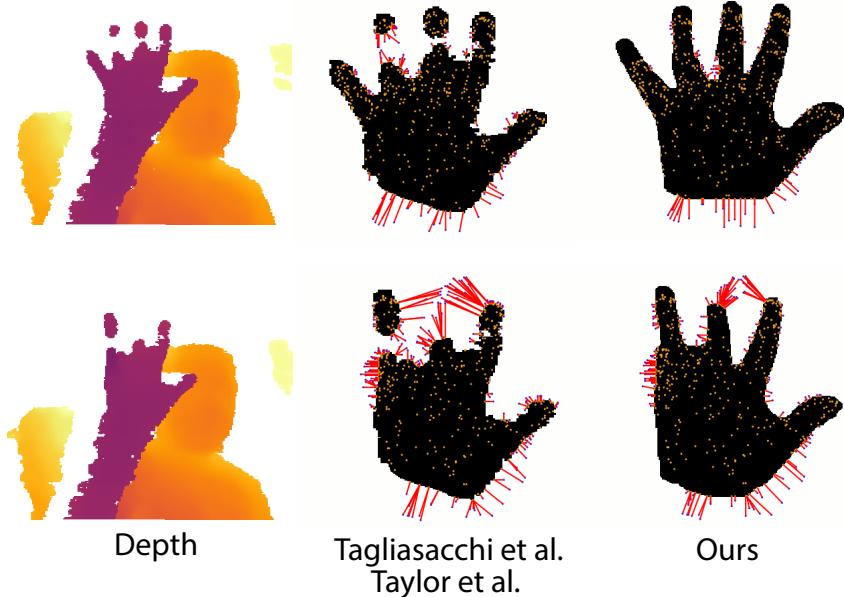


RGB-assisted Depth-based Registration for Accurate Hand Tracking

Submission Id 154



RGB



Depth

Tagliasacchi et al.
Taylor et al.

Ours

Figure 1: We present a method to use an RGB image to assist depth-based registration of a hand mesh. We highlight the limitation of existing methods that use a depth silhouette, where the correspondences (red lines) between the model (orange points) and the observation (black silhouette) are incorrect, leading to incorrect registration. Our RGBD silhouette is more reliable and thus the resulting correspondences are accurate.

ABSTRACT

Hand tracking is an essential component in modern virtual reality systems. Despite the presence of RGBD cameras in such systems, most existing methods have demonstrated hand tracking exclusively from either depth or RGB cameras. We present the first method to simultaneously leverage depth and RGB images for hand tracking. We formulate an energy minimization framework that registers a rigged hand mesh to RGB and depth data. To this end, we propose a novel method to extract a more accurate silhouette by utilizing information from both RGB and depth images. Our method outperforms state-of-the-art hand tracking methods using only RGB or depth information on publicly available datasets. We also demonstrate that our energy minimization framework is flexible enough to adapt to work with only RGB or only depth, or both.

We believe this work lays the foundation for simultaneously using RGB and depth cameras for more accurate shape registration.

CCS CONCEPTS

- Computing methodologies → Reconstruction; Image segmentation.

KEYWORDS

RGBD Hand Tracking, Shape and pose estimation

ACM Reference Format:

Submission Id 154. 2022. RGB-assisted Depth-based Registration for Accurate Hand Tracking. In *Proceedings of 13th Indian Conference on Computer Vision, Graphics and Image Processing (ICVGIP'22)*. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

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).

ICVGIP'22, December 2022, Gandhinagar, India

© 2022 Copyright held by the owner/author(s).

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM.

<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

Hands are the most natural way of interaction for humans. Understanding human hand motion from videos has been quite popular over the past few decades [18, 20, 23, 36, 40]. With the advent of virtual and augmented reality, hand tracking has gained a renewed interest in academia and industry [22, 24, 30, 42].

117 Although modern VR headsets are equipped with RGBD cameras
 118 capable of acquiring RGB and depth images of the user’s hand,
 119 most existing methods exclusively use RGB [2, 3, 6, 15, 39, 43] or
 120 depth [8, 11, 20, 22, 26, 27, 29, 30, 30–33, 35, 36, 38, 41] images as
 121 input. Depth-based methods are more accurate than RGB-based
 122 methods because depth information disambiguates the 3D pose
 123 that might appear similar in the 2D projected space. However, the
 124 quality of depth images often suffers from limited accuracy and
 125 stability due to holes at object boundaries or smooth and shiny
 126 surfaces and inconsistent depth values across frames [16]. As a
 127 result, the depth images cannot provide reliable data required for
 128 accurate registration.

129 We present a method that utilizes information from both RGB
 130 and depth images to achieve a more robust and accurate registration.
 131 (see Fig. 1) Specifically, we register the parametric hand mesh model,
 132 aMANO [9], to the point cloud obtained from the depth image and
 133 silhouette extracted from both RGB and depth images. Our novel
 134 silhouette extractor provides reliable data for accurate registration
 135 even from noisy depth data. Further, our registration method is
 136 adaptable to the available data source, i.e., only RGB, or only depth,
 137 or both.

138 Unlike existing methods [12, 15] that optimize the joint location
 139 alignment, our registration energy also minimizes surface-to-
 140 surface distance (whenever depth is available) and thus achieves a
 141 more accurate registration. When only the RGB image is present,
 142 our registration procedure achieves a more accurate and kinemat-
 143 ically valid pose than state-of-the-art iterative 2D joint location
 144 alignment [15].

145 We summarize our contributions below.

- We present the first method to utilize information from RGB and depth images to achieve a more accurate hand tracking.
- We also propose a novel silhouette extractor that provides reliable silhouette information for robust and accurate registration.

2 RELATED WORK

153 The problem of markerless hand tracking is closely associated with
 154 hand pose estimation and full 3D hand shape and pose reconstruc-
 155 tion. We refer the reader to Armanag et al. [1] for an in-depth
 156 overview of hand pose estimation. In this section, we discuss meth-
 157 ods that reconstruct the geometry of the hand from RGB and depth
 158 images.

2.1 Depth Hand Tracking

163 Oikonomidis et al. [20] introduced one of the early methods for
 164 reconstructing a primitive-based geometric hand model from depth
 165 images using particle swarm optimization (PSO). Qian et al. [22]
 166 demonstrate the use of ICP with PSO for tracking a sphere-based
 167 hand model. Tagliasacchi et al. [30] introduce the idea of using a
 168 gradient-based optimization, Levenberg-Marquardt (LM), to register
 169 a cylinder-based hand model to depth data. Tkach et al. [34] further
 170 extended the idea to a sphere-mesh model. We adapt the energy
 171 terms used in their method and use them in our registration frame-
 172 work. However, unlike these methods, we use a mesh-based model
 173 that can be easily integrable into an existing graphics pipeline.

175 Taylor et al. [33] introduced the first method to model a user-
 176 specific mesh from depth images, and Sharp et al. [26] used it in
 177 their robust hand tracker. Taylor et al. [32] use the hand model of
 178 Khamis et al. [11] to track a subdivided mesh. Further, Shen et al.
 179 [27] demonstrated that one could get away without subdividing the
 180 mesh using ideas from Phong shading. Our barycentric sampling of
 181 the mesh follows this approach. Unfortunately, unlike ours, these
 182 implementations are not publicly available, restricting their usage
 183 for scientific research.

184 Until recently, most methods used a fixed hand template model.
 185 Inspired by the recent advances in the human body shape model
 186 (e.g., SMPL [14]), Romero et al. [24] proposed a learned blend-shape
 187 hand mesh model, MANO. However, it cannot adapt to unseen
 188 hand shapes with substantially large deviations from the training
 189 set. Kalshetti and Chaudhuri [9] introduced adaptive MANO
 190 (aMANO), which extends MANO by incorporating local scale adap-
 191 tation parameters, to be used in a tracking framework Kalshetti
 192 and Chaudhuri [10]. However, these methods require depth images
 193 as input.

2.2 RGB Hand Tracking

196 de La Gorce et al. [5] introduced the idea of analysis-by-synthesis to
 197 track hands from monocular RGB video. However, their generative
 198 method is sensitive to initialization and not robust to any unseen
 199 background. Boukhayma et al. [3] introduce the first end-to-end
 200 deep learning method that regresses MANO shape and pose parame-
 201 ters from RGB images. Zhou et al. [46] propose a real-time method
 202 to regress joint rotations using an inverse kinematics module. Yang
 203 et al. [43] recover the mesh using a multi-stage hourglass network.
 204 Chen et al. [4] reconstruct the hand mesh using self-supervision
 205 from the texture and lighting in the input image along with detected
 206 2D keypoints. Kulon et al. [13] incorporate mesh convolutions for
 207 recovering the hand mesh. However, the accuracy of these discrim-
 208 inative methods suffers on unseen data.

209 Baek et al. [2] introduce the idea of iterative refinement using
 210 neural rendering. However, it assumes the availability of a 2D seg-
 211 mentation mask. Zhang et al. [45] also use a silhouette-based self-
 212 supervision loss. Panteleris et al. [21] use a non-linear least squares
 213 minimization to fit a 3D hand model to 2D keypoints detected us-
 214 ing OpenPose [28]. Mueller et al. [19] use 3D joint locations for
 215 fitting a hand model. Although these hybrid methods improve over
 216 discriminative-based methods, they are not as accurate as depth-
 217 based tracking methods. Our method uses a registration-based
 218 optimization method that leverages RGB and depth data.

3 METHOD

219 Given a sequence of RGB $\{I_c^{(f)}\}_{f=1}^{n_f}$ and corresponding depth $\{I_d^{(f)}\}_{f=1}^{n_f}$
 220 images acquired from a single RGBD camera, our goal is to reg-
 221 ister a hand model, M , to each of these frames. We formulate the
 222 registration problem as an energy minimization such that at each
 223 frame f , the hand model, $M^{(f)}$, explains the observed data in the
 224 RGB image, $I_c^{(f)}$, and the depth image, $I_d^{(f)}$. Specifically, we register
 225 the hand model to the point cloud obtained from the depth image
 226

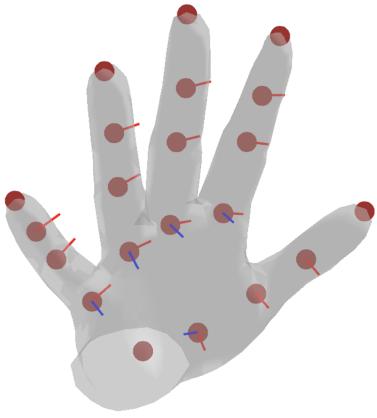


Figure 2: Degrees of freedom (DoF) at each joint: two at the metacarpophalangeal (MCP) joint, one at each proximal interphalangeal (PIP) and distal interphalangeal (DIP) joints of a finger. We plot the axis of rotation for each DoF.

and the silhouette obtained from both the RGB and the depth image. Our novel RGBD registration procedure outperforms existing state-of-the-art methods, as evident in Sec. 4.3.

We now describe the hand model used in our registration framework.

3.1 Hand Model

We use the parametric hand mesh model, aMANO [9], $M(\phi, \beta, \theta) \in \mathbb{R}^{n_v \times 3}$ with fixed triangulation among the n_v vertices. Here $\phi \in \mathbb{R}^{n_b}$ represents the user-specific local scale parameters for each of the n_b bones, calculated as

$$\phi_j = \frac{l_j^{(data)}}{l_j^{(template)}} \quad (1)$$

where $l_j^{(data)} \in \mathbb{R}$ is the length of the j^{th} bone in the observed data, obtained from keypoints (each bone connects two keypoints), and $l_j^{(template)} \in \mathbb{R}$ is the length of the j^{th} bone in the template mesh model. $\beta \in \mathbb{R}^{10}$ captures the user-specific shape parameters corresponding to the PCA shape blends of MANO [24]. The pose parameters are encoded by $\theta \in \mathbb{R}^{20}$ which represent the angles at each of the fixed axes of rotation for each degree of freedom (DoF) as shown in Fig. 2.

We now describe the procedure to deform the aMANO mesh. First, a template mesh vertex \bar{v}_i is offset as

$$v_i = \bar{v}_i + S_i \beta + P_i (r(\theta) - r(\bar{\theta})) \quad (2)$$

where $S_i \in \mathbb{R}^{3 \times 10}$ and $P_i \in \mathbb{R}^{3 \times 135}$ are the shape and pose blend shapes corresponding to the vertex v_i , $\beta \in \mathbb{R}^{10}$ is the shape parameter, $\theta \in \mathbb{R}^{15 \times 3}$ is the pose parameter capturing the axis angle rotation at each of the 15 joints, and $r(\theta) \in \mathbb{R}^{135}$ is the vectorized version of the stacked rotation matrices at each joint with pose θ ; $\bar{\theta}$ is the rest pose.

We then use the modified linear blend skinning [9] to pose the mesh. Let $a_j \in \mathbb{R}^3$ and $b_j \in \mathbb{R}^3$ be the start and end positions of

j^{th} bone in rest pose mesh (after applying MANO shape-blends) respectively, and $R_j \in \mathbb{R}^{3 \times 3}$ be the rotation matrix that takes bone j 's rest vector ($b_j - a_j$) to its pose vector ($b'_j - a'_j$). Now, the deformed vertex v'_i is given by

$$v'_i = \sum_{j=1}^{nb} W_{b_{ij}} \left\{ a'_j + R_j \left(W_{e_{ij}} s_j + (-a_j + v_i) \right) \right\} \quad (3)$$

where $s_j = (\phi_j - 1)(b_j - a_j)$, and $W_b \in \mathbb{R}^{n_v \times n_b}$ and $W_e \in \mathbb{R}^{n_v \times n_b}$ are the bone weight and endpoint weight matrices.

We also use a sparse regression matrix $K \in \mathbb{R}^{21 \times n_o}$ to calculate the 3D keypoints, $k \in \mathbb{R}^{21 \times 3}$, from vertices, $v \in \mathbb{R}^{n_v \times 3}$ as

$$k = Kv \quad (4)$$

3.2 RGBD Registration

We now detail our registration procedure that utilizes information from both RGB and depth images. Our RGBD registration procedure involves minimizing a registration energy to optimize the model parameters $\{\beta^{(f)}, \theta^{(f)}\}$ for each frame f . In the subsequent discussion, we drop the superscript f denoting the frame index for brevity. We use a sum of weighted terms in our registration energy as

$$E(\beta, \theta) = \sum_{\tau \in \mathcal{T}} \omega_\tau E_\tau(\beta, \theta) \quad (5)$$

where the terms E_τ are:

E_{data3D}	the model explains the depth point cloud
E_{data2D}	the model lies inside the observed sensor silhouette
E_{bound}	angle at each joint should respect the kinematic bounds
E_{pca}	hand pose lies in a low-dimensional manifold
E_{int}	fingers cannot inter-penetrate
E_{reinit}	model's fingertips are close to detected fingertips
E_{shape}	avoid drifting from human hand shape
E_{temp}	avoid jittery tracking

We refer the reader to Kalshetti and Chaudhuri [10] for details about these terms where the data terms only use the depth image. We adapt the E_{data2D} and E_{reinit} terms to leverage information from the RGB image too. In the absence of a depth image, we use an RGB-based hand pose estimation method [28, 37] to detect fingertips from the RGB image in the E_{reinit} term for achieving robust tracking.

We now focus our attention to the E_{data2D} term

$$E_{data2D}(\theta) = \sum_{i=1}^{n_{2D}} \|q_i - p_i(\theta)\|^2 \quad (6)$$

where $p_i \in \mathbb{R}^2$ is the 2D image space position of the evaluated predefined barycenter, and $q_i \in \mathbb{R}^2$ is the 2D image space position of the closest point on the silhouette, computed using the distance transform [17] of the cropped depth image. This term is similar to the E_{m2d} in Tkach et al. [34] and E_{bg} in Taylor et al. [32]. This term penalizes the projection of the model outside the observed silhouette. However, the silhouette obtained from the depth image is inaccurate because of missing regions, as shown in Fig. 3. Thus, the resulting distance transform is incorrect, leading to wrong correspondences, q_i .

349 Instead, we use the information from the RGB image to extract
 350 a more reliable silhouette described below.
 351

352 3.3 Silhouette Extraction

353 We present a new method to extract a reliable silhouette from RGB
 354 and the corresponding depth image. We assume the availability
 355 of a hand pose estimation network from either RGB [2, 28, 37] or
 356 depth [8, 38, 41] images, that provides detected keypoints, $k^{(data)}$.
 357

358 We use GrabCut [25] to segment the RGB image. Instead of man-
 359 ual user interaction, we automatically initialize the foreground and
 360 background Gaussian Mixture Models (GMMs) with the help of the
 361 depth image and the detected keypoints to obtain a more accurate
 362 segmentation of the hand region. Specifically, we preprocess the
 363 depth image to crop a bounding box around the hand using the
 364 detected keypoints. We then define the region outside the box as
 365 background, $\alpha = 0$, and the pixels around the detected keypoints
 366 as foreground, $\alpha = 1$.

367 We show the effect of our initialization on the resulting segmen-
 368 tation in Fig. 3. The correspondences obtained using the distance
 369 transform of our silhouette are more accurate than those obtained
 370 from the silhouette of the noisy depth image, as shown in Fig. 1.

372 3.4 Optimization

373 We use Levenberg-Marquardt to optimize the registration energy in
 374 Eq. 5 by linearizing each term written as a sum of squared residuals.
 375 To initialize the pose for the current frame, we use the previous
 376 frame’s optimized pose; for the first frame, we register the model to
 377 the keypoints obtained either from the dataset or marked manually.
 378

379 4 EVALUATION

380 4.1 Datasets

381 We evaluate our RGBD registration method on datasets that pro-
 382 vide RGB and corresponding depth frames. There are two publicly
 383 available datasets: Rendered Hand pose Dataset (RHD) [47] and
 384 Stereo hand pose Tracking Benchmark (STB) [44]. The RHD dataset
 385 is a synthetic dataset with 41258 images for training and 2728 im-
 386 ages for evaluation with a 320×320 pixels resolution. Each sample
 387 contains 21 3D keypoint annotations and 33 segmentation masks
 388 corresponding to palm, finger segments, human, and background.
 389 We combine the segmentation masks on the right hand and use
 390 them in our registration.
 391

392 The STB dataset contains 18000 stereo pairs and depth images
 393 captured from a Point Grey Bumblebee2 stereo camera and Intel
 394 Real Sense F200 active depth camera, respectively. The dataset
 395 contains 21 3D keypoint annotations corresponding to each image.
 396

397 We also capture a dataset comprising adult and child users using
 398 the Intel RealSense D435 stereo depth camera. This dataset con-
 399 tains 1000 RGB and corresponding depth frames. The user wears
 400 a wristband to crop out the wrist region from the hand, similar to
 401 Tagliasacchi et al. [30].

402 4.2 Metrics

404 We quantitatively evaluate the registration accuracy using a dense
 405 metric and a sparse metric.
 406

Method	E_{3D} (in mm)
HandTailor [15]	30.22
BiHand [43]	31.92
Our	21.12

407 **Table 1: Comparison of our registration framework with**
408 state-of-the-art methods on the RHD dataset.
 409
 410
 411
 412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431
 432
 433
 434
 435
 436
 437

Method	E_{3D} (in mm)
BiHand [43]	20
Our	14

411 **Table 2: Comparison of our registration framework with the**
412 best performing method on the STB dataset.
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431
 432
 433
 434
 435
 436
 437

Method	AUC_{5-20}
HandTailor [15]	0.632
Ours	0.775

438 **Table 3: AUC from 5mm to 20mm on the RHD dataset.**
 439
 440
 441
 442
 443
 444
 445
 446
 447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463
 464

450 The dense metric is the data-to-model error, E_{3D} . Given an ob-
 451 served depth image $I_d^{(data)}$ and a rendered depth image of the
 452 model $I_d^{(model)}$, E_{3D} is defined as
 453

$$454 E_{3D} = \frac{1}{|I_d^{(data)}|} \sum_{p \in I_d^{(data)}} \|p - \Pi_{I_d^{(model)}}(p)\| \quad (7)$$

455 where $\Pi_{I_d^{(model)}}$ denotes the closest point correspondence of the
 456 observed data point p to the rendered model point cloud.
 457

458 The sparse metric is the keypoint error, E_k . Given a set of ob-
 459 served keypoints $k^{(data)}$ and model keypoints $k^{(model)}$, the key-
 460 point error is given by
 461

$$462 E_k = \frac{1}{|k|} \sum_{i \in \{1 \dots |k|\}} \|k_i^{(data)} - k_i^{(model)}\| \quad (8)$$

463 where $k_i \in \mathbb{R}^3$ is the i^{th} keypoint. We report keypoint error wher-
 464 ever $k^{(data)}$ is available in the annotated dataset.
 465

466 4.3 Comparison with state-of-the-art

467 We compare our registration framework with state-of-the-art RGB
 468 3D hand shape and pose recovery methods using the dense data-
 469 to-model metric in Table 1 and Table 2. Our method clearly outper-
 470 forms state-of-the-art methods on each dataset.
 471

472 In hand pose estimation a challenging metric is the sparse key-
 473 point error metric. It is plotted as the fraction of test samples that
 474 have all predicted keypoints below a given maximum Euclidean
 475 distance from the ground truth. We compare the pose estimation
 476 accuracy of our method by using the area under curve (AUC) for a
 477 threshold range of 5mm to 20mm for the keypoints in Table 3. A
 478 higher area under the curve denotes more accurate results.
 479

480 The curves showing the percentage of correct keypoints (PCK)
 481 can be seen in Fig. 4.
 482

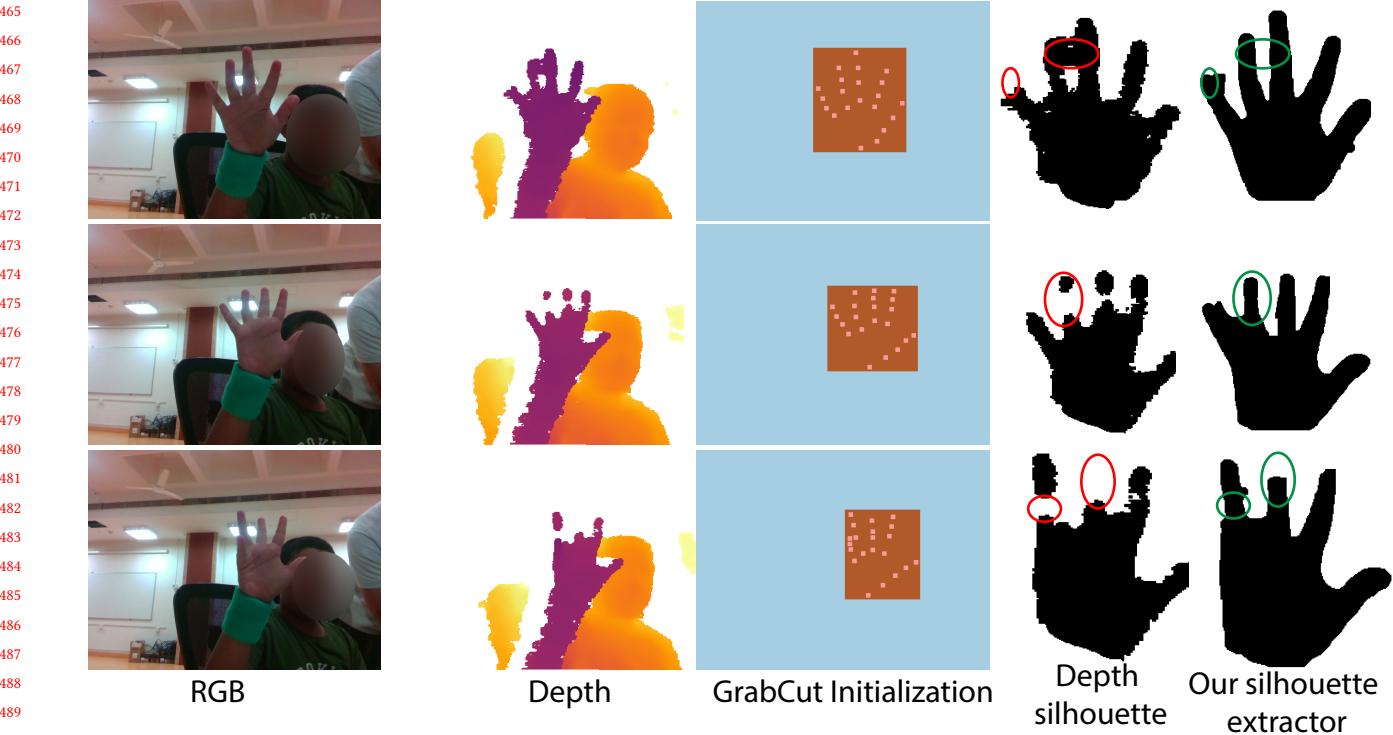


Figure 3: We use RGB and depth image along with the detected keypoints to initialize the GrabCut [25] segmentation algorithm. We define the region outside the bounding box around the keypoints as background (light blue) and the region inside the bounding box as probable foreground (brown). Further, we use the dilated keypoint locations to initialize the foreground (orange). Our silhouette extractor preserves more detail compared to the depth silhouette used in existing methods.

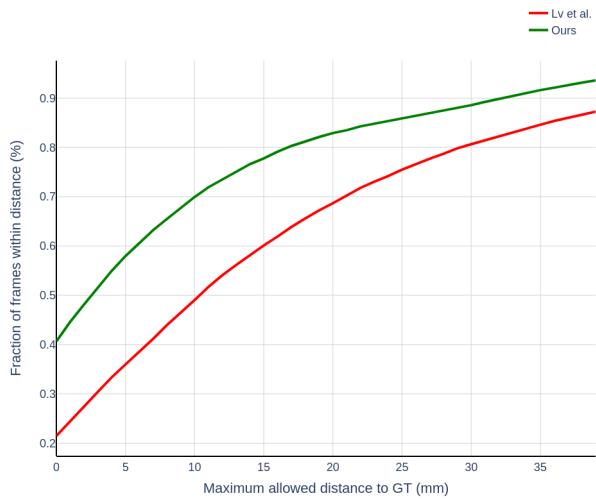


Figure 4: We calculate the percentage of correct frames within a threshold of the keypoint error on the RHD dataset. Our method outperforms state-of-the-art HandTailor method of Lv et al. [15] by a substantial margin.

We qualitatively compare the registered meshes on RHD and STB datasets in Fig. 5 and Fig. 6, respectively.

5 CONCLUSION

We present a new registration method to utilize information from RGB and depth images to achieve more accurate hand tracking than existing methods. To this end, we also propose a novel silhouette extractor that provides a reliable silhouette for the background data term in the optimization. Based on the available input data, RGB or depth or RGBD, our method is flexible enough to adapt to it and achieve state-of-the-art tracking accuracy. We lay the foundation for using the available RGBD sensors on virtual reality headsets for achieving accurate and reliable hand tracking.

Future work. Our method provides accurately registered meshes on RGB images and can serve as a dataset for training dense pose estimation networks for hands similar to dense human body pose estimation networks [7]. Further, our method can be extended to use texture and shading information for registration.

REFERENCES

- [1] Anil Armagan, Guillermo Garcia-Hernando, Seungryul Baek, Shreyas Hampali, Mahdi Rad, Zhaohui Zhang, Shipeng Xie, MingXiu Chen, Boshen Zhang, Fu Xiong, et al. 2020. Measuring Generalisation to Unseen Viewpoints, Articulations, Shapes and Objects for 3D Hand Pose Estimation under Hand-Object Interaction. In *ECCV*.

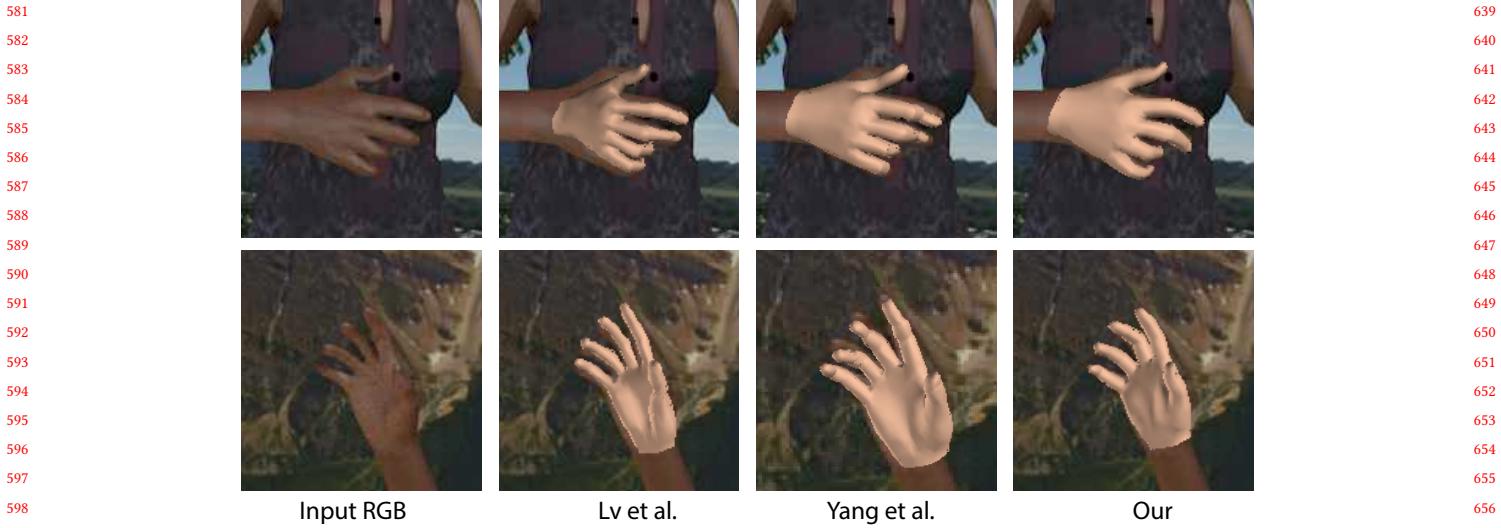


Figure 5: We qualitatively compare our method with Lv et al. [15] and Yang et al. [43] on the RHD dataset. Our method estimates shape and pose more accurately.

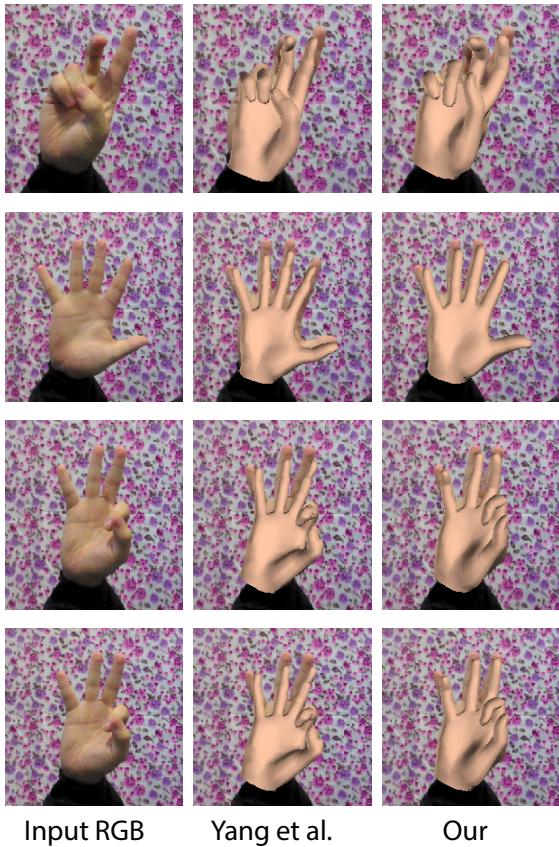


Figure 6: We qualitatively compare our method with Yang et al. [43] on the STB dataset. Our method aligns the finger articulations more accurately.

- [2] Seungryul Baek, Kwang In Kim, and Tae-Kyun Kim. 2019. Pushing the Envelope for RGB-Based Dense 3D Hand Pose Estimation via Neural Rendering. In *CVPR*.
- [3] Adnane Boukhayma, Rodrigo de Bem, and Philip HS Torr. 2019. 3d hand shape and pose from images in the wild. In *CVPR*. 10843–10852.
- [4] Yujin Chen, Zhigang Tu, Di Kang, Linchao Bao, Ying Zhang, Xuefei Zhe, Ruizhi Chen, and Junsong Yuan. 2021. Model-based 3D Hand Reconstruction via Self-Supervised Learning. In *CVPR*.
- [5] Martin de La Gorce, David J. Fleet, and Nikos Paragios. 2011. Model-Based 3D Hand Pose Estimation from Monocular Video. *IEEE TPAMI* 33, 9 (2011).
- [6] Liuhao Ge, Zhou Ren, Yuncheng Li, Zehao Xue, Yingying Wang, Jianfei Cai, and Junsong Yuan. 2019. 3D Hand Shape and Pose Estimation from a Single RGB Image. In *CVPR*.
- [7] Riza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. 2018. Densepose: Dense human pose estimation in the wild. In *CVPR*.
- [8] Weiting Huang, Pengfei Ren, Jingyu Wang, Qi Qi, and Haifeng Sun. 2020. AWR: Adaptive Weighting Regression for 3D Hand Pose Estimation. In *AAAI*.
- [9] Pratik Kalshetty and Parag Chaudhuri. 2022. Local Scale Adaptation for Augmenting Hand Shape Models. In *ACM SIGGRAPH 2022 Posters*.
- [10] Pratik Kalshetty and Parag Chaudhuri. 2022. Local Scale Adaptation to Hand Shape Model for Accurate and Robust Hand Tracking. *Comput. Graph. Forum* 41, 8 (2022). To appear.
- [11] Sameh Khamis, Jonathan Taylor, Jamie Shotton, Cem Keskin, Shahram Izadi, and Andrew Fitzgibbon. 2015. Learning an Efficient Model of Hand Shape Variation from Depth Images. In *CVPR*.
- [12] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. 2019. Learning to Reconstruct 3D Human Pose and Shape via Model-fitting in the Loop. In *ICCV*.
- [13] Dominik Kulon, Riza Alp Güler, Iasonas Kokkinos, Michael M. Bronstein, and Stefanos Zafeiriou. 2020. Weakly-Supervised Mesh-Convolutional Hand Reconstruction in the Wild. In *CVPR*.
- [14] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. 2015. SMPL: A Skinned Multi-Person Linear Model. *ACM TOG* 34, 6 (2015), 248:1–248:16.
- [15] Jun Lv, Wenqiang Xu, Lixin Yang, Sucheng Qian, Chongzhao Mao, and Cewu Lu. 2021. HandTailor: Towards High-Precision Monocular 3D Hand Recovery. In *BMVC*.
- [16] Tanwi Mallick, Partha Pratim Das, and Arun Kumar Majumdar. 2014. Characterizations of Noise in Kinect Depth Images: A Review. *IEEE Sensors Journal* 14, 6 (2014).
- [17] Calvin R Maurer, Rensheng Qi, and Vijay Raghavan. 2003. A linear time algorithm for computing exact Euclidean distance transforms of binary images in arbitrary dimensions. *IEEE TPAMI* 25, 2 (2003), 265–270.
- [18] Stan Melax, Leonid Keselman, and Sterling Orsten. 2013. Dynamics Based 3D Skeletal Hand Tracking. In *Proceedings of the ACM SIGGRAPH Symposium on Interactive 3D Graphics and Games*.

- 697 [19] Franziska Mueller, Florian Bernard, Oleksandr Sotnychenko, Dushyant Mehta,
698 Srinath Sridhar, Dan Casas, and Christian Theobalt. 2018. Ganerated hands for
699 real-time 3d hand tracking from monocular rgb. In *CVPR*.
700 [20] Iason Oikonomidis, Nikolaos Kyriazis, and Antonis A Argyros. 2011. Efficient
701 model-based 3D tracking of hand articulations using Kinect. In *BMVC*, Vol. 1:2.
702 3.
703 [21] Paschalis Panteleris, Iason Oikonomidis, and Antonis Argyros. 2018. Using a
704 single rgb frame for real time 3d hand pose estimation in the wild.
705 [22] Chen Qian, Xiao Sun, Yichen Wei, Xiaou Tang, and Jian Sun. 2014. Realtime
706 and Robust Hand Tracking from Depth. In *CVPR*. 1106–1113.
707 [23] James M. Rehg and Takeo Kanade. 1994. Visual tracking of high DOF articulated
708 structures: An application to human hand tracking. In *ECCV*.
709 [24] Javier Romero, Dimitrios Tzionas, and Michael J. Black. 2017. Embodied Hands:
710 Modeling and Capturing Hands and Bodies Together. *ACM TOG* 36, 6 (2017),
711 245:1–245:17.
712 [25] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. 2004. " GrabCut"
713 interactive foreground extraction using iterated graph cuts. *ACM TOG* 23, 3
714 (2004), 309–314.
715 [26] Toby Sharp, Cem Keskin, Duncan Robertson, Jonathan Taylor, Jamie Shotton,
716 David Kim, Christoph Rhemann, Ido Leichter, Alon Vinnikov, Yichen Wei, Daniel
717 Freedman, Eyal Krupka, Andrew Fitzgibbon, Shahram Izadi, and Pushmeet Kohli.
718 2015. Accurate, Robust, and Flexible Real-time Hand Tracking. In *CHI*. 3633–
719 3642.
720 [27] Jingjing Shen, Thomas J. Cashman, Qi Ye, Tim Hutton, Toby Sharp, Federica
721 Bogo, Andrew Fitzgibbon, and Jamie Shotton. 2020. The Phong Surface: Efficient
722 3D Model Fitting Using Lifted Optimization. In *ECCV*. 687–703.
723 [28] Tomas Simon, Hanbyul Joo, Iain Matthews, and Yaser Sheikh. 2017. Hand
724 Keypoint Detection in Single Images using Multiview Bootstrapping. In *CVPR*.
725 [29] Xiao Sun, Yichen Wei, Shuang Liang, Xiaou Tang, and Jian Sun. 2015. Cascaded
726 hand pose regression. In *CVPR*.
727 [30] Andrea Tagliasacchi, Matthias Schröder, Anastasia Tkach, Sofien Bouaziz, Mario
728 Botsch, and Mark Pauly. 2015. Robust articulated-icp for real-time hand tracking.
729 In *Comput. Graph. Forum*, Vol. 34:5. 101–114.
730 [31] David Joseph Tan, Tom Cashman, Jonathan Taylor, Andrew Fitzgibbon, Daniel
731 Tarlow, Sameh Khamis, Shahram Izadi, and Jamie Shotton. 2016. Fits Like a
732 Glove: Rapid and Reliable Hand Shape Personalization. In *CVPR*.
733 [32] Jonathan Taylor, Lucas Bordeaux, Thomas Cashman, Bob Corish, Cem Keskin,
734 Eduardo Soto, David Sweeney, Julien Valentin, Benjamin Luff, Arran Topalian, Er-
735 roll Wood, Sameh Khamis, Pushmeet Kohli, Toby Sharp, Shahram Izadi, Richard
736 Banks, Andrew Fitzgibbon, and Jamie Shotton. 2016. Efficient and Precise In-
737 teractive Hand Tracking through Joint, Continuous Optimization of Pose and
738 Correspondences. *ACM TOG* 35 (2016).
739 [33] Jonathan Taylor, Richard Stebbing, Varun Ramakrishna, Cem Keskin, Jamie
740 Shotton, Shahram Izadi, , and Andrew Fitzgibbon. 2014. User-Specific Hand
741 Modeling from Monocular Depth Sequences. In *CVPR*.
742 [34] Anastasia Tkach, Mark Pauly, and Andrea Tagliasacchi. 2016. Sphere-meshes
743 for real-time hand modeling and tracking. *ACM TOG* 35, 6 (2016), 1–11.
744 [35] Anastasia Tkach, Andrea Tagliasacchi, Edoardo Remelli, Mark Pauly, and Andrew
745 Fitzgibbon. 2017. Online generative model personalization for hand tracking.
746 *ACM TOG* 36, 6 (2017), 1–11.
747 [36] Jonathan Tompson, Murphy Stein, Yann Lecun, and Ken Perlin. 2014. Real-Time
748 Continuous Pose Recovery of Human Hands Using Convolutional Networks.
749 *ACM TOG* 33 (2014).
750 [37] Andrey Vakunov, Chuo-Ling Chang, Fan Zhang, George Sung, Matthias Grund-
751 mann, and Valentin Bazarevsky. 2020. MediaPipe Hands: On-device Real-time
752 Hand Tracking. In *Workshop on Computer Vision for AR/VR*.
753 [38] C. Wan, T. Probst, L. Gool, and A. Yao. 2018. Dense 3D Regression for Hand Pose
754 Estimation. In *CVPR*.
755 [39] Jiayi Wang, Franziska Mueller, Florian Bernard, Suzanne Sorli, Oleksandr Sotny-
756 chenko, Neng Qian, Miguel A. Otaduy, Dan Casas, and Christian Theobalt. 2020.
757 RGB2Hands: Real-Time Tracking of 3D Hand Interactions from Monocular RGB
758 Video. *ACM TOG* 39, 6 (2020).
759 [40] Robert Y. Wang and Jovan Popović. 2009. Real-Time Hand-Tracking with a Color
760 Glove. *ACM TOG* 28, 3 (2009).
761 [41] Fu Xiong, Boshen Zhang, Yang Xiao, Zhiguo Cao, Taidong Yu, Joey Zhou Tianyi,
762 and Junsong Yuan. 2019. A2J: Anchor-to-Joint Regression Network for 3D
763 Articulated Pose Estimation from a Single Depth Image. In *ICCV*.
764 [42] Hongyi Xu, Eduard Gabriel Bazavan, Andrei Zanfir, William T Freeman, Rahul
765 Sukthankar, and Cristian Sminchisescu. 2020. GHUM & GHUML: Generative 3D
766 Human Shape and Articulated Pose Models. In *CVPR*. 6184–6193.
767 [43] Lixin Yang, Jiasen Li, Wenqiang Xu, Yiqun Diao, and Cewu Lu. 2020. BiHand:
768 Recovering Hand Mesh with Multi-stage Bisected Hourglass Networks. In *BMVC*.
769 [44] Jiawei Zhang, Jianbo Jiao, Mingliang Chen, Liangqiong Qu, Xiaobin Xu, and
770 Qingxiang Yang. 2017. A hand pose tracking benchmark from stereo matching.
771 In *ICIP*.
772 [45] Xiong Zhang, Qiang Li, Hong Mo, Wenbo Zhang, and Wen Zheng. 2019. End-to-
773 end Hand Mesh Recovery from a Monocular RGB Image. In *ICCV*.
774