf{y}_i | z_i = k, \Theta)}{\sum_{l=1}^{N_p+1} P(\mathbf{y}_i | z_i = l, \Theta)} \quad -\ln(\cdot)$$

- M-step: minimize sum of negative log likelihood (energy).

$$E_{\text{data}}(\Theta) = \sum_{i=1}^{N_T} \sum_{k=1}^{N_p+1} w_i^k \left\| \mathbf{y}_i - \mathbf{x}_{v_i^k} \right\|^2$$

Rigidity energy [Cagniard *et al.* CVPR'10]

- $E_{\text{rigid}}(\Theta)$: smooth the motion of neighboring patches



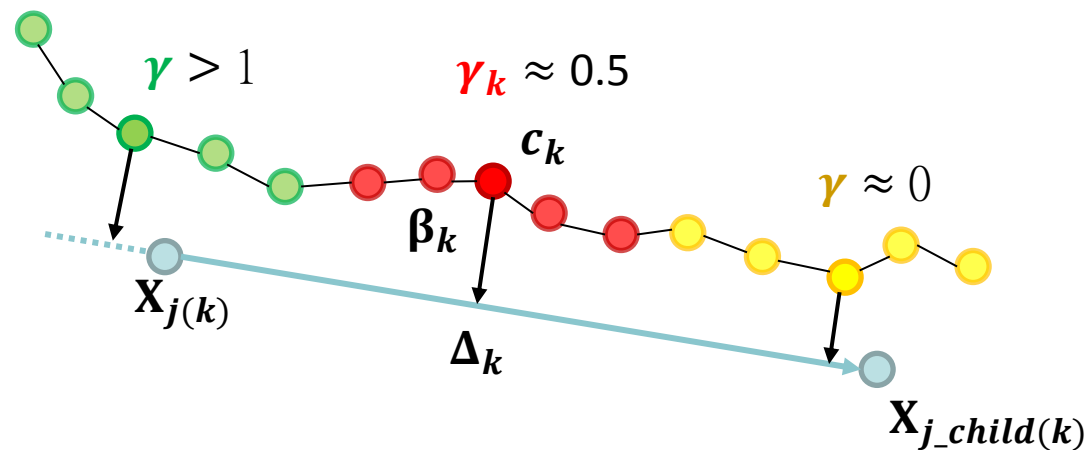
$$E_{\text{rigid}}(\Theta) = \sum_{k=1}^{N_P} \sum_{P_l \in N_k} \sum_{v \in P_k \cup P_l} \frac{\|\mathbf{x}_v^k - \mathbf{x}_v^l\|^2}{d}$$

For each patch, the real location and the predicted location should be consistent.

Θ is implicitly encoded in \mathbf{x}_v^k and \mathbf{x}_v^l

Bone-binding energy

- β coordinate : A **relative displacement** from patch to bone



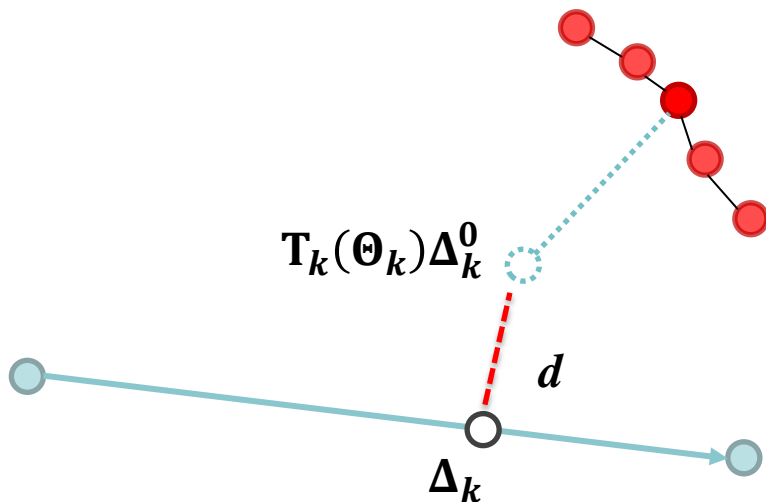
$$\beta_k = \Delta_k - c_k$$

$$\Delta_k = \gamma_k \mathbf{X}_{j(k)} + (1 - \gamma_k) \mathbf{X}_{j_child(k)}$$

γ_k is computed such that β_k is perpendicular to the bone.

Bone-binding energy

- $E_{\text{bone}}(\Theta, \mathbf{J})$: keep β consistent after transformation.



$$\begin{aligned}
 E_{\text{bone}}(\Theta, \mathbf{J}) &= \sum_{k=1}^{N_P} w_k \left\| \beta_k - T_k(\Theta_k) \beta_k^0 \right\|^2 \\
 &= \sum_{k=1}^{N_P} w_k \underbrace{\left\| \Delta_k - T_k(\Theta_k) \Delta_k^0 \right\|^2}_d
 \end{aligned}$$

For each patch, Δ_k predicted from the patch and Δ_k from the bone should be consistent.

w_k is weighed according to γ_k

Energy function

- $E_{\text{data}}(\Theta)$: how well the surface explains the observations.

$$E_{\text{data}}(\Theta) = \sum_{i=1}^{N_T} \sum_{k=1}^{N_P+1} w_i^k \left\| \mathbf{y}_i - \mathbf{x}_v^k \right\|^2$$

- $E_{\text{rigid}}(\Theta)$: smooth the motion of neighboring patches.

$$E_{\text{rigid}}(\Theta) = \sum_{k=1}^{N_P} \sum_{P_l \in N_k} \sum_{v \in P_k \cup P_l} \left\| \mathbf{x}_v^k - \mathbf{x}_v^l \right\|^2$$

- $E_{\text{bone}}(\Theta, \mathbf{J})$: keep the relationship between mesh

$$E_{\text{bone}}(\Theta, \mathbf{J}) = \sum_{k=1}^{N_P} w_k \left\| \Delta_k - T_k(\Theta_k) \Delta_k^0 \right\|^2$$

regularization terms
or
deformation prior

Minimizing the energy

$$E(\Theta, \mathbf{J}) = \lambda_d E_{\text{data}}(\Theta) + \lambda_r E_{\text{rigid}}(\Theta) + \lambda_b E_{\text{bone}}(\Theta, \mathbf{J})$$

- $\lambda_d = 10$, $\lambda_r = 1$, and $\lambda_b = 1$
- Each term is **quadratic** in terms of variables.
- Standard Gauss-Newton optimization is thus feasible.
- 3 - 4s per frame (including SVM training time).

Quantitative results

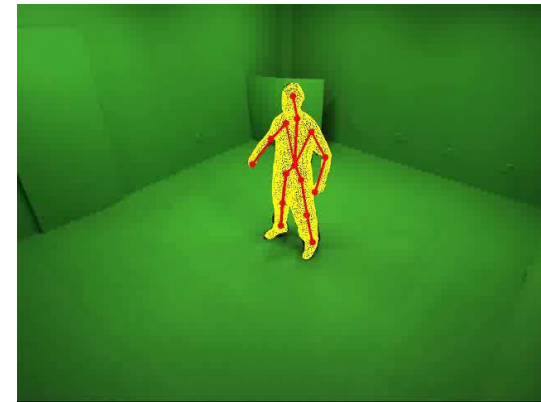
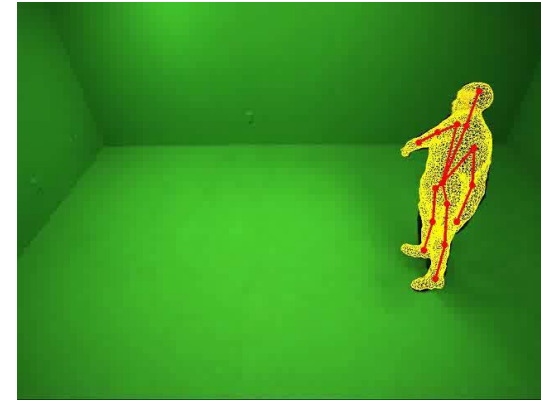
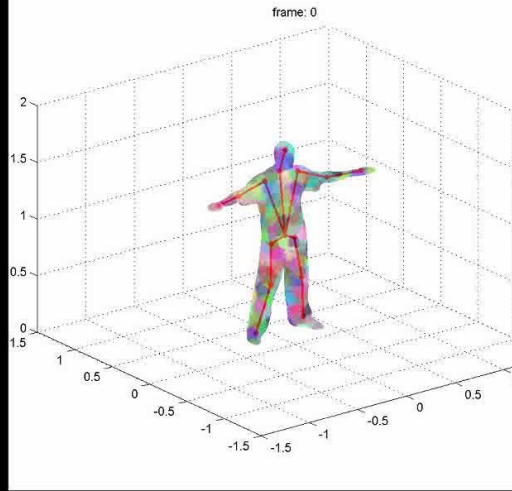
- **Pose:** **70.86mm** error in *Walking* sequence from HumanEvaII benchmark. (error < 80mm typically corresponds to a correct pose [Sigal *et al.* IJCV`12]).
- **Shape:** reprojection error (%)

Sequence	our	surface-based [4]	inverse kinematic
<i>Handstand 1</i> [1]	15.53	20.13	23.04
<i>Wheel</i> [1]	10.28	10.30	14.35
<i>Skirt</i> [1]	11.94	12.55	21.43
<i>Dance</i> [1]	9.95	9.90	15.01
<i>Crane</i> [2]	10.79	11.20	16.33
<i>Handstand 2</i> [2]	12.84	13.97	15.16
<i>Bouncing</i> [2]	9.87	9.95	14.64
<i>Free</i> [3]	14.12	14.69	-

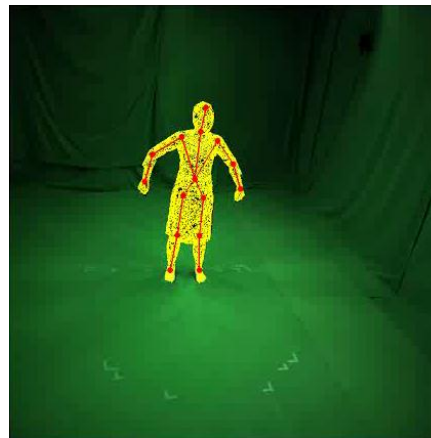
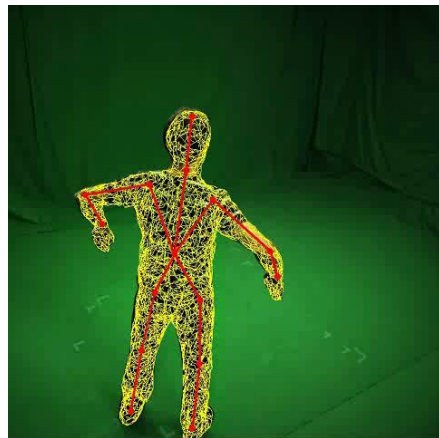
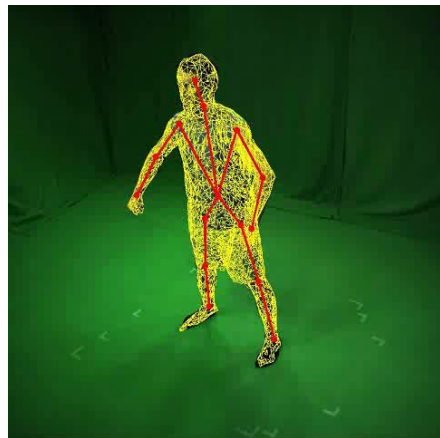
1. Gall *et al.* CVPR`09
2. Vlasic *et al.* ToG`08
3. Starck *et al.* CGA`07
4. Cagniard *et al.* ECCV`10

Qualitative results

[Starck *et al.* 2007]



[Vlasic *et al.* 2008]



[Gall *et al.* 2009]

Conclusion and future work

- A method that **jointly recovers the pose and the shape** of human body has been proposed.
- We introduce a novel **SVM-based classification** scheme that filters target point clouds and thus helps better correspondence search.
- Future directions include alleviating the requirement of **background subtraction**, and exploiting more **photometric information**.
- More experiment results in the poster session.

Thank you!

Questions?