

MonoPerfCap: Human Performance Capture from Monocular Video

INTRODUCTION

- MonoPerfCap tackles the problem of human performance capture with general clothing from monocular RGB videos.
- It can work even outdoors, with general background.
- It overcomes the limitation of depth cameras and multi-view setups.



KEY ASPECTS OF THE METHOD

- Marker-less approach
- Temporally coherent
- Reconstructs articulated human skeleton motion
- Reconstructs non-rigid surface deformation
- Pre-acquired person-specific template mesh is needed

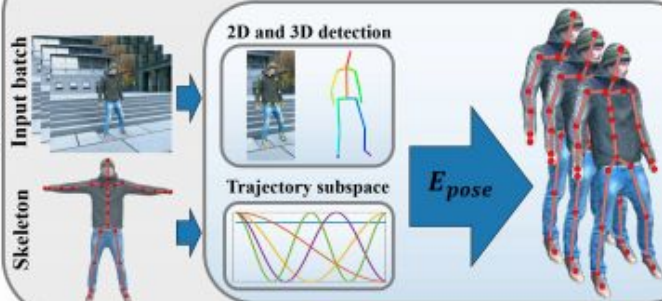
METHOD OVERVIEW

- Human motion is parameterized based on two level deformation hierarchy:
 - **Coarser level** - Articulated motion is captured in skeleton deformation space.
 - **Finer Level** - A deformation field is parameterized by an embedded deformation graph
- Motion capture is performed in a **coarse-to-fine** method, based on **two** steps:
 - **Batch based pose estimation**
 - **Silhouette based refinement**

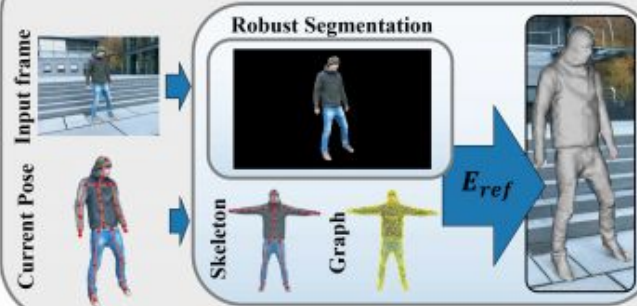
Surface Representation



Batch-based Pose Estimation (Sec. 4)



Silhouette-based Refinement (Sec. 5)



TEMPLATE MESH ACQUISITION

- High res. video of the person standing in a static T-pose is recorded using a handheld camera orbiting around the person.
- 60 images are uniformly sampled from the video for image based reconstruction of triangulated surface.



BATCH BASED POSE ESTIMATION

Articulated human motion is parameterized based on low dimensional skeleton subspace $\mathbf{S} = \{\mathbf{t}, \mathbf{R}, \Theta\}$

\mathbf{t} = Position of root joint ($\mathbf{t} \in \mathbf{R}^3$)

\mathbf{R} = Rotation of root joint ($\mathbf{R} \in \mathbf{SO}(3)$)

Θ = 27 joint angles stacked together ($\Theta \in \mathbf{R}^{27}$)

$\mathbf{N}_d = 16$ (No. of joints)

This leads to a **33 dimensional** deformation subspace.

BATCH BASED POSE ESTIMATION CONTD.

- Given monocular input video $\mathcal{V} = \{I_f\}_{f=1}^N$, with N image frames I_f .
- The goal is to estimate skeleton parameters S_f for all input frames.
- A novel batch-based approach is proposed that jointly recovers the motion for a continuous window in time:

$$\mathcal{B} = \{S_f \mid f_{start} \leq f \leq f_{end}\}$$

f_{start} = index of the first frame of the batch

f_{end} = index of the last frame of the batch

$$|\mathcal{B}| = 50$$

- The input video is partitioned into a series of overlapping frames (10 frames overlap), and a linear blending function is used in the overlap region.

BATCH BASED POSE ESTIMATION CONTD.

The problem of estimating the articulated motion of each batch \mathcal{B} is phrased as a constrained optimization problem-

$$\begin{aligned} \mathcal{B}^* = \operatorname{argmin}_{\mathcal{B}} \quad & E_{\text{pose}}(\mathcal{B}) , \\ \text{subject to} \quad & \Theta_{\min} \leq \Theta_f \leq \Theta_{\max} , \\ & \forall f \in [f_{\text{start}}, f_{\text{end}}] , \end{aligned}$$

Batch based pose estimation objective-

$$E_{\text{pose}}(\mathcal{B}) = \underbrace{E_{2d}(\mathcal{B}) + w_{3d}E_{3d}(\mathcal{B})}_{\text{data fitting}} + \underbrace{w_d E_d(\mathcal{B})}_{\text{regularization}}$$

BATCH BASED POSE ESTIMATION CONTD.

$E_{2d} \Rightarrow$ based on joint detections in image space

$E_{3d} \Rightarrow$ based on 3D joints

$E_d \Rightarrow$ based on DCT

For each input image I_f and each of the 16 joints J_i , we estimate -

$d_{f,i}^{2d} \Rightarrow$ 2D position of the joint in image space

$d_{f,i}^{3d} \Rightarrow$ 3D position of the joint

Resnet based CNN joint position regression method of **[Mehta et al. 2016]** is used to determine these terms.

BATCH BASED POSE ESTIMATION CONTD.

- The **2D** pose network is trained on the **MPII Human Pose** and **LSP** datasets, and the **3D** pose network is fine-tuned from the 2D pose network on the **H3.6M** and **3DHP** datasets.
- 2D joint alignment term is a re-projection constraint enforcing that the projected joint positions $\mathbf{J}_i(\mathbf{S}_f)$ closely match the corresponding 2D detections:

$$E_{2d}(\mathcal{B}) = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{S}_f \in \mathcal{B}} \frac{1}{N_d} \sum_{i=1}^{N_d} \left\| \Pi(\mathbf{J}_i(\mathbf{S}_f)) - \mathbf{d}_{f,i}^{2d} \right\|_2^2$$

($\Pi : \mathbf{R}^3 \rightarrow \mathbf{R}^2$ implements the full perspective camera projection)

BATCH BASED POSE ESTIMATION CONTD.

$$E_{3d}(\mathcal{B}) = \frac{1}{|\mathcal{B}|} \sum_{S_f \in \mathcal{B}} \frac{w_f}{N_d} \sum_{i=1}^{N_d} \left\| J_i(S_f) - (\mathbf{d}_{f,i}^{3d} + \mathbf{t}_f) \right\|_2^2$$

- Since the 3D joint detections are normalized for a skeleton with average bone length and are predicted relative to the root joint, rather than in camera space, they have to be rescaled to match the actor model, and mapped to their corresponding camera space position based on an unknown per-frame global translation \mathbf{t}_f .
- In order to prune frames with low 3D detection confidence, we measure the per-frame PCK error \mathbf{PCK}_f between the 2D joint detections and the projected 3D detections and apply a per-frame binary weight w_f to the 3D data term.

$$w_f = \begin{cases} 1 & \text{if } \mathbf{PCK}_f < \mathit{thres}_{pck}, \\ 0 & \text{else.} \end{cases}$$

BATCH BASED POSE ESTIMATION CONTD.

- Temporal smoothness is imposed by forcing the trajectory of each skeleton parameter to lie on a low dimensional linear subspace.
- All pose estimates $S_f \in \mathcal{B}$ are coupled by minimizing the distance to a $K = 8$ dimensional linear subspace $\text{DCT} \in \mathbb{R}^{K \times |\mathcal{B}|}$ [Park et al. 2015] spanned by the K lowest frequency basis vectors of the discrete cosine transform (DCT).

$$E_d(\mathcal{B}) = \frac{1}{|\mathcal{B}|} \left\| \Lambda S_{\mathcal{B}} \text{Null}(\text{DCT}) \right\|_F^2$$

$$S_{\mathcal{B}} = [S_{f_{start}}, \dots, S_{f_{end}}] \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{B}|}$$

$$\Lambda = \text{diag}([\lambda_t, \lambda_R, \lambda_{\Theta}])$$

- Finally, the constrained non-linear least squares optimization problem is solved using the Levenberg Marquardt (LM) algorithm.

SILHOUETTE-BASED REFINEMENT:

Why Refinement:

Batch-based pose optimization does not capture non-rigid surface deformation due to apparel and skin, and thus leads to misalignments between the skeleton-deformed template mesh and the input images, particularly at the boundaries.

This has 3 tasks

- 1) Automatic Silhouette Extraction
- 2) Silhouette-based Pose Refinement

SILHOUETTE-BASED REFINEMENT:

3)Silhouette-based Non-Rigid Surface Refinement

1)Automatic Silhouette Extraction:

To extract silhouette GrabCut was used.

GrabCut requires a user-specified initialization $T = \{T_b, T_{ub}, T_{uf}, T_f\}$.

T_b, T_f = known background and foreground masks.

T_{ub}, T_{uf} = uncertain background and foreground regions.

Segmentation is performed over T_{ub}, T_{uf} .

SILHOUETTE EXTRACTION:

To automate this we first rasterize the skeleton and the deformed dense actor template $V(S_f)$ to obtain two masks R and M respectively.

T is initialised as following:

$$T_f = R \cup \text{erosion}(M),$$

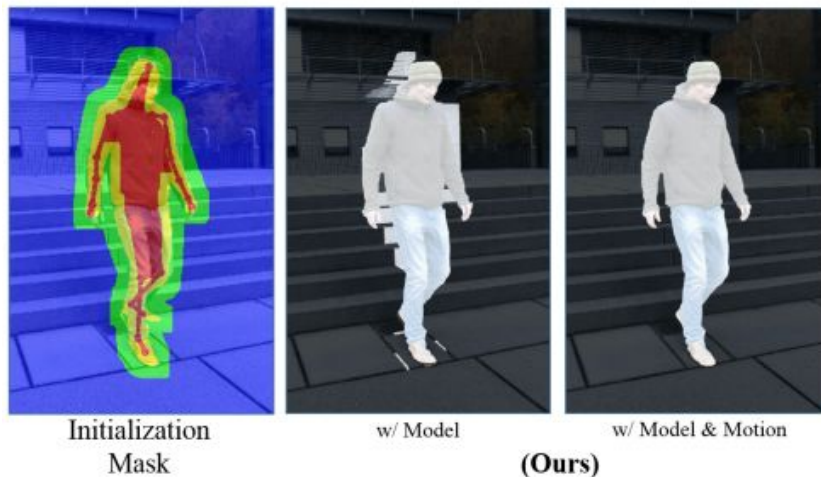
$$T_b = \text{dilation}(M),$$

$$T_{uf} = M - T_f ,$$

$$T_{ub} = \text{dilation}(M) - M.$$

SILHOUETTE EXTRACTION:

To improve robustness motion cues were incorporated in the objective. the temporal per-pixel color gradients between adjacent frames and encourage neighboring pixels with small temporal gradients to belong to the same region.



The masks $\{T_b, T_{ub}, T_{uf}, T_f\}$ are illustrated in red, blue, yellow and green.

2)Silhouette-based Pose Refinement

Refinement is performed in an Iterative Closest Point (ICP) manner.

In an iteration, for each boundary point of the projected surface model, search for its closest point on the image silhouette that shares a similar normal direction.

refine the pose by solving the following non-linear least squares optimization

$\mathbf{E}_{\text{ref}}(\mathbf{S}_f) = \mathbf{E}_{\text{con}}(\mathbf{S}_f) + w_{\text{stab}} \mathbf{E}_{\text{stab}}(\mathbf{S}_f)$ where \mathbf{E}_{con} aligns the mesh boundary with the input silhouette, \mathbf{E}_{stab} constrains the solution to stay close to the batch-based results and w_{stab} balances the importance of the two terms.

$\mathbf{E}_{\text{con}}(\mathbf{S}_f) = 1/|\mathbf{S}| \sum [\mathbf{n}_k^T \cdot (\mathbf{\Pi}(\mathbf{V}_k(\mathbf{S}_f) - \mathbf{s}_k))]^2 \quad \forall k \in \mathbf{S}$ where \mathbf{S} is the boundary of the actor model, \mathbf{v}_k the position of vertex k and $\mathbf{s}_k \in \mathbb{R}^2$ the corresponding silhouette point in the image with 2D normal \mathbf{n}_k .

Refinement Continued:

$E_{stab}(S_f) = 1/N_d \sum ||J_i(S_f) - J_i(S_f^*)|| \quad \forall i=1 \text{ to } N_d$ where S_f^* are the joint angles after batch-based pose estimation and $J_i(\cdot)$ computes the 3D position of joint J_i .



(a)



(b)



(c)



(d)



(e)

3) Silhouette-based Non-Rigid Surface Refinement

A deformation graph D , consisting of $M \approx 1000$ nodes, is generated from the template mesh using a uniform mesh decimation/simplification strategy.



A radius of influence is assigned for each deformation node by computing the maximum geodesic distance to its connected graph nodes.

Non-Rigid Surface Refinement Contd:

Each node defines a local warp field W_i that rotates $R_i \in SO(3)$ and translates $t_i \in R^3$ points $x \in R^3$ in the surrounding space.

$W_i(x) = R_i(x - \hat{g}_i) + \hat{g}_i + t_i$, where $\hat{g}_i \in R^3$ is the canonical position of node i , computed with the result of the pose refinement.

The graph and it's DoF are represented as $D = \{(R_i, t_i) | i \in [0, M)\}$.

The deformation is achieved by linear blending of the per-node warp fields i.e,

$v_i = W(\hat{v}_i) = \sum b_{i,k}(x) \cdot W_k(\hat{v}_i) \quad \forall k \in F_i$, where $v_i \in R^3$ is the deformed vertex position, $\hat{v}_i \in R^3$ is the canonical position of vertex i and F_i is the set of deformation nodes that influence vertex i .

Non-Rigid Surface Refinement Contd:

Given the embedded deformation graph, our silhouette-based surface refinement is expressed as the following optimization problem:

$E_{\text{surf}}(\mathbf{D}) = E_{\text{con}}(\mathbf{D}) + w_{\text{arap}} E_{\text{arap}}(\mathbf{D})$. Here E_{con} is the silhouette alignment term, E_{arap} is as-rigid-as-possible regularization term.

$E_{\text{con}}(\mathbf{D}) = 1/|S| \sum [n_k^T \cdot (\mathbf{I}(\mathbf{V}_k(\mathbf{D}) - \mathbf{s}_k))]^2 \quad \forall k \in S$ where S is the model silhouette, v_k the position of vertex k and $s_k \in \mathbb{R}^2$ its corresponding silhouette point with normal $n_k \in \mathbb{R}^2$.

$E_{\text{arap}}(\mathbf{D}) = 1/M \sum_{i=1}^m w_i \sum_{j \in N(i)} \|(\mathbf{g}_i - \mathbf{g}_j) - R_i(\hat{\mathbf{g}}_i - \hat{\mathbf{g}}_j)\|_2^2$ Here, $\mathbf{g}_i = W_i(\hat{\mathbf{g}}_i) = \hat{\mathbf{g}}_i + \mathbf{t}_i$ is the deformed

position of node $\hat{\mathbf{g}}_i$ and N_i is its 1-ring neighbourhood.

Contd.

Surface refinement is performed in an ICP like manner.

Initialize the optimization problem based on the pose refinement result and minimize E_{surf} using the LM algorithm.



Input



Before



After

Results from Paper

