

1 Using visual attention estimation on videos for automated prediction of autism spectrum disorder and symptom  
2 severity in preschool children

3

4 Ryan Anthony J. de Belen<sup>1\*</sup>, Valsamma Eapen<sup>2</sup>, Tomasz Bednarz<sup>3</sup> and Arcot Sowmya<sup>1</sup>

5 <sup>1</sup>School of Computer Science and Engineering, University of New South Wales, New South Wales, Australia

6 <sup>2</sup>School of Psychiatry, University of New South Wales, New South Wales, Australia

7 <sup>3</sup>School of Art & Design, University of New South Wales, New South Wales, Australia

8 \*Corresponding author

9 Email: [r.debelen@unsw.edu.au](mailto:r.debelen@unsw.edu.au) (RAJDB)

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

## 31 **Abstract**

32           Atypical visual attention in individuals with autism spectrum disorders (ASD) has been utilised as a  
33 unique diagnosis criterion in previous research. This paper presents a novel approach to the automatic and  
34 quantitative screening of ASD as well as symptom severity prediction in preschool children. We develop a  
35 novel computational pipeline that extracts learned features from a dynamic visual stimulus to classify ASD  
36 children and predict the level of ASD-related symptoms. Experimental results demonstrate promising  
37 performance that is superior to using handcrafted features and machine learning algorithms, in terms of  
38 evaluation metrics used in diagnostic tests. Using a leave-one-out cross-validation approach, we obtained an  
39 accuracy of 94.59%, a sensitivity of 100%, a specificity of 76.47% and an area under the receiver operating  
40 characteristic curve (AUC) of 96% for ASD classification. In addition, we obtained an accuracy of 94.74%, a  
41 sensitivity of 87.50%, a specificity of 100% and an AUC of 99% for ASD symptom severity prediction.

## 42 Introduction

43 Autism spectrum disorders (ASD) are currently being diagnosed through visual observation and  
44 analysis of children’s natural behaviours. While a gold standard observational tool is available, early screening  
45 of ASD in children still remains a complex problem. It is often expensive and time-consuming<sup>1</sup> to conduct  
46 interpretative coding of child observations, parent interviews and manual testing<sup>2</sup>. In addition, differences in  
47 professional training, resources and cultural context may affect the reliability and validity of the results obtained  
48 from a clinician’s observations<sup>3</sup>. Furthermore, the behaviours of children in their natural environments (e.g.,  
49 home) cannot be typically captured by clinical observation ratings. To reduce waiting periods for access to  
50 interventions, it is important to develop new methods of ASD diagnosis without compromising accuracy and  
51 clinical relevance. This is critical because early diagnosis and intervention can provide long-term improvements  
52 for the child and even have a greater effect on clinical outcomes<sup>4</sup>.

53 Recent advances in technology have allowed for the quantification of different biological and  
54 behavioural markers that are useful in ASD research (see <sup>5,6</sup> for reviews). Eye-tracking technology has shown  
55 promise in providing a non-invasive and objective tool for ASD research<sup>7,8</sup>. Several eye-tracking studies have  
56 identified unique visual attention patterns in ASD individuals. Gaze abnormalities in toddlers (<3-year-olds)  
57 include reduced attention to eye and head regions, reduced preference for biological motion, difficulties in  
58 response to joint attention behaviours<sup>9</sup> and scene monitoring challenges during explicit dyadic cues<sup>10</sup>. Pierce, et  
59 al.<sup>11</sup> Pierce, et al.<sup>12</sup> Moore, et al.<sup>13</sup> developed a geometric preference (“GeoPref”) test that contains both  
60 geometric and social videos. It was found that a subset of ASD participants exhibited a visual preference for  
61 geometric motion. This finding has already been leveraged by a growing number of studies that aim to leverage  
62 atypical visual attention to identify individuals with ASD<sup>14,15</sup> and predict symptom severity<sup>16</sup>.

63 Computational models that predict visual attention (i.e., saliency) have seen tremendous progress,  
64 starting from handcrafted features dating back to 1998<sup>17</sup> to a resurgence of deep neural networks (DNNs)<sup>18,19</sup>.  
65 This breakthrough has generated great interest in utilising saliency prediction as a diagnostic paradigm for ASD.  
66 For example, there is a growing collection of eye movements of ASD children recorded during image-<sup>20-22</sup> and  
67 video<sup>22</sup>- viewing tasks. Although the use of saliency detection models on image datasets has resulted in  
68 remarkable diagnostic performance, there is still a lack of diagnostic paradigms that utilise dynamic saliency  
69 detection. In fact, the most common approach of studies that utilise dynamic stimuli is to convert the eye-  
70 tracking data into an image and perform image classification to identify individuals with ASD. In this work, we

71 present a novel pipeline that leverages the dynamic visual attention of humans for ASD diagnosis, as well as  
72 symptom severity prediction.

73 This paper makes three major contributions to the field. First, we implement a data-driven approach to  
74 learn the dynamic visual attention of humans on videos and extract spatiotemporal features for downstream  
75 tasks (e.g., ASD classification and symptom severity prediction). Second, we develop a novel computational  
76 pipeline to diagnose ASD based on the learned features from dynamic visual stimuli. Finally, we use a similar  
77 method to predict the level of ASD-related symptoms from eye-tracking data of children obtained during a free-  
78 viewing task. In the next section, we discuss published works that are related to ours. Despite the growing  
79 literature, it is evident that the comparison of results is challenging due to the lack of publicly available datasets  
80 and open-source code repositories. This is even further complicated by the differences in the participants, age  
81 group and stimuli used in the experiments, making fair and straightforward performance comparisons more  
82 difficult. Nevertheless, we compare our work with a simple thresholding technique<sup>11-13</sup> and a machine learning  
83 (ML) classification approach using handcrafted features<sup>23,24</sup>.

## 84 **Related works**

85 Over the last decade, different behavioural and biological markers have already been quantified, to  
86 some extent, using computer vision methods (a comprehensive review<sup>5</sup> is available). Various data modalities,  
87 such as magnetic resonance imaging (MRI)/functional MRI<sup>25-30</sup>, eye-gaze data<sup>14,31-36</sup>, stereotyped behaviours<sup>37-42</sup>  
88 and multimodal data<sup>43</sup> have been utilised in autism diagnosis. We first provide a review of publicly available  
89 datasets that utilise the eye-tracking paradigm. Afterwards, related works that utilise eye-tracking data for the  
90 following purposes are reviewed: (i) saliency prediction in ASD, (ii) ASD diagnosis using static stimuli, (iii)  
91 ASD diagnosis using dynamic stimuli and (iv) ASD risk and symptom severity prediction. Each purpose has a  
92 corresponding table that includes the following information about the published research: mean age of the  
93 participants, gender distribution, stimuli and input used, methodology and conclusion. While not as exhaustive  
94 and rigorous in inclusion criteria as a systematic review, we hope that our discussion below will help the readers  
95 navigate the research landscape and better situate our work in the literature. Readers are also encouraged to read  
96 systematic reviews<sup>8,44</sup> for additional reference.

## 97 **Publicly available datasets**

98 There is a growing number of publicly available datasets that capture the eye-tracking data of ASD  
99 participants. In Table 1, we provide a summary of these datasets by providing descriptions of their target

100 application area, the mean age of the participants, sample size, stimuli used and data format provided by the

101 authors. There are two datasets for saliency estimation<sup>20,21</sup> and two datasets for ASD classification<sup>22,45</sup>.

102 *Table 1 List of publicly available datasets and their corresponding application area, mean age, sample size, stimuli and data format provided by the authors.*

Authors	Application area	Mean age (SD) in years	Sample size	Stimuli	Data format
Duan, et al. <sup>20</sup> Gutiérrez, et al. <sup>21</sup> (Saliency4ASD dataset)	Saliency estimation ASD classification	All participants: 8.00 (NR)	ASD: 14 TD: 14	300 images that depict diverse naturalistic scenes and may contain humans, animals, buildings or objects.	Image with the associated eye-tracking data of the participants
Le Meur, et al. <sup>22</sup> (MIE Fo and MIE No)	Saliency estimation	MIE Fo: ASD: 16.00 (2.00)  MIE No: 29.00 (7.00)	MIE Fo: ASD: 17  MIE No: ASD: 12	25 images with low semantic meaning and a low emotional arousal	Image with the associated eye-tracking data of the participants
Carette, et al. <sup>45</sup>	ASD classification	All participants: 7.88 (NR)	ASD: 29 TD: 30	Combination of static and dynamic stimuli that depict naturalistic scenes, initiate joint attention and static face or objects	Scanpath image that visualises the eye-tracking data of the participants. The visualised scanpath images are then converted to grayscale and rescaled for further processing.

103 ASD: Autism Spectrum Disorder, NR: Not reported, SD: Standard deviation, TD: Typically Developing

## 104 **Saliency prediction in ASD**

105           Accurately predicting the visual attention (i.e., saliency maps) of ASD individuals can boost prediction  
106 performance because classification models can better leverage the distinction between the visual attention of  
107 ASD and typically developing (TD) individuals. Table 2 shows the published research that aims to model the  
108 visual attention of ASD participants by developing different saliency models.

109           Duan, et al. <sup>46</sup> compared the performance of five state-of-the-art (SOTA) saliency prediction networks  
110 based on a deep neural network (DNN) architecture with pre-trained and fine-tuned weights on their dataset.  
111 Experimental results revealed that transfer learning provides a useful approach to modelling visual attention on  
112 images for individuals with ASD. Duan, et al. <sup>47</sup> combined high-level features (e.g., face size, facial features,  
113 face pose and facial expressions) and feature maps extracted from the SOTA saliency models to quantify visual  
114 attention on human faces in ASD. Their proposed approach reported higher performance when compared to  
115 other saliency models.

116           The remaining works used the Saliency4ASD dataset<sup>20,21</sup> for saliency estimation. For example, Fang, et  
117 al. <sup>48</sup> used U-net trained on a novel loss function for semantic feature learning, resulting in improved  
118 performance on some metrics. Wei, et al. <sup>49</sup> proposed a novel saliency prediction model for children with ASD.  
119 The fusion of multi-level features, deep supervision on attention maps and the single-side clipping operated on  
120 ground truths provided a boost in saliency prediction. Nebout, et al. <sup>50</sup> proposed a Convolutional Neural  
121 Network (CNN) with a coarse-to-fine architecture and trained using a novel loss function, achieving the best  
122 performance on most metrics when compared to general saliency models. Fang, et al. <sup>51</sup> proposed a model  
123 consisting of a spatial feature module and a pseudo-sequential feature module to generate an ASD-specific  
124 saliency map. Their model achieved the best performance on most metrics when compared to general saliency  
125 models and ASD-specific saliency models<sup>48-50</sup>. Finally, Wei, et al. <sup>52</sup> proposed a DNN architecture that enhances  
126 multi-level side-out feature maps using a scale-adaptive coarse-and-fine inception module. In addition, they  
127 designed a novel loss function to fit the atypical pattern of visual attention, resulting in SOTA performance.

128           This growing evidence suggests that researchers are starting to develop computational models that  
129 mimic the atypical visual attention on images of ASD individuals. However, there is still a huge gap in  
130 prediction performance as saliency prediction models trained on TD individuals do not generalise well on ASD  
131 individuals, as highlighted by Le Meur, et al. <sup>22</sup>. They revealed that current models trained on a TD dataset and  
132 fine-tuned on an ASD dataset perform well only on a small part of the ASD spectrum. To this end, they  
133 proposed two new eye-tracking datasets that cover a large part of the ASD spectrum.

Table 2 Saliency Prediction in ASD

Authors	Mean age (SD) in years	Sample size	Stimuli	Input used	Method	Conclusions
Duan, et al. <sup>46</sup>	7.8 (NR)	13	500 images	Image	They compared the performance of five different SOTA saliency models.	Transfer learning provides a useful approach to model the visual attention on images in individuals with ASD.
Duan, et al. <sup>47</sup>	ASD: 7.80 (2.10) TD: 8.00 (2.00)	ASD: 13 TD: 15	VAFA dataset: 300 images from open-source dataset <sup>53</sup> that depict various emotions and then classified into six expressions: (generally positive, very positive, neutral, generally negative, very negative and complex expressions)	Image	They computed fixation distributions on different pre-defined AOIs. Afterwards, statistical analyses were performed to identify differences in visual attention of ASD and TD participants while looking at effects of face pose and facial expressions. Afterwards, they compared six different SOTA deep learning-based saliency models on the VAFA dataset.	CASNet achieved the best performance in terms of the prediction of atypical visual attention of ASD individuals.
Fang, et al. <sup>48</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Image	They developed a saliency model based on the U-Net architecture. They also designed a new loss function called Positive and Negative Equilibrium Mean Square-Error that is used to determine model convergence.	Their model achieved higher performance on some metrics when compared to general saliency models.
Wei, et al. <sup>49</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Image	They first extracted multi-level features and combined these features using a fusion layer to output a saliency map. Deep supervision on the predicted saliency map was implemented to train the deeper layers of the network. They also utilised a single-side clipping approach to highlight regions that are mostly viewed by the participants.	Their model achieved the best performance on different metrics when compared to general saliency models.
Nebout, et al. <sup>50</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Image	They developed a two-stream network that extracts fine-scale	Their model achieved the best performance on most metrics



					and contextual information from the input image and the downscaled input image, respectively. Afterwards, a series of convolutional operations and concatenation is implemented to generate the saliency map.	when compared to general saliency models.
Fang, et al. <sup>51</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Image	They modelled the dynamic nature of human visual attention using a two-stream model that consists of a CNNs and a series of convolutional LSTM layers.	Their model achieved the best performance on most metrics when compared to general saliency models and ASD-specific saliency models <sup>48-50</sup> .
Wei, et al. <sup>52</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Image	They first extracted multi-level features from the input image. Afterwards, they passed it to a scale-adaptive coarse-and-fine inception module for a richer representation. These features are then combined using a feature fusion module and passed to a refinement and integration module. To better learn the atypical visual attention of ASD individuals, they developed a discriminative region enhancement loss.	Their approach achieved the best performance on different metrics when compared to general saliency models and ASD-specific saliency models <sup>48-50</sup> . Their experiments showed that their novel loss function improved the performance of other models in predicting atypical visual attention of ASD participants.
Le Meur, et al. <sup>22</sup>	Saliency4ASD MIE Fo and MIE No	Saliency4ASD MIE Fo and MIE No	Saliency4ASD MIE Fo and MIE No	Image	They compared six different saliency prediction models and analyse their saliency prediction performance in Saliency4ASD, MIE Fo and MIE No datasets.	Their results showed that current saliency models do not generalise well on ASD-specific dataset, hoping to raise awareness that researchers need different approaches to model the atypical visual attention of ASD people.

135 AOI: Area Of Interest, ASD: Autism Spectrum Disorder, LSTM: Long Short-Term Memory, NR: Not reported, SD: Standard deviation, SOTA: State-of-the-art, TD:  
136 Typically Developing

## 137 **Eye-tracking on static stimuli for ASD diagnosis**

138 As discussed in the previous section, it has been found that ASD participants exhibit atypical visual  
139 attention. As shown in Table 3, researchers explored the possibility of using the eye-tracking paradigm during  
140 image-viewing tasks to identify individuals with ASD. The earliest works explored different handcrafted  
141 features and ML models for ASD diagnosis. For example, Wang, et al.<sup>54</sup> used features extracted from images  
142 followed by a Support Vector Machine (SVM), while Yaneva, et al.<sup>55</sup> explored logistic-regression classification  
143 algorithms for detecting high-functioning ASD in adults. Liu, et al.<sup>34</sup> proposed a ML framework based on the  
144 frequency distribution of eye movements recorded during a face recognition task to identify individuals with  
145 ASD. The recent advances in deep learning (DL) also helped researchers better extract discriminative features  
146 from images. For example, Jiang and Zhao<sup>33</sup> used a DL approach followed by an SVM to distinguish  
147 individuals with ASD.

148 The succeeding works used the Saliency4ASD dataset<sup>20,21</sup>. Startsev and Dorr<sup>56</sup>, Arru, et al.<sup>57</sup> extracted  
149 features from the eye-tracking data and the input image and trained a random forest for ASD classification.  
150 Their analysis revealed that images that contain multiple faces provide significant differences in visual attention  
151 between ASD and TD individuals. Wu, et al.<sup>58</sup> proposed two machine learning approaches based on synthetic  
152 saccade generation and image classification with similar performance in terms of accuracy and AUC. Tao and  
153 Shyu<sup>59</sup> proposed a combination of CNN and long short-term memory (LSTM) networks to classify ASD and  
154 TD individuals. Exploiting a similar architecture, Chen and Zhao<sup>43</sup> proposed a multimodal approach to utilise  
155 information from behavioural modalities captured during photo-taking and image-viewing tasks, resulting in  
156 higher performance in both modalities. Using an additional dataset that contains people looking at other  
157 people/objects in the scene, Fang, et al.<sup>60</sup> proposed a DNN that achieved a higher accuracy when compared to a  
158 previous model<sup>33</sup>. Rahman, et al.<sup>61</sup> used several saliency prediction models and compared the performance of  
159 SVM and XGBoost. Observing that not all images highlight significant differences in visual attention between  
160 ASD and TD participants, Xu, et al.<sup>62</sup> used structural similarity between ASD and TD saliency maps to identify  
161 a subset of images in which a new bio-inspired metric was applied to identify ASD participants. Wei, et al.<sup>63</sup>  
162 proposed a dynamic filter and spatiotemporal feature extraction for ASD diagnosis, achieving the highest  
163 accuracy and similar specificity and AUC scores when compared to previous models<sup>56-59</sup>. Liaqat, et al.<sup>64</sup>  
164 proposed two ML approaches that include a branched MLP approach and an image-based approach for ASD  
165 classification and found that the latter approach resulted in slightly better performance. Mazumdar, et al.<sup>65</sup>  
166 extracted different handcrafted and DL features and compared 23 ML algorithms to identify individuals with

167 ASD. Their results were among the top 4 performing models across different metrics when compared to  
168 previous models<sup>56,59,64</sup>.

Table 3 Eye tracking on static stimuli for ASD diagnosis

Authors	Mean age (SD) in years	Sample size	Stimuli	Input used	Method	Conclusions
Wang, et al. <sup>54</sup>	ASD: 30.80 (11.1) TD: 32.30 (10.40)	ASD: 20 TD: 13	700 images from the OSIE dataset	Pixel-, object-, and semantic-level features extracted from the image. In addition, the image centre and background, as well as the ground-truth fixation maps were used.	Using the extracted features, they implemented an SVM to generate feature weights that were then combined to predict human fixation maps. They also conducted statistical analysis to investigate the atypical visual attention of ASD participants.	Their approach reported high performance in predicting the visual attention of both ASD and TD group. Their results showed that ASD group had increased biased towards the image centre, background and pixel-level, but reduced biased towards objects and semantic content of the image.
Yaneva, et al. <sup>55</sup>	Study 1: ASD: 37.00 (9.14) TD: 33.60 (8.60) Study 2: ASD: 41.00 (14.00) TD: 32.20 (9.90)	Study 1: ASD: 15 TD: 15 Study 2: ASD: 19 TD: 19	Study 1: 6 webpages with increasing visual complexity (e.g., low, medium, high) and 2 webpages in each category. Study 2: 8 randomly selected webpages from a list of top 100 websites, ensuring that there are 4 low visual complexity and 4 high visual complexity content.	Different computed eye-tracking variables (e.g., number of fixations, time to first look at an AOI) and non eye-tracking data-related variables (e.g., gender, visual complexity)	They computed eye-tracking related variables on different pre-defined AOIs. Afterwards, they trained several logistic regression classifiers using different combinations of the feature set for ASD classification.	Their results suggest that atypical visual attention of ASD individuals can be used as a biomarker for classification. They found differences in the information processing of ASD participants, regardless of specific information-location instructions across different time conditions. They also found that stimuli content and granularity have an impact on classification accuracy, while the stimuli complexity and gender do not exhibit the same effect.
Liu, et al. <sup>34</sup>	ASD: 7.90 (1.45) TD-Age Matched: 7.86 (1.38) TD-IQ Matched: 5.74 (1.01)	ASD: 29 TD-Age Matched: 29 TD-IQ Matched: 29	12 photos of adult Chinese female faces and 12 Caucasian female faces. 6 were used for memorisation task and 18 were used for a recognition task of	Frequency distribution of the visual attention of participants were computed.	They first quantised the fixation distribution of all participants using the k-means algorithm to generate cluster centroids. Afterwards, given a sequence of fixation locations, they assigned the cluster centroid closest to a	Their results showed a promising performance in classifying ASD participants based on visual attention on human faces.

			the 6 memorised faces.		participant's fixation location and counted the frequency of cluster assignments. This process was repeated on all the images and an SVM classifier was used for classification.	
Jiang and Zhao <sup>33</sup>	Same as Wang, et al. <sup>54</sup>	Same as Wang, et al. <sup>54</sup>	Same as Wang, et al. <sup>54</sup>	Images (and corresponding rescaled images) with the associated eye-tracking data of the participant	First, image selection using Fisher score ranking was implemented to reduce the number of input images from 700 to 100. Afterwards, each image and its corresponding rescaled image were passed to a two-branch VGG-16 network. The extracted features were then concatenated and used to predict the difference of fixation maps. Afterwards, a latent representation in the model was used for classification using SVM.	There was no direct comparison with other models since their model was one of the first to use eye-tracking for ASD classification. Nevertheless, the authors compared their approach with similar work that used different groups of subjects and input data and received the highest performance across different metrics.
Startsev and Dorr <sup>56</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Images with the associated eye-tracking data of the participant, including fixation durations.	First, they computed features extracted from the eye-tracking data and the input image. Afterwards, they trained a random forest for classification.	Their analysis revealed that images that contain multiple faces provide significant differences in visual attention between ASD and TD individuals.
Wu, et al. <sup>58</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Images with the associated eye-tracking data of the participant, including fixation durations.	They developed two networks: Synthetic saccade approach: a synthetic data generated by a scanpath model is aligned with the real eye-tracking data. Distance measures were then computed on these two data. Afterwards, different eye-tracking statistics were	Their experiments showed that both approaches resulted in similar classification performance in terms of accuracy and AUC.

					concatenated and used as features for MLP classification. Image-based approach: the real eye-tracking data were converted into an image. Afterwards, features were extracted from the input stimulus and the converted image and used as features for classification.	
Arru, et al. <sup>57</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Images with the associated eye-tracking data of the participant, including fixation durations.	First, they extracted features extracted from the image, eye-tracking data and bias towards the image centre. Afterwards, they trained a random forest that uses a bagging algorithm for classification.	Their results suggested that scene analysis, such as determining the objects attended by participants, could provide better results.
Tao and Shyu <sup>59</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Images with the associated eye-tracking data of the participant, including fixation durations.	First, they used a saliency model to generate a saliency map for a given image. Afterwards, square patches centred around the participant's fixations were extracted from the predicted saliency map. These patches were then passed to a CNN for feature extraction. The gaze duration associated with a patch location is concatenated with the extracted patch features and sequentially passed to an LSTM network followed by an FCL for classification.	Their results achieved an accuracy of 74.22% on the validation set and 57.90% on the test set.
Chen and Zhao <sup>43</sup>	Photo-taking task: NR Image-viewing task: NR	Photo-taking task: ASD: 22 TD: 23 Image-viewing	Photo-taking task: First-person photo taken by the participant	Photo-taking task: First-person photo taken by the participant	Photo-taking task: Given a sequence of photos taken by the participant, features are extracted using a CNN	Their results had the highest accuracy performance when compared to other models <sup>33,34</sup> .

	Saliency4ASD	task: ASD: 20 TD: 19  Saliency4ASD	Image-viewing task: 700 images from the OSIE dataset  Saliency4ASD	Image-viewing task: Images with the associated eye-tracking data of the participant.  Saliency4ASD	network and passed into a global average pooling layer. The sequence of image features is passed into an LSTM network and an FCL for classification. Image-viewing task: Given an image, features are extracted using a CNN network. Afterwards, using the associated eye-tracking data, features are extracted around the fixation location. The sequence of extracted features is then passed into an LSTM network and a FCL for classification. The authors also used multi-modal distillation to train both models.	
Fang, et al. <sup>60</sup>	Saliency4ASD GazeFollow4ASD: ASD: 9.60 (NR) TD: 8.90 (NR)	Saliency4ASD GazeFollow4ASD: ASD: 8 TD: 10	Saliency4ASD GazeFollow4ASD: Images that contain people looking at other people/objects in the scene	Saliency4ASD GazeFollow4ASD: Images with the gaze-following prior map indicating the eye locations of the people in the image and their gaze locations	First, they used a dilated CNN to extract coarse feature maps from the input image. Afterwards, these feature maps are passed to a convolutional LSTM network to generate enhanced features. A fusion layer is used to add the gaze-following prior map and a series of CNN layers is used to generate a difference of fixation maps. A latent representation in the model is passed to two FCLs for classification.	Their results had the highest accuracy performance when compared to a model <sup>33</sup> submitted to Saliency4ASD.
Rahman, et al. <sup>61</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Images with the associated eye-tracking data of the participant.	First, they used seven different saliency prediction models on a given image and computed evaluation metrics	Their model reported a higher performance compared to a previous SOTA model <sup>43</sup> for ASD

					between the predicted saliency and the recorded eye tracking data of the participant. This process is repeated for all the viewed images. The evaluation results for each saliency prediction model were concatenated. This feature representation was passed to an SVM and XGBoost for comparison of classification performance.	classification.
Xu, et al. <sup>62</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Images with the associated eye-tracking data of the participant.	Using structural similarity, they selected a subset of images that resulted into significant differences in visual attention of ASD and TD participants. Afterwards, they developed a bio-inspired metric that classifies ASD using the eye-tracking data.	Their results suggest that screening the photos to be viewed by participants and eventually used for classification is necessary to increase the model accuracy.
Wei, et al. <sup>63</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Images with the associated eye-tracking data of the participant.	First, an image encoder was used to extract visual features. Afterwards, the associated eye-tracking data of the participant was used as an input to three branches: (1) embedding layer to extract features (2) field of view maps generator layer that is composed of a spatial attention mechanism and LSTM network to extract spatiotemporal features (3) dynamic filters generator layer that uses CNNs. A final two FCLs were used for classification.	Their results had the highest accuracy and similar specificity and AUC scores when compared to models <sup>56-59</sup> submitted to Saliency4ASD.
Liaqat, et	Saliency4ASD	Saliency4ASD	Saliency4ASD	Images with the	They developed two	The image-based approach



al. <sup>64</sup>				associated eye-tracking data of the participant	networks: Branched MLP approach: it consists of a three-branch network that processes three different kinds of features: (1) a synthetic saccade is generated using a scanpath model, (2) a real scanpath and (3) statistical features. These features are passed to a series of MLPs for classification. Image-based approach: it consists of a two-branch network that extracts features from the input image and the eye tracking data and uses a final classification layer.	resulted in slightly better results than the branched MLP approach.
Mazumdar, et al. <sup>65</sup>	Saliency4ASD	Saliency4ASD	Saliency4ASD	Images with the associated eye-tracking data of the participant.	They computed features extracted from the image, eye-tracking data and centre bias of participants. Afterwards, they trained 23 different classifiers, such as decision trees, naïve bayes classifier, SVM, nearest neighbour classifier, and ensemble-based classifiers.	Their results were among the top 4 performing models across different metrics when compared to models <sup>56,59,64</sup> submitted to Saliency4ASD.

170 AOI: Area of Interest, ASD: Autism Spectrum Disorder, AUC: Area Under the Curve, CNN: Convolutional Neural Network, FCL: Fully-Connected Layer, IQ: Intelligence  
 171 Quotient, LSTM: Long Short-Term Memory, MLP: Multi-Layer Perceptron, NR: Not reported, SD: Standard deviation, SOTA: State-Of-The-Art, SVM: Support Vector  
 172 Machine, TD: Typically Developing

## 173 **Eye-tracking on dynamic stimuli for ASD diagnosis**

174 Prior research explored the possibility of using the eye-tracking paradigm during video-viewing tasks  
175 to identify specific neurological disorders. For example, Tseng, et al.<sup>66</sup> extracted low-level features from eye  
176 movement recorded from 15 minutes of videos and used an ML model to identify participants with attention  
177 deficit hyperactivity disorder, fetal alcohol spectrum disorder and Parkinson's disease. Although this work did  
178 not include ASD classification, it accentuates the efficacy of using eye-tracking on dynamic stimuli to identify  
179 the mental states of participants.

180 As shown in Table 4, there are recent works that utilise dynamic stimuli to differentiate ASD from TD  
181 subjects. Wan, et al.<sup>67</sup> investigated the difference in fixation times between ASD and TD children watching a  
182 10-second video of a female speaking. Their results revealed that fixation times at the mouth and body could  
183 significantly discriminate ASD from TD with a classification accuracy of 85.1%. Jiang, et al.<sup>68</sup> collected eye-  
184 tracking data during a dynamic affect recognition evaluation task, extracted handcrafted features and used a  
185 random forest classifier to identify ASD individuals. Zhao, et al.<sup>69</sup> collected eye-tracking data during a live  
186 interaction with an interviewer, extracted handcrafted features and employed four ML classifiers to identify  
187 individuals with ASD. These prior studies rely on handcrafted features that may provide less discriminative  
188 information between TD and ASD individuals.

189 Numerous studies employed an image classification approach based on a published dataset that  
190 contains the visualisation of eye-tracking data (i.e., scanpath images) of the participants during the experiment.<sup>45</sup>  
191 For example, Carette, et al.<sup>45,70</sup> used the raw pixel values as features and compared ML and DL algorithms for  
192 ASD classification. Their results revealed that DL algorithms achieved the highest performance when compared  
193 to ML models. Elbattah, et al.<sup>71</sup> trained a deep autoencoder and used a k-means clustering approach on the  
194 learned latent representation to identify clusters of participants. Their analysis revealed that an identified cluster  
195 contained a high percentage of ASD participants, suggesting that the algorithm can be used for ASD  
196 classification. Using a similar unsupervised learning approach, Akter, et al.<sup>72</sup> performed k-means clustering to  
197 divide the dataset into 4 groups and compared different ML models to identify participants with ASD. Cilia, et  
198 al.<sup>73</sup> used CNN and a fully-connected layer to predict ASD participants. Similarly, Kanhirakadavath and  
199 Chandran<sup>74</sup> compared Principal Component Analysis (PCA) and CNN for feature extraction and different ML  
200 and DL models for ASD classification. Gaspar, et al.<sup>75</sup> performed additional image augmentation to generate  
201 more training data. Afterwards, they used a kernel extreme learning machine optimised using the Giza Pyramids  
202 Construction metaheuristic algorithm to identify ASD individuals. Their approach achieved higher performance

203 when compared to ML approaches. Ahmed, et al. <sup>76</sup> compared ML, DL and a combination of both approaches  
204 for ASD diagnosis. The results in these prior studies suggest that DL models for feature extraction and ASD  
205 classification perform better when compared to traditional ML approaches.

206 There are also prior studies that explored the use of dynamic stimuli that are effective in evoking  
207 significant differences in visual attention of ASD and TD participants. For example, de Belen, et al. <sup>14</sup> used the  
208 GeoPref Test<sup>11,12</sup> in EyeXplain Autism, a system for eye-tracking data analysis, automated ASD prediction and  
209 interpretation of deep learning network predictions. Recently, Oliveira, et al. <sup>15</sup> used similar video stimuli,  
210 trained a visual attention model and utilised an ML model to identify individuals with ASD. Fan, et al. <sup>77</sup>, Fang,  
211 et al. <sup>78</sup> used biological motion stimuli and different ML classifiers for ASD diagnosis. Using a stimulus for  
212 initiating joint attention, Carette, et al. <sup>79</sup> extracted features related to saccadic movement (e.g., amplitude,  
213 velocity, acceleration) and trained an LSTM network to predict three diagnostic groups (i.e., ASD, TD,  
214 unclassified). Putra, et al. <sup>80</sup> collected eye-tracking data during Go/No-Go tasks, identified spatial and auto-  
215 regressive temporal gaze-related features that differ significantly between ASD and TD participants and applied  
216 an AdaBoost meta-learning algorithm to identify participants with ASD.

217 Although previous studies utilised dynamic stimuli, the most common approach was to convert the  
218 participant's eye-tracking data into an image, potentially losing spatiotemporal information that can be  
219 leveraged for classification. In addition, this approach disregards the pixel information around the fixation, a  
220 crucial insight into what part of the stimuli attracts human attention. In this paper, we propose a DNN approach  
221 that utilises dynamic saliency prediction to identify individuals with ASD.

222 While previous works have investigated the feasibility of leveraging visual attention in identifying  
223 individuals with ASD, limited research has been conducted to explore the effectiveness of exploiting the  
224 dynamic visual attention of the participant in ASD classification. Our approach utilises eye-tracking data  
225 captured during a dynamic stimulus viewing task. Our approach follows a similar deep learning framework  
226 reported in the literature<sup>33</sup>, however it provides an extension from static stimuli, widening the diagnostic  
227 paradigm to include dynamic stimuli.

Table 4 Eye Tracking on Dynamic Stimuli for ASD Diagnosis

Authors	Mean age (SD) in years	Sample size	Stimuli	Input used	Method	Conclusions
Wan, et al. <sup>67</sup>	ASD: 4.60 (0.70) TD: 4.80 (0.40)	ASD: 37 TD: 37	Dynamic, 10-second video of a female actor speaking	Eye-tracking data of the participant	They defined ten AOIs and computed different fixation time ratio. Afterwards, they used SVM to determine which AOI can be used for classification.	They found that using fixation times at the mouth and body results in an ASD classification accuracy of 85.1%, sensitivity of 86.5% and specificity of 83.8%.
Jiang, et al. <sup>68</sup>	ASD: 12.74 (2.45) TD: 14.11 (5.09)	ASD: 23 TD: 35	Combination of static and dynamic stimuli	Dynamic stimuli with the associated eye-tracking data of the participant	They computed eye-tracking variables (e.g., response time, fixation locations, length, frequency, duration, saccadic amplitude) and extracted face features using a DL model. They then used RF for classification.	The combination of all the handcrafted and extracted features resulted in a classification accuracy of 72.5%. Using a soft voting approach, the classification accuracy increased to 86.2% in identifying ASD participants.
Zhao, et al. <sup>69</sup>	ASD: 8.30 (2.09) TD: 9.07 (2.25)	ASD: 19 TD: 20	Dynamic, structured face-to-face conversation with a female interviewer	Dynamic stimuli with the associated eye-tracking data of the participant	They computed visual fixation ratios in four pre-defined AOIs across four sessions and added five features on session length, resulting in 21 features. Afterwards, they compared combinations of these features using different ML classifiers (e.g., SVM, LDA, DT and RF).	Their model that used the total session length, percentage of visual fixation time on the mouth AOI and the percentage of visual fixation time on the body as features achieved the highest classification accuracy. Looking at a single feature, the total session length was an effective discriminative feature.
Carette, et al. <sup>45</sup>	All participants: 7.88 (NR)	ASD: 29 TD: 30	Combination of static and dynamic stimuli that depict naturalistic scenes, initiate joint attention and static face or objects.	They visualised the eye-tracking data of a participant as a scanpath image. Using the scanpath images, they converted it to a grayscale image and rescaled for further processing.	They defined the ASD classification as an image classification problem using a logistic regression model.	Their result achieved an AUC of 0.819 based on 10-fold cross validation.
Carette, et al. <sup>70</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	They defined the ASD classification as an image classification problem using several ML and ANN models.	Their MLP achieved the best performance when compared to ML models. They noted that there was no performance increase as the complexity is increased.

Elbattah, et al. <sup>71</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	They trained an autoencoder for feature extraction. Afterwards, they implemented a k-means clustering algorithm and analysed the cluster qualities in terms of ASD classification.	They showed that by using a clustering technique on the latent space representation in the autoencoder bottleneck, they could get a cluster that contains a high percentage of ASD participants, suggesting that the algorithm can be used for ASD classification.
Akter, et al. <sup>72</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Using the scanpath images, they implemented a k-means clustering algorithm to divide the data into four groups. They trained different ML models in each cluster for classification.	Their results showed that the MLP achieved the best performance on different metrics when compared to ML models.
Cilia, et al. <sup>73</sup>	ASD: 7.58 (2.50) TD: 8.00 (2.67)	ASD: 29 TD: 30	Similar to Carette, et al. <sup>45</sup>	Scanpath images	They developed a four-layer CNN interspersed with pooling layers and a final FCLs for classification.	Their model achieved an accuracy of around 90%, sensitivity of around 83% and a precision of around 80%.
Kanhirakadavath and Chandran <sup>74</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	They compared two frameworks: (1) PCA for feature extraction and different ML techniques for classification. (2) CNN for feature extraction and different numbers of FCLs for classification.	Their results showed that the deep learning approach achieved higher performance across different metrics when compared to the different ML approaches.
Gaspar, et al. <sup>75</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Scanpath images	Their approach is a kernel extreme learning machine that uses giza pyramids construction metaheuristic algorithm for kernel parameters optimisation. They compared this technique to other optimisation algorithms, as well as ML algorithms, in terms of classification accuracy.	Their proposed pipeline achieved the highest performance on different metrics when compared to other optimisation algorithms. In addition, their model achieved the highest performance on difference metrics when compared to other ML algorithms.
Ahmed, et al. <sup>76</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Scanpath images	They developed three models that are based on ML, DL and hybrid techniques for classification.	The highest performing model was the ANN that uses the features extracted from the snake algorithm trained for image segmentation.

de Belen, et al. <sup>14</sup>	All participants: 4.60 (0.50)	ASD: 17 TD: 17	Same as Pierce, et al. <sup>11</sup> , Pierce, et al. <sup>12</sup> , Moore, et al. <sup>13</sup>	Dynamic stimuli with the associated eye-tracking data of the participant	They trained a VAM and used SVM for classification.	Using different number of fixations, their model achieved an accuracy of 68%-100%, sensitivity of 57% -100% and specificity of 65%-100%.
Oliveira, et al. <sup>15</sup>	Range: 3 to 18	ASD: 76 TD: 30	Dynamic, similar to GeoPref that contains biological and geometric movements	Dynamic stimuli with the associated eye-tracking data of the participant	For the entire video duration, they created two sets (one for each group) that contain the aggregated fixation locations on each frame. They created a group-specific fixation map which was then used to train VAMs. Afterwards, an individual classification was performed based on the VAMs.	Their model achieved an average precision of 90%, average recall of 69% and average specificity of 93%.
Fan, et al. <sup>77</sup>	All participants: Range: 3 to 13	ASD: 21 TD: 47	Point-light biological motion animation with upright/inverted persons that perform different actions.	They defined 5 'zones' where the visual attention of the participant is allocated. Afterwards, they computed data distribution within these zones.	They used the fixation distribution in different zones to identify zones helpful for classification. They trained an SVM for classification.	Their method achieved an AUC of 0.95.
Fang, et al. <sup>78</sup>	Age range: ASD: 4 to 10 TD: 2 to 15	ASD: 33 TD: 50	Same as Fan, et al. <sup>77</sup>	Same as Fan, et al. <sup>77</sup>	Using the extracted features, they compared kNN, Gaussian Naïve Bayes and Nonlinear SVM for ASD classification.	Their results showed that the nonlinear SVM achieved higher performance than the other MLP approaches.
Carette, et al. <sup>79</sup>	All participants: 8 to 10	ASD: 17 TD: 15	Dynamic, an actor initiating bids of joint attention	Eye-tracking data of the participants	Different saccadic movement variables were calculated as input to a two-layer LSTM network for classification.	Their model was able to identify ASD participants from TD participants with an accuracy of 83%.
Putra, et al. <sup>80</sup>	ASD: 5.00 (0.60) TD: 4.60 (0.40)	ASD: 21 TD: 31	Dynamic, CatChicken game	Eye-tracking data of the participants	They extracted different features and used the AdaBoost metalearning algorithm.	Their approach achieved an accuracy of 88.6%.

229 ANN: Artificial Neural Network, AOI: Area Of Interest, ASD: Autism Spectrum Disorder, AUC: Area Under the Curve, CNN: Convolutional Neural Networks, DL: Deep  
230 Learning, DT: Decision Tree, FCL: Fully-Connected Layer, kNN: k-Nearest Neighbour, LDA: Linear Discriminant Analysis, LSTM: Long Short-Term Memory, ML:  
231 Machine Learning, MLP: Multi-Layer Perceptron, NR: Not reported, PCA: Principal Component Analysis, RF: Random Forest, SD: Standard deviation, SVM: Support  
232 Vector Machine, TD: Typically Developing, VAM: Visual Attention Model

## 233 **Eye-tracking in ASD risk and symptom severity prediction**

234           Although there has been a great deal of research on the use of eye-tracking in ASD diagnosis, relatively  
235 little research focus on other applications, such as automatically predicting the risk of ASD (e.g., low, medium  
236 and high) and symptom severity, as shown in Table 5. Nevertheless, previous studies provide insights into the  
237 potential use of eye tracking in symptom severity prediction. For example, Kou, et al. <sup>81</sup> found that a reduction  
238 in visual preference for social scenes is significantly correlated with the ADOS social affect score, which may  
239 be useful in severity prediction. On the other hand, Bacon, et al. <sup>82</sup> found that a higher visual preference of  
240 toddlers for geometric scenes is significantly correlated with later symptom severity at school age, further  
241 suggesting the clinical utility of eye tracking for ASD symptom severity prediction.

242           Recently, Revers, et al. <sup>16</sup> trained two computational models<sup>83</sup> to generate saliency maps of severe and  
243 non-severe groups and used the RELIEFF algorithm<sup>84</sup> to select the most important features for classification.  
244 Afterwards, a neural network was trained to identify symptom severity for each fixation made by the participant.  
245 The final prediction is considered to be severe if more than 20 fixations were classified as severe by the trained  
246 neural network. Their approach obtained an average accuracy of 88%, precision of 70%, sensitivity of 87% and  
247 specificity of 60% in predicting symptom severity.

248           In a slightly different problem, Canavan, et al. <sup>23</sup>Fabiano, et al. <sup>24</sup> proposed a method for predicting  
249 ASD risk using eye gaze and demographic feature descriptors (e.g., age and gender). Handcrafted features, such  
250 as average fixation duration and average velocity, were tested on four different classifiers, namely random  
251 forests, decision trees, partial decision trees and a deep forward neural network. Although their results with a  
252 maximum classification rate of 93.45% are promising, it is crucial to compare their handcrafted features to  
253 features learned by modern deep learning models and determine if the latter improves the risk prediction  
254 accuracy. In this paper, we present the same DNN approach we used in ASD classification to predict the level of  
255 ASD-related symptoms.

Table 5 Eye Tracking in ASD Risk and Symptom Severity Prediction

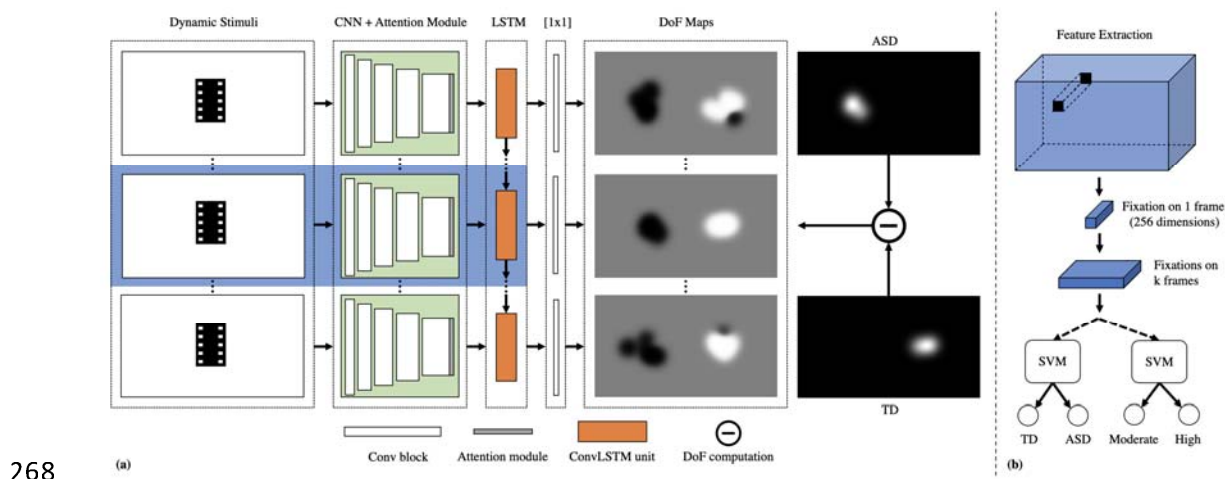
Authors	Mean age (SD) in years	Sample size	Stimuli	Input used	Method	Conclusions
Canavan, et al. <sup>23</sup> , Fabiano, et al. <sup>24</sup>	Two experiments: Experiment 1: Range: between 2 and 60 years old Experiment 2: Range: between 2 and 40 years old	Two experiments: Experiment 1: 257 with different risk types (low, medium, high and confirmed ASD) Experiment 2: 237 (subset of the first experiment)	Image and Video	They used the raw eye-tracking data (x and y locations), handcrafted features (e.g., average fixation duration, velocity), age and gender	They compared different ML and DL algorithms for ASD risk prediction.	Their approach achieved a maximum classification rate of 93.45%.
Revers, et al. <sup>16</sup>	Range: between 3 and 16 years old.	NSG: 49 SG: 39	Same as Pierce, et al. <sup>11</sup> , Pierce, et al. <sup>12</sup> , Moore, et al. <sup>13</sup>	They used the stimulus and the associated eye-tracking data of the participant.	They trained two computational models <sup>83</sup> to generate saliency maps of SG and NSG. Afterwards, they used RELIEFF algorithm to select features for classification. <sup>84</sup>	Their model achieved an average accuracy of 88%, precision of 70%, sensitivity of 87% and specificity of 60% for ASD symptom severity prediction.
Carette, et al. <sup>70</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	Same as Carette, et al. <sup>45</sup>	They defined the symptom severity prediction as an image classification problem using ANN models.	Their model achieved an average accuracy of around 83%. Their model was able to better identify TD participants compared to other ASD symptom severity. The prediction accuracy of symptom severity labels was 20% lower and worse for severe ASD participants.

ANN: Artificial Neural Network, ASD: Autism Spectrum Disorder, ML: Machine Learning, NSG: Non-Severe Group, SD: Standard Deviation, SG: Severe Group, TD: Typically Developing



## 259 Materials and methods

260 In this work, we used a data-driven approach to extract rich features learned from a dynamic stimulus  
261 to identify participants with autism and predict the level of ASD-related symptoms. In **Error! Reference source**  
262 **not found.**, an overview of the proposed approach is provided. The method is divided into different stages,  
263 including eye-tracking data collection, dynamic saliency detection trained on the difference of fixations between  
264 ASD and TD individuals, and SVM-based classification and severity prediction. This study was approved by the  
265 Human Research Ethics Committee of the University of New South Wales. Written informed consent was  
266 obtained from the parents/legally authorised representatives of the participants. All methods were carried out in  
267 accordance with relevant guidelines and regulations.



269 *Figure 1 Overview of the proposed feature learning/extraction, classification and symptom severity prediction*  
270 *approach. (a) Given a video input, per-frame features are learned using an end-to-end approach to predict the*  
271 *difference of fixation (DoF) maps; (b) Extracted features at fixated pixels from each fixation stage are cascaded*  
272 *and passed on to an SVM to identify individuals with ASD and predict the level of ASD-related symptoms.*

## 273 Eye-tracking

## 274 Participants

275 There were 57 children (9 females) in the ASD group and 17 children (9 females) in the TD group.  
276 Participants were matched by their age at the time of the study. 24 children in the ASD group were recruited  
277 from an Autism Specific Early Learning and Care Centre (ASELCC) and 33 children were recruited from the  
278 Child Development Unit (CDU) of a Children's Hospital. The TD children were recruited from a children's  
279 services preschool. All participants in the ASD group met the criteria for ASD based on the Diagnostic and  
280 Statistical Manual of Mental Disorders (DSM-5)<sup>85</sup> criteria and the diagnosis of ASD was confirmed using the

281 Autism Diagnostic Observation Schedule (ADOS), Second Edition<sup>86</sup>. Of the 57 ASD children, there were 24  
282 who showed high ASD-related symptoms and 33 had moderate symptoms. There are no specific exclusion  
283 criteria for the ASD group in this study. The TD group's exclusion criteria included known neurodevelopmental  
284 disorders, significant developmental delays and known visual/hearing impairments. No child had any visual  
285 acuity problems.

## 286 **Dynamic stimulus**

287 We used the GeoPref Test<sup>11,12</sup> dynamic stimulus, which has been shown to be an effective stimulus for  
288 detecting ASD subgroups. This stimulus consists of dynamic geometric images (DGIs) on one side and dynamic  
289 social images (DSIs) on the other. The DGIs were constructed from recordings of animated screen-saver  
290 programs. The DSIs were produced from a series of short sequences of children performing yoga exercises. It  
291 included images of children performing a wide range of movements (e.g., waving arms and appearing as if  
292 dancing). The stimulus contained a total of 28 different scenes and was presented in order, based on the  
293 originally published stimulus<sup>11,12</sup>. It has a resolution of 1281 x 720 pixels and contains a total of 1,488 frames,  
294 which is equivalent to 61 seconds of video playback.

## 295 **Eye-tracking apparatus and procedure**

296 Participants were tested using the Tobii X2-60 eye tracker and eye-tracking data was processed using  
297 Tobii Studio software to identify fixations and saccades. Eye movements were recorded at 60 Hz (with an  
298 accuracy of 0.5°) during the dynamic stimuli viewing. Each participant was seated approximately 60 cm in front  
299 of a 22" monitor with a video resolution of 1680 x 1050 pixels in a quiet room and shown dynamic visual  
300 stimuli in full-screen. A built-in five-point calibration in Tobii studio was completed before administering the  
301 task for accurate eye gaze tracking. The calibration procedure required gaze following on an image of an animal  
302 paired with auditory cues, starting with the centre of the screen, and moving across the four corners of the  
303 screen. The eye-tracking procedure was conducted during a clinical assessment or the intake assessment for  
304 entry to an early intervention program. Multiple attempts were made to ensure that the eye tracker has been  
305 calibrated properly for accurate data collection. Multiple attempts were also made to ensure that the participants  
306 were engaged during the experiment. As a result, depending on the capacity of the child, the procedure was  
307 conducted over 2 to 3 sittings or with smaller breaks in between. The overall clinical assessment and eye-  
308 tracking procedure were completed in approximately 2.5h per participant.

309

310

## 311 **Data processing and statistical analysis**

312 Tobii Studio's I-VT filter<sup>87</sup> was used to process the raw eye-tracking data, exclude random noise and  
313 define fixations for further analysis. More specifically, short fixations (<100ms) were discarded and adjacent  
314 fixations (75ms, 0.5°) were merged. Trials were excluded if the total fixation duration was less than 15 seconds.  
315 That is, to be included, the participant should be looking at the stimulus for approximately 25% of the entire  
316 video duration. Once included, the eye-tracking data captured during the entire length of the stimulus are used  
317 for training and evaluation.

318 A calibration quality assessment was performed to rule out the possibility of eye-tracking data quality  
319 as a confounding factor. In this assessment, a toy accompanied by a sound was used to attract the participants'  
320 gaze to the calibration point in the middle of the screen. The mean distance between the detected fixation  
321 locations and the calibration point was calculated as a measure of accuracy. A t-test showed no significant  
322 difference between the groups, suggesting that data quality did not differ between the two groups:  $t(64) = -$   
323  $0.445$ ,  $p = .658$ , ASD: 45.89 pixels (22.67), TD: 48.76 pixels (19.00).

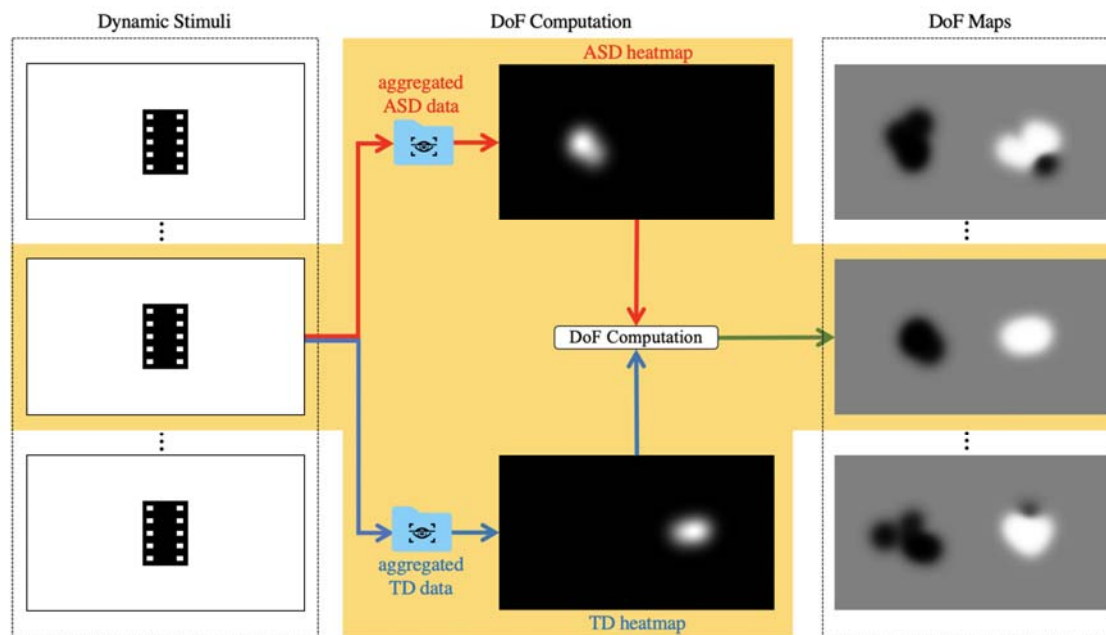
324 An additional data quality assessment was performed to determine the overall nature of the visual  
325 attention of the participants to the stimuli. A t-test showed no significant difference in visual attention between  
326 groups:  $t(72) = 0.011$ ,  $p = .991$ , ASD: 37.13 seconds (12.03), TD: 37.10 seconds (8.07). These analyses of  
327 quality suggest that it is unlikely that differences in data quality and general visual attention influenced the  
328 results.

329 An independent-samples t-test was used to investigate differences in visual attention across two groups  
330 for diagnosis (ASD vs. TD) and severity prediction (moderate vs. severe). All statistical analysis was performed  
331 in IBM SPSS Statistics Version 26.

## 332 **Computation of per-frame saliency maps**

333 Saliency detection models are typically optimised to detect salient features in a scene. They are trained  
334 on a probability distribution of eye fixations, called the fixation map. The per-frame fixation maps of each  
335 participant group were generated from the eye movement data collected in the study. For a given frame, all  
336 fixation points of the children in each group were overlaid in a binary map, in which the fixation points were set  
337 to 1 on a black background (value set to 0). The resulting per-frame fixation maps were smoothed with a

338 Gaussian kernel (bandwidth = 1°) and normalised by the sum to generate per-frame visual attention heatmaps  
339 (labelled ASD and TD heatmaps in **Error! Reference source not found.**).



340

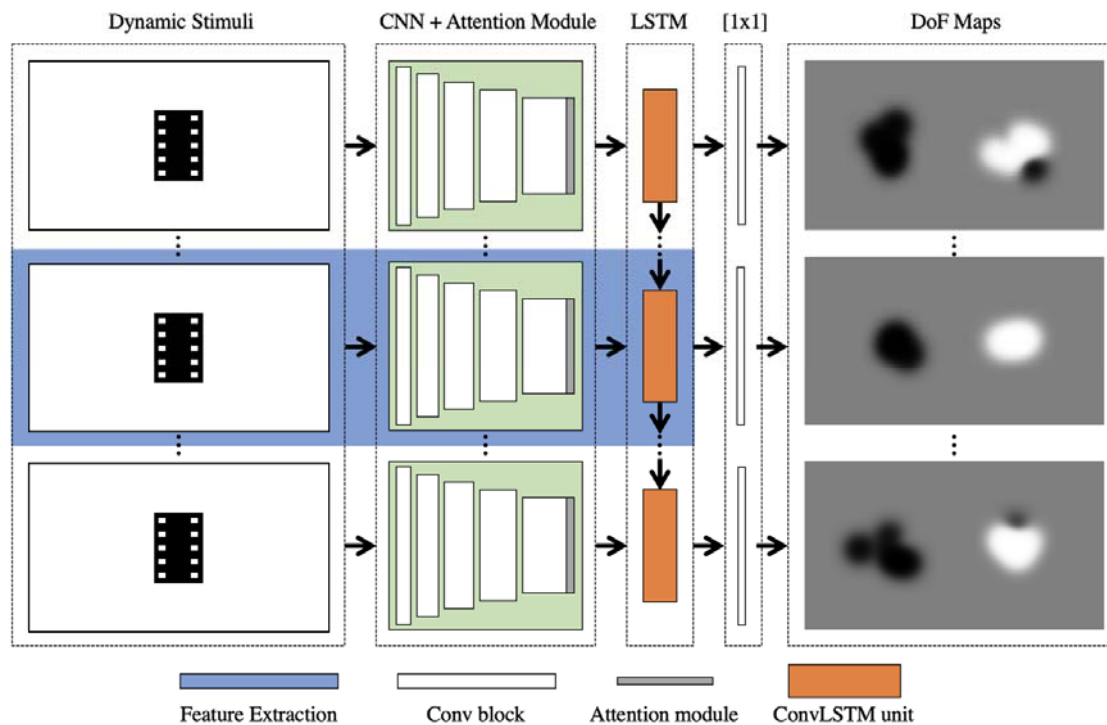
341 *Figure 2 Difference of Fixation (DoF) computation*

## 342 **Computation of per-frame difference of fixation (DoF) maps**

343 Similar to Jiang and Zhao<sup>33</sup>, our network was optimised on the difference of fixation (DoF) maps,  
344 highlighting the difference in visual attention between TD and ASD individuals. Since our approach uses a  
345 dynamic stimulus, we predict DoF maps on each frame. In particular, let  $I_{ASD}$  and  $I_{TD}$  be the fixation maps for the  
346 ASD and TD groups, respectively. The DoF map of a frame is computed as:

347 where  $I_{DoF}$  is a pixel-wise subtraction of fixation maps and  $\sigma$  represents the standard deviation of  $I$ .

348 The resulting DoF maps highlight the difference in visual attention between ASD and TD individuals (refer to  
349 **Error! Reference source not found.**). The white regions of the DoF map illustrate the visual attention of TD  
350 individuals while the black regions are for ASD individuals. Note that this is the opposite of the DoF  
351 computation elsewhere<sup>33</sup>. This also resulted in better training performance compared to DoF maps that highlight  
352 more fixations of the ASD group.



353

354 *Figure 3 Learning Difference of Fixation Maps*

### 355 **Per-frame prediction of difference of fixation maps**

356 As shown in **Error! Reference source not found.**, ACLNet<sup>88</sup>, one of the best models available for  
 357 dynamic saliency detection, is used for feature extraction. It consists of a CNN-LSTM network with an attention  
 358 mechanism to enable fast, end-to-end saliency prediction. Since ACLNet already contains an attention network  
 359 trained on TD individuals, we trained and fine-tuned our model with DoF maps that highlight more fixations of  
 360 the TD group.

361 Our model was optimised using the following loss function<sup>89</sup> which considers three different saliency  
 362 evaluation metrics instead of the binary-cross entropy loss used before<sup>33</sup>. We denote the predicted difference of  
 363 fixation map as  $\hat{Y}$  and the ground truth saliency map as  $Y$ . Our loss function  
 364 combines Kullback-Leibler (KL) divergence, the Linear Correlation Coefficient (CC) and the Normalised  
 365 Scanpath Saliency (NSS) similar to prior work<sup>88</sup>:

366  $\mathcal{L}$  is widely used for training saliency models and is computed by:

$$\mathcal{L} = \mathcal{L}_{KL} + \mathcal{L}_{CC} + \mathcal{L}_{NSS}$$

367  $\mathcal{L}_{CC}$  measures the linear relationship between  $Y$  and  $Q$ :

$$L_{CC}(Y, Q) = -\frac{cov(Y, Q)}{\sigma(Y)\sigma(Q)}$$

368 where  $cov(Y, Q)$  is the covariance of Y and Q while  $\sigma$  is the standard deviation.

369  $L_{NSS}$  is defined as:

$$L_{NSS}(Y, Q) = -\frac{1}{N} \sum_x \bar{Y}(x) \times Q(x)$$

370 where  $\bar{Y} = \frac{Y - \mu(Y)}{\sigma(Y)}$  and  $N = \sum_x Q(x)$ . It computes the mean of scores from the normalised saliency map  $\bar{Y}$  at

371 the predicted DoF maps Y.

## 372 **Training protocol**

373 Our classification and severity prediction models are iteratively trained with sequential DoF maps and  
374 image data. We train the model by using a loss defined over the predicted dynamic saliency maps from  
375 convLSTM. Let  $\{Y_t^d\}_{t=1}^T$  and  $\{Q_t^d\}_{t=1}^T$  denote the predicted dynamic saliency maps and continuous difference of  
376 fixation maps. We minimise the following loss:

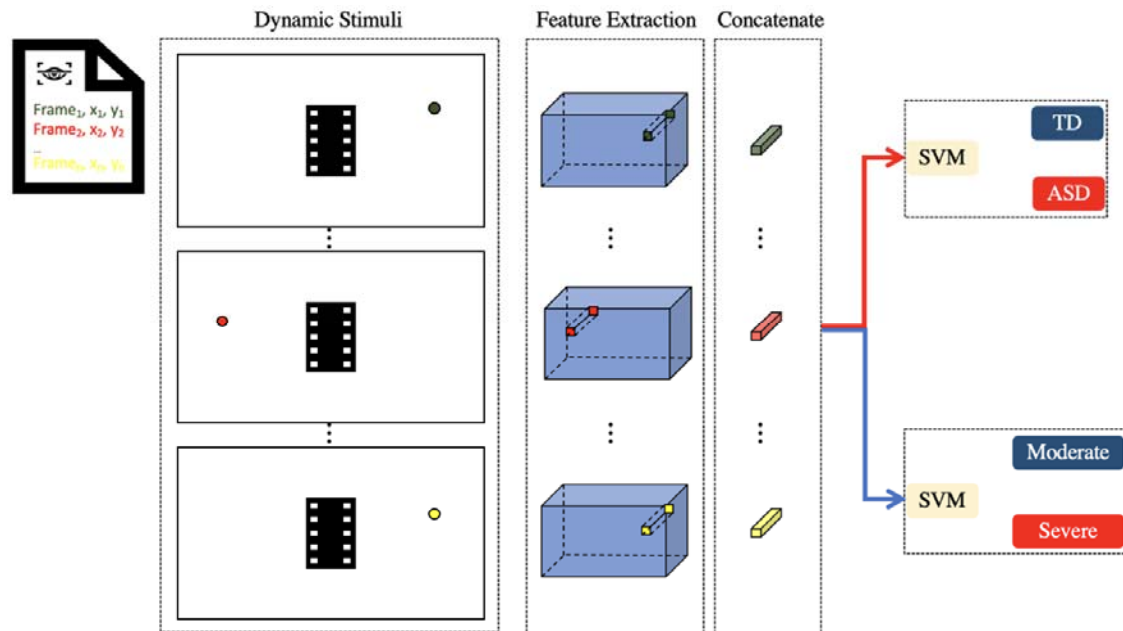
$$L^d = \sum_{t=1}^T L(Y_t^d, Q_t^d)$$

377 The parameters of ACLNet are initialised to the pre-trained parameters<sup>88</sup>. The network is then fine-tuned on the  
378 current dataset.

## 379 **ASD classification and symptom severity prediction**

380 Once the model has been trained to predict DoF maps of ASD and TD individuals from a given  
381 dynamic stimulus, feature extraction and classification are performed, with **Error! Reference source not**  
382 **found.** illustrating the process<sup>14</sup>. Based on the eye-tracking data, we determined the fixation positions and the  
383 corresponding frames in which they were recorded. Each saccade-fixation pair was considered a fixation stage.  
384 For each fixation stage, features were extracted from the corresponding fixation position on the feature map  
385 obtained from the convLSTM output (note that the convLSTM output is upsampled 4 times before extracting  
386 the feature map). More specifically, given a frame where a fixation has been identified, the feature map at the  
387 corresponding fixation is extracted, which results in a 256-dimensional feature vector at each fixation. For a  
388 corresponding number of fixation stages, feature vectors for all fixations are concatenated in their temporal  
389 order starting from the first fixation to the last fixation stage. This serves as the feature space in which

390 classification is performed. If there were fewer identified fixations, zeros are appended at the end. We explored  
391 the number of fixation stages that provided the best performance.



392

393 *Figure 4 Feature Extraction and Classification*

394 A linear decision boundary between ASD and TD individuals was determined by training an SVM on  
395 the extracted features. In addition, another SVM model was trained on the DoF maps of moderate and high ASD  
396 individuals to predict autism severity. We used the ADOS-2 calibrated severity scores (CSS) as ground truth to  
397 determine the ASD severity. Participants with ADOS CSS of 5-7 are considered to have moderate symptoms,  
398 while those with ADOS CSS of 8-10 are considered to have more severe (high) symptoms.

## 399 **Experimental setup**

### 400 **Training and testing protocols**

401 We report the model's performance on ASD classification and symptom severity prediction using  
402 leave-one-out cross-validation (LOOCV). Given the unbalanced nature and the limited number of samples in the  
403 dataset, LOOCV is used to provide an almost unbiased estimate of the probability of error<sup>90</sup>. In addition, it  
404 allows us to maximise the number of training samples per fold unlike in a k-fold validation approach. While a  
405 stratified k-fold cross-validation strategy may account for the group imbalance that is present in our dataset, it  
406 results in smaller training samples per fold. However, removing a single sample from the training set done in  
407 LOOCV also does not drastically change the class distribution. The combination of being able to use as much

408 training data as possible while also maintaining similar class distribution was the reason why we used LOOCV.  
409 The same evaluation approach has been employed in prior studies<sup>14,33,34,43,68,69</sup> in this application area.

## 410 **Implementation details**

411 We implemented our model in Tensorflow with Keras and Scikit-learn libraries. During the training  
412 phase, we fine-tuned the network with Adam optimizer and a batch size of one image for a total of 20 epochs.  
413 The learning rate was set to 0.0001. We did not perform any dropout and data augmentation. L2 regularisation  
414 with the penalty parameter  $C=1$  was used for SVM classification.

## 415 **Evaluation metrics**

416 We report on the performance of our model in terms of accuracy, sensitivity (i.e., true positive rate)  
417 and specificity (i.e., true negative rate) recorded at different numbers of fixations. Once the best number of  
418 fixations to be included in the classification was identified, the area under the receiver operating characteristic  
419 (ROC) curve and the confusion matrix were also computed. To obtain a meaningful area under the ROC curve  
420 (AUC) in an LOOCV, the output probability of the SVM for each fold (each consisting of just one subject) was  
421 saved and the AUC was computed on the set of these probability estimates. The computation of the confusion  
422 matrix was performed similarly using the predicted class to compare with the ground truth label.

## 423 **Computational load**

424 The entire training procedure for each video stimulus takes about 1 hour with two NVIDIA 2080 Super and a  
425 3.5GHz Intel processor (i7-7800X CPU). Once the model has been trained, feature extraction and SVM  
426 classification can be performed in less than 1 minute.

# 427 **Results**

## 428 **Datasets**

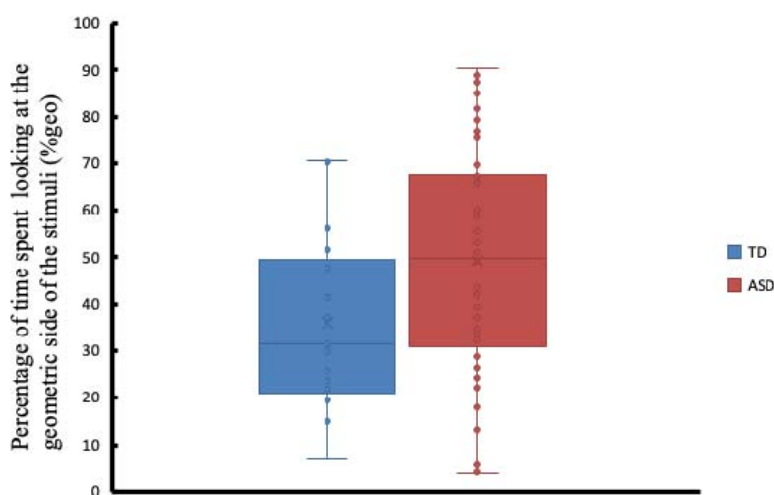
429 Children with ASD had a mean age of 4.63(standard deviation (SD) = 0.80) years and TD participants  
430 also had a mean age of 4.61 (SD = 0.47) years. There was no significant difference in age between the ASD and  
431 TD groups,  $t(72) = 0.009$ ,  $p = 0.993$ .

## 432 **Eye-tracking data analysis**

### 433 **ASD Classification**



434 It was previously shown that ASD individuals with severe symptoms tend to fixate more on the  
435 geometric stimuli than the social stimuli<sup>11,12</sup>. Shown in **Error! Reference source not found.** are the %Geo  
436 values, the percentage of time spent looking at the dynamic geometric stimuli. %Geo values are computed by  
437 dividing the total fixation duration on the geometric stimuli by the total fixation duration on both geometric and  
438 social stimuli. Independent-samples t-test was used to compare %Geo for each diagnostic group. Similar to  
439 published results elsewhere<sup>11-13</sup>, ASD participants in our study were significantly more attracted to dynamic  
440 geometric images when compared to TD participants ( $t = 2.11, p < .0386$ ). On average, the ASD group spent  
441 49.37% (standard deviation (SD) = 24.14%) of their attention looking at the dynamic geometric images, while  
442 the TD group spent 35.97% (SD = 18.58%) of their attention.

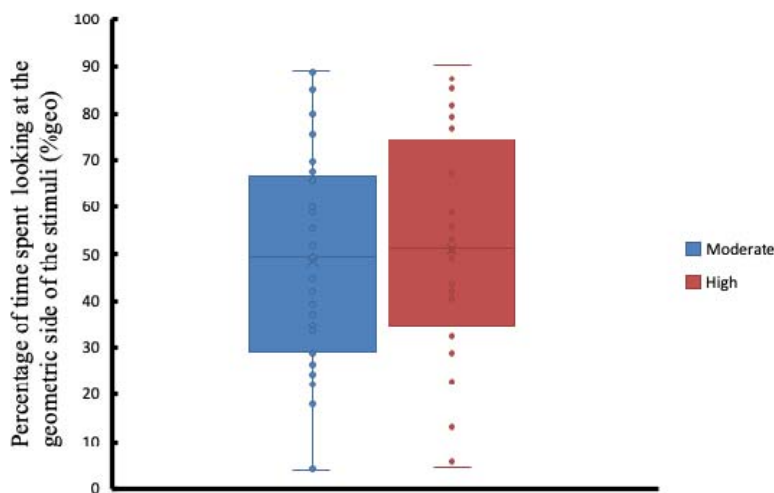


443  
444 *Figure 5 Comparison of the percentage of time spent looking at the dynamic geometric stimuli (%geo) between TD*  
445 *and ASD participants. Each box plot contains the interquartile range, the x marker corresponds to the mean value and the*  
446 *horizontal line inside correspond to the median. Each sample is also visualised using dot points.*

## 447 **ASD symptom severity prediction**

448 Shown in **Error! Reference source not found.** are the %Geo values, the percentage of time spent on  
449 looking at the dynamic geometric stimuli. There was no significant difference in the %Geo values between the  
450 moderate and severe ASD participants ( $t = 0.424, p < .6729$ ). On average, ASD participants with moderate  
451 symptoms fixated around 48.21% (SD = 23.82%) of their attention on the geometric stimuli. On the other hand,  
452 ASD participants with severe symptoms spent 50.98% (SD = 25.00%) of their attention looking at the geometric  
453 stimuli. We also performed pair-wise comparisons between the TD participants and the two ASD participant

454 groups (i.e., moderate and severe). There was a significant difference in the %Geo values between ASD  
455 participants with severe symptoms and TD participants ( $t = 2.096$ ,  $p < .0426$ ). On the other hand, there was only  
456 a trend toward a significant difference in the %Geo values between ASD participants with mild symptoms and  
457 TD participants ( $t = 1.846$ ,  $p < .0710$ ).



458

459 *Figure 6 Comparison of the percentage of time spent looking at the dynamic geometric stimuli (%geo) ASD participants*  
460 *with moderate and severe symptoms. Each box plot contains the interquartile range, the x marker corresponds to the mean*  
461 *value and the horizontal line inside correspond to the median. Each sample is also visualised using dot points.*

462

463 In recent years, it has been shown that stimuli that have both dynamic geometric and social images can  
464 reliably separate the visual attention of ASD and TD individuals. We contribute to the literature by showing that  
465 a DNN-based approach using dynamic stimuli can result in highly accurate ASD classification and even predict  
466 the level of ASD-related symptoms with promising performance.

## 467 **ASD classification performance**

468 In Figure 7, different performance metrics for ASD prediction on the GeoPref Test dynamic stimulus  
469 are shown. In Figure 7a the accuracy, sensitivity and specificity of the model as the number of fixations (i.e.,  
470 fixation length) increases are displayed. It can be observed that all measures generally increase as the number of  
471 fixations increases. In Figure 7b and Figure 7c, the receiver operating characteristics (ROC) curve and the  
472 confusion matrix of the model that reported the highest accuracy (i.e., using the optimal fixation length) in  
473 Figure 7a are shown. The area under the ROC curve (AUC) of our model is 0.96, significantly higher than

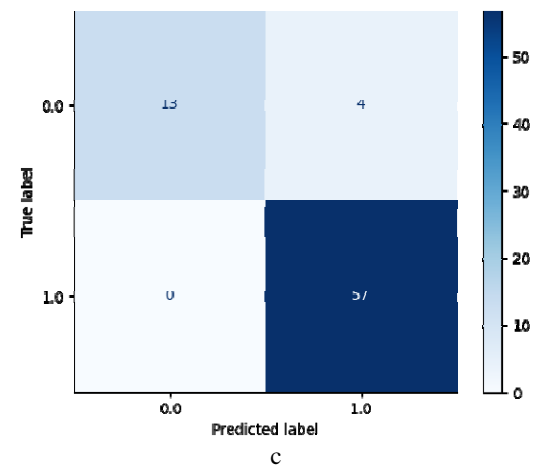
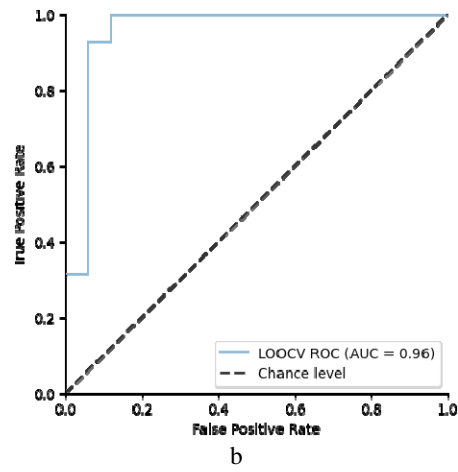
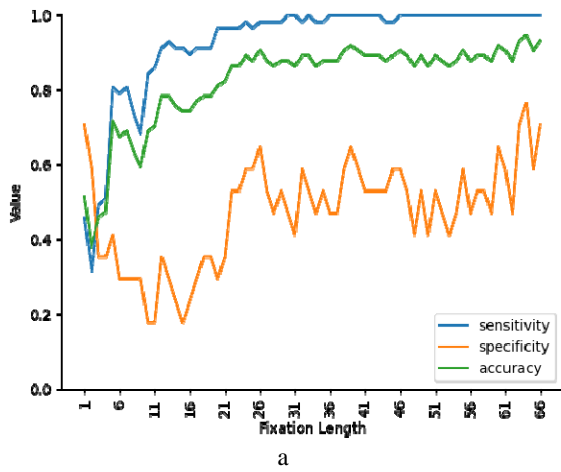
474 chance-level performance (AUC=0.5). Our model achieved the highest accuracy of 94.59% when 64 fixations  
475 were included in the analysis. The high sensitivity of our model (highest value = 100%) suggests that it can  
476 reliably identify ASD children. On the other hand, the specificity of our model (highest value = 76.47%)  
477 suggests that it can reliably identify children without the disorder. However, four (4) children were mistakenly  
478 flagged as having the disorder despite not having it.

### 479 **ASD severity prediction performance**

480 Similar to the results of the diagnosis prediction, it can be observed in Figure 8a that all performance  
481 measures for ASD severity prediction generally increase as the number of fixations (i.e., fixation length)  
482 increases. In Figure 8b and Figure 8c, the ROC curve and the confusion matrix of the model that reported the  
483 highest accuracy in Figure 8a are shown. Our model achieved the highest accuracy of 94.74% when 44  
484 fixations were included in the analysis. The area under the ROC curve (AUC) of our model is 0.99, significantly  
485 higher than chance-level performance (AUC=0.5). The high specificity of our model (highest value = 100%)  
486 suggests that it can reliably identify children with mild ASD. On the other hand, the high sensitivity of our  
487 model (highest value = 87.50%) suggests that it can reliably identify children with severe symptoms. However,  
488 three (3) children were mistakenly flagged as having severe diagnoses despite having milder symptoms.

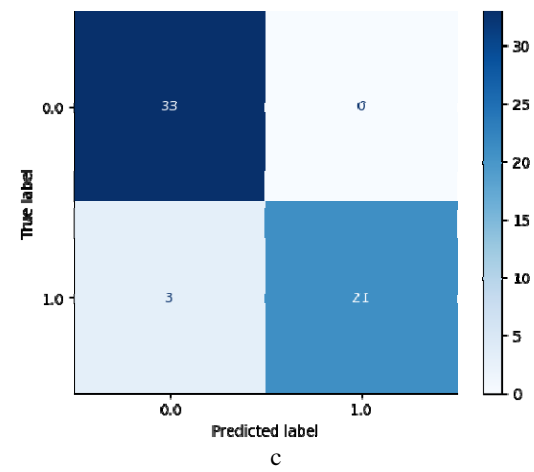
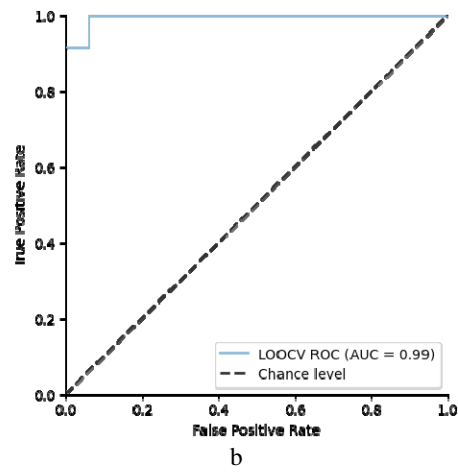
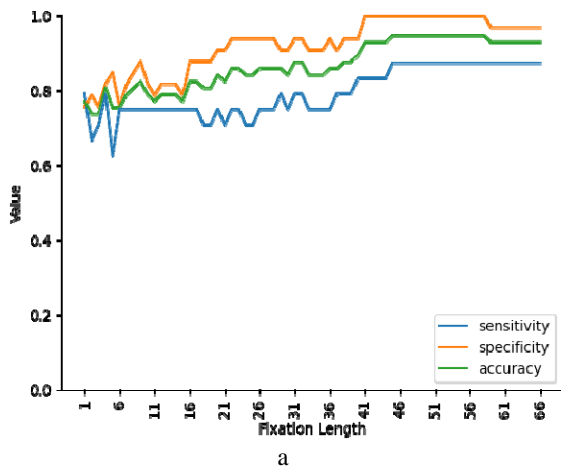
### 489 **Comparison with other approaches**

490 As outlined in the related work section, a straightforward comparison with previous approaches that  
491 utilise dynamic stimuli is not possible because the published dataset contains a visualisation of eye-tracking  
492 participants (i.e., scanpath images) rather than the stimuli used and the associated eye-tracking data that our  
493 model requires. Nevertheless, we compared our proposed approach with a simple thresholding method<sup>11-13</sup> and  
494 ML algorithms using handcrafted features<sup>23,24</sup>.



495  
496

Figure 7 Different performance metrics for ASD prediction. A.) the plot of the model's sensitivity, specificity and accuracy as the number of fixations (i.e., fixation length) increases. B.) the plot of the area under the receiving operating curve of the best-performing model. C.) the confusion matrix of the best-performing model.



497  
498

Figure 8 Different performance metrics for ASD symptom severity prediction. A.) the plot of the model's sensitivity, specificity and accuracy as the number of fixations (i.e., fixation length) increases. B.) the plot of the area under the receiving operating curve of the best-performing model. C.) the confusion matrix of the best-performing model.

## 499 **ASD Classification**

500 Following the cut-off of %Geo > 69% to determine ASD participants in a similar study<sup>11-13</sup>, we  
501 obtained a sensitivity of 22.80%, specificity of 88.23% and accuracy of 37.84%. The AUC obtained was 0.67.  
502 In comparison, our proposed model resulted in 77.2% higher sensitivity, 11.76% lower specificity and 56.75%  
503 higher accuracy when compared to solely utilising the %Geo values. Handcrafted features that include raw eye  
504 gaze points (x and y locations), average fixation duration, age and gender, were also used as input to a random  
505 forest regressor and a decision tree classifier for ASD classification similar to a previous study<sup>23,24</sup>. The random  
506 forest regressor achieved an accuracy of 72.97%, a sensitivity of 91.22% and a specificity of 0%. On the other  
507 hand, the decision tree classifier achieved an accuracy of 58.11%, a sensitivity of 70.18% and a specificity of  
508 17.65%.

509 Overall, our proposed model achieved the highest accuracy of 94.59%, the highest sensitivity of 100%  
510 and the second-best specificity of 76.47%. The comparison results in ASD classification suggest that our model  
511 better identified participants with ASD than the previous approaches.

512 *Table 6 ASD Classification Results Comparison with Prior Approaches*

Approach	Accuracy	Sensitivity	Specificity
Thresholding approach <sup>11-13</sup>	37.84%	22.80%	<b>88.23%</b>
Random forest regressor	72.97%	91.22%	0.00%
Decision tree classifier	58.11%	70.18%	17.65%
Ours	<b>94.59%</b>	<b>100%</b>	76.47%

## 513 **ASD symptom severity prediction**

514 We also used the same<sup>11-13</sup> cut-off of %Geo > 69% to identify ASD participants with severe symptoms  
515 and obtained a sensitivity of 25.00%, specificity of 78.79% and accuracy of 43.24%. The AUC obtained was  
516 0.54. Again, our proposed method showed promising results for severity prediction, resulting in a 62.50%  
517 increase in sensitivity, a 21.21% increase in specificity and a 51.5% increase in accuracy when compared to  
518 solely utilising the %Geo values. In comparison to our model, using handcrafted features and ML classifiers  
519 resulted in the same accuracy of 94.74%, slightly higher sensitivity of 91.67% and slightly lower specificity of  
520 96.97%.

521 Overall, our proposed model achieved the highest accuracy of 94.47%, the second-best sensitivity of  
522 87.50% and the highest specificity of 100%. The comparison results in ASD symptom severity prediction  
523 suggest that our model better identifies participants with moderate symptoms than the previous approaches.

524 *Table 7 ASD Symptom Severity Prediction Results Comparison with Prior Approaches*

Approach	Accuracy	Sensitivity	Specificity
Thresholding approach <sup>11-13</sup>	43.24%	25.00%	78.79%

Random forest regressor	<b>94.74%</b>	<b>91.67%</b>	96.97%
Decision tree classifier	<b>94.74%</b>	<b>91.67%</b>	96.97%
Ours	<b>94.74%</b>	87.50%	<b>100%</b>

## 525 Discussion

526 Over the past decade, eye-tracking studies have revealed significant differences in visual attention  
527 between ASD and TD individuals. This motivated researchers to leverage recent advances in saliency prediction  
528 when designing a more quantitative approach to ASD diagnosis, as well as risk and symptom severity  
529 prediction. In this context, researchers have explored the use of static and dynamic stimuli during free-viewing  
530 tasks. The most common approach in the literature comprised of a traditional two-stage method that consists of  
531 a feature extraction stage followed by a classification stage. Increasing evidence suggests that the DL-based  
532 approach produced more discriminative features when compared to ML-based approaches. Classification  
533 methods that utilise DL also resulted in better performance than ML models. The rapid advances in DL  
534 approaches and the increasing number of publicly available datasets may help further advance the literature and  
535 improve classification performance. In this paper, we utilised a combination of DL and ML approaches for ASD  
536 diagnosis and symptom severity prediction.

537 Unlike prior research that utilised dynamic stimuli and converted the participant's eye-tracking data  
538 into an image for classification, we propose a data-driven approach utilising a dynamic saliency model to extract  
539 discriminative features from the stimuli and an ML approach based on eye-tracking data to automatically  
540 identify individuals with ASD. In addition, we show that the same approach can predict the level of ASD-  
541 related symptoms in preschool children. Our approach to identifying children with ASD offers several  
542 advantages when compared to existing eye-tracking research. Most notably, our method only takes one minute  
543 of eye-tracking, a substantial decrease in recording time when compared to about 10 minutes required in  
544 previous studies<sup>33,34</sup>. While our method requires a substantially shorter amount of time, it is not a replacement  
545 for standard clinical assessments. Extensive experiments are necessary before the true clinical utility and  
546 usability of our proposed method can be realised.

547 Our results support other studies<sup>11-13</sup> that found a significant difference in the overall attention towards  
548 geometric stimuli between ASD and TD participants. This significant difference in visual attention was also  
549 found between ASD children with severe symptoms and TD children in our study. Despite these differences,  
550 using the ratio of visual attention towards the geometric stimuli and the total overall attention and implementing  
551 a thresholding technique employed previously<sup>11-13</sup> resulted in lower classification performance than our  
552 proposed model. Using an ML-based approach on handcrafted features<sup>23,24</sup> also resulted in lower accuracy in

553 ASD prediction and a similar accuracy in symptom severity prediction than our proposed model. Overall, our  
554 results demonstrate the feasibility of using our approach in accurately identifying ASD children and children  
555 with severe symptoms. Our model achieved promising performance with high accuracy, sensitivity and  
556 specificity.

557 Finally, most published research reviewed in this paper attempted to identify adults with ASD or older  
558 ASD children. In contrast, we investigated the possibility of diagnosing autism and predicting the level of ASD-  
559 related symptoms in preschool children (around 4 years old), an age range where diagnosis and assessment are  
560 typically performed. As a result, we provide an alternative to augment (and not replace) existing clinical  
561 observation tools with a more objective and efficient approach to ASD diagnosis. This takes us closer to an  
562 early ASD screening system and allows children to access intervention for better health outcomes. While our  
563 results are promising, our proposed approach needs to be trained and tested on a much larger dataset before it  
564 can be utilised in clinical settings.

565 From a clinical perspective, our findings suggest that eye-tracking technology could be used as a  
566 biomarker of the presence of ASD and symptom severity in preschool children. Initial findings already found  
567 significant correlations between changes in eye-tracking measures and changes in clinical measures captured  
568 before and after interventions, suggesting that eye-tracking can be utilised to quantify treatment response<sup>91</sup>.  
569 Given the rapid advances in technology supported by the promising performance of the classification models  
570 reviewed in this paper, it is not hard to imagine that future research would explore the use of a similar eye-  
571 tracking paradigm in predicting other clinical phenotypes and treatment response outcomes in preschool ASD  
572 children. This will have a tremendous impact on targeting interventions that maximise health outcomes in  
573 patients.

## 574 **Limitations**

575 Despite the utility of the current study, there are several limitations to keep in mind. First, there was a  
576 gender skew towards males in the ASD group, as would be clinically expected. Nevertheless, further studies  
577 with more female participants are required to clarify our results, as differences in autism presentation and  
578 diagnosis between males and females have been documented.<sup>92</sup> For example, studies have shown that girls on  
579 the spectrum behave similarly to neurotypical boys and girls on certain socially orientated tasks, such as  
580 enhanced attention to faces during scenes that do not have social interactions.<sup>93,94</sup> In addition, TD men with high  
581 ASD traits exhibit worse accuracy of gaze shifts, while TD women have similar gaze-following behaviour  
582 regardless of ASD traits.<sup>95</sup>

583 Further, the participant groups also differed in sample size, with the ASD group being three times as  
584 large as the TD group. The ASD participants in this study were recruited from an ASD-specific centre and there  
585 was good uptake to the study. Despite significant efforts of the team to recruit control participants, there was  
586 less interest from the families of neurotypical children to participate in the study, which is probably not  
587 surprising given the study is less meaningful for children without a developmental diagnosis. We also  
588 acknowledge that the dataset size is relatively small in comparison to the dataset required to train modern DL  
589 models. To aid our model training and leverage transfer learning, we utilised one of the best dynamic saliency  
590 detection model<sup>88</sup> and finetuned its weights to our dataset. This allowed our model to learn better and extract  
591 more robust and semantically meaningful features when compared to a model trained from scratch on our  
592 dataset. We believe that using the leave-one-out cross-validation approach to train and test the model addressed  
593 the class imbalance and small sample size in our study. This validation approach has been used extensively in  
594 prior research<sup>14,33,34,43,68,69</sup>.

595 It is also useful to note that the participant groups were matched on chronological age but not on  
596 developmental abilities. Further studies with larger sample sizes with a developmentally age-matched group are  
597 suggested to confirm our findings. As reported in the Materials and methods section, children with ASD were  
598 not excluded from the study if they had a comorbid diagnosis. Although this has implications for any strict  
599 interpretation of the findings reported here, the inclusion of comorbid conditions in ASD research is  
600 ecologically valid. Indeed, it is rare in clinical practice to encounter a young person who has a ‘pure’ autism  
601 spectrum diagnosis with no other psychiatric or developmental comorbidities.

602 Finally, we cannot report on the performance of the stimuli-based classification approaches and  
603 compare it with our dynamic stimuli-based classification approach since this study is part of a larger study that  
604 aimed to find differences in eye-tracking data between ASD and TD participants while watching dynamic  
605 stimuli. As such, no eye-tracking data from the same participants were collected while viewing static stimuli.

## 606 **Author contributions**

607 RAJDB conceptualised the methodology, conducted the pre-processing of the eye-tracking data,  
608 performed the statistical analysis, developed the deep neural network, and wrote the initial draft of the  
609 manuscript under the supervision of VE, TB and AS. All authors reviewed the manuscript and contributed to the  
610 revision of the article. All authors approved the final version of the manuscript.

## 611 **Additional information**



612           The authors have declared that no competing interests exist.

## 613    **Acknowledgements**

614           We extend our gratitude to the children and their families who participated in this study and to the staff  
615 where this study was conducted.

616    **Availability of data and materials:** The datasets generated and/or analysed during the current study are not  
617 publicly available but are available from the corresponding author upon reasonable request.

618

## References

619

- 620 1 Huerta, M., Bishop, S. L., Duncan, A., Hus, V. & Lord, C. Application of DSM-5 criteria  
621 for autism spectrum disorder to three samples of children with DSM-IV diagnoses of  
622 pervasive developmental disorders. *American Journal of Psychiatry* **169**, 1056-1064  
623 (2012).
- 624 2 Randall, M. *et al.* Diagnostic tests for autism spectrum disorder (ASD) in preschool  
625 children. *Cochrane Database of Systematic Reviews* (2018).
- 626 3 Taylor, L. J. *et al.* Brief Report: An Exploratory Study of the Diagnostic Reliability for  
627 Autism Spectrum Disorder. *Journal of Autism and Developmental Disorders* **47**, 1551-  
628 1558, doi:[10.1007/s10803-017-3054-z](https://doi.org/10.1007/s10803-017-3054-z) (2017).
- 629 4 Estes, A. *et al.* Long-term outcomes of early intervention in 6-year-old children with  
630 autism spectrum disorder. *Journal of the American Academy of Child & Adolescent*  
631 *Psychiatry* **54**, 580-587 (2015).
- 632 5 de Belen, R. A. J., Bednarz, T., Sowmya, A. & Del Favero, D. Computer vision in autism  
633 spectrum disorder research: a systematic review of published studies from 2009 to  
634 2019. *Translational Psychiatry* **10**, 333, doi:[https://doi.org/10.1038/s41398-020-](https://doi.org/10.1038/s41398-020-01015-w)  
635 [01015-w](https://doi.org/10.1038/s41398-020-01015-w) (2020).
- 636 6 Sapiro, G., Hashemi, J. & Dawson, G. Computer vision and behavioral phenotyping:  
637 an autism case study. *Current Opinion in Biomedical Engineering* **9**, 14-20,  
638 doi:<https://doi.org/10.1016/j.cobme.2018.12.002> (2019).
- 639 7 Ahmed, Z. A. T. & Jadhav, M. E. A Review of Early Detection of Autism Based on Eye-  
640 Tracking and Sensing Technology in 2020 *International Conference on Inventive*  
641 *Computation Technologies (ICICT)*. 160-166 (IEEE) (Year).
- 642 8 Kollias, K.-F., Syriopoulou-Delli, C. K., Sarigiannidis, P. & Fragulis, G. F. The  
643 Contribution of Machine Learning and Eye-Tracking Technology in Autism Spectrum  
644 Disorder Research: A Systematic Review. *Electronics* **10**, 2982 (2021).
- 645 9 de Belen, R. A. *et al.* Eye-tracking correlates of response to joint attention in  
646 preschool children with autism spectrum disorder. *BMC Psychiatry* **23**, 211,  
647 doi:[10.1186/s12888-023-04585-3](https://doi.org/10.1186/s12888-023-04585-3) (2023).
- 648 10 Osterling, J. & Dawson, G. Early recognition of children with autism: a study of first  
649 birthday home videotapes. *J Autism Dev Disord* **24**, 247-257,  
650 doi:[10.1007/bf02172225](https://doi.org/10.1007/bf02172225) (1994).
- 651 11 Pierce, K., Conant, D., Hazin, R., Stoner, R. & Desmond, J. Preference for Geometric  
652 Patterns Early in Life as a Risk Factor for Autism. *Archives of General Psychiatry* **68**,  
653 101-109, doi:[10.1001/archgenpsychiatry.2010.113](https://doi.org/10.1001/archgenpsychiatry.2010.113) (2011).
- 654 12 Pierce, K. *et al.* Eye Tracking Reveals Abnormal Visual Preference for Geometric  
655 Images as an Early Biomarker of an Autism Spectrum Disorder Subtype Associated  
656 With Increased Symptom Severity. *Biological Psychiatry* **79**, 657-666,  
657 doi:<https://doi.org/10.1016/j.biopsych.2015.03.032> (2016).
- 658 13 Moore, A. *et al.* The geometric preference subtype in ASD: identifying a consistent,  
659 early-emerging phenomenon through eye tracking. *Molecular autism* **9**, 1-13 (2018).
- 660 14 de Belen, R. A. J., Bednarz, T. & Sowmya, A. EyeXplain Autism: Interactive System for  
661 Eye Tracking Data Analysis and Deep Neural Network Interpretation for Autism  
662 Spectrum Disorder Diagnosis in *Extended Abstracts of the 2021 CHI Conference on*  
663 *Human Factors in Computing Systems*. Article 364 (Association for Computing  
664 Machinery), doi:<https://doi.org/10.1145/3411763.3451784> (Year).

- 665 15 Oliveira, J. S. *et al.* Computer-aided autism diagnosis based on visual attention  
666 models using eye tracking. *Scientific reports* **11**, 1-11 (2021).
- 667 16 Revers, M. C. *et al.* Classification of Autism Spectrum Disorder Severity Using Eye  
668 Tracking Data Based on Visual Attention Model in *2021 IEEE 34th International  
669 Symposium on Computer-Based Medical Systems (CBMS)*. 142-147 (IEEE) (Year).
- 670 17 Itti, L., Koch, C. & Niebur, E. A model of saliency-based visual attention for rapid  
671 scene analysis. *IEEE Transactions on pattern analysis and machine intelligence* **20**,  
672 1254-1259 (1998).
- 673 18 Borji, A. Saliency prediction in the deep learning era: Successes and limitations. *IEEE  
674 transactions on pattern analysis and machine intelligence* **43**, 679-700 (2019).
- 675 19 de Belen, R. A. J., Bednarz, T. & Sowmya, A. ScanpathNet: A Recurrent Mixture  
676 Density Network for Scanpath Prediction in *Proceedings of the IEEE/CVF Conference  
677 on Computer Vision and Pattern Recognition*. 5010-5020 (Year).
- 678 20 Duan, H. *et al.* A dataset of eye movements for the children with autism spectrum  
679 disorder in *Proceedings of the 10th ACM Multimedia Systems Conference*. 255–260  
680 (Association for Computing Machinery), doi:10.1145/3304109.3325818 (Year).
- 681 21 Gutiérrez, J., Che, Z., Zhai, G. & Le Callet, P. Saliency4ASD: Challenge, dataset and  
682 tools for visual attention modeling for autism spectrum disorder. *Signal Processing:  
683 Image Communication* **92**, 116092 (2021).
- 684 22 Le Meur, O., Nebout, A., Chereil, M. & Etchamendy, E. From Kanner Autism to  
685 Asperger Syndromes, the Difficult Task to Predict Where ASD People Look at. *IEEE  
686 Access* **8**, 162132-162140 (2020).
- 687 23 Canavan, S. *et al.* Combining gaze and demographic feature descriptors for autism  
688 classification in *2017 IEEE International Conference on Image Processing (ICIP)*.  
689 3750-3754 (IEEE) (Year).
- 690 24 Fabiano, D., Canavan, S., Agazzi, H., Hinduja, S. & Goldgof, D. Gaze-based  
691 classification of autism spectrum disorder. *Pattern Recognition Letters* **135**, 204-212  
692 (2020).
- 693 25 Chaddad, A., Desrosiers, C. & Toews, M. Multi-scale radiomic analysis of sub-cortical  
694 regions in MRI related to autism, gender and age. *Scientific Reports* **7**, 45639,  
695 doi:10.1038/srep45639 (2017).
- 696 26 Chaddad, A., Desrosiers, C., Hassan, L. & Tanougast, C. Hippocampus and amygdala  
697 radiomic biomarkers for the study of autism spectrum disorder. *BMC Neuroscience*  
698 **18**, 52, doi:10.1186/s12868-017-0373-0 (2017).
- 699 27 Chanel, G. *et al.* Classification of autistic individuals and controls using cross-task  
700 characterization of fMRI activity. *NeuroImage: Clinical* **10**, 78-88,  
701 doi:<https://doi.org/10.1016/j.nicl.2015.11.010> (2016).
- 702 28 Eslami, T. & Saeed, F. Auto-ASD-Network: A Technique Based on Deep Learning and  
703 Support Vector Machines for Diagnosing Autism Spectrum Disorder using fMRI Data  
704 in *Proceedings of the 10th ACM International Conference on Bioinformatics,  
705 Computational Biology and Health Informatics*. 646–651 (Association for Computing  
706 Machinery), doi:10.1145/3307339.3343482 (Year).
- 707 29 Zheng, W. *et al.* Multi-feature based network revealing the structural abnormalities  
708 in autism spectrum disorder. *IEEE Transactions on Affective Computing*, 1-1,  
709 doi:10.1109/TAFFC.2018.2890597 (2018).

- 710 30 Crimi, A., Dodero, L., Murino, V. & Sona, D. Case-control discrimination through  
711 effective brain connectivity in *2017 IEEE 14th International Symposium on*  
712 *Biomedical Imaging (ISBI 2017)*. 970-973, doi:10.1109/ISBI.2017.7950677 (Year).
- 713 31 Shukla, P., Gupta, T., Saini, A., Singh, P. & Balasubramanian, R. A Deep Learning  
714 Frame-Work for Recognizing Developmental Disorders in *2017 IEEE Winter*  
715 *Conference on Applications of Computer Vision (WACV)*. 705-714,  
716 doi:10.1109/WACV.2017.84 (Year).
- 717 32 Li, B. *et al.* A Facial Affect Analysis System for Autism Spectrum Disorder in *2019 IEEE*  
718 *International Conference on Image Processing (ICIP)*. 4549-4553,  
719 doi:10.1109/ICIP.2019.8803604 (Year).
- 720 33 Jiang, M. & Zhao, Q. Learning Visual Attention to Identify People with Autism  
721 Spectrum Disorder in *2017 IEEE International Conference on Computer Vision (ICCV)*.  
722 3287-3296, doi:10.1109/ICCV.2017.354 (Year).
- 723 34 Liu, W., Li, M. & Yi, L. Identifying children with autism spectrum disorder based on  
724 their face processing abnormality: A machine learning framework. *Autism Research*  
725 **9**, 888-898, doi:10.1002/aur.1615 (2016).
- 726 35 Liu, W. *et al.* Efficient autism spectrum disorder prediction with eye movement: A  
727 machine learning framework in *2015 International Conference on Affective*  
728 *Computing and Intelligent Interaction (ACII)*. 649-655,  
729 doi:10.1109/ACII.2015.7344638 (Year).
- 730 36 Vu, T. *et al.* Effective and efficient visual stimuli design for quantitative autism  
731 screening: An exploratory study in *2017 IEEE EMBS International Conference on*  
732 *Biomedical & Health Informatics (BHI)*. 297-300, doi:10.1109/BHI.2017.7897264  
733 (Year).
- 734 37 Vyas, K. *et al.* Recognition Of Atypical Behavior In Autism Diagnosis From Video Using  
735 Pose Estimation Over Time in *2019 IEEE 29th International Workshop on Machine*  
736 *Learning for Signal Processing (MLSP)*. 1-6, doi:10.1109/MLSP.2019.8918863 (Year).
- 737 38 Zunino, A. *et al.* Video Gesture Analysis for Autism Spectrum Disorder Detection in  
738 *2018 24th International Conference on Pattern Recognition (ICPR)*. 3421-3426,  
739 doi:10.1109/ICPR.2018.8545095 (Year).
- 740 39 Rajagopalan, S. S., Dhall, A. & Goecke, R. Self-Stimulatory Behaviours in the Wild for  
741 Autism Diagnosis in *2013 IEEE International Conference on Computer Vision*  
742 *Workshops*. 755-761, doi:10.1109/ICCVW.2013.103 (Year).
- 743 40 Rajagopalan, S. S. & Goecke, R. Detecting self-stimulatory behaviours for autism  
744 diagnosis in *2014 IEEE International Conference on Image Processing (ICIP)*. 1470-  
745 1474, doi:10.1109/ICIP.2014.7025294 (Year).
- 746 41 Wang, Z. *et al.* Screening Early Children with Autism Spectrum Disorder via  
747 Response-to-Name Protocol. *IEEE Transactions on Industrial Informatics*, 1-1,  
748 doi:10.1109/TII.2019.2958106 (2019).
- 749 42 Wang, Z., Xu, K. & Liu, H. Screening Early Children with Autism Spectrum Disorder via  
750 Expressing Needs with Index Finger Pointing in *Proceedings of the 13th International*  
751 *Conference on Distributed Smart Cameras*. Article 24 (Association for Computing  
752 Machinery), doi:10.1145/3349801.3349826 (Year).
- 753 43 Chen, S. & Zhao, Q. Attention-Based Autism Spectrum Disorder Screening With  
754 Privileged Modality in *2019 IEEE/CVF International Conference on Computer Vision*  
755 *(ICCV)*. 1181-1190, doi:10.1109/ICCV.2019.00127 (Year).

- 756 44 Minissi, M. E., Chicchi Giglioli, I. A., Mantovani, F. & Alcaniz Raya, M. Assessment of  
757 the autism spectrum disorder based on machine learning and social visual attention:  
758 A systematic review. *Journal of Autism and Developmental Disorders* **52**, 2187-2202  
759 (2022).
- 760 45 Carette, R., Elbattah, M., Dequen, G., Guérin, J.-L. & Cilia, F. Visualization of eye-  
761 tracking patterns in autism spectrum disorder: method and dataset in *2018*  
762 *Thirteenth International Conference on Digital Information Management (ICDIM)*.  
763 248-253 (IEEE) (Year).
- 764 46 Duan, H. *et al.* Learning to predict where the children with asd look in *2018 25th IEEE*  
765 *international conference on image processing (ICIP)*. 704-708 (IEEE) (Year).
- 766 47 Duan, H. *et al.* Visual attention analysis and prediction on human faces for children  
767 with autism spectrum disorder. *ACM Transactions on Multimedia Computing,*  
768 *Communications, and Applications (TOMM)* **15**, 1-23 (2019).
- 769 48 Fang, Y., Huang, H., Wan, B. & Zuo, Y. Visual attention modeling for autism spectrum  
770 disorder by semantic features in *2019 IEEE International Conference on Multimedia*  
771 *& Expo Workshops (ICMEW)*. 625-628 (IEEE) (Year).
- 772 49 Wei, W., Liu, Z., Huang, L., Nebout, A. & Le Meur, O. Saliency prediction via multi-  
773 level features and deep supervision for children with autism spectrum disorder in  
774 *2019 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*.  
775 621-624 (IEEE) (Year).
- 776 50 Nebout, A., Wei, W., Liu, Z., Huang, L. & Le Meur, O. Predicting saliency maps for asd  
777 people in *2019 IEEE International Conference on Multimedia & Expo Workshops*  
778 *(ICMEW)*. 629-632 (IEEE) (Year).
- 779 51 Fang, Y. *et al.* Visual attention prediction for Autism Spectrum Disorder with  
780 hierarchical semantic fusion. *Signal Processing: Image Communication* **93**, 116186  
781 (2021).
- 782 52 Wei, W. *et al.* Predicting atypical visual saliency for autism spectrum disorder via  
783 scale-adaptive inception module and discriminative region enhancement loss.  
784 *Neurocomputing* **453**, 610-622 (2021).
- 785 53 Min, X. *et al.* Visual attention analysis and prediction on human faces. *Information*  
786 *Sciences* **420**, 417-430 (2017).
- 787 54 Wang, S. *et al.* Atypical Visual Saliency in Autism Spectrum Disorder Quantified  
788 through Model-Based Eye Tracking. *Neuron* **88**, 604-616,  
789 doi:10.1016/j.neuron.2015.09.042 (2015).
- 790 55 Yaneva, V., Eraslan, S., Yesilada, Y. & Mitkov, R. Detecting high-functioning autism in  
791 adults using eye tracking and machine learning. *IEEE Transactions on Neural Systems*  
792 *and Rehabilitation Engineering* **28**, 1254-1261 (2020).
- 793 56 Startsev, M. & Dorr, M. Classifying autism spectrum disorder based on scanpaths and  
794 saliency in *2019 IEEE International Conference on Multimedia & Expo Workshops*  
795 *(ICMEW)*. 633-636 (IEEE) (Year).
- 796 57 Arru, G., Mazumdar, P. & Battisti, F. Exploiting visual behaviour for autism spectrum  
797 disorder identification in *2019 IEEE International Conference on Multimedia & Expo*  
798 *Workshops (ICMEW)*. 637-640 (IEEE) (Year).
- 799 58 Wu, C., Liaqat, S., Cheung, S.-c., Chuah, C.-N. & Ozonoff, S. Predicting autism  
800 diagnosis using image with fixations and synthetic saccade patterns in *2019 IEEE*  
801 *International Conference on Multimedia & Expo Workshops (ICMEW)*. 647-650 (IEEE)  
802 (Year).

- 803 59 Tao, Y. & Shyu, M.-L. SP-ASDNet: CNN-LSTM based ASD classification model using  
804 observer scanpaths in *2019 IEEE International conference on multimedia & expo*  
805 *workshops (ICMEW)*. 641-646 (IEEE) (Year).
- 806 60 Fang, Y., Duan, H., Shi, F., Min, X. & Zhai, G. Identifying children with autism  
807 spectrum disorder based on gaze-following in *2020 IEEE International Conference on*  
808 *Image Processing (ICIP)*. 423-427 (IEEE) (Year).
- 809 61 Rahman, S., Rahman, S., Shahid, O., Abdullah, M. T. & Sourov, J. A. Classifying eye-  
810 tracking data using saliency maps in *2020 25th International Conference on Pattern*  
811 *Recognition (ICPR)*. 9288-9295 (IEEE) (Year).
- 812 62 Xu, S., Yan, J. & Hu, M. A new bio-inspired metric based on eye movement data for  
813 classifying ASD and typically developing children. *Signal Processing: Image*  
814 *Communication* **94**, 116171 (2021).
- 815 63 Wei, W. *et al.* Identify autism spectrum disorder via dynamic filter and deep  
816 spatiotemporal feature extraction. *Signal Processing: Image Communication* **94**,  
817 116195 (2021).
- 818 64 Liaqat, S. *et al.* Predicting ASD diagnosis in children with synthetic and image-based  
819 eye gaze data. *Signal Processing: Image Communication* **94**, 116198 (2021).
- 820 65 Mazumdar, P., Arru, G. & Battisti, F. Early detection of children with autism spectrum  
821 disorder based on visual exploration of images. *Signal Processing: Image*  
822 *Communication* **94**, 116184 (2021).
- 823 66 Tseng, P.-H. *et al.* High-throughput classification of clinical populations from natural  
824 viewing eye movements. *Journal of neurology* **260**, 275-284 (2013).
- 825 67 Wan, G. *et al.* Applying eye tracking to identify autism spectrum disorder in children.  
826 *Journal of autism and developmental disorders* **49**, 209-215 (2019).
- 827 68 Jiang, M. *et al.* Classifying individuals with ASD through facial emotion recognition  
828 and eye-tracking in *2019 41st Annual International Conference of the IEEE*  
829 *Engineering in Medicine and Biology Society (EMBC)*. 6063-6068 (IEEE) (Year).
- 830 69 Zhao, Z. *et al.* Classification of Children With Autism and Typical Development Using  
831 Eye-Tracking Data From Face-to-Face Conversations: Machine Learning Model  
832 Development and Performance Evaluation. *J Med Internet Res* **23**, e29328,  
833 doi:10.2196/29328 (2021).
- 834 70 Carette, R. *et al.* Learning to Predict Autism Spectrum Disorder based on the Visual  
835 Patterns of Eye-tracking Scanpaths in *HEALTHINF*. 103-112 (Year).
- 836 71 Elbattah, M., Carette, R., Dequen, G., Guérin, J.-L. & Cilia, F. Learning clusters in  
837 autism spectrum disorder: Image-based clustering of eye-tracking scanpaths with  
838 deep autoencoder in *2019 41st Annual international conference of the IEEE*  
839 *engineering in medicine and biology society (EMBC)*. 1417-1420 (IEEE) (Year).
- 840 72 Akter, T., Ali, M. H., Khan, M. I., Satu, M. S. & Moni, M. A. Machine learning model to  
841 predict autism investigating eye-tracking dataset in *2021 2nd International*  
842 *Conference on Robotics, Electrical and Signal Processing Techniques (ICREST)*. 383-  
843 387 (IEEE) (Year).
- 844 73 Cilia, F. *et al.* Computer-aided screening of autism spectrum disorder: eye-tracking  
845 study using data visualization and deep learning. *JMIR Human Factors* **8**, e27706  
846 (2021).
- 847 74 Kanhirakadavath, M. R. & Chandran, M. S. M. Investigation of Eye-Tracking Scan Path  
848 as a Biomarker for Autism Screening Using Machine Learning Algorithms. *Diagnostics*  
849 **12**, 518 (2022).

- 850 75 Gaspar, A., Oliva, D., Hinojosa, S., Aranguren, I. & Zaldivar, D. An optimized Kernel  
851 Extreme Learning Machine for the classification of the autism spectrum disorder by  
852 using gaze tracking images. *Applied Soft Computing* **120**, 108654 (2022).
- 853 76 Ahmed, I. A. *et al.* Eye Tracking-Based Diagnosis and Early Detection of Autism  
854 Spectrum Disorder Using Machine Learning and Deep Learning Techniques.  
855 *Electronics* **11**, 530 (2022).
- 856 77 Fan, L. *et al.* Screening of Autism Spectrum Disorder Using Novel Biological Motion  
857 Stimuli. 371-384 (Springer Singapore) (Year).
- 858 78 Fang, H., Fan, L. & Hwang, J.-N. Auxiliary Diagnostic Method for Early Autism  
859 Spectrum Disorder Based on Eye Movement Data Analysis in *2021 IEEE 7th*  
860 *International Conference on Cloud Computing and Intelligent Systems (CCIS)*. 72-77  
861 (IEEE) (Year).
- 862 79 Carette, R. *et al.* Automatic autism spectrum disorder detection thanks to eye-  
863 tracking and neural network-based approach in *Internet of Things (IoT) Technologies*  
864 *for HealthCare: 4th International Conference, HealthyIoT 2017, Angers, France,*  
865 *October 24-25, 2017, Proceedings 4*. 75-81 (Springer) (Year).
- 866 80 Putra, P. U., Shima, K., Alvarez, S. A. & Shimatani, K. Identifying autism spectrum  
867 disorder symptoms using response and gaze behavior during the Go/NoGo game  
868 CatChicken. *Scientific reports* **11**, 1-12 (2021).
- 869 81 Kou, J. *et al.* Comparison of three different eye-tracking tasks for distinguishing  
870 autistic from typically developing children and autistic symptom severity. *Autism*  
871 *Research* **12**, 1529-1540, doi:<https://doi.org/10.1002/aur.2174> (2019).
- 872 82 Bacon, E. C. *et al.* Identifying prognostic markers in autism spectrum disorder using  
873 eye tracking. *Autism* **24**, 658-669 (2020).
- 874 83 Treisman, A. M. & Gelade, G. A feature-integration theory of attention. *Cognitive*  
875 *psychology* **12**, 97-136 (1980).
- 876 84 Kononenko, I., Šimec, E. & Robnik-Šikonja, M. Overcoming the myopia of inductive  
877 learning algorithms with RELIEFF. *Applied Intelligence* **7**, 39-55 (1997).
- 878 85 Association, A. P. *Diagnostic and statistical manual of mental disorders (DSM-5®)*.  
879 (American Psychiatric Pub, 2013).
- 880 86 Lord, C. *et al.* Autism diagnostic observation schedule, (ADOS-2) modules 1-4. *Los*  
881 *Angeles, California: Western Psychological Services* (2012).
- 882 87 Olsen, A. The Tobii I-VT fixation filter. *Tobii Technology* **21** (2012).
- 883 88 Wang, W., Shen, J., Guo, F., Cheng, M.-M. & Borji, A. Revisiting video saliency: A  
884 large-scale benchmark and a new model in *Proceedings of the IEEE Conference on*  
885 *Computer Vision and Pattern Recognition*. 4894-4903 (Year).
- 886 89 Huang, X., Shen, C., Boix, X. & Zhao, Q. Salicon: Reducing the semantic gap in saliency  
887 prediction by adapting deep neural networks in *Proceedings of the IEEE international*  
888 *conference on computer vision*. 262-270 (Year).
- 889 90 Vapnik, V. N. An overview of statistical learning theory. *IEEE transactions on neural*  
890 *networks* **10**, 988-999 (1999).
- 891 91 Bradshaw, J. *et al.* The Use of Eye Tracking as a Biomarker of Treatment Outcome in  
892 a Pilot Randomized Clinical Trial for Young Children with Autism. *Autism Research* **12**,  
893 779-793, doi:<https://doi.org/10.1002/aur.2093> (2019).
- 894 92 Lai, M.-C. & Szatmari, P. Sex and gender impacts on the behavioural presentation  
895 and recognition of autism. *Current Opinion in Psychiatry* **33**, 117-123,  
896 doi:10.1097/yco.0000000000000575 (2020).

- 897 93 Harrop, C. *et al.* Visual attention to faces in children with autism spectrum disorder:  
898 are there sex differences? *Molecular Autism* **10**, 28, doi:10.1186/s13229-019-0276-2  
899 (2019).
- 900 94 Harrop, C. *et al.* Social and Object Attention Is Influenced by Biological Sex and Toy  
901 Gender-Congruence in Children With and Without Autism. *Autism Research* **13**, 763-  
902 776, doi:<https://doi.org/10.1002/aur.2245> (2020).
- 903 95 Whyte, E. M. & Scherf, K. S. Gaze Following Is Related to the Broader Autism  
904 Phenotype in a Sex-Specific Way: Building the Case for Distinct Male and Female  
905 Autism Phenotypes. *Clinical Psychological Science* **6**, 280-287,  
906 doi:10.1177/2167702617738380 (2018).  
907