

# The Vitruvian Manifold: Inferring Dense Correspondences for One-Shot Human Pose Estimation

Jonathan Taylor<sup>†\*</sup>

Jamie Shotton<sup>†</sup>

Toby Sharp<sup>†</sup>

Andrew Fitzgibbon<sup>†</sup>

<sup>†</sup>Microsoft Research Cambridge

<sup>\*</sup>University of Toronto

The supplementary material for our CVPR 2012 paper *The Vitruvian Manifold: Inferring Dense Correspondences for One-Shot Human Pose Estimation* comprises this document and the accompanying video.

## 1. Supplementary Results

The accompanying video illustrates the inferred correspondences and convergence of the optimization algorithm for several poses, and additionally gives demonstrations of how the algorithm performs on extended sequences containing complex motion. In addition to this we include in Fig. 1 a qualitative comparison of the skeletons inferred by the algorithms of [1, 2] and the Microsoft Kinect for Windows software development kit (SDK) to accompany the quantitative results in the main paper.

## 2. Implementation Appendix

While the paper contains a full description of the algorithm, this section includes further details of the specific implementation we used for the interested reader, including parameter settings (Fig. 2) and many of the derivatives of the energy function required by the L-BFGS optimizer.

### 2.1. Prior

Our prior consists of: (i) a Gaussian on pose, learned from motion capture data; and (ii) a heuristic term that discourages self-intersection. These each provide a modest quantifiable increase in accuracy, and qualitatively encourage more realistic poses.

The Gaussian over pose has full covariance over the rotation degrees of freedom and a single independent element corresponding to the scale degree of freedom. There is no prior over global translation as the actor can appear at any point in the scene. See Fig. 3 for 16 random samples from the rotational component of the Gaussian pose prior (i.e. global scale, rotation and translation is fixed).

### 2.2. Model Parameterization

We explicitly write out the parameters of our model as  $\theta = (s^T, \mathbf{t}^T, \mathbf{q}_0^T, \dots, \mathbf{q}_L^T)^T \in \mathbb{R}^d$  with  $s \in \mathbb{R}$  defining global scale,  $\mathbf{t} \in \mathbb{R}^3$  global translation and each  $\mathbf{q}_l \in \mathbb{R}^4$  a quaternion defining the relative rotation of limb  $l$  if  $l \in \{1, \dots, L\}$  or global rotation if  $l = 0$ . To avoid overloading

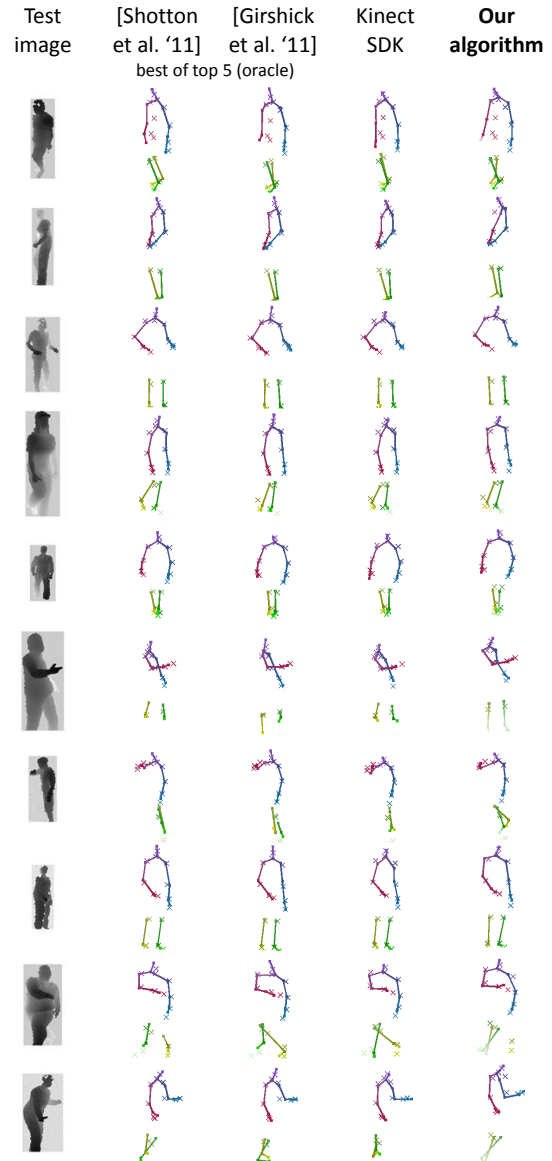


Figure 1. **Qualitative comparison with [2, 1] and the Kinect SDK.** All algorithms produce reasonable results on this random selection of test images. The quantitative results in the main paper, however, demonstrate the marked quantitative improvement in accuracy obtained by our algorithm. (Hips not drawn as [2, 1] do not predict hips, and the Kinect SDK hips are incompatible with our hips. Hips were not part of the metrics for this comparison). Note that the results from [2, 1] use the ‘best of top 5’ oracle and are not a realizable algorithm.

Parameter	Description	Value	Explanation
$\lambda_{vis}$	energy term weighting	$\frac{1}{n}$	$n$ is the number of image pixels
$\lambda_{prior}$	energy term weighting	$2 \times 10^{-4}$	
$\lambda_{int}$	energy term weighting	$\frac{1}{ S }$	
$\beta$	visibili sharpness	100	
$\gamma$	intersection sharpness	1000	
$r_s$	sphere radii	0.025m	
$m$	number of mesh vertices	$\sim 7000$	
$\eta$	Geman-McClure parameter	0.1m / 10 pixels	meters $\rightarrow$ depth; pixels $\rightarrow$ silhouettes
$\tau$	back-facing penalty	1.0	penalty paid
	number of trees	3	
	maximum depth of trees	20	
$b$	mean shift bandwidth	5 vits	used to cluster correspondences during tree training

Figure 2. Settings for parameters introduced in the main paper and the derivations below.

notation from the main paper, we use a fixed width font to denote these parameters and the transformations they specify. Explicitly, these specify a scaling transformation for  $s$

$$S(s) = \begin{pmatrix} s & 0 & 0 & 0 \\ 0 & s & 0 & 0 \\ 0 & 0 & s & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}, \quad (1)$$

a translation transformation for  $t = (t_1, t_2, t_3)$

$$T(t) = \begin{pmatrix} 1 & 0 & 0 & t_1 \\ 0 & 1 & 0 & t_2 \\ 0 & 0 & 1 & t_3 \\ 0 & 0 & 0 & 1 \end{pmatrix}, \quad (2)$$

and for each vector  $q_l$ , a rotation  $R(\vartheta(q_l))$  where  $\vartheta(q) = \frac{q}{\|q\|}$  normalizes  $q$  and for a unit quaternion  $p = (x, y, z, w)$

$$R(p) = \begin{pmatrix} 1 - 2(y^2 + z^2) & 2(xy - zw) & 2(xz + yw) & 0 \\ 2(xy + zw) & 1 - 2(x^2 + z^2) & 2(yz - xw) & 0 \\ 2(xz - yw) & 2(yz + xw) & 1 - 2(x^2 + y^2) & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}. \quad (3)$$

We can explicitly define the remaining quantities in equation (1) and (2) of the main paper by

$$R_{glob}(\theta) = T(t)R(\vartheta(q_0))S(s) \quad (4)$$

$$R_l(\theta) = T(t_l)R(\vartheta(q_l)) \quad (5)$$

where  $t_l$  for  $l \in \{1, \dots, L\}$  are fixed. It is readily seen that all transformations preserve unity in the 4th homogeneous coordinate, a property assumed in each vertex position  $p_i$ , and thus we can assume that  $\pi = P = I_{3 \times 4}$ , a matrix that drops the fourth coordinate. Further, dropping the

distinction between model vertices  $v_j$  and predicted correspondences  $u_i$ , we have

$$M(u_i; \theta) = P \sum_{k=1}^K \alpha_{ik} T_{l_{ik}}(\theta) T_{l_{ik}}^{-1}(\theta_0) p_i \quad (6)$$

For each, vertex  $u_i$  let  $n(u_i)$  denote it's surface normal in the base pose  $\theta_0$ , perhaps computed using the local mesh geometry. Then an appropriate way to define the surface normal in a new pose  $\theta$  is

$$N(u_i; \theta) = \frac{\hat{N}(u_i; \theta)}{\|\hat{N}(u_i; \theta)\|_2} \quad (7)$$

$$\hat{N}(u_i; \theta) = P \sum_{k=1}^K \alpha_{ik} \hat{T}_{l_{ik}}(\theta) \hat{T}_{l_{ik}}^{-1}(\theta_0) n(u_i), \quad (8)$$

where  $\hat{T}_l$  is composed only of the rotational components of  $T_l$ . That is,  $n(u_i)$  is simply rotated by each limb's transformation, and the resulting combination normalized.

### 2.3. Gradient

The L-BFGS optimizer requires the gradient of the energy function to be provided. This could be computed using finite differences, but this would be inexact and needlessly slow, especially considering that the Jacobian of vertex coordinates with respect to the optimization parameters is particularly sparse. Therefore it is of benefit to compute the derivatives analytically so that they can be calculated directly. Below, we provide a recipe for doing this with a simplified energy function consisting of just  $E'_{vis}$  with pixel weights and  $\beta$  set to unity:

$$E'_{\text{vis}} = \sum_{i=1}^n V_i(\theta) \cdot \rho(e_i(\theta)) + (1 - V_i(\theta)) \cdot \tau \quad (9)$$

$$= n\tau + \sum_{i=1}^n V_i(\theta) \cdot [\rho(e_i(\theta)) - \tau], \quad (10)$$

where  $e_i(\theta) = \|x_i - M(u_i; \theta)\|_2$ , Geman-McClure function  $\rho(e) = \frac{e^2}{e^2 + \eta^2}$ , and visibility weighting  $V_i(\theta) = \sigma(-N(u_i; \theta)^\top A)$  where  $\sigma(x)$  is the sigmoid function. Taking the derivative with respect to the  $j^{\text{th}}$  component<sup>1</sup> of the parameter vector  $\theta$  we get:

$$\frac{\partial E(\theta)}{\partial \theta_j} = \sum_{i=1}^n \frac{\partial \rho(e_i(\theta))}{\partial \theta_j} V_i(\theta) + \frac{\partial V_i(\theta)}{\partial \theta_j} \cdot [\rho(e_i(\theta)) - \tau]. \quad (11)$$

Working on the first term, observe that

$$\frac{\partial \rho(e_i(\theta))}{\partial \theta_j} = \frac{2e_i(\theta)\eta^2}{(e_i(\theta)^2 + \eta^2)^2} \frac{\partial e_i(\theta)}{\partial \theta_j}. \quad (12)$$

and

$$\frac{\partial e_i(\theta)}{\partial \theta_j} = \frac{1}{e_i(\theta)} (M(u_i; \theta) - x_i)^\top \frac{\partial M(u_i; \theta)}{\partial \theta_j}. \quad (13)$$

Further:

$$\frac{\partial M(u_i; \theta)}{\partial \theta_j} = P \sum_{k=1}^K \alpha_{ik} \frac{\partial T_{i_k}(\theta)}{\partial \theta_j} T_{i_k}^{-1}(\theta_0) p_i, \quad (14)$$

Note that  $T_{i_k}(\theta)$  is composed only of matrices of the form  $T(\mathbf{t})$ ,  $R(\mathbf{q})$  and  $S(\mathbf{s})$ , of which only one will depend on  $\theta_j$ . Thus, it is enough to be able to compute derivatives for these individually. Specifically:

$$\frac{\partial T(\mathbf{t})}{\partial \mathbf{t}_h} = \begin{pmatrix} 0 & 0 & 0 & \delta_{h1} \\ 0 & 0 & 0 & \delta_{h2} \\ 0 & 0 & 0 & \delta_{h3} \\ 0 & 0 & 0 & 0 \end{pmatrix}, \quad (15)$$

(using the Kronecker  $\delta$ ), and

$$\frac{\partial S(\mathbf{s})}{\partial \mathbf{s}} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}. \quad (16)$$

For a rotation parameter  $\mathbf{q}$ , we write the corresponding unit quaternion as  $\mathbf{p} = \vartheta(\mathbf{q}) = (\mathbf{x}, y, z, w)$ . Using the chain rule, we obtain the total derivative with respect to each component  $\mathbf{q}_j$  of  $\mathbf{q}$ :

$$\frac{\partial R(\vartheta(\mathbf{q}))}{\partial \mathbf{q}_j} = \sum_{m=1}^4 \frac{\partial R(\mathbf{p})}{\partial \mathbf{p}_m} \frac{\partial \vartheta_m(\mathbf{q})}{\partial \mathbf{q}_j} \quad (17)$$

<sup>1</sup>Here we use  $j$  to indicate a component of parameter vector  $\theta$ . In the main paper  $j$  refers to a model vertex index.

where the derivative of the normalization terms are given by

$$\frac{\partial \vartheta_m(\mathbf{q})}{\partial \mathbf{q}_j} = \frac{\delta_{mj}}{\|\mathbf{q}\|_2} - \frac{\mathbf{q}_m \mathbf{q}_j}{\|\mathbf{q}\|_2^3}. \quad (18)$$

and the other terms are easily calculated from equation Eq. 3. For example,

$$\frac{\partial R(\mathbf{p})}{\partial \mathbf{x}} = \begin{pmatrix} 0 & 2y & 2z & 0 \\ 2y & -4x & -2w & 0 \\ 2z & 2w & -4x & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}. \quad (19)$$

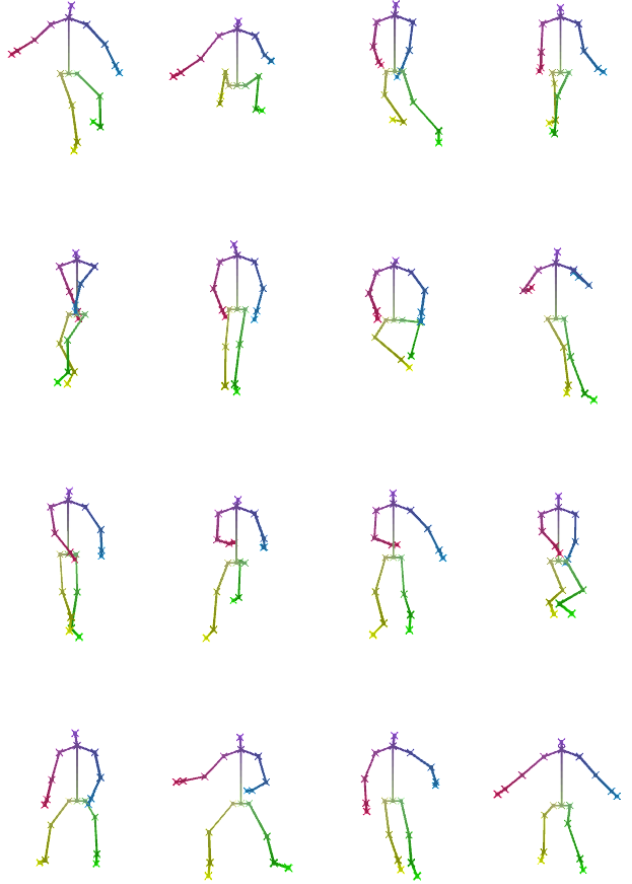


Figure 3. Samples from our Gaussian pose prior.

Turning to the second term on the right hand side of Eq. 11 we see

$$\frac{\partial V_i(\theta)}{\partial \theta_j} = -\sigma(-N_3(u_i; \theta))(1 - \sigma(-N_3(u_i; \theta))) \frac{\partial N_3(u_i; \theta)}{\partial \theta_j} \quad (20)$$

where subscript 3s denote the Z component of the normal vectors, and  $\frac{\partial N_3(u_i; \theta)}{\partial \theta_j}$  can be calculated easily from  $\frac{\partial \hat{N}(u_i; \theta)}{\partial \theta_j}$  as it is simply a normalization (similar to  $\frac{\partial \vartheta(\mathbf{q})}{\partial \mathbf{q}_j}$ )

above) and the calculation of  $\frac{\partial \tilde{N}(u_i; \theta)}{\partial \theta_j}$  is similar to that of  $\frac{\partial M(u_i; \theta)}{\partial \theta_j}$  above.

The remaining gradients of the full energy in the paper can be derived analytically in a similar manner using standard calculus.

## References

- [1] R. Girshick, J. Shotton, P. Kohli, A. Criminisi, and A. Fitzgibbon. Efficient regression of general-activity human poses from depth images. In *Proc. ICCV*, 2011.
- [2] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake. Real-time human pose recognition in parts from a single depth image. In *Proc. CVPR*, 2011.