Supplementary Materials Revitalizing Optimization for 3D Human Pose and Shape Estimation: A Sparse Constrained Formulation

Taosha Fan^{1,2}, Kalyan Vasudev Alwala¹, Donglai Xiang^{3,4}, Weipeng Xu³, Todd Murphey², Mustafa Mukadam¹

> ¹Facebook AI Research, ²Northwestern University, ³Facebook Reality Labs, ⁴Carnegie Mellon University

Abstract

In this supplementary material, we present the proofs of the propositions in the paper and a comprehensive complexity analysis of the dense unconstrained and sparse constrained formulations for 3D human pose and shape estimation, from which we derive an efficient algorithm to compute the Gauss-Newton direction. In addition, we present more results of qualitative comparisons and ablation studies to validate our work. Finally, we provide a more detailed description of our real-time motion capture framework, the prior loss of joint states, and how to implement our method on similar articulated tracking problems.

A. Proofs

A.1. Proof of Proposition 1

In this proof, we show the following two optimization problems are equivalent:

$$\min_{\mathbf{T}_0,\,\boldsymbol{\Omega},\,\boldsymbol{\beta}} \mathbf{E} = \sum_{i=0}^{K} \frac{1}{2} \|\mathbf{r}_i(\mathbf{T}_0,\,\boldsymbol{\Omega},\,\boldsymbol{\beta})\|^2, \tag{1}$$

and

$$\min_{\{\mathbf{T}_i,\,\boldsymbol{\Omega}_i,\,\boldsymbol{\beta}_i\}_{i=0}^K} \mathbf{E} = \sum_{i=0}^K \frac{1}{2} \|\mathbf{r}_i(\mathbf{T}_i,\,\boldsymbol{\Omega}_i,\,\boldsymbol{\beta}_i)\|^2 \qquad (2)$$

subject to

$$\mathbf{T}_{i} = \mathbf{F}_{i}(\mathbf{T}_{\mathrm{par}(i)}, \boldsymbol{\beta}_{\mathrm{par}(i)}, \boldsymbol{\Omega}_{i})$$

$$\triangleq \mathbf{T}_{\mathrm{par}(i)} \begin{bmatrix} \boldsymbol{\Omega}_{i} & \mathcal{S}_{i} \cdot \boldsymbol{\beta}_{\mathrm{par}(i)} + \mathbf{l}_{i} \\ \mathbf{0} & 1 \end{bmatrix}, \qquad (3a)$$

$$\boldsymbol{\beta}_i = \boldsymbol{\beta}_{\mathrm{par}(i)}.\tag{3b}$$

In Eqs. (1) and (2), $\mathbf{T}_i \in SE(3)$ is the rigid body transformation of body part *i*, and $\mathbf{\Omega}_i$ is the state of joint *i*, and $\Omega \triangleq (\Omega_1, \dots, \Omega_K) \in SO(3)^K$ are the joint states, and β and $\beta_i \in \mathbb{R}^P$ are the shape parameters, and $F_i(\cdot)$: $SE(3) \times \mathbb{R}^P \times SO(3) \to SE(3)$ is a function that maps $\mathbf{T}_{\text{par}(i)}, \beta_{\text{par}(i)}$ and Ω_i to \mathbf{T}_i . Note that Eqs. (1) and (2) are the dense unconstrained and sparse constrained formulations, respectively, for 3D human pose and shape estimation that are defined in the paper.

From Eq. (3b), if we let $\beta_0 = \beta$, then, $\beta_i = \beta$ for all $i = 1, \dots, K$. Thus, Eq. (2) is reduced to

$$\min_{\{\mathbf{T}_i,\,\boldsymbol{\Omega}_i,\,\boldsymbol{\beta}_i\}_{i=0}^K} \mathbf{E} = \sum_{i=0}^K \frac{1}{2} \|\mathbf{r}_i(\mathbf{T}_i,\,\boldsymbol{\Omega}_i,\,\boldsymbol{\beta})\|^2 \qquad (4)$$

subject to

$$\begin{aligned} \mathbf{\Gamma}_{i} = \mathbf{F}_{i}(\mathbf{T}_{\mathrm{par}(i)}, \boldsymbol{\beta}, \boldsymbol{\Omega}_{i}) \\ = \mathbf{T}_{\mathrm{par}(i)} \begin{bmatrix} \boldsymbol{\Omega}_{i} & \mathcal{S}_{i} \cdot \boldsymbol{\beta} + \mathbf{l}_{i} \\ \mathbf{0} & 1 \end{bmatrix}. \end{aligned}$$

$$(5)$$

Next, as mentioned in the paper, if we perform a top-down traversal of the kinematic tree of the SMPL model and recursively exploit Eq. (5) for each body part $i = 1, \dots, K$, then, all of $\mathbf{T}_i \in SE(3)$ can be represented as a function of the root pose $\mathbf{T}_0 \in SE(3)$, and the joint states $\mathbf{\Omega} \in SO(3)^K$, and the shape parameter $\boldsymbol{\beta} \in \mathbb{R}^P$, i.e.,

$$\mathbf{T}_{i} \triangleq \mathbf{T}_{i} \left(\mathbf{T}_{0}, \, \boldsymbol{\Omega}, \, \boldsymbol{\beta} \right) \tag{6}$$

If we use Eq. (6) to cancel out non-root rigid body transformations \mathbf{T}_i $(1 \le i \le K)$, then, each $\mathbf{r}_i(\cdot)$ in Eq. (4) is rewritten as a function of $\mathbf{T}_0 \in SE(3)$, and $\Omega \in SO(3)^K$, and $\boldsymbol{\beta} \in \mathbb{R}^P$, from which we obtain an optimization problem of a dense unconstrained formulation

$$\min_{\mathbf{T}_0,\,\boldsymbol{\Omega},\,\boldsymbol{\beta}} \mathsf{E} = \sum_{i=0}^{K} \frac{1}{2} \| \mathsf{r}_i(\mathbf{T}_0,\,\boldsymbol{\Omega},\,\boldsymbol{\beta}) \|^2$$

that is the same as Eq. (1). Therefore, it can be concluded that Eqs. (1) and (2) are equivalent. The proof is completed.

A.2. Proof of Proposition 2

The proof of proposition 2 is organized as follows: we present an overview of the steps to compute the Gauss-Newton direction in Section A.2.1, and show that the steps for the two formulations result in the same Gauss-Newton direction in Section A.2.2, and derive a dynamic programming algorithm to solve the quadratic program of the sparse constrained formulation in Section A.2.3, and analyze the complexity of the aforementioned steps to compute the Gauss-Newton direction in Section A.2.4.

A.2.1 Steps to Compute the Gauss-Newton Direction

With similar notation to the paper, we introduce $\mathbf{x} \triangleq (\mathbf{T}_0, \mathbf{\Omega}, \boldsymbol{\beta}) \in SE(3) \times SO(3)^K \times \mathbb{R}^P$ and $\mathbf{x}_i \triangleq (\mathbf{T}_i, \boldsymbol{\beta}_i) \in SE(3) \times \mathbb{R}^P$. Then Eqs. (1) and (2) can be rewritten as

$$\min_{\mathbf{x}} \mathbf{E} = \sum_{i=0}^{K} \frac{1}{2} \|\mathbf{r}_i(\mathbf{x})\|^2, \tag{7}$$

and

$$\min_{\{\mathbf{x}_i,\,\boldsymbol{\Omega}_i\}_{i=0}^{K}} \mathsf{E} = \sum_{i=0}^{K} \frac{1}{2} \|\mathbf{r}_i(\mathbf{x}_i,\,\boldsymbol{\Omega}_i)\|^2 \tag{8}$$

subject to

$$\mathbf{x}_{i} = \begin{bmatrix} \mathbf{F}_{i}(\mathbf{x}_{\mathrm{par}(i)}, \, \mathbf{\Omega}_{i}) \\ \boldsymbol{\beta}_{\mathrm{par}(i)} \end{bmatrix}, \tag{9}$$

respectively. For analytical clarity, we assume with no loss of generality that the residues $\mathbf{r}_i(\mathbf{x})$ and $\mathbf{r}_i(\mathbf{x}_i, \mathbf{\Omega}_i)$ are $N_i \times 1$ vectors for $i = 0, \dots, K$.

Following the procedure originally given in the paper, an overview of steps to compute the Gauss-Newton direction for the dense unconstrained and sparse constrained formulations is given in Table 1, which will be frequently used in the rest of this proof.

A.2.2 The Equivalence of the Gauss-Newton Direction

In Table 1, since Steps 2 and 3 are the reformulation of Step 1, we only need to show that the linearizations of dense unconstrained and sparse constrained formulations in Step 1, i.e., Eqs. (10) and (12), are equivalent. From Eq. (6), the rigid body transformation $\mathbf{T}_i \in SE(3)$ of body part *i* can be written as a function of \mathbf{T}_0 , $\boldsymbol{\Omega}$ and $\boldsymbol{\beta}$. Furthermore, it is by the definition of $\mathbf{r}_i(\cdot)$ that

$$\mathtt{r}_i(\mathbf{T}_0,\,oldsymbol{\Omega},\,oldsymbol{eta}) = \mathtt{r}_iig(\mathbf{T}_i(\mathbf{T}_0,\,oldsymbol{\Omega},\,oldsymbol{eta}),\,oldsymbol{\Omega}_i,\,oldsymbol{eta}ig).$$

From the equation above, $J_i \Delta x$ in Eq. (10) can be computed using $J_{i,1}$ and $J_{i,2}$ in Eq. (12):

$$\mathbf{J}_{i}\Delta\mathbf{x} = \mathbf{J}_{i,1} \begin{bmatrix} \frac{\partial \mathbf{T}_{i}}{\partial \mathbf{T}_{0}} \Delta \mathbf{T}_{0} + \frac{\partial \mathbf{T}_{i}}{\partial \boldsymbol{\Omega}} \Delta \boldsymbol{\Omega} + \frac{\partial \mathbf{T}_{i}}{\partial \boldsymbol{\beta}} \Delta \boldsymbol{\beta} \\ \boldsymbol{\beta} \end{bmatrix} + \mathbf{J}_{i,2}\Delta \boldsymbol{\Omega}_{i}.$$
(21)

Note that the partial derivatives $\frac{\partial \mathbf{T}_i}{\partial \mathbf{T}_0}$, $\frac{\partial \mathbf{T}_i}{\partial \Omega}$ and $\frac{\partial \mathbf{T}_i}{\partial \beta}$ in the right-hand side of Eq. (21) are obtained by the recursive implementation of Eq. (13). Therefore, it can be concluded that Eqs. (10) and (12) are equivalent to each other, which suggests that the dense unconstrained and sparse constrained formulations result in the same Gauss-Newton direction.

A.2.3 Algorithm to Solve Eq. (19)

In Table 1, it is straightforward to follow **Steps 1–2** of the sparse constrained formulation to compute the Gauss-Newton direction. Next, we need to solve the quadratic program of Eq. (19) in **Step 3**, which is nontrivial. In this subsection, we derive a dynamic programming algorithm that exploits the sparsity and constraints of Eq. (13) such that the Gauss-Newton direction can be exactly computed.

For notational simplicity, we let par(i), chd(i) and des(i) be the parent, children and descendants of body part i in the kinematics tree, and assume i > par(i) for all $i = 1, \dots, K$.

First, we define $\mathcal{E}_i(\cdot) : \mathbb{R}^{6+P} \to \mathbb{R}$ to be a function of $\Delta \mathbf{x}_{\text{par}(i)} \in \mathbb{R}^{6+P}$ in the form of an optimization problem of $\{\Delta \mathbf{x}_j, \Delta \Omega_j\}$ for $j \in \{i\} \cup \text{des}(i)$

$$\mathcal{E}_{i}(\Delta \mathbf{x}_{\text{par}(i)}) \triangleq \min_{\{\Delta \mathbf{x}_{j}, \Delta \mathbf{\Omega}_{j}\}_{j \in \{i\} \cup \text{des}(i)}} \sum_{j \in \{i\} \cup \text{des}(i)} \left[\frac{1}{2} \Delta \mathbf{x}_{j}^{\top} \mathbf{H}_{j,11} \Delta \mathbf{x}_{j} + \Delta \mathbf{\Omega}_{j}^{\top} \mathbf{H}_{j,21} \Delta \mathbf{x}_{j} + \frac{1}{2} \Delta \mathbf{\Omega}_{j}^{\top} \mathbf{H}_{j,22} \Delta \mathbf{\Omega}_{j} + \mathbf{g}_{j,1}^{\top} \Delta \mathbf{x}_{j} + \mathbf{g}_{j,2}^{\top} \Delta \mathbf{\Omega}_{j}\right] \quad (22)$$

subject to

$$\Delta \mathbf{x}_j = \mathbf{A}_j \Delta \mathbf{x}_{\text{par}(j)} + \mathbf{B}_j \Delta \mathbf{\Omega}_j, \ \forall j \in \{i\} \cup \text{des}(i), \ (23)$$

in which $\Delta \mathbf{x}_{\text{par}(i)} \in \mathbb{R}^{6+P}$ is given. Furthermore, if $\mathcal{E}_j(\cdot) : \mathbb{R}^{6+P} \to \mathbb{R}$ is defined for all $j \in \text{chd}(i)$, then, it is from Eq. (22) that $\mathcal{E}_i(\cdot)$ can be reduced to an optimization problem of $\Delta \mathbf{x}_i$ and $\Delta \Omega_i$

$$\mathcal{E}_{i}(\Delta \mathbf{x}_{\mathrm{par}(i)}) \triangleq \min_{\Delta \mathbf{x}_{i}, \Delta \Omega_{i}} \left[\frac{1}{2} \Delta \mathbf{x}_{i}^{\top} \mathbf{H}_{i,11} \Delta \mathbf{x}_{i} + \Delta \Omega_{i}^{\top} \mathbf{H}_{i,21} \Delta \mathbf{x}_{i} + \frac{1}{2} \Delta \Omega_{i}^{\top} \mathbf{H}_{i,22} \Delta \Omega_{i} + \mathbf{g}_{i,1}^{\top} \Delta \mathbf{x}_{i} + \mathbf{g}_{i,2}^{\top} \Delta \Omega_{i} + \sum_{j \in \mathrm{chd}(i)} \mathcal{E}_{j}(\Delta \mathbf{x}_{i}) \right]$$
(24)

subject to

$$\Delta \mathbf{x}_i = \mathbf{A}_i \Delta \mathbf{x}_{\mathrm{par}(i)} + \mathbf{B}_i \Delta \mathbf{\Omega}_i,$$

in which $\Delta \mathbf{x}_{\text{par}(i)} \in \mathbb{R}^{6+P}$ is given. Note that Eq. (24) is an intermediate procedure that is essential for our dynamic programming algorithm.

	Dense Unconstrained Formulation	Sparse Constrained Formulation
		The linearization of Eq. (8) results in
Step 1	The linearization of Eq. (7) results in $ \min_{\Delta \mathbf{x}} \Delta \mathbf{E} = \sum_{i=0}^{K} \frac{1}{2} \ \mathbf{J}_{i} \Delta \mathbf{x} + \mathbf{r}_{i} \ ^{2}, (10) $ in which $\Delta \mathbf{x} \triangleq (\Delta \mathbf{T}_{0}, \Delta \Omega, \Delta \beta) \in \mathbb{R}^{6+3K+P}, \Delta \mathbf{T}_{0} \in \mathbb{R}^{6}, \Delta \Omega \in \mathbb{R}^{3K} \text{ and } \Delta \beta \in \mathbb{R}^{P}$ are the Gauss-Newton directions of $\mathbf{x}, \mathbf{T}_{0}, \Omega$ and β , respectively, and $ \mathbf{J}_{i} \triangleq \left[\frac{\partial \mathbf{r}_{i}}{\partial \mathbf{T}_{0}} \frac{\partial \mathbf{r}_{i}}{\partial \Omega} \frac{\partial \mathbf{r}_{i}}{\partial \beta} \right] \in \mathbb{R}^{N_{i} \times (6+3K+P)} (11) $ is the Jacobian of $\mathbf{r}_{i}(\cdot)$, and $\mathbf{r}_{i} \in \mathbb{R}^{N_{i}}$ is the residue.	The intermation of Eq. (b) relation in $ \min_{\{\Delta \mathbf{x}_{i}, \Delta \Omega_{i}\}_{i=0}^{K}} \Delta \mathbf{E} = \sum_{i=0}^{K} \frac{1}{2} \ \mathbf{J}_{i,1} \Delta \mathbf{x}_{i} + \mathbf{J}_{i,2} \Delta \Omega_{i} + \mathbf{r}_{i} \ ^{2} $ (12) subject to $ \Delta \mathbf{x}_{i} = \mathbf{A}_{i} \Delta \mathbf{x}_{par(i)} + \mathbf{B}_{i} \Delta \Omega_{i}, (13) $ in which $\Delta \mathbf{x}_{i} \triangleq (\Delta \mathbf{T}_{i}, \Delta \beta_{i}) \in \mathbb{R}^{6+P}, \Delta \mathbf{T}_{i} \in \mathbb{R}^{6}, \Delta \Omega_{i} \in \mathbb{R}^{3} \text{ and } \Delta \beta_{i} \in \mathbb{R}^{P} \text{ are the Gauss-Newton directions of } \mathbf{x}_{i}, \mathbf{T}_{i}, \Omega_{i} \text{ and } \beta_{i}, \text{ respectively, and} $ $ \mathbf{J}_{i,1} \triangleq \left[\frac{\partial \mathbf{r}_{i}}{\partial \mathbf{T}_{i}} \frac{\partial \mathbf{r}_{i}}{\partial \beta_{i}} \right] \in \mathbb{R}^{N_{i} \times (6+P)} (14) $ and $ \mathbf{J}_{i,2} \triangleq \frac{\partial \mathbf{r}_{i}}{\partial \Omega_{i}} \in \mathbb{R}^{N_{i} \times 3} (15) $ are the Jacobians of $\mathbf{r}_{i}(\cdot)$, and $ \mathbf{A}_{i} \triangleq \left[\frac{\partial \mathbf{F}_{i}}{\partial \mathbf{T}_{par(i)}} \frac{\partial \mathbf{F}_{i}}{\partial \beta_{par(i)}} \right] \in \mathbb{R}^{(6+P) \times (6+P)} (16) $ and $ \mathbf{B}_{i} \triangleq \left[\frac{\partial \mathbf{F}_{i}}{\partial \Omega_{i}} \right] \in \mathbb{R}^{(6+P) \times 3} (17) $ are the partial derivatives of Eq. (13), and $\mathbf{r}_{i} \in \mathbb{R}^{N_{i}}$ is the residue.
Step 2	Reformulate Eq. (10) as $ \min_{\Delta \mathbf{x}} \Delta \mathbf{E} = \frac{1}{2} \Delta \mathbf{x}^{\top} \mathbf{H} \Delta \mathbf{x} + \mathbf{g}^{\top} \Delta \mathbf{x} (18) $ in which $\mathbf{H} \triangleq \sum_{i=0}^{K} \mathbf{J}_{i}^{\top} \mathbf{J}_{i} \in \mathbb{R}^{(6+3K+P) \times (6+3K+P)}$ is the Hessian, and $\mathbf{g} \triangleq \sum_{i=0}^{K} \mathbf{J}_{i}^{\top} \mathbf{r}_{i} \in \mathbb{R}^{(6+3K+P)}$ is the gradient.	Reformulate Eq. (12) as $ \min_{\{\Delta \mathbf{x}_{i}, \Delta \Omega_{i}\}_{i=0}^{K}} \Delta \mathbf{E} = \sum_{i=0}^{K} \left[\frac{1}{2} \Delta \mathbf{x}_{i}^{\top} \mathbf{H}_{i,11} \Delta \mathbf{x}_{i} + \Delta \Omega_{i}^{\top} \mathbf{H}_{i,21} \Delta \mathbf{x}_{i} + \frac{1}{2} \Delta \Omega_{i}^{\top} \mathbf{H}_{i,22} \Delta \Omega_{i} + \mathbf{g}_{i,1}^{\top} \Delta \mathbf{x}_{i} + \mathbf{g}_{i,2}^{\top} \Delta \Omega_{i} \right] (19) $ subject to $ \Delta \mathbf{x}_{i} = \mathbf{A}_{i} \Delta \mathbf{x}_{\text{par}(i)} + \mathbf{B}_{i} \Delta \Omega_{i}, $ in which $\mathbf{H}_{i,11} \triangleq \mathbf{J}_{i,1}^{\top} \mathbf{J}_{i,1} \in \mathbb{R}^{(6+P) \times (6+P)}, \mathbf{H}_{i,21} \triangleq \mathbf{J}_{i,2}^{\top} \mathbf{J}_{i,1} \in \mathbb{R}^{3 \times (6+P)}, \text{ and } \mathbf{H}_{i,22} \triangleq \mathbf{J}_{i,2}^{\top} \mathbf{J}_{i,2} \in \mathbb{R}^{3 \times 3} $ are the Hessians, and $\mathbf{g}_{i,1} \triangleq \mathbf{J}_{i,1}^{\top} \mathbf{r}_{i} \in \mathbb{R}^{6+P} \text{ and } \mathbf{g}_{i,2} \triangleq \mathbf{J}_{i,2}^{\top} \mathbf{r}_{i} \in \mathbb{R}^{6+P}$ are the gradients.
Step 3	Compute the Gauss-Newton direction from Eq. (18), which has a closed-form solution $\Delta \mathbf{x} = -\mathbf{H}^{-1}\mathbf{g}.$ (20)	Compute the Gauss-Newton direction from Eq. (19), which can be exactly solved by Algorithm 1.

Table 1: Steps to compute the Gauss-Newton direction for the dense unconstrained and sparse constrained formulations.

Algorithm 1 Solve Eq. (19) and compute the Gauss-Newton direction Input: $\{\mathbf{H}_{i,11}, \mathbf{H}_{i,21}, \mathbf{H}_{i,22}, \mathbf{g}_{i,1}, \mathbf{g}_{i,2}\}_{i=0}^{K}$ **Output**: $\{\Delta \mathbf{x}_i, \Delta \mathbf{\Omega}_i\}_{i=0}^K$ and $\Delta \mathbf{E}_0$ 1: for $i = K \rightarrow 1$ do 2: $\mathbf{N}_{i,11} = \mathbf{H}_{i,11} + \sum_{j \in \operatorname{chd}(i)} \mathbf{M}_{j,11}$ $N_{i,21} = H_{i,21}$ 3: $\mathbf{N}_{i,22} = \mathbf{H}_{i,22}$ 4: $\mathbf{n}_{1,i} = \mathbf{g}_{i,1} + \sum_{j \in \mathrm{chd}(i)} \mathbf{m}_{j,1}$ 5: $\mathbf{n}_{i,2} = \mathbf{g}_{i,2}$ 6: $\Delta \overline{\mathbf{E}}_i = \sum_{i \in \mathrm{chd}(i)} \Delta \mathbf{E}_j$ 7: $\mathbf{Q}_{i,11} = \mathbf{A}_i^\top \mathbf{N}_{i,11} \mathbf{A}_i$ 8: $\mathbf{Q}_{i,21} = \mathbf{B}_i^\top \mathbf{N}_{i,11} \mathbf{A}_i + \mathbf{N}_{i,21} \mathbf{A}_i$ 9: $\mathbf{Q}_{i,22} = \mathbf{B}_i^{\top} \mathbf{N}_{i,11} \mathbf{B}_i + \mathbf{N}_{i,21} \mathbf{B}_i + \mathbf{B}_i^{\top} \mathbf{N}_{i,21}^{\top} + \mathbf{N}_{i,22}$ 10: $\mathbf{q}_{i,1} = \mathbf{A}_i^\top \mathbf{n}_{i,1}$ 11: $\mathbf{q}_{i,2} = \mathbf{B}_i^\top \mathbf{n}_{i,1} + \mathbf{n}_{i,2}$ 12: $\mathbf{K}_i = -\mathbf{Q}_{i,22}^{-1}\mathbf{Q}_{i,21}$ 13: $\mathbf{k}_{i} = -\mathbf{Q}_{i \ 22}^{-1}\mathbf{q}_{i,2}$ 14: $\mathbf{M}_{i,11} = \mathbf{Q}_{i,11} - \mathbf{Q}_{i,21}^{\top} \mathbf{Q}_{i,22}^{-1} \mathbf{Q}_{i,21}$ 15: $\mathbf{m}_{1,i} = \mathbf{q}_{i,1} - \mathbf{Q}_{i,21}^{\top} \mathbf{Q}_{i,22}^{-1} \mathbf{q}_{i,2}$ 16: $\Delta \mathbf{E}_i = \Delta \overline{\mathbf{E}}_i - \frac{1}{2} \mathbf{q}_{i,2}^{\top} \mathbf{Q}_{i,22}^{-1} \mathbf{q}_{i,2}$ 17: 18: end for 19: $\Delta \mathbf{\Omega}_0 = \mathbf{0}$ 20: $\mathbf{M}_0 = \mathbf{H}_{0,11} + \sum_{j \in \text{chd}(0)} \mathbf{M}_{j,11}$ 21: $\mathbf{m}_0 = \mathbf{g}_{0,1} + \sum_{j \in \text{chd}(0)} \mathbf{m}_{j,1}$ 22: $\Delta \overline{\mathbf{E}}_0 = \sum_{i \in \mathrm{chd}(0)} \Delta \mathbf{E}_i$ 23: $\mathbf{x}_0 = -\mathbf{M}_0^{-1}\mathbf{m}_0$ 24: $\Delta \mathbf{E}_0 = \Delta \overline{\mathbf{E}}_0 - \frac{1}{2} \mathbf{m}_0^\top \mathbf{M}_0^{-1} \mathbf{m}_0$ 25: for $i = 1 \rightarrow K$ do 26: $\Delta \mathbf{\Omega}_i = \mathbf{K}_i \Delta \mathbf{x}_{\mathrm{par}(i)} + \mathbf{k}_i$ 27: $\Delta \mathbf{x}_i = \mathbf{A}_i \Delta \mathbf{x}_{\mathrm{par}(i)} + \mathbf{B}_i \Delta \mathbf{\Omega}_i$ 28: end for

Next, suppose that there exists $\mathbf{M}_j \in \mathbb{R}^{(6+P)\times(6+P)}$, $\mathbf{m}_j \in \mathbb{R}^{6+P}$ and $\Delta \mathbf{E}_j \in \mathbb{R}$ for all $j \in \text{chd}(i)$ such that $\mathcal{E}_j(\Delta \mathbf{x}_i)$ can be written as

$$\mathcal{E}_j(\Delta \mathbf{x}_i) = \frac{1}{2} \Delta \mathbf{x}_i^{\mathsf{T}} \mathbf{M}_j \Delta \mathbf{x}_i + \mathbf{m}_j^{\mathsf{T}} \Delta \mathbf{x}_i + \Delta \mathbf{E}_j.$$
(25)

Applying Eq. (25) to Eq. (24), we obtain

$$\mathcal{E}_{i}(\Delta \mathbf{x}_{\text{par}(i)}) = \min_{\Delta \mathbf{x}_{i}, \Delta \Omega_{i}} \frac{1}{2} \Delta \mathbf{x}_{i} \mathbf{N}_{i,11} \Delta \mathbf{x}_{i} + \Delta \Omega_{i}^{\top} \mathbf{N}_{i,21} \Delta \mathbf{x}_{i} + \frac{1}{2} \Delta \Omega_{i}^{\top} \mathbf{N}_{i,22} \Delta \Omega_{i} + \mathbf{n}_{i,1}^{\top} \Delta \mathbf{x}_{i} + \mathbf{n}_{i,2}^{\top} \Delta \Omega_{i} + \Delta \overline{\mathbf{E}}_{i} \quad (26)$$

subject to

$$\Delta \mathbf{x}_i = \mathbf{A}_i \Delta \mathbf{x}_{\mathrm{par}(i)} + \mathbf{B}_i \Delta \mathbf{\Omega}_i,$$

in which

$$\mathbf{N}_{i,11} = \mathbf{H}_{i,11} + \sum_{j \in \mathrm{chd}(i)} \mathbf{M}_i,$$
(27a)

$$\mathbf{N}_{i,21} = \mathbf{H}_{i,21},\tag{27b}$$

$$\mathbf{N}_{i,22} = \mathbf{H}_{i,22},\tag{27c}$$

$$\mathbf{n}_{i,1} = \mathbf{g}_{i,1} + \sum_{j \in \mathrm{chd}(i)} \mathbf{m}_j, \qquad (27\mathrm{d})$$

$$\mathbf{n}_{i,2} = \mathbf{g}_{i,2},\tag{27e}$$

$$\Delta \overline{\mathsf{E}}_i = \sum_{j \in \mathrm{chd}(i)} \Delta \mathsf{E}_j.$$
(27f)

Substitute Eq. (13) into Eq. (26) to cancel out $\Delta \mathbf{x}_i$ and simplify the resulting equation to an unconstrained optimization problem on $\Delta \Omega_i \in \mathbb{R}^3$:

$$\mathcal{E}_{i}(\Delta \mathbf{x}_{\mathrm{par}(i)}) = \min_{\Delta \mathbf{\Omega}_{i}} \frac{1}{2} \Delta \mathbf{x}_{\mathrm{par}(i)} \mathbf{Q}_{i,11} \Delta \mathbf{x}_{\mathrm{par}(i)} + \Delta \mathbf{\Omega}_{i}^{\top} \mathbf{Q}_{i,21} \Delta \mathbf{x}_{\mathrm{par}(i)} + \frac{1}{2} \Delta \mathbf{\Omega}_{i}^{\top} \mathbf{Q}_{i,22} \Delta \mathbf{\Omega}_{i} + \mathbf{q}_{i,1}^{\top} \Delta \mathbf{x}_{\mathrm{par}(i)} + \mathbf{q}_{i,2}^{\top} \Delta \mathbf{\Omega}_{i} + \Delta \overline{\mathbf{E}}_{i}, \quad (28)$$

in which

$$\mathbf{Q}_{i,11} = \mathbf{A}_i^\top \mathbf{N}_{i,11} \mathbf{A}_i, \tag{29a}$$

$$\mathbf{Q}_{i,21} = \mathbf{B}_i^{\top} \mathbf{N}_{i,11} \mathbf{A}_i + \mathbf{N}_{i,21} \mathbf{A}_i, \qquad (29b)$$

$$\mathbf{Q}_{i,22} = \mathbf{B}_i^{\top} \mathbf{N}_{i,11} \mathbf{B}_i + \mathbf{N}_{i,21} \mathbf{B}_i + \mathbf{B}_i^{\top} \mathbf{N}_{i,21}^{\top} + \mathbf{N}_{i,22}, \quad (29c)$$

$$\mathbf{q}_{i,1} = \mathbf{A}_i^\top \mathbf{n}_{i,1},\tag{29d}$$

$$\mathbf{q}_{i,2} = \mathbf{B}_i^{\top} \mathbf{n}_{i,1} + \mathbf{n}_{i,2}. \tag{29e}$$

It is obvious that Eq. (28) has a closed-form solution

$$\Delta \mathbf{\Omega}_i = \mathbf{K}_i \Delta \mathbf{x}_{\text{par}(i)} + \mathbf{k}_i, \qquad (30)$$

in which

$$\mathbf{K}_{i} = -\mathbf{Q}_{i,22}^{-1}\mathbf{Q}_{i,21}$$
 and $\mathbf{k}_{i} = -\mathbf{Q}_{i,22}^{-1}\mathbf{q}_{i,2}$. (31)

If we use Eq. (30) to eliminate $\Delta \Omega_i$ in Eq. (28), there exists $\mathbf{M}_i \in \mathbb{R}^{(6+P) \times (6+P)}$, $\mathbf{m}_i \in \mathbb{R}^{6+P}$ and $\Delta \mathbf{E}_i \in \mathbb{R}$ such that

$$\mathcal{E}_{i}(\Delta \mathbf{x}_{\mathrm{par}(i)}) = \frac{1}{2} \Delta \mathbf{x}_{\mathrm{par}(i)}^{\top} \mathbf{M}_{i} \Delta \mathbf{x}_{\mathrm{par}(i)} + \mathbf{m}_{i}^{\top} \Delta \mathbf{x}_{\mathrm{par}(i)} + \Delta \mathbf{E}_{i}, \quad (32)$$

in which

$$\mathbf{M}_{i} = \mathbf{Q}_{i,11} - \mathbf{Q}_{i,21}^{\top} \mathbf{Q}_{i,22}^{-1} \mathbf{Q}_{i,21}, \qquad (33a)$$

	Dense Unconstrained Formulation	Sparse Constrained Formulation
Step 1	O(N(6+3K+P))	O(K(9+P)) + O(N(9+P))
Step 2	$O(N(6+3K+P)^2)$	$O(N(9+P)^2)$
Step 3	$O\bigl((6+3K+P)^3\bigr)$	$O(K(9+P)^2) + O((6+P)^3)$
Total	$O((6+3K+P)^3) + O(N(6+3K+P)^2)$	$O(K(9+P)^2) + O((6+P)^3) + O(N(9+P)^2)$

	· · ·
	<u>م</u>
	<u> </u>
۰.	u,

	Dense Unconstrained Formulation	Sparse Constrained Formulation	
Step 1	O(KN)	O(K) + O(N)	
Step 2	$O(K^2N)$	O(N)	
Step 3	$O(K^3)$	O(K)	
Total	$O(K^3) + O(K^2N)$	O(K) + O(N)	
(b)			

Table 2: The summary of the computational complexities for the steps to compute the Gauss-Newton direction for the dense unconstrained and sparse constrained formulations, in which K is the number of joints, P is the number of shape parameters, N is the number of measurements for all the body parts. Note that the number of shape parameters P is assumed to be varying in (a) and constant in (b).

$$\mathbf{m}_i = \mathbf{q}_{i,1} - \mathbf{Q}_{i,21}^\top \mathbf{Q}_{i,22}^{-1} \mathbf{q}_{i,2}, \qquad (33b)$$

$$\Delta \mathbf{E}_i = \Delta \overline{\mathbf{E}}_i - \frac{1}{2} \mathbf{q}_{i,2}^{\top} \mathbf{Q}_{i,22}^{-1} \mathbf{q}_{i,2}.$$
 (33c)

Therefore, if there exists $\mathbf{M}_j \in \mathbb{R}^{(6+P)\times(6+P)}$, $\mathbf{m}_j \in \mathbb{R}^{6+P}$ and $\Delta \mathbf{E}_j \in \mathbb{R}$ for all $j \in \operatorname{chd}(i)$ such that Eq. (25) holds, we might further obtain $\mathbf{M}_i \in \mathbb{R}^{(6+P)\times(6+P)}$, $\mathbf{m}_i \in \mathbb{R}^{6+P}$ and $\Delta \mathbf{E}_i \in \mathbb{R}$ with which $\mathcal{E}_i(\Delta \mathbf{x}_{\operatorname{par}(i)})$ can be written as Eq. (32).

In the kinematic tree, a body part *i* at the leaf node has no children, for which Eq. (27) is simplified to $\mathbf{N}_{i,11} = \mathbf{H}_{i,11}$, $\mathbf{N}_{i,21} = \mathbf{H}_{i,21}$, $\mathbf{N}_{i,22} = \mathbf{H}_{i,22}$, $\mathbf{n}_{i,1} = \mathbf{g}_{i,1}$, $\mathbf{n}_{i,2} = \mathbf{g}_{i,2}$ and $\Delta \mathbf{\bar{E}}_i = 0$, then, it is possible to recursively compute $\mathbf{M}_i \in \mathbb{R}^{(6+P)\times(6+P)}$, $\mathbf{m}_i \in \mathbb{R}^{6+P}$ and $\Delta \mathbf{E}_i \in \mathbb{R}$ for each $i = 1, \dots, K$ following Eqs. (27), (29) and (33) through the bottom-up traversal of kinematic tree.

It is by definition that Ω_0 is a dummy variable and $\Delta \Omega_0 = 0$. Thus, if $\mathcal{E}_i(\Delta \mathbf{x}_0)$ in Eq. (32) is known for each $i \in \text{chd}(0)$, Eq. (19) is equivalent to an unconstrained optimization problem on $\Delta \mathbf{x}_0 \in \mathbb{R}^{6+P}$:

$$\min_{\Delta \mathbf{x}_0} \frac{1}{2} \Delta \mathbf{x}_0^\top \mathbf{H}_{0,11} \Delta \mathbf{x}_0 + \mathbf{g}_{0,1}^\top \Delta \mathbf{x}_0 + \sum_{j \in \text{chd}(0)} \mathcal{E}_i(\Delta \mathbf{x}_0).$$

From Eq. (32), the equation above is equivalent to

$$\min_{\Delta \mathbf{x}_0} \frac{1}{2} \Delta \mathbf{x}_0^\top \mathbf{M}_0 \Delta \mathbf{x}_0 + \mathbf{m}_0^\top \mathbf{x}_0 + \Delta \overline{\mathbf{E}}_0 \qquad (34)$$

in which

$$\mathbf{M}_0 = \mathbf{H}_{0,11} + \sum_{j \in \mathrm{chd}(0)} \mathbf{M}_j, \tag{35a}$$

$$\mathbf{m}_0 = \mathbf{g}_{0,1} + \sum_{j \in \text{chd}(0)} \mathbf{m}_j, \qquad (35b)$$

$$\Delta \overline{\mathsf{E}}_0 = \sum_{i \in \mathrm{chd}(0)} \Delta \mathsf{E}_i. \tag{35c}$$

It is straightforward to show that

$$\Delta \mathbf{x}_0 = -\mathbf{M}_0^{-1} \mathbf{m}_0 \tag{36}$$

solves Eq. (34) with

$$\Delta \mathbf{E}_0 = \Delta \overline{\mathbf{E}}_0 - \frac{1}{2} \mathbf{m}_0^\top \mathbf{M}_0^{-1} \mathbf{m}_0$$
(37)

to be the expected cost reduction as well as the minimum objective value of Eq. (19).

At last, we recursively compute $\{\Delta \mathbf{x}_i, \Delta \mathbf{\Omega}_i\}_{i=1}^K$ using Eqs. (13), (30) and (31) through a top-down traversal of the kinematics tree, from which the Gauss-Newton direction is exactly retrieved.

From our analysis, the resulting algorithm to solve Eq. (19) and compute the Gauss-Newton direction is summarized in Algorithm 1. In the next subsection, we show that Algorithm 1 scales linearly with respect to the number of joints.

	Dense Unconstrained Formulation	Sparse Constrained Formulation
Step 1	 (a) It takes O(N_i(6 + 3K + P)) time to compute J_i ∈ ℝ^{N_i×(6+3K+P)} in Eq. (11) for each i = 0,, K. (b) In total, it takes O(N(6+3K+P)) time to compute J_i ∈ ℝ^{N_i×(6+3K+P)} for all i = 0,, K. 	 (a) It takes O(9 + P) time to compute A_i ∈ R^{(6+P)×(6+P)} and B_i ∈ R^{(6+P)×3} in Eqs. (16) and (17) for each i = 0,, K. Note that the bottom of A_i and B_i in Eqs. (16) and (17) are either zero or identity matrices, which simplifies the computation. (b) It takes O(N_i(9 + P)) time to compute J_{i,1} ∈ R^{N_i×(9+P)} and J_{i,2} ∈ R^{N_i×3} in Eqs. (14) and (15) for each i = 0,, K. (c) Note that J_{i,1}, J_{i,2}, A_i and B_i are intermediates to compute J_i in Eq. (11) using the chain rule. (d) In total, it takes O(K(9 + P) + O(N(9 + P))) time to compute J_{i,1}, J_{i,2}, A_i and B_i for all i = 0,, K.
Step 2	(a) It takes $O(N_i(6+3K+P)^2)$ to compute $\mathbf{J}_i^{\top} \mathbf{J}_i \in \mathbb{R}^{(6+3K+P)\times(6+3K+P)}$ for each $i = 0, \dots, K$. (b) In total, it takes $O(N(6+3K+P)^2)$ time to compute $\mathbf{H} = \sum_{i=0}^{K} \mathbf{J}_i^{\top} \mathbf{J}_i \in \mathbb{R}^{(6+3K+P)\times(6+3K+P)}$ in Eq. (18).	 (a) It takes O (N_i(9 + P)²) time to compute H_{i,11} ∈ ℝ^{(6+P)×(6+P)}, H_{i,21} ∈ ℝ^{3×(6+P)} and H_{i,22} ∈ ℝ^{3×3} in Eq. (19) for each i = 0,, K. (b) In total, it takes O(N(9 + P)²) time to compute H_{i,11} ∈ ℝ^{(6+P)×(6+P)}, H_{i,21} ∈ ℝ^{3×(6+P)} and H_{i,22} ∈ ℝ^{3×3} for all i = 0,, K.
Step 3	(a) In total, it takes $O((6+3K+P)^3)$ to compute the matrix inverse of $\mathbf{H} \in \mathbb{R}^{(6+3P+K)\times(6+3P+K)}$ and solve Eq. (20).	 (a) It takes O((9 + P)²) time to run lines 2-17 and lines 26-27 in Algorithm 1 for each i = 1,, K. Note that A_i and B_i in Eqs. (16) and (17) are zero and identity matrices at the bottom, which can be exploited to simplify the computation. (b) It takes O((6 + P)³) time to compute the matrix inverse of M₀ ∈ ℝ^{(6+P)×(6+P)} in line 23 of Algorithm 1. (c) In total, it takes O(K(9+P)²) + O((6+P)³) to compute the Gauss-Newton direction.
Total	The overall complexity is $O((6+3K+P)^3)+O(N(6+3K+P)^2)$.	The overall complexity is $O(K(9+P)^2) + O((6+P)^3) + O(N(9+P)^2)$.

Table 3: The analysis of the computational complexities for the steps to compute the Gauss-Newton direction for the dense unconstrained and sparse constrained formulations. In this table, K is the number of joints, P is the number of shape parameters, N is the number of measurements for all the body parts, and N_i is the number of measurements associated with body part i.

A.2.4 Complexity Analysis

In Table 2, we present a short summary of the computational complexities for each step to compute the Gauss-Newton direction, and in Table 3, we present a comprehensive analysis of the computational complexities that leads to results in Table 2. The analysis also proves the complexity conclusions in Proposition 2.

In Tables 2 and 3, it can be concluded that our sparse constrained formulation is O(K) times faster for Step 1,



Figure 1: The CPU time ratio of the SMPL+H and SMPL models to compute the Gauss-Newton direction with (a) different numbers of measurements and no shape parameters, (b) different numbers of measurements and 10 shape parameters, and (c) different numbers of shape parameters. The SMPL and SMPL+H models have K = 23 and K = 51 joints, respectively. In Figs. 1 (a) to 1(c), the solid lines denote the actual CPU time ratio of the SMPL+H and SMPL models that is obtained from the experiments, whereas the dashed lines denote the expected CPU time ratio that is approximated from the complexity analysis in Tables 2 and 3. It can be seen the impact of the number of joints is around two orders of magnitude less on our method.

and $O(K^2)$ times for Steps 2 and 3 than the dense unconstrained formulation in terms of the number of joints K. In total, our sparse constrained formulation scales linearly with respect to the number of joints instead of cubically as the dense unconstrained formulation.

Furthermore, in terms of the number of measurements N, Tables 2 and 3 indicate that the complexity of our sparse constrained formulation is $O(N(9 + P)^2)$ or O(N), whereas that of the dense constrained formulation is $O(N(6 + 3K + P)^2)$ or $O(K^2N)$. This suggests that our sparse constrained formulation has the the number of joints K and measurements N decoupled in the computation, and as a result, is much more efficient to handle optimization problems with more measurements. Note that it is common in [3, 9, 14, 15, 16, 21] to introduce extra measurements to improve the estimation accuracy.

B. Ablation Studies

In addition to the results of ablation studies in the paper, we present a more complete analysis on the impact of the number of joints K, the number of measurements N, and the number of shape parameters P on the computation of the Gauss-Newton direction.

B.1. Experiments

As mentioned in the paper, the CPU time to compute the Gauss-Newton direction w/ and w/o our method is recorded for the SMPL and SMPL+H models in the following experiments.

Experiment 1. The number of shape parameters P is 0 and the number of measurements N increases from 120 to 600 for both of the SMPL and SMPL+H models.

Experiment 2. The number of shape parameters P is 10 and the number of measurements N increases from 120 to

600 for both of the SMPL and SMPL+H models.

Experiment 3. The number of shape parameters P increases from 0 to 10, and each joint of the SMPL and SMPL+H models is assigned with a 2D keypoint, a 3D keypoint, and a part orientation field as measurements.

B.2. Number of the Joints

The CPU time ratio of the SMPL+H and SMPL models to compute the Gauss-Newton direction is used as the metric to evaluate the impact of the number of joints K. Note that the SMPL and SMPL+H models have K = 23 and K = 51 joints, respectively. The CPU time ratio reflects the additional time induced as a result of the more joints on the SMPL+H model. The CPU time ratios of the three experiments are reported in Fig. 1 and discussed as follows:

 In Experiment 1, there are no shape parameters and the computation of the Gauss-Newton direction is dominated by the number of measurements N. From Tables 2 and 3, it is known that our method has O(N) complexity, which is not related with the number of joints K, and thus, the expected CPU time ratio with our method should be

$$\frac{1}{1} = 1.$$

In contrast, the CPU time without our method is approximately $O((3K+6)^2)$, which suggests an expected CPU time ratio of

$$\left(\frac{3\times51+6}{3\times23+6}\right)^2 = 4.49.$$

The numbers of 1 and 4.49 in the two equations above are consistent with the results in Fig. 1(a).



Figure 2: The computation of the Gauss-Newton direction with different numbers of measurements and no shape parameters. The results are (a) the CPU time with and without our method on the SMPL and SMPL+H models, and (b) the speedup of our method on the SMPL and SMPL+H models, and (c) the speed up of our method on the SMPL model, and (d) the speed up of our method on the SMPL+H model.



Figure 3: The computation of the Gauss-Newton direction with different numbers of measurements and 10 shape parameters. The results are (a) the CPU time with and without our method on the SMPL and SMPL+H models, and (b) the speedup of our method on the SMPL and SMPL+H model, and (d) the speedup of our method on the SMPL+H model.

 In Experiment 2, there are 10 shape parameters. However, the analysis is still similar to that of Experiment
 From Tables 2 and 3, the expected CPU time ratio of the SMPL+H and SMPL models w/ and w/o our method should be around

$$\frac{1}{1} =$$

and

$$\left(\frac{3\times51+6+10}{3\times23+6+10}\right)^2 = 3.95$$

1

respectively, which is consistent with the results in Fig. 1(b).

3. In Experiment 3, the number of measurements N is proportional to the number of joints of the SMPL and SMPL+H models. Then, as a result of Tables 2 and 3, the CPU time w/ and w/o our method to compute the Gauss-Newton direction should be around O(K) and $O((3K+6)^3)$, respectively, and the corresponding expected CPU time can be also approximated by

and

$$\left(\frac{3\times51+6}{3\times23+6}\right)^3 = 9.53,$$

 $\frac{51}{23} = 2.22$

which is consistent with the results in Fig. 1(c).

4. From Fig. 1 and the discussions above, it can be further concluded that the number of joints has around $O(K^2)$ times less impact on our method, which suggests that our sparse constrained formulation is more suitable for human models with more joints.

B.3. Number of the Measurements

The CPU time w/ and w/o our method to compute the Gauss-Newton direction and the corresponding speedup in Experiments 1 and 2 are reported in Figs. 2 and 3. It can be seen from Figs. 2 and 3 that our method has $4.73 \sim 13.91x$ speedup on the SMPL model and a $12.17 \sim 43.24x$ speedup on the SMPL+H model. Furthermore, no matter whether there are shape parameters or not, the speedup of our method is greater as the number of measurements increases, which means that our sparse constrained formulation is more efficient to solve optimization problems with more more measurements.

B.4. Number of the Shape Parameters

The CPU time w/ and w/o our method to compute the Gauss-Newton direction and the corresponding speedup in Experiment 3 are reported in Fig. 4. It can be seen from Fig. 4 that our method has a $4.92 \sim 7.78x$ speedup on the SMPL model and a $18.63 \sim 34.18x$ speedup on the SMPL+H model, which is consistent with the analysis that our sparse constrained formulation has better scalability on



Figure 4: The computation of the Gauss-Newton direction with different number of shape parameters. The results are (a) the CPU time with and without our method on the SMPL and SMPL+H models, and (b) the speedup of our method on the SMPL and SMPL+H models, and (c) the speed up of our method on the SMPL model, and (d) the speed up of our method on the SMPL+H model.

human models with more joints. On the SMPL+H model, the CPU time taken to compute the Gauss-Newton direction without our method is as many as 2.5 ms, which is difficult to be used in real time considering that most optimization methods need around $20 \sim 30$ iterations to converge. As a comparison, our method is significantly faster on both of the SMPL and SMPL+H models, for which the CPU time is $0.027 \sim 0.13$ ms. In particular, note that if there are no shape parameters, our method has a further acceleration of the computation—this has is important for real-time video tracking of 3D human pose and shape, in which the shape parameters that are estimated from the first few frames can be reused.

C. Qualitative Results

In this section, we present more qualitative comparisons with SPIN [8] and SMPLify [3] on the Human3.6M, MPI-INF-3DHP and 3DPW datasets. The results are shown in Figs. 5 to 7.

D. Real-Time Motion Capture Framework

D.1. Human Detection

The YOLOv4-CSP [2,20] is used for human detection to make a balance between accuracy and efficiency. The size of input images for YOLOv4-CSP is 512×512 .

D.2. 2D Keypoint Estimation

The AlphaPose [5] is used for 2D keypoint estimation with 256×192 input images. The following datasets are used to train AlphaPose.

Human3.6M [4, 6] is a popular dataset for 3D human pose estimation. Following the standard training-testing protocol in [17], we use subjects S1, S5-S8 for training.

MPI-INF-3DHP [13] is a multi-view markerless dataset with 8 training subjects and 6 test subjects. We use subjects S1-S8 that are downsampled to 10 FPS for training. **COCO** [10] is a large-scale dataset for 2D joint detection. We use the COCO training datasets for training.

MPII [1] is a 2D human pose dataset that is extracted from online videos. We use the MPII training datasets for training.

D.3. 3D Keypoint Regression

In our real-time motion capture framework, we use a light-weight fully connected neural network for 2D-to-3D lifting. The 3D Keypoint regression network can be regarded as a modification of VideoPose3D [18]. From the 3D keypoint regression network, we further obtain the part orientation field [21] for each body part. We use the training datasets of Human3.6M [6] and MPI-INF-3DHP [13] that are downsampled to 10 FPS to train the 3D keypoint regression network.

E. Prior Loss of Joint States

We use the normalizing flow [7] to describe the joint state prior loss $E_{\Omega,i}$. The normalizing flow is trained on the AMASS dataset [12] and has the structure of FC6 \rightarrow PReLU \rightarrow FC6 whose input is the 6D representation of rotation. We remark that the normalizing flow structure above to learn admissible joint states is inspired by the work of [22].

F. Implementation

F.1. Overview

While originally designed for 3D human pose and shape estimation, we emphasize that our method can be extended to any types of articulated tracking problems in computer vision and robotics [19]. The only requirement is that the objective can be written as

$$\mathbf{E} = \sum_{0 \le i \le K} \frac{1}{2} \| \mathbf{r}_i(\mathbf{T}_i, \mathbf{\Omega}_i, \boldsymbol{\beta}) \|^2,$$
(38)



Figure 5: Qualitative comparisons of our method (second row in pink), SPIN [8] (third row in gray), and SMPLify [3] (fourth row in purple) on the Human3.6M dataset.

in which K is the number joints, \mathbf{T}_i is the pose of body part i, Ω_i is the joint state and β is the shape parameters. Empirically, such a requirement can be satisfied with ease, e.g., we might assume that the keypoints selected to calculate the losses are rigidly attached to a single body part. As a matter of fact, as long as the objective is in the form of Eq. (38), the steps to compute the Gauss-Newton direction in Table 1 and the complexity analysis in Tables 2 and 3 hold as well. Thus, there are no difficulties to implement our method on practical articulated tracking problems.

F.2. Extract S_i and l_i from the SMPL Model

At the rest pose of the SMPL model [11], it is known that the joint positions linearly depend on the vertex positions, and the vertex positions also linearly depend on the shape parameters $\beta \in \mathbb{R}^{P}$. Thus, we conclude that the joint positions $\overline{\mathbf{t}}_{i} \in \mathbb{R}^{3}$ at the rest pose linearly depend on the shape parameters, i.e., there exists $\mathcal{J}_{i} \in \mathbb{R}^{3 \times P}$ and $\mathbf{c}_{i} \in \mathbb{R}^{3}$ in the SMPL model such that $\overline{\mathbf{t}}_{i}$ at the rest pose takes the form of

$$\overline{\mathbf{t}}_i = \mathcal{J}_i \cdot \boldsymbol{\beta} + \mathbf{c}_i. \tag{39}$$

Note that joint position $\mathbf{t}_i \in \mathbb{R}^3$ is also the translation of pose $\mathbf{T}_i = \begin{bmatrix} \mathbf{R}_i & \mathbf{t}_i \\ \mathbf{0} & 1 \end{bmatrix} \in SE(3)$ where $\mathbf{R}_i \in SO(3)$ is the

rotation. Moreover, the relative joint position $\Delta \overline{\mathbf{t}}_i \in \mathbb{R}^3$ between any connected body parts is constant, and thus, we obtain $\Delta \overline{\mathbf{t}}_i = \overline{\mathbf{t}}_i - \overline{\mathbf{t}}_{\text{par}(i)}$, in which par(i) denotes the index of the parent of body part *i*. Then, joint position $\mathbf{t}_i \in \mathbb{R}^3$ at any poses satisfies

$$\mathbf{t}_{i} = \mathbf{R}_{\mathrm{par}(i)} \Delta \bar{\mathbf{t}} + \mathbf{t}_{\mathrm{par}(i)} = \mathbf{R}_{\mathrm{par}(i)} \left(\bar{\mathbf{t}}_{i} - \bar{\mathbf{t}}_{\mathrm{par}(i)} \right) + \mathbf{t}_{\mathrm{par}(i)}.$$
(40)

In the equation above, $\mathbf{R}_{\text{par}(i)}$ is rotation of pose $\mathbf{T}_{\text{par}(i)} \in SE(3)$. Substituting Eq. (39) into Eq. (40) to cancel out $\mathbf{\bar{t}}_i$ and $\mathbf{\bar{t}}_{\text{par}(i)}$, we obtain

$$\mathbf{t}_{i} = \mathbf{R}_{\mathrm{par}(i)} \left(\mathcal{S}_{i} \cdot \boldsymbol{\beta} + \mathbf{l}_{i} \right) + \mathbf{t}_{\mathrm{par}(i)}, \tag{41}$$

in which

and

$$S_i = \mathcal{J}_i - \mathcal{J}_{\text{par}(i)} \in \mathbb{R}^{3 \times P}$$
 (42)

$$\mathbf{l}_i = \mathbf{c}_i - \mathbf{c}_{\mathrm{par}(i)} \in \mathbb{R}^3.$$
(43)

It is immediate to show that $S_i \cdot \beta + l_i$ is the relative joint position between body parts *i* and par(*i*), and thus, the corresponding relative pose $\mathbf{T}_{\text{par}(i),i}$ is

$$\mathbf{T}_{\mathrm{par}(i),i} \triangleq \begin{bmatrix} \mathbf{\Omega}_i & \mathcal{S}_i \cdot \boldsymbol{\beta} + \mathbf{l}_i \\ \mathbf{0} & 1 \end{bmatrix}, \quad (44)$$

in which $\Omega_i \in SO(3)$ is the state of joint *i*.



Figure 6: Qualitative comparisons of our method (second row in pink), SPIN [8] (third row in gray), and SMPLify [3] (fourth row in purple) on the MPI-INF-3DHP dataset.



Figure 7: Qualitative comparisons of our method (second row in pink), SPIN [8] (third row in gray), and SMPLify [3] (fourth row in purple) on the MPI-INF-3DHP dataset.

References

- Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2D human pose estimation: New benchmark and state of the art analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2014. 9
- [2] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. YOLOv4: Optimal speed and accuracy of object detection. arXiv preprint arXiv:2004.10934, 2020. 9
- [3] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J. Black. Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. In *European conference on computer vision (ECCV)*, 2016. 7, 9, 10, 11
- [4] Cristian Sminchisescu Catalin Ionescu, Fuxin Li. Latent structured models for human pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2220–2227, 2011. 9
- [5] Hao-Shu Fang, Shuqin Xie, Yu-Wing Tai, and Cewu Lu. RMPE: Regional multi-person pose estimation. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 2334–2343, 2017.
- [6] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6M: Large scale datasets and predictive methods for 3D human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7):1325–1339, 2013. 9
- [7] Ivan Kobyzev, Simon Prince, and Marcus Brubaker. Normalizing flows: An introduction and review of current methods. *IEEE Transactions on Pattern Analysis* and Machine Intelligence, 2020. 9
- [8] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In *Proceedings of the IEEE/CVF Conference on Computer Vision*, 2019. 9, 10, 11
- [9] Christoph Lassner, Javier Romero, Martin Kiefel, Federica Bogo, Michael J. Black, and Peter V. Gehler. Unite the People: Closing the loop between 3D and 2D human representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017. 7
- [10] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 9

- [11] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. ACM Trans. Graphics (Proc. SIGGRAPH Asia), 34(6):248:1– 248:16, Oct. 2015. 10
- [12] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. AMASS: Archive of motion capture as surface shapes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5442–5451, Oct. 2019. 9
- [13] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. Monocular 3D human pose estimation in the wild using improved cnn supervision. In *Proceedings of the International Conference on 3D Vision*. IEEE, 2017. 9
- [14] Dushyant Mehta, Oleksandr Sotnychenko, Franziska Mueller, Weipeng Xu, Mohamed Elgharib, Pascal Fua, Hans-Peter Seidel, Helge Rhodin, Gerard Pons-Moll, and Christian Theobalt. XNect: Real-time multi-person 3D motion capture with a single RGB camera. ACM Transactions on Graphics (TOG), 39(4):82–1, 2020. 7
- [15] Dushyant Mehta, Srinath Sridhar, Oleksandr Sotnychenko, Helge Rhodin, Mohammad Shafiei, Hans-Peter Seidel, Weipeng Xu, Dan Casas, and Christian Theobalt. VNect: Real-time 3D human pose estimation with a single RGB camera. ACM Transactions on Graphics (TOG), 36(4):1–14, 2017. 7
- [16] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3D hands, face, and body from a single image. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019. 7
- [17] Georgios Pavlakos, Xiaowei Zhou, Konstantinos G Derpanis, and Kostas Daniilidis. Coarse-to-fine volumetric prediction for single-image 3D human pose. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 7025–7034, 2017. 9
- [18] Dario Pavllo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3D human pose estimation in video with temporal convolutions and semisupervised training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019. 9
- [19] Tanner Schmidt, Richard Newcombe, and Dieter Fox. DART: dense articulated real-time tracking with consumer depth cameras. *Autonomous Robots*, 39(3):239–258, 2015. 9

- [20] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Scaled-YOLOv4: Scaling cross stage partial network. arXiv preprint arXiv:2011.08036, 2020. 9
- [21] Donglai Xiang, Hanbyul Joo, and Yaser Sheikh. Monocular total capture: Posing face, body, and hands in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10965–10974, 2019. 7, 9
- [22] Andrei Zanfir, Eduard Gabriel Bazavan, Hongyi Xu, William T Freeman, Rahul Sukthankar, and Cristian Sminchisescu. Weakly supervised 3d human pose and shape reconstruction with normalizing flows. In *European Conference on Computer Vision*, 2020. 9