

# In-the-Wild Learning from Demonstration for Therapies for Autism Spectrum Disorder

Ala'aldin Hijaz<sup>1</sup>, Jessica Korneder<sup>2</sup>, and Wing-Yue Geoffrey Louie<sup>1</sup>, *Member, IEEE*

**Abstract**—Current studies have demonstrated that Socially Assistive Robots (SARs) delivering Applied Behavior Analysis (ABA) based interventions can teach individuals with Autism Spectrum Disorder (ASD) valuable social, emotional, communication and academic skills. These robot-mediated interventions (RMIs) are typically delivered via teleoperation, which places additional or similar workloads on therapists as administering interventions directly. The autonomous delivery of ABA therapies to individuals with ASD by a robot could significantly reduce workload and improve the usability as well as acceptance of this technology. However, pre-programming the autonomy of a SAR with a limited set of interventions is not sufficient for clinical practice due to the rapidly changing and different learning needs of individuals with ASD. In order to be applicable in clinical settings, therapists must be capable of customizing and personalizing interventions to the needs of each individual. Towards this goal, in this paper we present the initial development and deployment of a proof-of-concept Learning from Demonstration (LfD) system in-the-wild to learn the verbal behavior of therapists during the delivery of an ABA-based intervention to children with ASD. We also present preliminary data on the results of a policy trained on data collected from demonstrations provided during this in-the-wild deployment of our LfD system.

## I. INTRODUCTION

The Center for Disease Control estimates that 1 in every 54 individuals are diagnosed with Autism Spectral Disorder (ASD) [1]. ASD is a condition that affects an individual's behavior, including their social, emotional and communication skills [1]. Early childhood intervention for an individual with ASD has been shown to positively impact his/her outcomes [2]. Applied Behavioral Analysis (ABA) therapy is an evidenced-based practice for delivering therapeutic interventions to teach individuals with ASD new valuable cognitive, emotional and behavioral skills [3]. More recently, there has been interest in utilizing Socially Assistive Robots (SARs) to deliver ABA-based therapies to individuals with ASD because Robot Mediated Interventions (RMIs) which utilize the principles of ABA have been effective in helping children with ASD learn new beneficial skills [4]–[7].

The majority of existing research in RMIs have an ABA therapist teleoperating a SAR to deliver an ABA session whether by pre-scripted social behaviors [7]–[10] or a combination of motion tracking and real-time streaming of the teleoperator's voice [11]. Teleoperation of a robot by a therapist during an intervention requires significant cognitive load

while also affecting the treatment integrity of the intervention [11]. Meanwhile, therapists would still be responsible for delivering the interventions which reduces the robot's overall utility as it would require the same or additional workload as when therapists directly deliver an intervention [12]. Hence, the autonomous delivery of ABA therapies by SARs can potentially be a more effective, efficient, usable, and acceptable approach because therapists would only need to monitor the SAR behaviors and intervene only when necessary [13].

In order for SARs to autonomously deliver an intervention they must be capable of adapting to the non-deterministic human-robot interactions within an ABA intervention and rapidly changing as well as different learning needs of children with ASD. These non-deterministic human-robot interactions includes the different responses of a child during an intervention, ambient and background noise, unexpected behaviors from the child and differences in responses between different children. Hence, a SAR must be capable of learning a policy which maps the current state of an intervention to the appropriate robot behavior to generalize to different human-robot interaction scenarios within an ABA intervention session. Furthermore, children with ASD have rapidly changing and different learning needs due to different behaviors, preferences, and traits of learners with ASD [14]. Hence, therapists desire the ability to personalize interventions delivered by SARs and do not believe they can be pre-programmed with a set of interventions to meet all the learning needs of individuals with ASD [14]. Learning from demonstration is an approach for non-experts to teach a robot a policy to autonomously complete a task [15]. LfD refers to the process of transferring a new skill to a robot by having a human user demonstrate the skill to a robot [16].

In this paper, we present a proof-of-concept LfD system capable of learning the verbal behavior of an ABA therapist to enable a SAR to autonomously deliver an ABA intervention. First, the learning data was captured during a teleoperated RMI where therapists remotely controlled the robot. The robot was teleoperated using motion tracking and voice streaming via a Virtual Reality (VR) interface where the verbal responses of the therapists and the children were recorded concurrently via microphones [11]. The verbal utterances of the therapists during the intervention are clustered using unsupervised learning. The clusters are then used as labels for the children's utterances. Namely, the system takes as input the raw audio spectrogram signal from the child's verbal behavior and outputs the corresponding ABA therapist verbal response accordingly. We designed an entirely unsupervised learning method that allows the SAR

This work was supported by the National Science Foundation grant #1948224

<sup>1</sup>Intelligent Robotics Laboratory, Oakland University, Michigan, USA (e-mail: louie@oakland.edu)

<sup>2</sup>Applied Behavior Analysis Clinic, Oakland University, Michigan, USA

to automatically learn the verbal structure of an ABA intervention session from demonstrations provided by a therapist in-the-wild.

## II. RELATED WORK

There is currently a handful of research focusing on having robots learn from human demonstrations social human-robot interaction (HRI) tasks [17]–[19]

In [17], a deep reinforcement learning-based system was developed for robots to learn to deliver interventions to individuals with ASD from demonstrations provided via teleoperation control of a SAR. Namely, the researchers teleoperated a robot using three pre-scripted robot behaviors during the delivery of mock greeting interventions to healthy adult participants. A DeepQ network was then trained using the observed RGB video stream from the robot’s camera and audio spectrogram from the robot’s microphone as the inputs to the model and the teleoperated robot behaviors as the expected behaviors the robot should take during the intervention. The learned DeepQ network enabled the robot to autonomously respond with the appropriate scripted robot behavior in response to a participant’s actions during an intervention.

In [18], a LfD system was developed to learn to deliver group recreational activities from demonstrations provided by healthcare professionals teleoperating a SAR. Namely, healthcare professionals pre-scripted robot behaviors (i.e. speech and arm motions) and demonstrated the structure of an activity by teleoperating the robot behaviors. The demonstrations were then utilized to learn a random forest classifier which mapped the appropriate robot behavior to the state of an activity. The pre-scripted robot behaviors and learned random forest classifier was then utilized by the robot to autonomously deliver a group recreational activity to older adults.

In [19], a LfD system was developed to learn appropriate shopkeeper social behavior by observing human-human interactions between a customer and a shopkeeper. Namely, speech as well as physical locations of participants within the space were recorded during mock shopkeeper and customer interactions within a lab setting. Speech was collected by the participants using a handheld audio recorder which they manually indicated the beginning and ending of their speech so the speech-to-text API could translate the utterance to text. A dynamic hierarchical clustering approach was used to group speech utterances of the shopkeeper into a set of discrete speech behaviors. These clusters were then used to label the customers utterances to train a naïve Bayesian classifier to output robot speech according to input customer utterances. Although the system could learn the appropriate shopkeeper verbal behaviors from human-human interactions, the approach required participants to manually identify when they were speaking during the demonstrations and contained significant speech recognition errors due to limitations in the speech-to-text API.

Current research has demonstrated that LfD approaches can be utilized for learning policies for social tasks from

human demonstrations. However, existing approaches still either require a demonstrator to pre-script or utilize pre-scripted robot behaviors to demonstrate a social task to a robot [17] [18]. Demonstrating a task with pre-scripted behaviors can be unnatural because users have a limited set of behaviors they can utilize during a social interaction and can also eventually result in repetitive interactions. Such repetitive behavior can be limited when teaching individuals with ASD social interaction skills, since the goal of ABA therapy is to teach these individuals to generalize to human-human interactions rather than teaching them how to interact with a robot [20] [21]. Existing LfD approaches are also evaluated by collecting demonstration data within laboratory settings [17] [19]. Such laboratory settings may not reflect real-world interaction data because adult participants are providing mock demonstrations of social interactions and the demonstrations are in a controlled environment. If such approaches are to be utilized for the delivery of therapies to children with ASD, it is necessary that intervention demonstrations be collected within these settings by the actual users that will be interacting with these systems (i.e., therapists and children with ASD). Namely, existing approaches must account for variable responses from children and, consequently, the therapists’ responses while demonstrating an intervention. It is also necessary to address the noisy and unstructured environments where these interventions take place because current speech recognition techniques are not suitable for children or the noisy environments where interventions are held (e.g., classrooms) [22].

In this paper, we introduce a proof-of-concept LfD system that learns the verbal behavior of a therapist from their demonstrations of an ABA-based intervention within a real clinical setting. Namely, our approach is applied to data collected on therapists delivering an emotion recognition intervention within a classroom of an ABA clinic to children with ASD by teleoperating the speech and motion of a SAR [11].

## III. LEARNING FROM DEMONSTRATION USER STUDY

We conducted a user study to collect therapist demonstrations on the delivery of an ABA-based emotion recognition intervention to children with ASD. The primary objective was to collect demonstrations of intervention delivery in a real-world clinical setting. Namely, all demonstrations were collected from practicing therapists delivering an intervention through a robot to a child with ASD in a uncontrolled classroom setting the children typically receive therapies. In addition to the robot-mediated intervention being delivered by the robot as it was teleoperated by a therapist, there was always one to three additional one-on-one therapies between a human therapist and a child occurring within the classroom.

### A. *The Emotion Recognition Intervention*

A board-certified behavior analyst-doctoral (Dr. Korneder) designed the emotion recognition intervention which the therapists delivered through the robot. The main goal of the intervention was to teach the children to recognize emotions

only from body language and without any sound effects nor facial expressions. Recognizing emotions from body language was a useful skill to learn during the COVID-19 pandemic due to facial expressions being hidden by masks [23]. The robot’s lack of facial expressions was especially useful for simulating a scenario where a human’s facial expressions are hidden by a face mask.

The emotions that were being taught to the children were: happy, surprised, tired, sad, scared and angry. The intervention was based on standard ABA procedures which includes three steps. First, a Discriminative Stimulus (SD) is presented by the therapist teleoperator asking the child how they (i.e the robot) are feeling and presenting the emotion using only the robot’s body language. Second, the child is given an opportunity to respond to the question. Third, social praise is provided by the therapist teleoperator if the child answered correctly, a prompt if the child does not answer within a predefined time, or an error correction if the child answers incorrectly. Herein, a complete trial is defined as this three step process. All interventions were one-on-one between the robot and one child. During the interventions a therapist teleoperator was remotely controlling the robot while located in a different room than the classroom where the robot and child were located. Fig. 1 depicts an intervention being delivered to a child while the robot is teleoperated.

### B. Therapist ABA Intervention Demonstrations

Therapists remotely teleoperated the humanoid Pepper robot using a VR-based interface to deliver the emotion recognition intervention to a child with ASD through the robot. Namely, the therapists equipped a virtual reality headset, earphones, a microphone, and hand-held controllers. The virtual reality headset and earphones enabled the therapist to view a video stream of the robot’s camera and hear an audio stream of the robot’s microphone while the robot was interacting with a child. The therapist could then deliver the intervention by controlling the robot’s joints by naturally moving their arms with the hand-held controllers while the VR interface tracked their body motions and mapped it to the robots motions. The therapist could also control the robot’s speech by speaking through a microphone. Please refer to [11] for more details on the VR teleoperation interface, intervention design, and learning outcomes of the study.

### C. Data Collection

A total of eight therapist participants demonstrated the emotion recognition intervention by teleoperating the robot to deliver the intervention to a child with ASD. There were a total of four children participants and the child participating in the demonstration was assigned according to their schedule availability. Each therapist conducted one intervention session with a child. Each intervention session consisted of 9 trials with three different emotions presented three times each. A total of 72 trials of demonstration data was collected from all the participating therapists. Namely, in this study we collected the audio from the microphone utilized by the therapist to control the robot’s speech during

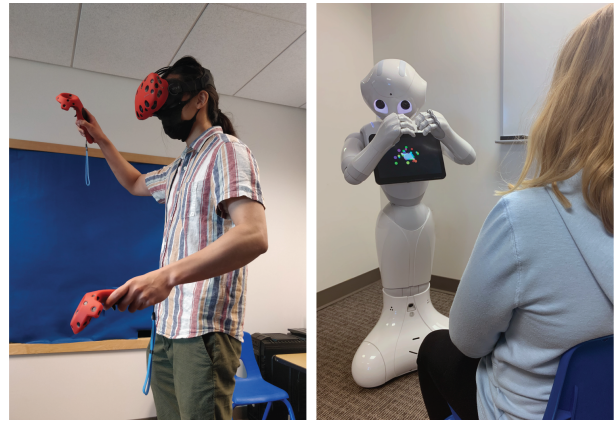


Fig. 1. User teleoperating a robot to demonstrate the delivery of an intervention to a child

the intervention and the audio from the robot’s microphone which recorded all the audio within the classroom where the child participants were working with the robot.

## IV. LEARNING FROM DEMONSTRATION SYSTEM ARCHITECTURE

Our LfD system consists of two steps: 1) identifying therapist behaviors during an intervention from the speech utterances of the therapist teleoperators, and 2) learning a policy which maps the appropriate therapist behavior to execute according to the children’s verbal responses. The first step consists of an unsupervised learning approach where the therapists’ speech utterances were grouped utilizing a K-Means clustering algorithm. These groups of utterances, herein referred to as therapist behaviors, are the different classes of verbal responses therapist have towards the children’s current verbal responses. For example, if the current emotion being taught to the child is “scared” and the child’s verbal response is incorrect, the therapist’s behavior should correspond to a “correction prompt” based on the current feeling (e.g. “I am feeling scared”). If the child’s response is correct, the therapist’s behavior should correspond to “social praise” (e.g. “Good job!”). Once these therapist behaviors are identified, they are utilized to label the child’s verbal responses automatically. The child’s verbal responses and the therapist behaviors used to label them are then used to train a Deep Neural Network (DNN). Namely, the input data for training the DNN is the raw spectrogram of the child’s verbal behavior and the therapist’s behavior labeled to the child’s verbal behavior is the expected output of the DNN. The final learned DNN can then be utilized to determine the appropriate therapist behavior a robot should execute given a child’s verbal behavior.

### A. Identifying Therapist Behaviors During an Intervention

The first step in the LfD system consists of identifying the discrete verbal behaviors therapists have during a demonstrated intervention using an unsupervised learning technique on the captured therapist audio data. Each therapist had his/her own audio data. Google’s speech-to-text API was applied to each therapist’s audio data to obtain their

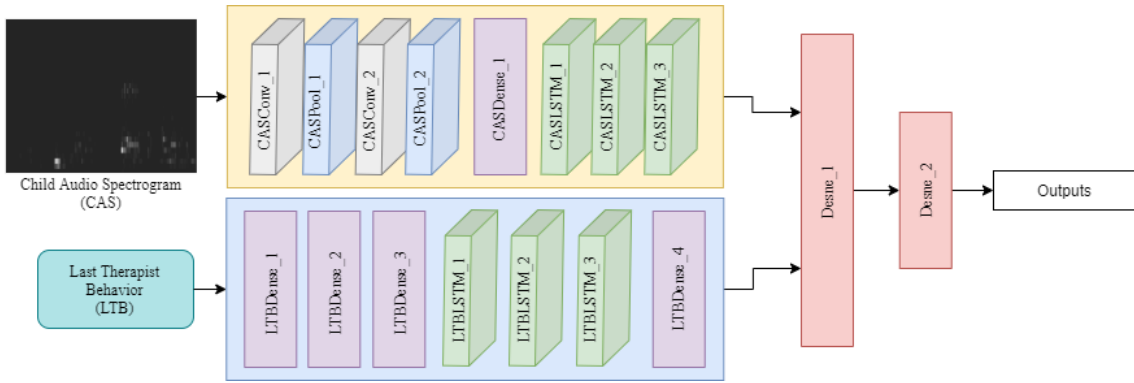


Fig. 2. Deep neural network model architecture

utterances during an intervention. Due to limitations in current speech-to-text technologies, there were some therapist speech utterances that were incorrectly translated or were undetected. Such errors in utterances were manually fixed by a researcher. Latent Semantic Analysis (LSA) was then applied to vectorize the utterances and K-Means clustering was used to cluster these utterances into discrete therapist behaviors. The Elbow method was used to identify the optimal number of clusters (i.e., therapist behaviors) for the utterances [24]. The Elbow method consists of plotting the number of clusters against data distortion and the optimum number of cluster for the data was considered where an “elbow” appears in the plot. Namely, the “elbow” refers to when the distortion converges and increasing the number of clusters does not significantly affect the distortion of the data. According to the elbow method, we identified the optimal number of clusters was 10. After examining the contents of each cluster, two clusters had similar utterances such as “How am I feeling” and “How do I feel”, which correspond to the same discriminative stimulus therapist behavior. Consequently, we combined the two clusters and had a final total of nine clusters. The final nine clusters represented the following therapist behaviors: discriminative stimulus, prompts corresponding to each emotion (i.e. scared, sad, happy, tired, surprised and angry), social praise (e.g. “Good job!” or “Nice work!”), and random instructions (e.g. “clap your hands”, “touch your head”, or “what’s your name”). These random instructions did not pertain to the emotion recognition intervention but were a technique utilized by therapist’s to maintain the motivation and engagement of the children. These clusters (i.e., therapist intervention behaviors) were then used to label the children’s responses.

### B. Extracting and Labeling Child Verbal Responses

The children’s verbal response data during the intervention session was extracted automatically by utilizing the start and end times of therapists’ behaviors. Namely, the children’s verbal responses were considered as the duration between the start of the current therapist behavior and end of the therapist’s last executed behavior. The rationale for utilizing the duration between therapist speech behaviors to obtain children’s verbal responses is due to Google’s speech-to-text API performing poorly with children’s speech and noisy

classroom environments with multiple other inhabitants. Each extracted child verbal response was extracted from the audio data of the intervention session and converted to a spectrogram. The spectrogram was combined with the therapist’s last executed behavior to form a data tuple which was labeled with the current therapist behavior. This data tuple was then utilized to learn a policy to map the appropriate therapist behavior to execute according to the children’s verbal responses.

### C. Modeling Therapist Behaviors During an Intervention

A multi-input DNN was utilized to model therapist behaviors during the delivery of an intervention, Fig. 2. The inputs to the model are the last therapist behavior (LTB) and the Child’s Audio Spectrogram (CAS) after the execution of the last therapist behavior. The spectrogram and the last therapist behavior are input into two separate networks and the network outputs are then combined with two dense layers to estimate the next appropriate therapist behavior to be executed. Hence, the DNN defines a policy which maps the appropriate therapist behavior to execute given a child’s verbal response and the previous executed therapist behavior.

The first network is a Convolutional Neural Network with LSTM layers (CNN-LSTM) and takes as input the child’s audio spectrogram which is a 55x55x1 image. The image is connected to two convolution layers with RELU activation functions, two max-pooling layers, two dense layers, and three LSTM layers. The two convolutional layers have kernel sizes 4x4 and 3x3 with 32 and 64 filters respectively. The 32-filter convolution layer (CASConv\_1) is downsampled with a 3x3 max-pooling layer (CASPool\_1) and connected to the 64-filter convolution layer (CASConv\_2). The CASConv\_2 is downsampled with a 3x3 max-pooling layer (CASPool\_2) and connected to the first dense layer (CASDense\_1). The CASDense\_1 consists of 256 fully connected neurons with a dropout of 0.25. The CASDense\_1 is then connected to 3 LSTM layers (CASLSTM\_1, CASLSTM\_2 and CASLSTM\_3) with 100 units each.

The second network is a dense network with LSTM layers and takes as input the last therapist behavior which is defined by the enumeration for its cluster. The last therapist behavior is passed through four dense layers (LTBDense\_1, LTBDense\_2, LTBDense\_3 and LTBDense\_4) with RELU

activation functions and three LSTM layers (LTBLSTM\_1, LTBLSTM\_2 and LTBLSTM\_3). The number of neurons in the first three dense layers increases progressively as we go deeper into the network with 4, 8 and 16 fully connected neurons respectively. The first three layers are then followed by three LSTM layers. The LSTM layers are then connected to the a dense layer consisting of 3 fully connected neurons.

The outputs of LTBDense\_4 and CASLSTM\_3 are then combined in one dense layer (Dense\_1) with 300 fully connected neurons and a dropout of 0.25. The Dense\_1 is then connected to Dense\_2 which consists of 150 fully connected neurons. The Dense\_2 layer is then connected with the output layer consisting of 9 fully connected neurons, with each neuron representing a possible therapist behavior to be executed.

#### D. Training the Model

A total of 100 data tuples were obtained from the learning from demonstration user study. Table I provides 6 example data tuples collected from the therapist demonstrations. We utilized 75 data tuples for training the DNN model and reserved 25 data tuples for evaluating our approach. As previously mentioned, each data tuple consisted of a spectrogram of a child’s response, the therapist’s last executed behavior, and the current therapist behavior which is utilized to label the data tuple. The model was trained with a batch size of 5, 100 epochs and 0.0001 initial learning rate. The loss was calculated using categorical cross-entropy, and was minimized using the Adam optimizer, which is an implementation of the gradient descent method [25].

### V. RESULTS AND DISCUSSION

We evaluated our model on the remaining 25 data tuples from the demonstration data collected during our user study. The overall classification accuracy of the model was 43.48%. The confusion matrix for these results is provided in Table II.

From these results we can see that the model can correctly identify "SD" and "Social Praise" behaviors, which usually indicate the beginning and end of a trial. However, the model was not able to detect different prompt behaviors for each emotion. We hypothesize this is due to the limited amount of demonstration data currently available for each prompt cluster, and the data imbalance present in the current training dataset. Namely, therapists demonstrated prompts for each emotion less than other behaviors. Prompts for emotions such as "Sad", "Surprised" and "Tired" had only an average of 3 data tuples. On the other hand, the "SD" and "Social Praise" behaviors were demonstrated more frequently during the interventions with the children, which explains the model’s ability to better classify these two classes. Regarding prompts for each emotion, we believe the amount of data compared to the number of classes and the unbalanced amount of data instances had the largest impact on the model performance when classifying a prompt. The input spectrograms also had minimal variance when therapists provided the different prompts for each emotion to the child. Moreover, we can

TABLE I  
SAMPLES OF DIALOGUE FROM THE TRAINING DATA

Sample \ Action	Last Therapist Behavior	Child Response	Therapist Behavior
Sample 1	Sad Prompt	"Happy!"	Sad Prompt: "I Feel Sad"
Sample 2	Sad Prompt	"Sad!"	Social Praise: "Good Job!"
Sample 3	Social Praise	"....."	SD: "How Do I Feel? .. Tired"
Sample 4	SD	"Tired!"	Social Praise: "Great Work!"
Sample 5	Social Praise	"....."	SD: "How Do I Feel? .. Angry"
Sample 6	SD	"Tired!"	Angry Prompt: "I Feel Angry"

observe from the confusion matrix that even with more data instances, the "SD" and "Social Praise" classes had the most false positives. This may indicate overfitting of the data which we believe is also due to an unbalanced dataset. This was further emphasized in the confusion matrix by some classes not being included in the test data (i.e. "Sad" and "Random"). Nevertheless, the concept has shown potential for future investigation and we expect that with additional as well as balanced demonstration data of therapists performing these behaviors that our model could be trained to handle these instances.

### VI. CONCLUSIONS

In this paper, we developed and deployed a LfD system in-the-wild to have SARs learn from therapists the delivery of an ABA-based intervention to children with ASD. The system was applied in a real-world ABA clinical setting with practicing therapists demonstrating the delivery of an emotion recognition intervention to children with ASD in a classroom. The preliminary demonstration data collected from the therapists and results of our work demonstrate that the LfD system architecture is capable of learning the discrete verbal behaviors used by therapists during an intervention. Furthermore, preliminary results suggest that with limited demonstration data the system can learn to apply the appropriate therapist behavior within some intervention scenarios with a child. Such results are promising and we expect that as we continue to collect demonstration data in-the-wild the model behavior will significantly improve. Hence, in the future we plan on continuing to collect therapist demonstrations and improving upon our model development.



TABLE II

CONFUSION MATRIX SUMMARIZING THE PERFORMANCE OF THE MODEL LEARNED FROM THE THERAPIST DEMONSTRATIONS

Predicted \ Actual	"Angry"	"Surprised"	"Tired"	"Random"	"Scared"	"Happy"	"Social Praise"	"SD"	"Sad"
"Angry"	0	0	0	0	0	0	1	0	0
"Surprised"	1	0	0	0	0	0	0	0	0
"Tired"	0	0	0	0	0	0	1	1	0
"Random"	0	0	0	0	0	0	0	0	0
"Scared"	0	0	0	0	0	0	0	2	0
"Happy"	0	0	0	0	0	0	0	0	0
"Social Praise"	1	0	0	0	0	0	5	2	0
"SD"	0	0	0	0	0	1	2	5	0
"Sad"	0	0	0	0	0	0	0	0	0

## REFERENCES

- [1] K. A. Shaw, M. J. Maenner, J. Baio, *et al.*, "Early identification of autism spectrum disorder among children aged 4 years—early autism and developmental disabilities monitoring network, six sites, united states, 2016," *MMWR Surveillance Summaries*, vol. 69, no. 3, p. 1, 2020.
- [2] D. Granpeesheh, D. R. Dixon, J. Tarbox, A. M. Kaplan, and A. E. Wilke, "The effects of age and treatment intensity on behavioral intervention outcomes for children with autism spectrum disorders," *Research in Autism Spectrum Disorders*, vol. 3, no. 4, pp. 1014–1022, 2009.
- [3] R. M. Foxx, "Applied behavior analysis treatment of autism: The state of the art," *Child and Adolescent Psychiatric Clinics of North America*, vol. 17, no. 4, pp. 821–834, 2008, treating Autism Spectrum Disorders.
- [4] E. S. Kim, *et al.*, "Potential clinical impact of positive affect in robot interactions for autism intervention," in *2015 International Conference on Affective Computing and Intelligent Interaction (ACII)*, 2015, pp. 8–13.
- [5] Z. Warren, *et al.*, "Brief report: development of a robotic intervention platform for young children with asd," *Journal of autism and developmental disorders*, vol. 45, no. 12, pp. 3870–3876, 2015.
- [6] M. Begum, *et al.*, "Measuring the efficacy of robots in autism therapy: How informative are standard hri metrics?," in *Proceedings of the Tenth Annual ACM/IEEE International Conference on Human-Robot Interaction*, ser. HRI '15. New York, NY, USA: Association for Computing Machinery, 2015, p. 335–342.
- [7] W. Y. G. Louie, J. Korneder, I. Abbas, and C. Pawluk, "A study on an applied behavior analysis-based robot-mediated listening comprehension intervention for asd," *Paladyn, Journal of Behavioral Robotics*, vol. 12, no. 1, pp. 31–46, 2021.
- [8] M. Hirokawa, A. Funahashi, Y. Itoh, and K. Suzuki, "Adaptive behavior acquisition of a robot based on affective feedback and improvised teleoperation," *IEEE Transactions on Cognitive and Developmental Systems*, vol. 11, no. 3, pp. 405–413, 2019.
- [9] B. A. English, A. Coates, and A. Howard, "Recognition of gestural behaviors expressed by humanoid robotic platforms for teaching affect recognition to children with autism - a healthy subjects pilot study," in *Social Robotics*, A. Kheddar, *et al.*, Eds. Cham: Springer International Publishing, 2017, pp. 567–576.
- [10] A. Taheri, A. Meghdari, M. Alemi, and H. Pourtemad, "Teaching music to children with autism: A social robotics challenge," *Scientia Iranica*, vol. 26, no. Special Issue on: Socio-Cognitive Engineering, pp. 40–58, 2019.
- [11] R. Kulikovskiy, *et al.*, "Can therapists design robot-mediated interventions and teleoperate robots using VR to deliver interventions for ASD?" *2021 IEEE International Conference on Robotics and Automation (ICRA)*, to be published.
- [12] P. G. Esteban, *et al.*, "How to build a supervised autonomous system for robot-enhanced therapy for children with autism spectrum disorder," *Paladyn, Journal of Behavioral Robotics*, vol. 8, no. 1, pp. 18–38, 2017.
- [13] J.-J. Cabibihan, H. Javed, M. Ang, and S. M. Aljunied, "Why robots? a survey on the roles and benefits of social robots in the therapy of children with autism," *International journal of social robotics*, vol. 5, no. 4, pp. 593–618, 2013.
- [14] A. M. Alcorn, *et al.*, "Educators' Views on Using Humanoid Robots With Autistic Learners in Special Education Settings in England," *Frontiers in Robotics and AI*, vol. 6, pp. 1–15, 2019.
- [15] S. Chernova and A. L. Thomaz, "Robot learning from human teachers," *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 8, no. 3, pp. 1–121, 2014.
- [16] S. Calinon, "Learning from demonstration (programming by demonstration)," *Encyclopedia of Robotics*, pp. 1–8, 2018.
- [17] M. Clark-Turner and M. Begum, "Deep reinforcement learning of abstract reasoning from demonstrations," in *Proceedings of the 2018 ACM/IEEE International Conference on Human-Robot Interaction*, 2018, pp. 160–168.
- [18] W. Y. G. Louie and G. Nejat, "A Social Robot Learning to Facilitate an Assistive Group-Based Activity from Non-expert Caregivers," *International Journal of Social Robotics*, vol. 12, pp. 1159–1176, 2020.
- [19] P. Liu, D. F. Glas, T. Kanda, and H. Ishiguro, "Data-driven hri: Learning social behaviors by example from human–human interaction," *IEEE Transactions on Robotics*, vol. 32, no. 4, pp. 988–1008, 2016.
- [20] S. Calinon and A. Billard, "Incremental learning of gestures by imitation in a humanoid robot," in *Proceedings of the ACM/IEEE international conference on Human-robot interaction*, 2007, pp. 255–262.
- [21] K. Vogeley and G. Bente, "'artificial humans': Psychology and neuroscience perspectives on embodiment and nonverbal communication," *Neural Networks*, vol. 23, no. 8, pp. 1077–1090, 2010, social Cognition: From Babies to Robots.
- [22] G. Yeung and A. Alwan, "On the difficulties of automatic speech recognition for kindergarten-aged children," *Interspeech 2018*, pp. 1661–1665, 2018.
- [23] C.-C. Carbon, "Wearing face masks strongly confuses counterparts in reading emotions," *Frontiers in Psychology*, vol. 11, p. 2526, 2020.
- [24] P. Bholowalia and A. Kumar, "Ebk-means: A clustering technique based on elbow method and k-means in wsn," *International Journal of Computer Applications*, vol. 105, no. 9, 2014.
- [25] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2017.