

Avatars for the Masses: Smartphone-Based Reconstruction of Humans for Virtual Reality

Timo Menzel^{1,†}, Erik Wolf^{2,†}, Stephan Wenninger¹, Niklas Spinczyk¹, Lena Holderrieth², Carolin Wienrich³, Ulrich Schwanecke⁴, Marc Erich Latoschik², Mario Botsch^{1,*}

¹Computer Graphics Group, TU Dortmund University, Dortmund, Germany

²Human-Computer Interaction Group, Julius-Maximilians-Universität Würzburg, Würzburg, Germany

³Psychology of Intelligent Interactive Systems Group, Julius-Maximilians-Universität Würzburg, Würzburg, Germany

⁴Computer Vision and Mixed Reality Group, RheinMain University of Applied Sciences, Wiesbaden, Germany

Correspondence*:

Mario Botsch

mario.botsch@tu-dortmund.de

† These authors contributed equally to this work

2 ABSTRACT

3 Realistic full-body avatars play a key role in representing users in virtual environments, where
4 they have been shown to considerably improve important effects of immersive experiences such
5 as body ownership and presence. Consequently, the demand for realistic virtual humans – and
6 methods for creating them – is rapidly growing. However, despite extensive research into 3D
7 reconstruction of avatars from real humans, an *easy and affordable* method for generating *realistic*
8 *and VR-capable* avatars is still lacking: Existing methods are either limited to complex capture
9 hardware and/or controlled lab environments, do not provide sufficient visual fidelity, or cannot be
10 rendered at sufficient frame rates for multi-avatar VR applications. To make avatar reconstruction
11 widely available, we developed *Avatars for the Masses* – a client-server-based online service
12 for scanning real humans with an easy-to-use smartphone application that empowers even non-
13 expert users to capture photorealistic and VR-ready avatars. The data captured by the smartphone
14 are transferred to a reconstruction server, where the avatar is generated in a fully automated
15 process. Our advancements in capturing and reconstructing allow for higher-quality avatars even
16 in less controlled in-the-wild environments. Extensive qualitative and quantitative evaluations show
17 our method's avatars to be on par with the ones generated by expensive expert-operated systems.
18 It also generates more accurate replicas in comparison to the current state-of-art in smartphone-
19 based reconstruction, produces much less artifacts and provides a much higher rendering
20 performance in VR in comparison to three representative neural methods. A comprehensive
21 user study confirms similar perception results compared to avatars reconstructed with expensive
22 expert-operated systems, and it underscores a sufficient usability of the overall system. To truly
23 bring avatars to the masses, we will make our smartphone application publicly available for
24 research purposes.

25 **Keywords:** Avatar generation, 3D scanning, virtual human, embodiment, virtual body ownership

1 INTRODUCTION

26 Avatars are digital representations of users that can be dynamically rendered in virtual environments in real
27 time to reflect the behavior of their users (Bailenson and Blascovich, 2004). While avatars can be of almost
28 any conceivable shape and appearance, in this research, we specifically refer to humanoid representations
29 that vary from stylized to realistically reconstructed 3D models. Such avatars may appear generic, lacking
30 distinctive or individual features, or they can be personalized to closely resemble the appearance of their
31 respective user. With the recent surge in virtual reality (VR) research (Skarbez and Jiang, 2024) and
32 the increasing availability of mature head-mounted displays (HMDs) (Sutherland, 1968), avatars have
33 become increasingly important as faithful self-representations of users in almost countless scenarios. These
34 scenarios include metaverse-like social VR environments (Latoschik et al., 2019; Yoon et al., 2019; Aseeri
35 and Interrante, 2021; Mystakidis, 2022) or VR applications to support mental health (Sampaio et al., 2021;
36 Döllinger et al., 2022). Among them are critical applications for which maintaining user identity and
37 conveying realistic emotions are crucial for authentic interactions and a sophisticated user experience
38 (UX). Prior work has shown that realistically personalized full-body avatars, which can look deceptively
39 similar to the user, are superior for the outlined scenarios by increasing the user's sense of presence and
40 embodiment or self-identification with the avatar (Waltemate et al., 2018; Salagean et al., 2023; Fiedler
41 et al., 2024; Kim et al., 2023), or to increase emotional response (Gall et al., 2021; Waltemate et al., 2018).

42 Unfortunately, many approaches for scanning-based full-body avatar generation rely on complex and
43 expensive multi-camera rigs for photogrammetric reconstruction, such as (Achenbach et al., 2017; Shetty
44 et al., 2024; Ma et al., 2021). Methods for generating avatars from monocular video input make avatar
45 generation more affordable, but early approaches (Alldieck et al., 2018a,b) suffered from insufficient
46 quality, as shown in (Wenninger et al., 2020). Recent avatar reconstructions adapt NeRFs (Mildenhall
47 et al., 2020) or Gaussian Splatting (Kerbl et al., 2023) as underlying representations, for instance (Jiang
48 et al., 2023; Moreau et al., 2024). Although this is an exciting and very promising research direction, our
49 experiments in Section 4 clearly demonstrate that these approaches are not (yet) capable of providing
50 sufficient visual quality and rendering performance for VR applications. So far, the method of Wenninger
51 et al. (2020), which reconstructs mesh-based avatars from smartphone videos, seems to be the most suitable
52 for the affordable reconstruction of photorealistic and VR-capable full-body avatars. However, while the
53 low hardware requirements make avatar reconstruction more affordable, the scanning process requires
54 sufficient experience, the reconstruction process involves commercial products, and the system's operation
55 requires expert knowledge. Consequently, there is still no approach for fast, affordable, and easy-to-operate
56 reconstruction of photorealistic and VR-capable full-body avatars. This prevents the full potential of
57 photorealistic avatars from being realized for many applications.

58 To bridge this gap and make avatar reconstruction both affordable and widely available to non-expert
59 users, we present *Avatars for the Masses*, an easy-to-use system for smartphone-based person scanning and
60 server-based avatar reconstruction. In particular, our contributions are:

- 61 • An easy-to-use smartphone application that visually guides the user through the scanning process,
62 enabling even non-expert users to achieve high-quality results;
- 63 • A server-based pipeline that fully automatically reconstructs a photorealistic avatar from smartphone-
64 captured data in about 20 minutes, without relying on commercial components;
- 65 • Technical improvements in the capture and reconstruction processes that result in high quality results
66 even in uncontrolled outdoor environments;

- 67 • Qualitative and quantitative technical evaluations and comparisons with several state-of-the-art
68 approaches that clearly demonstrate the advantages of our system;
- 69 • A user-centric evaluation through a user study that evaluates and confirms both our smartphone app's
70 usability and the resulting avatars' quality (*captured by non-expert first-time users!*).

71 Our evaluations demonstrate that the proposed system is indeed fast, affordable, and easy to use, and that it
72 achieves avatar quality almost on par with that of complex camera rigs – even in challenging “in-the-wild”
73 capture scenarios. As such, and due to the lack of commercial components, it has the potential to bring
74 avatars to the masses. We will make our system publicly available for research to encourage this.

2 RELATED WORK

75 In this section, we describe the mechanisms and implications of representing oneself through an avatar in
76 virtual reality (Section 2.1), before discussing different approaches to generate realistic avatars (Section 2.2).
77 In the following, we restrict our discussion to avatars *personalized* (as opposed to generic), *realistic* (as
78 opposed to stylized), and *full-body* (as opposed to head-only or upper-body-only), because these are the
79 most challenging with regard to the outlined desiderata.

80 2.1 Avatars for Self-Representation in Virtual Reality

81 The egocentric embodiment of avatars for self-representation in VR (Slater et al., 2010) can positively
82 impact the UX of virtual environments (Mottelson et al., 2023). This includes improving the key
83 psychometric properties of VR, such as the sense of presence (Waltemate et al., 2018; Wolf et al., 2021;
84 Skarbez et al., 2017), or intensifying emotional responses to virtual content (Waltemate et al., 2018; Gall
85 et al., 2021). Other advantages may include improved spatial perception (Mohler et al., 2010; Leyrer et al.,
86 2011), reduced cognitive load (Steed et al., 2016), or higher performance and accuracy (Jung and Hughes,
87 2016; Pastel et al., 2020) when performing tasks in VR.

88 A crucial aspect in evaluating the effectiveness of avatar embodiment is the sense of embodiment (SoE),
89 consisting of the feeling of owning (ownership), controlling (agency), and being located within (self-
90 location) a virtual body in a virtual environment (Kilteni et al., 2012; de Vignemont, 2011). Previous work
91 has shown that realistic and personalized avatars increase the SoE towards the avatar (Waltemate et al.,
92 2018; Fiedler et al., 2023; Salagean et al., 2023) and thus contribute to an overall plausible VR experience
93 (Latoschik and Wienrich, 2022).

94 Photorealistic and personalized avatars are particularly valuable for maintaining the user's identity, which
95 is beneficial in social VR experiences (Yoon et al., 2019; Aseeri and Interrante, 2021; Mystakidis, 2022)
96 or applications supporting mental health (Sampaio et al., 2021; Döllinger et al., 2022; Turbyne et al.,
97 2021). Previous work has also shown that self-related cues through avatar embodiment and personalization
98 significantly increase self-identification with the avatar (Fiedler et al., 2024), potentially maintaining a
99 more accurate self-perception in VR, even in body-swap paradigms (Döllinger et al., 2024). However, a
100 realistic personalization of avatars can also harm UX, as their human-like realism combined with their high
101 affinity to the user can potentially trigger Uncanny Valley effects, leading to negative emotional responses
102 such as eeriness towards the avatars (Mori et al., 2012; Döllinger et al., 2023).

103 Overall, comprehensive evidence exists for notable effects of photorealistic personalized avatars on
104 important user states. Consequently, we will employ representative psychometric measures for a prominent
105 selection of the aforementioned effects of avatars to evaluate the 3D reconstruction quality achieved with

106 our developed system. Therefore, Section 5 reports on a user study evaluating our avatars with respect to
107 the sense of embodiment, plausibility, and a potential uncanny valley effect. In addition, this user study
108 also evaluates the general usability and user satisfaction of the smartphone front-end to ensure appropriate
109 ease of use and user satisfaction.

110 **2.2 Generation of Realistic Personalized Avatars**

111 The growing demand for virtual avatars has triggered a lot of research in scanning-based avatar
112 reconstruction in the recent years. We restrict ourselves to realistic full-body avatars and discuss related
113 approaches with respect to our target application requirements: The avatar generation should be affordable
114 and easy to use, the resulting avatars should accurately resemble the scanned person, and the avatars should
115 be suitable for VR applications – meaning they can be rendered from arbitrary camera views and at a
116 sufficiently high frame rate for multi-avatar (social) VR applications.

117 Many approaches employ complex and expensive rigs of 50–100 cameras to capture high-quality photos
118 or videos of the person to be scanned (Feng et al., 2017; Achenbach et al., 2017; Ma et al., 2021; Kwon
119 et al., 2023; Salagean et al., 2023; Morgenstern et al., 2024; Shetty et al., 2024; Pang et al., 2024). While
120 these methods achieve highly accurate reconstructions, they are restricted to dedicated capture laboratories
121 whose operation requires expert knowledge.

122 Instead of simultaneously taking images with multiple cameras, other approaches use a single monocular
123 camera (or smartphone) to capture a sequence of images or videos. While this design choice considerably
124 reduces hardware cost and complexity, it increases the capture time, which inevitably causes small
125 movements of the scanned person and reduces geometric accuracy. Early approaches suffer from
126 considerably lower quality compared to camera rigs (Alldieck et al., 2018a,b, 2019), most visible in
127 the face region. Wenninger et al. (2020) address this problem by incorporating close-ups of the head into
128 the avatar reconstruction, producing a quality that is objectively quite close and subjectively very similar to
129 those of camera rigs (Bartl et al., 2021). On the downside, their method is rather complicated to operate,
130 is intended for controlled indoor environments, and relies on commercial components, which prevents
131 widespread use by researchers and non-experts.

132 More recently, neural geometry representations, such as Neural Radiance Fields (NeRFs) (Mildenhall
133 et al., 2020) or 3D Gaussian Splatting (3DGS) (Kerbl et al., 2023) have been extensively adapted to avatar
134 reconstruction. Neural avatar representations (Peng et al., 2021b; Zhao et al., 2022; Jiang et al., 2024;
135 Guo et al., 2023; Xiao et al., 2024; Lin et al., 2024) are capable of reconstructing fine details, since they
136 are not restricted to a fixed mesh topology. Avatars based on NeRFs (Liu et al., 2021; Peng et al., 2021a;
137 Jiang et al., 2023, 2022; Wang et al., 2023; Yu et al., 2023; Wang et al., 2024; Zheng et al., 2022, 2023) or
138 3DGS (Hu et al., 2024; Shao et al., 2024; Moreau et al., 2024; Li et al., 2024; Habermann et al., 2023) are
139 therefore better suited for reconstructing clothing and hair. However, as our experiments with recent neural
140 avatars show (see Section 4), their generation from image and video data can be very time-consuming
141 (from hours up to days), their rendering is not fast enough for multi-avatar VR applications (where 90 fps
142 at 2k resolution for left/right eye is desired), and their visual fidelity is not sufficient (when viewed from
143 directions not covered by training data). This last point is a particularly challenging limitation, since in
144 multi-user social VR applications there is no control over viewing directions and avatar poses, which can
145 quickly lead to visual artifacts.

146 We therefore employ a traditional mesh-based representation for virtual avatars and build on the approach
147 of Wenninger et al. (2020), which we advance in several important aspects. First, our smartphone
148 application visually guides the user through the capture process, thereby ensuring high-quality input

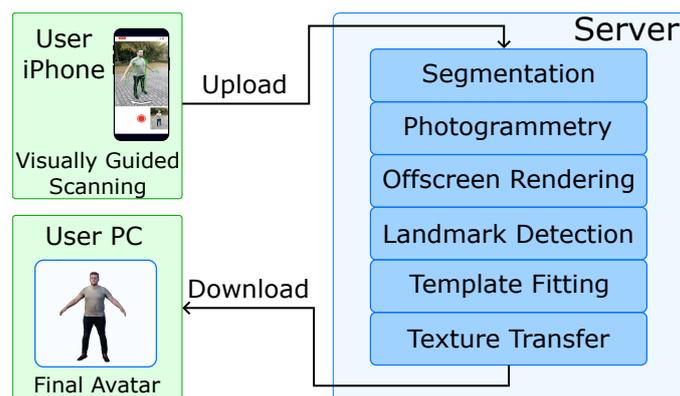


Figure 1. The user scans a subject with our smartphone application (top left). The captured images are uploaded to our processing server (right), where a fully automatic reconstruction pipeline generates an avatar in about 20 minutes. The user can then download his/her avatar into any VR application (bottom left).

149 data. Second, we technically improve the data acquisition, image pre-processing, and template fitting,
 150 leading to more accurate and more robust avatar reconstructions. Third, we replace the commercial
 151 components of Wenninger et al. (2020) with carefully selected non-commercial alternatives, allowing us to
 152 make our system publicly available. Finally, we evaluate our approach (i) by qualitative and quantitative
 153 comparisons to state-of-the-art avatar reconstruction methods, and (ii) in terms of a carefully designed user
 154 study, in which first-time users successfully reconstruct and evaluate avatars.

3 AVATAR RECONSTRUCTION

155 Our approach extends and improves the work by Wenninger et al. (2020). We start with a brief overview of
 156 their method in order to point out our specific technical improvements later on. Wenninger et al. (2020)
 157 record two videos of the to-be-scanned person: The *body video* circles around the scanning subject twice
 158 to capture both the lower and the upper body. The *head video* circles around the face/head in a close-up
 159 manner to capture facial details. From these two videos, individual frames are extracted and fed into
 160 Agisoft Metashape (Agisoft, 2023), a commercial photogrammetry reconstruction tool, resulting in two
 161 point clouds for the body and head. A template mesh is then fitted to the point clouds in a two-step process:
 162 The template is first fitted to the body point cloud (for the overall shape) and then to the head point cloud
 163 (for fine-tuning facial features). Landmarks detected by OpenPose (Cao et al., 2019) guide the template
 164 fitting process. In a final step, the avatar texture is generated from the input images.

165 Our approach, as outlined in Figure 1, introduces guided smartphone-based data capturing (Section 3.1)
 166 and a fully automatic server-based reconstruction pipeline (Section 3.2). In the following, we describe the
 167 components of both phases and point out the main contributions and technical improvements compared to
 168 Wenninger et al. (2020).

169 3.1 Smartphone-Based Data Acquisition

170 Analogous to Wenninger et al. (2020), we capture people by performing (i) a full-body scan in A-pose
 171 and (ii) a close-up head scan using a smartphone (see Figure 1, top left). However, our approach differs in
 172 the kind of data that is captured (Section 3.1.1) and how the user performs the scanning (Section 3.1.2).

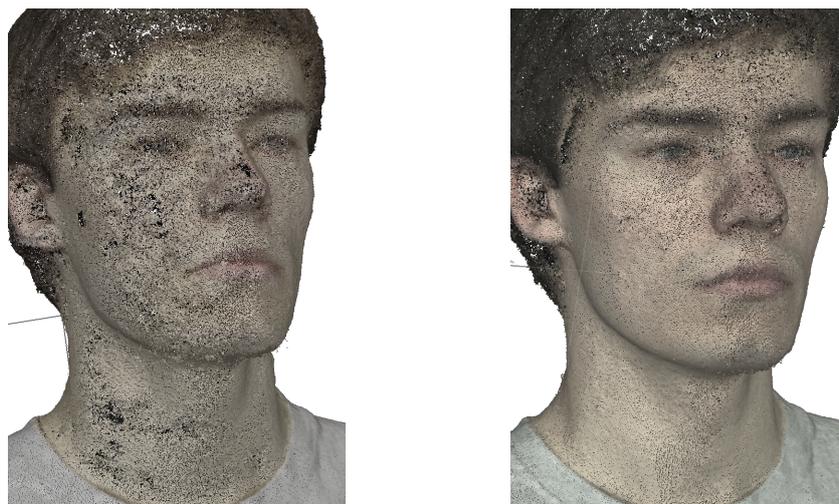


Figure 2. Point clouds reconstructed via photogrammetry from video frames suffer from compression artifacts (left). The higher quality of individually captured images yields more accurate point clouds (right).

173 3.1.1 Capturing Videos vs. Images

174 Videos captured with current smartphone cameras are compressed using H.264 or H.265. These algorithms
175 are optimized for viewing each video frame for a fraction of a second only, hence allowing for rather
176 aggressive per-frame compression. In addition, inter-frame compression exploits blockwise similarity of
177 consecutive frames, which further degrades image quality (Wiegand et al., 2003). As photogrammetry
178 algorithms use image gradients to detect feature points, the block edges can negatively influence the quality
179 of the resulting 3D point cloud (Figure 2, left). Furthermore, extracted video frames can be affected by
180 motion blur, which Wenninger et al. (2020) had to handle explicitly.

181 In contrast, individual photographs can be captured at considerably higher quality, since they suffer
182 much less from motion blur, avoid the inter-frame block compression artifacts, and allow to use less
183 aggressive compression in general. Higher-quality images in turn yield more accurate photogrammetry
184 results (Figure 2, right), which will eventually result in more accurate avatars with fine geometric details
185 and higher-quality textures. Our capture process (described next) therefore records individual photographs
186 instead of videos, at a resolution of 3024×4032 pixels.

187 In addition to the high-resolution RGB images, we also capture coarse depth images (576×768 pixels)
188 using the smartphone's depth sensor and the phone's orientation (resp. gravity vector). The former helps to
189 determine the reconstructed subject's correct size/scaling, and the latter to determine its correct orientation.
190 Our scanning application is designed for Apple iOS devices, and this additional information is conveniently
191 included in the meta-data of Apple's HEIC image format.

192 3.1.2 Scanning UI

193 Our extensive experience with the approach of Wenninger et al. (2020) revealed that the quality of the
194 photogrammetry reconstruction strongly depends on the correct distance and orientation of the camera to
195 the subject. If the camera is too close, some regions (e.g. feet, hands) might not be captured. If the camera
196 is too far away, valuable image resolution is wasted and the point cloud becomes less dense and more noisy.
197 Those errors typically occurred to unexperienced users – despite detailed previous instructions.

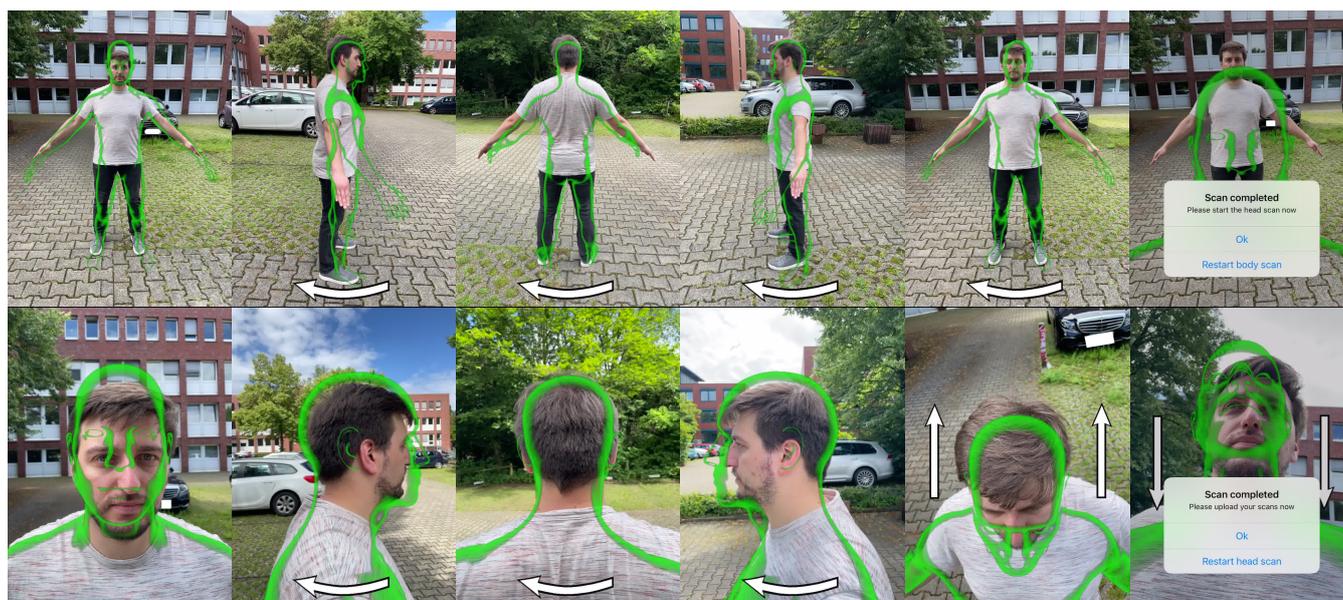


Figure 3. The guided scanning procedure of the smartphone application. Top row, from left to right: Initial overlay before starting the body scan; overlay during the body scan; end of body scan. Bottom row, from left to right: Initial overlay before starting the head scan; overlay during the head scan; end of head scan.

198 To avoid these errors, our smartphone application visually guides the user through the scanning process:
 199 On top of the camera feed, we overlay in green the silhouette of a virtual human model (of average size
 200 and shape) as seen from the intended camera position (see Figure 3). The user adjusts the phone/camera to
 201 roughly match the silhouette of the subject and the model. It is not necessary to precisely fit the overlay to
 202 the subject. The purpose of the overlay only is to guide the user to maintain proper distance and orientation.
 203 To also guide the user's movement around the scanned subject, the virtual camera moves around the green
 204 virtual model in the same way that the user should move around the scanned subject. In addition, the
 205 direction of the movement is indicated by a white arrow.

206 During the scanning process, the app captures images at a frequency of 1 Hz and a resolution of
 207 3024×4032 pixels. The speed of the virtual camera's movement is chosen to result in 105 images for the
 208 entire scanning process, as this number experimentally turned out to be the best compromise: Fewer images
 209 degraded the point cloud quality, more images did not improve the results but increased the computation
 210 time. Thanks to this well-controlled capture process, we require just one circle around the subject for
 211 the full-body scan and one for the head scan – thus reducing the scanning time to about two-thirds of
 212 Wenninger et al. (2020). Since a shorter scanning time reduces artifacts caused by subject movement, it
 213 also improves geometric accuracy.

214 A dialog informs the user when the full-body and head scans are complete (Figure 3, right column), after
 215 which the captured images are uploaded to the reconstruction server. All further user instructions or hints
 216 during the scanning procedure are displayed in Figure 3. The entire scanning process can also be seen
 217 in the accompanying video. To further minimize scanning errors, the app displays a step-by-step tutorial
 218 before the scanning process, covering subject preparation (hairstyle, accessories, shoes, and clothes), scan
 219 pose requirements (A pose), and scan process explanations (body/head scan procedure).

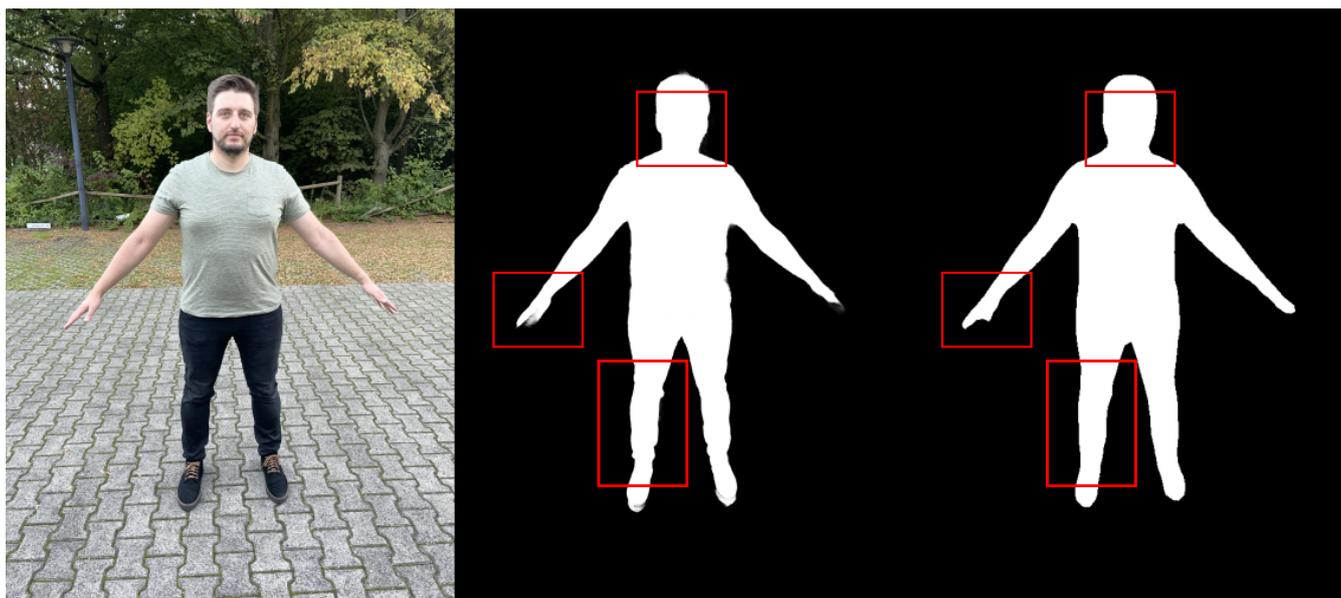


Figure 4. When segmenting foreground and background of the input images, Apple's person segmentation (center) better preserves small details, such as hair and clothing creases, compared to DeepLabV3 (right).

220 3.2 Server-Based Reconstruction Pipeline

221 The avatar reconstruction pipeline, whose individual tasks are described in this section, includes several
222 computationally expensive tasks. To speed up the avatar generation process and to reduce the load on
223 the smartphone's resources, the captured images are uploaded to a compute server, where the avatar is
224 automatically reconstructed and can be downloaded by the user.

225 3.2.1 Image Preprocessing

226 Our experiments with Wenninger et al. (2020) revealed that their method works well in controlled indoor
227 environments, but in outdoor environments it often gives noticeably worse results. This is due to non-static
228 background, such as leaves moving in the wind or cars driving by. These background movements violate
229 the photogrammetry assumption of a static scene, leading to incorrect extrinsic camera parameters and
230 consequently to errors in the reconstructed 3D point cloud.

231 To eliminate these problems and thereby make the reconstruction process much more robust with respect
232 to "in-the-wild" capture environments, we segment the input images into foreground and background
233 and mask out the background before passing the images to the photogrammetry process. To this end, we
234 compared DeepLabV3 (Chen et al., 2017) and Apple's person segmentation (Apple Inc., 2023d) (on macOS
235 15.3.1), and decided for the latter since it produced slightly more accurate and more detailed masks in our
236 experiments (see Figure 4). Moreover, as the image background is excluded from the reconstruction, the
237 number of image features to be matched by the photogrammetry is significantly reduced (accelerating this
238 process by 30%), the resulting point cloud contain considerably fewer points (accelerating later template
239 fitting), and the point clouds contain significantly less noise and outliers (improving accuracy and robustness
240 of the template fitting). Overall, this image preprocessing leads to faster computations and much cleaner
241 point clouds – in particular in uncontrolled outdoor environments (see Figure 5).

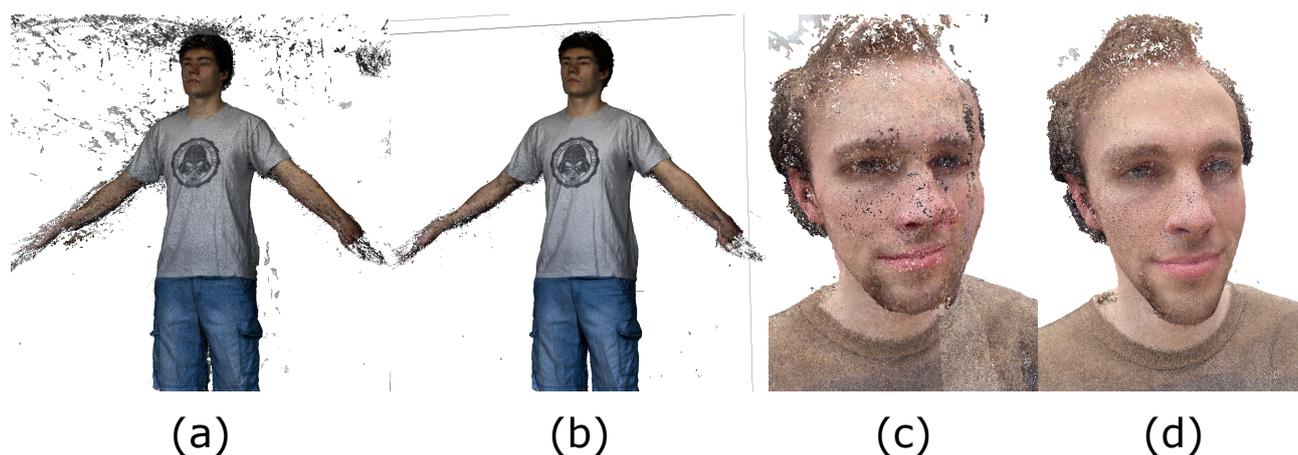


Figure 5. Point clouds reconstructed from unprocessed input images captured in outdoor environments suffer from significant noise, outliers, and misalignments in the face region (a, c). By processing the images through person segmentation and background removal, these artifacts are largely eliminated (b, d).



Figure 6. The (meshed) photogrammetry reconstructions of Agisoft Metashape (left) and Apple's RealityKit (right) are very similar in geometric fidelity, texture quality, and computational performance.

242 3.2.2 Photogrammetry

243 The captured and segmented images are passed to the photogrammetry stage, which reconstructs a dense
 244 3D point cloud (see Figure 1). Wenninger et al. (2020) employ Metashape (Agisoft, 2023) for this task, a
 245 widely used commercial photogrammetry software. Unfortunately, its license restrictions explicitly prohibit
 246 the use in server-based reconstruction scenarios (which we aim for). In order to make our system publicly
 247 available for research purposes, we compared several non-commercial alternatives for photogrammetric
 248 reconstruction, including MeshRoom (Alice Vision, 2025), COLMAP (Schönberger and Frahm, 2016), and
 249 Apple's RealityKit (Apple Inc., 2023b) (on macOS 15.3.1). From these frameworks, Apple's RealityKit
 250 consistently produced the highest-quality results, which are very similar in terms of geometric fidelity,
 251 texture quality, and processing time to those of Agisoft Metashape (see Figure 6).

252 Apple's high-level photogrammetry API allows one to specify different configuration options, where
253 we found the settings *raw* detail, *high* feature sensitivity, and *unordered* samples to yield best results. In
254 addition to the high-resolution RGB images, we also provide segmentation masks, coarse depth images,
255 gravity vectors, and EXIF data, which significantly improves the photogrammetry results compared to
256 processing only the RGB images. The depth data and gravity vectors help to determine the correct scale
257 and orientation of the reconstructed object, which make the upcoming steps more reliable.

258 3.2.3 Landmark Detection

259 The previous photogrammetry stage produces a (static) high-resolution textured triangle mesh (see
260 Figure 6). This mesh can suffer from artifacts in insufficiently scanned regions and is lacking animation
261 controls (body skeleton, facial blendshapes). The well-established approach is to fit a high-quality template
262 mesh with all required animation controls to the photogrammetry output (point cloud or mesh). This results
263 in a reconstructed avatar mesh that inherits its triangulation, UV texture layout, and animation controls
264 from the template model, while closely resembling the geometric shape and the texture/material of the
265 photogrammetry scan. This template fitting process (described in the next subsection) has to be initialized
266 and guided by a set of landmarks on both the template model (where they are pre-selected once) and the
267 photogrammetry output (where they have to be manually selected or automatically detected).

268 Wenninger et al. (2020) detect body and face landmarks in the 2D input images using OpenPose (Cao
269 et al., 2019) and (Dlib, 2022), respectively, where they select the best suited (e.g., most frontal) input
270 images based on heuristics. The detected 2D landmarks are re-projected onto the photogrammetry point
271 cloud. While working well in most cases, their approach can fail if wrong images are selected for landmark
272 detection or if a 2D landmark in sparsely sampled regions back-projects to the wrong surface part. We
273 avoid both problems by rendering the photogrammetry mesh (Figure 6) from several camera positions and
274 performing landmark detection on the resulting synthetic images. Our controlled capturing process enables
275 the straightforward selection of suitable camera views, and the back-projection onto the rendered 3D mesh
276 is well-defined for any detected 2D landmark.

277 We detect 4 hand landmarks (two knuckles on the left and right hands) and 37 face landmarks (eye
278 contours, tip of the nose, and mouth features), which are passed on to the template-fitting stage. To ensure
279 accurate and reliable landmark detection, we compared on a wide range of examples the face/hand landmark
280 detection of OpenPose (Cao et al., 2019), Dlib (2022), Apple's Vision Framework (Apple Inc., 2023c,
281 2025), and Google's MediaPipe (Lugaresi et al., 2019). Since MediaPipe produced the most reliable results
282 in our experiments, we chose this landmark detector for our reconstruction pipeline.

283 3.2.4 Template Fitting

284 The previous two stages result in two high-resolution textured photogrammetry meshes from the body
285 and head scans of the subject standing in A-pose with neutral facial expression, as well as a set of 37 face
286 and 4 body landmarks. We perform curvature-adaptive point sampling on the two photogrammetry meshes
287 to convert them into two point clouds for body and head, respectively. To reconstruct the avatar, we fit a
288 fully rigged statistical template mesh to the photogrammetry data, guided by the landmarks.

289 Our template mesh was designed by a skilled artist (to be free of license restrictions) and has a slightly
290 higher resolution (23752 vertices) than the template from the Autodesk Character Generator in Wenninger
291 et al. (2020). Its animation rig consists of a full-body skeleton with 59 joints, as well as 52 facial blends
292 that are compatible with ARKit (Apple Inc., 2023a). The template was fit to 1700 scans of the CAESAR
293 database (Robinette et al., 1999) to derive a 30-dimensional PCA subspace of human body shapes.

294 In a first step, the template is coarsely fitted to the point clouds by iteratively optimizing *alignment*
295 (position, orientation, scale), overall *body shape* (PCA weights), and *body pose* (skeleton joint angles).
296 In the second step, the initial template fit is refined by optimizing all individual *vertex positions*. Both
297 optimization phases minimize the sum of squared distances of photogrammetry points to their closest
298 points on the template mesh in a non-rigid ICP manner, guided by the landmark points. Both steps are
299 regularized to prevent overfitting: the first step by Tikhonov regularization on the PCA weights, the second
300 step by a discrete bending energy (see Achenbach et al. (2017); Wenninger et al. (2020) for details).

301 Our method differs from Wenninger et al. (2020) in two aspects: Wenninger et al. first fit the template to
302 the body point cloud and then refine the result by fitting it to the head point cloud. Since in their approach
303 the absolute scaling of these point clouds is unknown, the proportions of body to head can be slightly
304 wrong. In contrast, our coarse depth images determine the absolute scale. We also pre-align the body and
305 head point clouds using landmark-guided ICP and then fit the template to both point clouds *simultaneously*.
306 In this process, closest-point correspondences to the head/body regions of the template mesh are computed
307 from the head/body point clouds only, respectively. This approach effectively avoids the wrong body-head
308 proportions (see Figure 9). In addition, since our advanced scanning process yields more accurate point
309 clouds, we require less regularization in the fine-scale fitting step, resulting in more geometric details.

310 3.2.5 Texture Transfer

After reconstructing the geometric shape of the avatar in the previous step, the final step reconstructs the texture image. The two photogrammetry meshes already feature high-quality textures generated from the input images, but with a rather poor UV texture layout. To have a uniform texture layout for all avatars, we transfer texture_P of the photogrammetry mesh to texture_A of the avatar mesh (having the high-quality texture layout of the template). For each texel \mathbf{u}_A in the avatar's UV layout we determine the corresponding 3D point \mathbf{x}_A on the avatar mesh (based on texture coordinates), find its closest point \mathbf{x}_P on the photogrammetry mesh, and copy its color by its texture coordinate \mathbf{u}_P :

$$\text{texture}_A[\mathbf{u}_A] \leftarrow \text{texture}_P[\mathbf{u}_P]$$

311 Note that we actually fill two textures, from the body and head scan, respectively. These two texture
312 images are then combined into one using Poisson Image Editing (Pérez et al., 2003). This final step of the
313 reconstruction pipeline results in a textured avatar mesh (see Figure 7 for some examples).

4 QUANTITATIVE AND QUALITATIVE EVALUATION

314 Reconstructing an avatar with our approach starts by capturing a person using our iPhone application
315 (iPhone 12 Pro and iPhone 13 Pro Max in our experiments). The app visually guides the user through
316 the scanning procedure and takes 105 images (45 full-body and 60 head images), which takes about
317 two minutes and is shown in the accompanying video. The captured image data (about 320 MB) is then
318 uploaded to our server, which takes less than one minute over WiFi. The reconstruction is performed in a
319 fully automatic manner on the server (Mac Studio, M1-Max 10-Core CPU, 32-Core integrated GPU, 64 GB
320 RAM) and takes about 19 minutes (1 min segmentation, 7 min photogrammetry, 5 min offscreen rendering,
321 2 min landmark detection, 4.5 min template fitting and texture generation). The whole process, therefore,
322 takes about 22 minutes only, after which the avatar can then be downloaded in file formats compatible with
323 VR and game engines.



Figure 7. Avatars reconstructed with our approach, all being scanned in uncontrolled outdoor settings.

324 In the following, we compare our results to those of a complex multi-camera rig Achenbach et al. (2017),
325 to the smartphone-based method of Wenninger et al. (2020), and to three recent neural avatar techniques
326 based on NeRFs or 3D Gaussian Splatting (Müller et al., 2022; Shao et al., 2024; Lei et al., 2024).

327 **4.1 Quantitative Comparisons**

328 Following Wenninger et al. (2020), we evaluate the accuracy of our avatar reconstruction by reporting
329 reprojection errors. To this end, we render the resulting textured avatar onto the images captured during the
330 scan process using the camera calibration data from the photogrammetry process (see Figure 8), and then
331 compute the root-mean-square errors over all rendered pixels in CIELab color space, averaged over all
332 images. This metric allows us to measure errors resulting from inaccuracies in both color and geometry.

333 We perform this evaluation on the 33 subjects that were scanned during the user study described in
334 Section 5. These participants were scanned by (i) another 33 non-expert first-time users of our smartphone
335 application, as well as (ii) an expert using the multi-camera rig at the Embodiment Lab of JMU Würzburg
336 (106 Canon EOS 1300D DSLR cameras, based on Achenbach et al. (2017)). Generating the avatars using
337 the expert-operated multi-camera rig took about 15 minutes. Our pipeline, on the other hand, took around 22
338 minutes. Despite the tremendous difference in expertise of the scanning person and in cost and complexity
339 of the scanner setup, the results obtained with the camera rig are only slightly better than our smartphone
340 scans (see Figure 8). Averaging over all 33 scans and comparing our RSME ($M = 32.83$, $SD = 4.88$) with
341 that of the multi-camera rig ($M = 32.29$, $SD = 6.36$) reveals that the error increases by less than 2%, while
342 the financial cost decreases by more than 98%.

343 **4.2 Qualitative Comparisons**

344 Besides the easy-to-use visually guided scanning procedure, our method improves the approach of
345 Wenninger et al. (2020) by several technical contributions, as described in Section 3.2. To evaluate
346 the effect of these contributions, we compare with their method in Figure 9. The two subjects were
347 captured in an outdoor environment by recording videos (for their method) and images (for our method)
348 on the same iPhone 12 Pro. The avatar generation took 15 minutes (their method) and 22 minutes
349 (our method). Our method produces noticeably more accurate results, with more geometric detail and
350 higher-quality textures. This results from more accurate photogrammetry point clouds (due to recording
351 higher-resolution images instead of videos and due to background removal) and from better template

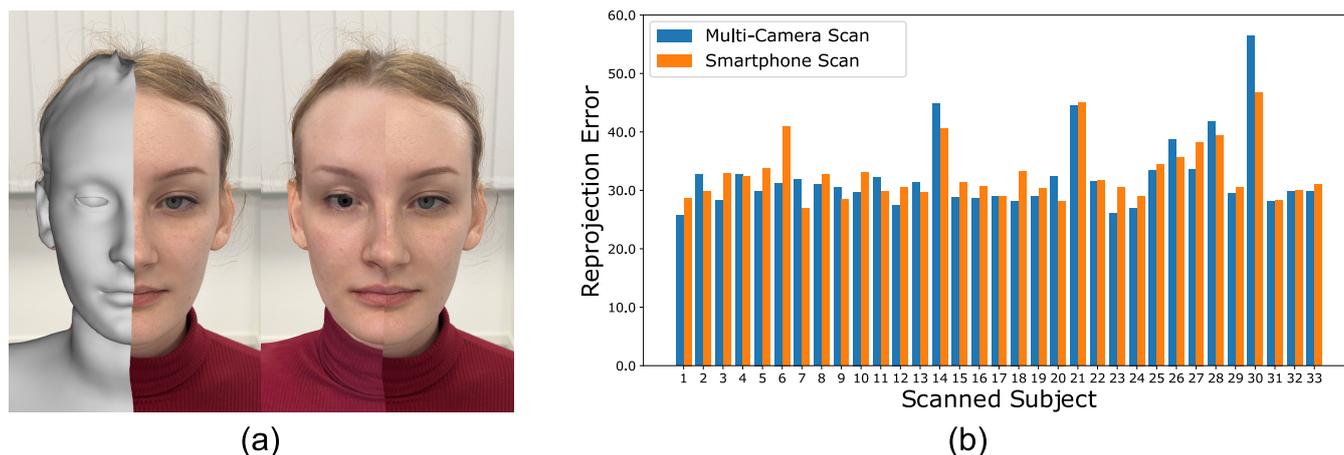


Figure 8. (a): We evaluate the accurate of our approach by reprojecting our avatars (left half of images) into the captured images (right half of images). (b): Despite the significant difference in hardware complexity, the reprojection errors of our smartphone scans are only slightly worse than those of the multi-camera rig.

352 fitting (due to simultaneously fitting to body and head point clouds, and requiring less regularization). The
 353 results of Wenninger et al. (2020) suffer from considerably less geometric details and wrong body-to-head
 354 proportions. These differences are even more prevalent for the lower subject, where the head is unnaturally
 355 deformed due to camera misalignments in the photogrammetry step (see Figure 5 c).

356 In order to evaluate whether neural avatars are a viable alternative to mesh-based avatars for VR
 357 applications, we experimented with three recent approaches: the NeRF-based InstantAvatar (Jiang et al.,
 358 2023) and the 3DGS-based methods SplattingAvatar (Shao et al., 2024) and GART (Lei et al., 2024) – since
 359 those methods reported fast training times and high rendering performance. To achieve optimal results, we
 360 followed the recommendations of these projects and use their training scripts. Since all three methods can
 361 reconstruct avatars from the People Snapshot format (Alldieck et al., 2018b), we recorded equivalent videos
 362 (1080×1080 pixels, subject rotating in A-pose) using the same iPhone as for our scans. These videos
 363 are then converted to the People Snapshot format using (Alldieck et al., 2018b), and the per-frame SMPL
 364 poses are refined using Anim-NeRF (Chen et al., 2021). On this prepared data we ran InstantAvatar using
 365 their provided scripts. The data resulting from InstantAvatar then act as input for running SplattingAvatar
 366 and GART. We used batch size 4 for InstantAvatar and 2min training of GART, as these produced the best
 367 results. The results from different configurations are shown in the supplementary material.

368 The (required) video pre-processing (landmark detection, segmentation, VideoAvatars pipeline, and
 369 Anim-NeRF) took about 17 hours on a compute server with three Nvidia RTX 6000 having 48 GB GPU
 370 memory each. Training of InstantAvatar, SplattingAvatar, and GART took another 2–5 minutes, 25 minutes,
 371 and 2 minutes, respectively, on a different server with Nvidia RTX A5000 and 24 GB GPU memory. With
 372 more than 17 hours, the overall reconstruction time of these methods is 45 times longer than ours.

373 Figure 10 shows the resulting avatars in a training pose and from a training camera view. While
 374 InstantAvatar produces visual clutter, both SplattingAvatar and GART are visually appealing – although
 375 more blurry than our reconstruction. However, as shown in Figure 11 and the accompanying video, when it
 376 comes to novel poses and/or novel viewpoints, the quality of neural avatars quickly degrades to a level not
 377 acceptable for social VR applications. In particular, for multi-avatar VR applications, where users/avatars
 378 can take on arbitrary poses and be viewed from arbitrary camera positions, the artifacts shown in Figure 11
 379 are more the rule rather than the exception. Avoiding these generalization problems would require much

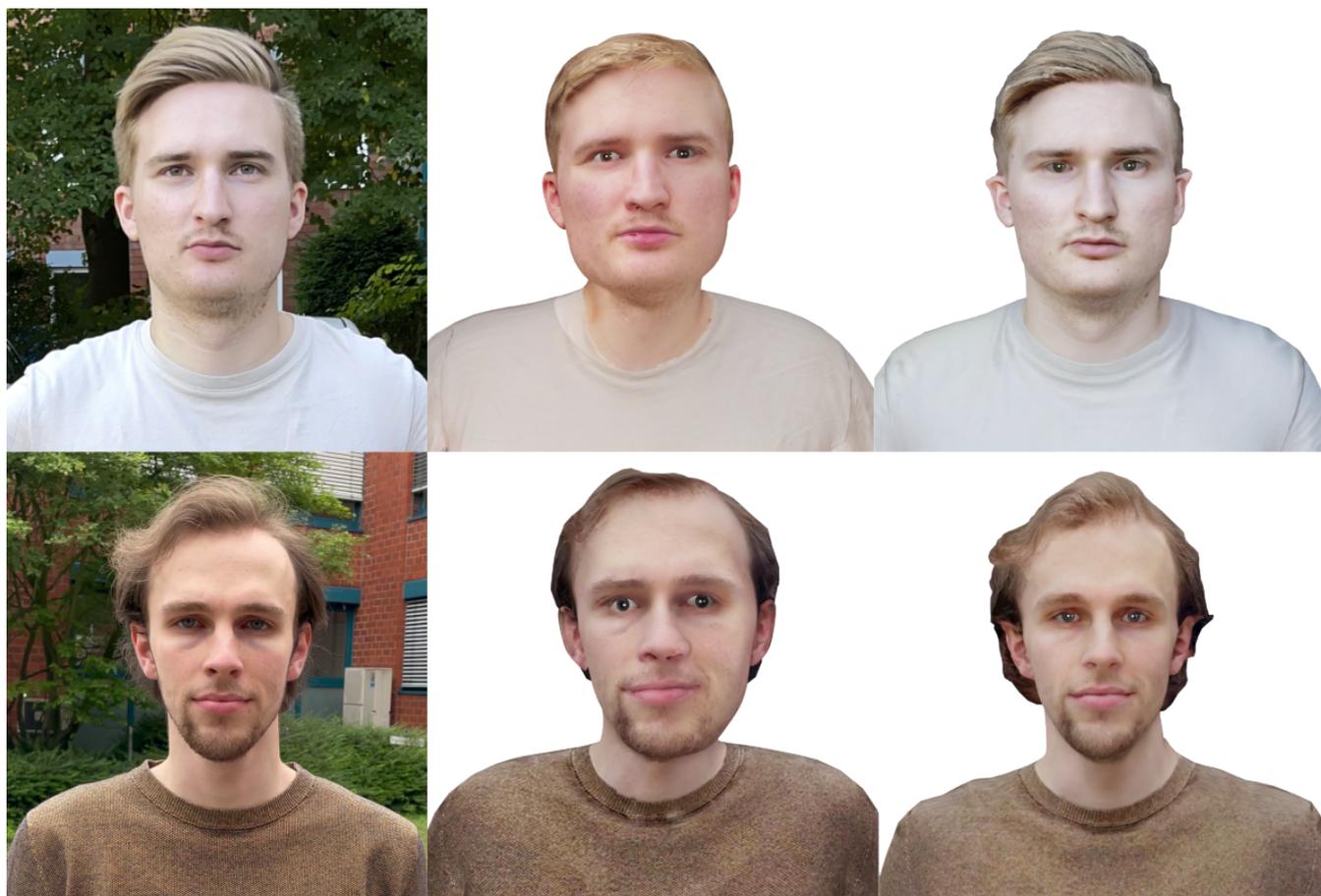


Figure 9. Two subjects captured in an outdoor environment (left), with avatars reconstructed using (Wenninger et al., 2020) (center) and our approach (right). Our avatars are considerably more accurate in terms of geometry and texture, while those of Wenninger et al. (2020) suffer from photogrammetry misalignments, (required) strong regularization, and wrong body-to-head proportions.

380 more training data, i.e., capturing the subject in significantly more poses and from significantly more
381 camera views, which in turn would make the scanning significantly more complex and the reconstruction
382 significantly more expensive – therefore requiring a complex multi-camera video recording setup. Our
383 mesh-based avatars, in contrast, are sufficiently regularized by the statistical human body template and its
384 animation controls to enable generalization to novel views and novel poses, even when captured from 105
385 smartphone images only.

386 In addition to long reconstruction times and suboptimal visual quality, the rendering performance
387 of neural avatars is not (yet) sufficient for VR applications, where currently HMDs require around
388 90 fps stereoscopic rendering (i.e., 180 fps monoscopic rendering) at about $2k$ resolution per eye. We
389 therefore evaluated the rendering performance at a resolution of 2160×2160 pixels on a VR workstation
390 (AMD Ryzen 9 7950X CPU, RTX 4090 24GB GPU, 64GB RAM). For this monoscopic rendering,
391 InstantAvatar achieved 1.22 fps (819 ms/frame), SplattingAvatar 171 fps (5.9 ms/frame), and GART 245 fps
392 (4.1 ms/frame). While InstantAvatar was far from the required frame rate for VR application, rendering
393 performance of SplattingAvatar and GART were just on the edge of VR applicability. However, contrary to
394 these rather low performance results, our mesh-based representation can be rendered at 4615 fps – therefore
395 comfortably allowing even multiple avatars in the same virtual environment at the same time.

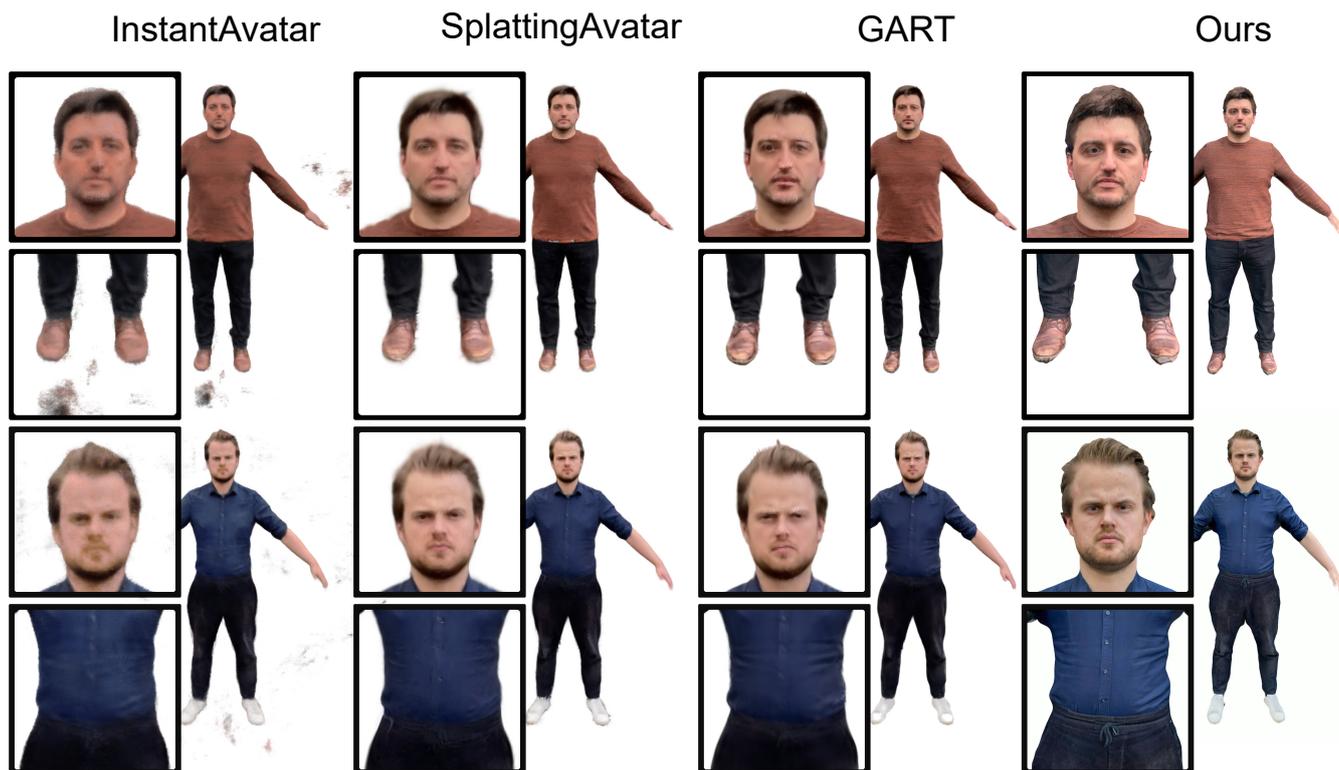


Figure 10. Reconstructed avatars rendered in a pose and view from the training data. From left to right: InstantAvatar (Jiang et al., 2023), SplattingAvatar (Shao et al., 2024), GART (Lei et al., 2024), and ours.

5 USER STUDY

396 We conducted a comprehensive user study following a multi-method approach. The study evaluated
397 the quality of our generated avatars (called below *smartphone avatars*) by measuring their impact on a
398 prominent selection of well-known and often studied avatar effects. In addition, it also evaluated the general
399 usability and user satisfaction of the smartphone front-end. The purpose was to assess the quality of the
400 avatars subjectively and to improve the user experience of scanning and being scanned with the smartphone
401 app.

402 To this end, we arranged the participants into dyads, where one participant had to perform a smartphone
403 app scan of another participant. While the scanning participant evaluated the app's usability afterward (in
404 the following called *smartphone app evaluation*), the scanned participant assessed the perception of the
405 scanning processes and the generated avatar (in the following called *avatar evaluation*).

406 For the smartphone app evaluation, participants performing the smartphone scan were asked to assess
407 the app's usability using standardized questionnaires, allowing for comparison with validated benchmarks.
408 Additionally, we conducted semi-structured interviews to gather more feedback on the user experience of
409 both scanning and being scanned with the smartphone app. The results are used as part of a user-centered
410 design process to improve the app.

411 For the avatar evaluation, we adopted and extended the approach from Bartl et al. (2021) and utilized a
412 counterbalanced within-subject design comparing our generated smartphone avatars to (a) photorealistically
413 reconstructed personalized avatars from a state-of-the-art expert system (in the following called *camera rig
414 avatar*, see Section 4.1) and (b) gender- and ethnicity-matched generic avatars. We chose condition (a) to
415 compare the quality of our *smartphone avatars* to the quality of personalized avatars frequently used in

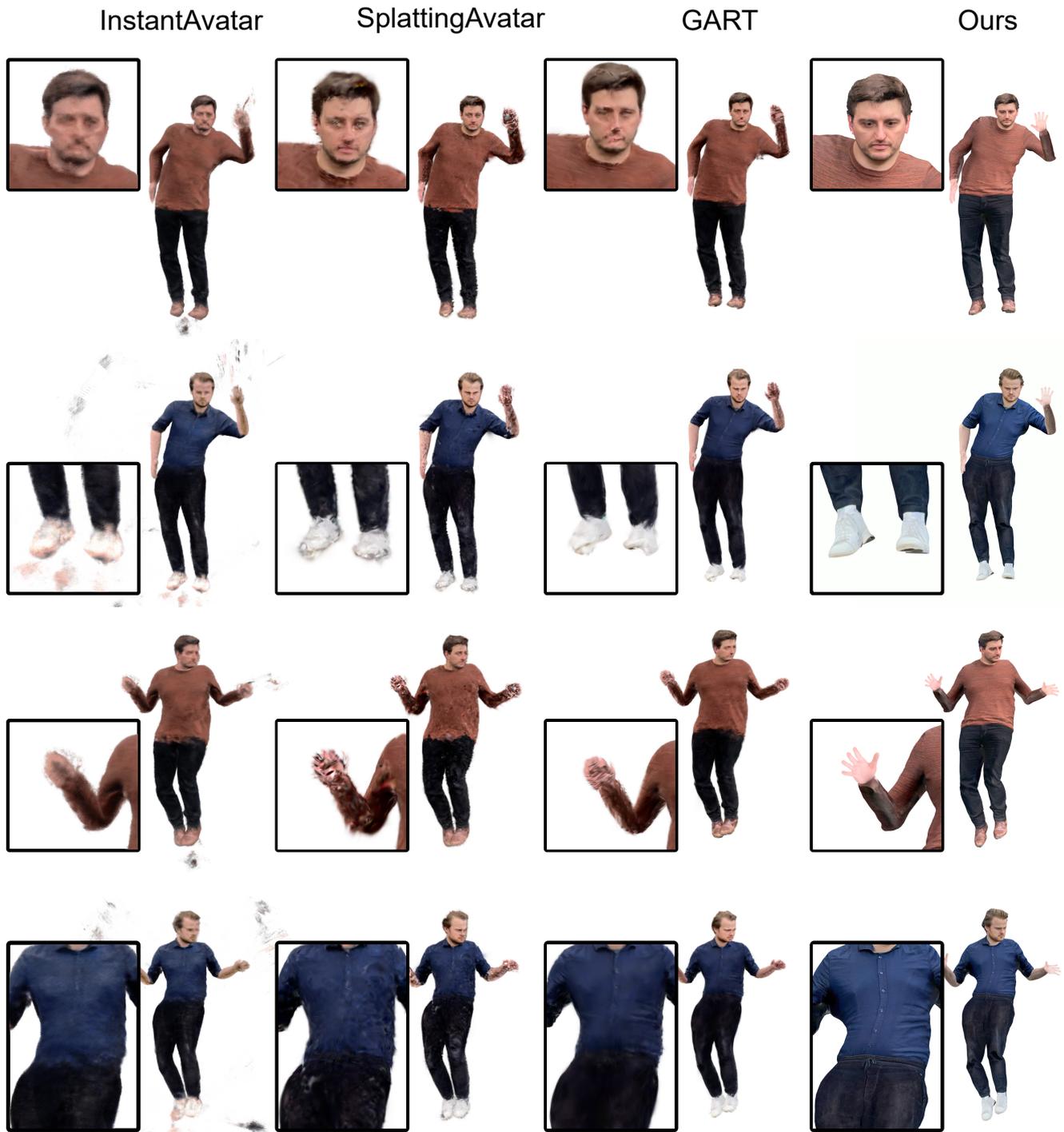


Figure 11. Avatar reconstructions animated in novel poses. From left to right: InstantAvatar (Jiang et al., 2023), SplattingAvatar (Shao et al., 2024), GART (Lei et al., 2024), and our result. Although InstantAvatar, SplattingAvatar, and GART produced visually appealing results in the training poses, the reconstructions get noisier and blurrier in novel poses. Details, e.g., hands and faces, are hardly recognizable anymore and the renderings get even blurrier. Our results in contrast are as sharp and detailed as in the training poses.

416 recent avatar research but reconstructed by a rather costly and complex technical setup. This allowed us to
 417 assess the impact of the proposed method's reconstruction quality on typical well-known and often studied
 418 avatar effects in relation to the much lower technical requirements of our method. We chose condition (b)

419 to assess the quality of both reconstruction methods in terms of self-similarity and self-attribution of the
420 resulting avatars and to measure the overall effect of personalization. An additional in-VR comparison to
421 any of the neural avatars (see Section 4.2) unfortunately was still unreasonable. The given state of the art
422 was inadequate for a VR evaluation due to the low rendering performance but most importantly due to the
423 significant artifacts (see Figure 11 and video), specifically given the arbitrary poses and camera positions
424 typical for the VR exposure.

425 During individual *one-by-one exposures*, the scanned participants embodied each of the three avatar
426 types successively while engaging in various body-centered movement tasks in front of a virtual mirror
427 within a VR environment. Afterward, they evaluated the avatars regarding (a) sense of embodiment
428 and self-identification, (b) plausibility, and (c) uncanny valley effects. In a final *side-by-side exposure*,
429 participants simultaneously embodied each type of avatar while observing them exclusively from an
430 allocentric perspective in three different virtual mirrors (one for each type) and answering different
431 preference questions. Afterward, we asked the participants why they preferred their chosen avatars.

432 5.1 Apparatus

433 5.1.1 Avatars

434 In the following, we explain the integration of the three different avatar types utilized in our study.

435 Each participant attending the smartphone app evaluation (in the following called *scanning participant*)
436 used our smartphone app to create a personalized avatar for the corresponding participant attending the
437 avatar evaluation (in the following called *scanned participant*). We maintained uniform lighting conditions
438 to enhance the avatars' comparability with the camera rig avatars. The scanning participant received
439 instructions from the smartphone app tutorial and was directed to guide the scanned participant accordingly.
440 No further post-processing was performed on the smartphone avatars.

441 We created a personalized camera rig avatar for each participant in the avatar evaluation using the
442 expert body scanner of the Embodiment Lab at the University of Würzburg (see Section 4.1). No further
443 post-processing was performed on the camera rig avatars.

444 Since avatars that do not match the user's gender and ethnicity have been shown to impact SoE particularly
445 negatively (Do et al., 2024) and consequently would lead to an unequal comparison with personalized
446 avatars that are matched in gender and ethnicity, we decided to match both between user and generic
447 avatars. To this end, we chose the Validated Avatar Library for Inclusion and Diversity (VALID) (Do et al.,
448 2023). Through a LimeSurvey questionnaire, each participant in the avatar evaluation was asked to select
449 the VALID avatar that most closely matched their own gender and ethnicity. As the participants typically
450 attend studies dressed casually, they could choose between 42 casually dressed VALID avatars, consisting
451 of three male and three female avatars, each of seven different ethnicities.

452 5.1.2 Virtual Reality System

453 The VR system was realized using Unity 2020.3.25f1 (Unity Technologies, 2020). We utilized a Valve
454 Index head-mounted display (HMD) featuring a resolution of 1440×1600 px per eye and a total field
455 of view of $114.1 \times 109.4^\circ$ (Wolf et al., 2022a). Its refresh rate was set to 90 Hz. Participants' hand and
456 finger movements were tracked through two Index controllers and their built-in proximity sensors. Four
457 SteamVR base stations covered the 3×3 m tracking area. All mentioned components were integrated into
458 the VR system using SteamVR version 2.3 (Valve Corporation, 2024a) and its corresponding Unity plug-in
459 version 2.7.3 (Valve Corporation, 2024b). We routed the HMD's cable to a VR-capable workstation (Intel

460 Core i7-7700K CPU, NVIDIA GeForce GTX 1080, 16 GB RAM) running the VR system on Windows
461 10. For body tracking, we utilized the markerless body tracking system from Captury. Body poses were
462 captured using eight FLIR Blackfly S BFS-PGE-16S2C RGB cameras running at 100 Hz, which have been
463 connected via two 4-port 1 GBit/s ethernet frame-grabber to a high-end workstation (NVIDIA GeForce
464 RTX 3080 Ti, 32 GB RAM, AMD Ryzen 9 5900x) running Captury Live in version 259 (Captury, 2023a)
465 on Ubuntu 18 LTS. The body poses were continuously integrated into the VR system using Captury's
466 corresponding Unity plug-in (Captury, 2023b).

467 5.1.3 Avatar Embodiment

468 We realized avatar embodiment by retargeting the participant's tracked body pose to the used avatar in
469 real-time following the joint approaches described in previous work (Döllinger et al., 2022; Wolf et al.,
470 2022b). During a short calibration process, in which the participant had to stand rigidly and upright, the
471 embodied avatar was calibrated to continuously follow the position of the HMD and scaled to match
472 the participant's eye height. To avoid sliding feet and inaccuracies in hand and feet positions caused by
473 variations in skeletal structure, segment lengths, or insufficient hand tracking, we utilized an IK-supported
474 end-effector optimization using FinalIK version 2.1. Due to a higher accuracy and sampling rate, hand
475 positions and finger poses were taken from the Index controllers, while elbow, knee, and foot positions
476 were taken from Captury.

477 5.1.4 Virtual Environment and Tasks

478 Our virtual environment was based on different Unity assets, which we adapted to create a realistically
479 rendered setting. Figure 12 depicts the virtual environment, accommodating up to three virtual mirrors.
480 Following the guidelines for self-observation mirror placement by Wolf et al. (2022a), each virtual mirror
481 was placed at a distance of 1.5 m from the participant during the study.

482 During each **one-by-one exposure**, participants embodied one of the three avatars in the virtual
483 environment, where only the middle virtual mirror was shown. They could either observe their embodied
484 avatar directly from an egocentric perspective or look into the virtual mirror to receive an allocentric
485 perspective. Participants were asked to perform various body movement tasks in front of the virtual mirror
486 to promote visuomotor coupling and induce SoE (Slater et al., 2010; González-Franco et al., 2010). The
487 body movement tasks adhered to a structured protocol adapted from Roth and Latoschik (2020) and can be
488 found in the supplements of this work.

489 During the **side-by-side exposure**, participants embodied all three avatars simultaneously in the virtual
490 environment, where all three virtual mirrors were shown. While they received no egocentric perspective
491 on the avatars, they could observe each avatar through a virtual mirror. The mirrors were labeled with
492 small numbers, and participants responded to four different preference questions by identifying the mirror
493 number displaying their preferred avatar. The assignment of avatars to mirrors changed randomly after
494 each question. The preference questions can be found in the supplements of this work. Figure 12 depicts
495 the side-by-side exposure.

496 5.2 Measures

497 5.2.1 Quantitative Measures

498 We assessed all quantitative measures using previously published questionnaires. When available, we
499 used validated translated German versions of the utilized questionnaires. Otherwise, we used back-and-forth



Figure 12. The three mirrors showing the expert (left), smartphone (middle), and generic (right) avatar of a female participant during the side-by-side exposure.

500 translations to translate the items into German. Participants answered all questionnaires on a MacBook Pro
501 using LimeSurvey (Limesurvey GmbH, 2024).

502 We captured the **usability** of the smartphone app using the System Usability Scale (SUS) (Brooke, 1996).
503 It provides a fast and simple way to assess a system's usability using ten questionnaire items each answered
504 on a 5-point Likert scale. The calculated overall score ranges between 0 and 100 (*100 = highest usability*)
505 and can be compared with benchmarks provided by previous work (Bangor et al., 2009; Sauro and Lewis,
506 2016; Kortum and Sorber, 2015).

507 For assessing **Sense of Embodiment and Self-Identification** (SoE) towards the avatars, we captured
508 virtual body ownership (VBO) and agency (AG) utilizing the corresponding items of the Virtual
509 Embodiment Questionnaire (VEQ) (Roth and Latoschik, 2020) and self-location (SL) using the additional
510 items introduced by Fiedler et al. (VEQ+) (Fiedler et al., 2023). We used the items capturing self-similarity
511 (SS) and self-attribution (SA) from the VEQ+ to assess self-identification towards the avatars. Each factor
512 measured comprises four items rated on a 7-point Likert scale (*7 = highest VBO, AG, SL, SS, and SA*).

513 We captured the avatars' **plausibility** utilizing the Virtual Human Plausibility Questionnaire (VHPQ)
514 (Mal et al., 2022, 2024). It consists of seven items that assess the avatars' appearance and behavior
515 plausibility (ABP) and four items for matching the virtual environment (MVE). Each item is rated on a
516 7-point Likert scale (*7 = highest ABP and MVE*).

517 We captured tendencies of the avatars' appearance towards the **uncanny valley** using the revised version
518 of the Uncanny Valley Index (UVI) (Ho and MacDorman, 2017). It comprises four items each to assess

519 the avatars' humanness (HU) and attractiveness (AT) and eight items to capture the avatars' eeriness (EE).
520 While the items are answered on a range between -3 and 3, we report them on a range between 1 and 7 (7
521 = *highest HU, AT, EE*).

522 As a control measure, we captured participants' physical symptoms associated with **VR sickness** in
523 a pre-post comparison using the Virtual Reality Sickness Questionnaire (VRSQ) (Kim et al., 2018). It
524 consists of nine items, each of which represents a typical symptom of VR sickness and is answered on a
525 scale between 0 and 3 (*3 = highest symptomatology*). The total score of the VRSQ ranges between 0 and
526 100 (*100 = highest VR sickness*).

527 5.2.2 Qualitative Measures

528 We conducted semi-structured interviews to assess the user experiences related to both scanning and
529 being scanned with the smartphone app. The interview protocols incorporated a retrospective thinking-
530 aloud approach (Bowers and Snyder, 1990; Simon and Ericsson, 1993) to comprehensively analyze the
531 interactions with the smartphone app while not influencing the scan experiences. We further included
532 predefined questions to query positive and negative feelings experienced during the use of the app and while
533 being scanned, the app's functionality and its intended purpose, the impact of the scanning participant on
534 comfort or discomfort when being scanned, and the clarity and comprehensibility of the scanning process.
535 Additionally, participants described aspects of the process they found efficient or challenging and reported
536 any problematic incidents they faced. Finally, participants could suggest improvements to both the scan
537 app functionality and the scanning process and were asked about their scan preferences and if they would
538 participate in a body scan again. Participants in the avatar evaluation were further asked which avatar they
539 preferred regarding self-representation similarity, fidelity, plausibility, and suitability, along with reasons
540 behind their choices. The complete interview protocols and exact phrasing of the preference questions can
541 be found in the supplements of this work.

542 5.3 Procedures

543 In the following, we describe the standardized experimental procedures of our smartphone app and
544 avatar evaluations. Figure 13 visualizes both procedures and highlights their intersection during the
545 smartphone app scan. Initially, participants in both procedures received information about the study and
546 privacy, consented to participate, and generated two pseudonymization codes to store personal (i.e., voice
547 recordings and avatars) and evaluation data separately. Subsequently, they proceeded with their respective
548 evaluation procedures.

549 5.3.1 Smartphone App Evaluation

550 Each participant in the smartphone app evaluation first completed a tutorial on how to perform a body
551 scan using the smartphone app. As soon as the other participant arrived for the scan in the laboratory,
552 both participants were introduced to each other. The participant performing the scan verified that all
553 requirements for the scan were met and instructed the scanned participant not to speak or move during the
554 scan. To ensure that an assessable avatar was generated, the scanning participant performed two successive
555 scans. After scanning, the scanned participant left the laboratory, and the scanning participant answered
556 the SUS questionnaire using LimeSurvey. Following that, the participant was interviewed and completed
557 demographics. On average, the entire smartphone app evaluation took approximately 41 min.

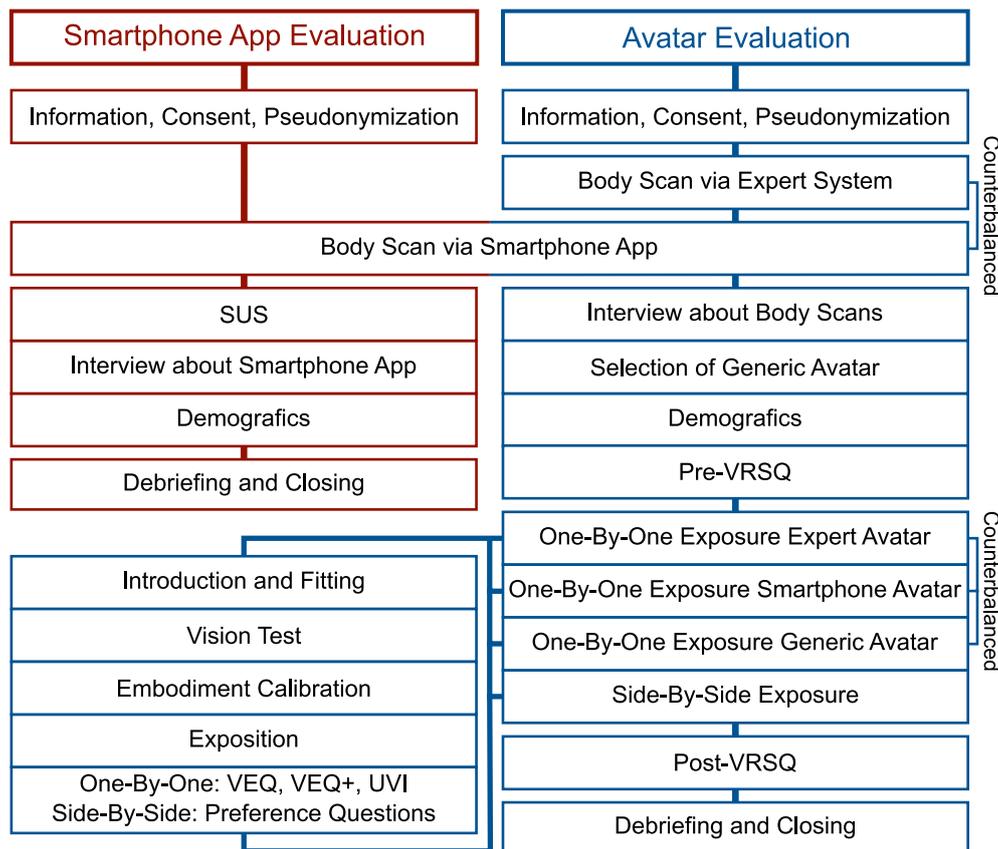


Figure 13. Experimental procedure of a dyad, illustrating the process of evaluating the smartphone app (left) and the avatars (right).

5.3.2 Avatar Evaluation

Each participant in the avatar evaluation first participated in a smartphone and expert scan conducted in a counterbalanced order. After the scans, the participant was interviewed about the scan processes, chose a generic avatar as described above, completed the demographics, and answered the pre-VRSQ. The one-by-one exposures followed in a counterbalanced order, each lasting on average 7.6 min. After each exposure, the participant answered the VEQ, VEQ+, and UVI. The following side-by-side exposure averaged 4.2 min and was accompanied by the preference questions answered verbally in VR. For each exposure, a vision test and the avatar embodiment calibration were performed following the instructions on a virtual whiteboard. In addition, the participant received audio instructions for all tasks. Finally, the participant completed the post-VRSQ. On average, the entire avatar evaluation lasted 103 min.

5.4 Participants

Adhering to the ethical standards of the Declaration of Helsinki, our study received approval from the ethics review board of the Institute Human-Computer-Media (MCM) at the University of Würzburg¹. We recruited a total of 66 participants organized into 33 dyads using the local participant management system and compensated them either by course credits or cash, both depending on the duration of their participation. In none of the dyads, participants knew each other before the study. All participants had normal or corrected vision and no hearing impairment. Participants evaluating the smartphone app (19

¹ <https://www.mcm.uni-wuerzburg.de/forschung/ethikkommission/>

575 female, 14 male) were aged between 19 and 41 ($M = 26.60, SD = 5.48$). None of them had used the
576 smartphone app before. Participants evaluating the avatars (25 female, 8 male) were aged between 20 and
577 49 ($M = 27.64, SD = 6.90$). While none of them had been scanned with the smartphone app before, nine
578 participants had previously taken part in an expert scan. Most participants in the avatar evaluation (29
579 White, 2 Asian, 1 MENA) chose a generic avatar that matched their ethnicity. Only one White participant
580 chose a Hispanic avatar. Ten participants used VR for the first time, 20 up to ten times, one more than ten
581 times, and two more than 20 times.

582 We excluded one dyad from our statistical analysis as one participant used the smartphone app contrary
583 to the instructions, resulting in an unusable avatar. While all participants stated that they had more than
584 five years of experience with the German language, we had to exclude another participant from the avatar
585 evaluation as the experimenter felt that the participant did not understand the questions and instructions
586 correctly, which was confirmed by implausible answers and outliers in the data. Hence, 32 datasets
587 remained for the smartphone app and 31 for the avatar evaluation.

588 5.5 Data Analysis

589 We conducted all quantitative analyses using SPSS version 29.0.2.0 (IBM, 2022). Before running the
590 statistical tests, we checked whether our data met the assumption of normality and sphericity for parametric
591 testing. Shapiro-Wilk tests showed clear violations of the normality assumption for both dimensions of
592 the VHPQ and minor violations for VEQ agency and VEQ+ self-location. Mauchly's test for sphericity
593 confirmed homoscedasticity between the groups for all of our measures. Since variance analysis shows
594 robustness to slight violations of normality for groups with $N \geq 30$ (Wilcox, 2022), we decided to perform
595 parametric tests for all measures except those from the VHPQ. All main tests have been performed against
596 an α of .05, while post-hoc tests have been Bonferroni adjusted.

597 The qualitative feedback has been analyzed following the principles of thematic analysis (Braun and
598 Clarke, 2006). Due to space restrictions, we decided to report the results mainly based on the frequency of
599 certain feedback while mostly refraining from direct quotes.

600 5.6 Results

601 5.6.1 Smartphone App Evaluation

602 The quantitative evaluation of the smartphone app's usability resulted in a reasonably high SUS score
603 ($M = 78.83, SD = 12.23$). We compared the results to absolute benchmarks from existing literature.
604 According to Sauro and Lewis (2016), our smartphone app shows above-average usability. While a score
605 between 77.2 and 78.8 leads to a usability grade of *B+*, a score between 78.9 and 80.7 relates to an *A-*.
606 This grade matches the classifications of the adjective rating scale of Bangor et al. (2009), where a score
607 above 71.4 is considered *good*, while a score above 85.5 would be *excellent*. According to the work of
608 Kortum and Sorber (2015), our smartphone app's usability can almost keep up with the usability of the ten
609 most-used iPhone apps, which have an average SUS score of 79.3.

610 When analyzing interviews about the usability of the smartphone app, the majority of the 32 participants
611 performing the smartphone app scan found it highly usable. Twenty-nine participants found the app's
612 functionality and purpose easy to understand, while 26 reported that they constantly knew how to
613 use it. As particularly useful features, 20 participants highlighted the overlay for controlling scan
614 distance and movement, 16 participants the initial tutorial, and five participants the arrows indicating
615 the movement direction. Nonetheless, challenges were also noted. Twenty-three participants reported

Table 1. Exact descriptive values for each measure of the avatar evaluation per group and statistical results of the group comparisons.

	Smartphone	Camera Rig	Generic	Group Comparisons
	<i>M (SD)</i>	<i>M (SD)</i>	<i>M (SD)</i>	
Sense of Embodiment				
VEQ Ownership (VBO)	4.10 (1.50)	4.79 (1.28)	3.75 (1.44)	$F(2, 60) = 11.011, p < .001, \eta_p^2 = .268$
VEQ Agency (AG)	5.57 (0.96)	5.97 (0.79)	5.79 (0.81)	$F(2, 60) = 2.845, p = .066, \eta_p^2 = .087$
VEQ+ Self-Location (SL)	4.11 (1.07)	4.14 (1.03)	3.73 (1.13)	$F(2, 60) = 3.502, p = .036, \eta_p^2 = .275$
VEQ+ Self-Similarity (SS)	5.69 (1.11)	5.85 (0.68)	2.69 (1.23)	$F(2, 60) = 82.651, p < .001, \eta_p^2 = .734$
VEQ+ Self-Attribution (SA)	4.69 (1.21)	4.99 (0.99)	3.37 (1.16)	$F(2, 60) = 31.390, p < .001, \eta_p^2 = .511$
Plausibility				
VHPQ Appearance/Behaviour (ABP)	4.84 (0.82)	5.33 (0.78)	5.27 (0.78)	$\chi^2(2) = 4.581, p = .101, W = .074$
VHPQ Match to VE (MVE)	5.15 (0.98)	5.47 (1.14)	5.80 (0.65)	$\chi^2(2) = 5.782, p = .056, W = .093$
Uncanny Valley				
UVI Humanness (HU)	3.46 (1.16)	3.84 (1.08)	3.36 (0.90)	$F(2, 60) = 2.444, p = .095, \eta_p^2 = .075$
UVI Eeriness (EE)	4.05 (0.79)	3.87 (0.97)	3.15 (0.79)	$F(2, 60) = 19.313, p < .001, \eta_p^2 = .392$
UVI Attractiveness (AT)	3.90 (1.20)	4.13 (0.84)	4.52 (0.69)	$F(2, 60) = 3.264, p = .045, \eta_p^2 = .098$

616 difficulties maintaining an appropriate moving pace while scanning, with six participants emphasizing this
617 problem, especially for the head scan. Similarly, seven and six participants reported issues with aligning
618 the overlay while moving and keeping the correct distance, respectively. Six participants mentioned the
619 need for high concentration, and 18 felt a bit uncomfortable due to the close proximity to the scanned
620 participant. Six participants considered the relatively long duration of the scan process as unpleasant. To
621 address the mentioned aspects, eight participants suggested a more detailed tutorial, and another four
622 suggested an initial overlay mapping to the height of the scanned participant. To improve the scan process,
623 five participants recommended more interaction with the scanned person, five more additional feedback on
624 pacing their movement during the scan, and another five stressed the need to shorten the scan duration.

625 In addition to feedback on performing the scan, we obtained reports from the 32 scanned participants
626 on their scanning experience. Overall, the process was clear and manageable, with 30 participants
627 completely understanding the required actions. All participants confirmed their willingness to participate in
628 a smartphone app scan again. However, compared to expert scans, 21 participants noted the smartphone
629 app scan was slower, and 22 found it less comfortable. Prolonged posing discomfort was mentioned
630 by twelve participants, while wardrobe and hairstyle constraints were issues for another four. Fourteen
631 participants anticipated a difference between an expert and a beginner performing the smartphone scan,
632 with four believing the expert would be faster. When asked about suggestions for improvement, four
633 participants indicated that they would accelerate the process to reduce the discomfort of holding the scan
634 pose. Regarding the head scan, four participants suggested a fixation to aid focus, and three to increase the
635 distance between the camera and the head.

636 5.6.2 Avatar Evaluation

637 To perform group comparisons on our avatar evaluation data, we calculated either a repeated-measures
638 ANOVA for measures that met the requirements for parametric analysis or Friedman tests as a non-
639 parametric alternative. The descriptive data and the results of the group comparisons can be found in
640 Table 1. For all tests revealing significant differences between groups, we calculated Bonferroni-corrected
641 pairwise post-hoc comparisons that are reported in Figure 14.

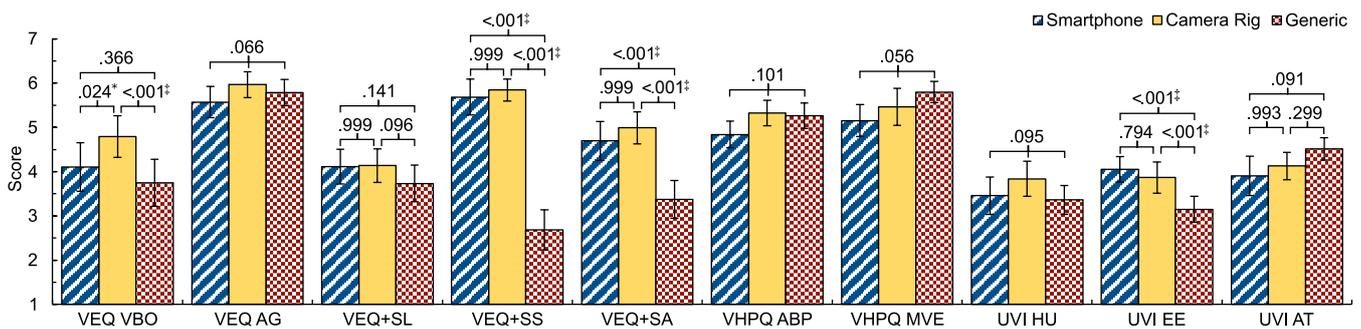


Figure 14. Bar charts for each measure and each group of the avatar evaluation, including statistical test results of the group comparisons and post-hoc tests where applicable. Error bars represent 95 % confidence intervals. Statistical significance indicators: * $p < .05$; † $p < .01$; ‡ $p < .001$.

642 During the side-by-side exposure, we asked participants about their preferences regarding self-
 643 representation similarity, fidelity, plausibility, and suitability, along with reasons behind their choices. Out
 644 of the 31 participants included in the analysis, 16 perceived the smartphone avatars to be more similar to
 645 themselves, while 13 preferred the camera rig avatars. Regarding self-representation fidelity, 11 participants
 646 preferred the smartphone avatars, 19 chose the camera rig avatars, and one favored the generic one. To feel
 647 most plausibly represented in VR, 12 participants chose the smartphone avatars, 18 the camera rig avatars,
 648 and one the generic one. When asked which avatar the participants would prefer to be represented in VR,
 649 10 chose the smartphone avatars, 17 the camera rig avatars, and four the generic ones. When asked for
 650 their reasoning, participants favoring smartphone avatars mostly mentioned a detailed facial reconstruction
 651 and realism as key factors. Those participants who preferred camera rig avatars highlighted the accuracy
 652 of body shape reconstruction, noting issues with smartphone avatars' body proportions, particularly the
 653 arms. Participants who chose generic avatars consistently did so because of overall dissatisfaction with
 654 their personal appearance rather than avatar quality.

6 DISCUSSION

655 In this section, we discuss the results of the comparisons with the different avatar reconstruction methods
 656 and the results of our user study and present the limitations of our work.

657 6.1 Smartphone App Evaluation

658 We evaluated the usability of our smartphone app quantitatively using the SUS questionnaire and
 659 qualitatively using semi-structured interviews, including a retrospective thinking-aloud approach. The
 660 SUS results showed that our smartphone app is already well usable. The qualitative feedback confirmed
 661 this impression and highlighted the overlay and tutorial as particularly positive features. However, the
 662 qualitative feedback also revealed areas for improvement.

663 As part of the user-oriented design process, we already incorporated suggested improvements. To address
 664 comments regarding the duration of the scan and the pace, we added the option to shorten or extend the
 665 scan speed using technical means. The unclear parts of the tutorial have been improved to prepare users
 666 for the scan better. Furthermore, we have also added warnings if the scanned person is not sufficiently
 667 centered. Other feedback could not be implemented due to technical limitations or requires further research.
 668 For example, the distance between the smartphone and the scanned person, especially during the head
 669 scan, could only be increased by the loss of detail in the reconstructed avatars. However, since the high

Table 2. Comparison of the different avatar reconstruction methods regarding our requirements. The symbols represent: ✓ completely, ● partially, ✗ not fulfilled.

Method	Easy	Fast	Affordable	Realistic	Full-Body	VR-Ready
Achenbach et al. (2017)	✗	✓	✗	✓	✓	✓
Wenninger et al. (2020)	✗	✓	✓	✓	✓	✓
Jiang et al. (2023)	●	✗	✓	●	✓	✗
Shao et al. (2024)	●	✗	✓	●	✓	✗
Lei et al. (2024)	●	✗	✓	●	✓	✗
Ours	✓	✓	✓	✓	✓	✓

670 quality of the faces is a significant advantage of our system, we decided to keep the required distance.
 671 Furthermore, the interaction between the scanning and scanned person and visual aids (e.g., fixation point)
 672 for the scanned person lies outside the influence of our smartphone application.

673 6.2 Avatar Evaluation

674 Our quantitative and qualitative comparisons in Section 4 demonstrate that our smartphone application
 675 enables even non-experts to reconstruct avatars of a similar quality and accuracy as those produced with an
 676 expert-operated multi-camera rig – at a fraction of the price, complexity, and required expertise. Compared
 677 to the previous smartphone-based reconstruction (Wenninger et al., 2020), our proposed method is easier to
 678 use and gives higher-quality results even in more challenging in-the-wild scenarios.

679 Our experiments with InstantAvatar (Jiang et al., 2023), SplattingAvatar (Shao et al., 2024), and GART
 680 Lei et al. (2024) revealed that neural avatars generalize rather poorly to poses and camera views far
 681 from training data – a situation that cannot be avoided in multi-avatar VR applications. Although these
 682 generalization problems can be reduced with more training data, this is beyond the capabilities of a simple
 683 smartphone-based scanning solution. Also in terms of reconstruction times and rendering performance are
 684 neural avatars not yet suitable for VR applications, such that classical mesh-based avatars appear to still be
 685 the preferred representation. Table 2 summarizes the fulfillment of our requirements with respect to avatar
 686 reconstruction. Although many methods show their strengths in a subset of the criteria, only our system
 687 fulfills all of them.

688 Compared to the work of Waltemate et al. (2018), our user study confirmed that realistic avatars still offer
 689 substantial benefits over generic avatars for self-representation, even when the generic avatars are also
 690 personalized in gender and ethnicity (Do et al., 2023, 2024). With regard to the comparison, some further
 691 notable findings need to be addressed. The statistically significant difference in virtual body ownership
 692 between the smartphone and camera rig avatars can potentially be attributed to observed motion artifacts,
 693 which can degrade the avatars' appearance. However, smartphone avatars still perform descriptively
 694 better than generic avatars. Regarding self-identification, the smartphone and camera rig avatars both
 695 show significant advantages to generic avatars, although the smartphone avatars were generated using
 696 a significantly cheaper method than the camera rig avatars. For the smartphone avatars, participants
 697 emphasized particularly the high similarity of the head. However, results also showed that the eeriness of
 698 realistic avatars was significantly higher than generic avatars. This is likely attributable to an Uncanny
 699 Valley effect originating from the emotional relatedness to self-personalized avatars, which has also been
 700 observed in other research (Mori et al., 2012; Döllinger et al., 2023). When considering the plausibility of

701 the avatars, it is noticeable that the reconstruction described most realistically had the lowest match with
702 the perceived plausibility. This discrepancy might be attributed to the incongruence between the virtual
703 environment's realistic style and the avatars' photorealistic style (Latoschik and Wienrich, 2022).

704 6.3 Limitations

705 Since our method uses photogrammetry software to generate point clouds from images, the input images
706 must contain as little movement as possible. If movement occurs in the background, the segmentation
707 significantly improves the photogrammetry results. However, the motions of the scanned subject violate the
708 photogrammetry assumption, i.e., that the scanned object is rigid and not moving, leading to less accurate
709 point clouds and, therefore, geometric deformations in the final avatar. Figure 7 shows this problem in
710 more detail, as the arms of the second avatar (from left) have visible differences in thickness.

711 We use a mesh-based representation for our avatars. On the one hand, this enables high-performance
712 rendering and novel pose generation. On the other hand, we represent cloth, hair, and skin by a single
713 textured mesh, which can lead to visual artifacts. An interesting direction for future work would be
714 to combine mesh-based avatars (potentially with multiple layers for skin and cloth) with volumetric
715 details (such as hair) represented by Gaussian Splatting – as this would combine the strengths of both
716 representations.

717 Our system uses image segmentation to preprocess the input and mask out regions that do not contain
718 people. For that reason, people in the background are a challenging task, as they are not removed. We want
719 to explore the capabilities of the depth sensor to remove people in the background from the masks.

720 We rely on Apple frameworks for both our scanning client and reconstruction server. For the client, an
721 Android-based app would be possible, but we chose iOS because Apple's photogrammetry yields correctly
722 scaled results thanks to the LiDAR sensor of the pro-level iPhones. On the reconstruction server, only the
723 photogrammetry and the segmentation frameworks are Apple-specific, all other parts of the reconstruction
724 pipeline are cross-platform compatible. We chose Apple's RealityKit and Vision frameworks since in our
725 extensive tests this platform produced the best results while not being limited by restrictive licensing.

726 The sample in our study consisted of white participants only. As this potentially limits the generalizability
727 of our results, in future work a larger population sample with greater variability in age, sex, and ethnicity
728 should be tested.

7 CONCLUSION

729 We presented *Avatars for the Masses*, a system that allows non-expert users to scan people and automatically
730 reconstruct realistic VR-ready full-body avatars that achieve similar perception results compared to avatars
731 reconstructed with expensive expert-operated systems. Inspired by the approach of Wenninger et al. (2020),
732 we presented methods to resolve present obstacles that prevent the wide accessibility of realistic full-body
733 avatars. Our custom smartphone application enables laypeople to easily and quickly capture high-quality
734 input images, which, together with background segmentation and an improved template fitting algorithm,
735 result in more convincing reconstructions while reducing restrictions on scanning locations. Our end-to-end
736 solution computes VR-ready avatars that can be easily integrated into existing VR pipelines. To further
737 empower people to create realistic full-body avatars and encourage more avatar-related studies, we will
738 make *Avatars for the Masses* publicly available for research purposes.

REFERENCES

- 739 Achenbach, J., Waltemate, T., Latoschik, M. E., and Botsch, M. (2017). Fast Generation of Realistic
740 Virtual Humans. In *Proceedings of the ACM Symposium on Virtual Reality Software and Technology*.
741 doi:10.1145/3139131.3139154
- 742 [Dataset] Agisoft (2023). Agisoft Metashape. [visited on 2023-09-29]
- 743 [Dataset] Alice Vision (2025). Meshroom. [visited on 2025-02-12]
- 744 Alldieck, T., Magnor, M., Bhatnagar, B. L., Theobalt, C., and Pons-Moll, G. (2019). Learning to
745 Reconstruct People in Clothing From a Single RGB Camera. In *Proceedings of the IEEE/CVF Conference
746 on Computer Vision and Pattern Recognition (CVPR)*
- 747 Alldieck, T., Magnor, M., Xu, W., Theobalt, C., and Pons-Moll, G. (2018a). Detailed Human Avatars from
748 Monocular Video. In *International Conference on 3D Vision (3DV)*. 98–109. doi:10.1109/3DV.2018.
749 00022
- 750 Alldieck, T., Magnor, M., Xu, W., Theobalt, C., and Pons-Moll, G. (2018b). Video Based Reconstruction of
751 3D People Models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition
752 (CVPR)*
- 753 [Dataset] Apple Inc. (2023a). ARKit Framework - Blendshapes. [visited on 2023-09-27]
- 754 [Dataset] Apple Inc. (2023b). RealityKit Framework - Object Capture - PhotogrammetrySession. [visited
755 on 2024-07-05]
- 756 [Dataset] Apple Inc. (2023c). Vision Framework - VNFaceLandmarks2D. [visited on 2024-07-05]
- 757 [Dataset] Apple Inc. (2023d). Vision Framework - VNGeneratePersonSegmentationRequest. [visited on
758 2024-07-05]
- 759 [Dataset] Apple Inc. (2025). Vision Framework - VNDetectHumanHandPosesRequest. [visited on
760 2025-02-12]
- 761 Aseeri, S. and Interrante, V. (2021). The Influence of Avatar Representation on Interpersonal
762 Communication in Virtual Social Environments. *IEEE Transactions on Visualization and Computer
763 Graphics* 27, 2608–2617. doi:10.1109/TVCG.2021.3067783
- 764 Bailenson, J. N. and Blascovich, J. (2004). Avatars. In *Encyclopedia of Human-Computer Interaction*
765 (Great Barrington, MA, USA: Berkshire Publishing Group). 64–68
- 766 Bangor, A., Kortum, P., and Miller, J. (2009). Determining what individual SUS scores mean: Adding an
767 adjective rating scale. *Journal of usability studies* 4, 114–123
- 768 Bartl, A., Wenninger, S., Wolf, E., Botsch, M., and Latoschik, M. E. (2021). Affordable but not Cheap: A
769 Case Study of the Effects of Two 3D-Reconstruction Methods of Virtual Humans. *Frontiers in Virtual
770 Reality* 2. doi:10.3389/frvir.2021.694617
- 771 Bowers, V. A. and Snyder, H. L. (1990). Concurrent versus Retrospective Verbal Protocol for Comparing
772 Window Usability. *Proceedings of the Human Factors Society Annual Meeting* 34, 1270–1274. doi:10.
773 1177/154193129003401720
- 774 Braun, V. and Clarke, V. (2006). Using thematic analysis in psychology. *Qualitative Research in Psychology*
775 3, 77–101. doi:10.1191/1478088706qp063oa
- 776 Brooke, J. (1996). SUS: A quick and dirty usability scale. In *Usability evaluation in industry* (Taylor &
777 Francis). 4–7
- 778 Cao, Z., Hidalgo Martinez, G., Simon, T., Wei, S., and Sheikh, Y. A. (2019). OpenPose: Realtime
779 Multi-Person 2D Pose Estimation using Part Affinity Fields. *IEEE Transactions on Pattern Analysis and
780 Machine Intelligence*
- 781 [Dataset] Captury (2023a). CapturyLive 259
- 782 [Dataset] Captury (2023b). Unity plugin

- 783 [Dataset] Chen, J., Zhang, Y., Kang, D., Zhe, X., Bao, L., Jia, X., et al. (2021). Animatable Neural
784 Radiance Fields from Monocular RGB Videos
- 785 [Dataset] Chen, L.-C., Papandreou, G., Schroff, F., and Adam, H. (2017). Rethinking atrous convolution
786 for semantic image segmentation
- 787 de Vignemont, F. (2011). Embodiment, ownership and disownership. *Consciousness and Cognition* 20,
788 82–93. doi:<https://doi.org/10.1016/j.concog.2010.09.004>
- 789 [Dataset] Dlib (2022). Dlib C++ Library
- 790 Do, T. D., Isabella Protko, C., and McMahan, R. P. (2024). Stepping into the Right Shoes: The Effects
791 of User-Matched Avatar Ethnicity and Gender on Sense of Embodiment in Virtual Reality. *IEEE*
792 *Transactions on Visualization and Computer Graphics*, 1–10doi:10.1109/TVCG.2024.3372067
- 793 Do, T. D., Zelenty, S., Gonzalez-Franco, M., and McMahan, R. P. (2023). VALID: A perceptually validated
794 virtual avatar library for inclusion and diversity. *Frontiers in Virtual Reality* 4. doi:10.3389/frvir.2023.
795 1248915
- 796 Döllinger, N., Mal, D., Keppler, S., Wolf, E., Botsch, M., Israel, J. H., et al. (2024). Virtual body swapping:
797 A vr-based approach to embodied third-person self-processing in mind-body therapy. In *Proceedings of*
798 *the 2024 CHI Conference on Human Factors in Computing Systems*. 1–18
- 799 Döllinger, N., Beck, M., Wolf, E., Mal, D., Botsch, M., Latoschik, M. E., et al. (2023). "If it's not me it
800 doesn't make a difference" – The impact of avatar customization and personalization on user experience
801 and body awareness in virtual reality. In *2023 IEEE International Symposium on Mixed and Augmented*
802 *Reality (ISMAR)*. doi:10.1109/ISMAR59233.2023.00063
- 803 Döllinger, N., Wolf, E., Mal, D., Wenninger, S., Botsch, M., Latoschik, M. E., et al. (2022). Resize Me!
804 Exploring the User Experience of Embodied Realistic Modulatable Avatars for Body Image Intervention
805 in Virtual Reality. *Frontiers in Virtual Reality* 3. doi:10.3389/frvir.2022.935449
- 806 Feng, A., Suma, E., and Shapiro, A. (2017). Just-in-Time, Viable, 3D Avatars from Scans. In *ACM*
807 *SIGGRAPH 2017 Talks*. doi:10.1145/3084363.3085045
- 808 Fiedler, M. L., Wolf, E., Döllinger, N., Botsch, M., Latoschik, M. E., and Wienrich, C. (2023). Embodiment
809 and personalization for self-identification with virtual humans. In *IEEE Conference on Virtual Reality and*
810 *3D User Interfaces Abstracts and Workshops (VRW)*. 799–800. doi:10.1109/VRW58643.2023.00242
- 811 Fiedler, M. L., Wolf, E., Döllinger, N., Mal, D., Botsch, M., Latoschik, M. E., et al. (2024). From avatars
812 to agents: Self-related cues through embodiment and personalization affect body perception in virtual
813 reality. *IEEE Transactions on Visualization and Computer Graphics* Accepted for publication
- 814 Gall, D., Roth, D., Stauffert, J.-P., Zarges, J., and Latoschik, M. E. (2021). Embodiment in virtual reality
815 intensifies emotional responses to virtual stimuli. *Frontiers in Psychology* 12, 3833. doi:10.3389/fpsyg.
816 2021.674179
- 817 González-Franco, M., Pérez-Marcos, D., Spanlang, B., and Slater, M. (2010). The contribution of real-time
818 mirror reflections of motor actions on virtual body ownership in an immersive virtual environment. In
819 *2010 IEEE Virtual Reality Conference (VR)*. 111–114. doi:10.1109/VR.2010.5444805
- 820 Guo, C., Jiang, T., Chen, X., Song, J., and Hilliges, O. (2023). Vid2Avatar: 3D Avatar Reconstruction
821 From Videos in the Wild via Self-Supervised Scene Decomposition. In *Proceedings of the IEEE/CVF*
822 *Conference on Computer Vision and Pattern Recognition (CVPR)*. 12858–12868
- 823 Habermann, M., Liu, L., Xu, W., Pons-Moll, G., Zollhoefer, M., and Theobalt, C. (2023). HDHumans:
824 A Hybrid Approach for High-fidelity Digital Humans. *Proc. ACM Comput. Graph. Interact. Tech.* 6.
825 doi:10.1145/3606927
- 826 Ho, C.-C. and MacDorman, K. F. (2017). Measuring the uncanny valley effect. *International Journal of*
827 *Social Robotics* 9, 129–139. doi:10.1007/s12369-016-0380-9

- 828 Hu, L., Zhang, H., Zhang, Y., Zhou, B., Liu, B., Zhang, S., et al. (2024). GaussianAvatar: Towards Realistic
829 Human Avatar Modeling from a Single Video via Animatable 3D Gaussians. In *Proceedings of the*
830 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 634–644
- 831 [Dataset] IBM (2022). SPSS Statistics. <https://www.ibm.com/products/spss-statistics>
- 832 Jiang, B., Hong, Y., Bao, H., and Zhang, J. (2022). SelfRecon: Self Reconstruction Your Digital Avatar
833 from Monocular Video. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*
- 834 Jiang, T., Chen, X., Song, J., and Hilliges, O. (2023). InstantAvatar: Learning Avatars From Monocular
835 Video in 60 Seconds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
836 *Recognition (CVPR)*. 16922–16932
- 837 Jiang, Z., Guo, C., Kaufmann, M., Jiang, T., Valentin, J., Hilliges, O., et al. (2024). Multiply:
838 Reconstruction of multiple people from monocular video in the wild. In *Proceedings of the IEEE/CVF*
839 *Conference on Computer Vision and Pattern Recognition (CVPR)*. 109–118
- 840 Jung, S. and Hughes, C. E. (2016). The effects of indirect real body cues of irrelevant parts on virtual body
841 ownership and presence. In *Proceedings of the 26th International Conference on Artificial Reality and*
842 *Telexistence and the 21st Eurographics Symposium on Virtual Environments*. 107–114
- 843 Kerbl, B., Kopanas, G., Leimkühler, T., and Drettakis, G. (2023). 3d gaussian splatting for real-time
844 radiance field rendering. *ACM Transactions on Graphics* 42
- 845 Kilteni, K., Groten, R., and Slater, M. (2012). The sense of embodiment in virtual reality. *Presence: Teleoperators & Virtual Environments* 21, 373–387. doi:10.1162/PRES.a.00124
- 847 Kim, D. Y., Lee, H. K., and Chung, K. (2023). Avatar-mediated experience in the metaverse: The
848 impact of avatar realism on user-avatar relationship. *Journal of Retailing and Consumer Services* 73.
849 doi:<https://doi.org/10.1016/j.jretconser.2023.103382>
- 850 Kim, H. K., Park, J., Choi, Y., and Choe, M. (2018). Virtual reality sickness questionnaire (VRSQ):
851 Motion sickness measurement index in a virtual reality environment. *Applied Ergonomics* 69, 66–73.
852 doi:<https://doi.org/10.1016/j.apergo.2017.12.016>
- 853 Kortum, P. and Sorber, M. (2015). Measuring the Usability of Mobile Applications for Phones and Tablets.
854 *International Journal of Human–Computer Interaction* 31, 518–529. doi:10.1080/10447318.2015.
855 1064658
- 856 Kwon, Y., Liu, L., Fuchs, H., Habermann, M., and Theobalt, C. (2023). DELIFFAS: Deformable Light
857 Fields for Fast Avatar Synthesis. In *Advances in Neural Information Processing Systems*. vol. 36,
858 40944–40962
- 859 Latoschik, M. E., Kern, F., Stauffert, J.-P., Bartl, A., Botsch, M., and Lugrin, J.-L. (2019). Not alone here?!
860 scalability and user experience of embodied ambient crowds in distributed social virtual reality. *IEEE*
861 *Transactions on Visualization and Computer Graphics (TVCG)* 25, 2134–2144
- 862 Latoschik, M. E. and Wienrich, C. (2022). Congruence and plausibility, not presence: Pivotal conditions
863 for XR experiences and effects, a novel approach. *Frontiers in Virtual Reality* 3. doi:10.3389/frvir.2022.
864 694433
- 865 Lei, J., Wang, Y., Pavlakos, G., Liu, L., and Daniilidis, K. (2024). GART: Gaussian Articulated Template
866 Models. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 19876–
867 19887. doi:10.1109/CVPR52733.2024.01879
- 868 Leyrer, M., Linkenauger, S. A., Bühlhoff, H. H., Kloos, U., and Mohler, B. (2011). The influence of eye
869 height and avatars on egocentric distance estimates in immersive virtual environments. In *Proceedings*
870 *of the ACM SIGGRAPH Symposium on Applied Perception in Graphics and Visualization*. 67–74.
871 doi:10.1145/2077451.2077464

- 872 Li, Z., Zheng, Z., Wang, L., and Liu, Y. (2024). Animatable Gaussians: Learning Pose-dependent Gaussian
873 Maps for High-fidelity Human Avatar Modeling. In *Proceedings of the IEEE/CVF Conference on*
874 *Computer Vision and Pattern Recognition (CVPR)*. 19711–19722
- 875 [Dataset] Limesurvey GmbH (2024). LimeSurvey: An Open Source survey tool
- 876 Lin, W., Zheng, C., Yong, J.-H., and Xu, F. (2024). Relightable and Animatable Neural Avatars from
877 Videos. *Proceedings of the AAAI Conference on Artificial Intelligence* 38, 3486–3494. doi:10.1609/aaai.
878 v38i4.28136
- 879 Liu, L., Habermann, M., Rudnev, V., Sarkar, K., Gu, J., and Theobalt, C. (2021). Neural actor: neural
880 free-view synthesis of human actors with pose control. *ACM Trans. Graph.* 40. doi:10.1145/3478513.
881 3480528
- 882 [Dataset] Lugaresi, C., Tang, J., Nash, H., McClanahan, C., Uboweja, E., Hays, M., et al. (2019).
883 MediaPipe: A Framework for Building Perception Pipelines
- 884 Ma, S., Simon, T., Saragih, J., Wang, D., Li, Y., Torre, F. L., et al. (2021). Pixel Codec Avatars. In *2021*
885 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 64–73. doi:10.1109/
886 CVPR46437.2021.00013
- 887 Mal, D., Wolf, E., Döllinger, N., Botsch, M., Wienrich, C., and Latoschik, M. E. (2022). Virtual Human
888 Coherence and Plausibility – Towards a Validated Scale. In *IEEE Conference on Virtual Reality and 3D*
889 *User Interfaces Abstracts and Workshops (VRW)*. 788–789. doi:10.1109/VRW55335.2022.00245
- 890 Mal, D., Wolf, E., Döllinger, N., Botsch, M., Wienrich, C., and Latoschik, M. E. (2024). From 2D-screens
891 to VR: Exploring the effect of immersion on the plausibility of virtual humans. In *CHI 24 Conference*
892 *on Human Factors in Computing Systems Extended Abstracts*. 8
- 893 Mildenhall, B., Srinivasan, P. P., Tancik, M., Barron, J. T., Ramamoorthi, R., and Ng, R. (2020). NeRF:
894 Representing scenes as neural radiance fields for view synthesis. In *ECCV*
- 895 Mohler, B. J., Creem-Regehr, S. H., Thompson, W. B., and Bühlhoff, H. H. (2010). The effect of viewing
896 a self-avatar on distance judgments in an HMD-based virtual environment. *Presence* 19, 230–242.
897 doi:10.1162/pres.19.3.230
- 898 Moreau, A., Song, J., Dharmo, H., Shaw, R., Zhou, Y., and Pérez-Pellitero, E. (2024). Human Gaussian
899 Splatting: Real-time Rendering of Animatable Avatars. In *Proceedings of the IEEE/CVF Conference on*
900 *Computer Vision and Pattern Recognition (CVPR)*. 788–798
- 901 Morgenstern, W., Bagdasarian, M. T., Hilsmann, A., and Eisert, P. (2024). Animatable virtual humans:
902 Learning pose-dependent human representations in uv space for interactive performance synthesis.
903 *IEEE Transactions on Visualization and Computer Graphics* 30, 2644–2650. doi:10.1109/TVCG.2024.
904 3372117
- 905 Mori, M., MacDorman, K. F., and Kageki, N. (2012). The uncanny valley [from the field]. *IEEE Robotics*
906 *Automation Magazine* 19, 98–100. doi:10.1109/MRA.2012.2192811
- 907 Mottelson, A., Muresan, A., Hornbæk, K., and Makransky, G. (2023). A systematic review and meta-
908 analysis of the effectiveness of body ownership illusions in virtual reality. *ACM Trans. Comput.-Hum.*
909 *Interact.* doi:10.1145/3590767
- 910 Müller, T., Evans, A., Schied, C., and Keller, A. (2022). Instant Neural Graphics Primitives with a
911 Multiresolution Hash Encoding. *ACM Trans. Graph.* 41. doi:10.1145/3528223.3530127
- 912 Mystakidis, S. (2022). Metaverse. *Encyclopedia* 2, 486–497. doi:10.3390/encyclopedia2010031
- 913 Pang, H., Zhu, H., Kortylewski, A., Theobalt, C., and Habermann, M. (2024). ASH: Animatable Gaussian
914 Splats for Efficient and Photoreal Human Rendering. In *2024 IEEE/CVF Conference on Computer*
915 *Vision and Pattern Recognition (CVPR)*. 1165–1175. doi:10.1109/CVPR52733.2024.00117

- 916 Pastel, S., Chen, C.-H., Petri, K., and Witte, K. (2020). Effects of body visualization on performance in
917 head-mounted display virtual reality. *PLOS ONE* 15, 1–18. doi:10.1371/journal.pone.0239226
- 918 Peng, S., Dong, J., Wang, Q., Zhang, S., Shuai, Q., Zhou, X., et al. (2021a). Animatable neural radiance
919 fields for modeling dynamic human bodies. In *2021 IEEE/CVF International Conference on Computer
920 Vision (ICCV)*. 14294–14303. doi:10.1109/ICCV48922.2021.01405
- 921 Peng, S., Zhang, Y., Xu, Y., Wang, Q., Shuai, Q., Bao, H., et al. (2021b). Neural Body: Implicit
922 Neural Representations with Structured Latent Codes for Novel View Synthesis of Dynamic Humans
923 . In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 9050–9059.
924 doi:10.1109/CVPR46437.2021.00894
- 925 Pérez, P., Gangnet, M., and Blake, A. (2003). Poisson image editing. *ACM Transactions on Graphics* 22,
926 313–318
- 927 Robinette, K., Daanen, H., and Paquet, E. (1999). The CAESAR project: a 3-D surface anthropometry
928 survey. In *Second International Conference on 3-D Digital Imaging and Modeling (Cat. No.PR00062)*.
929 380–386. doi:10.1109/IM.1999.805368
- 930 Roth, D. and Latoschik, M. E. (2020). Construction of the virtual embodiment questionnaire (VEQ).
931 *IEEE Transactions on Visualization and Computer Graphics* 26, 3546–3556. doi:10.1109/TVCG.2020.
932 3023603
- 933 Salagean, A., Crellin, E., Parsons, M., Cosker, D., and Stanton Fraser, D. (2023). Meeting Your Virtual
934 Twin: Effects of Photorealism and Personalization on Embodiment, Self-Identification and Perception of
935 Self-Avatars in Virtual Reality. In *Proceedings of the CHI Conference on Human Factors in Computing
936 Systems*. doi:10.1145/3544548.3581182
- 937 Sampaio, M., Navarro Haro, M. V., De Sousa, B., Vieira Melo, W., and Hoffman, H. G. (2021). Therapists
938 Make the Switch to Telepsychology to Safely Continue Treating Their Patients During the COVID-19
939 Pandemic. Virtual Reality Telepsychology May Be Next. *Frontiers in Virtual Reality* 1. doi:10.3389/
940 frvir.2020.576421
- 941 Sauro, J. and Lewis, J. R. (2016). *Quantifying the user experience: Practical statistics for user research*
942 (Morgan Kaufmann)
- 943 Schönberger, J. L. and Frahm, J.-M. (2016). Structure-from-Motion Revisited. In *Conference on Computer
944 Vision and Pattern Recognition (CVPR)*
- 945 Shao, Z., Wang, Z., Li, Z., Wang, D., Lin, X., Zhang, Y., et al. (2024). SplattingAvatar: Realistic Real-Time
946 Human Avatars with Mesh-Embedded Gaussian Splatting. In *Proceedings of the IEEE/CVF Conference
947 on Computer Vision and Pattern Recognition (CVPR)*. 1606–1616
- 948 Shetty, A., Habermann, M., Sun, G., Luvizon, D., Golyanik, V., and Theobalt, C. (2024). Holoported
949 Characters: Real-time Free-viewpoint Rendering of Humans from Sparse RGB Cameras. In *Proceedings
950 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 1206–1215
- 951 Simon, H. A. and Ericsson, K. A. (1993). *Protocol analysis: Verbal reports as data* (The MIT Press)
- 952 Skarbez, R., Frederick P. Brooks, J., and Whitton, M. C. (2017). A survey of presence and related concepts.
953 *ACM Computing Surveys* 50, 96. doi:10.1145/3134301
- 954 Skarbez, R. and Jiang, D. (2024). A Scientometric History of IEEE VR. In *2024 IEEE Conference Virtual
955 Reality and 3D User Interfaces (VR)*. 990–999. doi:10.1109/VR58804.2024.00118
- 956 Slater, M., Spanlang, B., Sanchez-Vives, M. V., and Blanke, O. (2010). First person experience of body
957 transfer in virtual reality. *PLOS ONE* 5, e10564. doi:10.1371/journal.pone.0010564
- 958 Steed, A., Pan, Y., Zisch, F., and Steptoe, W. (2016). The impact of a self-avatar on cognitive load in
959 immersive virtual reality. In *2016 IEEE Virtual Reality (VR)*. 67–76. doi:10.1109/VR.2016.7504689

- 960 Sutherland, I. E. (1968). A head-mounted three-dimensional display. In *Proceedings of the December*
961 *9-11, 1968, Fall Joint Computer Conference, Part I*. 757–764. doi:10.1145/1476589.1476686
- 962 Turbyne, C., Goedhart, A., de Koning, P., Schirmbeck, F., and Denys, D. (2021). Systematic Review and
963 Meta-Analysis of Virtual Reality in Mental Healthcare: Effects of Full Body Illusions on Body Image
964 Disturbance. *Frontiers in Virtual Reality* 2, 39. doi:10.3389/frvir.2021.657638
- 965 [Dataset] Unity Technologies (2020). Unity 2020.3.25f1
- 966 [Dataset] Valve Corporation (2024a). Steam VR 2.3
- 967 [Dataset] Valve Corporation (2024b). Steam VR Plugin 2.7.3
- 968 Waltemate, T., Gall, D., Roth, D., Botsch, M., and Latoschik, M. E. (2018). The Impact of Avatar
969 Personalization and Immersion on Virtual Body Ownership, Presence, and Emotional Response. *IEEE*
970 *Transactions on Visualization and Computer Graphics* 24, 1643–1652. doi:10.1109/TVCG.2018.
971 2794629
- 972 Wang, S., Antic, B., Geiger, A., and Tang, S. (2024). IntrinsicAvatar: Physically Based Inverse Rendering
973 of Dynamic Humans from Monocular Videos via Explicit Ray Tracing. In *Proceedings of the IEEE/CVF*
974 *Conference on Computer Vision and Pattern Recognition (CVPR)*. 1877–1888
- 975 Wang, T., Zhang, B., Zhang, T., Gu, S., Bao, J., Baltrusaitis, T., et al. (2023). RODIN: A Generative Model
976 for Sculpting 3D Digital Avatars Using Diffusion. In *2023 IEEE/CVF Conference on Computer Vision*
977 *and Pattern Recognition (CVPR)*. 4563–4573. doi:10.1109/CVPR52729.2023.00443
- 978 Wenninger, S., Achenbach, J., Bartl, A., Latoschik, M. E., and Botsch, M. (2020). Realistic Virtual
979 Humans from Smartphone Videos. In *Proceedings of the ACM Symposium on Virtual Reality Software*
980 *and Technology*. doi:10.1145/3385956.3418940
- 981 Wiegand, T., Sullivan, G., Bjontegaard, G., and Luthra, A. (2003). Overview of the H.264/AVC video
982 coding standard. *IEEE Transactions on Circuits and Systems for Video Technology* 13, 560–576.
983 doi:10.1109/TCSVT.2003.815165
- 984 Wilcox, R. R. (2022). *Introduction to robust estimation and hypothesis testing* (Academic Press)
- 985 Wolf, E., Döllinger, N., Mal, D., Wenninger, S., Andrea, B., Botsch, M., et al. (2022a). Does Distance
986 Matter? Embodiment and Perception of Personalized Avatars in Relation to the Self-Observation Distance
987 in Virtual Reality. *Frontiers in Virtual Reality* 3. doi:10.3389/frvir.2022.1031093
- 988 Wolf, E., Fiedler, M. L., Döllinger, N., Wienrich, C., and Latoschik, M. E. (2022b). Exploring Presence,
989 Avatar Embodiment, and Body Perception with a Holographic Augmented Reality Mirror. In *2022 IEEE*
990 *Conference on Virtual Reality and 3D User Interfaces (VR)*. 350–359. doi:10.1109/VR51125.2022.00054
- 991 Wolf, E., Merdan, N., Döllinger, N., Mal, D., Wienrich, C., Botsch, M., et al. (2021). The embodiment
992 of photorealistic avatars influences female body weight perception in virtual reality. In *IEEE Virtual*
993 *Reality and 3D User Interfaces (VR)*. 65–74. doi:10.1109/VR50410.2021.00027
- 994 Xiao, J., Zhang, Q., Xu, Z., and Zheng, W.-S. (2024). NECA: Neural Customizable Human Avatar.
995 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
996 20091–20101
- 997 Yoon, B., Kim, H.-i., Lee, G. A., Billinghamurst, M., and Woo, W. (2019). The Effect of Avatar Appearance
998 on Social Presence in an Augmented Reality Remote Collaboration. In *IEEE Conference on Virtual*
999 *Reality and 3D User Interfaces (VR)*. 547–556. doi:10.1109/VR.2019.8797719
- 1000 Yu, W., Fan, Y., Zhang, Y., Wang, X., Yin, F., Bai, Y., et al. (2023). NOFA: NeRF-Based One-Shot
1001 Facial Avatar Reconstruction. In *ACM SIGGRAPH 2023 Conference Proceedings*. doi:10.1145/3588432.
1002 3591555

- 1003 Zhao, H., Zhang, J., Lai, Y.-K., Zheng, Z., Xie, Y., Liu, Y., et al. (2022). High-fidelity human avatars
1004 from a single rgb camera. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
1005 (*CVPR*). 15883–15892. doi:10.1109/CVPR52688.2022.01544
- 1006 Zheng, Z., Huang, H., Yu, T., Zhang, H., Guo, Y., and Liu, Y. (2022). Structured Local Radiance Fields for
1007 Human Avatar Modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
1008 *Recognition (CVPR)*. 15893–15903
- 1009 Zheng, Z., Zhao, X., Zhang, H., Liu, B., and Liu, Y. (2023). AvatarReX: Real-time Expressive Full-body
1010 Avatars. *ACM Trans. Graph.* 42. doi:10.1145/3592101