SPECIAL ISSUE PAPER

# Rapid avatar capture and simulation using commodity depth sensors

Ari Shapiro[1]*, Andrew Feng[1], Ruizhe Wang[2], Hao Li[2], Mark Bolas[1], Gerard Medioni[2] and Evan Suma[1]

[1] Institute for Creative Technologies, University of Southern California, Playa Vista, CA 90094, USA
[2] University of Southern California, Los Angeles, CA 90089, USA

## ABSTRACT

We demonstrate a method of acquiring a 3D model of a human using commodity scanning hardware and then controlling that 3D figure in a simulated environment in only a few minutes. The model acquisition requires four static poses taken at 90° angles relative to each other. The 3D model is then given a skeleton and smooth binding information necessary for control and simulation. The 3D models that are captured are suitable for use in applications where recognition and distinction among characters by shape, form, or clothing is important, such as small group or crowd simulations or other socially oriented applications. Because of the speed at which a human figure can be captured and the low hardware requirements, this method can be used to capture, track, and model human figures as their appearances change over time. Copyright © 2014 John Wiley & Sons, Ltd.

**\*Correspondence**

Ari Shapiro, Institute for Creative Technologies, University of Southern California, Playa Vista, CA 90094, USA.
E-mail: shapiro@ict.usc.edu

## 1. INTRODUCTION

Recent advances in low-cost scanning have enabled the capture and modeling of real-world objects into a virtual environment in 3D. For example, a table, a room, or work of art can be quickly scanned, modeled, and displayed within a virtual world with a handheld, consumer scanner. There is great value to the ability to quickly and inexpensively capture real-world objects and create their 3D counterparts. While numerous generic 3D models are available for low-cost or no-cost for use in 3D environments and virtual worlds, it is unlikely that such acquired 3D model matches the real object to a reasonable extent without individually modeling the object. In addition, the ability to capture specific objects that vary from the generic counterparts is valuable for recognition, interaction, and comprehension within a virtual world. For example, a real table could have a noticeable scratch, design, imperfection, or size that differs greatly from a stock 3D model of a table. These individual markers can serve as landmarks for people interacting with the virtual scene.

The impact of recognizing living objects in a virtual environment can be very powerful, such as the effect of seeing a relative, partner, or even yourself in a simulation.

However, living objects present simulation challenges due to their dynamic nature. Organic creatures, such as plants, can be difficult to scan because of their size and shape, which require high levels of details and stable scanning environments. Similarly, other living objects, such as people or animals, can be scanned but require much more complex models to model motion and behavior. In addition, the particular state of the living object can vary tremendously; an animal may grow, a plant can blossom flowers, and a person can wear different clothes, inhale or exhale, and gain or lose weight. Thus, capturing a moment in time of a living object is usually not sufficient for its representation in dynamic environments, where the 3D representation of that living object is expected to breath, move, grow, and respond to interaction in non-trivial ways.

In this work, we demonstrate a process for capturing human subjects and generating digital characters from those models using commodity scanning hardware. Our process is capable of capturing a human subject using still four poses, constructing a 3D model, then registering it, and controlling it within an animation system within minutes. The digital representation that our process is able to construct is suitable for use in simulations, games, and other applications that use virtual characters. Our

technique is able to model many dynamic aspects of human behavior (Figure 1). As shown in Figure 2, our main contribution in this work is a near-fully automated, rapid, low-cost end-to-end system for capturing, modeling, and simulation of a human figure in a virtual environment that requires no expert intervention.

## 2. RELATED WORK

### 2.1. 3D Shape Reconstruction

A 3D shape reconstruction has been extensively explored, among which the 3D shape reconstruction of human subjects is of specific interest to computer vision and computer graphics, with its potential applications in recognition, animation, and apparel design. With the availability of low-cost 3D cameras (e.g., Kinect and Primesense), many inexpensive solutions for 3D human shape acquisition have been proposed. The work by Tong *et al*. [1] employs three Kinect devices and a turntable. As the turntable rotates, multiple shots are taken with the three precalibrated Kinect sensors to cover the entire body. All frames are registered in a pairwise non-rigid manner using the Embedded Deformation Model [2], and loop-closure is explicitly addressed at the final stage. The work carried out in [3] utilizes two Kinect sensors in front of the self-turning subject. The subject stops at several key poses, and the captured frame is used to update the online model. Again, the dynamic nature of the turning subject is considered under the same non-rigid registration framework [2], and the loop is implicitly closed.

More recently, solutions that utilize only a single 3D sensor have been proposed, and this allows for home-based scanning and applications. The work in [4] asks the subject to turn in front of a fixed 3D sensor, and four key poses are uniformly sampled to perform shape reconstruction. The four key poses are registered in a top-bottom-top fashion, assuming an articulated tree structure of human body. Their reconstructed model, however, suffers from a low-resolution issue at a distance. To overcome the resolution issue, KinectAvatar [5] considers color constraints among consecutive frames for super-resolution. They register all super-resolution frames under a probabilistic framework. More recently, the work in [6] asks the

subject to come closer and obtain a super-resolution scan at each of eight key poses. The eight key poses are then aligned in a multi-view non-rigid manner to generate the final model. Inspired by their work, we follow the same idea of asking the subject to get closer but employ a different super-resolution scheme. Unlike the work in [6] where they merge all range scans using the Iterative Closest Point algorithm [7] along with the Poisson Surface Reconstruction algorithm [8], we use the KinectFusion algorithm [9], which incrementally updates an online volumetric model.

All these works capture the static geometry of human subjects, and additional efforts are necessary to convert the static geometry into an animated virtual character. The research works [10,11] focus on capturing the dynamic shapes of an actor's full body performance. The capturing sessions usually require a dedicated setup with multiple cameras and are more expensive than capturing only the static geometry. The resulting dynamic geometries can be played back to produce the animations of the scanned actor. The work in [12] combines dynamic shapes from multiple actors to form a shape space. The novel body deformations are driven by motion capture markers and can be synthesized based on an actor's new performance.
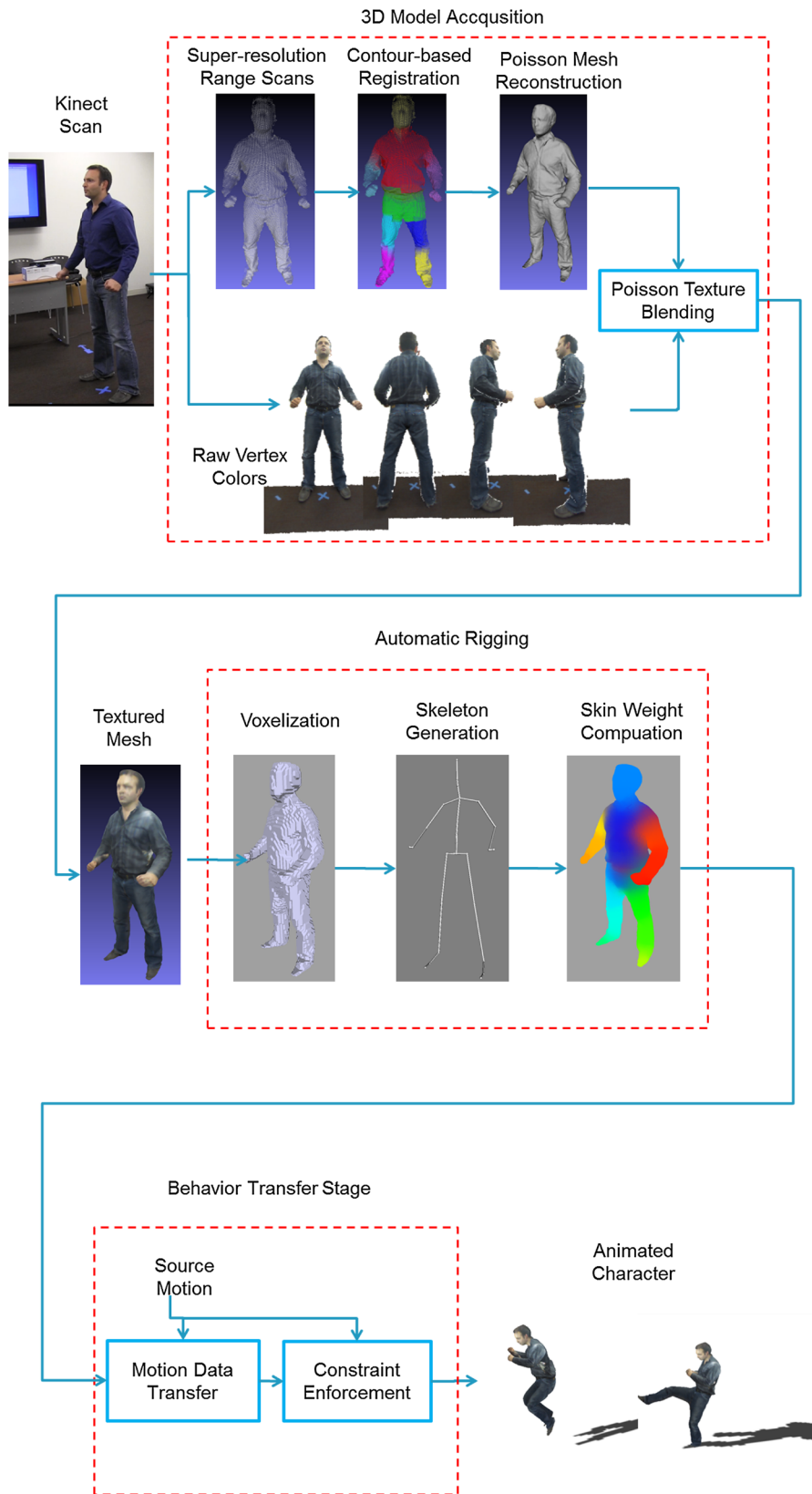
Other research has created a database of people that show the diversity of shape, size, and posture in a small population of shape, size, and posture [12]. The data set has been employed for human body modeling by fitting the model to input range scans of subject of interest [13]. This data set has also been used to manipulate a scanned human model by modifying the model proportions according to the data [14].

### 2.2. Automatic Rigging and Retargeting

While it is relatively easy to obtain static 3D character models, either from the Internet or through 3D scanning, it requires much more efforts to create an animated virtual character. A 3D model needs to be rigged with a skeleton hierarchy and appropriate skinning weights. Traditionally, this process needs to be performed manually and is time consuming even for an experienced animator. An automatic skinning method is proposed in [15] to reduce the manual efforts of rigging a 3D model. The method produces reasonable results but requires a connected and



**Figure 1.** The 3D models captured in our system can be readily applied in real-time simulation to perform various behaviors such as jumping and running with the help of auto-rigging and animation retargeting.

**Figure 2.** The overall work flow of our fast avatar capture system.

watertight mesh to work. The method proposed by Bharaj *et al.* [16] complements the previous work by automatically skinning a multi-component mesh. It works by detecting the boundaries between disconnected components to find potential joints. Thus, the method is suitable for rigging the mechanical characters that usually consist of many components. Other rigging algorithms can include manual annotation to identify important structures such as wrists, knees, and neck [17].

Recent work has shown the capability of capturing a human figure and placing that character into a simulation using 48 cameras with processing time on the order of 2 h [18]. Our method differs in that we use a single commodity camera and scanner and our processing time takes a few minutes. While this introduces a trade-off in visual quality, the minimal technical infrastructure required makes our approach substantially more accessible to a widespread audience. In addition, our method requires no expert intervention during the rigging and animation phases.

# 3. 3D MODEL RECONSTRUCTION

We propose a convenient and fast way to acquire accurate static 3D human models of different shapes by the use of a single commodity hardware, for example, Kinect. The subject turns in front of the Kinect sensor in a natural motion while staying static at four key poses, namely, front, back, and two profiles, for approximately 10 s each. For each key pose, a super-resolution range scan is generated as the Kinect device, controlled by a built-in motor, moves up and down (Section 3.1). The four super-resolution range scans are then aligned in a multi-view piecewise rigid manner, assuming small articulations between them. Traditional registration algorithms (e.g., Iterative Closest Point [7]), which are based on the *shape coherence*, fail in this scenario because the overlap between consecutive frames is very small. Instead, we employ *contour coherence* (Section 3.2) and develop a contour-based registration method [19], which iteratively minimizes the distance between the closest points on the predicted and observed contours (Section 3.3). For more details on using *contour coherence* for multi-view registration of range scans, please refer to [19]. In this paper, we summarize their method and give a brief introduction. At the final stage, the four aligned key poses are processed to generate a watertight mesh model using the Poisson Surface Reconstruction algorithm [8]. The corresponding texture information of the four super-resolution range scans is inferred using the Poisson Texture Blending algorithm [20] (Section 3.4).

## 3.1. Super-resolution Range Scan

Given the field of view of the Kinect sensor, the subject must stand 2 m away in order to cover the full body while turning in front of the device. The data are heavily
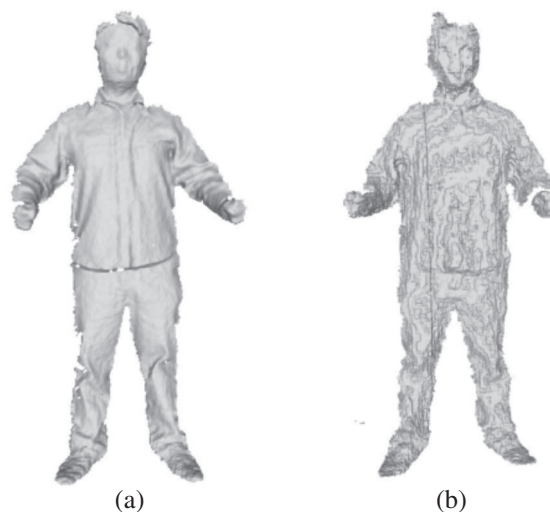
quantized at that distance (Figure 3(b)), thus produces a poor quality scan, which results in a coarse model after integration. Here, instead, we ask the subject to come closer and stay as rigid as possible at the four key poses, while the Kinect device scans up and down to generate a super-resolution range scan. Each pose takes 10 s, and approximately 200 frames are merged using the KinectFusion algorithm [9] (Figure 3(a)). This process greatly improves the quality of the input and allows us to capture more details, such as wrinkles of clothes and face as shown in Figure 3. It is worth mentioning that the ground is removed by using the RANSAC algorithm [21], assuming that the subject of interest is the only thing in the sensor's predefined capture range.
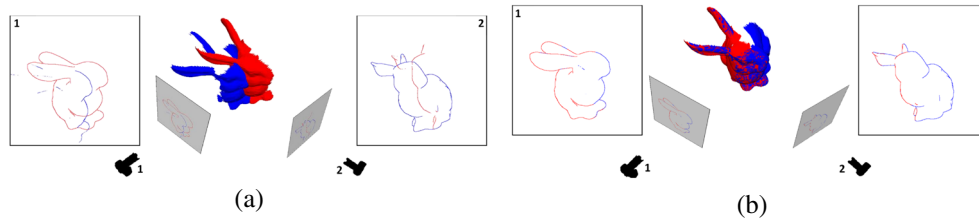
## 3.2. Contour Coherence as a Clue

The amount of overlap between two consecutive super-resolution range scans is limited as they are 90° apart (i.e., wide baseline). As such, traditional *shape coherence*-based methods (e.g., Iterative Closest Point and its variants [22]) fail, as it is hard to establish the point-to-point correspondences on two surfaces with small overlap.

An example of two wide baseline range scans of the Stanford bunny with approximately 35% overlap is given in Figure 4(a). Traditional methods fail, as most closest-distance correspondences are incorrect.

While the traditional notion of *shape coherence* fail, we propose the concept of *contour coherence* for wide baseline range scan registration. *Contour coherence* is defined as the agreement between the observed apparent contour and the predicted apparent contour. As shown in Figure 4(a), the observed contours extracted from the original 2.5D range scans, that is, red lines in image 1 and blue



(a)          (b)

**Figure 3.** (a) Super-resolution range scans after integrating approximately 200 frames using the KinectFusion algorithm and (b) low-resolution single range scan at the distance of 2 m.

**Figure 4.** (a) Two roughly aligned wide baseline 2.5D range scans of the Stanford bunny with the observed and predicted apparent contours extracted. The two meshed points cloud are generated from the two 2.5D range scans, respectively. (b) Registration result after maximizing the contour coherence.

lines in image 2, do not match the corresponding predicted contours extracted from the projected 2.5D range scans, that is, blue lines in image 1 and red lines in image 2. We maximize *contour coherence* by iteratively finding closest correspondences among apparent contours and minimizing their distances. The registration result is shown in Figure 4(b) with the *contour coherence* maximized and two wide baseline range scans well aligned. The *contour coherence* is robust in the presence of wide baseline in the sense that no matter the amount of overlap between two range scans, only the shape area close to the predicted contour generator is considered when building correspondences on the contour, thus avoiding the search for correspondences over the entire shape.

### 3.3. Contour Coherence-based Registration Method

We apply the notion of *contour coherence* to solve the registration problem of four super-resolution range scans with small articulations. For simplicity, we start the discussion with the contour-based rigid registration of two range scans. As shown in Figure 4(a), the observed contour and the predicted contour do not match. In order to maximize the *contour coherence*, we iteratively find the closest pairs of points on two contours and minimize their distances. Assume that point $\mathbf{u} \in \mathbb{R}^2$ is on predicted contour in image 1 of Figure 4(a) (i.e., blue line) and point $\mathbf{v} \in \mathbb{R}^2$ is its corresponding closest point on the observed contour in image 1 (i.e., red line), we minimize their distance as

$$\left\| \mathbf{v} - \mathcal{P}_1 \left( T_1^{-1} T_2 \mathcal{V}_2 \left( \tilde{\mathbf{u}} \right) \right) \right\| \tag{1}$$

where $\tilde{u}$ is the corresponding pixel location in image 2 of $\mathbf{u}$, $\mathcal{V}_2$ maps the pixel location $\tilde{u}$ to its 3D location in the coordinate system of camera 2, $T_1$ and $T_2$ are the camera to world transformation matrices of camera 1 and 2, respectively, and $\mathcal{P}_1$ is the projection matrix of camera 1. Assuming known $\mathcal{P}_1$ and $\mathcal{P}_2$, we iterate between finding all closest contour points on images 1 and 2 and minimizing the sum of their distances (Eq. 1) to update the camera poses $T_1$ and $T_2$ until convergence. We use quaternion to represent the rotation part of $T$ and Levenberg–Marquardt algorithm to solve for the minimization as it is nonlinear in
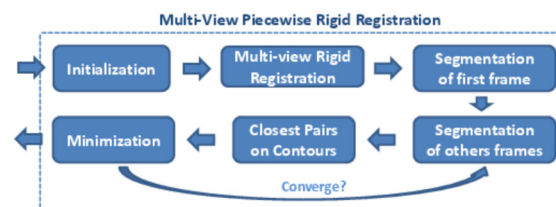
parameters. It is worth mentioning that minimizing Eq. 1 updates $T_1$ and $T_2$ at the same time, and this enables us to perform multi-view rigid registration in the case of three or more frames.

The extension from rigid registration to piecewise rigid registration is quite straightforward. Each segment (i.e., segmented body part) is considered rigid, and all the rigid segments are linked by a hierarchical tree structure in the case of body modeling. We again iteratively find the closest pairs on contours between all corresponding body segments and minimize the sum of their distances.
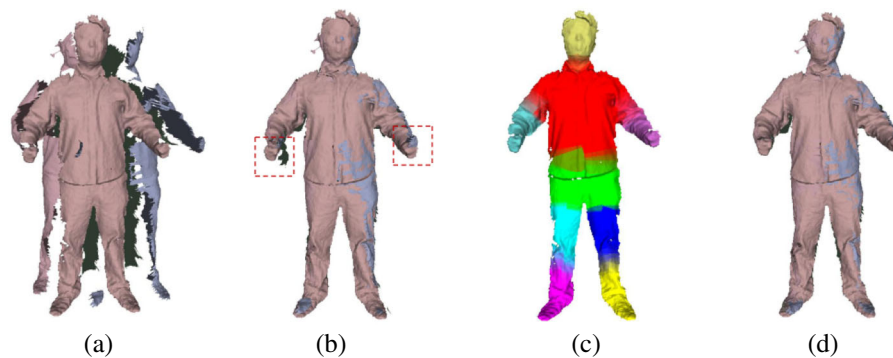
A complete pipeline of our registration method is given in Figure 5. First, the four super-resolution range scans are initialized by assuming a 90° rotation between consecutive frames (Figure 6(a)). Second, they are further aligned by the multi-view rigid registration method considering the whole body as rigid (Figure 6(b)). While the translation part of the camera pose is not well estimated by the initialization procedure, it is corrected by the multi-view rigid registration step. As indicated by the red boxes, however, the small articulations between frames still remain unresolved under the rigid assumption. Third, the front pose is roughly segmented into nine body parts in a heuristic way (Figure 6(c)). Fourth, we iteratively propagate the segmentation to other frames, find closest pairs on contours between corresponding rigid body parts, and minimize their distances to update the camera poses, as well as the human poses of each frame (Figure 6(d)).

### 3.4. Watertight Mesh Model with Texture

At this point, we have aligned all four super scans to produce a point cloud with normal vectors. Poisson mesh reconstruction [8] is used to obtain a watertight mesh from



**Figure 5.** General pipeline of our registration method.

**Figure 6.** (a) Four super-resolution range scans after initialization; (b) four super-resolution range scans after multi-view rigid registration, with red boxes indicating unresolved small articulations under the rigid assumption; (c) rough segmentation of the front pose; and (d) four super-resolution range scans after multi-view piecewise rigid registration.

the point clouds. The Kinect camera also captures the color information from the scanned person when generating the superscans at each pose. For each superscan, we also store a color image corresponding to the range scan and combine the color images to produce the texture for the watertight mesh. We follow a similar procedure as in [6] to corrode the color images and remove unreliable pixels. The corroded color images are then transferred onto the superscans as vertex colors to produce color meshes before going through the registration process. Finally, these aligned color meshes are used to texture the watertight mesh generated from Poisson reconstruction. We apply the Poisson texture blending algorithm in [20] to fill out the gaps and holes in the texture and produce the final color mesh.

## 4. RESOLUTION INDEPENDENT AUTOMATIC RIGGING

Animating a 3D character model usually requires a skeletal structure to control the movements. Our system automatically builds and adapts a skeleton to the 3D scanned character. Thus, it can later apply the rich sets of behavior on the character through motion retargeting.

The auto-rigging method in our system is similar to the one proposed in [15]. The method builds a distance field from the mesh and uses the approximate medial surface to extract the skeletal graph. The extracted skeleton is then matched and refined based on the template skeleton. The method is automatic and mostly robust, but it requires a watertight and single component mesh to work correctly. This poses a big restriction on the type of 3D models the method can be applied to. For example, the production meshes usually come with many props and thus have multiple components. On the other hand, the mesh produced from range scans tend to contain holes, non-manifold geometry, or other topological artifacts that require additional cleanup. Moreover, the resulting mesh produced through the super-resolution scans usually consists of hundreds of thousands of vertices. Such high-resolution

meshes would cause the auto-rigging method to fail during optimization process to build the skeleton. To alleviate this limit, we proposed a modified method that works both for generic production models and large meshes.

Our key idea is that the mesh could be approximated by a set of voxels and the distance field could be computed using the voxels. The voxels are naturally free from any topological artifacts and are easy to processed. It is carried out by first converting the mesh into voxels using depth buffer carving in all positive and negative $x$, $y$, and $z$ directions. This results in six depth images that can be used to generate the voxelization of the original mesh. Although most small holes in the original mesh are usually removed in the resulting voxels because of discretization, some holes could still remain after the voxelization. In removing the remaining holes, we perform the image hole filling operation in the depth images to fill up the small empty pixels. After voxelization, we select the largest connected component and use that as the voxel representation for the mesh. The resulting voxels are watertight and connected and can be converted into distance field to construct the skeleton. Figure 2 demonstrates the process of converting the original mesh into voxel representation to produce the skeleton hierarchy and skinning weights.

The voxel representation is only an approximation of the original mesh. Therefore, the resulting distance field and, consequently, the skeleton could be different from the one generated with the original mesh. In our experiments, we found that the resulting skeletons tend to be very similar as shown in Figure 7 and do not impact the overall animation quality in the retargeting stage. Once we obtain the skeleton, the skinning weights can be computed using the original mesh instead of the voxels because the weight computation in [15] does not rely on the distance field. Alternatively, the skinning weights can be computed using the techniques in [23], which use voxels to approximate the geodesic distance for computing bone influence weights. Thus, we can naturally apply their algorithm using our resulting voxels and skeleton to produce higher-quality smooth bindings.

**Figure 7.** The voxelization produces the skeleton similar to the one extracted from original mesh. Left: original mesh and its skeleton. Right: voxel representation of original mesh and its corresponding skeleton.

## 5. BEHAVIOR TRANSFER

The behavior transfer stage works by retargeting an example motion set from our canonical skeleton to the custom skeleton generated from automatic rigging. Here, we use the method from [24] to perform motion retargeting. The retargeting process can be separated into two stages. The first stage is to convert the joint angles encoded in a motion from our canonical skeleton to the custom skeleton. This is carried out by first recursively rotating each bone segment in target skeleton to match the global direction of that segment in source skeleton at default pose so that the target skeleton is adjusted to have the same default pose as the source skeleton. Once the default pose is matched, we address the discrepancy between their local frames by adding suitable pre-rotation and post-rotation at each joint in target skeleton. These pre-rotation and post-rotation are then used to convert the joint angles from source canonical skeleton to the target skeleton.

The second stage is using inverse kinematics to enforce various positional constraints such as foot positions to remove motion artifacts such as foot sliding. The inverse kinematic method we use is based on damped Jacobian Pseudo-Inverse algorithm [25]. We apply this inverse kinematic method at each motion frame in the locomotion sequences to ensure that the foot joint is in the same position during the foot plant stage. After the retargeting stage, the acquired 3D skinned character can be incorporated into the animation simulation system to execute a wide range of common human-like behaviors such as walking and gesturing.

## 6. APPLICATIONS

### 6.1. 3D Capture for Use in Games and Simulation

We demonstrate our method by showing the capture and processing, registration, and subsequent simulation of a human figure in our accompanying video and in Figure 8

in the following text. The construction of a 3D model takes approximately 4 min, and the automatic rigging, skinning, and registration of a deformable skeleton take approximately 90 s. Models typically contain between 200K and 400K vertices and 400K and 800K faces. Simulation and control of the character are performed in real time using various animations and procedurally based controllers for gazing and head movement. The 3D models captured in this way are suitable for use in games where characters need to be recognizable from a distance but do not require face-to-face or close interactions.

### 6.2. Temporal Avatar Capture

Because our method enables the capture of a 3D character without expert assistance and uses commodity hardware, it is economically feasible to perform 3D captures of the same subject over a protracted period of time. For example, a 3D model could be taken every day of the same subject,



**Figure 8.** A representative captured character from scan containing 306K vertices and 613K faces. Note that the distinguishing characteristics are preserved in the capture and simulation, such as hair color, clothing style, height, and skin tone.

**Figure 9.** Models generated from captures over a period of 4 days. Note changes and commonality in clothing, hair styles, and other elements of appearance.

which would reflect their differences in appearance over time. Such captures would reflect changes in appearance such as hair style or hair color, clothing, or accessories worn. In addition, such temporal captures could reflect personal changes such as growth of facial hair, scars, and weight changes. Such temporal information could be analyzed to determine clothing preferences or variations in appearance (Figure 9).

Note that our method will generate a skeleton for each 3D model. Thus, avatars of the same subject will share the same topology but have different bone lengths.

### 6.3. Crowds

Many applications that use virtual crowds require tens, hundreds, or thousands of characters to populate the virtual space. Research has experimented with saliency to show the needed variation in traditionally modeled characters to model a crowd [26] and the number of variations needed [27]. By reducing the cost of constructions of 3D characters, crowd members can be generated from a population of captured subjects rather than through traditional 3D means.

## 7. DISCUSSION

We have demonstrated a technique that allows the capture and simulation of a human figure into a real-time simulation without expert intervention in a matter of few minutes.

### 7.1. Limitations

The characters generated are suitable for applications where recognizability and distinction among the virtual characters are important. In the course of our experiments, we have found the virtual characters to be recognizable to those familiar with the subjects. The characters are not suitable for close viewing or in simulations where face details are needed, such as conversational agent or talking head applications. Higher levels of detail are needed for areas such as the face and hands before other models of synthetic motion, such as emotional expression, lip syncing, or gesturing could be used. Additionally, our method makes no distinction between the body of the captured subject and their clothing. Thus, bulky clothing or accessories could change the skeletal structure of the virtual character. Also, the behaviors associated with the characters are retargeted from sets of motion data and control algorithms but are not generated from movements or motion gleaned from the subject itself. Thus, motion transferred to all captured subjects shares the same characteristics, differing only by the online retargeting algorithm, which accommodates differently sized characters. This homogeneity can be partially circumvented by including variations in the set of motion data, such as differing locomotion or gesturing sets for male and female characters. For future work, we plan on extracting movement models from the captured subjects in order to further personalize their virtual representation.

## REFERENCES

1. Tong J, Zhou J, Liu L, Pan Z, Yan H. Scanning 3D full human bodies using kinects. *IEEE Transactions on Visualization and Computer Graphics* 2012; **18**(4): 643–650.
2. Sumner RW, Schmid J, Pauly M. Embedded deformation for shape manipulation. In *ACM*

*Transactions on Graphics (TOG)*, Vol. 26. ACM: New York, NY, USA, 2007; 80.

3. Zeng M, Zheng J, Cheng X, Liu X. Templateless quasi-rigid shape modeling with implicit loop-closure. In *2013 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE: Portland, Oregon, USA, 2013; 145–152.

4. Wang R, Choi J, Medioni G. Accurate full body scanning from a single fixed 3d camera. In *2012 Second International Conference on 3D Imaging, Modeling, Processing, Visualization and Transmission (3DIMPVT)*. IEEE: Zurich, Switzerland, 2012; 432–439.

5. Cui Y, Chang W, Nöll T, Stricker D. KinectAvatar: fully automatic body capture using a single kinect. In *Computer VISION-ACCV 2012 Workshops*. Springer: Daejeon, Korea, 2013; 133–147.

6. Li H, Vouga E, Gudym A, Luo L, Barron JT, Gusev G. 3D self-portraits. *ACM Transactions on Graphics (Proceedings SIGGRAPH Asia 2013)* 2013; **32**(6): 187.

7. Chen Y, Medioni G. Object modelling by registration of multiple range images. *Image and Vision Computing* 1992; **10**(3): 145–155.

8. Kazhdan M, Bolitho M, Hoppe H. Poisson surface reconstruction. In *Proceedings of the Fourth Eurographics Symposium on Geometry Processing*, Sardinia, Italy, 2006.

9. Newcombe RA, Davison AJ, Izadi S, Kohli P, Hilliges O, Shotton J, Molyneaux D, Hodges S, Kim D, Fitzgibbon A. KinectFusion: real-time dense surface mapping and tracking. In *2011 10th IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE: Basel, Switzerland, 2011; 127–136.

10. Vlasic D, Peers P, Baran I, Debevec P, Popovic J, Rusinkiewicz S, Matusik W. Dynamic shape capture using multi-view photometric stereo. In *In ACM Transactions on Graphics*, Yokohama, Japan, 2009.

11. Wu C, Stoll C, Valgaerts L, Theobalt C. On-set performance capture of multiple actors with a stereo camera. In *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia 2013),* volume 32, Hong Kong, November 2013.

12. Anguelov D, Srinivasan P, Koller D, Thrun S, Rodgers J, Davis J. Scape: shape completion and animation of people. In *ACM SIGGRAPH 2005 Papers*. SIGGRAPH '05. ACM: New York, NY, USA, 2005; 408–416.

13. Weiss A, Hirshberg D, Black MJ. Home 3D body scans from noisy image and range data. In *2011 IEEE International Conference on Computer Vision (ICCV)*. IEEE: Barcelona, Spain, 2011; 1951–1958.

14. Jain A, Thormählen T, Seidel HP, Theobalt C. MovieReshape: tracking and reshaping of humans in videos. *ACM Transactions on Graphics (Proceedings SIGGRAPH Asia 2010)* 2010; **29**(5): 148.

15. Baran I, Popović J. Automatic rigging and animation of 3D characters. *ACM Transactions on Graphics* 2007; **26**(3): 72.

16. Bharaj G, Thormählen T, Seidel HP, Theobalt C. Automatically rigging multi-component characters. *Computer Graphics Forum (Proceedings Eurographics 2012)* 2011; **30**(2): 755–764.

17. Mixamo auto-rigger, 2013. Availvable from: http://www.mixamo.com/c/auto-rigger [10 March 2014].

18. xxarray demo at ces, 2014. Availvable from: http://gizmodo.com/nikon-just-put-me-in-a-video-game-and-it-was-totally-in-1497441443 [10 March 2014].

19. Wang R, Choi J, Medioni G. 3D modeling from wide baseline range scans using contour coherence. In *2014 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE: Columbus, Ohio, USA, 2014.

20. Chuang M, Luo L, Brown BJ, Rusinkiewicz S, Kazhdan M. Estimating the Laplace–Beltrami operator by restricting 3D functions. In *Computer Graphics Forum*, Vol. 28. Wiley Online Library: New York, NY, USA, 2009; 1475–1484.

21. Fischler MA, Bolles RC. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM* 1981; **24**(6): 381–395.

22. Rusinkiewicz S, Levoy M. Efficient variants of the ICP algorithm. In *2001 Proceedings of the Third International Conference on 3-D Digital Imaging and Modeling*. IEEE: Quebec City, Quebec, Canada, 2001; 145–152.

23. Dionne O, de Lasa M. Geodesic voxel binding for production character meshes. In *Proceedings of the 12th ACM SIGGRAPH/Eurographics Symposium on Computer Animation*. SCA '13. ACM: New York, NY, USA, 2013; 173–180.

24. Feng A, Huang Y, Xu Y, Shapiro A. Fast, automatic character animation pipelines. *Computer Animation and Virtual Worlds* 2013; **25**(1): 3–16.

25. Buss SR. Introduction to inverse kinematics with jacobian transpose, pseudoinverse and damped least squares methods. Technical report, *IEEE Journal of Robotics and Automation* 2004.

26. McDonnell R, Larkin M, Hernández B, Rudomin I, O'Sullivan C. Eye-catching crowds: saliency based selective variation. *ACM Transactions on Graphics* 2009; **28**(3): 55:1–55:10.

27. McDonnell R, Larkin M, Dobbyn S, Collins S, O'Sullivan C. Clone attack! perception of crowd variety. *ACM Transactions on Graphics* 2008; **27**(3): 26:1–26:8.

# AUTHORS' BIOGRAPHIES

**Ari Shapiro** is a research scientist at the USC Institute for Creative Technologies where he leads the Character Animation and Simulation research group. His research focuses on synthesizing controllable models of movement and behavior for virtual characters. For several years, he worked on character animation tools and algorithms in the research and development and software departments of visual effects and video games companies such as Industrial Light and Magic, LucasArts and Rhythm & Hues Studios.

Shapiro has published many academic articles in the field of computer graphics and animation for virtual characters. He completed his PhD in Computer Science at UCLA in 2007 in the field of computer graphics with a dissertation on character animation using motion capture, physics and machine learning. He also holds an MS in Computer Science from UCLA and a BA in Computer Science from the University of California, Santa Cruz.

**Andrew Feng** is currently a research associate at the Institute for Creative Technologies. He received the PhD and MS degrees in Computer Science from the University of Illinois at Urbana-Champaign. His research interests include character animation, mesh deformation, mesh skinning and real-time rendering.

**Ruizhe Wang** received a BS degree in Electrical Engineering from Tsinghua University in June 2010 and an MS degree in Electrical Engineering from Caltech in December 2011. He is currently pursuing a PhD degree in Computer Science and doing research in the Computer Vision Lab at the University of Southern California. His research interest includes computer vision, computer graphics and machine learning. He is a member of the IEEE.

**Hao Li** is an assistant professor of Computer Science at the University of Southern California since 2013. Prior to joining USC, he spent at year at Industrial Light & Magic/Lucasfilm as a research lead, developing next generation real-time performance capture technologies. He spent a year as a postdoctoral fellow at Columbia and Princeton Universities in 2011. His research lies in geometry processing, 3D reconstruction, and realtime performance capture. His algorithms are widely deployed in the industry ranging from leading visual effects studios to manufacturers of radiation therapy systems. Hao has extensive publications at ACM SIGGRAPH and was named top 35 innovator under 35 by MIT Technology Review in 2013 and NextGen 10: Innovators under 40 by CSQ magazine in 2014. He was also awarded the Swiss National Science Foundation fellowship for prospective researchers in 2011 and obtained the best paper award at the Symposium of Computer Animation in 2009. Hao obtained his PhD at ETH Zurich and his MSc degree at the University of Karlsruhe (TH).

**Mark Bolas** is an Associate Director at USC's Institute of Creative Technologies where he is the director of the Mixed-Reality Lab; he is also an Associate Professor in the Interactive Media Division of the USC School of Cinematic Arts. His work focuses on creating virtual environment transducers 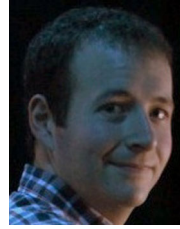and experiences that fully engage one's perception and cognition in order to form a visceral memory of the experience. Mark's work at USC has included the development of a range of new imaging technologies, including a 150° head-mounted display, a nonobtrusive retroreflective display for mixed reality applications, and contributions to USC's 360° light field display. Mark holds more than 20 patents, is the recipient of numerous product awards and the 2005 IEEE Virtual Reality Technical Achievement Award. Bolas' thesis work at Stanford's School of Engineering, 'Design and Virtual Environments', was among the first efforts to map the breadth of virtual reality as a new medium.

**Gerard Medioni** received the Diplôme d'Ingénieur from the École Nationale Supérieure des Télécommunications (ENST), Paris, in 1977 and the MS and PhD degrees from the University of Southern California (USC) in 1980 and 1983, respectively. He has been with USC since then and is currently a Professor of Computer Science and Electrical Engineering, a Co-Director of the Institute for Robotics and Intelligent Systems (IRIS), and a Co-Director of the USC Games Institute. He served as the chairman of the Computer Science Department from 2001 to 2007. He has made significant contributions to the field of computer vision. His research covers a broad

spectrum of the field, such as edge detection, stereo and motion analysis, shape inference and description, and system integration. He has published three books, more than 50 journal papers, and 150 conference articles. He is a holder of eight international patents. He is an associate editor of the Image and Vision Computing Journal, Pattern Recognition and Image Analysis Journal, and the International Journal of Image and Video Processing. He served as a Program Co-Chair of the 1991 IEEE Computer Vision and Pattern Recognition (CVPR) Conference and the 1995 IEEE International Symposium on Computer Vision, a General Co-Chair of the 1997 IEEE CVPR Conference, a Conference Co-Chair of the 1998 International Conference on Pattern Recognition, a General Co-Chair of the 2001 IEEE CVPR Conference, a General Co-Chair of the 2007 IEEE CVPR Conference, and a General Co-Chair of the 2009 IEEE CVPR Conference. He is a fellow of the IEEE, IAPR, and AAAI and a member of the IEEE Computer Society.

**Evan Suma** is a Research Assistant Professor at the Institute for Creative Technologies and the Computer Science Department at the University of Southern California. As one of the co-directors of the MxR Lab, his work focuses on techniques and technologies that enhance immersive virtual experiences and address important challenges in the domains of training, education, and rehabilitation. Dr Suma has written or co-authored over 50 published papers in the areas of virtual environments, 3D user interfaces, and human-computer interaction. Additionally, he also created the Flexible Action and Skeleton Toolkit (FAAST), a gestural interaction software framework that has been widely adopted by the research and hobbyist communities. He received his PhD in Computer Science from the University of North Carolina at Charlotte in 2010.