

1 **Title**

2

3 A computational probe into the behavioral and neural markers of atypical facial emotion
4 processing in autism.

5

6 **Authors**

7

8 Kohitij Kar^{1,2*}

9

10 1. McGovern Institute for Brain Research and Department of Brain and Cognitive
11 Sciences, Massachusetts Institute of Technology, Cambridge, MA, 01239, USA

12 2. Center for Brains, Minds, and Machines, Massachusetts Institute of Technology,
13 Cambridge, MA, 01239, USA

14

15 *Correspondence should be addressed to Kohitij Kar.

16

17 **Contact Info**

18

19 McGovern Institute for Brain Research
20 Massachusetts Institute of Technology,
21 77 Massachusetts Institute of Technology, 46-6161,
22 Cambridge, MA 02139

23 E-mail: kohitij@mit.edu

24

25

26 **Abstract**

27
28

29 Despite ample behavioral evidence of atypical facial emotion processing in individuals
30 with autism (IwA), the neural underpinnings of such behavioral heterogeneities remain
31 unclear. Here, I have used brain-tissue mapped artificial neural network (ANN) models of
32 primate vision to probe candidate neural and behavior markers of atypical facial emotion
33 recognition in IwA at an image-by-image level. Interestingly, the ANNs' image-level
34 behavioral patterns better matched the neurotypical subjects' behavior than those
35 measured in IwA. This behavioral mismatch was most remarkable when the ANN
36 behavior was decoded from units that correspond to the primate inferior temporal (IT)
37 cortex. ANN-IT responses also explained a significant fraction of the image-level
38 behavioral predictivity associated with neural activity in the human amygdala — strongly
39 suggesting that the previously reported facial emotion intensity encodes in the human
40 amygdala could be primarily driven by projections from the IT cortex. Furthermore, in
41 silico experiments revealed how learning under noisy sensory representations could lead
42 to atypical facial emotion processing that better matches the image-level behavior
43 observed in IwA. In sum, these results identify primate IT activity as a candidate neural
44 marker and demonstrate how ANN models of vision can be used to generate neural
45 circuit-level hypotheses and guide future human and non-human primate studies in
46 autism.

47
48
49
50
51
52
53
54

55 **Keywords**

56

57 Autism, Amygdala, Inferior Temporal Cortex, Artificial Neural Networks, Facial emotion
58 recognition

59

60 Introduction

61
62 The ability to recognize others' mood, emotion, and intent from facial expressions lie at
63 the core of human interpersonal communication and social engagement. This relatively
64 automatic, visuocognitive feature that neurotypically developed human adults take for
65 granted shows significant differences in children and adults with autism¹⁻⁴. A mechanistic
66 understanding of the underlying neural correlates of such behavioral mismatches is key
67 to designing efficient cognitive therapies and other approaches to help individuals with
68 autism.

69
70 There is a growing body of work on how facial identity is encoded in the primate brain,
71 especially in the Fusiform Face Areas (FFA) in humans^{5,6} and in the topographically
72 specific "face patch" systems of the inferior temporal (IT) cortex of the rhesus macaques
73⁷⁻⁹. Also, previous research has linked human amygdala neural responses with
74 recognizing facial emotions¹⁰⁻¹². For instance, subjects who lack a functional amygdala
75 often exhibit selective impairments in recognizing fearful faces^{13,14}. Wang et al.¹⁵ also
76 demonstrated that the human amygdala parametrically encodes the intensity of specific
77 facial emotions (e.g., fear, happiness) and their categorical ambiguity. A critical question,
78 however, is whether the atypical facial emotion recognition broadly reported in individuals
79 with autism (IwA) arises purely from differences in sensory representations (i.e., purely
80 perceptual alterations^{16,17}) or is due to a primary (but not mutually exclusive) variation in
81 the development and function of specialized affect processing regions (e.g., atypical
82 amygdala development leading to specific differences in encoding emotion). There are
83 two main roadblocks toward answering this question. First, heterogeneity and
84 idiosyncrasies are commonplace across behavioral reports in autism, including facial
85 affect processing (for a formal meta-analysis of recognition of emotions in autism see:
86^{18,19}). The inability to parsimoniously explain such heterogeneous findings prevent us from
87 designing more efficient follow-up experiments to probe such questions further. Second,
88 in the absence of neurally mechanistic models of behavior, it remains challenging to infer
89 neural mechanisms from behavioral results and generate testable neural circuit level
90 predictions that can be validated or falsified using neurophysiological approaches.
91 Therefore, we need brain-mapped computational models that can predict at an image-
92 by-image level how primates represent facial emotions across different parts of their brain
93 and how such representations are linked to their performance in facial emotion judgment
94 tasks (like the one used in⁴).

95
96 The differences in facial emotion judgments between neurotypical adults and individuals
97 with autism are often interpreted with inferential models (e.g., psychometric functions)
98 that base their predictions on high-level categorical descriptors of the stimuli (e.g., overall
99 facial expression levels of "happiness", "fear" and other primary emotions²⁰). Such
100 modeling efforts are likely to ignore an important source of variance produced by the
101 image-level sensory representations of each stimuli being tested. To interpret this source
102 of variance, it is necessary to develop models that are image computable. Recent
103 progress in computer vision and computational neuroscience has led to the development
104 of artificial neural network (ANN) models that can both perform human-like object
105 recognition^{21,22} as well as contain internal components that match human and macaque

106 visual systems^{23,24}. Such image-computable ANNs can generate testable neural
107 hypotheses^{25,26} and help design experiments that leverage on the image-level variance
108 to guide us beyond the standard parametric approaches.

109
110 In this study, I have used a family of brain-tissue mapped ANN models of primate vision
111 to generate testable hypotheses and identify candidate neural and behavior markers of
112 atypical facial emotion recognition in IwA. Specifically, I have compared the predictions
113 of ANN models with behavior measured in neurotypical adults and people with autism⁴,
114 and facial emotion decodes from neural activity measured in the human amygdala¹⁵.
115 Furthermore, I performed in silico perturbation experiments to simulate and test autism-
116 relevant hypotheses of underlying neural mechanisms. I observed that the ANNs could
117 accurately predict the human facial emotion judgments at an image-by-image level.
118 Interestingly, the models' image-level behavioral patterns better matched the neurotypical
119 human subjects' behavior than those measured in individuals with autism. This behavioral
120 mismatch was most remarkable when the model behavior was constructed from units that
121 correspond to the primate IT cortex. Interestingly, I also observed this behavioral
122 mismatch when comparing neural decodes from a distinct population of visually facilitated
123 neurons in the human amygdala with *Control* and IwA behavior. However, ANN-IT
124 activation patterns could fully account for the image-level behavioral predictivity of the
125 human amygdala population responses that has been previously implicated in autism-
126 related facial emotion processing differences^{12,15}. Furthermore, in silico experiments
127 revealed that learning the emotion discrimination task with noisier ANN-IT representations
128 (i.e., with higher response variability per unit) result in weaker synaptic connections
129 between the model-IT and the downstream decision unit that improve the model's match
130 to the image-level behavioral patterns measured in the IwA. In sum, these results argue
131 that noisier sensory representations in the primate inferior temporal cortex that drive a
132 distinct population of neurons in the human amygdala is a key candidate mechanism of
133 atypical facial emotion processing in individuals with autism — a testable neural
134 hypothesis for future human and nonhuman primate studies.

135

136 Results

137
138 As outlined above, I reasoned that the ability to predict the image-level differences in
139 facial emotion judgments between individuals with autism (IwA) and neurotypical adults
140 (*Controls*) allow us to 1) design more efficient experiments to study the atypical facial
141 processing observed in IwA, 2) efficiently probe the underlying neural correlates. In this
142 study, I first took a data-driven approach to discover such image-level differences in
143 behavior across *Controls* and IwA in a facial emotion discrimination task ⁴. I then used
144 brain-mapped computational models of primate vision to probe the underlying neural
145 mechanisms that could drive such differences.

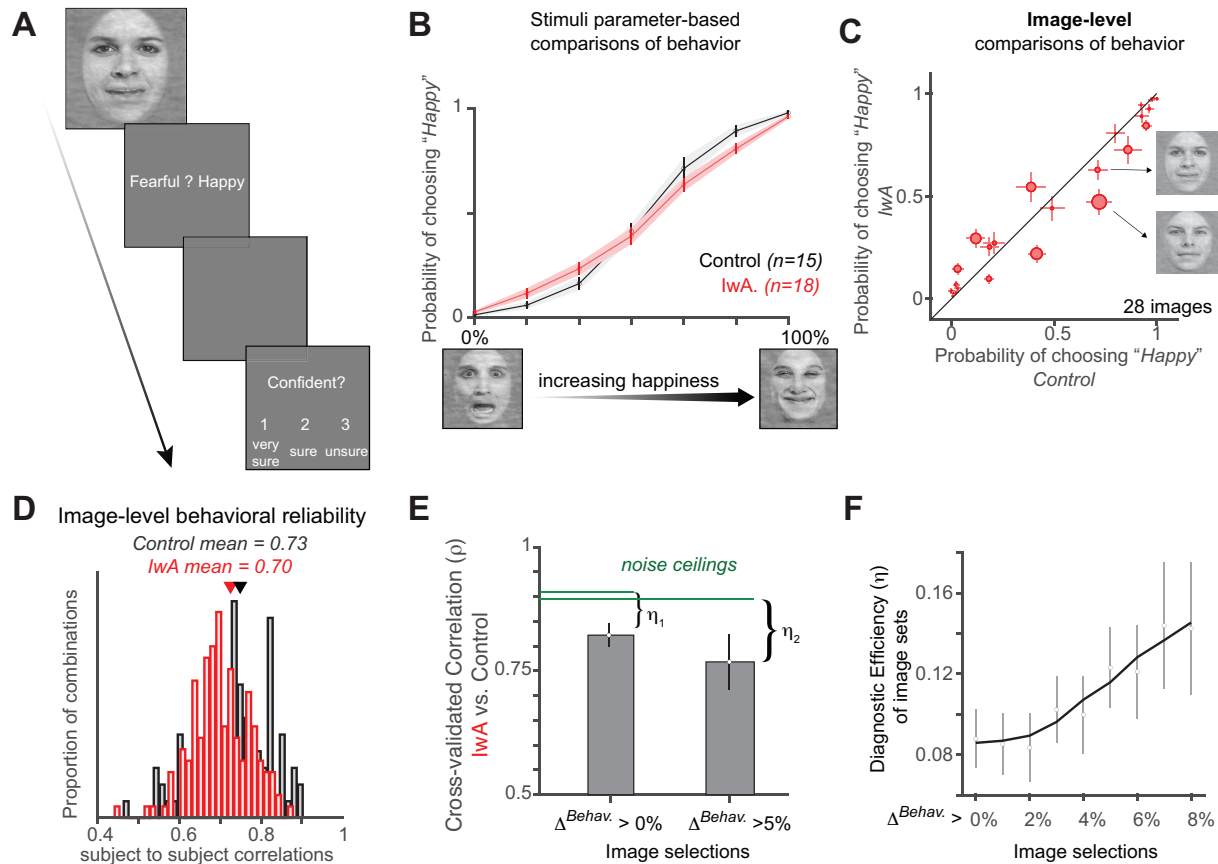
146
147 The behavioral and neural measurements analyzed in this study were performed by
148 Wang et al. ^{4,15}. During the task, participants were shown images of individual faces with
149 specific levels of morphed emotions (for 1 sec) and asked to discriminate between two
150 emotions, fear and happiness (Figure 1A; see Methods for details). The authors observed
151 a reduced specificity in facial emotion judgment among individuals with autism (IwA)
152 compared to neurotypical *Controls* (Figure 1B). Notably, the study controlled for low-level
153 image confounds, and eye movement patterns across the two groups did not explain the
154 reported behavioral differences. Therefore, the behavioral results significantly narrowed
155 the space of neural hypotheses to sensory and affect-processing circuits.

156 157 **Image-level differences can be leveraged to produce** 158 **stronger behavioral markers of atypical facial emotion** 159 **judgments in autism**

160
161 Wang and Adolphs ⁴ primarily investigated the differences in behavior of IwA and
162 *Controls*, across parametric variations of facial emotion levels (e.g., levels of happiness
163 and fear). Here, I first examined whether the image-by-image behavioral patterns
164 (irrespective of their facial identity or emotion levels), across the IwA and *Control* groups
165 could be reliably estimated. Therefore, I computed the individual subject-to-subject
166 correlations in image-level behavior (Figure 1D) which show that both of the groups
167 exhibit highly reliable image-level behavior. The internal reliability (see Methods) for
168 *Control* and IwA groups are 0.73 and 0.70, respectively. A visual inspection of the
169 comparison of behavioral patterns across the two groups (Figure 1C) show that there are
170 pairs of images (two such examples are shown in Figure 1C) for which the *Control* group
171 exhibited very similar behavior, but the IwA made very different behavioral responses.
172 This further confirms that diagnostic image-level variations in behavior could be further
173 utilized to gain more insight into the mechanisms that drive the atypical facial emotion
174 responses in IwA. Next, I quantified how stimuli selection based on high image-level
175 differences can be leveraged to design more efficient behavioral experiments. To do this,
176 I selected images based on the difference in behavior between the two groups (Δ^{Behav} :
177 using data from four randomly selected individual subjects from each group) and tested
178 the resulting correlation between the two groups' behavior (using the held-out subject
179 population). This was repeated several times to get a mean measure of the cross-

180 validated raw correlation (y-axis in Figure 1E). A noise-ceiling was measured for each
181 image-set selection based on image-level internal reliability of the held-out test population
182 (see Methods). The difference between the noise ceiling and the raw correlation is
183 referred to as the diagnostic efficiency η of the image-set, which is a measure of how
184 efficient the image-set is in discriminating between the *lwA* and *Control* behavior. Figure
185 1F shows how η varies across more and more efficient selection of image-sets (based on
186 higher differences in image-level behavior with *Controls* and *lwA*). These results suggest
187 that one reasonable goal of the field should be to find more efficient ways to predict which
188 images will produce the highest η values. Focusing human behavioral testing on such
189 images is likely going to yield stronger inferences and lead to a better understanding of
190 the behavioral and neural markers driving the difference in behavior.

191
192
193
194
195
196
197
198
199
200
201
202
203



204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229

Figure 1. Behavioral task and image-level assessment of behavioral markers. **A.** Subjects, both neurotypical (*Control*; $n=15$) population and individuals with autism (*IwA*; $n=18$) viewed a face for 1 sec in their central ~ 12 deg, followed by a question asking them to identify the facial emotion (fearful or happy). After a blank screen of 500 ms, subjects were then asked to indicate their confidence in their decision ('1' for 'very sure', '2' for 'sure' or '3' for 'unsure'). **B.** The psychometric curves show the proportion of trials judged as "happy" as a function of facial emotion morph levels (ranging from 0% happy (100% fearful; left) to 100% happy (0% fearful; right)). *IwA* (red curve), on average, showed lower specificity (slope of the psychometric curve) compared to the *Controls* (black curve). The shaded area and errorbars denotes SEM across participants. **C.** Image-level differences in behavior between *Controls* vs. *IwA*. Each red dot corresponds to an image. The size of the dot is scaled by the difference in behavior between the *Controls* and *IwA*. Errorbars denote SEM across subjects. Two example images are highlighted that show similar emotional ("happiness") judgments by the *Controls* but drive significantly different behaviors in *IwA* — demonstrating the importance of investigating individual image-level differences. **D.** The estimated image-by-image happiness judgments were highly reliable as demonstrated by comparisons across individuals (estimated separately for each group). The mean reliability (average of the individual subject to subject correlations) was 0.73 and 0.70 for the *Controls* (black histogram) and *IwA* (red histogram), respectively. **E.** Correlation between image-by-image behavioral patterns measured in *Controls* vs. *IwA*, with two different selections of images (cross-validated image selections with held-out subjects). Noise ceilings were calculated based on measured behavioral (split-half) reliability across populations within each group (see Methods). The difference between the noise ceiling and the mean raw correlation is referred to as the diagnostic efficiency of the image-set (η) **F.** Diagnostic efficiency (η) as a function of image selection criteria. Errorbars denote bootstrap confidence intervals. Facial images shown in this figure are morphed and processed version of the original face images. These images have full re-use permission.

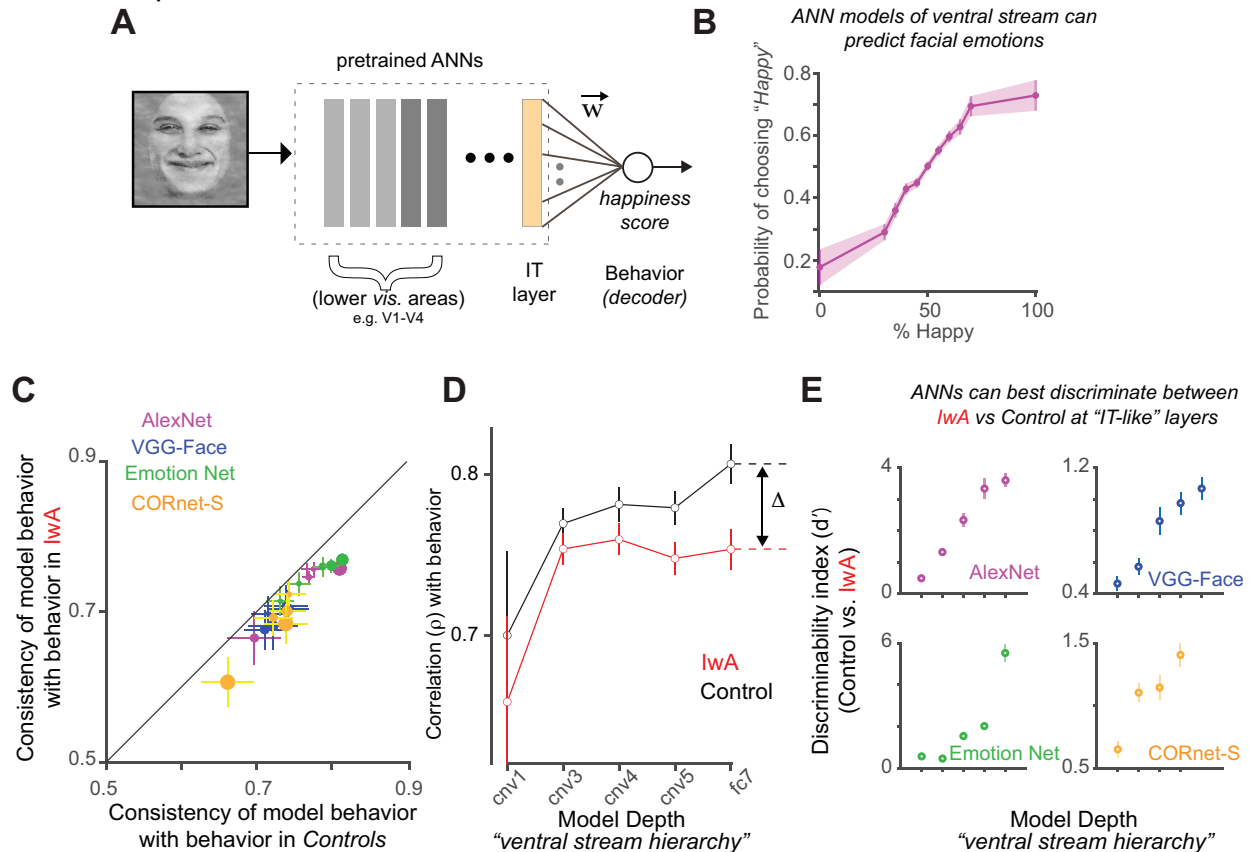
230 **ANN models of primate vision trained on varied objectives** 231 **can perform facial emotion judgment tasks**

232 To investigate how one can predict the image-level facial emotion judgments, I first tested
233 how accurately current ANN models of primate vision can be trained to perform such
234 tasks. One advantage of using these ANNs is that there are significant correspondences
235 between their architectural components and the areas in the primate ventral visual cortex
236 ^{24,25,27} (as shown in the schematic Figure 2A). Also, there is a significant match in the
237 predicted behavioral patterns of such models with primate behavior (including face-
238 related tasks) measured during multiple object recognition tasks^{21,22}. Taken together,
239 these models are great candidates for generating testable hypotheses regarding both
240 neural and behavioral markers of specific visual tasks. I selected four different ANNs to
241 test their behavioral predictions with respect to the facial emotional judgement task.
242 These ANNs were pretrained to perform image classification (AlexNet²⁸, CORnet-S²⁹),
243 face recognition (VGGFace³⁰) and emotion recognition (Emotion-Net³¹). I observed that,
244 a 10-fold cross validated partial least square regression model (see Methods for details)
245 could be used to train each model to perform the task. The variation of the behavioral
246 responses of the model with parametric changes in the level of happiness in the faces
247 qualitatively matched the patterns observed in the human data (Figure 2B).
248

249 **ANN model predictions better match the behavioral patterns** 250 **measured in neurotypical adults compared to individuals** 251 **with autism**

252 Next, I quantified how well the ANNs can predict the human image-level behavioral
253 responses (across both *Controls* and *lWA*). Interestingly, ANN models significantly
254 better predicted the image-level behavior measured in *Control* compared to the
255 behavior measured in *lWA* (Figure 2C; 20 models tested; paired t-test; $p < 0.00001$; $t(19) =$
256 10.99). To dissect which layer of the ANN best discriminated between the behavior of
257 *Controls* and *lWA*, I compared individual models constructed from different layers of the
258 same pretrained ANN architectures. This revealed two critical points. First, the
259 correlation between model behavior and the *Control* group behavior increased as a
260 function of model depth (black line; e.g. AlexNet shown in Figure 2D), which
261 corresponds to the ventral visual hierarchy as reported in many studies^{23,24}. Second, the
262 difference in the model's predictivity of behavior measured in *Controls* vs. *lWA* across
263 layers is also highest at deeper layers, which corresponds to primate IT (comparison of
264 the black and the red line for AlexNet shown in Figure 2D). This overall qualitative
265 observation was consistent across all four tested models (Figure 2E). Given the high
266 discriminability index (see Methods), established mappings between the layers and
267 primate brain, as well as wide usage among researchers, I have used AlexNet for the
268 subsequent analysis presented in this study. Therefore, these results suggest that
269 population neural activity in primate IT could play a significant role in the atypical facial
270 emotion processing in people with autism, and the image-level differences in sensory
271 representations in IT might explain the difference in behavior observed across the
272 images. However, such a role has been previously attributed to the human amygdala

273 responses¹⁵. Therefore, I next tested whether the human amygdala responses can
 274 predict the image-level behavior and how well this predictivity could be explained by the
 275 ANN-IT representations.



276

277 **Figure 2. Testing ANN-models on facial emotion recognition tasks.** **A.** ANN models of the primate
 278 ventral stream (typically comprising V1, V2, V4 and IT like layers) can be trained to predict human facial
 279 emotion judgments. This involves building a regression model, i.e., determining the weights \vec{w} based on
 280 the model layer activations (as the predictor) to predict the image ground truth ("level of happiness") on a
 281 set of training images, and then testing the predictions of this model on held-out images. **B.** An ANN model's
 282 predicted psychometric curves (e.g., AlexNet, shown here) show the proportion of trials judged as "happy"
 283 as a function of facial emotion morph levels ranging from 0% happy (100% fearful; left) to 100% happy (0%
 284 fearful; right). This curve demonstrates that activations of ANN layers (layer 'fc7' that corresponds to the
 285 "model- IT" layer) can be successfully trained to predict facial emotions. **C.** Comparison of ANN's image-
 286 level behavioral patterns with the behavior measured in *Controls* (x-axis) and IwA (y-axis). Four ANNs (with
 287 5 models each generated from different layers of the ANNs are shown here in different colors. ANN
 288 predictions better match the behavior measured in the *Controls* compared to IwA. The correlation values
 289 (x and y axes) were corrected by the noise estimates per human population so that the differences are not
 290 due to differences in noise-levels in measurements across the IwA and *Control* subject pools. The dot size
 291 refers to the degree of discrepancy between ANN predictivity of *Controls* vs. IwA. **D.** A comparison of the
 292 ANN predictivity (results from AlexNet shown here) of behavior measured in IwA vs. *Controls* as function
 293 of model layers (convolutional (cnv) layers 1,3,4, and 5 and the fully connected layer 7, 'fc7' -- that
 294 approximately corresponds to the ventral stream cortical hierarchy). The difference between the ANN's
 295 predictivity of behavior in IwA and *Controls* increases with depth and is referred to as Δ . **E.** Discriminability
 296 index (d' ; ability to discriminate between image-level behavioral patterns measured in IwA vs. *Controls*;
 297 see Methods) as a function of model layers (all four tested models shown separately in individual panels).
 298 The difference in ANN predictivity between *Controls* and IwA was largest at the deeper (more IT-like) layers of
 299 the models instead of earlier (more V1, V2, and V4-like) layers. Errorbars denote bootstrap confidence

300 intervals. Facial images shown in this figure are morphed and processed version of the original face images.
301 These images have full re-use permission.

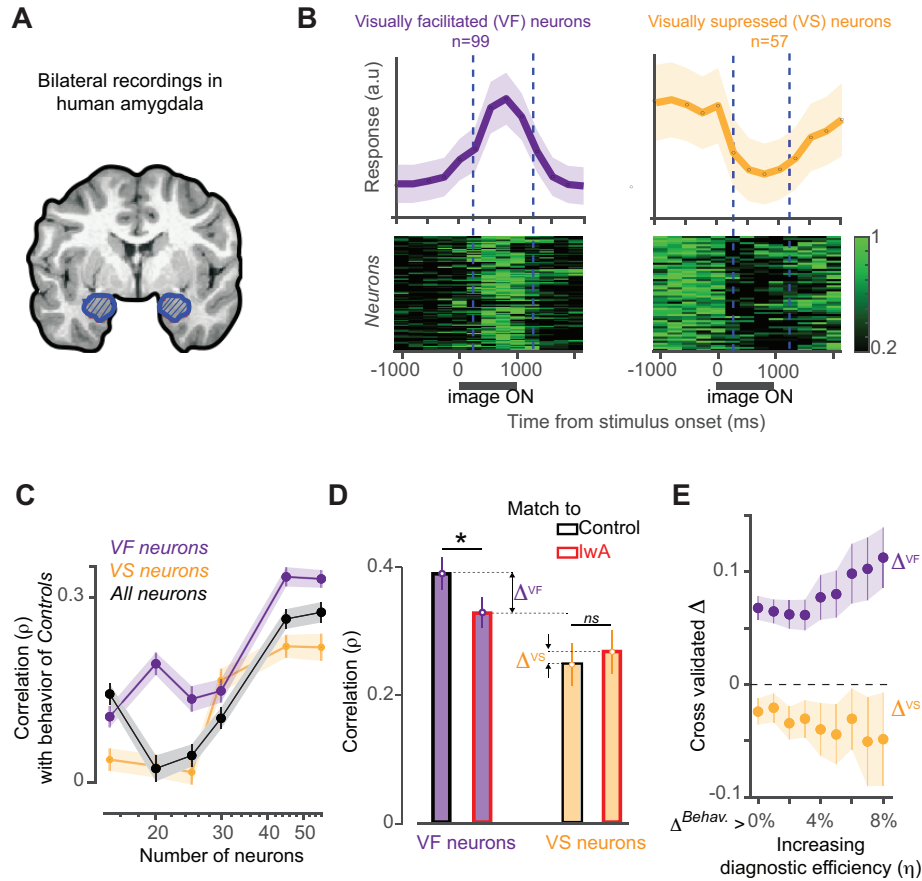
302

303 **Two distinct neural population coding schemes in the human** 304 **amygdala**

305 Wang et al.¹⁵ recorded bilaterally from implanted depth electrodes in the human
306 amygdala (schematic shown in Figure 3A) from patients with pharmacologically
307 intractable epilepsy. Subjects were presented each image for 1s (same as the task
308 description above⁴) to discriminate between two emotions, fear and happiness. Similar
309 to previous reports¹⁵, I observed two distinct population of neurons in the human
310 amygdala. These two populations were marked by significant response suppression
311 (visually suppressed (VS); 57 neurons; Figure 3B, right panel) and facilitation (visually
312 facilitated (VF); 99 neurons; Figure 3B, left panel) respectively, after the onset of the facial
313 image stimulus. I first tested how well the population-level activity (250-1500 ms post
314 image onset) of three specific subsamples of the amygdala neurons (VS only, VF only
315 and VS + VS neurons) predicted the behavioral patterns measured in human subjects. I
316 observed that each of these populations of VF, VS, and mixed (equal number of VS and
317 VF neurons) could significantly ($p < 0.0001$; permutation test for significance of
318 correlation) predict the image-level facial emotion judgments measured in *Controls*.
319 Figure 3C shows how these three populations predict the image-level behavior
320 measured in *Controls* as a function of the number of neurons sampled to build the neural
321 population decoders. Given that all of these groups exhibit an increase in behavioral
322 predictivity with the number of neurons, it is difficult to reject any of these decoding
323 models (with the current neural dataset). Therefore, in the following analyses I have
324 examined the VF and VS units separately. Next, I estimated how well the VS and VF
325 population predicted the behavioral patterns measured in the *Control* and IwA
326 respectively. Interestingly, I observed that similar to the ANN-IT behavior, neural
327 decodes out of the VF neurons in the human amygdala better match the *Control* group
328 behavior compared to the ones measured in IwA (Figure 3C; Δ^{VF} is significantly greater
329 than 0; permutation test of correlation: $p < 0.05$). However, the VS neurons did not show
330 this trend (Figure 3D; Δ^{VS} is not significantly different from 0; permutation test of
331 correlation; $p > 0.05$). Figure 3E shows how VF (and not VS) neurons become more
332 discriminatory of the IwA vs. *Control* behavior (i.e., Δ^{VF} increases) as we choose image-
333 sets with higher diagnostic efficiencies (η). Consistent with prior work, these results
334 provide evidence that neural responses in the human amygdala are implicated in atypical
335 facial processing in people with autism. However, the results presented here also
336 critically identify the VF neurons as a stronger candidate neural marker of the differences
337 in facial emotion processing observed in IwA.

338

339



340

341 **Figure 3. Facial emotion representation in the population neural activity of human amygdala. A.**

342 Schematic of bilateral amygdala (blue patch) recordings performed by Wang et al. **B.** Two distinct

343 population of neurons observed in the human amygdala. The visually facilitated (VF; shown in purple)

344 neurons ($n=99$) increased their responses after the onset of the face stimuli (top left panel: averaged

345 normalized spike rate across time; 250 ms time bins). The bottom left panel shows the normalized firing

346 rate across time for each VF neuron. The visually suppressed (VS; shown in yellow) neurons ($n=57$)

347 decreased their responses after the onset of the face stimuli (top right panel: averaged normalized spike

348 rate across time; 250 ms time bins). The bottom right panel shows the normalized firing rates across time

349 for each VS neuron. Errorbars denote SEM across neurons. **C.** An estimate (correlation) of how three

350 subsamples of neural populations, VS (yellow), VF (purple) and VS+VF ('All', black) predict the image-level

351 behavior measured in *Controls* as a function of the number of neurons sampled to build the neural decoders.

352 Errorbars denote bootstrapped CI. **D.** Comparison of how well the VS (yellow bars) and VF (purple bars)

353 neurons predict the behavior measured in *Controls* vs. IwA. The red and black edges denote the predictivity

354 of IwA and *Controls* respectively. Δ^{VF} and Δ^{VS} are the differences in the human amygdala (neural decode)

355 predictivity of facial emotion judgments measured in *Controls* and IwA from the VF and VS neurons

356 respectively. Errorbars denote bootstrap CI. **E.** Δ^{VF} and Δ^{VS} as function of image selection (which is

357 proportional to the diagnostic efficiency η estimated per image-set). The cross validation was done at the

358 level of subjects for each image selection. Errorbars denote bootstrap CI.

359

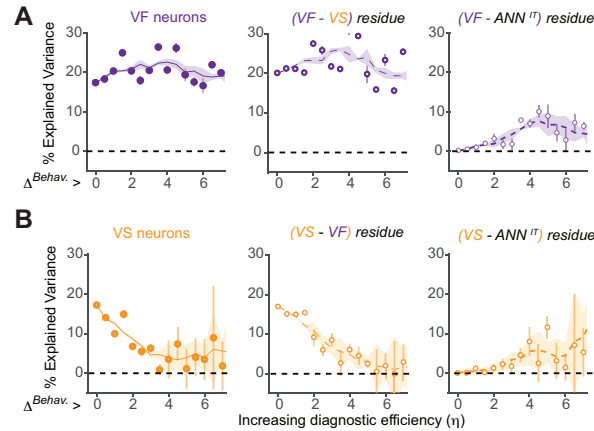
360

361 **ANN-IT features can explain a significant fraction of the**
362 **image-level behavioral predictivity of the human amygdala**
363 **population**

364 Given the significant predictivity of facial emotion judgments observed in the ANN IT
365 layers and the presence of strong anatomical connections between primate IT and
366 amygdala³², I further asked how much of the image-level predictivity estimated from the
367 amygdala activity is likely driven by input projections from the IT cortex. To test this, I first
368 asked (with a linear regression analyses; see Methods) how well the image-by-image
369 behavioral predictions from the ANN-IT models (AlexNet-fc7 tested here) can explain the
370 image-by-image neural decoding patterns estimated from the amygdala neurons
371 (separately for VS and VF neurons). The residue of this analyses (see Methods)
372 contained the variance in the amygdala decodes that was not explained by the predictions
373 of the ANN-IT models. Therefore, the amount of variance in the measured behavioral
374 patterns explained by this residue provides an estimate of how much of the behavior is
375 purely driven by the amygdala responses independent of the image-driven sensory
376 representations. Assuming a feedforward hierarchical circuit whereby the IT cortex drives
377 the human amygdala and not the other way around, a lower percentage of explained
378 variance (%EV) obtained after such an analysis should indicate that the source of the
379 signal in amygdala is at least partially coming from the IT cortex. Interestingly, this
380 analysis revealed that the behavioral predictivity (%EV) of the human amygdala is
381 significantly reduced once I regressed out the variance that is driven by the ANN-IT
382 responses. For instance, when considering all images (i.e., very low diagnostic efficiency
383 of the imageset), I observed that VS and VF neurons could explain approximately 17.24%
384 and 17.39% (a lower bound of the %EV since neural noise has not been accounted for)
385 of the behavioral variance (Figure 4A, B; left panel). However, once the ANN-IT driven
386 variance was regressed out these values significantly dropped to 0.06% and 0.2%
387 respectively (Figure 4A, B; right panel). Overall, VF neural residuals (after regressing out
388 ANN-IT predictions) explained significantly less variance at all tested η levels. VS neural
389 residuals explained significantly less variance only at lower η levels ($\Delta^{Behav} < 2.5\%$).
390 Given that VS neurons showed a drop in %EV for higher η levels, it is not surprising that
391 I did not observe any differences with the residual predictivity at those levels. Interestingly,
392 there was no significant change in %EV across the image selections when VS activity
393 was regressed out of VF activity (and vice versa; Figure 4A, B; middle panel), providing
394 further evidence that they largely support a complimentary coding scheme for facial
395 emotions within the amygdala. In sum, these results suggest that input projections from
396 the IT cortex into the amygdala³² might be the primary carrier of the facial emotion related
397 signals. Furthermore, the results also suggest a likely difference in how VS and VF
398 neurons are affected in lWA – with VF neurons being more diagnostic of the atypical
399 behavior observed in lWA.

400

401



402

403

404 **Figure 4. Amount of behavioral variance (measured in *Controls*) explained by different neural**
405 **markers. A.** Left panel: Percentage of behavioral variance explained by the human amygdala (VF) neural
406 activity as a function of the overall differences in image-level behavior between IwA and *Controls*. As
407 demonstrated in Figure 1F the x-axis is proportional to the diagnostic efficiency (η). Middle panel:
408 Percentage of variance explained by the residual (VS-based predictions regressed out of the predictions
409 from VF-based neural decodes). There was no significant change in %EV across the image selections
410 when VS was regressed out, suggesting a complimentary coding scheme. Right panel: Percentage of
411 behavioral variance explained by the residual (ANN-IT predictions regressed out of the predictions from
412 VF-based neural decodes). There was a significant difference (reduction in %EV) between the two cases
413 for all levels of tested η . **B.** Left panel: Percentage of behavioral variance explained by the human amygdala
414 (VS) neural activity as a function of the overall differences in image-level behavior between IwA and
415 *Controls*. Middle panel: Percentage of variance explained by the residual (VF-based predictions regressed
416 out of the predictions from VS-based neural decodes). There was no significant change in %EV across the
417 image selections when VF was regressed out, suggesting a complimentary coding scheme. Right panel:
418 Percentage of variance explained by the residual (ANN-IT predictions regressed out of the predictions from
419 VS-based neural decodes). There was a significant difference (reduction in %EV) between the two cases
420 while $\Delta^{Behav.}$ was less than 2. All %EV values were estimated in a cross validated way, wherein the image
421 selections and the final estimates were done based on different groups of subjects. Errorbars denote
422 bootstrapped CI.

423

424

425

426

427

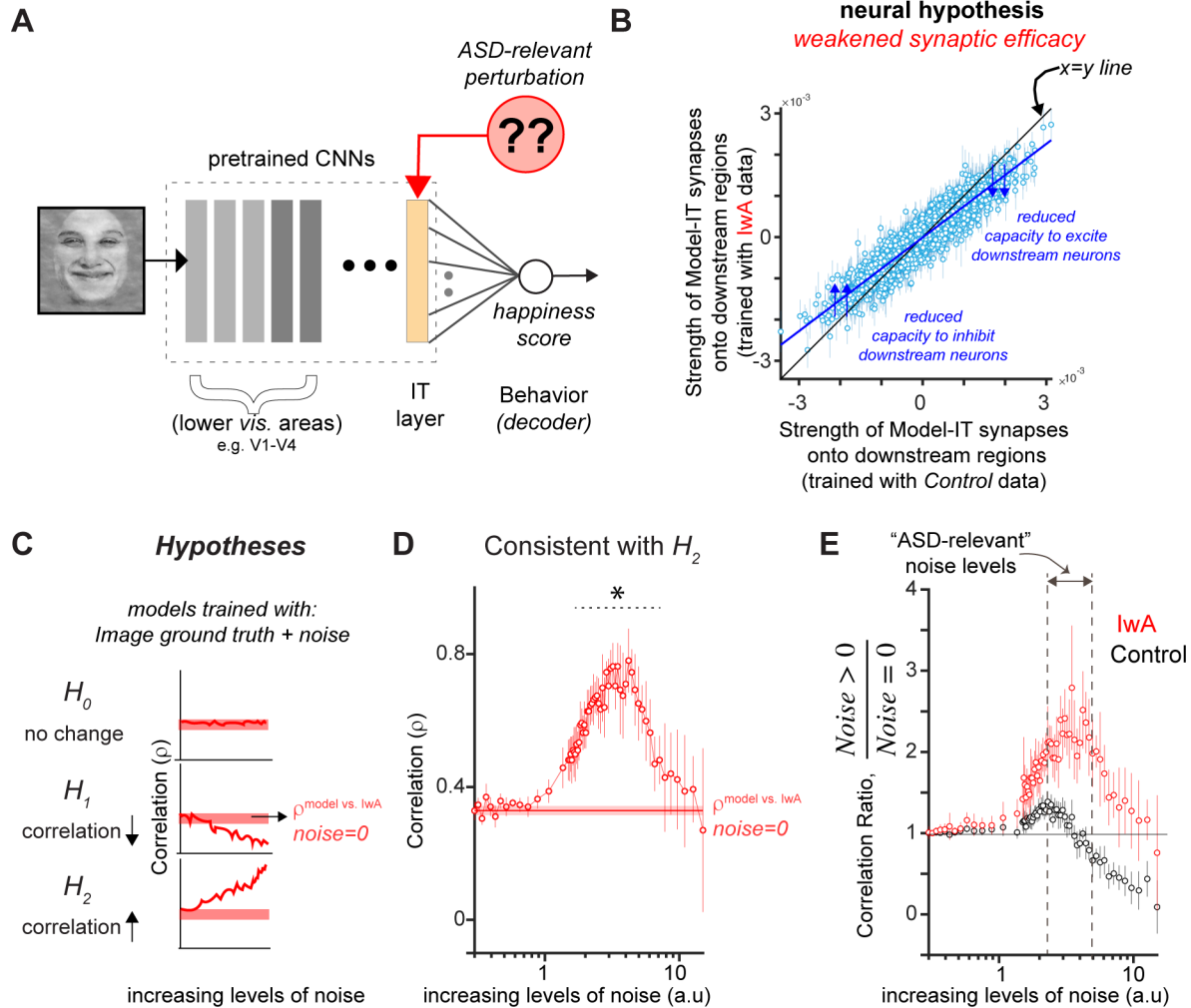
428

429

430 **In silico perturbations with additional noise in ANN-IT layers**
431 **improves the model's match with the behavior of individuals**
432 **with autism**

433 To further probe how IT representations might be different in lwA compared to *Controls*
434 (Figure 5A), I compared ANNs independently trained to predict the behavior of *Controls*
435 and lwA. I directly compared the learned weights, that is the synaptic strengths between
436 the model-IT layer and the behavioral output node in the two cases. I observed that are
437 model trained on the behavior measured in lwA yielded weaker synaptic strengths for
438 both excitatory (positively weighted) and inhibitory (negatively weighted) connections
439 (Figure 5B), compared to models trained to reproduce the behavior measured in *Controls*.
440 I further explored how this modest difference in the models could be simulated such that
441 an ANN trained on ground truth labels of human facial emotions could be transformed
442 into behaving more like what we observe in lwA. Based on previous studies^{33,34}, I
443 hypothesized that increased noise (scaled according to overall responsiveness of the
444 model units) in the sensory representations during learning could potentially yield weaker
445 synaptic strengths between the model-IT layer and the trained behavioral output node. Of
446 note, although a noisy representation likely yields a reduced specificity in behavioral
447 performance, an addition of specific amounts of noise does not necessarily guarantee a
448 stronger or weaker correlation with the image-level behavioral patterns observed in lwA.
449 Therefore, such in silico perturbations could produce three primary outcomes. First,
450 adding noise might produce no effects in the model's behavioral match with the behavior
451 of lwA (Figure 5C, top panel, H_0). Second, the added noise might weaken the correlation
452 achieved by a noiseless model (Figure 5C, middle panel, H_1). Third, and consistent with
453 an Autism Spectrum Disorder (ASD)-relevant mechanism, addition of noise could improve
454 the correlation with the image-level behavior measured in lwA (Figure 5C, bottom panel,
455 H_2). I observed that at specific levels of added noise (Figure 5D; dashed black line) during
456 the model training (transfer learning), the model's behavioral match with lwA significantly
457 improved (assessed by permutation test of correlation) beyond the levels noted with a
458 noise-free model (Figure 5D). In addition, this increase in the predictivity of lwA behavior
459 with addition of noise is significantly higher than that observed when compared to the
460 model's predictivity of the behavior measured in the *Controls* (as shown in Figure 5E).
461 Within the dashed black lines (Figure 5E), noise added to each model unit were drawn
462 from a normal distribution with zero mean and standard deviation equal to 2 to 5 times
463 the width of the response distribution of that unit across all tested images. Taken together,
464 this strongly suggests that additional noise in sensory representations is a very likely
465 candidate mechanism implicated in atypical facial emotion processing in adult with
466 autism.

467



468

469 **Figure 5. In silico experiments on ANNs to probe neural mechanisms underlying atypical facial**
 470 **emotion judgments in individuals with autism. A.** What changes can one induce in the model-IT layer
 471 to simulate the behavioral patterns measured in lWA? **B.** Comparison of synaptic strengths (weights)
 472 between ANN-IT and the behavioral node when models are independently trained with the behavior
 473 measured in lWA vs. *Controls*. ANN fits to behavior of lWA yielded weaker synaptic strengths for both
 474 excitatory (positively weighted) and inhibitory (negatively weighted) connections. Each blue dot refers to
 475 the weights in the connection between an individual model unit in the IT-layer and the decision ("level of
 476 happiness") node. **C.** Hypotheses and corresponding predictions H_0 : Addition of noise could lead to no
 477 differences in how it affects the model's match to behavior measured in lWA. H_1 : Addition of noise could
 478 reduce the models' match to behavior measured in lWA compared to the noise-free model. H_2 : Addition of
 479 noise could improve the models' match to the behavior measured in lWA compared to the noise-free model.
 480 H_2 supports the "high IT variability in autism" hypotheses. **D.** Correlation of ANN behavior with lWA as a
 481 function of levels of added noise. The results show that at specific noise regimes ANNs are significantly
 482 more predictive of the behavior measured in lWA compared to the noiseless model. Errorbars denote
 483 bootstrapped CI. **E.** Ratio of ANN behavioral predictivity of noisy vs. noise-free ANNs. At specific levels of
 484 noise, referred to as the Autism Spectrum Disorder (ASD)-relevant noise levels, the ANNs trained with
 485 noise show much higher predictivity for behavior measured in lWA while suffering a reduction in predictivity
 486 of the *Controls*. Errorbars denote bootstrapped CI. Facial images shown in this figure are morphed and
 487 processed version of the original face images. These images have full re-use permission.

488 Discussion

489
490 The overall goal of this study was to identify candidate neural and behavioral markers of
491 atypical facial emotion judgments observed in individuals with autism. Based on
492 discovering reliable image-by-image differences between the behavior of *Controls* and
493 IwA that could not be explained by categorical ambiguity in the stimuli, I reasoned that
494 such image-level variance could be leveraged to probe the neural mechanisms of
495 behavioral differences observed in IwA. Therefore, I used image-computable, brain-tissue
496 mapped artificial neural network models of primate vision to further probe the issue. By
497 using computational models (that have established brain tissue correlates) to explain
498 experimental data, I hereby demonstrate how such an approach could be used to probe
499 the neural mechanisms that underlie the differences in facial emotion processing
500 observed in individuals with autism. Below, I discuss the findings with their relevance to
501 future experiments and candidate mechanisms implicated in atypical facial emotion
502 recognition in IwA.
503

504 **ANN based predictions can be used to efficiently screen** 505 **images and provide neural hypotheses for more powerful** 506 **experiments**

507
508 A family of ANN models can currently predict a significant amount of variance measured
509 in various object recognition related behaviors and neural circuits³⁵. Given that the results
510 presented here demonstrate the ability of such ANNs to discriminate between the
511 behavior measured in *Controls* and IwA, we can further leverage the ANNs to screen
512 facial image stimuli and select images where the predicted behavioral differences are
513 maximum. Further, such models can be reverse engineered^{25,36} to synthesize images that
514 could achieve maximum differences to optimize behavioral testing and diagnosis. Such
515 deep image synthesis methods could also modify the facial images such that the
516 differences in the observed behavior between the *Controls* and IwA are minimized.
517 Although clearly at an early stage, such methods have a significant potential to improve
518 future cognitive therapies. Unlike many machine learning approaches that are not closely
519 tied to the computation and architecture of the primate brain, the ANNs used in this study
520 have established homologies with the primate brain and behavior³⁵. As demonstrated in
521 this study, these links allow us to relate the ANN predictions to distinct brain areas directly.
522 Specifically, the ANN results presented here suggest that population activity patterns in
523 areas like the human and macaque inferior temporal cortex are vital candidates for neural
524 markers of atypical facial processing in autism. The modeling results provide further
525 insights into the most affected aspects of the population responses, implicating noisier
526 sensory representations (see below) as a source of the differences in sensory
527 representation, learning and subsequent decision making. Besides the specific
528 hypotheses generated in this study, it is essential to note that ANN models of primate
529 vision are an active area of research, and we are witnessing the gradual emergence of
530 better brain-matched models^{29,37-39}. Therefore, this study establishes a critical link
531 between atypical face processing in autism and how to leverage ANNs to study this.

532

533 **Modeling results imply the need for more fine grain neural** 534 **measurements in the primate IT cortex and amygdala**

535 The ANN-based computational analyses in this study provide specific neural hypotheses
536 that can be tested using macaque electrophysiology and human fMRI experiments. First,
537 I observed that the ANN-IT layers could best discriminate between the behavior of
538 *Controls* vs. IwA. Therefore, such signals are likely also measurable in the primate IT
539 cortex and are key candidates for neural markers of atypical facial emotion processing
540 in autism. Given that most ANN models are feedforward-only or have minimal dynamics,
541 it will be critical to test how the different temporal components of IT population
542 responses carry the facial emotion signal. Similar to predictions of ANN-IT layers, I
543 observed that population activity in the human amygdala also better matches behavior
544 measured in the *Controls* than IwA. There can be multiple reasons for the observed
545 differences in behavioral predictivity. First, it is possible that due to the atypical
546 development of the human amygdala in IwA, the behavior they exhibit does not match
547 well with the neural decodes out of the neurotypical amygdala. Second, the lack of
548 predictivity might be carried forward from responses in the IT cortex -- as predicted by
549 the ANNs. The current study attempted to disambiguate between these two factors. I
550 asked how well ANN-IT predictions can account for the amygdala activity's behavioral
551 patterns. Indeed, the image-level predictivity of facial emotion judgments observed in
552 the human amygdala's population activity (both VF and VS neurons) was significantly
553 explained away by the ANN-IT features (Figure 4A, B; left panel). This result is consistent
554 with the hypothesis that the higher-level visual cortices (like IT) primarily drive the facial
555 affect signal observed in the human amygdala. Simultaneous neural recordings in IT and
556 amygdala or finer grain causal perturbation experiments need to be conducted to test
557 this hypothesis more directly. Notably, the behavioral mismatch (neural decodes
558 vs. *Control/IwA* behavior) was specific to the decodes constructed from the VF neurons
559 (and not VS neurons). Therefore, future experimental investigations should dissect the
560 role of IT cortex and how it functionally influences the VF and VS neurons, which are
561 likely part of a complimentary coding scheme. Furthermore, it will be essential to
562 examine how the IT cortical activity is driven by feedback projections from the amygdala,
563 given that evidence for the importance of such connections from ventrolateral PFC has
564 been demonstrated for object recognition⁴⁰.

565

566 **High variability in sensory representation can lead to weaker** 567 **efferent synaptic strengths during learning and development**

568 In a psychophysical discrimination task, the typical consequence of having a noisy
569 detector is a reduction in the sensitivity of performance, which manifests as a reduced
570 estimated slope of the psychometric function. This is consistent with what Wang and
571 Adolphs⁴ had observed. Given that the idea of higher sensory variability in autism is also
572 consistent with previous findings³⁴, I considered this as a potential neural mechanism
573 that could explain the image-level differences I have observed in the facial emotion

574 discrimination behavior in IwA. Therefore, I tested the “increased sensory noise
575 hypothesis” to test whether such a perturbation could simulate the weaker efferent
576 synaptic connections from IT-like layers as revealed by the ANN based analyses (Figure
577 5B). Indeed, addition of noise during learning made the ANN behavior more matched
578 with that observed in IwA. First, this could suggest that perhaps the behavior measured
579 in IwA results from additional noise in the sensory representations that affects the
580 subjects’ behavior during the task. However, this could also be the result of executing
581 an inference engine (in the brain) that learned its representations under high sensory
582 noise during development (as a child). An estimate of noise levels (sensory cortical signal
583 variability) in children with autism and a quantitative probe into how that could potentially
584 interact with learning new tasks is essential to test this hypothesis. As demonstrated in
585 this study, the ANN models provide a very efficient framework to generate more
586 diagnostic image-sets for these future studies given that we can simulate any level (and
587 type) of noise under different learning regimes and make predictions on effect sizes.
588 Such model-driven hypotheses are likely to play a vital role in guiding future experimental
589 efforts and inferences.
590

591 **High variability in sensory representation can qualitatively** 592 **explain other ASD-specific behavioral reports**

593
594 Addition of noise during the transfer learning procedure of the ANN models made the
595 model’s behavioral output more consistent with the behavior measured in IwA (Figure
596 5D). Such a mechanism can indeed qualitatively explain other previous behavioral
597 observations made in individuals with autism. For example, Behrmann et al.⁴¹ observed
598 that reaction times measured during object discrimination tasks, in adults with autism
599 were significantly higher than the *Control* subjects. This difference was especially high
600 during more fine-grained discrimination tasks. Such a behavioral phenomenon can be
601 explained by an increase in sensory noise in IwA that leads to longer time requirements
602 during integration of information⁴², and weaker performances on finer discrimination
603 tasks. The ANN based approach demonstrated in this study, however, provides guidance
604 beyond the qualitative predictions of overall effect types. Specific image-level predictions
605 provided by ANNs will help researchers to design more diagnostic behavioral experiments
606 and make measurements that can efficiently discriminate among competing models of
607 brain mechanisms.
608

609 **Potential underlying mechanisms behind increased neural** 610 **variability**

611 An imbalance in the ratio of the excitatory and inhibitory processes in cortical circuits has
612 been proposed as an underlying mechanism for various atypical behaviors observed in
613 autism⁴³. I speculate that such an E/I imbalance could arise due to lower inhibition in the
614 cortical networks. This could lead to larger neural variability and a subsequent noisier,
615 less efficient sensory processing. Therefore, the results observed in the in-silico
616 experiments are not biologically implausible. In fact, genetic mutations that impact the
617 generation and function of interneurons have been previously linked with autism^{44,45}.

618 Therefore, cell-type specific causal perturbation approaches are necessary to test
619 whether a decreased inhibition in the visuocortical pathway (especially in the primate IT
620 cortex) leads to noisier sensory representations and can reproduce the specific image-
621 level differences in facial emotion processing reported in this study. The image-level
622 behavioral measurements and ANN predictions reported here will enable such stronger
623 forms of hypothesis testing during the interpretation of such experimental results.

624 **Methods and Materials**

625

626 **Human Behavior**

627 In this study, I have re-analyzed behavioral data that was previously collected and used
628 in a study by Wang and Adolphs⁴. The raw behavioral dataset was kindly shared via
629 personal communication.

630

631 **Participants**

632 In the original study (for further details see⁴), eighteen high-functioning participants with
633 ASD (15 male) were recruited. All ASD participants met DSM-V/ICD-10 diagnostic criteria
634 for autism spectrum disorder (ASD) and met the cutoff scores for ASD on the Autism
635 Diagnostic Observation Schedule-2 (ADOS-2) revised scoring system for Module 4, and
636 the Autism Diagnostic Interview-Revised (ADI-R) or Social Communication Questionnaire
637 (SCQ) when an informant was available. The ASD group had a full-scale IQ (FSIQ) of
638 105 ± 13.3 (from the Wechsler Abbreviated Scale of Intelligence-2), a mean age of
639 30.8 ± 7.40 years, a mean Autism Spectrum Quotient (AQ) of 29.3 ± 8.28 , a mean SRS-2
640 Adult Self Report (SRS-A-SR) of 84.6 ± 21.5 , and a mean Benton score of 46.1 ± 3.89
641 (Benton scores 41–54 were in the normal range). ADOS item scores were not available
642 for two participants, so we were unable to utilize the revised scoring system. But these
643 individuals' original ADOS algorithm scores all met the cutoff scores for ASD.

644

645 Fifteen neurologically and psychiatrically healthy participants with no family history of
646 ASD (11 male) were recruited as *Controls*. *Controls* had a comparable FSIQ of 107 ± 8.69
647 (two-tailed t-test, $P=0.74$) and a comparable mean age of 35.1 ± 11.4 years ($P=0.20$), but
648 a lower AQ (17.7 ± 4.29 , $P=4.62 \times 10^{-5}$) and SRS-A-SR (51.0 ± 30.3 , $P=0.0039$) as expected.
649 Participants gave written informed consent, and all original experiments were approved
650 by the Caltech Institutional Review Board. All participants had normal or corrected-to-
651 normal visual acuity. No enrolled participants were excluded for any reasons.

652

653 **Facial emotion judgment task**

654 During the task, Wang and Adolphs⁴ asked participants to discriminate between two
655 emotions, fear and happiness. The image-set includes faces of four individuals (2 female)
656 each posing fear and happiness expressions from the STOIC database (Roy et al. 2007),
657 which are expressing highly recognizable emotions. To generate the morphed expression
658 continua for the experiments, the authors interpolated pixel value and location between
659 fearful exemplar faces and happy exemplar faces using a piece-wise cubic-spline
660 transformation over a Delaunay tessellation of manually selected control points. They
661 created 5 levels of fear-happy morphs, ranging from 30% fear/70% happy to 70%
662 fear/30% happy in steps of 10% (Figure 1B). Low-level image properties were equalized
663 using the SHINE toolbox⁴⁶. In each trial, a face was presented for 1 second followed by
664 a question prompt asking participants to make the best guess of the facial emotion (Figure
665 1A). After stimulus offset, participants had 2 seconds to respond, otherwise the trial was
666 aborted and discarded. Participants were instructed to respond as quickly as possible,
667 but only after stimulus offset. No feedback message was displayed, and the order of faces
668 was completely randomized for each participant. Images were presented approximately

669 in the central 12° of visual angle. A subset of the participants (11 participants with autism
670 and 11 *Controls*) also performed confidence ratings after emotion judgment and a 500
671 ms blank screen, participants were asked to indicate their confidence by pushing the
672 button '1' for 'very sure', '2' for 'sure' or '3' for 'unsure'. This question also had 2
673 seconds to respond. All images used in this study has free re-use permission as set
674 here¹⁵.

675

676 **Estimating image-level behavioral reliability**

677

678 To estimate the image-level behavioral reliability (Figure 1D), I first estimated the
679 probability of choosing "Happy" per image in each subject (15 Controls, 18 IwA) -- referred

680 to as the P_C and the P_{IwA} vectors. Then, for each possible combination of selecting 2
681 subjects from the subject pools, I estimated the subject-to-subject Kendall rank correlation
682 coefficient. This was done separately for the Controls and IwA, leading to the red and
683 black histograms in Figure 1D respectively. These correlations scores are not corrected
684 by the individual subjects' internal reliability (across trials). Therefore, they represent the
685 lower bound of the inter subject correlations.

686

687 **Estimating noise ceilings for IwA vs. Control correlations**

688

689 I define the noise ceiling of a correlation as the highest possible value of correlation
690 expected given the noise measured independently in the two variables that are being

691 tested. To estimate this, first I individually estimate the split half reliability of the P_C and

692 the P_{IwA} vectors. Each split is constructed with a random sampling of half of the subjects
693 and taking the average across them and doing same for the other half of the subjects.
694 For each iteration, such splits were made, and the correlation between the resulting
695 vectors was computed. This correlation score was corrected by the Spearman-Brown
696 correction procedure to account for the halving of subject numbers. I then computed the
697 average across 100 such iterations, referred to as $\rho_{P_C^1, P_C^2}$ and $\rho_{P_{IwA}^1, P_{IwA}^2}$ for the *Controls*

698 and IwA respectively. The noise ceiling was then estimated as,

699

700

$$\sqrt{\rho_{P_C^1, P_C^2} * \rho_{P_{IwA}^1, P_{IwA}^2}}$$

701

702 Intuitively, if both groups provided noiseless data, then these reliabilities should be each
703 at 1, and therefore the noise ceiling shall also be set at 1. Noisy data will lead to <1 values
704 for the individual $\rho_{P_C^1, P_C^2}$ and $\rho_{P_{IwA}^1, P_{IwA}^2}$ reliabilities, and hence the noise ceiling shall

705 also be <1. Of note, each selection of image with result in a different P vector and
706 therefore will result in a slightly different noise ceiling estimate, as demonstrated in Figure
707 1E (two green lines).

708

709 **Estimating cross-validated diagnostic efficiency (η) of image-sets**

710 Diagnostic Efficiency (η); shown in Figure 1E, and 1F) of an image-set is defined as the
711 cross-validated estimate of the difference between the noise ceiling and the raw
712 correlation between the P_C and the P_{IwA} vectors. The cross validation is achieved by the
713 choosing the images based on a specific subset of subjects and then measuring the noise
714 ceiling and the raw correlation on a different held-out set of subjects. For efficient
715 collection of human subject data that could optimally discriminate between the behavior
716 measured in *Controls* and *IwA*, one must aspire for the highest η values for image-sets.
717

718 **Depth recording in human amygdala**

719
720 In this study I have re-analyzed the neural data that was previously collected and used in
721 a study by Wang et al.¹⁵. The raw neural dataset was kindly shared via personal
722 communication. Wang and colleagues recorded bilaterally from implanted depth
723 electrodes in the amygdala from patients with pharmacologically intractable epilepsy.
724 Target locations in the amygdala were verified using post-implantation structural MRIs.
725 At each site, they recorded from eight 40 μm microwires inserted into a clinical electrode.
726 Bipolar wide-band recordings (0.1–9 kHz), using one of the eight microwires as reference,
727 were sampled at 32 kHz and stored continuously for off-line analysis with a Neuralynx
728 system (Digital Cheetah; Neuralynx, Inc.). The raw signal was filtered with a zero-phase
729 lag 300-3000 Hz bandpass filter and spikes were sorted using a semiautomatic template
730 matching algorithm. Units were carefully isolated and spike sorting quality were assessed
731 quantitatively. Subjects were presented each image for 1s (similar to the task description
732 above) to discriminate between two emotions, fear and happiness.
733

734 **Selection of neurons for analyses**

735
736 In the original study, only units with an average firing rate of at least 0.2 Hz (entire task)
737 were considered. Only single units were considered. In addition to that, in this study I
738 have further restricted the neural dataset to neurons that have a significant visual
739 response (both increase and decrease). To estimate that I compared the neural firing
740 rates (per image) averaged across two specific time bins, [-1000 0] and [250 1250], where
741 0 is the onset of the image. If the paired Wilcoxon Signed Rank test between these two
742 firing rate vectors were significant, the site was considered for further analyses. Thus, I
743 considered 156 total neurons: 99 visually facilitated (VF) neurons and 57 visually
744 suppressed (VS) neurons.
745
746
747

748 **Decoding facial emotion judgment from neural population activity**

749
750 To decode facial emotion judgments from the neural responses per image, I used a linear
751 model that linked the neural responses to the levels of happiness (ground truth from
752 image generation). Building the model, essentially involves solving a regression problem

753 estimating the weights (\vec{w}) per neuron and a *bias* term. I used a partial least squares
754 (MATLAB command: *plsregress*) regression procedure, using 15 retained components. I
755 also used 10-fold cross validation. For each fold, the model was trained (i.e., \vec{w} and *bias*
756 were estimated) using the data from the other 9 folds (training data), and predictions were
757 generated for the held-out fold (test images). This was repeated for each of the folds and
758 the entire procedure was repeated 100 times. The predictions of the trained neural model
759 on the held-out test images were used for future correlation analyses. Given the training
760 scheme, every image was assigned as the test-image once per iteration.
761

762 **ANN models of primate vision**

763
764 The term "model" in this study always refer to a specific modification of a pre-trained ANN.
765 For instance, I have used an Image-Net pretrained deep neural network, AlexNet to build
766 multiple models. Each model was constructed by deleting all layers succeeding a given
767 layer. For instance, the '*cnv5*' model was built by removing all layers of AlexNet that
768 followed the output of its fifth convolutional layer. The feature activations from the fifth
769 convolutional layer output were then trained with the linear regression procedure (similar
770 to the neural decodes).
771

772 **Estimating model facial emotion judgment behavior**

773
774 To decode facial emotion judgments from the model responses per image, I used the
775 same linear modeling approach as the neural data (see above), that linked the model
776 feature activations to the level of happiness (ground truth from image generation). The
777 model features, per layer, were extracted using the MATLAB command *activations* for
778 AlexNet²⁸, VGGFace³⁰ and EmotionNet³¹ in MATLAB-R 2020b. For the CORnet-S²⁹
779 model, I used the code from: <https://github.com/dicarololab/CORnet>.
780

781 **Estimation of discriminatory index (d')**

782 The discrimination index was computed to quantify the difference between the match of
783 the ANNs' (models per layer) behavioral predictions to the behavior measured in *Controls*
784 and *IwA* (as shown in Figure 2E). It was calculated as:

$$785 \frac{\rho^{Control} - \rho^{IwA}}{\sqrt{\left\{\frac{1}{2} * (\sigma_{Control}^2 + \sigma_{IwA}^2)\right\}}}$$

786 where $\rho^{Control}$ and ρ^{IwA} was the correlation between ANN predictions and behavior
787 measured in *Controls* and *IwA* respectively. $\sigma_{Control}$ and σ_{IwA} was the standard deviation
788 of the bootstrap estimates of the correlations with random subsampling features from
789 the model layers. To make the comparisons fair across all layers, 1000 features were
790 randomly subsampled (without repetition) 100 times to estimate the ANN predictions.
791

792 **Estimation of residuals between ANN-IT and human amygdala's** 793 **behavioral predictions**

794 I first estimated the cross-validated test predictions (ANN^{Pred}) of behavioral patterns from
795 an ANN-IT layer (e.g., AlexNet 'fc7' model used in the study) using the partial least
796 squares regression method. The ground truth values of image-level facial happiness were
797 used as the dependent variable in this analysis. Next, I used the same algorithm but with
798 the human amygdala neural features (instead of the ANN-IT features) as the predictors
799 to estimate the neurally decoded behavioral patterns ($Amygdala^{Pred}$). I then used a
800 generalized linear regression model (MATLAB: *glmfit*) to estimate the residues while
801 using ANN^{Pred} as the predictor and $Amygdala^{Pred}$ as the dependent variable. The
802 square of the Pearson correlation (%EV) between this residue vector (one value per
803 image) and the image-level behavioral vector (Probability of choosing "Happy" per image)
804 measured in the *Controls* is plotted in the y-axis of Figure 4 (left panels). These %EV
805 values were corrected by the noise estimates in the behavioral data per image selection.
806 In addition, all %EV values were estimated in a cross validated way, wherein the image
807 selections and the final estimates were done based on different groups of subjects.

808 809 **In silico model perturbation and training**

810
811 *Generation of activity scaled additive noise values:* To estimate how much noise shall be
812 added to each unit (feature) of the model layer, I used the following procedure. First, I
813 estimated the standard deviation (σ , across all 28 images) of the activation distribution
814 per unit in a noise-free model. The addition of noise was made proportional to this value.
815 To vary noise levels, a scalar factor (C ; x-axis in Figure 5D and 5E) was multiplied with σ
816 per unit. For each unit, the noise added was drawn from a normal distribution that had a
817 standard deviation of $C*\sigma$.

818
819 *Training the model with and without noise:* To simulate a learning scheme with noise, I
820 modified the model feature activations in the following way. During training of the
821 regression model (i.e., estimating \vec{w} and *bias*), the noisy version of the model was
822 generated by concatenating 1000 randomly drawn features (which were fixed for each
823 iteration of the procedure), with ten repetitions of the same features but with the added
824 noise on top of it. This procedure was repeated several times to estimate the variance in
825 the model predictions per noise level. For the noise free model, the same 1000 randomly
826 drawn features were repeated without addition of any noise.

827 828 **Statistics**

829
830 All correlation scores reported in this study are Kendall rank coefficients (unless otherwise
831 mentioned). For significance tests of correlations (between two variables of interest), I
832 have used a bootstrapped permutation test. To do this, I first constructed a null hypothesis
833 by mixing the two variables and then randomly drew (as many times as the number of
834 elements in the original variable) with replacements two elements from the mixed dataset

835 to create two vectors. These two vectors can be constructed multiple times (typically
836 >100) and correlated. The resulting correlation distribution was considered as the null
837 hypothesis. Then the true raw correlation was compared to this distribution to determine
838 a p-value of rejecting the null distribution.

839 **Data and Code Availability**

840

841 All the data and code used in this study will be freely available to download and use
842 during the time of journal publication from [https://github.com/kohitij-](https://github.com/kohitij-kar/2021_faceEmotion_ASD)
843 [kar/2021_faceEmotion_ASD](https://github.com/kohitij-kar/2021_faceEmotion_ASD).

844

845 **Acknowledgments**

846

847 I thank R. Adolphs, P. Sinha (and Sinha Lab members), and J.J. DiCarlo for helpful
848 comments and discussions. I thank S. Wang for sharing the behavioral and neural
849 datasets used in this study. I thank S. Wang, S. Sanghavi, A. Peter, and Y. Bai for
850 helpful comments on the manuscript.

851 Bibliography

- 852
- 853 1 Adolphs, R., Sears, L. & Piven, J. Abnormal processing of social information from faces
854 in autism. *J Cogn Neurosci* **13**, 232-240, doi:10.1162/089892901564289 (2001).
- 855 2 Golarai, G., Grill-Spector, K. & Reiss, A. L. Autism and the development of face
856 processing. *Clin Neurosci Res* **6**, 145-160, doi:10.1016/j.cnr.2006.08.001 (2006).
- 857 3 Kennedy, D. P. & Adolphs, R. Perception of emotions from facial expressions in high-
858 functioning adults with autism. *Neuropsychologia* **50**, 3313-3319,
859 doi:10.1016/j.neuropsychologia.2012.09.038 (2012).
- 860 4 Wang, S. & Adolphs, R. Reduced specificity in emotion judgment in people with autism
861 spectrum disorder. *Neuropsychologia* **99**, 286-295,
862 doi:10.1016/j.neuropsychologia.2017.03.024 (2017).
- 863 5 Kanwisher, N., McDermott, J. & Chun, M. M. The fusiform face area: a module in human
864 extrastriate cortex specialized for face perception. *J Neurosci* **17**, 4302-4311 (1997).
- 865 6 Tsao, D. Y. & Livingstone, M. S. Mechanisms of face perception. *Annu Rev Neurosci* **31**,
866 411-437, doi:10.1146/annurev.neuro.30.051606.094238 (2008).
- 867 7 Tsao, D. Y., Freiwald, W. A., Knutsen, T. A., Mandeville, J. B. & Tootell, R. B. Faces and
868 objects in macaque cerebral cortex. *Nat Neurosci* **6**, 989-995, doi:10.1038/nn1111 (2003).
- 869 8 Tsao, D. Y., Moeller, S. & Freiwald, W. A. Comparing face patch systems in macaques
870 and humans. *Proc Natl Acad Sci U S A* **105**, 19514-19519, doi:10.1073/pnas.0809662105
871 (2008).
- 872 9 Freiwald, W. A., Tsao, D. Y. & Livingstone, M. S. A face feature space in the macaque
873 temporal lobe. *Nat Neurosci* **12**, 1187-1196, doi:10.1038/nn.2363 (2009).
- 874 10 Adolphs, R., Tranel, D., Damasio, H. & Damasio, A. Impaired recognition of emotion in
875 facial expressions following bilateral damage to the human amygdala. *Nature* **372**, 669-672,
876 doi:10.1038/372669a0 (1994).
- 877 11 Adolphs, R. Fear, faces, and the human amygdala. *Curr Opin Neurobiol* **18**, 166-172,
878 doi:10.1016/j.conb.2008.06.006 (2008).
- 879 12 Rutishauser, U., Mamelak, A. N. & Adolphs, R. The primate amygdala in social
880 perception - insights from electrophysiological recordings and stimulation. *Trends Neurosci* **38**,
881 295-306, doi:10.1016/j.tins.2015.03.001 (2015).
- 882 13 Broks, P. *et al.* Face processing impairments after encephalitis: amygdala damage and
883 recognition of fear. *Neuropsychologia* **36**, 59-70, doi:10.1016/s0028-3932(97)00105-x (1998).
- 884 14 Adolphs, R. *et al.* Recognition of facial emotion in nine individuals with bilateral
885 amygdala damage. *Neuropsychologia* **37**, 1111-1117, doi:10.1016/s0028-3932(99)00039-1
886 (1999).
- 887 15 Wang, S. *et al.* The human amygdala parametrically encodes the intensity of specific
888 facial emotions and their categorical ambiguity. *Nat Commun* **8**, 14821,
889 doi:10.1038/ncomms14821 (2017).
- 890 16 Behrmann, M., Thomas, C. & Humphreys, K. Seeing it differently: visual processing in
891 autism. *Trends Cogn Sci* **10**, 258-264, doi:10.1016/j.tics.2006.05.001 (2006).
- 892 17 Robertson, C. E. & Baron-Cohen, S. Sensory perception in autism. *Nat Rev Neurosci*
893 **18**, 671-684, doi:10.1038/nrn.2017.112 (2017).
- 894 18 Uljarevic, M. & Hamilton, A. Recognition of emotions in autism: a formal meta-analysis.
895 *J Autism Dev Disord* **43**, 1517-1526, doi:10.1007/s10803-012-1695-5 (2013).
- 896 19 Lozier, L. M., Vanmeter, J. W. & Marsh, A. A. Impairments in facial affect recognition
897 associated with autism spectrum disorders: a meta-analysis. *Dev Psychopathol* **26**, 933-945,
898 doi:10.1017/S0954579414000479 (2014).

- 899 20 Ekman, P. & Keltner, D. Universal facial expressions of emotion. Segerstrale U, P.
900 Molnar P, eds. Nonverbal communication: Where nature meets culture, 27-46 (1997).
- 901 21 Rajalingham, R., Schmidt, K. & DiCarlo, J. J. Comparison of Object Recognition
902 Behavior in Human and Monkey. *J Neurosci* **35**, 12127-12136, doi:10.1523/JNEUROSCI.0573-
903 15.2015 (2015).
- 904 22 Rajalingham, R. *et al.* Large-scale, high-resolution comparison of the core visual object
905 recognition behavior of humans, monkeys, and state-of-the-art deep artificial neural networks.
906 *bioRxiv*, 240614 (2018).
- 907 23 Khaligh-Razavi, S. M. & Kriegeskorte, N. Deep supervised, but not unsupervised,
908 models may explain IT cortical representation. *PLoS Comput Biol* **10**, e1003915,
909 doi:10.1371/journal.pcbi.1003915 (2014).
- 910 24 Yamins, D. L. *et al.* Performance-optimized hierarchical models predict neural
911 responses in higher visual cortex. *Proc Natl Acad Sci U S A* **111**, 8619-8624,
912 doi:10.1073/pnas.1403112111 (2014).
- 913 25 Bashivan, P., Kar, K. & DiCarlo, J. J. Neural population control via deep image
914 synthesis. *Science* **364**, doi:10.1126/science.aav9436 (2019).
- 915 26 Kar, K., Kubilius, J., Schmidt, K., Issa, E. B. & DiCarlo, J. J. Evidence that recurrent
916 circuits are critical to the ventral stream's execution of core object recognition behavior. *Nat*
917 *Neurosci* **22**, 974-983, doi:10.1038/s41593-019-0392-5 (2019).
- 918 27 Cadena, S. A. *et al.* Deep convolutional models improve predictions of macaque V1
919 responses to natural images. *PLoS Comput Biol* **15**, e1006897,
920 doi:10.1371/journal.pcbi.1006897 (2019).
- 921 28 Krizhevsky, A., Sutskever, I. & Hinton, G. E. in Proceedings of the 25th International
922 Conference on Neural Information Processing Systems - Volume 1 1097-1105 (Curran
923 Associates Inc., Lake Tahoe, Nevada, 2012).
- 924 29 Kubilius, J. *et al.* in Advances in Neural Information Processing Systems. 12785-12796.
- 925 30 Parkhi, O. M., Vedaldi, A. & Zisserman, A. Deep face recognition. (2015).
- 926 31 García, L. Face, Age and Emotion Detection. *MATLAB Central File Exchange* (2021).
- 927 32 Webster, M. J., Ungerleider, L. G. & Bachevalier, J. Connections of inferior temporal
928 areas TE and TEO with medial temporal-lobe structures in infant and adult monkeys. *J*
929 *Neurosci* **11**, 1095-1116 (1991).
- 930 33 MacDonald, S. W., Nyberg, L. & Backman, L. Intra-individual variability in behavior: links
931 to brain structure, neurotransmission and neuronal activity. *Trends Neurosci* **29**, 474-480,
932 doi:10.1016/j.tins.2006.06.011 (2006).
- 933 34 Haigh, S. M., Heeger, D. J., Dinstein, I., Minshew, N. & Behrmann, M. Cortical variability
934 in the sensory-evoked response in autism. *Journal of autism and developmental disorders* **45**,
935 1176-1190 (2015).
- 936 35 Schrimpf, M. *et al.* Brain-score: Which artificial neural network for object recognition is
937 most brain-like? *BioRxiv*, 407007 (2018).
- 938 36 Xiao, W. & Kreiman, G. XDream: Finding preferred stimuli for visual neurons using
939 generative networks and gradient-free optimization. *PLoS Comput Biol* **16**, e1007973,
940 doi:10.1371/journal.pcbi.1007973 (2020).
- 941 37 Nayebi, A. *et al.* in Advances in Neural Information Processing Systems. 5290-5301.
- 942 38 Lee, H. *et al.* Topographic deep artificial neural networks reproduce the hallmarks of the
943 primate inferior temporal cortex face processing network. *bioRxiv* (2020).
- 944 39 Zhuang, C. *et al.* Unsupervised neural network models of the ventral visual stream. *Proc*
945 *Natl Acad Sci U S A* **118**, doi:10.1073/pnas.2014196118 (2021).

- 946 40 Kar, K. & DiCarlo, J. J. Fast Recurrent Processing via Ventrolateral Prefrontal Cortex Is
947 Needed by the Primate Ventral Stream for Robust Core Visual Object Recognition. *Neuron* **109**,
948 164-176 e165, doi:10.1016/j.neuron.2020.09.035 (2021).
- 949 41 Behrmann, M. *et al.* Configural processing in autism and its relationship to face
950 processing. *Neuropsychologia* **44**, 110-129, doi:10.1016/j.neuropsychologia.2005.04.002
951 (2006).
- 952 42 Ratcliff, R., Smith, P. L., Brown, S. D. & McKoon, G. Diffusion Decision Model: Current
953 Issues and History. *Trends Cogn Sci* **20**, 260-281, doi:10.1016/j.tics.2016.01.007 (2016).
- 954 43 Rubenstein, J. L. & Merzenich, M. M. Model of autism: increased ratio of
955 excitation/inhibition in key neural systems. *Genes Brain Behav* **2**, 255-267, doi:10.1034/j.1601-
956 183x.2003.00037.x (2003).
- 957 44 Chao, H. T. *et al.* Dysfunction in GABA signalling mediates autism-like stereotypies and
958 Rett syndrome phenotypes. *Nature* **468**, 263-269, doi:10.1038/nature09582 (2010).
- 959 45 Sohal, V. S. & Rubenstein, J. L. R. Excitation-inhibition balance as a framework for
960 investigating mechanisms in neuropsychiatric disorders. *Mol Psychiatry* **24**, 1248-1257,
961 doi:10.1038/s41380-019-0426-0 (2019).
- 962 46 Willenbockel, V. *et al.* Controlling low-level image properties: the SHINE toolbox. *Behav*
963 *Res Methods* **42**, 671-684, doi:10.3758/BRM.42.3.671 (2010).