

**Title:** Leveraging video data from a digital smartphone autism therapy to train an emotion detection classifier

**Authors:**

Cathy Hou<sup>1</sup>, Haik Kalantarian<sup>2</sup>, Peter Washington<sup>3</sup>, Kaiti Dunlap<sup>2</sup>, Dennis P. Wall<sup>2,4,5</sup>

**Affiliations:**

<sup>1</sup> Department of Computer Science, Stanford University, Stanford, California, USA.

<sup>2</sup> Department of Pediatrics (Systems Medicine), Stanford University, Stanford, California, USA.

<sup>3</sup> Department of Bioengineering, Stanford University, Stanford, California, USA

<sup>4</sup> Department of Biomedical Data Science, Stanford University, Stanford, California, USA

<sup>5</sup> Department of Psychiatry and Behavioral Sciences (by courtesy), Stanford University, Stanford, California, USA.

**Key words:** emotion recognition, autism, machine learning, precision therapy

## **ABSTRACT:**

Autism spectrum disorder (ASD) is a neurodevelopmental disorder affecting one in 40 children in the United States and is associated with impaired social interactions, restricted interests, and repetitive behaviors. Previous studies have demonstrated the promise of applying mobile systems with real-time emotion recognition to autism therapy, but existing platforms have shown limited performance on videos of children with ASD. We propose the development of a new emotion classifier designed specifically for pediatric populations, trained with images crowdsourced from an educational mobile charades-style game: *Guess What?*. We crowdsourced the acquisition of videos of children portraying emotions during remote game sessions of *Guess What?* that yielded 6,344 frames from fifteen subjects. Two raters manually labeled the frames with four of the Ekman universal emotions (happy, scared, angry, sad), a “neutral” class, and “n/a” for frames with an indeterminable label. The data were pre-processed, and a model was trained with a transfer-learning and neural-architecture-search approach using the Google Cloud AutoML Vision API. The resulting classifier was evaluated against existing approaches (Microsoft’s Azure Face API and Amazon Web Service’s Rekognition) using the standard metrics of F1 score. The resulting classifier demonstrated superior performance across all evaluated emotions, supporting our hypothesis that a model trained with a pediatric dataset would outperform existing emotion-recognition approaches for the population of interest. These results suggest a new strategy to develop precision therapy for autism at home by integrating the model trained with a personalized dataset to the mobile game.

## INTRODUCTION:

Autism spectrum disorder (ASD) is a neurodevelopmental disorder affecting one in 40 children in the United States and is associated with impaired social interactions, restricted interests, and repetitive behaviors [1-2]. While there is no cure, studies have shown the efficacy of Applied Behavioral Analysis (ABA) therapy, if administered at a young age and customized to address the child's unique deficits [3-4]. However, caring for a child with ASD can generate a financial burden on the family [5]. Additionally, the increasing prevalence of the condition is resulting in a short supply of certified specialists, further hindering treatment options [6].

We developed *Guess What?* [7-9], a charades-style mobile game that delivers social training to children with ASD at home, to mitigate the high costs and shortage of traditional interventions. To play the game, the child interprets and acts out prompts that are displayed on the screen while the caregiver is tasked with guessing the prompt. Multiple decks and prizes tailor to the child's preferences and help increase engagement.

*Guess What?* incorporates two teaching methods based on ABA principles: Discrete Trial Training (DTT) and Pivotal Response Treatment (PRT). DTT breaks down the skill into discrete trials that build up the skill step by step [10]. Each trial follows a specific set of steps consisting of an antecedent, prompt, response, reinforcement, and brief pause [10]. PRT is less structured and initiated by the child, emphasizing natural reinforcement and targeting pivotal areas of a child's development instead of specific behaviors [11]. Multiple studies suggest that DTT helps improve emotion recognition and expression [10] and PRT enhances communication skills in children with ASD [11].

Previous studies have demonstrated the promise of applying mobile systems with real-time emotion recognition to ABA therapy [12-18]. Integrating an automatic emotion

classifier into *Guess What?* will provide supplemental reinforcement to the caregiver and allow for the development of additional features integral to ABA therapy: adapting prompts and difficulty to target the child's specific deficits and offering appropriate visual cues to assist the child [3].

However, existing emotion recognition platforms are not optimized for research on children [19-20] as a result of being trained on datasets in which pediatric populations are highly underrepresented such as the CIFAR-100, ImageNet [21], Cohn-Kanade Database [22] and Belfast-Induced Natural Emotion Databases [23]. *Guess What?* can serve as a data acquisition tool and aggregate emotive videos for autism research that can be used to train a more effective automatic emotion recognition platform. The use of data collected from mobile devices, such as the built-in camera, allow for continuous phenotyping and repeat diagnoses in home settings [24-39]. This motivates the development of a new emotion classifier designed specifically for pediatric populations, trained with images crowdsourced from *Guess What?*.

## METHODS:

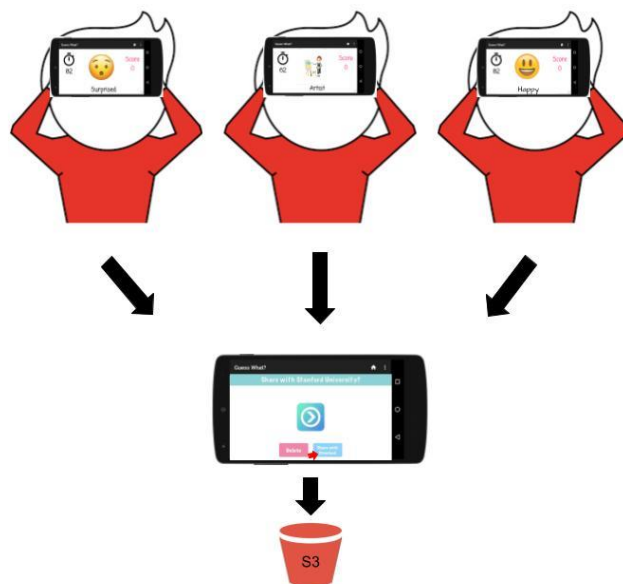


Figure 1. Crowdsourced videos taken during game sessions are stored in an Amazon S3 bucket (with participant's consent).

### Game Design

*Guess What?* is available for both Android and iOS platforms [7-9]. The child begins by selecting one of the following themed decks to play: animals, emoji, faces, gestures, jobs, objects, sports, chores, and a special deck created for toddlers. The caregiver will hold the device outwards with the screen facing the child. During the 90-second game session, the child acts out the prompt displayed on the screen while the caregiver guesses. If the caregiver's guess is correct, the child prompts the caregiver, who tilts the phone. The game then rewards the child with a point, resulting in another image appearing on the screen. This process repeats for a fixed amount of time. The entire game session is recorded using the front-facing camera on the device, focusing on the child's actions. If the user grants permission to share this footage, the video is uploaded to a secure and encrypted Amazon Web Services S3 bucket. This data upload and

storage process is fully compliant with the Stanford University's High-Risk Application security standards. Additional metadata included with the video includes the prompts used in the session, timing logs, and the number of points awarded.

### Data Collection

Figure 1 illustrates the data aggregation process. The two decks that are the most closely associated with emotion recognition and expression are the *emoji* and *faces* decks. These decks contain emoticons (cartoon representations of facial emotions) and real images of children expressing various emotions, respectively. Using crowdsourced videos from fifteen subjects remotely playing these two decks subsampled at 5 frames per second (FPS), a dataset consisting of 6,344 frames was built.

### Data Processing

To establish ground truth, two raters manually labeled each of the 6,344 frames with either one of four Ekman universal emotions (happy, sad, scared, angry) [40] or a neutral label. In cases where there were no faces in the frame or the face was too blurry to discern, the raters labeled the frame with "n/a." To filter the data, all frames with rater disagreement or labeled as "n/a" were discarded. Faces were then extracted from the remaining frames using the OpenCV library [41], yielding 757 frames. Figure 2 is a confusion matrix illustrating the distribution of the raters' labels. The Cohen's Kappa statistic for inter-rater reliability [42], a metric which accounts for agreements due to chance, was 0.8, indicating a high level of reliability between the two raters.

**Distribution of Frames Between Two Raters**

	HAPPY	NEUTRAL	SCARED	ANGRY	SAD
Rater 2 HAPPY	868	201	38	14	14
NEUTRAL	100	1151	59	9	49
SCARED	1	32	59	2	4
ANGRY	13	4	0	30	2
SAD	2	9	5	12	60
Rater 1					

Figure 2. Confusion matrix of the two raters' emotion labels.

### Classifier Training

The proposed emotion classifier was trained using Google Cloud's AutoML pipeline, which leverages Google's transfer learning and neural architecture search technologies to automate the determination of the strongest network architecture and optimal hyperparameter configurations to minimize the loss functions [43-44]. Due to an uneven distribution of emotions, data augmentation methods from the Imgaug library [45] were performed to increase the number of viable frames to 989, with roughly 200 corresponding to each of the five emotions. The specific methods performed included horizontally flipping the image, cropping the image, blurring the image, improving or worsening contrast, adding Gaussian noise to the image, brightening or darkening the image, and applying affine transformations to the image, all performed in random order and of varied magnitudes [45]. 861 frames were used for training and

validation, and the remaining 128 frames were used for testing. Figure 3 illustrates the data processing and classifier training procedure.

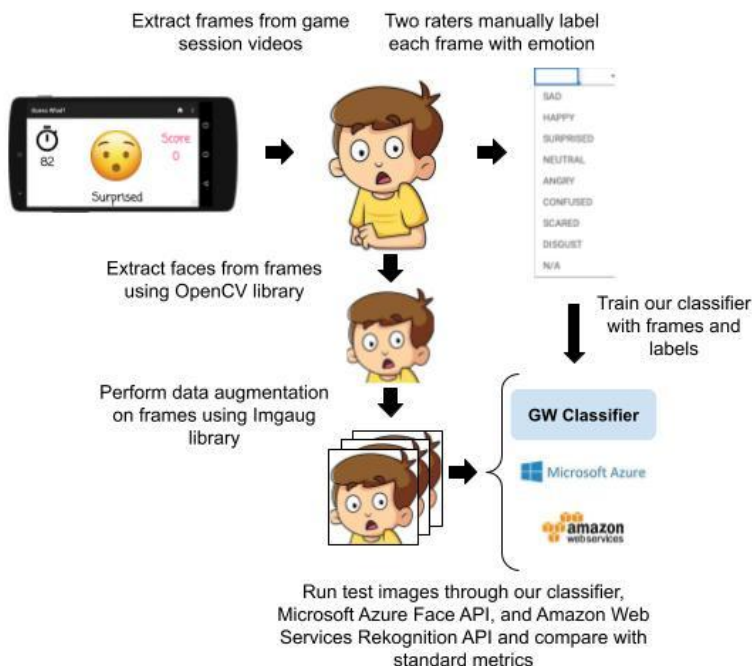


Figure 3. The classifier's training procedure.

### Data Analysis

To evaluate the performance of the models, the F1 score was calculated as follows:

$$P = \frac{t_p}{t_p + f_p} \quad R = \frac{t_p}{t_p + f_n} \quad F1 = 2 \times \frac{P \times R}{P + R}$$

P stands for precision, R stands for recall, F1 stands for F1 score,  $t_p$  stands for true positive,  $f_p$  stands for false positive, and  $f_n$  stands for false negative.



## RESULTS AND DISCUSSION:

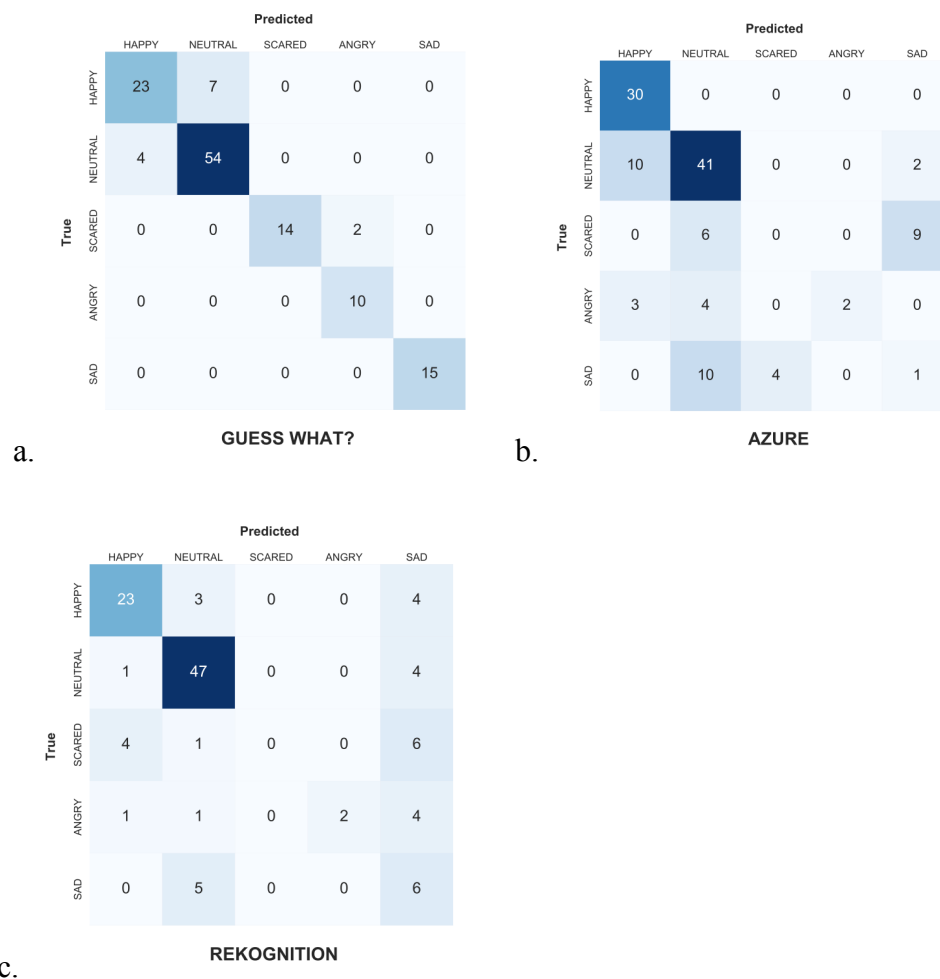


Figure 4. Confusion matrices of the proposed classifier, Azure, and Rekognition.

The proposed classifier was compared to two existing emotion recognition platforms: Microsoft’s Azure Face API [46] and Amazon Web Services’ Rekognition [47]. The same 128 frames that were tested on the proposed classifier were tested on these two classifiers. Figure 4b and 4c show that these classifiers can recognize happy and neutral but perform poorly on angry,

sad, and scared classes. These results suggest the need for a new classifier that demonstrates stronger performance across all emotions for pediatric populations.

The performance of the proposed classifier is illustrated in Figure 4a. Figure 4a shows that the most discrepancies occurred between differentiating neutral from happy, which contrasts with the performance of the other classifiers. However, this proposed classifier generally showed a very strong performance for all five emotions, especially when compared to the other two existing classifiers.

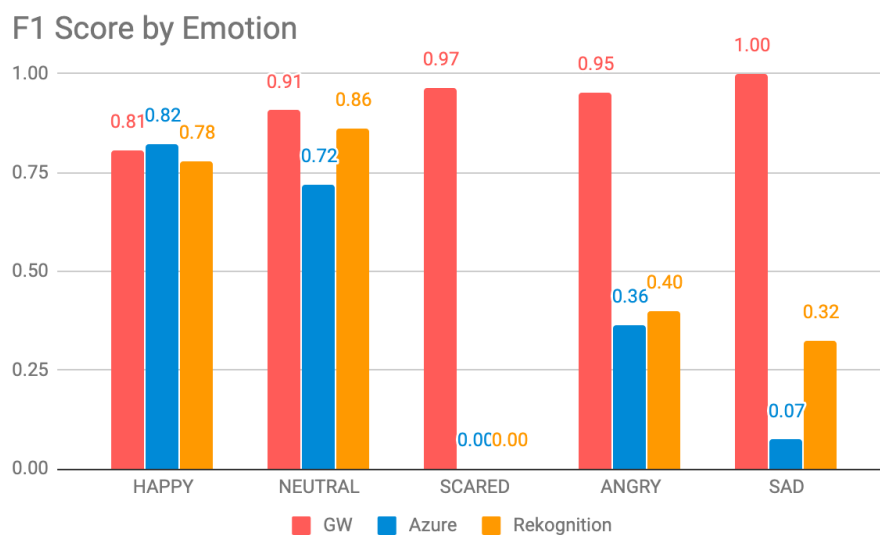


Figure 5. F1 score by emotion by classifier

Figure 5 shows the F1 scores [48] of each classifier separated by emotion. The proposed classifier displayed the highest F1 scores across all but one emotion in comparison to the existing commercial emotion recognition platforms. The one deviation occurred when the Azure classifier had an F1 score of 0.82 for happy, while the proposed classifier had an F1 score of 0.81. However, the Azure classifier had the lowest F1 scores for all of the other classes.

## CONCLUSION AND FUTURE DIRECTIONS:

Both the Azure and Rekognition classifiers performed relatively well on happy and neutral frames but failed with other emotions. In addition, the proposed classifier performed better on scared, angry, and sad frames, while the two existing classifiers demonstrated stronger performance on happy and neutral frames. Because heavy data augmentation procedures had to be performed on the scared, angry, and sad classes to evenly distribute the training set, these results suggest that overfitting may have occurred with the proposed classifier. Generating a more diverse dataset of frames to begin with may alleviate this issue and will be addressed in future work. However, due to the personalized nature of classifiers, perhaps the overfitting is useful in this case. As a result of the limited training set and contrasting performance of the classifiers, in future work, a transfer learning approach will be taken to improve the performance of the proposed classifier across all emotions: the five emotions addressed in this study as well as disgust and surprise which are two other Ekman emotions. Additionally, this emotion classifier generalized to children with ASD will be integrated into *Guess What?* to provide supplemental reinforcement and allow for the development of new features including adapting the game to target specific deficits and providing appropriate guiding feedback.

## **ACKNOWLEDGEMENTS:**

These studies were supported by awards to DW by the National Institutes of Health (1R21HD091500-01 and 1R01EB025025-01). Additionally, we acknowledge the support of grants to DW from The Hartwell Foundation, the David and Lucile Packard Foundation Special Projects Grant, Beckman Center for Molecular and Genetic Medicine, Coulter Endowment Translational Research Grant, Berry Fellowship, Spectrum Pilot Program, Stanford's Precision Health and Integrated Diagnostics Center (PHIND), Wu Tsai Neurosciences Institute Neuroscience: Translate Program, and Stanford's Institute of Human Centered Artificial Intelligence as well as philanthropic support from Mr. Peter Sullivan. HK would also like to acknowledge support from the Thrasher Research Fund and Stanford NLM Clinical Data Science program (T-15LM007033-35). PW would like to acknowledge support from Mr. Schroeder and the Stanford Interdisciplinary Graduate Fellowship (SIGF) as the Schroeder Family Goldman Sachs Graduate Fellow. Finally, we acknowledge the Stanford Institutes of Medicine Summer Research Program for funding CH.

## REFERENCES:

- [1] Kogan, M. D., Vladutiu, C. J., Schieve, L. A., Ghandour, R. M., Blumberg, S. J., Zablotsky, B., ... & Lu, M. C. (2018). The prevalence of parent-reported autism spectrum disorder among US children. *Pediatrics*, *142*(6), e20174161.
- [2] American Psychiatric Association. (2013). *Diagnostic and statistical manual of mental disorders (DSM-5®)*. American Psychiatric Pub.
- [3] Virués-Ortega, J. (2010). Applied behavior analytic intervention for autism in early childhood: Meta-analysis, meta-regression and dose–response meta-analysis of multiple outcomes. *Clinical psychology review*, *30*(4), 387-399.
- [4] Dawson, G. (2008). Early behavioral intervention, brain plasticity, and the prevention of autism spectrum disorder. *Development and psychopathology*, *20*(3), 775-803.
- [5] Horlin, C., Falkmer, M., Parsons, R., Albrecht, M. A., & Falkmer, T. (2014). The cost of autism spectrum disorders. *PloS one*, *9*(9), e106552.
- [6] Ning, M., Daniels, J., Schwartz, J., Dunlap, K., Washington, P., Kalantarian, H., ... & Wall, D. P. (2019). Identification and Quantification of Gaps in Access to Autism Resources in the United States: An Infodemiological Study. *Journal of medical Internet research*, *21*(7), e13094.
- [7] Kalantarian, H., Washington, P., Schwartz, J., Daniels, J., Haber, N., & Wall, D. P. (2019). Guess What?. *Journal of Healthcare Informatics Research*, *3*(1), 43-66.
- [8] Kalantarian, H., Washington, P., Schwartz, J., Daniels, J., Haber, N., & Wall, D. (2018, June). A Gamified Mobile System for Crowdsourcing Video for Autism Research. In *2018 IEEE International Conference on Healthcare Informatics (ICHI)* (pp. 350-352). IEEE.
- [9] Kalantarian, H., Jedoui, K., Washington, P., & Wall, D. P. (2018). A Mobile Game for Automatic Emotion-Labeling of Images. *IEEE Transactions on Games*.
- [10] Sigafos, J., Carnett, A., O'Reilly, M. F., & Lancioni, G. E. (2019). Discrete trial training: A structured learning approach for children with ASD.

- [11] Koegel, R. L., & Koegel, L. K. (2006). *Pivotal response treatments for autism: Communication, social, & academic development*. Paul H Brookes Publishing.
- [12] Voss, C., Schwartz, J., Daniels, J., Kline, A., Haber, N., Washington, P., ... & Feinstein, C. (2019). Effect of Wearable Digital Intervention for Improving Socialization in Children With Autism Spectrum Disorder: A Randomized Clinical Trial. *JAMA pediatrics*, *173*(5), 446-454.
- [13] Daniels, J., Schwartz, J. N., Voss, C., Haber, N., Fazel, A., Kline, A., ... & Wall, D. P. (2018). Exploratory study examining the at-home feasibility of a wearable tool for social-affective learning in children with autism. *npj Digital Medicine*, *1*(1), 32.
- [14] Daniels, J., Schwartz, J., Haber, N., Voss, C., Kline, A., Fazel, A., ... & Wall, D. (2017). 5.13 design and efficacy of a wearable device for social affective learning in children with autism. *Journal of the American Academy of Child & Adolescent Psychiatry*, *56*(10), S257.
- [15] Washington, P., Voss, C., Kline, A., Haber, N., Daniels, J., Fazel, A., De, T., Feinstein, C., Winograd, T., & Wall, D.P. (2017). SuperpowerGlass: A Wearable Aid for the At-Home Therapy of Children with Autism. *IMWUT*, *1*, 112:1-112:22.
- [16] Voss, C., Washington, P., Haber, N., Kline, A., Daniels, J., Fazel, A., ... & Wall, D. (2016, September). Superpower glass: delivering unobtrusive real-time social cues in wearable systems. In *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing: Adjunct* (pp. 1218-1226). ACM.
- [17] Washington, P., Voss, C., Haber, N., Tanaka, S., Daniels, J., Feinstein, C., ... & Wall, D. (2016, May). A wearable social interaction aid for children with autism. In *Proceedings of the 2016 CHI Conference Extended Abstracts on Human Factors in Computing Systems* (pp. 2348-2354). ACM.
- [18] Porayska-Pomsta, K., Frauenberger, C., Pain, H., Rajendran, G., Smith, T., Menzies, R., ... & Avramides, K. (2012). Developing technology for autism: an interdisciplinary approach. *Personal and Ubiquitous Computing*, *16*(2), 117-127.

- [19] Kalantarian, H., Jedoui, K., Washington, P., Tariq, Q., Dunlap, K., Schwartz, J., & Wall, D. P. (2019). Labeling images with facial emotion and the potential for pediatric healthcare. *Artificial intelligence in medicine*, 98, 77-86.
- [20] Kalantarian, H., Khaled J., Washington, P., ... & Wall, D. (2018, December). The Limitations of Real-Time Emotion Recognition for Autism Research.
- [21] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).
- [22] Cohn, J. F. (1999). Cohn-Kanade AU-coded facial expression database. *Pittsburgh University*.
- [23] Douglas-Cowie, E., Cowie, R., & Schröder, M. (2000). A new emotion database: considerations, sources and scope. In *ISCA tutorial and research workshop (ITRW) on speech and emotion*.
- [24] Abbas, Halim, Ford Garberson, Eric Glover, and Dennis P. Wall. "Machine learning for early detection of autism (and other conditions) using a parental questionnaire and home video screening." In 2017 IEEE International Conference on Big Data (Big Data), pp. 3558-3561. IEEE, 2017.
- [25] Abbas, Halim, Ford Garberson, Stuart Liu-Mayo, Eric Glover, and Dennis P. Wall. "Multi-modular Ai Approach to Streamline Autism Diagnosis in Young children." *Scientific reports* 10, no. 1 (2020): 1-8.
- [26] Duda, M., J. A. Kosmicki, and D. P. Wall. "Testing the accuracy of an observation-based classifier for rapid detection of autism risk." *Translational psychiatry* 4, no. 8 (2014): e424-e424.
- [27] Duda, M., R. Ma, N. Haber, and D. P. Wall. "Use of machine learning for behavioral distinction of autism and ADHD." *Translational psychiatry* 6, no. 2 (2016): e732-e732.

- [28] Duda, M., N. Haber, J. Daniels, and D. P. Wall. "Crowdsourced validation of a machine-learning classification system for autism and ADHD." *Translational psychiatry* 7, no. 5 (2017): e1133-e1133.
- [29] Fusaro, Vincent A., Jena Daniels, Marlena Duda, Todd F. DeLuca, Olivia D'Angelo, Jenna Tamburello, James Maniscalco, and Dennis P. Wall. "The potential of accelerating early detection of autism through content analysis of YouTube videos." *PLOS one* 9, no. 4 (2014): e93533.
- [30] Levy, Sebastien, Marlena Duda, Nick Haber, and Dennis P. Wall. "Sparsifying machine learning models identify stable subsets of predictive features for behavioral detection of autism." *Molecular autism* 8, no. 1 (2017): 65.
- [31] Leblanc, Emilie, Peter Washington, Maya Varma, Kaitlyn Dunlap, Yordan Penev, Aaron Kline, and Dennis P. Wall. "Feature replacement methods enable reliable home video analysis for machine learning detection of autism." *Scientific reports* 10, no. 1 (2020): 1-11.
- [32] Stark, David E., Rajiv B. Kumar, Christopher A. Longhurst, and Dennis P. Wall. "The quantified brain: a framework for mobile device-based assessment of behavior and neurological function." *Applied clinical informatics* 7, no. 2 (2016): 290.
- [33] Tariq, Qandeel, Jena Daniels, Jessey Nicole Schwartz, Peter Washington, Haik Kalantarian, and Dennis Paul Wall. "Mobile detection of autism through machine learning on home video: A development and prospective validation study." *PLoS medicine* 15, no. 11 (2018): e1002705.
- [34] Tariq, Qandeel, Scott Lanyon Fleming, Jessey Nicole Schwartz, Kaitlyn Dunlap, Conor Corbin, Peter Washington, Haik Kalantarian, Naila Z. Khan, Gary L. Darmstadt, and Dennis Paul Wall. "Detecting developmental delay and autism through machine learning models using home videos of Bangladeshi children: Development and validation study." *Journal of medical Internet research* 21, no. 4 (2019): e13822.



- [35] Washington, Peter, Haik Kalantarian, Qandeel Tariq, Jessey Schwartz, Kaitlyn Dunlap, Brianna Chrisman, Maya Varma et al. "Validity of online screening for autism: crowdsourcing study comparing paid and unpaid diagnostic tasks." *Journal of medical Internet research* 21, no. 5 (2019): e13668.
- [36] Washington, Peter, Emilie Leblanc, Kaitlyn Dunlap, Yordan Penev, Aaron Kline, Kelley Paskov, Min Woo Sun et al. "Precision Telemedicine through Crowdsourced Machine Learning: Testing Variability of Crowd Workers for Video-Based Autism Feature Recognition." *Journal of personalized medicine* 10, no. 3 (2020): 86.
- [37] Washington, Peter, Emilie Leblanc, Kaitlyn Dunlap, Yordan Penev, Maya Varma, Jae-Yoon Jung, Brianna Chrisman et al. "Selection of trustworthy crowd workers for telemedical diagnosis of pediatric autism spectrum disorder." *PSB*, 2021.
- [38] Washington, Peter, Kelley Marie Paskov, Haik Kalantarian, Nathaniel Stockham, Catalin Voss, Aaron Kline, Ritik Patnaik et al. "Feature selection and dimension reduction of social autism data." In *Pac Symp Biocomput*, vol. 25, pp. 707-718. 2020.
- [39] Washington, Peter, Natalie Park, Parishkrita Srivastava, Catalin Voss, Aaron Kline, Maya Varma, Qandeel Tariq et al. "Data-driven diagnostics and the potential of mobile artificial intelligence for digital therapeutic phenotyping in computational psychiatry." *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging* (2019).
- [40] Ekman, P., Friesen, W. V., O'sullivan, M., Chan, A., Diacoyanni-Tarlatzis, I., Heider, K., ... & Scherer, K. (1987). Universals and cultural differences in the judgments of facial expressions of emotion. *Journal of personality and social psychology*, 53(4), 712.
- [41] Wagner, P. (2012). Face recognition with opencv. *Order A J. Theory Ordered Sets Its Appl*, 1-26.
- [42] McHugh, M. L. (2012). Interrater reliability: the kappa statistic. *Biochemia medica: Biochemia medica*, 22(3), 276-282.

- [43] Li, F. F., & Li, J. (2018). Cloud AutoML: Making AI accessible to every business. *Internet: [https://www. blog. google/topics/google-cloud/cloud-automl-making-ai-accessible-everybusiness](https://www.blog.google/topics/google-cloud/cloud-automl-making-ai-accessible-everybusiness)*.
- [44] Wong, C., Houlsby, N., Lu, Y., & Gesmundo, A. (2018). Transfer learning with neural automl. In *Advances in Neural Information Processing Systems* (pp. 8356-8365).
- [45] Jung, A. (2017). Imgaug: a library for image augmentation in machine learning experiments.
- [46] Face API - Facial Recognition Software: Microsoft Azure. (n.d.). Retrieved from <https://azure.microsoft.com/en-us/services/cognitive-services/face/>
- [47] Amazon Rekognition. (n.d.). Retrieved from <https://aws.amazon.com/rekognition/>
- [48] Goutte, C., & Gaussier, E. (2005, March). A probabilistic interpretation of precision, recall and F-score, with implication for evaluation. In *European Conference on Information Retrieval* (pp. 345-359). Springer, Berlin, Heidelberg.