# Toward an Automated Measure of Social Engagement for Children with Autism Spectrum Disorder – a Personalized Computational Modeling Approach

1  **Hifza Javed[1], WonHyong Lee[2], Chung Hyuk Park[1*]**

2  [1]Assistive Robotics and Telemedicine Laboratory, Department of Biomedical Engineering, School of
3  Engineering and Applied Science, The George Washington University, Washington, DC, USA
4  [2]School of Computer Science and Electrical Engineering, Handong Global University, South Korea

5  **\* Correspondence:**
6  Chung Hyuk Park
7  chpark@gwu.edu

10  **Abstract**

11  Social engagement is a key indicator of an individual's socio-emotional and cognitive states. For a
12  child with Autism Spectrum Disorder (ASD), this serves as an important factor in assessing the
13  quality of the interactions and interventions. So far, qualitative measures of social engagement have
14  been used extensively in research and in practice, but a reliable, objective, and quantitative measure
15  is yet to be widely accepted and utilized. In this paper, we present our work on the development of a
16  framework for the automated measurement of social engagement in children with ASD that can be
17  utilized in real-world settings for the long-term clinical monitoring of a child's social behaviors as
18  well as for the evaluation of the intervention methods being used. We present a computational
19  modeling approach to derive the social engagement metric based on a user study with children
20  between the ages of 4 and 12 years. The study was conducted within a child-robot interaction setting
21  that targets sensory processing skills in children. We collected video, audio and motion-tracking data
22  from the subjects and used them to generate personalized models of social engagement by training a
23  multi-channel and multi-layer convolutional neural network. We then evaluated the performance of
24  this network by comparing it with traditional classifiers and assessed its limitations, followed by
25  discussions on the next steps towards finding a comprehensive and accurate metric for social
26  engagement in ASD.

27  **1    Introduction**

28  Social engagement of a child is an indicator of his/her socioemotional and cognitive states. It is the
29  interaction of a child with the environment in a contextually appropriate manner and reflects a
30  complex internal state that signifies the occupation of the child with a person or a task. Much of the
31  research so far has relied on the perceptual evaluation of engagement, utilizing questionnaires and
32  behavioral assessments administered by trained professionals, which typically attempt to identify key
33  behavioral traits that serve as important indicators of social engagement. Automatic quantification of
34  engagement is still limited but can allow not only for an objective interpretation of engagement and
35  the contributing target behaviors, but also help to identify methods to improve engagement in
36  different settings, especially when targeting a specific health condition. Therefore, it serves both as
37  an outcome measure and as an objective measure of the quality of an activity, interaction, or
38  intervention [1].

39  Social engagement has often been reported to be particularly deficient in children with Autism
40  Spectrum Disorder (ASD). ASD is a neurodevelopmental disorder that causes significant impairment
41  in three broad areas of functioning: communication, social interaction, and restricted and repetitive
42  behaviors [2]. This means that children interact with their peers infrequently, thus preventing the
43  formation of lasting and meaningful social relationships and resulting in social withdrawal. These
44  children often feel isolated from or rejected by peers and are more likely to develop behavioral
45  problems [3] as well as anxiety and depression [4][5].
46
47  Behavioral and physiological cues can provide insight into the engagement state of a child, with
48  gestures, subtle body language changes, facial expressions, vocal behaviors, and various
49  physiological signals, all carrying significant indications of a child's level of interest and engagement
50  in an interaction. Eye gaze focus, smiling, vocalizations, joint-attention, imitation, self-initiated
51  interactions and triadic interactions are among the important behavioral cues that can be utilized to
52  assess engagement [6-17]. Heart rate, electrodermal activity, electrocardiography, electromyography,
53  blood pressure etc. are among the key physiological indicators of engagement state [18-20]. A
54  combination of these multi-modal behavioral and physiological features can present a comprehensive
55  feature set for effective engagement evaluation.
56
57  A major hurdle in the path toward automated measurement of social engagement is of the
58  identification and classification of these key behaviors. While it may be a simple task for trained
59  professionals to identify these high-level behaviors and infer a fairly accurate engagement state from
60  real-time observations of a child's interactions, it remains a considerable challenge for the state-of-
61  the-art algorithms and machines. Instead, the current technologies are better equipped to extract
62  lower-level behaviors that can be used as a rough estimation of the target behaviors.
63
64  This paper presents our first step toward an automated quantifiable measure of social engagement
65  derived from behavioral data collected from two groups of children, one typically developing (TD)
66  and one with ASD. Research from our team thus far has focused on child-robot interaction scenarios
67  that target several ASD symptoms, including sensory processing [21], imitation [22], emotion
68  recognition and emotion regulation skills [23]. In these studies, we collected multi-modal interaction
69  data, including video and audio recordings, as well as motion tracking data. The overall goal of our
70  work is to develop a framework for personalized child-robot interactions for ASD. To this end, our
71  framework aims to 1) sense important features of a child's interaction with a robot, 2) interpret and
72  derive meaningful deductions about a child's engagement in the interaction, 3) identify target
73  behaviors that may be lacking in the detected interaction pattern, 4) reassess the current robot
74  behavior strategy and modulate it to elicit a higher level of engagement from the child. This paper
75  focuses on step 2 of the above approach by processing the multimodal behavioral data collected from this
76  study through a deep learning-based multi-label classification model in order to contribute towards
77  deriving an automated measure of social engagement.
78
79  This paper is organized as follows. Section 2 discusses the previous studies that have designed
80  methods to formulate an automated measure of social engagement. Section 3 describes the child-
81  robot interaction scenario we used in this study. Sections 4 and 5 present the modalities of the data
82  we collected during our experiments and the methods we employed to label these data. Sections 6
83  and 7 discuss our feature extraction methods and design of our convolutional neural network for
84  multi-label classification. Sections 8, 9 and 10 describe the user study, its results and a comparison of
85  the proposed network with other classical algorithms. Section 11 presents a discussion on these
86  findings while Section 12 concludes this paper with comments on the future work.

87  **2    Related work**
88  Several studies in the past have contributed to this area of research with each method typically
89  varying in terms of the feature set, number of engagement classes and computational model that were
90  used, as well as the demographics of the participants from whom the data were collected. Rajgopalan
91  et al. [24] showed the feasibility of utilizing low-level behavioral features in the absence of accurate
92  high-level features, and used a two-stage approach to first find hidden structures in the data (using
93  Hidden Conditional Random Fields) and then learn them through a Support Vector Machine (SVM).
94  Only head pose orientation estimates were used to assess engagement and the approach was
95  evaluated by conducting experiments on labeled child interaction data from the Multimodal Dyadic
96  Behavior Dataset [25], obtaining an accuracy of around 70%.

98  Gupta et. al. [26] designed an engagement prediction system that utilized only the prosodic features
99  of a child's speech as observed during a structured interaction between a child and a psychologist
100  involving several tasks from the Rapid ABC database. Three engagement classes and two levels of
101  prosodic features (local for short-term and global for task-wide patterns) were defined. The system
102  achieved an unweighted average recall of 55.8%, where the best classification results were obtained
103  by using an SVM that utilized both categories of the prosodic features. Another study by Lala et. al.
104  [27] used several verbal and non-verbal behavioral features, including nodding, eye gaze, laughing
105  and verbal backchannels. The authors collected their own dataset comprising audio and video
106  recordings based on conversational scenarios between a human user and a humanoid robot, while
107  human annotators provided labels to establish ground truth. A Bayesian binary classifier was used to
108  classify the user as engaged or not engaged and obtained an AUC (area under the precision-recall
109  curve) score of 0.62.

111  A study from Castellano et.al. [28] used both behavioral features from the user (gaze focus and
112  smiling) and contextual information from the activity in order to train a Bayesian classifier to detect
113  engagement in users for a child-robot interaction scenario. The labels generated from human coding
114  were based only on the two user behaviors. The authors reported only a slight improvement in the
115  classifier recognition rate when using both behavioral and contextual features (94.79%) versus when
116  only behavioral features were utilized (93.75%), highlighting the key importance of the behavioral
117  information.

119  Kim et. al. [29] investigated the use of vocal/acoustic features in determining child engagement in
120  group interaction scenarios. The annotation scheme involves the giving and receiving of attention
121  from other group members. They used a combination of ordinal regression and ranking with SVM to
122  detect engagement in children and found this technique to outperform classification, simple
123  regression and rule-based approaches. Such a system may be acceptable to use with typically-
124  developing children, but since children with ASD may often be non-verbal and/or shy or unwilling to
125  communicate using speech/vocalizations, the exclusive use of acoustic features may not be suited to
126  research involving the ASD population.

128  Another study from Parekh et. al. [30] developed a video system for measuring engagement in
129  patients with dementia, which uses deep-learning based computer vision algorithms to evaluate their
130  engagement in an activity to provide behavior analytics based on facial expression and gaze analysis.
131  Ground truth was extracted through scoring performed by human annotators by classifying
132  engagement states in terms of attention and attitude. The video system presented in this study was
133  exclusively tested with elderly patients with dementia who were required to participate in a digital
134  interaction while seated directly in front of the camera. Additionally, since only facial expressions

135   and gaze features were utilized, the proximity of the participants to the camera was important, hence,
136   limiting their physical movements.
137
138   Oertel et. al. [31] studied the relation between group involvement and individual engagement using
139   several features of eye gaze patterns defined as presence, entropy, symmetry and maxgaze. They
140   utilized the Stockholm Werewolf Corpus, which is a video dataset of participants engaging in a game
141   that involved the use of speech and eye gaze. Once again, since only eye gaze patterns were used as
142   features to train a classifier, participants were required to remain seated in front of the cameras.
143
144   A study that specifically tested their system on the ASD population was from Anzalone et. al. [32]
145   that used a combination of static (focus of attention, head stability and body posture stability) and
146   dynamic (joint attention, synchrony, and imitation) metrics within two distinct use cases including
147   one where the robot attempts to learn the colors in its environment with the help of a human, and
148   another that elicits joint attention from participating children with ASD. The features were extracted
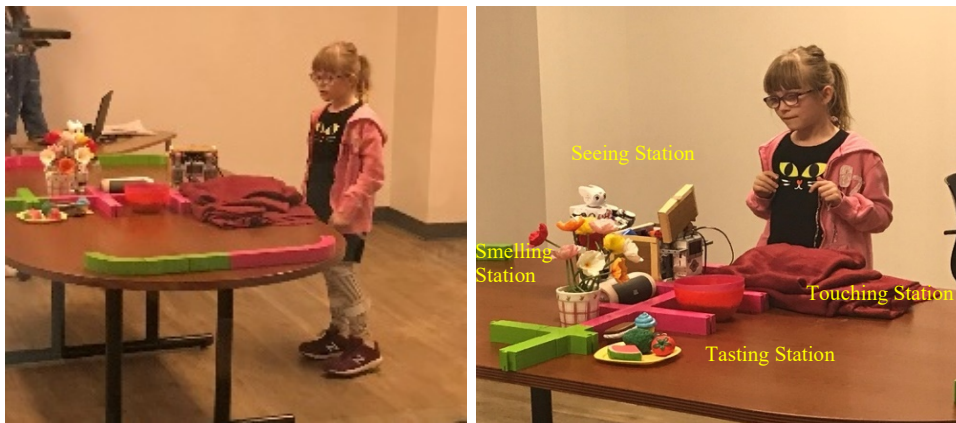149   using histogram heatmaps and clustered using the K-means algorithm.



Figure 1. Station set up for the sensory maze game (the child's photo rights reserved).

150   In [33] Rudovic et. al. also targeted the automated measurement of engagement for ASD children
151   with multimodal data collection including features from video (facial expressions, head movements,
152   body movements, poses, and gestures), audio, and physiological (heart rate, electrodermal activity
153   and heart rate) data. The child-robot interaction setting involved an emotion recognition activity with
154   a humanoid robot that required children to be seated in front of the robot [34]. Participating children
155   belonged to one of two cultures (Eatsern European and Asian) and the cultural differences were also
156   taken into account during engagement estimation. The authors generated ground truth through expert
157   human labelers who marked changes in engagement on a 0-5 Likert scale that is based on the
158   different levels of prompting required from the therapist during the interaction with the robot. In fact,
159   in this work, child engagement is considered to be a function of task-driven behavioral engagement
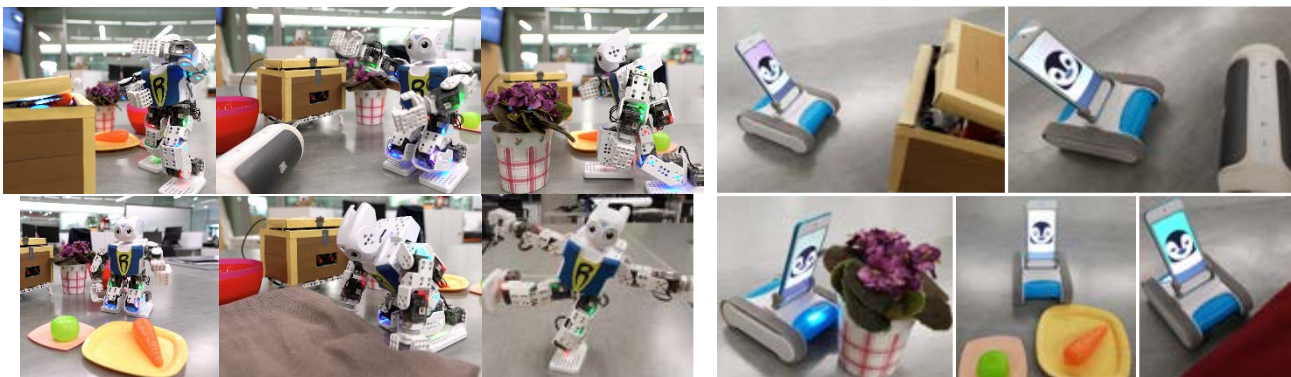160   and affective engagement.
161
162   Despite the overlap, this approach is significantly different from the one proposed in this paper in
163   several ways. Firstly, we define engagement as a function of several key behavioral indicators that
164   provide an insight into an individual's internal engagement state [21], which generates a novel
165   measure to estimate social engagement state i.e. the engagement index. Additionally, our methods do
166   not restrict the movement of the subjects by requiring them to be seated in front of a camera or a
167   robot, and the interaction design allows for free, naturalistic movement in order to closely resemble
168   real-world social settings as opposed to other restrictive experimental approaches. Importantly, this

169   approach toward engagement estimation can be easily generalized to any child, with or without ASD,
170   and to a variety of different, interactive experimental settings that may or may not involve a robot.
171
172   The work described in this paper presents a social engagement prediction system for children. It
173   utilizes a combination of features extracted from facial expressions and upper body motion tracking
174   data to train a deep convolutional neural network that can then classify the engagement state of a
175   child. We intentionally designed the experiments to not be strictly structured in order to encourage
176   naturalistic and unguided child-robot interactions during data collection that impose no restrictions
177   on the movement of a child. The nature of the features used in our approach allow for independence
178   of interaction context and can easily be extended to a variety of scenarios within laboratory or home
179   settings. In addition, a unique engagement model is obtained for every individual participant to
180   ensure personalized interaction with the robot, giving it potential to be used as an intervention tool
181   for ASD.

182   ## 3      Interaction Scenario Design
183   For this work, we used socially assistive robots to design a child-robot interaction that targeted the
184   sensory processing difficulties in ASD, as detailed in our previous work [21]. In this pedagogical
185   setting, two different mobile robots were used to model socially acceptable responses to potentially
186   overwhelming sensory stimulation that a child is likely to encounter in everyday experiences. The
187   humanoid robot, Robotis Mini (from Robotis) and the iPod-based robot, Romo (from Romotive) both
188   had their unique set of capabilities. While Mini used gestures and speech to communicate, Romo relied
189   mostly on its large set of emotional expressions and some movements.



Mini at the stations.
Top (L-R): Seeing station, Hearing station, Smelling station
Bottom (L-R): Tasting station, Touching station, Celebration station

Romo at the stations.
Top (L-R): Seeing station, Hearing station
Bottom (L-R): Smelling station, Tasting station, Touching station

Figure 2. The two robots at each sensory station.

190
191   A maze-like setup consisting of a station for each of the visual, auditory, olfactory, gustatory, tactile
192   and vestibular senses was used, as shown in Figure 1. Though one of the goals of the interaction was
193   to leverage the relationship between a robot and a child with ASD, as established by a plethora of
194   previous research [35-38], the focus of this work [21] was to assess the potential of this setup as a
195   tool to socially engage children with ASD and to use the collected data to contribute towards deriving
196   an automated measure of social engagement. Each sensory station simulated an everyday experience,
197   such as encountering bright lights at the *Seeing station*, loud music at the *Hearing station*, scented
198   flowers at the *Smelling station*, different food items at the *Tasting station*, materials with different
199   textures at the *Touching station* and summersaulting to celebrate at the vestibular station (Figure 2).
200   These scenarios were chosen to incorporate everyday stimulation that all children experience in

201  uncontrolled environments like malls, playgrounds, cinemas etc. and in the activities of daily living
202  such as eating meals and dressing. This interaction was designed to be highly interactive and
203  engaging, and required the child to participate actively by answering questions from the robots,
204  following their instructions, and 'helping' them complete the maze. Details of this study, including
205  the nature of interaction between the children and the robots, can be found in [21].

206  ## 4    Multimodal Data Collection
207  A high-quality measure for social engagement estimation must take into account all behavioral and
208  physiological cues that can serve as quantifiers of social motivation and social interaction. As
209  discussed in Section 1, a number of behavioral traits and physiological signals can be used effectively
210  to this end. However, when designing an interaction for autistic children, their unique needs and
211  sensitivities must be taken into account. For this study, this meant that only non-contact sensors
212  could be used in order to limit tactile disturbances to the children and enable free movement to allow
213  for naturalistic interaction. The combination of sensors also needed to provide a wholistic and
214  accurate representation of a child's engagement changes over the length of the interaction.
215
216  We collected video recordings of the child-robot interactions with a camcorder placed in one corner
217  of the room, which was repositioned by an instructor as the child moved during the interaction. From
218  these recordings, we were able to extract audio data as well as 2-D motion tracking data with the
219  OpenPose library [39]. While OpenPose provides full body motion tracking (Figure 3), we were only
220  able to utilize upper body data since the chosen experimental setting meant that children were often
221  standing in front of the table that hosted the maze setup, preventing a full-body view from being
222  captured. In addition, OpenPose also allowed for the extraction of facial expression datapoints from
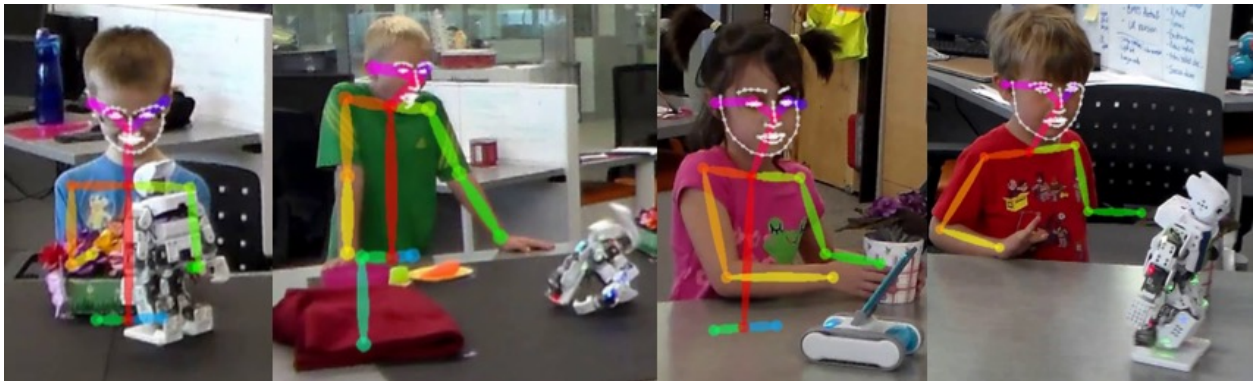223  the same video data.



Figure 3. Upper body and facial keypoints generated by OpenPose.

224  ## 5    Extracting Ground Truth
225  Unlike some of the previous studies described in Section 2, we did not use any existing video
226  datasets to test our methods. Since our goal was to derive an engagement measure specific to the
227  interactions that we designed for children with ASD, we opted to test our methods on the relatively
228  limited data available from our user study. To extract ground truth for a child's engagement in the
229  interaction with the robots, we defined six target behaviors that have been found to be key behavioral
230  indicators of social engagement [40-51]. These included eye gaze focus, vocalizations, smiling, self-
231  initiated interactions, triadic interactions and imitation.
232
233  Three raters then coded these videos using the Behavioral Observation Research Interactive Software
234  (BORIS) [52] to annotate the start and stop times of each target behavior as it was identified in the

235  video recordings. An inter-coder correlation (ICC) score of $0.8752 \pm 0.145$ was achieved for the 18
236  participants, which was used to evaluate the quality of the annotations. Details of the evaluation
237  criteria are reported in [21].
238
239  An eye gaze event was tagged each time the child's gaze moved to the robots or the setup and
240  stopped when the gaze focus was lost. Vocalizations comprised of any verbal expression from the
241  child, including but not limited to a shriek of excitement while interacting with the robots or the
242  utterance of words to communicate sentiments or queries regarding the robots. Smiling recorded all
243  events where a child was observed to visibly express joy in the form of a smile or laugh. Self-
244  initiated interactions involved all interactions with the robots or setup that are initiated by the child.
245  Triadic interactions comprised of an interaction where a child voluntarily involved a third entity in
246  the interaction with the robot, such as sharing their excitement with the parent. Lastly, imitations
247  included all events of voluntary imitation the robot's actions by the child. An in-depth report on the
248  inclusion criteria of the target behaviors, their significance and annotations in video data can be
249  found in [21].
250
251  Based on these annotations, multiple analytics were derived to quantify the social engagement with
252  respect to each robot and target behavior, and across stations to obtain a fine-grained analysis of the
253  child's interaction preferences [21]. However, for the current work, we have only used the raw time
254  series data of every child's changing engagement state as determined by the chosen target behaviors.
255  These overall engagement changes are shown in Figure 4, along with the subplots of each
256  contributing key behavior.
257
258  Therefore, each instance of time was mapped to an engagement state. Every behavior contributed a
259  factor of 1/6 to the engagement value, thus resulting in a metric with seven distinct values that ranged
260  from 0 (no target behavior observed) to 1 (all target behaviors observed).

261  **6   Feature Extraction**
262  An ideal automated engagement measure in this case would incorporate all of the above behaviors,
263  but also necessitates the automated classification of these behaviors. This is no trivial task, and
264  involves contributions from multiple disciplines including computer vision, speech analysis and
265  machine learning. As a part of a more practical approach that is fitting of a first step toward the
266  derivation of an automated measure of social engagement in ASD, we decided to extract low-level
267  behavioral components from our video data as indicators of engagement in the interactions with the
268  robots. For this purpose, we utilized the 2D body tracking and facial expression data generated by
269  OpenPose [39].
270
271  Using the body tracking data, we derived three new features based on Laban Movement Analysis
272  (LMA), a method for describing and interpreting all types of human movement [53] used frequently
273  in a variety of fields including dance, acting, music, and physical therapy etc. LMA categorizes all
274  body movements into the categories of body effort, space and shape. Out of the four categories, effort
275  represents the dynamics of human movement and provides an insight into the subtle characteristics of
276  movements with respect to inner intention. This makes it an important feature to use in studies
277  involving the estimation of affect, intention and engagement states. Effort itself is classified into
278  space, weight and time, which are the three features that we incorporated in our current work. Space
279  represents the area taken up over the course of a movement, weight indicates the power or impact of
280  movement, and time conveys the speed of an action, including a sense of urgency or a lack thereof in
281  a movement. The equations [55,56] for each of these features are as shown in Table 1.
282

283 OpenPose generates 50 keypoints for skeletal tracking as described in [39]. In addition to the skeletal
284 data, we also recorded facial keypoints to incorporate the changes in a child's facial expressions in
285 our feature set. Figure 5 (taken from [54]) depicts these datapoints. While a total of 69 facial
286 keypoints is available, we only used the lip and eye keypoints shown on the right. Including the x and
287 y coordinates for each of the 34 facial keypoints and and the three Laban features derived from the
288 upper body skeletal keypoints created a total of 71 features in the dataset. A moving window of 1
289 second, i.e. 30 frames, was used to compute the Laban features in order to incorporate the sequential
290 nature of the movement data. A 1 second interval was chosen to capture meaningful, yet rapidly
291 changing movement patterns in response to the actions of the robot during the child-robot interaction.
292 The number of available datapoints per participant depended on the length of interaction of each
293 participant and ranged between 9300 and 30508 datapoints. Further details are listed in Table 3.
294
295 We initially attempted to use some derived features from the raw skeletal keypoints based on Laban
296 Movement Analysis (LMA), which is a method for describing and interpreting all types of human
297 movement [53], mainly used to represent the dynamics of human movement and provide an insight
298 into the subtle characteristics of movements with respect to inner intention. However, with some
299 preliminary tests, we found that the classifier trained on raw keypoints outperformed one trained on
300 derived features, and hence dropped the Laban features from the dataset. We also used Principal
301 Component Analysis (PCA) to reduce the dimensionality of the dataset.
302
303

Table 1. Equations for the derived Laban features adopted from [55,56].

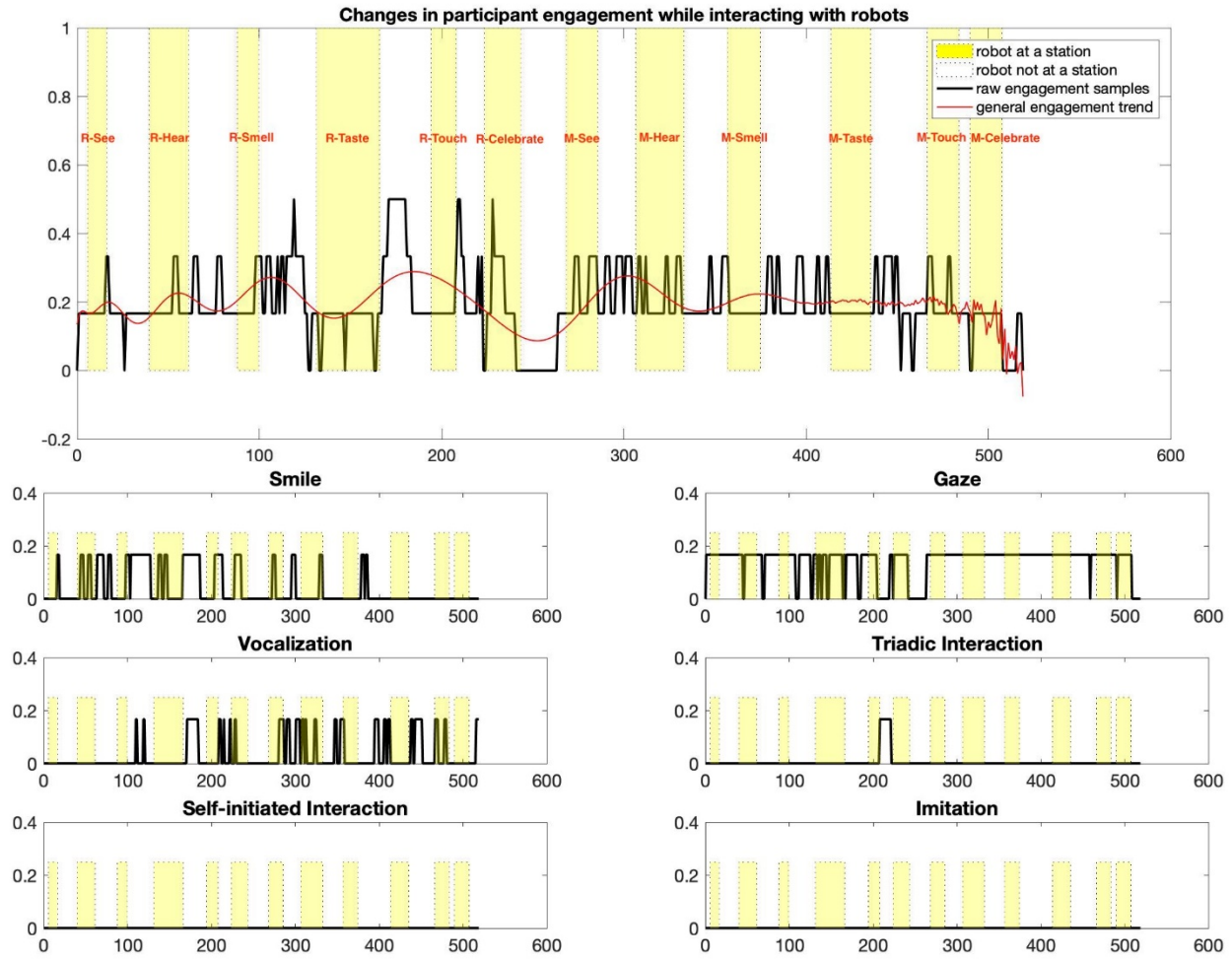| Feature | Equation |
|---|---|
| Space | $$Space = \left(0.5|\vec{a}||\vec{d}|\sin(\theta_1)\right) + \left(0.5|\vec{c}||\vec{b}|\sin(\theta_2)\right)$$ where $$\vec{a} = Position\ vector\ from\ left\ shoulder\ to\ left\ hand$$ $$\vec{b} = Position\ vector\ from\ right\ shoulder\ to\ left\ shoulder$$ $$\vec{c} = Position\ vector\ from\ right\ hand\ to\ right\ shoulder$$ $$\vec{d} = Position\ vector\ from\ left\ hand\ to\ right\ hand$$ $$\theta_1 = Angle\ between\ \vec{a}\ \&\ \vec{d}$$ $$\theta_2 = Angle\ between\ \vec{c}\ \&\ \vec{b}$$ |
| Weight | $$Weight = \sum_i \tau_i \omega_i(t)$$ where $$\tau_i = L^2 \omega_i^2 \sin(\theta) * mass$$ $$\omega_i = \frac{d\theta}{dt}$$ $$L = distance\ between\ joints$$ $$i = Joint\ Number$$ $$\dot{\omega}_i = Angular\ Velocity\ for\ Joint\ i$$ |
| Time | $$Time_i = \sum_i \dot{\omega}_i(t)$$ where $$i = Joint\ Number$$ $$\dot{\omega}_i = Angular\ Velocity\ for\ Joint\ i$$ |

304
305

Figure 4. Plots depicting changes in the overall engagement level of a child during an interaction, along with subplots of the target behaviors contributing to this engagement [20].



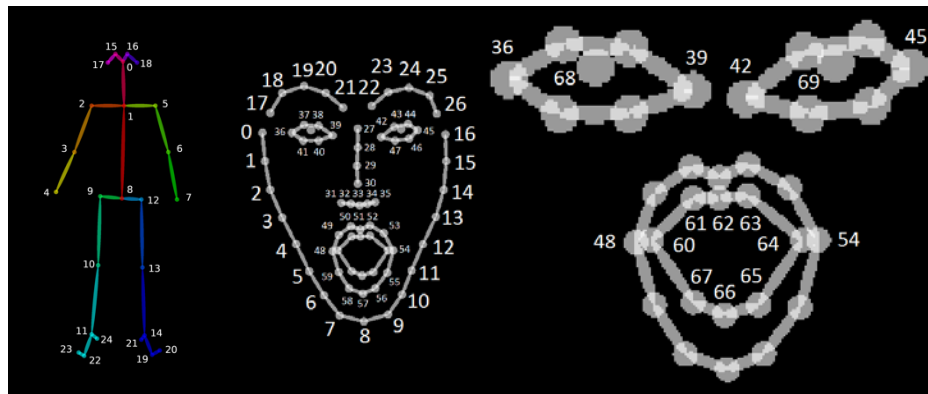Figure 5. Illustrations of the skeletal and facial keypoints extracted by OpenPose [54] (permission acquired from the author for using these image with citation).

306

## 7    Network Architecture

308    We used a multi-channel and multi-layer convolutional neural network (CNN) for this temporal
309    multi-label classification problem. The network was composed of two Conv1D layers to identify

310  temporal data patterns (with 5 channels with 64 and128 filters respectively and a kernel size of 3 with
311  20% dropout) and three dense layers for classification (kernel sizes 256, 256, and 7 (number of
312  output labels: value ranges of engagement level)). This is illustrated in more detail in Figure 6. A 10-
313  fold cross-validation (train/test split of 0.8/0.2) was used for every subject's individual dataset and
314  optimization was performed using the Adam optimizer.
315
316  The two Conv1D layers are meant to extract high-level features from the temporal data since the
317  dataset being used has a high input dimension and a relatively small number of datapoints. Since the
318  data have a non-linear structure, the first two dense layers are used to spread the feature dimension,
319  whereas the last one generates the output dimension. The dropout layers are used to avoid overfitting.
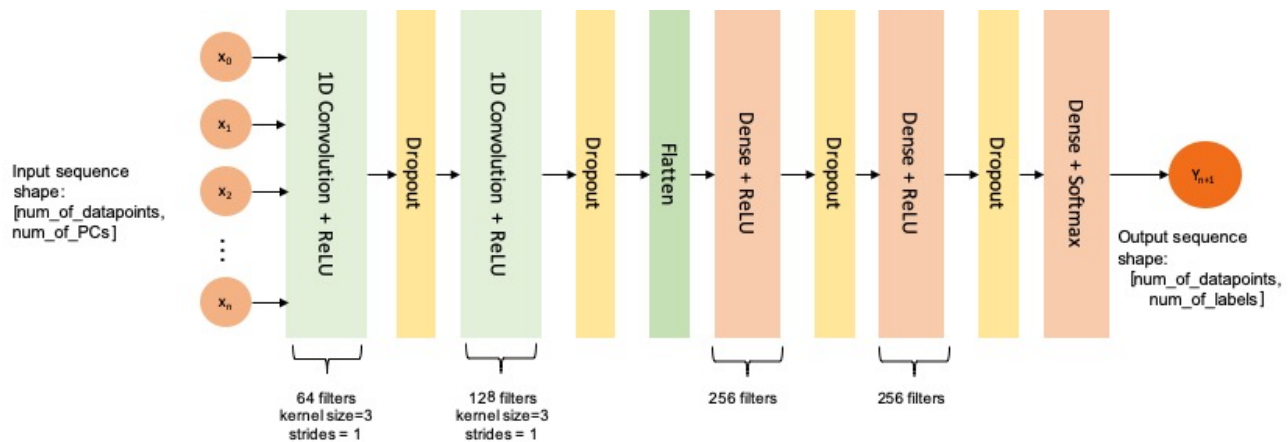320



321
322  Figure 4. Architecture of the CNN used for multi-label classification.

323  **8    User Study**
324  We conducted a user study with a total of 18 children, 13 TD and 5 with ASD between the ages of 4
325  and 12 years who participated in a one-time interaction with our robots within the setting of a sensory
326  maze game. The average age of the TD group was $7.07\pm2.56$ years and that of the ASD group was
327  $8.2\pm1.10$ years. The TD group consisted of 5 females and 8 males, whereas the ASD group was
328  composed of all male participants. These details are presented in Table 2.
329
330  The participants were allowed to participate for the entire course of the interaction as designed with
331  the two robots, one after another. The data presented in this study is for one-time interactions
332  between each subject and the robots. The length of the interaction for each participant is listed in
333  Table 2. The average TD interaction length was 464.92 seconds whereas that of the ASD group was
334  620 seconds. Individual engagement prediction models were generated for each participant and their
335  performances were evaluated.
336
337          Table 2. Demographic details of the subjects
338

339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355

| ID | Age | Gender | Group |
|----|-----|--------|-------|
| 1 | 10 | M | TD |
| 2 | 4 | F | TD |
| 3 | 5 | F | TD |
| 4 | 11 | F | TD |
| 5 | 9 | M | TD |
| 6 | 10 | F | TD |
| 7 | 9 | M | TD |
| 8 | 5 | M | TD |
| 9 | 5 | F | TD |
| 10 | 5 | M | TD |
| 11 | 5 | M | TD |
| 12 | 5 | M | TD |
| 13 | 9 | M | TD |
| 14 | 7 | M | ASD |
| 15 | 8 | M | ASD |
| 16 | 10 | M | ASD |
| 17 | 8 | M | ASD |
| 18 | 8 | M | ASD |

356 **9    Results**
357 Table 3 presents the detailed results produced by training, validation and testing our network for
358 every subject in the study. The length of interaction is important and provides an insight into the
359 number of video frames, and hence, the datapoints that would be available to the network. The
360 datapoint count is also affected by the processing performed by OpenPose, which can drop some
361 frames where processing could not be completed. This is particularly evident in the case of
362 participant 6 and 12, where the number of available datapoints are far fewer than expected.
363
364 Before presenting the results, it must be highlighted that the metrics shown in this work are all
365 weighted metrics, so as to address the impact of the imbalance in engagement level samples within
366 the dataset. The network has an average accuracy of 0.7985 for the TD group and 0.8061for the ASD
367 group in the training stage. For the test data, the performance remains steady with an average
368 accuracy of 0.7767 for the ASD group and 0.7918 for the TD group.
369
370 Figure 7 depicts the accuracy and loss plots for training and validation data for a participant from
371 each group illustrating the changes in accuracy with respect to the number of epochs. Figure 8 shows
372 the timeseries plots of the changing engagement states for the participants. The red line shows the
373 true engagement as determined by the annotations [21]. Predictions made by the network are marked
374 in blue. Since the dataset was randomly partitioned into test and training data, the predictions on the
375 test set appear as a scatter plot.
376
377 Table 3. Performance metrics for the individual classifiers (TD Group: ID1 – ID13, ASD Group: ID14 – ID18)

| ID | Interaction length (s) | No. of datapoints (frames) | Train | | Validation | | Test |
|----|----|----|----|----|----|----|----|
| | | | Accuracy | Loss | Accuracy | Loss | Accuracy |
| 1 | 315 | 9444 | 0.8101 | 0.5028 | 0.7790 | 0.6681 | 0.7946 |
| 2 | 519 | 15357 | 0.6499 | 0.7278 | 0.6398 | 0.7797 | 0.6393 |
| 3 | 540 | 16412 | 0.6703 | 0.8723 | 0.6407 | 1.0095 | 0.6526 |
| 4 | 658 | 10933 | 0.8302 | 0.4189 | 0.8131 | 0.4923 | 0.8240 |
| 5 | 797 | 22996 | 0.9255 | 0.1903 | 0.9198 | 0.2484 | 0.9159 |
| 6 | 696 | 9300 | 0.9200 | 0.2850 | 0.8925 | 0.3856 | 0.9124 |

| 7 | 316 | 9388 | 0.7821 | 0.5423 | 0.7417 | 0.7946 | 0.7338 |
| 8 | 457 | 13725 | 0.7561 | 0.6065 | 0.7418 | 0.6796 | 0.7483 |
| 9 | 574 | 10463 | 0.6671 | 0.8486 | 0.6535 | 0.9333 | 0.6364 |
| 10 | 780 | 16627 | 0.9104 | 0.2253 | 0.8831 | 0.3907 | 0.8698 |
| 11 | 726 | 12726 | 0.8390 | 0.3843 | 0.8303 | 0.4039 | 0.8283 |
| 12 | 685 | 9723 | 0.8118 | 0.5162 | 0.7715 | 0.6980 | 0.7720 |
| 13 | 540 | 12879 | 0.8084 | 0.4296 | 0.7812 | 0.5858 | 0.7702 |
| 14 | 517 | 15502 | 0.8163 | 0.4417 | 0.7952 | 0.5621 | 0.7907 |
| 15 | 578 | 14624 | 0.9204 | 0.2276 | 0.8923 | 0.3390 | 0.9108 |
| 16 | 679 | 15950 | 0.6810 | 0.7582 | 0.6501 | 0.9095 | 0.6398 |
| 17 | 610 | 16401 | 0.8306 | 0.3946 | 0.8232 | 0.4923 | 0.8366 |
| 18 | 1058 | 30508 | 0.7822 | 0.5467 | 0.7759 | 0.6323 | 0.7812 |

Table 4. Average metrics to compare classifier performance

| ID | Average interaction length (s) | Train | | Validation | | Test |
| | | Accuracy | Loss | Accuracy | Loss | Accuracy |
|---|---|---|---|---|---|---|
| TD | 584.8 | 0.7985 | 0.5038 | 0.7760 | 0.6207 | 0.7767 |
| ASD | 688.4 | 0.8061 | 0.4738 | 0.7873 | 0.5870 | 0.7918 |

In addition to the individual models described above, we also trained a group model for each of the two groups by using all the datapoints collected from the participants from each group. The ASD classifier was able to achieve a training accuracy of 0.6389 and a test accuracy of 0.6524, while the TD classifier achieved a slightly higher training accuracy of 0.6733 and a test accuracy of 0.6803. The slightly superior performance of the classifiers on the test data as opposed to the training data can be attributed to the use of regularization techniques used when constructing the classifier structure, in this case, the Dropout layers, which are only applied during the training phase.

We also trained a combined classifier on the data collected from all the participants. This model underperformed slightly compared to the group-specific classifiers, indicating that a group-specific classifier may be better suited for generalization to all participants within the group rather than a single classifier for all participants (Table 5). Accuracy and loss plots for the training and validating processes for all three grouped conditions are shown in Figure 9.

Table 5. Performance metrics for group classifiers.

| Classifier | Train | | Validation | | Test |
| | Accuracy | Loss | Accuracy | Loss | Accuracy |
|---|---|---|---|---|---|
| TD | 0.6733 | 0.8472 | 0.6800 | 0.8263 | 0.6803 |
| ASD | 0.6389 | 0.9320 | 0.6512 | 0.8858 | 0.6524 |
| Combined | 0.6733 | 0.8472 | 0.6800 | 0.8263 | 0.6803 |

## 10   Comparison with Other Machine Learning Classifiers

A number of standard Machine Learning (ML) classifiers were also trained for all the scenarios described above as a way to situate the performance of the CNN, which included Support Vector Classification (SVC), Random Forest (RF), Decision Trees (DT) and K-Nearest Neighbors (KNN). The reported metrics were also averaged across all participants to compare the overall performance of the classifiers. As before, each classifier was trained and tested on entire group datasets to compare performance as a generalized group classifier. These results are shown in Table 6.

405
406

Table 6. Performance metrics for all classifiers under individual and group conditions.

| ID | CNN | | SVC | | RF | | DT | | KNN | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 |
| 1 | 0.79 | 0.77 | 0.77 | 0.72 | 0.80 | 0.78 | 0.77 | 0.75 | 0.81 | 0.79 |
| 2 | 0.64 | 0.62 | 0.58 | 0.55 | 0.75 | 0.75 | 0.65 | 0.64 | 0.72 | 0.71 |
| 3 | 0.65 | 0.59 | 0.66 | 0.55 | 0.67 | 0.61 | 0.65 | 0.58 | 0.67 | 0.61 |
| 4 | 0.82 | 0.79 | 0.82 | 0.76 | 0.83 | 0.81 | 0.82 | 0.79 | 0.83 | 0.81 |
| 5 | 0.92 | 0.91 | 0.89 | 0.87 | 0.93 | 0.92 | 0.90 | 0.89 | 0.93 | 0.93 |
| 6 | 0.91 | 0.89 | 0.92 | 0.90 | 0.90 | 0.89 | 0.91 | 0.89 | 0.92 | 0.90 |
| 7 | 0.73 | 0.73 | 0.61 | 0.59 | 0.80 | 0.80 | 0.72 | 0.71 | 0.80 | 0.80 |
| 8 | 0.75 | 0.74 | 0.51 | 0.47 | 0.82 | 0.82 | 0.66 | 0.66 | 0.82 | 0.81 |
| 9 | 0.64 | 0.57 | 0.63 | 0.56 | 0.65 | 0.60 | 0.63 | 0.57 | 0.67 | 0.61 |
| 10 | 0.87 | 0.87 | 0.79 | 0.77 | 0.88 | 0.87 | 0.82 | 0.82 | 0.85 | 0.85 |
| 11 | 0.77 | 0.76 | 0.69 | 0.65 | 0.78 | 0.77 | 0.72 | 0.71 | 0.76 | 0.74 |
| 12 | 0.83 | 0.78 | 0.81 | 0.74 | 0.84 | 0.81 | 0.82 | 0.79 | 0.84 | 0.80 |
| 13 | 0.77 | 0.77 | 0.73 | 0.69 | 0.79 | 0.80 | 0.77 | 0.77 | 0.79 | 0.80 |
| 14 | 0.79 | 0.79 | 0.70 | 0.69 | 0.82 | 0.81 | 0.73 | 0.73 | 0.81 | 0.81 |
| 15 | 0.91 | 0.90 | 0.87 | 0.83 | 0.92 | 0.90 | 0.90 | 0.88 | 0.92 | 0.91 |
| 16 | 0.64 | 0.62 | 0.61 | 0.57 | 0.67 | 0.65 | 0.62 | 0.60 | 0.68 | 0.66 |
| 17 | 0.84 | 0.84 | 0.70 | 0.69 | 0.88 | 0.88 | 0.76 | 0.75 | 0.84 | 0.84 |
| 18 | 0.78 | 0.78 | 0.63 | 0.60 | 0.79 | 0.78 | 0.61 | 0.58 | 0.78 | 0.78 |
| Average | 0.78 | 0.76 | 0.72 | 0.68 | 0.81 | 0.79 | 0.75 | 0.73 | 0.80 | 0.79 |
| TD | 0.68 | 0.65 | 0.63 | 0.58 | 0.74 | 0.74 | 0.64 | 0.61 | 0.74 | 0.73 |
| ASD | 0.72 | 0.71 | 0.60 | 0.58 | 0.77 | 0.76 | 0.61 | 0.60 | 0.76 | 0.76 |
| Combined | 0.65 | 0.62 | 0.59 | 0.54 | 0.74 | 0.71 | 0.60 | 0.56 | 0.71 | 0.71 |

407
408 After averaging over the metrics for all participants, RF is seen to have the best performance
409 followed by KNN and CNN respectively. A similar trend is seen for grouped classifiers, where RF
410 once again outperforms all other classifiers in terms of both the accuracy and the F1 score, followed
411 again by KNN and CNN respectively. All classifier performances drop slightly when data from the
412 two groups are combined, suggesting that a single classifier may not be as useful for generalization
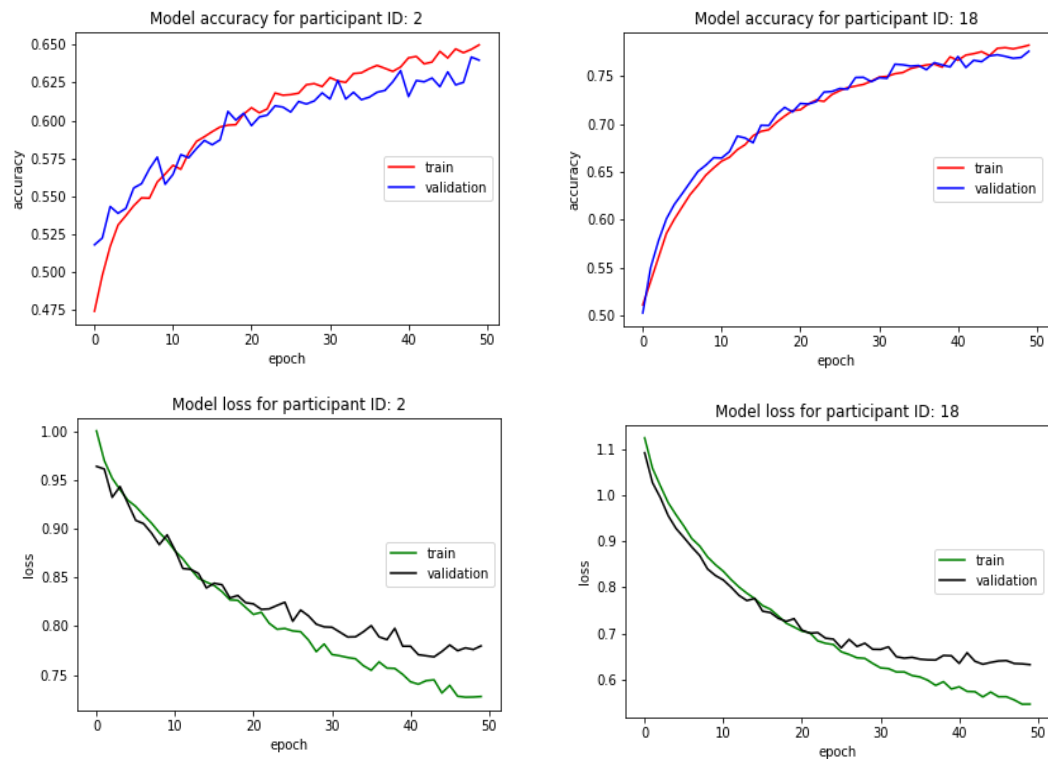413 as a group-specific classifier.

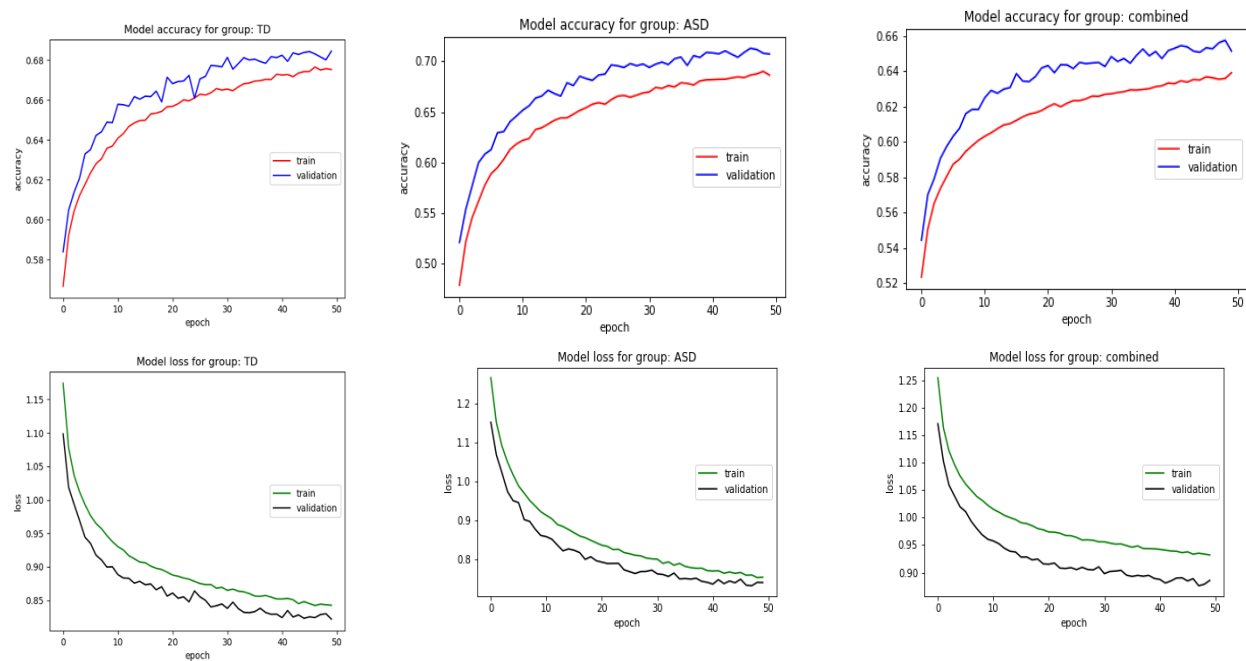Figure 7. Classifier accuracy and loss with respect to the number of epochs for two different participants.



Figure 8. Plots showing the ground truth labels in red and the classifier predictions in blue.

414

415     **11     Discussion**

416    In this work, we propose the use of a Deep Learning Convolutional Neural Network to model and
417    predict child social engagement as a part of our larger goal to personalize child-robot interactions.
418    We utilized key social behaviors as indicators of engagement in an interaction, which formed the
419    criterion for the human-generated labels that serves as the ground truth for this engagement
420    classification approach.
421
422    We found that the proposed CNN was able to achieve a performance that was comparable to the
423    highest performing classical ML approaches in this work. The RF and KNN classifiers only slightly
424    outperform the CNN in the case of both individual classifiers and grouped classifiers. The individual
425    classifiers serve as personalized engagement prediction networks for the unique behavioral
426    expressions of each individual participant, whereas the grouped classifiers were used to evaluate the
427    potential for a single classifier to generalize the learnt patterns to all the participants within a group.
428
429    On the individual level, the CNN was able to attain a best case accuracy of 0.92 (participant 5) and a
430    worst case accuracy of 0.64 (participant 2). On the other hand, the RF classifier reached a highest
431    accuracy of 0.93 (participant 5) and lowest accuracy of 0.65 (participant 9). For the averaged metrics
432    as well as the grouped metrics, the RF accuracy is no more than 2% higher than that of the CNN.
433    The individual ASD and TD classifiers were generally found to achieve a higher accuracy than the
434    single classifier trained on data from all the participants. This points the possibility of a generalized
435    group classifier that can be used effectively to classify social engagement for all the children in each
436    group while providing a high level of personalization in the interaction.
437
438    The CNN is a complex structure with a large number of tunable parameters that generally requires
439    much larger datasets to fully exploit the potential of deep networks. Given the number of input
440    features, the number of output classes and the size of the dataset (generated by single session child-
441    robot interactions only) used in this study, the CNN was able to achieve a performance comparable to
442    simpler ML classifiers but not exceed them. We anticipate that as we continue to collect interaction
443    data from additional participants for a long-term study involving multiple sessions, the proposed deep
444    learning network will likely become a more suitable choice for social engagement classification.
445    It must also be pointed out that in terms of deployment to a robotic platform, a CNN may also be a
446    more suitable option since the traditional algorithms require expensive resources when deployed to
447    mobile platform in real-world applications, whereas deep learning algorithms can fully take
448    advantage of the scalable computing platforms with GPUs that have low-cost modules (like the
449    NVidia Jetson Nano) while retaining the capacity to handle much larger datasets.
450
451    The current work is limited in that it only utilizes single session data for each participant based on
452    which the classifiers are trained. Classifier performance is likely to improve as subsequent sessions
453    are conducted and larger datasets are collected. Another limitation of this work is that the datasets for
454    the two groups are unbalanced, with 13 participants in the TD group and only 5 in the ASD group
455    generating much larger training dataset for the TD classifier than ASD. Conducting long-term studies
456    with a population such as ASD remains a considerable challenge for all researchers in the field and
457    explains the lack of open multi-modal datasets to benefit the ASD research community.
458
459    Since our focus in this work was to evaluate social engagement in a naturalistic interaction setting,
460    the video recordings of the sessions mainly focused on the participant but also included other
461    members of the research team and/or parent in several segments of the videos as the child moved
462    around the room to interact with the robots. OpenPose was chosen to process the movements of the
463    participants particularly because it offers a feature to track multiple persons by assigning each a fixed
464    ID. In practice, however, this ID assignment was found to lack reliability, which we discovered by

465    visualizing the participant's skeletal tracking data. In addition, we also found that the number of
466    frames in the input video and the number of frames generated as output by OpenPose were often
467    inconsistent, contributing to the loss of data.
468
469    It would be interesting to see how the classifier performance changes over long-term interactions
470    between the children and robots. Child engagement is likely to vary with continued exposure to the
471    robots and inclusion of additional temporal features in the dataset may become important. We also
472    aim to incorporate additional modalities to our dataset, including physiological signals like heart rate,
473    electrodermal activity, body temperature and blood pressure, as well as audio features. For this
474    complex feature set, we foresee a deep learning network to be a more suitable classifier choice
475    capable of identifying patterns relating to different levels of social engagement in children.

476    **12    Conclusion**
477    In this paper, we presented a multi-label convolutional neural network classifier to formulate an
478    automated measure of social engagement for children. To provide a personalized metric that is the
479    best representation of the unique expression of emotion, interest and intention of each individual, we
480    trained a separate classifier for each subject and then evaluated its performance. We designed the
481    study to ensure the participants were not restricted in their movements at all in order to closely mimic
482    naturalistic interactions in the real world. The use of this setting increases the complexity of data
483    collection and analysis but enables the generalization of the presented analysis techniques to other
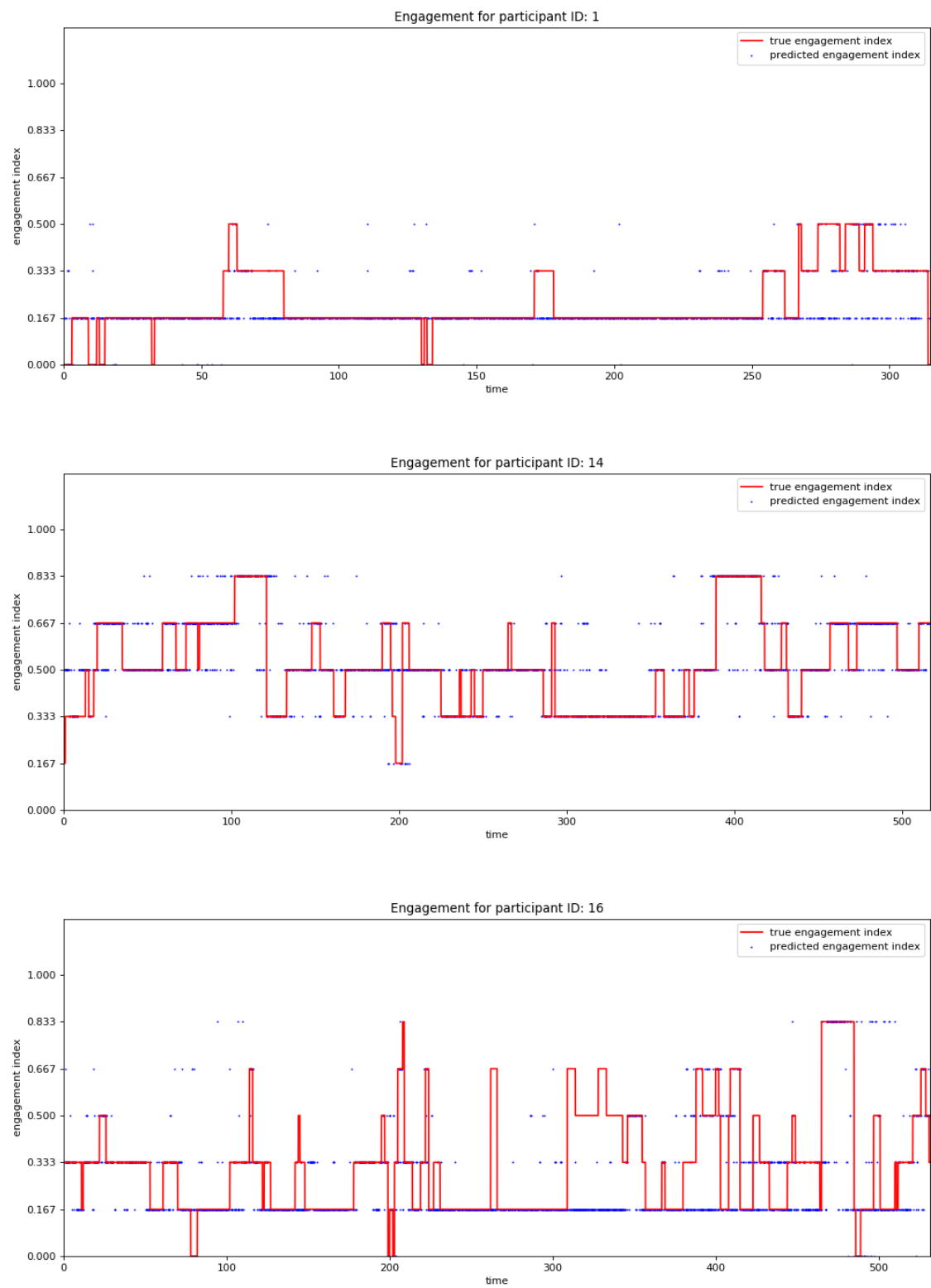484    interaction scenarios and populations, which sets this work apart from other research studies in this
485    domain.
486

Figure 9. Classifier accuracy and loss for training and test datasets for three grouped conditions.

487
488

489                                                                      **References**
490    [1] Kishida, Y., & Kemp, C. (2006). Measuring child engagement in inclusive early childhood
491    settings: Implications for practice. Australasian Journal of Early Childhood, 31(2), 14-19.
492    [2] American Psychiatric Association. (2013). Diagnostic and statistical manual of mental disorders
493    (DSM-5®). American Psychiatric Pub.
494    [3] Ollendick, T. H., Weist, M. D., Borden, M. C., & Greene, R. W. (1992). Sociometric status and
495    academic, behavioral, and psychological adjustment: a five-year longitudinal study. Journal of
496    consulting and clinical psychology, 60(1), 80.
497    [4] Tantam, D. (2000). Psychological disorder in adolescents and adults with Asperger syndrome.
498    Autism, 4(1), 47-62.
499    [5] Bellini, S. (2006). The development of social anxiety in adolescents with autism spectrum
500    disorders. Focus on autism and other developmental disabilities, 21(3), 138-145.
501    [6] Dubey, I., Ropar, D., & de C Hamilton, A. F. (2015). Measuring the value of social engagement
502    in adults with and without autism. Molecular autism, 6(1), 35.
503    [7] Sanefuji, Wakako, and Hidehiro Ohgami. "Imitative behaviors facilitate communicative gaze in
504    children with autism." Infant Mental Health Journal 32, no. 1 (2011): 134-142.
505    [8] Ingersoll, Brooke. "The social role of imitation in autism: Implications for the treatment of
506    imitation deficits." Infants & Young Children 21, no. 2 (2008): 107-119.
507    [9] Tiegerman, Ellenmorris, and Louis H. Primavera. "Imitating the autistic child: Facilitating
508    communicative gaze behavior." Journal of autism and developmental disorders 14, no. 1 (1984): 27-
509    38.
510    [10] Slaughter, Virginia, and Su Sen Ong. "Social behaviors increase more when children with ASD
511    are imitated by their mother vs. an unfamiliar adult." Autism Research 7, no. 5 (2014): 582-589.
512    [11] Tiegerman, Ellenmorris, and Louis Primavera. "Object manipulation: An interactional strategy
513    with autistic children." Journal of Autism and Developmental Disorders 11, no. 4 (1982): 427-438.
514    [12] Katagiri, Masatoshi, Naoko Inada, and Yoko Kamio. "Mirroring effect in 2-and 3-year-olds with
515    autism spectrum disorder." Research in Autism Spectrum Disorders 4, no. 3 (2010): 474-478.
516    [13] Contaldo, Annarita, Costanza Colombi, Antonio Narzisi, and Filippo Muratori. "The social
517    effect of "being imitated" in children with autism spectrum disorder." Frontiers in psychology 7
518    (2016): 726.
519    [14] Stanton, Cady M., Peter H. Kahn Jr, Rachel L. Severson, Jolina H. Ruckert, and Brian T. Gill.
520    "Robotic animals might aid in the social development of children with autism." In Proceedings of the
521    3rd ACM/IEEE international conference on Human robot interaction, pp. 271-278. ACM, 2008.
522    [15] Tapus, Adriana, Andreea Peca, Amir Aly, Cristina Pop, Lavinia Jisa, Sebastian Pintea, Alina S.
523    Rusu, and Daniel O. David. "Children with autism social engagement in interaction with Nao, an
524    imitative robot: A series of single case experiments." Interaction studies 13, no. 3 (2012): 315-347.
525    [16] Wimpory, Dawn C., R. Peter Hobson, J. Mark G. Williams, and Susan Nash. "Are infants with
526    autism socially engaged? A study of recent retrospective parental reports." Journal of Autism and
527    Developmental Disorders 30, no. 6 (2000): 525-536.
528    [17] Nadel, Jacqueline. "Imitation and imitation recognition: Functional use in preverbal infants and
529    nonverbal children with autism." The imitative mind: Development, evolution, and brain bases 4262
530    (2002).
531    [18] Hernandez, J., Riobo, I., Rozga, A., Abowd, G. D., & Picard, R. W. (2014). Using electrodermal
532    activity to recognize ease of engagement in children during social interactions. In Proceedings of the
533    2014 ACM International Joint Conference on Pervasive and Ubiquitous Computing (pp. 307-317).
534    ACM.
535    [19] Lahiri, U., Bekele, E., Dohrmann, E., Warren, Z., & Sarkar, N. (2012). Design of a virtual
536    reality based adaptive response technology for children with autism. IEEE Transactions on Neural
537    Systems and Rehabilitation Engineering, 21(1), 55-64.

538 [20] Kushki, A., Andrews, A. J., Power, S. D., King, G., & Chau, T. (2012). Classification of activity
539 engagement in individuals with severe physical disabilities using signals of the peripheral nervous
540 system. PloS one, 7(2), e30373.
541 [21] Javed, H., Burns, R., Jeon, M., Howard, A. M., and Park, C. H., An Interactive Framework to
542 Facilitate Sensory Experiences for Children with ASD, 2019 (Accepted), ACM Transactions on
543 Human-Robot Interaction. arXiv preprint arXiv:1901.00885.
544 [22] Bevill, R., Cowin, S., & Park, C. H. (2017). Motion Learning for Emotional Interaction and
545 Imitation Of Children With Autism Spectrum Disorder.
546 [23] Javed, H., Jeon, M., Howard, A., & Park, C. H. (2018). Robot-assisted socio-emotional
547 intervention framework for children with Autism Spectrum disorder. In Companion of the 2018
548 ACM/IEEE International Conference on Human-Robot Interaction (pp. 131-132). ACM.
549 [24] Rajagopalan, S. S., Murthy, O. R., Goecke, R., & Rozga, A. (2015, May). Play with me—
550 Measuring a child's engagement in a social interaction. In 2015 11th IEEE International Conference
551 and Workshops on Automatic Face and Gesture Recognition (FG) (Vol. 1, pp. 1-8). IEEE.
552 [25] Rehg, J., Abowd, G., Rozga, A., Romero, M., Clements, M., Sclaroff, S., ... & Rao, H. (2013).
553 Decoding children's social behavior. In Proceedings of the IEEE conference on computer vision and
554 pattern recognition (pp. 3414-3421).
555 [26] Gupta, R., Bone, D., Lee, S., & Narayanan, S. (2016). Analysis of engagement behavior in
556 children during dyadic interactions using prosodic cues. Computer speech & language, 37, 47-66.
557 [27] Lala, D., Inoue, K., Milhorat, P., & Kawahara, T. (2017). Detection of social signals for
558 recognizing engagement in human-robot interaction. arXiv preprint arXiv:1709.10257.
559 [28] Castellano, G., Pereira, A., Leite, I., Paiva, A., & McOwan, P. W. (2009, November). Detecting
560 user engagement with a robot companion using task and social interaction-based features.
561 [29] Kim, J., Truong, K. P., & Evers, V. (2016, September). Automatic detection of children's
562 engagement using non-verbal features and ordinal learning. In WOCCI (pp. 29-34).
563 [30] Parekh, V., Foong, P. S., Zhao, S., & Subramanian, R. (2018, March). AVEID: Automatic
564 Video System for Measuring Engagement in Dementia. In 23rd International Conference on
565 Intelligent User Interfaces (pp. 409-413). ACM.
566 [31] Oertel, C., & Salvi, G. (2013, December). A gaze-based method for relating group involvement
567 to individual engagement in multimodal multiparty dialogue. In Proceedings of the 15th ACM on
568 International conference on multimodal interaction (pp. 99-106). ACM.
569 [32] Anzalone, S. M., Boucenna, S., Ivaldi, S., & Chetouani, M. (2015). Evaluating the engagement
570 with social robots. International Journal of Social Robotics, 7(4), 465-478.
571 [33] Rudovic, O., Lee, J., Dai, M., Schuller, B., & Picard, R. W. (2018). Personalized machine
572 learning for robot perception of affect and engagement in autism therapy. Science Robotics, 3, 19.
573 [34] Rudovic, O., Lee, J., Mascarell-Maricic, L., Schuller, B. W., & Picard, R. W. (2017). Measuring
574 engagement in robot-assisted autism therapy: A cross-cultural study. Frontiers in Robotics and AI, 4,
575 36.
576 [35] Dautenhahn, K., & Werry, I. (2004). Towards interactive robots in autism therapy: Background,
577 motivation and challenges. Pragmatics & Cognition, 12(1), 1-35.
578 [36] Diehl, J. J., Schmitt, L. M., Villano, M., & Crowell, C. R. (2012). The clinical use of robots for
579 individuals with autism spectrum disorders: A critical review. Research in autism spectrum disorders,
580 6(1), 249-262.
581 [37] Cabibihan, J. J., Javed, H., Ang, M., & Aljunied, S. M. (2013). Why robots? A survey on the
582 roles and benefits of social robots in the therapy of children with autism. International journal of
583 social robotics, 5(4), 593-618.
584 [38] Scassellati, B. (2007). How social robots will help us to diagnose, treat, and understand autism.
585 In Robotics research (pp. 552-563). Springer, Berlin, Heidelberg.

586    [39] Cao, Z., Hidalgo, G., Simon, T., Wei, S. E., & Sheikh, Y. (2018). OpenPose: realtime multi-
587    person 2D pose estimation using Part Affinity Fields. arXiv preprint arXiv:1812.08008.
588    [40] Nadel, Jacqueline. "Imitation and imitation recognition: Functional use in preverbal infants and
589    nonverbal children with autism." The imitative mind: Development, evolution, and brain bases 4262
590    (2002).
591    [41] Dubey, Indu, Danielle Ropar, and Antonia F. de C Hamilton. "Measuring the value of social
592    engagement in adults with and without autism." Molecular autism 6, no. 1 (2015): 35.
593    [42] Sanefuji, Wakako, and Hidehiro Ohgami. "Imitative behaviors facilitate communicative gaze in
594    children with autism." Infant Mental Health Journal 32, no. 1 (2011): 134-142.
595    [43] Ingersoll, Brooke. "The social role of imitation in autism: Implications for the treatment of
596    imitation deficits." Infants & Young Children 21, no. 2 (2008): 107-119.
597    [44] Tiegerman, Ellenmorris, and Louis H. Primavera. "Imitating the autistic child: Facilitating
598    communicative gaze behavior." Journal of autism and developmental disorders 14, no. 1 (1984): 27-
599    38.
600    [45] Slaughter, Virginia, and Su Sen Ong. "Social behaviors increase more when children with ASD
601    are imitated by their mother vs. an unfamiliar adult." Autism Research 7, no. 5 (2014): 582-589.
602    [46] Tiegerman, Ellenmorris, and Louis Primavera. "Object manipulation: An interactional strategy
603    with autistic children." Journal of Autism and Developmental Disorders 11, no. 4 (1982): 427-438.
604    [47] Katagiri, Masatoshi, Naoko Inada, and Yoko Kamio. "Mirroring effect in 2-and 3-year-olds with
605    autism spectrum disorder." Research in Autism Spectrum Disorders 4, no. 3 (2010): 474-478.
606    [48] Contaldo, Annarita, Costanza Colombi, Antonio Narzisi, and Filippo Muratori. "The social
607    effect of "being imitated" in children with autism spectrum disorder." Frontiers in psychology 7
608    (2016): 726.
609    [49] Stanton, Cady M., Peter H. Kahn Jr, Rachel L. Severson, Jolina H. Ruckert, and Brian T. Gill.
610    "Robotic animals might aid in the social development of children with autism." In Proceedings of the
611    3rd ACM/IEEE international conference on Human robot interaction, pp. 271-278. ACM, 2008.
612    [50] Tapus, Adriana, Andreea Peca, Amir Aly, Cristina Pop, Lavinia Jisa, Sebastian Pintea, Alina S.
613    Rusu, and Daniel O. David. "Children with autism social engagement in interaction with Nao, an
614    imitative robot: A series of single case experiments." Interaction studies 13, no. 3 (2012): 315-347.
615    [51] Wimpory, Dawn C., R. Peter Hobson, J. Mark G. Williams, and Susan Nash. "Are infants with
616    autism socially engaged? A study of recent retrospective parental reports." Journal of Autism and
617    Developmental Disorders 30, no. 6 (2000): 525-536.
618    [52] Friard, Olivier, and Marco Gamba. "BORIS: a free, versatile open-source event-logging
619    software for video/audio coding and live observations." Methods in Ecology and Evolution 7, no. 11
620    (2016): 1325-1330.
621    [53] Groff, E. (1995). Laban movement analysis: Charting the ineffable domain of human movement.
622    Journal of Physical Education, Recreation & Dance, 66(2), 27-30.
623    [54] CMU-Perceptual-Computing-Lab (2019). OpenPose Demo – Output. Retrieved from
624    https://github.com/CMU-Perceptual-Computing-Lab/openpose/blob/master/doc/output.md.
625    [55] Masuda, M., Kato, S., & Itoh, H. (2009, December). Emotion detection from body motion of
626    human form robot based on laban movement analysis. In International Conference on Principles and
627    Practice of Multi-Agent Systems (pp. 322-334). Springer, Berlin, Heidelberg.
628    [56] Wakayama, Y., Okajima, S., Takano, S., & Okada, Y. (2010, September). IEC-based motion
629    retrieval system using Laban movement analysis. In International Conference on Knowledge-Based
630    and Intelligent Information and Engineering Systems (pp. 251-260). Springer, Berlin, Heidelberg.
631