

Northumbria Research Link

Citation: Chan, Jacky C. P. and Ho, Edmond (2021) Emotion Transfer for 3D Hand and Full Body Motion using StarGAN. Computers, 10 (3). p. 38. ISSN 2073-431X

Published by: MDPI

URL: <https://doi.org/10.3390/computers10030038>
<<https://doi.org/10.3390/computers10030038>>

This version was downloaded from Northumbria Research Link:
<http://nrl.northumbria.ac.uk/id/eprint/45755/>

Northumbria University has developed Northumbria Research Link (NRL) to enable users to access the University's research output. Copyright © and moral rights for items on NRL are retained by the individual author(s) and/or other copyright owners. Single copies of full items can be reproduced, displayed or performed, and given to third parties in any format or medium for personal research or study, educational, or not-for-profit purposes without prior permission or charge, provided the authors, title and full bibliographic details are given, as well as a hyperlink and/or URL to the original metadata page. The content must not be changed in any way. Full items must not be sold commercially in any format or medium without formal permission of the copyright holder. The full policy is available online: <http://nrl.northumbria.ac.uk/policies.html>

This document may differ from the final, published version of the research and has been made available online in accordance with publisher policies. To read and/or cite from the published version of the research, please visit the publisher's website (a subscription may be required.)



**Northumbria
University**
NEWCASTLE



UniversityLibrary

Emotion Transfer for 3D Hand and Full Body Motion using StarGAN

Jacky C. P. Chan ¹ and Edmond S. L. Ho ^{2*} 

¹ Department of Computer Science, Hong Kong Baptist University, Hong Kong; cpchan@comp.hkbu.edu.hk

² Department of Computer and Information Sciences, Northumbria University, United Kingdom; e.ho@northumbria.ac.uk

* Correspondence: e.ho@northumbria.ac.uk;

Abstract: In this paper, we propose a new data-driven framework for 3D hand and full-body motion emotion transfer. Specifically, we formulate the motion synthesis task as an image-to-image translation problem. By presenting a motion sequence as an image representation, the emotion can be transferred by our framework using StarGAN. To evaluate our proposed method's effectiveness, we first conducted a user study to validate the perceived emotion from the captured and synthesized hand motions. We further evaluate the synthesized hand and full body motions qualitatively and quantitatively. Experimental results show that our synthesized motions are comparable to the captured motions and those created by an existing method in terms of naturalness and visual quality.

Keywords: hand animation; body motion; skeletal motion; emotion; motion capture; generative adversarial network; style transfer; user study

1. Introduction

Effectively expressing emotion is crucial to improve the realism of 3D character animation. While animating facial expressions to reflect the character's emotional states has been an active research area [1–4], less attention has been paid to expressing emotion by other body parts, practically 'the body language'. In this work, we propose a general framework for synthesizing new hand and full body motions from an input motion, namely *emotion transfer*, by specifying the target emotion label. Our objective is to create motions for the character to present four emotions: anger, sadness, fear, and joy. Building upon our pilot study [5], we found that hand motion plays a vital role in computer animation since the subtle hand gestures can express a lot of different meanings and are useful for understanding a person's personality [6]. A classic example would be the character *Thing T. Thing* of the "*The Addams Family*" which is a hand, and it can 'act' and express a lot of different emotions solely by the fingers and hand movements. It is not surprising to see researchers proposing frameworks [7,8] for synthesizing hand and finger movements based on the given full-body motion to improve the expressiveness of the animation.

However, synthesizing hand motion is not a trivial task. Capturing hand motion using an optical motion capture system is not easy as the fingers are in proximity, and the labeling of the markers can be mixed up easily. As a result, most of the previous hand motion synthesis frameworks are based on physics-based motion generation models [8–12]. Recently, more effective hand motion capturing approaches are proposed. Alexanderson et al. introduce a new system for a passive optical motion capture system that can better obtain correct markers labels of fingers in real-time [13]. Han et al. [14] improve the difficulties in marker labeling for the optical MOCAP system using convolutional neural networks. While hand motion can be synthesized or captured using the approaches mentioned above, those motions are always challenging to be

Citation: Chan, J. C. P.; Ho, E. S. L. Emotion Transfer for 3D Hand and Full Body Motion using StarGAN. *Computers* **2021**, *1*, 0. <https://doi.org/>

Received:

Accepted:

Published:

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Copyright: © 2021 by the authors. Submitted to *Computers* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

38 reused because of the difficulties in transferring the styles to improve the expressiveness
39 in different scenes.

40 This paper focuses on validating the effectiveness of the StarGAN-based emotion
41 transfer framework proposed in our pilot study [5], which consist of three main compo-
42 nents as illustrated in Figure 1: 1) converting motion data into image representation, 2)
43 synthesizing new image by specifying the target emotion label using StarGAN, and 3)
44 converting the synthesized image into motion data for animating 3D mesh models. In
45 particular, we conducted a user study on evaluating how users perceived the emotion
46 from the hand motions captured in [5] and those synthesized by our method. The natu-
47 ralness and visual quality of the motions synthesized by our method are also evaluated
48 and compared with exiting work [15]. We further demonstrate the generality of the
49 proposed framework by transferring emotion on full-body 3D skeletal motions.

50 The contributions in this work can be summarized as follows:

- 51 • We proposed a new framework for transferring emotions in synthesizing hand and
52 full body skeletal motions, which is built upon the success of our pilot study [5].
- 53 • We conducted a user study to validate the perceived emotion on the dataset we
54 captured and open-sourced in our pilot study [5].
- 55 • We provide qualitative and quantitative results on the hand and full-body motion
56 emotion transfer using the proposed framework, showing its validity by comparing
57 them to captured motions.

58 2. Related Work

59 2.1. Hand animation

60 Examples of hand animation can be easily found in various applications such as
61 movies, games, and animations. However, capturing hand motion using existing motion
62 capture systems is not a trivial task. There was no commercial or academic real-time
63 vision-based hand motion capture solution until a recent work presented by Han et
64 al. [14]. As a result, most of the previous work focused on synthesizing hand motions
65 based on physics-based models. Liu [9] proposed an optimization-based approach for
66 synthesizing hand-object manipulations. Given the initial hand pose for the desired
67 initial contact, the properties of the object to interact and the kinematics goals, the 2-stage
68 physics-based framework will synthesize the reaching and object manipulation motions
69 accordingly. Andrews and Kry [10] proposed a hand motion synthesis framework for
70 object manipulation. The method divides a manipulation task into 3 phases (approach,
71 actuate, and release), and each phase is associated with a control policy for generating
72 physics-based hand motion.

73 Liu et al. [11] introduced an optimization-based approach to hand manipulation of
74 grasping pose. A physically plausible hand animation will be created by providing the
75 grasping pose and the partial trajectory of the object. Ye and Liu [8] proposed a physics-
76 based hand motion synthesis framework to generate detailed hand-object manipulations
77 which match seamlessly with the full-body motion with wrist movements provided as
78 the input. With the initial hand pose driven by the wrist movements, feasible contact
79 trajectories (i.e. contacts between the hand/fingers and the object) will be found by
80 random sampling at every time-step. Finally, detailed hand motion will be computed by
81 using the contact trajectories as constraints in spacetime optimization. Bai and Liu [12]
82 presented a solution to manipulate the orientation of a polygonal object using both the
83 palm and fingers of a robotic hand. Their method considers the physical properties such
84 as collisions, gravitational, and contact forces.

85 A PCA-based framework is proposed for data-driven approaches in [16] for gener-
86 ating detailed hand animation from a set of sparse markers. Jörg et al [7] proposed a
87 data-driven framework for synthesizing the finger movements for an input full-body
88 motion. The methods employ a simple approach for searching for an appropriate finger
89 motion from the pre-recorded motion database based on the input wrist and body mo-
90 tion. While a wide range of approaches for hand motion synthesis are presented in the

91 literature, less attention has been paid to synthesizing expressive hand motions using
92 high-level and intuitive control. Irimia et al. [15] proposed a framework for generating
93 hand motion with different emotion by interpolation in the latent space. For every hand
94 motion captured with different emotions, the hand poses are collected and projected to
95 the latent space using PCA. New motion can be created by interpolating the hand poses
96 using the latent representation. In contrast, our framework enables emotion transfer
97 between different hand motions.

98 2.2. *Style transfer for motion*

99 Motion style transfer is a technique used to convert the style of a motion to another
100 style, thus creating new motions without losing the primitive content of the original
101 one. An early work by Unuma et al. [17] proposed using Fourier principles to create an
102 interactive and real-time control of locomotion with emotion, and include cartoon-ish
103 exaggerations and expressions. Amaya et al. [18] introduced a model that could emulate
104 emotional animation using signal processing techniques. The emotional transform is
105 based on the speed and spatial amplitude of the movements.

106 Brand et al. [19] proposed a learned statistical model to synthesize styled motion
107 sequences by interpolating and extrapolating the motion data. Urtasun et al. [20]
108 lowered the dimension of the motion data by principal component analysis (PCA) and
109 model the style transfer as the difference between the features. Ikemoto et al. [21] edited
110 the motion using Gaussian process models of dynamics and kinematics for motion style
111 transfer. Xia et al. [22] proposed a time-varying mixture of autoregressive models to
112 represent the style difference between two motions. Their method learns such models
113 automatically from unlabeled heterogeneous motion data.

114 Hsu et al. [23] presented a solution for translating the style of a human motion by
115 comparing the difference of the behaviour of the aligned input and output motions using
116 a linear time-invariant model. Shapiro et al. introduced a novel method of interactive
117 motion data editing based on motion decomposition, which separates the style and
118 expressiveness from the main motion [24]. The method uses Independent Component
119 Analysis (ICA) to separate the style from the motion data.

120 Machine learning is applied to learn the style transfer from samples. Holden et
121 al. leveraged a Convolution neural network to learn the style transfer from unaligned
122 motion clips [25]. Smith et al. proposed a compact network architecture for learning the
123 style transfer, which focuses on pose, foot contact, and timing [26]. To learn a motion
124 controller with behavior styles applicable to unseen environment, Lee and Popović
125 proposed an inverse reinforcement learning-based approach that works with a small set
126 of motion samples [27].

127 Until now, the only research of style transfer was for a full-body character. This
128 paper will propose a general method to full body character and the human hand. While
129 we share a similar interest with the pilot study [15] on synthesizing hand motion with
130 emotion, the previous work is technically interpolating emotion strength instead of
131 emotion transfer.

132 2.2.1. *Image style transfer*

133 Inspired by the encouraging results in image style transfer, we proposed formulat-
134 ing the emotion transfer for motion synthesis as an image-to-image translation problem.
135 In this section, we review the recently proposed approaches in image style transfer.

136 Selim et al. [28] presented a new technique of style transfer that uses Convolutional
137 Neural Networks (CNN) for extracting features from the input images. The method
138 uses style transfer to transfer the features from a portrait painting to a portrait image.
139 The method is generic and different kinds of styles can be transferred given the training
140 data contains the required styles. To maintain the integrity of the facial structure and
141 to capture the texture from the painting, the method uses spatial constraints such as
142 the local transfer of the colour distribution. Elad et al. [29] presented a method for

143 transferring the style from painting to image indifferently of the portrayed subject. The
 144 method uses a fast style transfer process that gradually changes the resolution of the
 145 output image. To obtain the result, the creators applied multi-patch sizes and different
 146 resolution scales of the input source. The method is also able to control the colour
 147 pallet for the output image depending on the desire of the developer. Matsuo et al.
 148 [30] presented another style transfer method that uses CNN by combining a neural
 149 style transfer method with segmentation to obtain a partial texture style transfer. The
 150 method uses a CNN-based weakly supervised semantic segmentation technique and
 151 transfers the style to selected areas of the picture while maintaining the image's structure.
 152 The method uses neural style transfer to change the style of the selected part of the
 153 image. Unfortunately, a problem appears when the sources fail to map the style transfer,
 154 changing the background of the image even when the user does not select that area. In
 155 this work, we will represent the motion features as an image and CNN will be used in
 156 the core network.

157 3. Methodology

158 In this section, the proposed emotion transfer framework will be presented, and
 159 the overview is shown in Figure 1. Firstly, we introduce two datasets, hand motion
 160 database captured using Senso VR Glove and a full-body motion database captured using
 161 the MOCAP system (Section 3.1). These two databases contain motions with various
 162 emotional states and types. Next, the captured motions are standardized (Section
 163 3.2) as a pre-processing step for the learning process. The motion data will then be
 164 transformed into an RGB image representation for learning the emotion transfer model
 165 using StarGAN. The StarGAN model learns how to generate a new image given a target
 166 domain label and the input image (Section 3.3). Finally, the synthesized new image will
 167 be converted to the joint angle/position space for generating the final 3D animation
 168 (Section 3.3.3). The details of each step will be explained in the following subsections.

169 3.1. Motion Datasets

170 To learn how the motion features are mapped to emotion status, motion data is
 171 collected for training the models to be proposed in this paper. In particular, we used
 172 two datasets, which include the hand motion dataset collected in our pilot study [5] for
 173 hand motion synthesis, and the 3D skeletal motions in Body Movement Library [31] for
 174 full-body motion synthesis.

175 3.1.1. Hand Motion with Emotion Label

176 We start with the details of hand motion dataset which was captured in our pilot
 177 study [5]. High-quality 3D skeletal hand motions were captured using the Senso VR
 178 glove (<https://senso.me/>). There are 35 motions in total, with 7 different action types,
 179 including *Crawling*, *Gripping*, *Patting*, *Impatient*, *Hand on Mouse*, *Pointing*, and *Pushing*.
 180 Each motion types are captured under 5 different types of emotions and their character-
 181 istics are listed on Table 1. Readers are referred to [5] for the details of the data capturing
 182 process.

In this dataset, each hand motion at each frame is represented by a vector P_j

$$P_j = [p_{j,x}^0, p_{j,y}^0, p_{j,z}^0, \dots, p_{j,x}^{n-1}, p_{j,y}^{n-1}, p_{j,z}^{n-1}] \quad (1)$$

183 where j is the index of the frame, n is the joint number in the 3D hand skeletal structure
 184 and $n = 27$ in all of the data we capture, and p contains the joint rotations on the x,y and
 185 z axes, respectively. Therefore, each keyframe is a 81-dimensional feature vector. The
 186 hand translation was discarded as in [5,15] due to the inconsistent global locations of the
 187 hand in the captured motions.

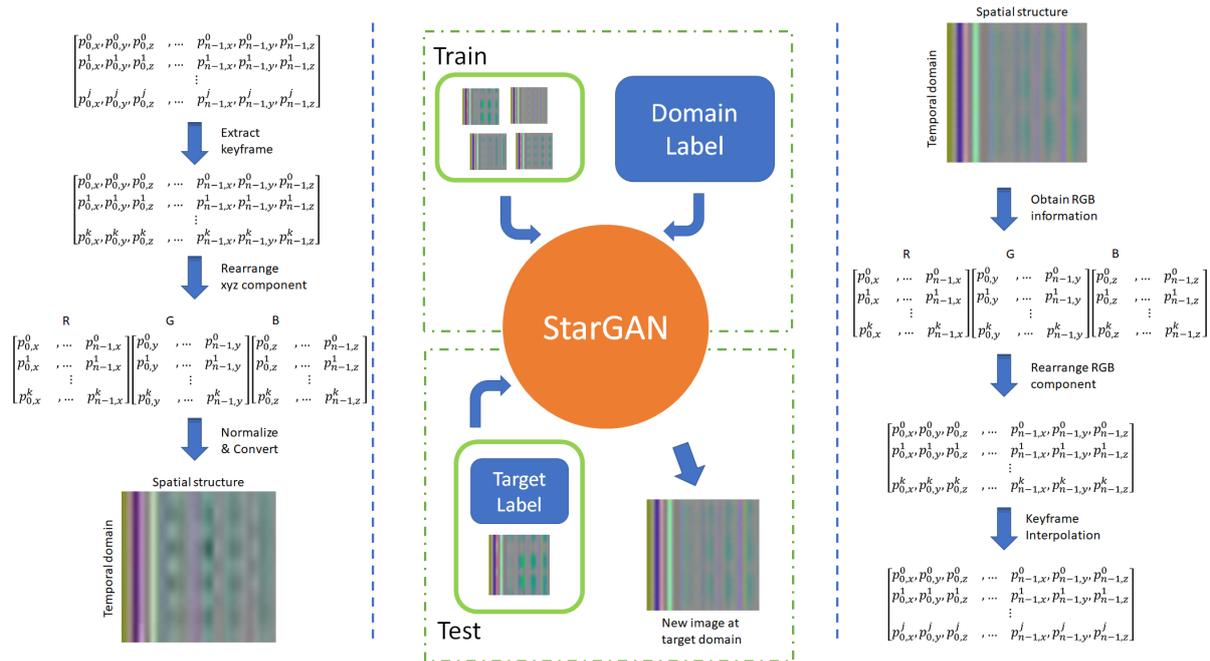


Figure 1. The overview of the proposed emotion transfer framework. (left) Convert motion data to image format. (middle) StarGAN learn how to generate realistic fake image given a sample and target domain label. (right) Obtain new motion data from the generated image

Emotion	Characteristics
Angry	exaggerated, fast, large range of motion
Happy	energetic, large range of motion
Neutral	normal, styleless motion
Sad	sign of tiredness, small range of motion
Fearful	asynchronous finger movements, small range of motion

Table 1: The 5 types of emotions used in the hand motion dataset [5] and their characteristics.

3.1.2. Full Body Motion with Emotion Label

Here, the details of the full-body motion dataset are presented. The Body Movement Library [31] were captured with the Falcon Analog optical motion capture system. In order to capture the emotion expressions naturally from the subjects, 30 nonprofessional participants (15 females and 15 males) with an age range from 17 to 29 years old were recruited. Three motion types, *knocking*, *lifting*, and *throwing*, are included in our experiment. A skeletal structure with 33 joints was used in all of the captured motions. Note that only the 3D joint positions in Cartesian coordinates are available. There are 3 motion types in the dataset, including *knocking*, *lifting* and *throwing*. Each subject performed each motion type with four different emotion status: *Neutral*, *Angry*, *Happy*, and *Sad*. The dataset contains 4,080 motions in total.

3.2. Standardizing Motion Feature

Due to the environmental setting and personal style, the captured hand motion data varies significantly representation both spatially and temporally. Data standardization (or normalization) is used to facilitate the learning process in the later stage. While some advanced techniques such as Recurrent Neural Network (RNN) can be used to model

204 data sequences with variations in length, such a method requires a significant amount of
 205 data to train the model, which is not feasible with the dataset that have been collected.

206 To handle the temporal difference, keyframes are extracted from the motion by
 207 curve simplification to facilitate the learning process. By considering the reconstruction
 208 errors when interpolating the in-between motion using spline interpolation, we found
 209 the optimal numbers of keyframes for every type of motion. For hand motion, good
 210 performance can be achieved by extracting 9 keyframes as in [5].

As a result, each hand motion sequence is represented by a vector M in the joint angle space:

$$M = [Pk_0, \dots, Pk_{k-1}, Pk_k] \quad (2)$$

211 where k is the total number of keyframes and $k = 9$, Pk_i is the i -th keyframe and Pk is
 212 having the same representation as in P (Eq. 1).

For full-body motions, we empirically found that the optimal number of keyframes of knocking, throwing, lifting are 13, 20, and 13, respectively. In [32], Chan et al. observed that people express different emotions by using different speeds and rhythms. Such an observation aligns well with the characteristics we found in hand motions as presented in Table 1. Specifically, there is a significant difference in the speed of body movements. For example, the arm of the subject swings faster in the *happy* throwing motion than the *sad* one (see Figure 8). To better represent this key characteristic, we compute the joint velocity between adjacent keyframes as follows:

$$v_{i+1} = (k_{i+1} - k_i) / \Delta t \quad (3)$$

213 where Δt is the duration between the two adjacent keyframes. Therefore, each keyframe
 214 is represented by 3D joint positions and velocities and results in a $99 + 99 = 198$ -
 215 dimensional feature vector. The full body motion sequence can then be represented as a
 216 sequence of keyframes as in Eq. 1.

217 3.3. Emotion transfer

218 Generative adversarial network (GAN) based framework gain much attention in the
 219 area of Computer Graphics. Encouraging results are found in style transfer frameworks
 220 such as CycleGAN [33] and DualGAN [34] for image-to-image translation. These results
 221 inspire us to adopt such kind of framework for emotion transfer for skeletal motion
 222 synthesis tasks.

223 In the rest of this section, we will first explain how to represent a motion in the
 224 format of an image in Section 3.3.1. Next, the justifications on adopting the StarGAN [35]
 225 framework will be given in Section 3.3.2. Finally, a new motion will be reconstructed
 226 from the synthesized image (Section 3.3.3).

227 3.3.1. Representing motion as an image

In order to use the Image-to-Image domain translation framework for motion emotion transfer, we will show how to represent a motion sequence as an image. The x , y and z components (i.e. angles for hand motions and positions for full-body motions) are arranged chronologically as the RGB components of the image. Each frame of a motion is represented as a row of an image, while each joint of a motion is represented as a column of an image. Hence, each keyframed hand motion M will be arranged as a 3-channel ($27 \times k$) matrix. On the other hand, each keyframed full-body motion M_{full} will be arranged as a 3-channel ($66 \times k$) matrix. The values in the 3-channel matrix are then re-scaled into the range of $[0, 255]$, which is the typical range of RGB value, as follows:

$$v_{i,c}^m = \text{round}\left(255 \times \frac{(p_{i,c}^m - p_{min})}{(p_{max} - p_{min})}\right) \quad (4)$$

228 where m is the joint index, i is the keyframe index, $c \in \{x, y, z\}$ represents the channel
 229 index, $V_{i,c}^m$ is the normalized pixel value, p_{max} and p_{min} are the maximum and minimum

230 values among all the joint angles/positions/velocities existed in the dataset. Noted that
 231 the images are saved in Bitmap format to avoid data loss during compression. Examples
 232 of the image representation of the motions are illustrated in Figure 2 and 3. It can be
 233 seen that the different motions are represented by different image patterns, which will
 234 be useful for extracting discriminative patterns in the learning process. From Figure 3,
 235 we can see that the main difference between the four emotions is at the right-hand side
 236 of the images, which is about joint velocity.

237 When converting the motion into an image, sequential ordering is used. Such an
 238 approach is commonly used for arranging the data of each joint [36]. In this study,
 239 our main focus is to evaluate the performance of adapting the StarGAN network for
 240 emotion transfer in motion synthesis tasks. As a result, we directly utilize the original
 241 StarGAN network architecture which contains 2D Convolutional layers. Since 2D
 242 Convolution focuses on local neighbours (i.e. image pixels nearby) only, using the
 243 sequential ordering method can result in sub-optimal results when representing the tree-
 244 like skeletal structure for full-body and hand motions as the neighbouring joints are not
 245 necessarily close-by after converting into the image representation. While encouraging
 246 results are obtained in this study, we will explore the use of other approaches to better
 247 represent motions in the StarGAN framework in the future, such as Graph Convolutional
 248 Networks (GCNs) [37] and its variants [38] which demonstrated a better performance in
 249 modelling human-like skeletal motions.

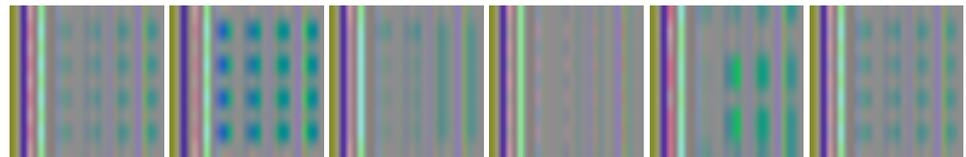


Figure 2. Examples of the image representation of neutral hand motions. From left to right: Crawling, Gripping, Impatient, Patting, Pointing and Pushing.



Figure 3. Examples of the image representation of full body motions (throwing) under different emotions. From left to right: Angry, Happy, Neutral and Sad.

250 3.3.2. Emotion transfer as Image-to-Image domain translation

251 One of the potential applications of the proposed system is to create new motion
 252 by controlling the *emotion labels*. To support the translation between multi-domain
 253 and considering the robustness and scalability, StarGAN [35] is adapted to translate
 254 motion from one emotion to another emotion while preserving the basic information
 255 of the input motion. Comparing to typical GANs with cycle consistency losses such as
 256 CycleGAN [33] and DualGAN [34] for style transfer, StarGAN [35] can perform image-
 257 to-image translations for multiple domains using only a single model, which is suitable
 258 for transferring different types of emotions. Readers are referred to our pilot study [5]
 259 and [35] for the technical details.

260 3.3.3. Reconstructing Hand Motion from Generated Images

Since the output of StarGAN is an image, we need to reconstruct it to obtain the
 new motion (i.e., joint position/angle space). The first step is to re-scale the RGB values:

$$p_{i,c}^m = \left(\frac{v_{i,c}^m}{255} \times (p_{max} - p_{min}) \right) + p_{min} \quad (5)$$

where $v_{i,j}^m$ is the pixel value for the c -th channel of the m -th row at i -th column on the synthesized image, p_{max} and p_{min} are same values as in Equation 4. Next, we rearrange the pixel values to convert the image representation back to the keyframed motion M . Then the duration between two adjacent keyframes can be approximated by using joint velocity:

$$\Delta t = (k_{i+1} - k_i) / v_{i+1} \quad (6)$$

261 Finally, the new motion is produced by applying spline interpolation on those keyframes.

262 4. Experimental Results

263 To evaluate the effectiveness of the proposed emotion transfer framework, a wide
 264 range of experiments are conducted to assess the performance qualitatively and quan-
 265 titatively. In particular, we first carried out a user study (Section 4.1) to understand
 266 how users perceive the emotions from the captured and synthesized hand motions.
 267 Next, a series of hand and full body animations are synthesized to compare different
 268 emotions visually (Section 4.2). To demonstrate the framework’s practicality, we employ
 269 a leave-one-out cross-validation approach to split the data into training, and testing
 270 sets. The results presented in the section are synthesized from unseen samples. We
 271 employed a leave-one-out cross-validation approach for hand motions as in the existing
 272 approaches [5,15]. For each action type, we only captured 1 sample for each emotion
 273 type in the hand motion dataset. As a result, the leave-one-out cross-validation is the
 274 best data-split approach we can use in order to 1) maximize the number of training
 275 data, while 2) keeping the testing sample to be ‘unseen’ during the training process. For
 276 full-body motions, the motions of 26 subjects are used for training, while the motions of
 277 the remaining 4 subjects are used for synthesizing the results. The readers are referred to
 278 the video demo submitted as supplementary materials for more results.

279 4.1. User Study on Hand Animations

280 A user study was conducted to evaluate how users perceived the emotion from each
 281 hand animation created by either the captured motions or synthesized by our proposed
 282 method. The study was published online using Google Form, and the animations are
 283 embedded in the online form to facilitate the side-by-side comparison. We invited a
 284 group of final year undergraduate students who are studying in the Computer Science
 285 programme. In particular, the students completed a course on Computer Graphics and
 286 Animation recently, although the course does not cover any specific knowledge on hand
 287 animation and emotion-based motion synthesis. The range of age in the group is 20-25.
 288 At the end of the study, 30 completed sets of the survey were received. The study is
 289 divided into several parts, including the emotion recognition from the captured dataset
 290 (Section 4.1.1) and synthesized animations (Section 4.1.2) for validating whether the
 291 participants perceived the emotion correctly, and the evaluation of naturalness and
 292 visual quality of the captured motions, synthesized motions and results created by Irimia
 293 et al. [15] (Section 4.1.3). The details are explained in the following subsections.

294 4.1.1. Evaluating the emotion perceived from the captured hand motions

295 To evaluate how users perceive emotion from hand animation, we first analyze
 296 the accuracy of emotion recognition from the captured hand motions. Specifically, we
 297 randomly select 4 hand animations from the capture motions, which include one motion
 298 in each emotion category (i.e. *Angry*, *Happy*, *Sad*, and *Afraid*). Note that no briefing, such
 299 as the characteristics for each emotion class as listed on Table 1, is provided to the users.
 300 As a result, the users labeled each hand animation based on their interpretation of the
 301 emotions. For each animation, the *neutral* emotion of the hand animation is shown to
 302 the user first. Next, the user was asked to choose the most suitable emotion label for
 303 the hand animation with emotion. The averaged and class-level emotion recognition
 304 accuracies are reported in Table 2 (see the ‘Captured’ column).

Emotion	Captured	Synthesized
Angry	70.00%	73.33%
Happy	60.00%	60.00%
Sad	70.00%	60.00%
Fearful	63.33%	70.00%
Average	65.83%	65.83%

Table 2: Emotion recognition accuracy (%) on the captured and synthesized hand motions in the user study.

		Perceived Emotion			
		Angry	Happy	Sad	Feared
Ground Truth Labels	Angry	70.00%	15.00%	3.33%	11.67%
	Happy	23.33%	60.00%	6.67%	10.00%
	Sad	6.67%	10.00%	70.00%	13.33%
	Feared	13.33%	10.00%	13.33%	63.33%

Table 3: Confusion matrix of the emotion recognition test on the captured hand motions in the user study.

305 The averaged recognition accuracy is 65.83%, which shows that the majority of the
 306 users perceived the correct emotion from the captured hand motions. The class-level
 307 recognition accuracies further show the consistency among all classes. In particular,
 308 a good recognition accuracy at 70.00% was achieved in both Angry and Sad classes.
 309 To further analyze the results, the confusion matrix of the emotion recognition test is
 310 presented in Table 3. It can be seen that the recognition accuracy of the Happy class is
 311 lower than other classes. This is mainly caused by the inter-class similarity between
 312 happy and angry since the motions in both classes are having a large range of motion.
 313 Although ambiguity can be found between motions from different classes, the results
 314 highlight the correct recognition is dominating the results.

315 4.1.2. Evaluating the emotion perceived from the hand motions synthesized by our
 316 method

317 Similar to Section 4.1.1, we also analyze the accuracy of emotion recognition from
 318 the hand motions synthesized using our method. Again, each user was asked to label 4
 319 synthesized hand animations. The averaged and class-level emotion recognition accura-
 320 cies are reported in Table 2 (see the ‘Synthesized’ column). The averaged recognition
 321 accuracy is 65.83%, which is the same as the accuracy obtained from the captured dataset.
 322 This result highlights no significant difference in the expressiveness of emotion between
 323 the captured and synthesized hand motions. While the averaged recognition accuracies
 324 are the same, it can be seen that the class-level accuracies are different, and the confusion
 325 matrix is presented in Table 4. The results indicate that the synthesized motions demon-
 326 strated a lower level of inter-class similarity. The perceived emotion only spread across 3
 327 classes for Angry, Sad and Feared (i.e. with one class having 0% of voting) instead of 4
 328 classes as in the results obtained from captured motions.

329 While the inter-class similarity is reduced in general, it can be observed that some
 330 pairs of classes are having higher similarity in the motions synthesized by our method.
 331 For example, the ‘Sad-Happy’ pair (i.e. misperceiving *Sad* as *Happy*) increases from
 332 10% to 21.67%. this can be caused by the small changes in the speed of the motions
 333 as the synthesized motions tend to have a slightly smaller range of speed difference
 334 across different emotions. This affects the *Sad* and *Happy* motions since the major
 335 difference between these 2 classes is the speed. The ‘Feared-Sad’ and ‘Sad-Feared’
 336 also demonstrated an increase in ambiguity. This can be caused by the averaging
 337 effect by training the StarGAN model using all different types of motions. The most
 338 discriminative characteristics of the *Feared* class is the asynchronous finger movement.

		Perceived Emotion			
		Angry	Happy	Sad	Feared
Ground Truth Labels	Angry	73.33%	16.67%	0.00%	10.00%
	Happy	23.33%	60.00%	6.67%	10.00%
	Sad	0.00%	21.67%	60.00%	18.33%
	Feared	6.67%	0.00%	23.33%	70.00%

Table 4: Confusion matrix of the emotion recognition test on the synthesized hand motions in the user study.

Synthesized is more pleasant	46.67%
Do not know	13.33%
Captured is more pleasant	40.00%

Table 5: Emotion recognition accuracy (%) on the captured and synthesized hand motions in the user study.

339 Learning a generic model for all different emotions and action types has reduced the
 340 discriminative characteristics in the synthesized motions. It will be an interesting future
 341 direction to explore a better way to separate the emotion, action and personal style into
 342 different components in the motion to better preserve the motion characteristics.

343 4.1.3. Comparing the naturalness and visual quality of the synthesized animations with
 344 the captured motions and baseline

345 To evaluate the visual quality and naturalness of the synthesized hand animations,
 346 we follow the design of the user study conducted in [39] by playing back the captured
 347 and synthesized motions side-by-side, and ask the user which one is 'more pleasant'
 348 or the user can answer 'do not know' if it is difficult to judge. In this experiment, 4
 349 pairs of animations were randomly selected from our database and shown to each user.
 350 We randomly selected *non-neutral* motions from the captured data and paired them up
 351 with the corresponding synthesized motions, which are emotionally transferred from
 352 *neutral* to the target emotion state. The results are summarized in Table 5. From the
 353 results, it can be seen that both of the synthesized and captured motions have received
 354 a similar percentage of users rating as 'more pleasant', while the motion produced by
 355 our method is 6.67% higher than the captured motions. There are 13.3% of users who
 356 cannot decide which motion is better. To validate the user study results, A/B testing
 357 is used to find out the statistical significance of the results. By treating the captured
 358 motions as variant A and our synthesized motions as variant B, the p-value is 0.2301.
 359 This suggests the conclusion on '*the synthesized motions are visually more pleasant than the*
 360 *captured motions*' is not statistically significant at the 95% confidence interval. On the
 361 other hand, when including our synthesized motions and the 'do not know' option in
 362 variant B, the p-value becomes 0.0127. This suggests the conclusion on '*the synthesized*
 363 *motions are not visually less pleasant than the captured motions*' is statistically significant
 364 at the 95% confidence interval. In summary, the visual quality between the captured and
 365 synthesized hand motions are similar, and arguably our method will not degrade the
 366 visual quality of the input hand motion, as the results indicate.

367 We further compared the motions synthesized by our method with those created
 368 using a PCA-based method proposed by Irimia et al. [15]. Again, we follow the user
 369 study explained above to evaluate the difference in the visual quality and naturalness
 370 of the synthesized motions. A side-by-side comparison will be given to the user to
 371 rate whether ours or Irimia et al. [15]'s method produces 'more pleasant' animation or
 372 no decision can be made. Each user was asked to rate 4 pairs of animations, and the
 373 results are presented in Table 6. The results show that the animations synthesized by our
 374 methods have better visual quality than those generated by Irimia et al. [15] with a more
 375 significant margin of 8.33% more users rated our results as 'more pleasant', although

Synthesized is more pleasant	45.00%
Do not know	18.33%
Irimia et al. [15] is more pleasant	36.67%

Table 6: Emotion recognition accuracy (%) on the captured and synthesized hand motions in the user study.

376 the p-value ($p = 0.1757$) computed in the A/B test suggests that the results are not
 377 statistically significant at the 95% confidence interval. Similar to the comparison between
 378 the captured and synthesized motions, we group our synthesized motions and the ‘do
 379 not know’ option and compare with the results created generated by Irimia et al. [15].
 380 The p-value becomes 0.0012 which suggests ‘our synthesized motions are not visually less
 381 pleasant than those generated by Irimia et al. [15]’. Again, the results highlight the methods
 382 compared in this study are producing motions with similar visual quality. In addition
 383 to the visual quality, the capability of multi-class emotion transfers in our proposed
 384 method is another advantage over Irimia et al. [15] since their approach is essentially
 385 interpolating the motion between two different emotional states.

386 4.2. Evaluation on Emotion Transfer

387 In this section, a wide range of results synthesized by our method are presented.
 388 We will first present the synthesized hand animations in Section 4.2.1 and 4.2.2. Next,
 389 the synthesized full-body animations will be discussed in Section 4.2.3. The animations
 390 are also included in the accompanying video demo.

391 4.2.1. Hand Animation

392 Here, we demonstrate the effectiveness of the proposed method by showing some
 393 of the synthesized hand animations. Like the experiments mentioned earlier, we used
 394 an unseen neutral hand motion as input and synthesized the animations by specifying
 395 the emotion labels. Due to the limited space, we visualize the results (Figure 4a to 4d)
 396 on 4 motion sets including *crawling*, *patting*, *impatient* and *pushing*. In each figure, each
 397 row contains 5 hand models which are animated by motions with different emotions
 398 (from left to right): *angry*, *happy*, *input (neutral)*, *sad*, *fearful*. The 4 rows in each figure are
 399 referring to the keyframes of the 4 progression stages (i.e. 0%, 33%, 67%, and 100%) in
 400 each animation.

401 The experimental results are consistent. To assess the correctness of the synthesized
 402 motion, we can compare the changes of the motion between keyframes in each column
 403 in each figure and evaluate if the changes align with the characteristics listed in Table
 404 1. Specifically, input motions become more exaggerated by transferring to *angry*. The
 405 range of motion increases, and the motion becomes faster. This is highlighted by the
 406 movement of the thumb in the video demo. By transferring to the *happy* emotion, the
 407 motion becomes more energetic with a larger range of motion when compared with the
 408 input (*neutral*) motion. The motion’s speed is getting higher as well, although the motion
 409 is less exaggerated than those transferred to *angry*. With the *sad* emotion, the synthesized
 410 motions show the sign of tiredness, which results in slower movement. Finally, the
 411 *fearful* emotion brings the asynchronous finger movements to the neutral motion as those
 412 characteristics can be found in the captured data. In summary, the consistent observation
 413 of the synthesized motions highlighted the effectiveness of our framework.

414 4.2.2. Comparing Emotion Transferred Motions with Captured Data

415 Next, we compare the synthesized motions with the captured data. Recall that
 416 leave-one-out cross-validation is used in defining the training and testing data sets. As a
 417 result, the motion type of the input (i.e., testing) motion is not included in the training
 418 data. It is possible that the action of the synthetic motion looks slightly different from

419 the captured motion. Having said that, an effective emotion transfer framework should
420 be able to transfer the characteristics of the corresponding emotion to the new motion.

421 The results are illustrated in Figure 5a to 5e. Similar to Section 4.2.1, we extracted
422 3 keyframes (i.e. each row in each figure) at the different progression stages (0%, 50%
423 and 100%) of the animation. The hand models in each column (from left to right) were
424 animated by the input (*neutral*), synthesized (i.e. emotion transferred) and captured
425 motions, respectively. Here, readers can focus on whether the synthesized motions
426 (middle column in each figure) contains the characteristics of the corresponding emotion
427 as in the captured motions (right column in each figure). Readers can also compare the
428 difference between the input (*neutral*, left column in each figure) and the synthesized
429 motion to evaluate the changes made by the proposed framework.

430 It can be seen that the motions synthesized by the proposed framework have the
431 characteristics of the corresponding emotion. For example, the motions are exaggerated
432 in the *angry* crawling and pushing motions in Figure 5a and 5d, respectively. On the
433 other hand, the *sad* crawling motion shows the sign of tiredness. The *fearful* emotion can
434 again transfer the asynchronous finger movements to the pushing motion, as illustrated
435 in Figure 5b. Finally, a larger range of motion can be seen in the *happy* impatient motion
436 (Figure 5c).

437 4.2.3. Body motion synthesis results

438 To further demonstrate the generality of the proposed framework, we trained the
439 proposed framework using 3D skeletal full-body motion in this experiment. Again,
440 unseen motions with the *neutral* emotion are used as input, and new motions are syn-
441 thesized by specifying the emotion labels. Three types of motions, including knocking,
442 lifting, and throwing, are included in the test and the screenshots of some examples are
443 illustrated in Figure 6 to 8. We selected 3 key moments from each animation representing
444 the early, middle, and late stages of the motion in the screenshots. To facilitate the side
445 by side comparison, we show the input motion (blue), synthesized motions (green) and
446 the captured motions (purple) in Figure 6 to 8.

447 In general, there is a consistent trend in terms of the difference in the speed between
448 motions with different emotions. Specifically, the *angry* motions are the fastest, with
449 *happy* motions are slightly slower, *neutral* motions are the third, and *sad* motions are the
450 slowest in most of the samples. Such pattern can be seen in the motions synthesized
451 by our method. Another observation from the results is the small difference between
452 the synthesized motions and the captured ones (i.e., ground truth). We believe the
453 small difference is mainly caused by the proposed method emphasized learning emo-
454 tion transfer without explicitly modeling the personal style differences from motions
455 performed by different subjects. As a result, small differences can be introduced when
456 synthesizing new motions by our method. This is an interesting further direction to
457 incorporate personal style in the motion modeling process to further strengthen the
458 proposed method.

459 5. Conclusion and Discussions

460 In this paper, we propose a new framework for synthesizing motion by emotion
461 transfer. We demonstrated the generality of the framework by modeling hand and full
462 body motions in a wide range of experiments. A user study is conducted to verify
463 the perceived emotion from the hand motions as well as evaluating the visual quality
464 and naturalness of the animations. Experimental results show that our method can 1)
465 generate different styles of motions according to the emotion type, 2) the characteristics
466 in each emotion type can be transferred to new motions, and 3) achieving similar or
467 better visual quality when comparing the hand motions synthesized by our method
468 with those captured motions and created by [15].

469 In the future, we are interested in incorporating the personal style into the motion
470 modeling framework. In addition to specifying the emotion labels to synthesize different

471 motions, de-tangling the ‘base’ motion and personal style can further increase the
472 variations in the synthesized motions. Another further direction will be evaluating the
473 feasibility of having multiple emotion labels with different levels of strengths for motion
474 representation and synthesis. Such a direction is inspired by the emotion recognition
475 results of the user study in which not all users are agreed on a single type of emotion
476 being associated with each hand motions in the study. It is also an interesting future
477 direction to quantitatively evaluate the results by comparing the differences between
478 the synthesized and ground-truth motion numerically. To achieve this goal, more hand
479 motions have to be captured. As in our pilot study, we have difficulties capturing the
480 global translation and rotation in high quality. As a result, the global transformation is
481 discarded, which limits the expression of emotion. One possible solution is to capture
482 the hand motions using state-of-the-art MOCAP solutions such as [14].

483 Acknowledgements

484 We gratefully acknowledge the support of NVIDIA Corporation with the donation
485 of the Titan Xp GPU used for this research.

486 References

- 487 1. Karras, T.; Aila, T.; Laine, S.; Herva, A.; Lehtinen, J. Audio-Driven Facial Animation by
488 Joint End-to-End Learning of Pose and Emotion. *ACM Trans. Graph.* **2017**, *36*. doi:
489 10.1145/3072959.3073658.
- 490 2. Tinwell, A.; Grimshaw, M.; Nabi, D.A.; Williams, A. Facial expression of emotion and
491 perception of the Uncanny Valley in virtual characters. *Computers in Human Behavior* **2011**,
492 *27*, 741–749. Web 2.0 in Travel and Tourism: Empowering and Changing the Role of Travelers,
493 doi:<https://doi.org/10.1016/j.chb.2010.10.018>.
- 494 3. Courgeon, M.; Clavel, C. MARC: a framework that features emotion models for facial
495 animation during human–computer interaction. *Journal on Multimodal User Interfaces* **2013**,
496 *7*, 311–319. doi:10.1007/s12193-013-0124-1.
- 497 4. Ruttkay, Z.; Noot, H.; Ten Hagen, P. Emotion Disc and Emotion Squares: Tools
498 to Explore the Facial Expression Space. *Computer Graphics Forum* **2003**, *22*, 49–
499 53, [<https://onlinelibrary.wiley.com/doi/pdf/10.1111/1467-8659.t01-1-00645>]. doi:
500 <https://doi.org/10.1111/1467-8659.t01-1-00645>.
- 501 5. Chan, J.C.P.; Irimia, A.S.; Ho, E.S.L. Emotion Transfer for 3D Hand Motion using Star-
502 GAN. *Computer Graphics and Visual Computing (CGVC)*; Ritsos, P.D.; Xu, K., Eds. The
503 Eurographics Association, 2020. doi:10.2312/cgvc.20201146.
- 504 6. Wang, Y.; Tree, J.E.F.; Walker, M.; Neff, M. Assessing the Impact of Hand Motion on Virtual
505 Character Personality. *ACM Trans. Appl. Percept.* **2016**, *13*. doi:10.1145/2874357.
- 506 7. Jörg, S.; Hodgins, J.; Safonova, A. Data-Driven Finger Motion Synthesis for Gesturing
507 Characters. *ACM Trans. Graph.* **2012**, *31*. doi:10.1145/2366145.2366208.
- 508 8. Ye, Y.; Liu, C.K. Synthesis of Detailed Hand Manipulations Using Contact Sampling. *ACM*
509 *Trans. Graph.* **2012**, *31*, 41:1–41:10. doi:10.1145/2185520.2185537.
- 510 9. Liu, C.K. Synthesis of Interactive Hand Manipulation. *Proceedings of the 2008 ACM*
511 *SIGGRAPH/Eurographics Symposium on Computer Animation*; Eurographics Association:
512 Aire-la-Ville, Switzerland, Switzerland, 2008; SCA '08, pp. 163–171.
- 513 10. Andrews, S.; Kry, P.G. Policies for Goal Directed Multi-Finger Manipulation. *VRIPHYS*,
514 2012.
- 515 11. Liu, C.K. Dexterous Manipulation from a Grasping Pose. *ACM SIGGRAPH 2009 Papers*; ACM:
516 New York, NY, USA, 2009; SIGGRAPH '09, pp. 59:1–59:6. doi:10.1145/1576246.1531365.
- 517 12. Bai, Y.; Liu, C.K. Dexterous manipulation using both palm and fingers. *2014 IEEE In-*
518 *ternational Conference on Robotics and Automation (ICRA)*, 2014, pp. 1560–1565. doi:
519 10.1109/ICRA.2014.6907059.
- 520 13. Alexanderson, S.; O’Sullivan, C.; Beskow, J. Robust online motion capture labeling of finger
521 markers. *Proceedings of the 9th International Conference on Motion in Games*. ACM, 2016,
522 pp. 7–13.
- 523 14. Han, S.; Liu, B.; Wang, R.; Ye, Y.; Twigg, C.D.; Kin, K. Online Optical Marker-Based Hand
524 Tracking with Deep Labels. *ACM Trans. Graph.* **2018**, *37*. doi:10.1145/3197517.3201399.

- 525 15. Irimia, A.S.; Chan, J.C.P.; Mistry, K.; Wei, W.; Ho, E.S.L. Emotion Transfer for Hand Animation.
526 Motion, Interaction and Games; ACM: New York, NY, USA, 2019; MIG '19, pp. 41:1–41:2.
527 doi:10.1145/3359566.3364692.
- 528 16. Wheatland, N.; Jörg, S.; Zordan, V. Automatic Hand-Over Animation Using Principle Com-
529 ponent Analysis. Proceedings of Motion on Games; Association for Computing Machinery:
530 New York, NY, USA, 2013; MIG '13, p. 197–202. doi:10.1145/2522628.2522656.
- 531 17. Unuma, M.; Anjyo, K.; Takeuchi, R. Fourier Principles for Emotion-based Human Figure
532 Animation. Proceedings of the 22Nd Annual Conference on Computer Graphics and
533 Interactive Techniques; ACM: New York, NY, USA, 1995; SIGGRAPH '95, pp. 91–96. doi:
534 10.1145/218380.218419.
- 535 18. Amaya, K.; Bruderlin, A.; Calvert, T. Emotion from Motion. Proceedings of the Conference
536 on Graphics Interface '96; Canadian Information Processing Society: Toronto, Ont., Canada,
537 Canada, 1996; GI '96, pp. 222–229.
- 538 19. Brand, M.; Hertzmann, A. Style machines. Proceedings of the 27th annual conference on
539 Computer graphics and interactive techniques, 2000, pp. 183–192.
- 540 20. Urtasun, R.; Glardon, P.; Boulic, R.; Thalmann, D.; Fua, P. Style-based motion synthesis.
541 Computer Graphics Forum. Wiley Online Library, 2004, Vol. 23, pp. 799–812.
- 542 21. Ikemoto, L.; Arikian, O.; Forsyth, D. Generalizing motion edits with gaussian processes.
543 *ACM Transactions on Graphics (TOG)* **2009**, *28*, 1–12.
- 544 22. Xia, S.; Wang, C.; Chai, J.; Hodgins, J. Realtime style transfer for unlabeled heterogeneous
545 human motion. *ACM Transactions on Graphics (TOG)* **2015**, *34*, 1–10.
- 546 23. Hsu, E.; Pulli, K.; Popović, J. Style Translation for Human Motion. ACM SIGGRAPH
547 2005 Papers; ACM: New York, NY, USA, 2005; SIGGRAPH '05, pp. 1082–1089. doi:
548 10.1145/1186822.1073315.
- 549 24. Shapiro, A.; Cao, Y.; Faloutsos, P. Style Components. Proceedings of Graphics Interface 2006;
550 Canadian Information Processing Society: Toronto, Ont., Canada, Canada, 2006; GI '06, pp.
551 33–39.
- 552 25. Holden, D.; Habibie, I.; Kusajima, I.; Komura, T. Fast Neural Style Transfer for Motion Data.
553 *IEEE Computer Graphics and Applications* **2017**, *37*, 42–49. doi:10.1109/MCG.2017.3271464.
- 554 26. Smith, H.J.; Cao, C.; Neff, M.; Wang, Y. Efficient Neural Networks for Real-Time Motion
555 Style Transfer. *Proc. ACM Comput. Graph. Interact. Tech.* **2019**, *2*. doi:10.1145/3340254.
- 556 27. Lee, S.J.; Popović, Z. Learning Behavior Styles with Inverse Reinforcement Learning. *ACM*
557 *Trans. Graph.* **2010**, *29*. doi:10.1145/1778765.1778859.
- 558 28. Selim, A.; Elgharib, M.; Doyle, L. Painting Style Transfer for Head Portraits Using
559 Convolutional Neural Networks. *ACM Trans. Graph.* **2016**, *35*, 129:1–129:18. doi:
560 10.1145/2897824.2925968.
- 561 29. Elad, M.; Milanfar, P. Style Transfer Via Texture Synthesis. *IEEE Transactions on Image*
562 *Processing* **2017**, *26*, 2338–2351. doi:10.1109/TIP.2017.2678168.
- 563 30. Matsuo, S.; Shimoda, W.; Yanai, K. Partial style transfer using weakly supervised seman-
564 tic segmentation. 2017 IEEE International Conference on Multimedia Expo Workshops
565 (ICMEW), 2017, pp. 267–272. doi:10.1109/ICMEW.2017.8026228.
- 566 31. Ma, Y.; Paterson, H.M.; Pollick, F.E. A motion capture library for the study of identity, gender,
567 and emotion perception from biological motion. *Behavior Research Methods* **2006**, *38*, 134–141.
- 568 32. Chan, J.C.P.; Shum, H.P.H.; Wang, H.; Yi, L.; Wei, W.; Ho, E.S.L. A generic framework for edit-
569 ing and synthesizing multimodal data with relative emotion strength. *Computer Animation*
570 *and Virtual Worlds* **2019**, *30*, e1871, [<https://onlinelibrary.wiley.com/doi/pdf/10.1002/cav.1871>].
571 e1871 cav.1871, doi:10.1002/cav.1871.
- 572 33. Zhu, J.; Park, T.; Isola, P.; Efros, A.A. Unpaired Image-to-Image Translation Using Cycle-
573 Consistent Adversarial Networks. 2017 IEEE International Conference on Computer Vision
574 (ICCV), 2017, pp. 2242–2251.
- 575 34. Yi, Z.; Zhang, H.; Tan, P.; Gong, M. DualGAN: Unsupervised Dual Learning for Image-to-
576 Image Translation. 2017 IEEE International Conference on Computer Vision (ICCV), 2017,
577 pp. 2868–2876.
- 578 35. Choi, Y.; Choi, M.; Kim, M.; Ha, J.; Kim, S.; Choo, J. StarGAN: Unified Generative Adversarial
579 Networks for Multi-domain Image-to-Image Translation. 2018 IEEE/CVF Conference on
580 Computer Vision and Pattern Recognition, 2018, pp. 8789–8797.
- 581 36. Holden, D.; Saito, J.; Komura, T. A Deep Learning Framework for Character Motion Synthesis
582 and Editing. *ACM Trans. Graph.* **2016**, *35*. doi:10.1145/2897824.2925975.

-
- 583 37. Kipf, T.N.; Welling, M. Semi-Supervised Classification with Graph Convolutional Networks.
584 International Conference on Learning Representations (ICLR), 2017.
- 585 38. Men, Q.; Ho, E.S.L.; Shum, H.P.H.; Leung, H. A Quadruple Diffusion Convolutional
586 Recurrent Network for Human Motion Prediction. *IEEE Transactions on Circuits and Systems
587 for Video Technology* **2020**, pp. 1–1. doi:10.1109/TCSVT.2020.3038145.
- 588 39. Aristidou, A.; Zeng, Q.; Stavrakis, E.; Yin, K.; Cohen-Or, D.; Chrysanthou, Y.; Chen, B.
589 Emotion Control of Unstructured Dance Movements. Proceedings of the ACM SIGGRAPH /
590 Eurographics Symposium on Computer Animation; Association for Computing Machinery:
591 New York, NY, USA, 2017; SCA '17. doi:10.1145/3099564.3099566.

**(a)** *crawling***(b)** *impatient***(c)** *pushing***(d)** *patting*

Figure 4. Screenshots (one frame per row) of transferring the input (neutral) motion to different types of emotions. Columns from left to right: *angry, happy, input (neutral), sad, fearful*.



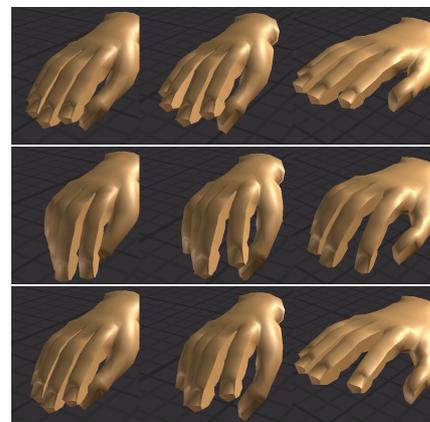
(a) *crawling*, transferred to *angry*

(b) *pushing*, transferred to *fearful*

(c) *impatient*, transferred to *happy*



(d) *pushing*, transferred to *angry*



(e) *crawling*, transferred to *sad*

Figure 5. Screenshots (one frame per row) of the comparison between the input (*neutral*, left), synthesized (i.e. emotion transferred, middle) and captured (right) motions.

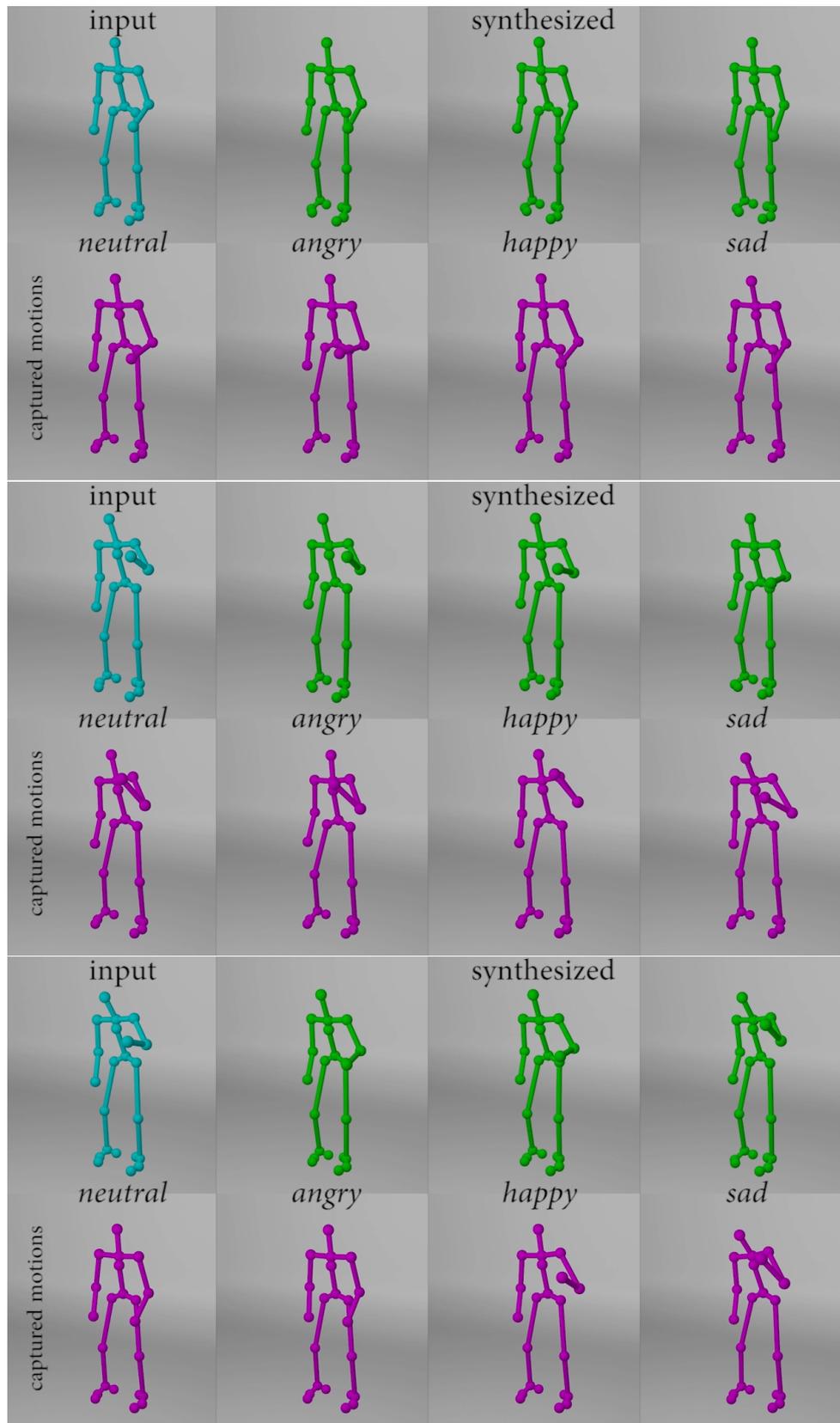


Figure 6. Screenshots of the *knocking* motion extracted from different stages - early (top row), middle (middle row), late (bottom row).

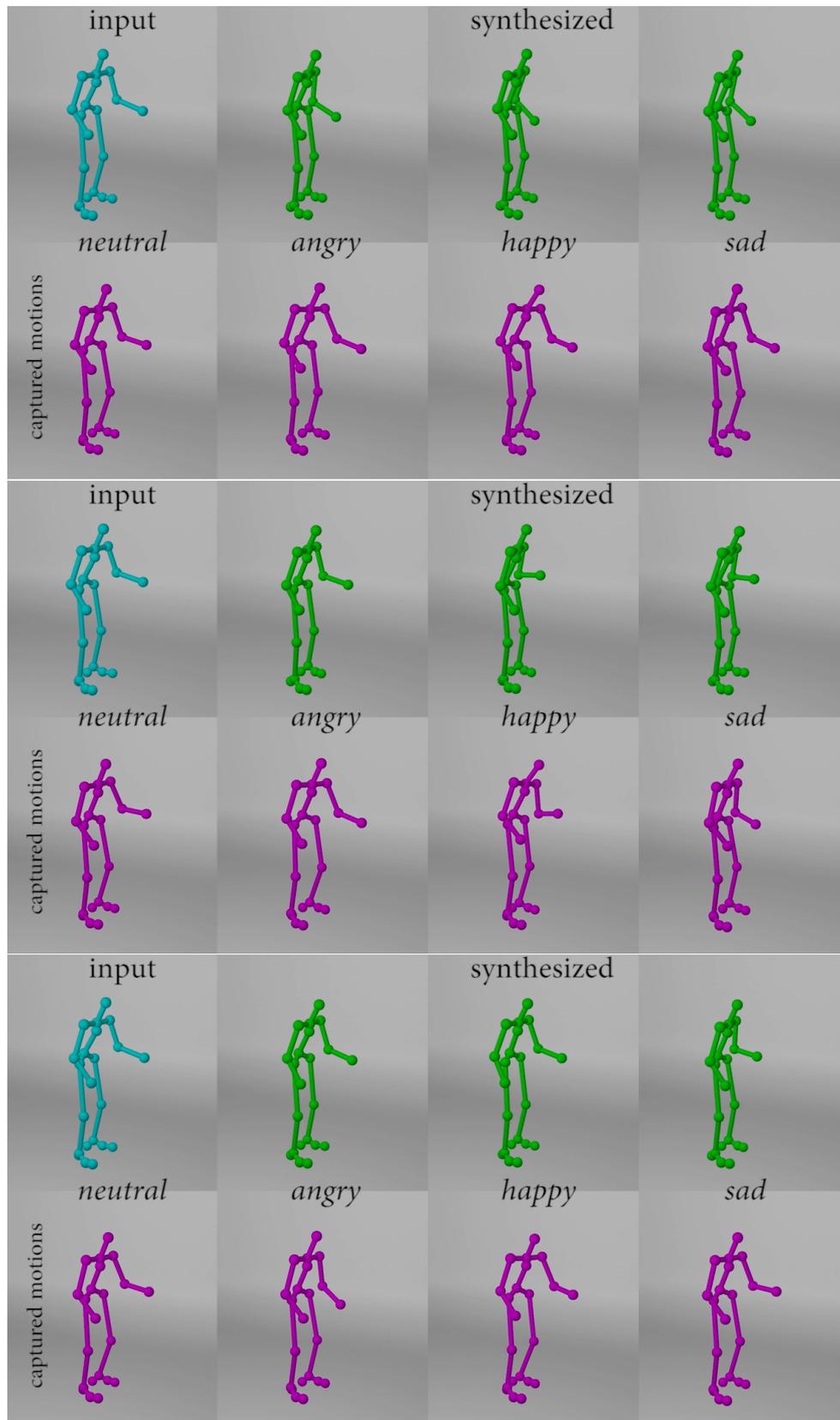


Figure 7. Screenshots of the *lifting* motion extracted from different stages - early (top row), middle (middle row), late (bottom row).

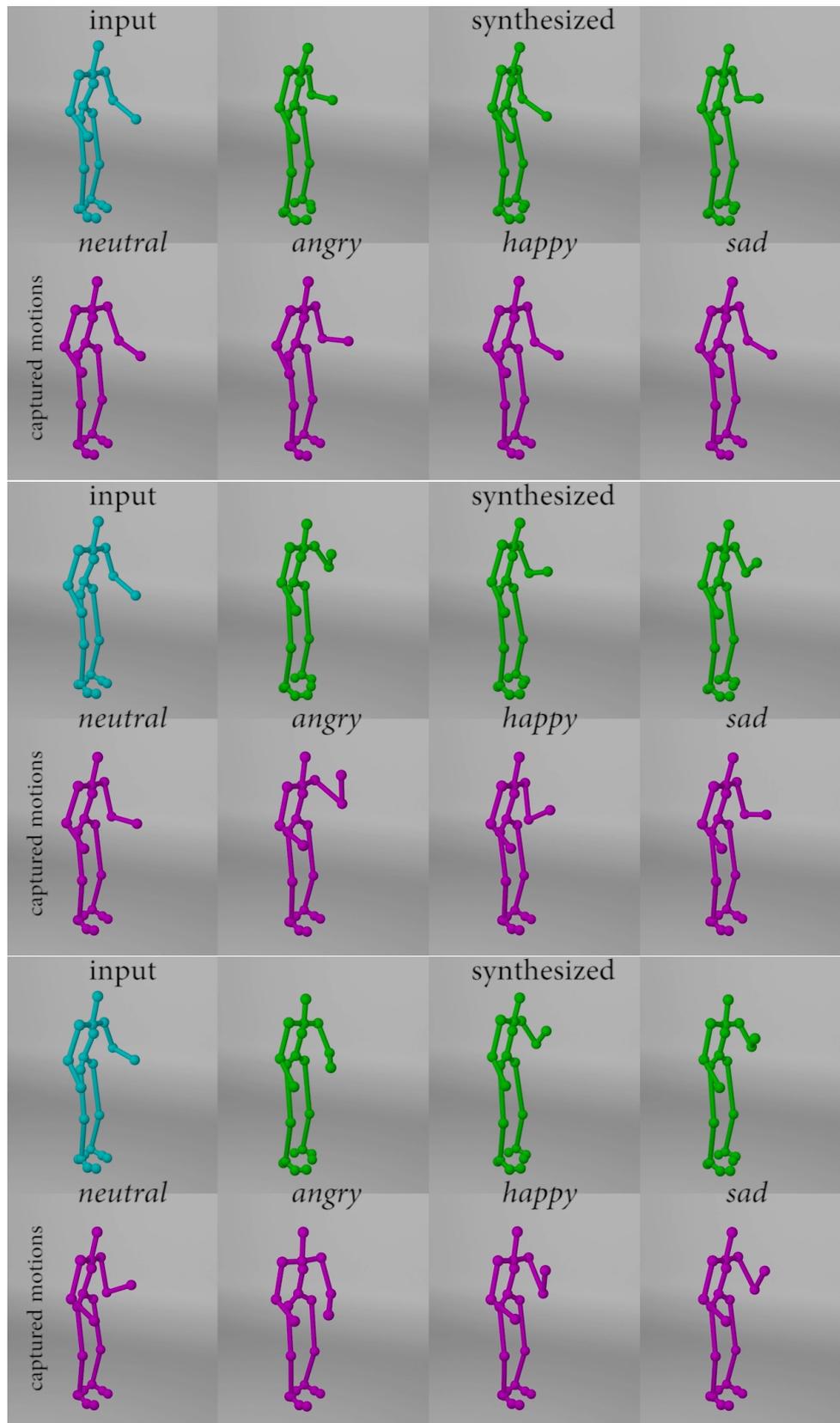


Figure 8. Screenshots of the *throwing* motion extracted from different stages - early (top row), middle (middle row), late (bottom row).

