

1 **Using Computer-vision and Machine Learning to Automate Facial Coding of Positive and**
2 **Negative Affect Intensity**

3
4 Nathaniel Haines

5 The Ohio State University

6
7 Matthew W. Southward

8 The Ohio State University

9
10 Jennifer S. Cheavens

11 The Ohio State University

12
13 Theodore Beauchaine

14 The Ohio State University

15
16 Woo-Young Ahn

17 Seoul National University

18
19 Address correspondence and reprint requests to Woo-Young Ahn, Department of Psychology,

20 Seoul National University, Seoul, Korea 08826. Email: wahn55@snu.ac.kr. Tel: +82-2-880-2538

21

22

23

24

25

26

27

28 **Abstract**

29 Facial expressions are fundamental to interpersonal communication, including social interaction,
30 and allow people of different ages, cultures, and languages to quickly and reliably convey
31 emotional information. Historically, facial expression research has followed from discrete
32 emotion theories, which posit a limited number of distinct affective states that are represented
33 with specific patterns of facial action. Much less work has focused on dimensional features of
34 emotion, particularly positive and negative affect intensity. This is likely, in part, because
35 achieving inter-rater reliability for facial action and affect intensity ratings is painstaking and
36 labor-intensive. We use computer-vision and machine learning (CVML) to identify patterns of
37 facial actions in 4,648 video recordings of 125 human participants, which show strong
38 correspondences to positive and negative affect intensity ratings obtained from highly trained
39 coders. Our results show that CVML can both (1) determine the importance of different facial
40 actions that human coders use to derive positive and negative affective ratings, and (2) efficiently
41 automate positive and negative affect intensity coding on large facial expression databases.
42 Further, we show that CVML can be applied to individual human judges to infer which facial
43 actions they use to generate perceptual emotion ratings from facial expressions.

44

45 **Keywords:** emotion, facial expressions, positive and negative affect, computer-vision, machine
46 learning

47

48

49

50

51 **Using Computer-vision and Machine Learning to Automate Facial Coding of Positive and**
52 **Negative Affect Intensity**

53 The ability to effectively communicate emotion is essential for adaptive human function. Of
54 all the ways that we communicate emotion, facial expressions are among the most flexible,
55 allowing us to rapidly convey information to people of different ages, cultures, and languages.
56 Further, facial expressions signal complex action tendencies including threat and cooperative
57 intent (1-3). Unsurprisingly, the ability to produce and recognize facial expressions of emotion is
58 of interest to researchers throughout the social and behavioral sciences.

59 Facial expressions can be interpreted using either message- or sign-based approaches (4).
60 Message-based approaches describe the meaning conveyed by a facial expression (e.g.,
61 happiness), whereas sign-based approaches describe observable facial actions that
62 embody/comprise messages (e.g., cheek raising may indicate happiness). Although message-
63 based approaches are used effectively by psychologists to measure facial expression messages
64 (e.g., happiness), they do not describe facial behavior comprehensively. Instead, they rely on
65 expert judgments of holistic facial expressions—provided by highly trained coders—rather than
66 on facial movements themselves. This renders message-based approaches susceptible to
67 individual differences (unreliability) among human coders, which can impede valid comparisons
68 of results across studies and research sites—even when the same construct is measured.

69 In comparison, multiple comprehensive, standardized sign-based protocols have been
70 developed and used to answer a variety of research questions (4). Among these protocols, the
71 Facial Action Coding System (FACS; 5) may be the most widely used. FACS comprises
72 approximately 33 anatomically-based facial actions (termed action units [AUs]), which interact
73 to generate different facial expressions.

74 Originally developed from a basic emotion theory perspective, the relation between FACS-
75 based AUs and discrete emotions is an active research topic (6). Distinct patterns of AUs reliably
76 map onto each basic emotion category (happiness, sadness, anger, fear, surprise, and disgust),
77 and the existence of distinct patterns of AUs that people use to label different emotional
78 expressions is often used as evidence to support discrete theories of emotion (see 7). For
79 example, oblique lip-corner contraction (AU12), together with cheek raising (AU6) reliably
80 signals enjoyment (8), while brow furrowing (AU4) tends to signal negative emotions like anger
81 and sadness (e.g., 9). Recently, research on how people perceive discrete emotions from AUs has
82 revealed up to 21 discrete categories composed of compound basic emotions (e.g., happily-
83 surprised; 10). Together, these studies suggest that people use the presence of distinct AUs to
84 evaluate emotional content from facial expressions (11), a hypothesis supported by neuroimaging
85 studies showing that differential patterns of BOLD responding in the posterior superior temporal
86 sulcus discriminate between AUs (12).

87 Despite the clear links between AUs and discrete emotion perception, little is known about
88 how AUs map onto dimensional features of emotion (7), especially positive and negative affect
89 (i.e., valence). This is a potentially important oversight given the centrality of valence to
90 dimensional theories of emotion (e.g., 13-15), of which valence is the most consistently
91 replicated dimension (16). Early work using facial electromyography (EMG) showed that
92 zygomatic (AU12) and corrugator (AU4) activity may indicate more positive and more negative
93 subjective intensity, respectively (e.g., 9). However, later studies found that interactions between
94 multiple AUs better describe valence intensity (e.g., 17), and in follow-up work, researchers have
95 proposed that the face may represent positive and negative affect simultaneously with
96 independent sets of AUs (e.g., 18). In one of the few studies directly linking AUs to perceived

97 valence intensity, Messinger et al. (19) found that cheek raising (AU6) was common to
98 perceptual judgments of both intense positive and negative affect, which challenges the idea that
99 a single AU can adequately capture the entire range of valence intensity. Altogether, current
100 evidence suggests that zygomatic (AU12) and corrugator (AU4) activity indicate perceived
101 positive and negative affect, but the extent to which these and other discrete facial actions map
102 onto the entire range of perceived positive or negative affect intensity is unclear.

103 Comprehensive follow-up investigations have been difficult to pursue, in part, because
104 measuring AUs is labor- and time-intensive and requires highly skilled annotators. Indeed, FACS
105 training requires an average of 50-100 hours, and minutes of video can take expert coders
106 multiple hours to rate reliably (20). These characteristics limit sample sizes, reduce feasibility of
107 replication efforts, and discourage researchers from coding facial expressions. Instead,
108 researchers tend to rely on measures of emotional responding that are not observable in social
109 interactions (e.g., heart rate variability). Recently, automated computer-vision and machine
110 learning (CVML) based approaches have emerged that make it possible to scale AU annotation
111 to larger numbers of participants (e.g., 21-23) thus making follow-up studies more feasible. In
112 fact, inter-disciplinary applications of CVML have allowed researchers to automatically identify
113 pain severity (e.g., 24), depressive states (e.g., 25), and discrete emotions from facial expressions
114 (e.g., 26).

115 Work using CVML to detect valence intensity from facial expressions is ongoing (see 27). In
116 fact, there are annual competitions held to develop CVML models that best characterize
117 dimensional features of emotions such as valence and arousal (e.g., 28). Currently, basic
118 emotions can be coded automatically with accuracy comparable to human coders, but valence
119 intensity models show lower concurrent validity. For example, state-of-the-art CVML models

120 show correlations between human- and computer-coded valence ranging from $r = .60-.71$
121 (29,30). While impressive, these CVML models are often constructed to detect valence directly
122 from frame-by-frame video input without intermediately capturing AUs, so it is unclear if
123 successful valence detection depends on prior detection of specific AUs. Furthermore, these
124 results have been collected from relatively small samples or on continuously collected valence
125 ratings (human ratings collected in real-time using dials or joysticks), and it is unclear if these
126 models generalize to other research settings where participants' emotional expressions to
127 evocative stimuli are coded within discrete, trial-by-trial time intervals (e.g., 31). Indeed,
128 contemporary work using CVML has shifted from evaluating facial expressions in controlled
129 laboratory settings toward accurately capturing continuous facial expressions of emotion “in the
130 wild”, which is a much more difficult task (e.g., 29,32). However, to be useful for the majority of
131 applications within the cognitive, decision, and other social and behavioral sciences—where
132 controlled laboratory settings are the norm—it is only necessary that CVML models are
133 optimized for performances within laboratory settings as opposed to in the real-world. Lastly,
134 most valence-detecting CVML models assume a unidimensional valence continuum as opposed
135 to separable continua for positive and negative affect—to our knowledge, there are few
136 opensource datasets used in CVML research that characterize valence as multi-dimensional (see
137 33), and very little work has been done with CVML to separate positive and negative affect (cf.
138 34). Notably, positive and negative affect can vary independently and have different predictive
139 values (10,15,35), suggesting that CVML models designed to account for each dimension
140 separately may be most beneficial for behavioral science applications.

141 Using a well-validated method of emotion induction and both computer-vision measurement
142 of discrete facial actions and continuous measures of positive and negative affect intensity, we

143 (1) identified specific correspondences between perceived emotion intensity and discrete facial
144 AUs, and (2) developed a reliable, valid, and efficient method of automatically measuring the
145 separable dimensions of positive and negative affect intensity. Importantly, data used to train and
146 validate our CVML models were collected from a commonly-used psychological task and
147 contained 4,648 video-recorded, evoked facial expressions from 125 human subjects across
148 multiple task instructions. Our findings shed light on the mechanisms of valence recognition
149 from facial expressions and point the way to novel research applications of large-scale emotional
150 facial expression coding.

151 **Method**

152 *Participants*

153 Video recordings and human coder data were collected as part of a larger study (31). The current
154 study included 125 participants (84 females), ages 18-35 years. All participants gave informed
155 consent prior to the study, and the study protocol (#2011B0071) was approved by The Ohio State
156 Behavioral and Social Sciences Institutional Review Board. Self-reported ethnicities of
157 participants were as follows: Caucasian ($n=96$), East Asian ($n=14$), African-American ($n=5$),
158 Latino ($n=3$), South Asian ($n=3$), and unspecified ($n=4$).

159 *Measures*

160 ***Emotion-evoking task.*** We used an emotion-evoking task, depicted in Figure 1, that has been
161 used in several previous studies to elicit facial expressions of emotion across multiple task
162 instructions (31,36). Participants viewed 42 positive and negative images selected from the
163 International Affective Picture System (IAPS) to balance valence and arousal. Selections were
164 based on previously reported college-student norms (37). Images were presented in 6 blocks of 7
165 trials each, whereby each block consisted of all positive or all negative images. For each block,

166 participants were asked to either *enhance*, *react normally*, or *suppress* their naturally evoked
167 emotional expressions to the images. These instructions effectively increased variability in facial
168 expressions within participants, and they reflect common social situations where emotions must
169 be enhanced or suppressed to achieve desired outcomes. Given known individual differences in
170 suppression and enhancement of facial expressions (31,36), we expected that these task
171 instructions would allow us to create a more generalizable CVML model than with no
172 instructions at all. Block order was randomized across participants. Instructions were given so
173 that each valence was paired once with each condition. All images were presented for 10 s, with
174 4 s between each image presentation. Participants' reactions to each image were video-recorded
175 with a 1080p computer webcam (Logitech HD C270). Due to experimenter error, 1 participant's
176 videos were not recorded correctly, and 7 participants were shown only 41 recordings, resulting
177 in 6,293 usable recordings. Among these, 3 were corrupted and could not be viewed. Thus, 6,290
178 10-s recordings were potentially available.

179

180 **Figure 1. Emotion-evoking task**

181 Participants ($N=125$) viewed a total of 42 images each, divided into 6 blocks of 7 trials. Images
182 were presented for 10 s, with a 4 s inter-trial-interval. Each block of images consisted of either
183 positive or negative image content. In each of the 3 blocks containing positive and negative
184 image content, participants were asked to either *enhance*, *react normally*, or *suppress* their
185 emotional expressions, so that each valence type (i.e., positive or negative) was paired once with
186 each task instruction (enhance, react normally, suppress). All images were selected from the
187 International Affective Picture System (37). Participants' reactions to the images were video
188 recorded and their facial expressions were subsequently rated for positive and negative emotion

189 intensity by a team of trained coders. The same recordings were then analyzed by FACET, a
190 computer vision tool which automatically identifies facial Action Units (AUs).

191

192 ***Manual coding procedure.*** A team of three trained human coders, unaware of participants'
193 task instructions, independently viewed and rated each 10-s recording for both positive and
194 negative emotion intensity. Presentation of recordings was randomized for each coder. Ratings
195 were collected on a 7-point Likert scale ranging from 1 (*no emotion*) to 7 (*extreme emotion*).
196 Coders completed an initial training phase during which they rated recordings of pre-selected
197 non-study cases and discussed specific facial features that influenced their decisions (see the
198 Supplementary Text for the coding guide). The goal of this training was to ensure that all coders
199 could reliably agree on emotion intensity ratings. In addition, coders participated in once-
200 monthly meetings throughout the coding process to ensure reliability and reduce drift.
201 Agreement between coders across all usable recordings (6,290 recordings) was high, with
202 intraclass correlation coefficients (ICCs(3); 38) of .88 and .94 for positive and negative ratings,
203 respectively. The ICC(3) measure reported above indicates absolute agreement of the average
204 human-coder rating within each condition (*enhance, react normally, suppress*) for each of the
205 125 participants. To prepare data for CVML analysis, we performed an additional quality check
206 to screen out videos in which participants' faces were off-camera or covered. Any recording in
207 which a participant's face was covered, obscured, or off-camera for 1 s or more was removed
208 from analysis. If 50% or more of a participant's recordings were excluded, we excluded all of
209 his/her recordings. This resulted in a total of 4,648 usable recordings across 125 participants.
210 With over 4,000 individually-coded recordings, our sample size is in the typical range for
211 machine learning applications (39).

212 ***Automated coding procedure.*** We then analyzed each of the 4,648 recordings with FACET
213 (23). FACET is a computer-vision tool that automatically detects 20 FACS-based AUs (see
214 Supplementary Table 1 for descriptions and depictions of FACET-detected AUs). FACET
215 outputs values for each AU indicating the algorithm's confidence in the AU being present.
216 Confidence values are output at a rate of 30 Hz, resulting in a time-series of confidence values
217 for each AU being present with each frame of a video-recording. Each point in the time-series is
218 a continuous number ranging from about -16 to 16, whereby more positive and more negative
219 numbers indicate increased and decreased probability of the presence of a given AU,
220 respectively. We refer to this sequence of numbers as an AU evidence time-series.

221 Each AU evidence time-series was converted to a point estimate by taking the area under the
222 curve (AUC) of the given time-series and dividing the AUC by the total length of time that a face
223 was detected throughout the clip. This creates a normalized measure that does not render biased
224 weights to clips of varying quality (e.g., clips in which participants' faces are occasionally not
225 detected). Point-estimates computed this way represent the expected probability that a participant
226 expressed a given AU across time. We used the AU evidence time-series point estimates as
227 predictor (independent) variables to train a machine learning model to predict human valence
228 intensity ratings. It took FACET less than 3 days to extract AU evidence time-series data from
229 all recordings (running on a standard 8-core desktop computer). Note that we did not use a
230 baseline correction for each subject, which would require human annotation of a neutral facial
231 expression segment for each participant. Therefore, the models reported here may be applied to
232 novel facial recordings with no human judgment.

233

234 ***Machine Learning Procedure***

235 Figure 2 depicts the machine learning procedure. We trained a random forest (RF) model to
236 predict human-coded valence ratings from the AU evidence time-series point estimates described
237 above (see Supplementary Text for details on training). RFs are constructed by generating
238 multiple decision trees and averaging predictions of all trees together. We chose the RF model
239 because (1) it can automatically capture interactions between independent variables, and we
240 know that humans use multiple AUs simultaneously when evaluating facial expressions; (2) the
241 importance of each independent variable can be extracted from the RF to make inferences
242 regarding which AUs human coders attended to while rating valence intensity (analogous to
243 interpreting *beta* weights from a multiple regression; 39); and (3) RFs have previously shown
244 robust representations of the mapping from facial features (e.g., AUs) to discrete emotions and
245 valence intensity (40,41). Given high agreement among coders and a large literature showing
246 that aggregating continuous ratings from multiple, independent coders leads to reliable estimates
247 despite item-level noise (i.e., ratings for each recording; see 42), we used the average of all
248 coders' ratings for each recording as the outcome (dependent) variable to train the RF.

249

250 **Figure 2. Machine learning procedure**

251 The goal of our first analysis was to determine whether or not CVML could perform similarly to
252 humans in rating facial expressions of emotion. For each AU evidence time-series, we computed
253 the normalized (i.e., divided by the total time that FACET detected a face) Area Under the Curve
254 (AUC), which captures the probability that a given AU is present over time. All AUC values (20
255 total) were entered as predictors into the random forest (RF) model to predict the average coder
256 rating for each recording. To test how similar the model ratings were to human ratings, we
257 separated the data into training (3,060 recordings) and test (1,588 recordings) sets. We fit the RF

258 to the training set and made predictions on the unseen test set. Model performance was assessed
259 by comparing the Pearson and intraclass correlations between computer- and human-generated
260 ratings in the test sets.

261

262 The RF model contains 2 tuning parameters, namely: 1) *ntrees*—the number of decision trees
263 used in the forest, and 2) *mtry*—the number of predictors to sample from at each decision node
264 (i.e., “split”) in a tree. A grid search over $ntrees \in \{100, 200, 300, \dots, 1000\}$ showed that out-
265 of-bag prediction accuracy converged by 500 trees for both positive and negative datasets (not
266 reported). A grid search over $mtry \in \{1, 2, 3, \dots, 20\}$ revealed negligible differences in out-of-
267 bag prediction accuracy for values ranging from 5 to 20. Because RFs do not over-fit the data
268 with an increasing number of trees (39), we set $ntrees = 500$ for models presented in all reported
269 analyses to ensure convergence. Because initial grid searches over *mtry* failed to improve the
270 model, we set *mtry* heuristically (39) as $mtry = p/3$, where *p* represents the number of predictors
271 (i.e., 1 for each AU) in an $n \times p$ matrix ($n =$ number of cases) used to train the model. We fit the
272 RF model using the *easymf* R package (43), which provides a wrapper function for the
273 *randomForest* R package (44). All R codes used for model fitting along with the trained RF
274 models will be made available on our GitHub repository upon publication
275 (<https://github.com/CCS-Lab>).

276 ***Correspondence between human coders and model predictions.*** Model performance refers
277 to how similar the model- and human-generated valence intensity rating are. To assess model
278 performance, we split the 4,648 recordings into training ($n=3,060$; 65.8%) and test ($n=1,588$;
279 34.2%) sets, trained the model on the training set (see the Supplementary Text for details), and
280 then made predictions on the unseen test set to assess how well the RF predicted valence

281 intensity ratings on new data. The data were split randomly with respect to participants so that
282 the training and test data contained 66% and 34% of each participant's recordings, respectively.
283 This separation ensured that training was conducted with all participants, thus creating a more
284 generalizable final model. We fit a separate RF model to positive and negative human ratings. To
285 see if the way we split the training and test data influenced our results, we made 1,000 different
286 training/test-set splits and assessed model performance across all splits (45,46). We used Pearson
287 correlations and ICC coefficients to check model performance on training- and test-sets. Pearson
288 correlations measure the amount of variance in human ratings captured by the model, whereas
289 ICCs measure absolute agreement between human- and model-predicted ratings at the item level
290 (i.e., per recording). Therefore, high correlations and ICCs indicate the model is capturing a large
291 amount of variance in human coder ratings and generating ratings using a similar scale as human
292 coders, respectively. We used McGraw and Wong's ICC(1), as opposed to other ICC methods
293 (38), because we were interested in absolute agreement across all clips, regardless of
294 condition/participant. One-way models were used to compute ICCs in all cases. In general, ICCs
295 between .81 and 1.00 are considered "almost perfect" (i.e., excellent) and ICCs between .61 and
296 .80 are considered "substantial" (i.e., good; 47). We also checked model performance using a
297 different folding scheme for separating training and test sets which ensured that participants'
298 recordings were not shared across splits. This analysis revealed negligible differences in
299 prediction accuracy for positive ratings and a decrease in accuracy for negative ratings, which
300 suggests that more training data may be necessary to capture negative as opposed to positive
301 affect intensity (see Supplementary Text).

302 *Importance of AUs for positive and negative affect.* To identify the specific AUs that human
303 coders were influenced most by when making affective ratings, we fit the RF model to the entire

304 dataset (all 4,648 recordings) without splitting into training and test sets. We used this method to
305 identify independent variables that were robust across all samples (45,46). After fitting the RF
306 models, the importance of each independent variable was estimated using *increase in node*
307 *purity*, a measure of change in residual squared error (increases in prediction accuracy)
308 attributable to the independent variable across all decision trees (39). Importance scores for each
309 independent variable extracted from the RF then represent the magnitude—but not direction of—
310 the effect a given AU has on human coders' valence ratings.

311 To identify potential individual differences in emotion recognition between human coders,
312 we also fit the RF to each coder's ratings independently. We used randomization tests to
313 determine the minimum number of ratings necessary to accurately infer which AUs the coders
314 attended to while generating emotion ratings. For each of the 3 coders, we performed the
315 following steps: (1) randomly sample n recordings rated by coder i , (2) fit the RF model to the
316 subset of n recordings/ratings according to the model fitting procedures outlined above, (3)
317 compute the ICC(2) of the extracted RF feature importances (i.e., *increase in node purity*)
318 between the permuted model and the model fit to all recordings/ratings from coder i , and (4)
319 iterate steps 1-3 twenty times for each value of n (note that different subsets of n
320 recordings/ratings were selected for each of these twenty iterations). We varied $n \in \{30, 40, 50,$
321 $60, 70, 80, 90, 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 1200, 1400, 1600, 1800,$
322 $2000, 2500, 3000\}$.

323 **Results**

324 *Model performance across participants.* Table 1 shows correlations between the model-
325 predicted and the average of the human coders' ratings per recording across both training and
326 test sets. Overall, the RF showed good to excellent performance across both training and test sets

327 for positive and negative ratings. Notably, these results were supported by both the Pearson
328 correlations and the ICCs, suggesting that the RF produced ratings that not only captured
329 variance in, but also showed high agreement with, human ratings. Sensitivity analyses (see
330 Figure 3) indicated that model performance was robust across different training and test splits of
331 the data. These results suggest that variance in human-coded valence intensity can be captured
332 by the presence of discrete AUs.

333

334 **Table 1. Correlations between human- and computer-generated valence ratings**

335

Data Set	Correlation [95% CI]			
	<i>r</i>	ICC(1)		
	(+)	(-)	(+)	(-)
Training	.89 [.88, .90]	.77 [.75, .78]	.88 [.87, .89]	.71 [.69, .72]
Test	.88 [.87, .89]	.74 [.72, .77]	.87 [.86, .88]	.68 [.65, .71]

336

337 *Notes.* (+) = positive valence ratings; (-) = negative valence ratings; *r* = Pearson's correlation;

338 ICC = Intraclass correlation coefficient. Training and test sets contained 3,060 and 1,588

339 recordings, respectively.

340

341 **Figure 3. Sensitivity of model performance to different training/test splits**

342 Results of sensitivity analyses across different splits of the training and test sets. We created

343 1,000 different splits of the training and test sets, fit the RF to each training set, and then made

344 predictions on each respective test set. We stored the Pearson correlations between human- and

345 model-generated ratings for each iteration. Distributions therefore represent uncertainty in

346 prediction accuracy. Means of the distributions (superimposed on respective graphs) are

347 represented by dashed red lines.

348

349 ***Model performance within participants.*** We also checked model performance for each of
350 the 125 participants by computing correlations between human- and model-generated ratings for
351 each participant separately (Figure 4). Although the RF model performed well for many
352 participants in the positive (median $r = .91$, $ICC(1) = .80$) and negative (median $r = .73$, $ICC(1) =$
353 $.51$) affect test sets, 5 participants within the positive and 7 participants within the negative affect
354 test-set yielded negative correlations between human- and computer-generated emotion ratings
355 (Figure 4). Further analyses of within-participant model performance revealed significant
356 positive associations between within-subject variance in model-predicted ratings and within-
357 participant prediction accuracy (all $r_s \geq .54$, $p_s < .001$; see Supplementary Figure 1A). We found
358 the same relation between human-assigned ratings and within-participant variance (see
359 Supplementary Figure 1B). This suggests that the RF model was more accurate in predicting
360 human-rated emotion if participants expressed a wider range of emotional intensity.

361

362 **Figure 4. Model performance within participants**

363 Distributions of within-participant Pearson correlations for positive and negative ratings in the
364 training (all 125 participants) and test (122 participants; correlations could not be computed for 3
365 participants who had 0 variance in human ratings) sets. Red dashed lines represent median
366 within-participant Pearson correlations for each distribution. Intraclass correlations for
367 corresponding figures are reported in text.

368

369 ***Importance of AUs across task instructions.*** To identify which facial expressions human
370 coders may have used to generate positive and negative emotion ratings, we examined the
371 importance of all AUs in predicting human emotion ratings (Figure 5). Note that importance

372 values for the RF do not indicate directional effects, but instead reflect relative importance of a
373 given AU in predicting human-coded positive/negative affect intensity. The RF identified AUs
374 12 (*lip corner pull*), 6 (*cheek raiser*), and 25 (*lips part*) as three of the four most important AUs
375 for predicting positive emotion. Together, AUs 12 and 6 accounted for 50% of the total
376 importance (analogous to proportion of variance accounted for) of all AUs for positive ratings. In
377 contrast to positive ratings, relative importance values for AUs of negative ratings were
378 distributed more evenly across AUs, a trend which was also found when the RF was fit
379 individually to each coder (see *Individual differences in emotion recognition* below). In fact, the
380 5 most important AUs for negative ratings (Figure 5) accounted for 50% of the total importance
381 of all AUs, compared to only AU12 and AU6 accounting for the same proportion of importance
382 in positive ratings. Notably, the importance of AUs for positive and negative emotion ratings
383 were largely independent. In fact, when the ICC(3) is computed by treating positive and negative
384 importance weights for each AU as averaged ratings from two “coders”, the ICC(3) is negative
385 and non-significant ($ICC(3) = -.11, p = .58$), which would only be expected if different facial
386 expressions were important for the coders to rate positive versus negative valence.

387

388 **Figure 5. Relative importance of each AU for positive and negative ratings**

389 Relative importance of each AU for positive and negative human-coder ratings. Relative
390 importance (normalized *increase in node purity* from the RF model) is a measure of change in
391 residual squared error, (increase in prediction accuracy) attributable to each AU being included
392 in the model, and it can be interpreted as how important a given AU is with respect to all other
393 AUs. Note that relative importance is not directional, but does capture interactions among AUs.
394 Visual depictions of the 5 most important AUs for predicting positive and negative ratings are

395 shown on the graphs. Error bars indicate 95% confidence intervals (CIs) estimated from the
396 sensitivity analysis (i.e., 1,000 train/test splits).

397

398 ***Sensitivity of AUs to task instructions.*** To determine if task instructions (*enhance, react*
399 *normally, suppress*) affected model performance or our interpretation of which AUs map onto
400 positive and negative affect, we fit the RF model to all recordings from each condition separately
401 and then compared model performance and AU importance scores across conditions. Table 2
402 shows correlations between human- and computer-generated valence ratings within the different
403 conditions. For positive ratings, correlations were consistently high ($r_s > .80$) across all
404 conditions. In contrast, for negative ratings, correlations were highest in the enhance condition,
405 followed by the react normally and suppress conditions. Of note, all correlations between
406 human- and computer-generated ratings were lower when data were separated by condition
407 compared to when condition was ignored (cf., Table 2 to Table 1). This suggests the lower
408 number of recordings included in the training samples may be partially responsible for lower
409 model performance, but also that CVML performs best when trained on a wider range of
410 emotional intensity.

411

412 **Table 2. Correlations between human- and computer-generated ratings within conditions**

413

Condition	Correlation [95% CI]				Number of recordings	
	<i>r</i> (+)	<i>r</i> (-)	ICC(1) (+)	ICC(1) (-)	Training	Test
Enhance	.81 [.78, .84]	.64 [.59, .68]	.79 [.76, .82]	.61 [.55, .66]	1,047	569
Normal	.81 [.78, .84]	.55 [.49, .61]	.79 [.76, .82]	.49 [.42, .55]	880	516
Suppress	.85 [.83, .87]	.44 [.38, .51]	.83 [.80, .85]	.35 [.28, .42]	1,040	596

414

415 *Notes.* (+) = positive valence ratings; (–) = negative valence ratings; r = Pearson’s correlation;
416 ICC = Intraclass correlation coefficient. All results reported are on test sets.

417

418 Despite only moderate correlations for negative ratings in these conditions, relative
419 importance values for AUs across conditions showed only small differences (Figure 6). In fact,
420 ICCs between AU importance values across conditions were excellent for both positive and
421 negative ratings (Figure 6). This suggests that the task instructions did not strongly influence
422 which AUs were most important to coders.

423

424 **Figure 6. AU relative importance values across task instructions**

425 Relative importance of each AU for positive valence and negative valence human-coder ratings
426 within each of the three task instructions (*enhance*, *react normally*, *suppress*). Intraclass
427 correlation coefficients—both treating importance values as average [ICC(3)] and single [ICC(1)]
428 units—are superimposed. We show ICC(3) here because the AU importance scores could be
429 interpreted as “averages” across all recordings.

430

431 *Individual differences in emotion recognition.* All three coders showed similarly-ordered
432 importance profiles, indicating that they attended to similar AUs while generating emotion
433 ratings (Figure 7). Agreement between all three individual coders’ importance profiles supported
434 this claim—ICC(3)s were high for both positive (ICC(3) = 0.95) and negative (ICC(3) = 0.93)
435 importance profiles. Figure 8 shows the results of the randomization test. For positive ratings,
436 ICC(2)s for all 3 coders reached 0.75 (regarded as “excellent” agreement; see 47) after just 60
437 recordings/ratings. For negative ratings, ICC(2)s for all 3 coders reached 0.75 after 200

438 recordings/ratings. Because the recordings in our task were 10 s long and coders rated
439 positive/negative emotion intensity after each recording, the task used in the current study could
440 be condensed to 200 recordings (~33 minutes) and still reveal individual differences in AU
441 importances between coders with good accuracy. Future studies may be able to shorten the task
442 even further by testing shorter video recordings (i.e., less than 10 s per recording).

443

444 **Figure 7. Individual differences in AU importances between coders**

445 Feature importances (not normalized to show relative differences) extracted from the RF model
446 fit separately to each coder. Coders all show similarly ordered importance profiles, suggesting
447 that they attended to similar facial expressions while generating emotion ratings. Note that
448 positive importance estimates are distributed across few predictors (i.e., AUs 6, 12, and 18),
449 whereas negative importance estimates are more spread out throughout all predictors. Agreement
450 between all three individual coders' importance profiles was high, with ICC(3)s of .95 and .93
451 and ICC(1)s of .86 and .83 for positive and negative ratings, respectively.

452

453 **Figure 8. Number of recordings necessary to accurately estimate AU importance**

454 Grid searches over the number of recordings/ratings necessary to achieve reliable estimates of
455 AU importances for each valence-coder pair (coders appear in the same order as in Figure 7).
456 Reliability is indexed by the ICC(2) between AU importance profiles (i.e. *increase in node*
457 *purity*) extracted from the model fit to all the recordings that coders rated versus the model fit to
458 subsets of recordings that they rated. Note that the ICC(2) assumes that importance estimates are
459 “average” units (similar to ICC(3)s in Figures 6 & 7). The RF model was fit to each sample of
460 size n along the x -axis, AU importance profiles were extracted from the model, and ICC(2)s

461 were then calculated between the given sample and full-data AU importance profile scores. We
462 iterated this procedure 20 times within each different sample size to estimate the variation in
463 estimates across recordings. Shading reflects the 2 standard errors from the mean ICC within
464 each sample across all 20 iterations. The red-dashed line indicates an ICC(2) of .75, which is
465 considered “excellent”. For positive ratings, the ICC(2) reached .75 after only 60
466 recordings/ratings for each coder. For negative ratings, coders 1 and 3 reached an ICC(2) of .75
467 by 120 recordings/ratings, whereas coder 3 reached an ICC(2) of .75 by 200 recordings/ratings.

468

469 **Discussion**

470 Our study offers strong evidence that people use discrete AUs to make judgments regarding
471 positive and negative affect intensity from facial expressions, indicating that patterns of discrete
472 AUs reliably represent dimensions of facial expressions of emotion (analogous to how specific
473 patterns of AUs map to the basic emotions). Our CVML analysis identified AU12, AU6, and
474 AU25 as especially important features for positive affect intensity ratings. Together, these AUs
475 represent the core components of a genuine smile (48). Note that AU12 and AU6 interact to
476 signify a *Duchenne smile*, which can indicate genuine happiness (8), and previous research
477 demonstrates that accurate observer-coded enjoyment ratings rely on AU6 (49). Additionally, the
478 five most important AUs we identified for negative affect intensity map on to those found in
479 negative, discrete emotions such as anger (AUs 4 and 5), disgust (AU9), sadness (AU4), and fear
480 (AUs 2, 4, and 5). Together, the AUs that we identified for positive and negative affect are
481 consistent with prior studies suggesting that positive and negative facial expressions occupy
482 separate dimensions (15,50). Notably, the AUs accounting for the majority of the variance in
483 positive affect had no overlap with those for negative affect, evidenced by near-zero ICCs,

484 indicating that our human coders used distinct patterns of facial expressions to evaluate positive
485 versus negative intensity ratings. The existence of distinct patterns of AUs which represent
486 positive and negative affect intensity explains paradoxical findings that facial expressions can be
487 simultaneously evaluated as both positive and negative (e.g., happily-disgusted; 10). Further, our
488 results suggest that the use of CVML to identify individual differences in positive and negative
489 affect recognition from dynamic facial expressions is a potentially important avenue for future
490 research. While the current study only observed differences between three trained coders (see
491 Figure 8), future studies may collect emotion ratings from naïve groups of participants and
492 perform similar analyses.

493 Our results also provide support for the use of CVML as a valid, efficient alternative to
494 human coders, and with further validation we expect CVML to expand the possibilities of future
495 facial expression research in the social and behavioral sciences. For example, adoption of
496 automatic facial coding tools will allow researchers to more easily incorporate facial expressions
497 into models of human decision making. Decades of research show clear links between facial
498 expressions of emotion and cognitive processes in aggregate (see 51,52), yet the dynamics
499 between cognitive mechanisms and facial expressions are poorly understood in part due to
500 difficulties accompanying manual coding. In fact, we are currently using computational
501 modeling to explore cognition-expression relationships with the aid of CVML (53), which would
502 be infeasible with manual coding of facial expressions. For example, in the current study it took
503 less than three days to automatically extract AUs from 4,648 video recordings and train ML
504 models to generate valence intensity ratings (using a standard desktop computer). In stark
505 contrast, it took six months for three human coders to be trained and then code *affect intensity*

506 across our 125 subjects—FACS coding would have taken much longer, rendering the scale of
507 this project infeasible.

508 Models used in this study predicted positive emotion intensity with greater accuracy than
509 negative emotion intensity, which may be due to the number of discrete facial actions associated
510 with negative compared to positive emotional expressions. To support this claim, we found that
511 importance scores for negative, but not positive, emotion ratings were spread across many
512 different AUs and showed more variation across task instructions (Figures 5 and 6). This
513 suggests that a wider range of facial expressions were used by coders when generating negative
514 rather than positive emotion ratings. Future studies might address this with CVML models that
515 can detect more than the 20 AUs used here. Additionally, our results suggest that negative affect
516 intensity requires more training data for CVML than positive affect, as evidenced by large
517 discrepancies in model performance between our CVML model that ignored the task instructions
518 compared to those that we fit to data from each task instruction separately. Future studies might
519 address this by devoting more time to collecting and coding negative, rather than positive,
520 affective facial expressions.

521 Our interpretation of the computer-vision coded AUs in this study is potentially limited
522 because we did not compare reliability of AU detection between FACET and human FACS
523 experts. Additionally, FACET only detects 20 of the approximately 33 AUs described by FACS,
524 so it is possible that there were other important AUs to which the human coders attended when
525 generating valence ratings that we were unable to capture. However, our models showed
526 excellent prediction accuracy on new data (i.e., capturing ~80% of the variance in human ratings
527 of positive affect intensity), and we identified theoretically meaningful patterns of AUs for
528 positive and negative emotion intensity that are consistent with prior studies (e.g., components of

529 the *Duchenne smile*). It is unlikely that we would achieve these results if FACET did not reliably
530 detect similar, important AUs which represented the intensity of positive and negative facial
531 expressions produced by our 125 participants. Finally, as computer vision advances, we expect
532 that more AUs will be easier to detect. CVML provides a scalable method that can be re-applied
533 to previously collected facial expression recordings as technology progresses.

534 Although this study investigated positive and negative affect, our method could easily be
535 extended to identify facial actions that are associated with other emotional constructs (e.g.,
536 arousal). The ability to identify specific AUs responsible for facial expression recognition has
537 implications for various areas within the social and behavioral sciences. Opportunities may be
538 particularly pronounced for psychopathology research, where deficits and/or biases in
539 recognizing facial expressions of emotion are associated with a number of psychiatric disorders,
540 including autism, alcoholism, and depression (54-56). CVML provides a framework through
541 which both normal and abnormal emotion recognition can be studied efficiently and
542 mechanistically, which could lead to rapid and cost-efficient markers of emotion recognition in
543 psychopathology (57).

544

545

546

547

548

549

550

551

552

553

554

555 **Author Contributions**

556 W.-Y. Ahn and N. Haines developed the study concept. Data were collected by M. Southward
557 and J. Cheavens. N. Haines and W.-Y. Ahn conducted all statistical and machine learning
558 analyses. All authors interpreted the results. N. Haines drafted the manuscript, and all authors
559 contributed to critical manuscript revisions. All authors approved the final version of the
560 manuscript for submission.

561

562 **Acknowledgements**

563 We thank S. Bowman-Gibson for aiding in the manual quality check for all recordings, and J.
564 Haaser, J. Borden, S. Choudhury, S. Okey, T. St. John, M. Stone, and S. Tolliver for manually
565 coding videos. We also thank J. Cohn, J. Myung, A. Rogers, and H. Hahn for their comments
566 and suggestions on previous drafts of the manuscript.

567

568

569

570

571

572

573

574

575

576

577

578

References

- 579 1. Krumhuber E, Manstead ASR, Cosker D, Marshall D, Rosin PL, Kappas A. Facial
580 dynamics as indicators of trustworthiness and cooperative behavior. *Emotion*.
581 2007;7(4):730-735.
- 582 2. Reed L, DeScioli P, Pinker S. The Commitment Function of Angry Facial Expressions.
583 *Psychological Science*. 2014;25(8):1511-1517.
- 584 3. Reed L, Zeglen K, Schmidt K. Facial expressions as honest signals of cooperative intent
585 in a one-shot anonymous Prisoner's Dilemma game. *Evol Hum Behav*. 2012;33(3):200-
586 209.
- 587 4. Cohn JF, Ekman P. Measuring facial action. In: Harrigan JA, Rosenthal R, Scherer KR,
588 editors. *The new handbook of nonverbal behavior for research methods in the affective*
589 *sciences*. New York: Oxford University Press; 2005. p. 9-64.
- 590 5. Ekman P, Friesen W, Hager JC. Facial action coding system: The manual on CD ROM.
591 [CD-ROM]. Salt Lake City; 2002.
- 592 6. Ekman P, Rosenberg EL. *What the face reveals: basic and applied studies of spontaneous*
593 *expression using the facial action coding system (FACS)*. Oxford: Oxford University
594 Press; 2005.
- 595 7. Keltner D, Ekman P. Facial expression of emotion. In: Lewis M, Haviland-Jones JM,
596 editors. *Handbook of emotions*. 2nd ed. New York: Guilford Press; 2000. P. 236-249.
- 597 8. Ekman P, Davidson R, Friesen W. The Duchenne smile: Emotional expression and brain
598 physiology: II. *Journal of Personality and Social Psychology*. 1990;58(2):342-353.
- 599 9. Brown SL, Schwartz GE. Relationships between facial electromyography and subjective
600 experience during affective imagery. *Biol Psychol*. 1980;11(1):49-62.

- 601 10. Du S, Tao Y, Martinez AM. Compound facial expressions of emotion. Proc Natl Acad
602 Sci USA. 2014;111(15):E1454-E1462.
- 603 11. Martinez AM. Visual Perception of facial expressions of emotion. Curr Opin Psychol.
604 2017;17:27-33.
- 605 12. Srinivasan R, Golomb JD, Martinez AM. A neural basis of facial action recognition in
606 humans. J Neurosci. 2017;36(16): 4434-4442.
- 607 13. Russell JA. A circumplex model of affect. J Pers Soc Psychol. 1980;39(6):1161-1178.
- 608 14. Schlosberg H. Three dimensions of emotion. Psychol Rev. 1954;61(2):81-88.
- 609 15. Watson D, Tellegen A. Toward a consensual structure of mood. Psychol Bull. 1985;
610 98(2):2918-235.
- 611 16. Smith CA, Ellsworth PC. Patterns of cognitive appraisal in emotion. J Pers Soc Psychol.
612 1985;48(4):813-838.
- 613 17. Cacioppo JT, Petty RE, Losch ME, Kim HS. Electromyographic activity over facial
614 muscle regions can differentiate the valence and intensity of affective reactions. J Pers
615 Soc Psychol. 1986;50(2):260-268.
- 616 18. Larsen JT, Norris CJ, Cacioppo JT. Effects of positive and negative affect on
617 electromyographic activity over zygomaticus major and corrugator supercili.
618 Psychophysiology. 2003;40:776-785.
- 619 19. Messinger DS, Mattson WI, Mahoor MH, Colm JF. The eyes have it: making positive
620 expressions more positive and negative expressions more negative. Emotion.
621 2012;12(3):430-436.
- 622 20. Bartlett MS, Hager JC, Ekman P, Sejnowski TJ. Measuring facial expressions by
623 computer image analysis. Psychophysiology. 2003;36(2):253-263.

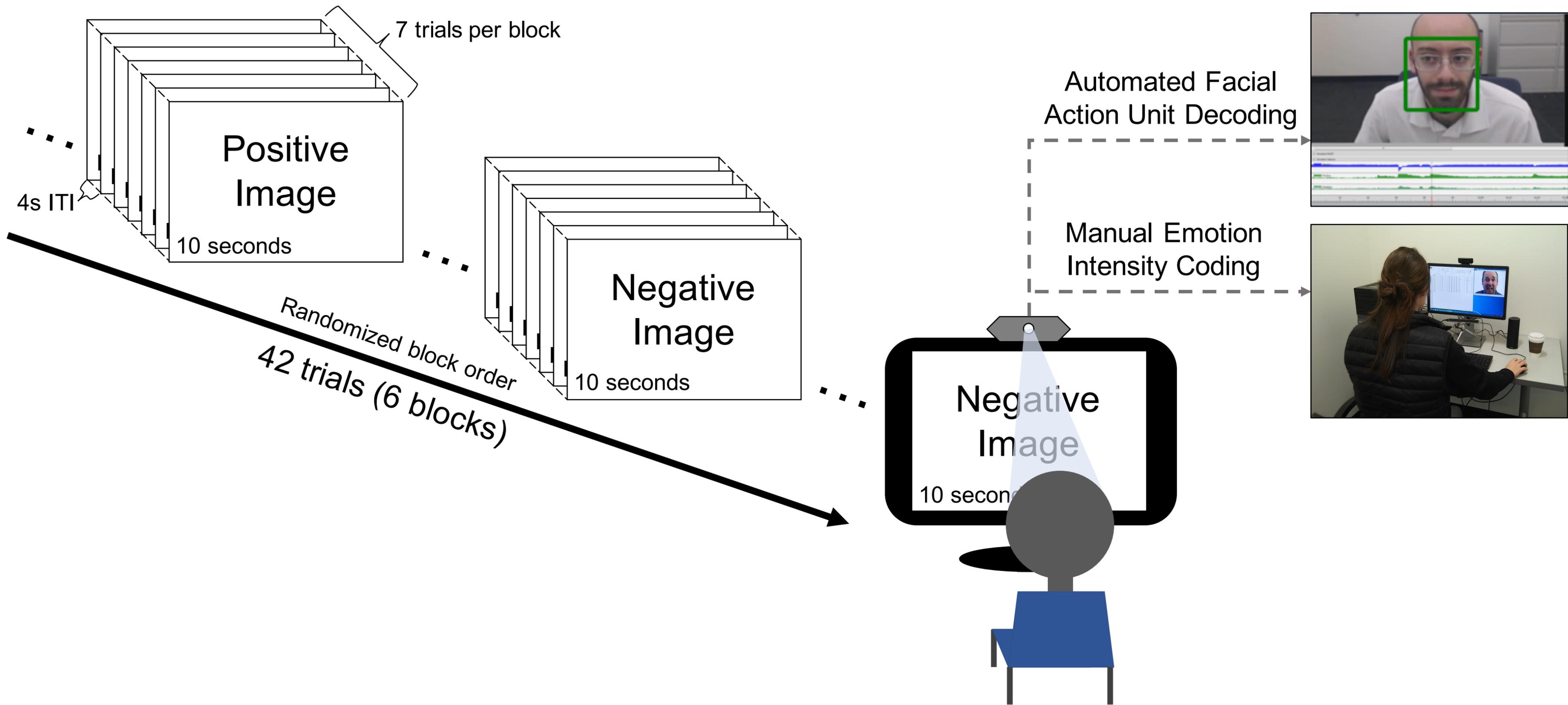
- 624 21. Baltrusaitis T, Robinson P, Morency LP. Openface: an open source facial behavior
625 analysis toolkit. In: 2016 IEEE Winter conference on Applications of Computer Vision;
626 2016 March 7-9; Lake Placid, NY.
- 627 22. Lewinski P, den Uyl TM, Butler C. Automated facial coding: Validation of basic
628 emotions and FACS AUs in FaceReader. *J Neurosci Psychol Econ*. 2014; 7(4):227-236.
- 629 23. Stöckli S, Schulte-Mecklenbeck M, Borer S, Samson AC. Facial expression analysis with
630 AFFDEX and FACET: a validation study. *Behav Res Methods*. 2017;26(5):1-15.
- 631 24. Sikka K, Ahmed AA, Diaz D, Goodwin MS, Craig KD, Bartlett MS, et al. Automated
632 assessment of children's postoperative pain using computer vision. *Pediatrics*.
633 2015;136(1):e124-e131.
- 634 25. Dibeklioglu , Hammal Z, Yang Y, Cohn JF. Multimodal detection of depression in
635 clinical interviews. In: 2015 ACM on International Conference on Multimodal
636 Interaction - ICMI '15; 2015 Nov 9-13; New York, New York, USA. ACM Press; 2015.
- 637 26. Kotsia I, Pitas I. Facial expression recognition in image sequences using geometric
638 deformation features and support vector machines. *IEEE Trans Image Process*.
639 2007;16(1):172-187.
- 640 27. Gunes H, Pantic M. Automatic, dimensional, and continuous emotion recognition.
641 *International Journal of Synthetic Emotions*. 2010; 1(1):68-99.
- 642 28. Ringeval F, Schuller B, Valstar M, Jaiswal S, Marchi E, Lalanne D, et al. AV+EC 2015:
643 The First Affect Recognition Challenge Bridging Across Audio, Video, and
644 Physiological Data. In: 5th International Workshop on Audio/Visual Emotion Challenge;
645 2015 Oct 26-30; Brisbane, Australia.

- 646 29. Mollahosseini A, Hasani B, Mahoor MH. AffectNet: A Database for Facial Expression,
647 Valence, and Arousal Computing in the Wild. *IEEE Trans Affect Comput.* 2017.
- 648 30. Nicolaou MA, Gunes H, Pantic M. (2011). Continuous prediction of spontaneous affect
649 from multiple cues and modalities in valence-arousal space. *IEEE Trans Affect Comput.*
650 2011;2:92-105.
- 651 31. Southward MW, Cheavens JS. (2017). Assessing the relation between flexibility in
652 emotional expression and symptoms of anxiety and depression: The roles of context
653 sensitivity and feedback sensitivity. *J Soc Clin Psychol.* Feb 2017;36(2):142-157.
- 654 32. Kossaifi J, Tzimiropoulos G, Todorovic S, Pantic M. (2017). AFEW-VA database for
655 valence and arousal estimation in-the-wild. *Image Vis Comput.* Sept 2017;65:23-26.
- 656 33. Haamer E, Rusadze E, Lüsi I, Ahmed T, Sergio, Escalera, et al. Review on Emotion
657 Recognition Databases | IntechOpen [Internet]. Intech open. IntechOpen; 2018 [cited
658 2018Sep23]. Available from: [https://www.intechopen.com/books/human-robot-](https://www.intechopen.com/books/human-robot-interaction-theory-and-application/review-on-emotion-recognition-databases)
659 [interaction-theory-and-application/review-on-emotion-recognition-databases](https://www.intechopen.com/books/human-robot-interaction-theory-and-application/review-on-emotion-recognition-databases)
- 660 34. Bailenson JN, Pontikakis ED, Mauss IB, Gross JJ, Jabon ME, Hutcherson CAC, et al.
661 Real-time classification of evoked emotions using facial feature tracking and
662 physiological responses. *Int J Hum Comput Stud.* 2008;66(5):303-317.
- 663 35. Watson D, Clark LA, Tellegen A. Development and validation of brief measures of
664 positive and negative affect: The PANAS scales. *J Pers Soc Psychol.* 1998;54(6):1063-
665 1070.
- 666 36. Bonanno GA, Papa A, Lalande K, Westphal M, Coifman K. The importance of being
667 flexible: The ability to both enhance and suppress emotional expression predicts long-
668 term adjustment. *Psychol Sci.* 2004;15(7):482-487.

- 669 37. Lang PJ, Bradley MM, Cuthbert BN. International Affective Picture System (IAPS):
670 Technical manual and affective ratings (Technical Report A-4). Gainesville, FL; 1995.
- 671 38. McGraw KO, Wong SP. (1996). Forming inferences about some intraclass correlation
672 coefficients. *Psychol Methods*. 1996;1(4):30-46.
- 673 39. Hastie T, Tibshirani R, Friedman J. The elements of statistical learning. New York(NY):
674 Springer New York; 2009.
- 675 40. Amirian M, Kächele M, Thiam P, Kessler V, Schwenker F. Continuous Multimodal
676 Human Affect Estimation using Echo State Networks: Proceedings of the 6th
677 International Workshop on Audio/Visual Emotion Challenge; 2016 Oct 16; Amsterdam,
678 Netherlands. Proceedings of the 6th International Workshop on Audio/Visual Emotion
679 Challenge; 2016. p. 67-74.
- 680 41. Pu X, Fan K, Chen X, Ji L, Zhou Z. Facial expression recognition from image sequences
681 using twofold random forest classifier. *Neurocomputing*. 2015;168:1173-1180.
- 682 42. Rosenthal R. Conducting judgment studies: some methodological issues. In: Harrigan JA,
683 Rosenthal R, Scherer KR, editors. Series in Affective Science. The new handbook of
684 methods in nonverbal behavior research. New York: Oxford University Press; 2005. p.
685 199-234.
- 686 43. Ahn WY, Hendricks P, Haines N. Easyml: Easily Build and Evaluate Machine Learning
687 Models. *bioRxiv*. 2017.
- 688 44. Liaw A, Wiener M. *R News*. 2002Dec;
- 689 45. Ahn, WY, Ramesh D, Moeller FG, Vassileva J. Utility of machine-learning approaches
690 to identify behavioral markers for substance use disorders: impulsivity dimensions as
691 predictors of current cocaine dependence. *Front in Psychiatry*. 2016; 7:290.

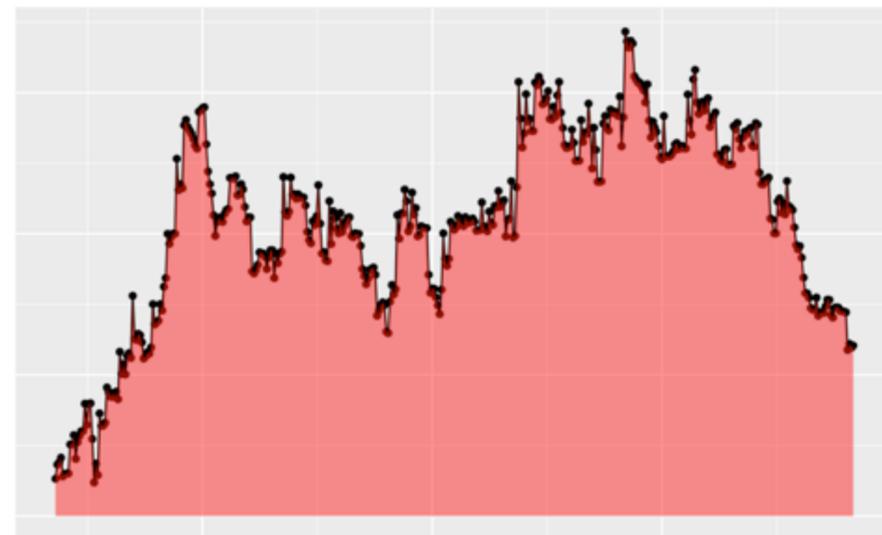
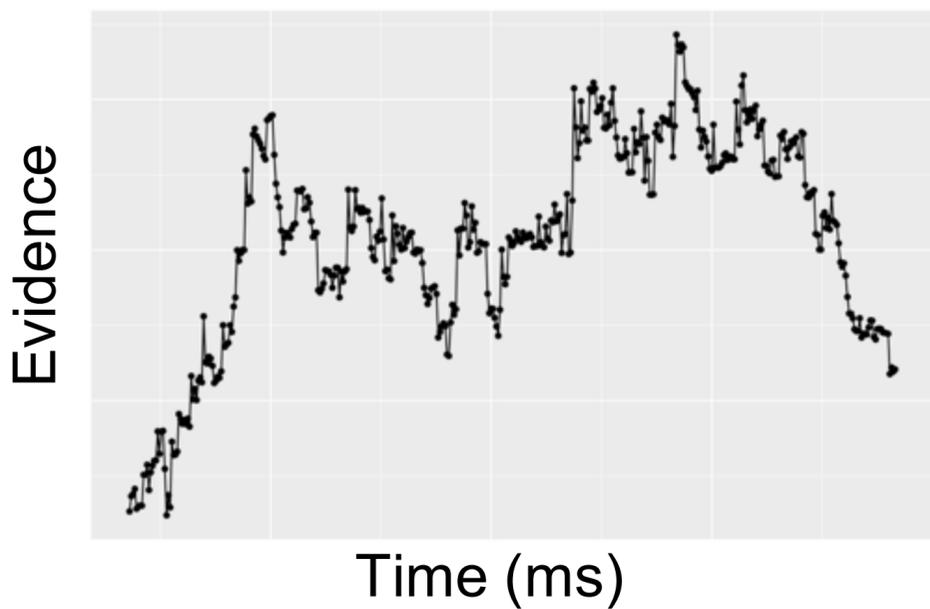
- 692 46. Ahn WY, Vassileva J. Machine-learning identifies substance-specific behavioral markers
693 for opiate and stimulant dependence. *Drug Alcohol Depend.* 2016;161:247-257.
- 694 47. Landis JR, Koch GG. The measurement of observer agreement for categorical data.
695 *Biometrics.* 1997;33(1):159-174.
- 696 48. Korb S, With S, Niedenthal P, Kaiser S, Grandjean D. The perception and mimicry of
697 facial movements predict judgments of smile authenticity. *PLoS ONE.* 2014;
698 9(6):e99194.
- 699 49. Frank MG, Ekman P, Friesen WV. Behavioral markers and recognizability of the smile
700 of enjoyment. *J Pers Soc Psychol.* 1993; 64(1):83-93.
- 701 50. Belsky J, Hsieh KH, Crnic K. Infant positive and negative emotionality: One dimension
702 or two? *Dev Psychol.* 1996;32(2):289-298.
- 703 51. Erickson K, Schulkin J. Facial expressions of emotion: A cognitive neuroscience
704 perspective. *Brain Cog.* 2003;52(1):52-60.
- 705 52. Izard CE. Basic emotions, relations among emotions, and emotion-cognition relations.
706 *Psychol Rev.* 1992; 99(3):561-565.
- 707 53. Haines N, Rass O, Shin YW, Busemeyer JR, Brown JW, O'Donnell B, et al. (in
708 preparation). Regret induces rapid learning from experience-based decisions: A model-
709 based facial expression analysis approach.
- 710 54. Celani G, Battacchi MW, Arcidiacono L. The understanding of the emotional meaning of
711 facial expressions in people with autism. *J Autism Dev Disord.* 1999;29(1):57-66.
- 712 55. Philippot P, Kornreich C, Blairy S, Baert I, Dulk AD, Bon OL, et al. (1999). Alcoholics'
713 deficits in the decoding of emotional facial expression. *Alcohol Clin Exp Res.*
714 1999;23(6):1031-1038.

- 715 56. Rubinow D R, Post RM. Impaired recognition of affect in facial expression in depressed
716 patients. *Biol Psychiatry*. 1992;31(9): 947-953.
- 717 57. Ahn, W-Y, Busemeyer, JR. Challenges and promises for translating computational tools
718 into clinical practice. *Curr Opin Behav Sci*. 2016; 1:1-7.



AU Evidence Time-Series

Area Under the Curve



0.94

For each AU
20 predictors

Time-normalized Area Under the Curve (AUC) score is computed for each AU

Predictors = AUC for each AU evidence time-series
Outcome = Average of human coder ratings per recording

All data
(100%)

Divided into Training/Test sets

Training set
(66%)

Test set
(34%)

Train the random forest model

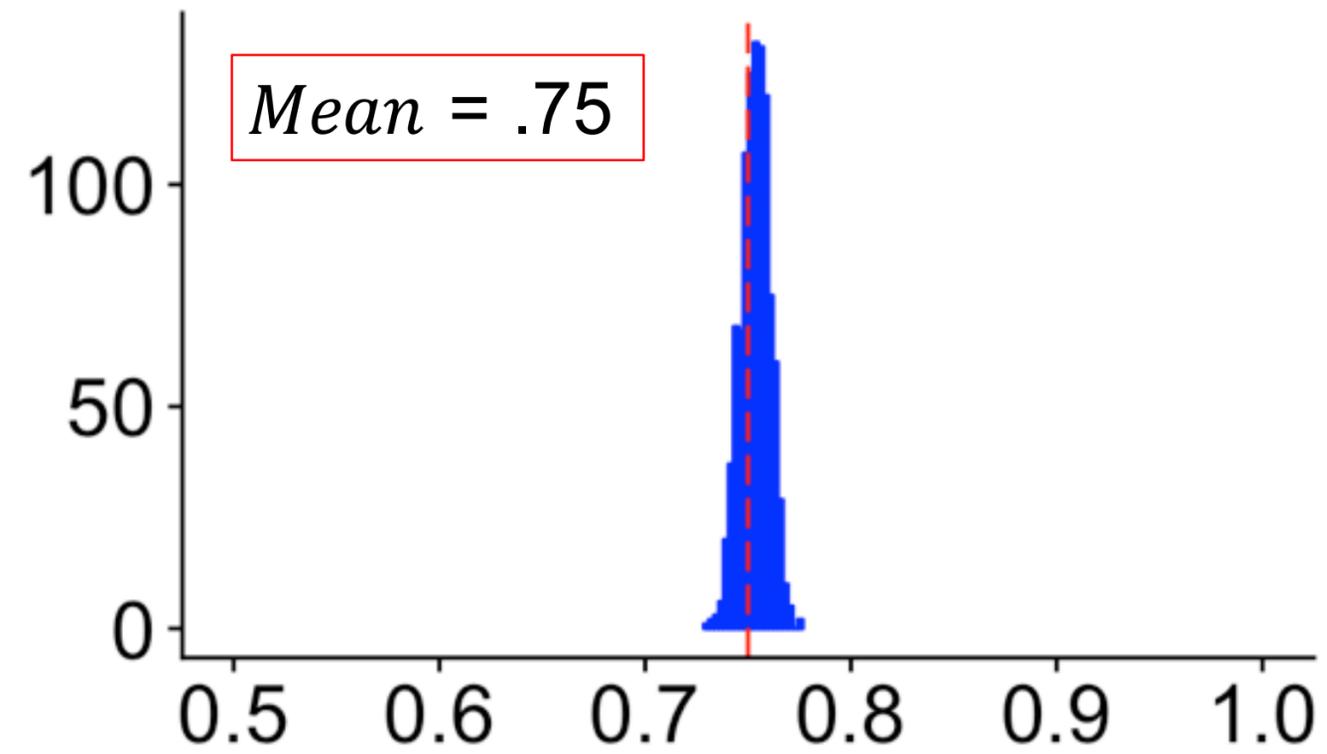
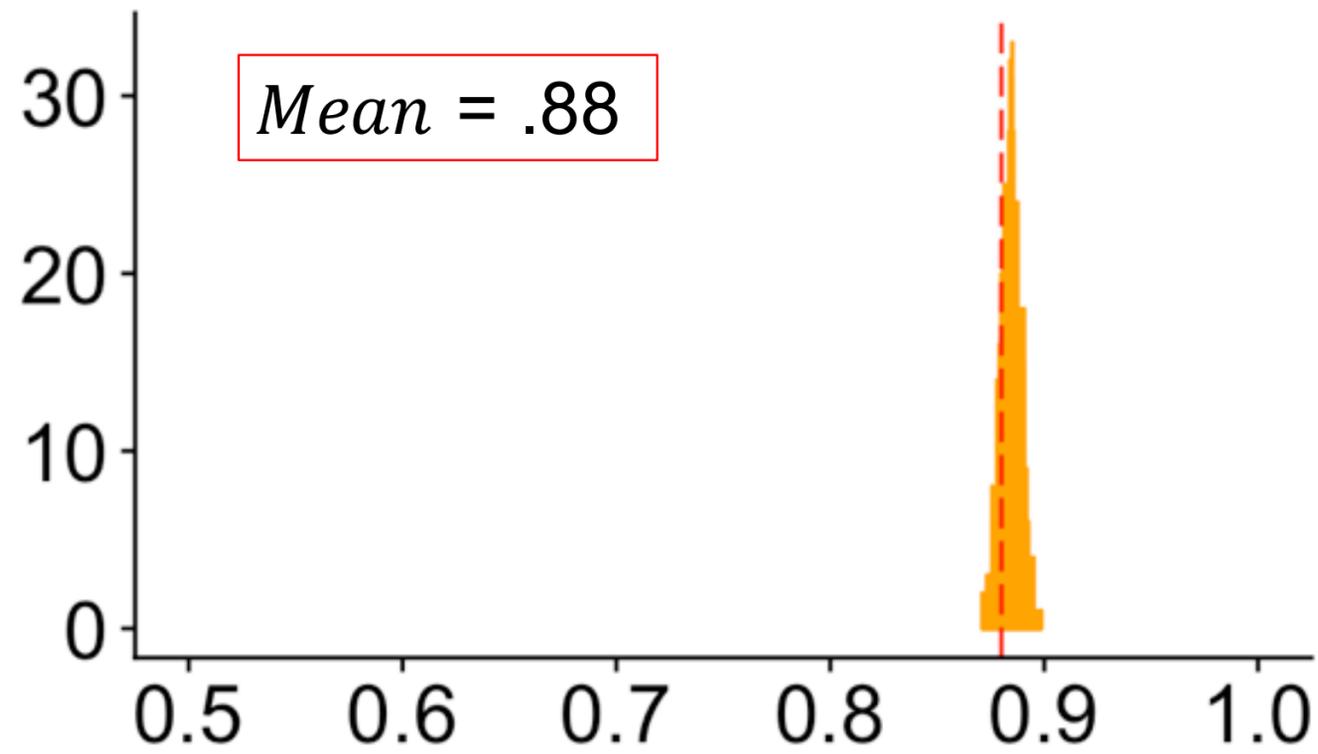
Out-of-sample predictions

Positive

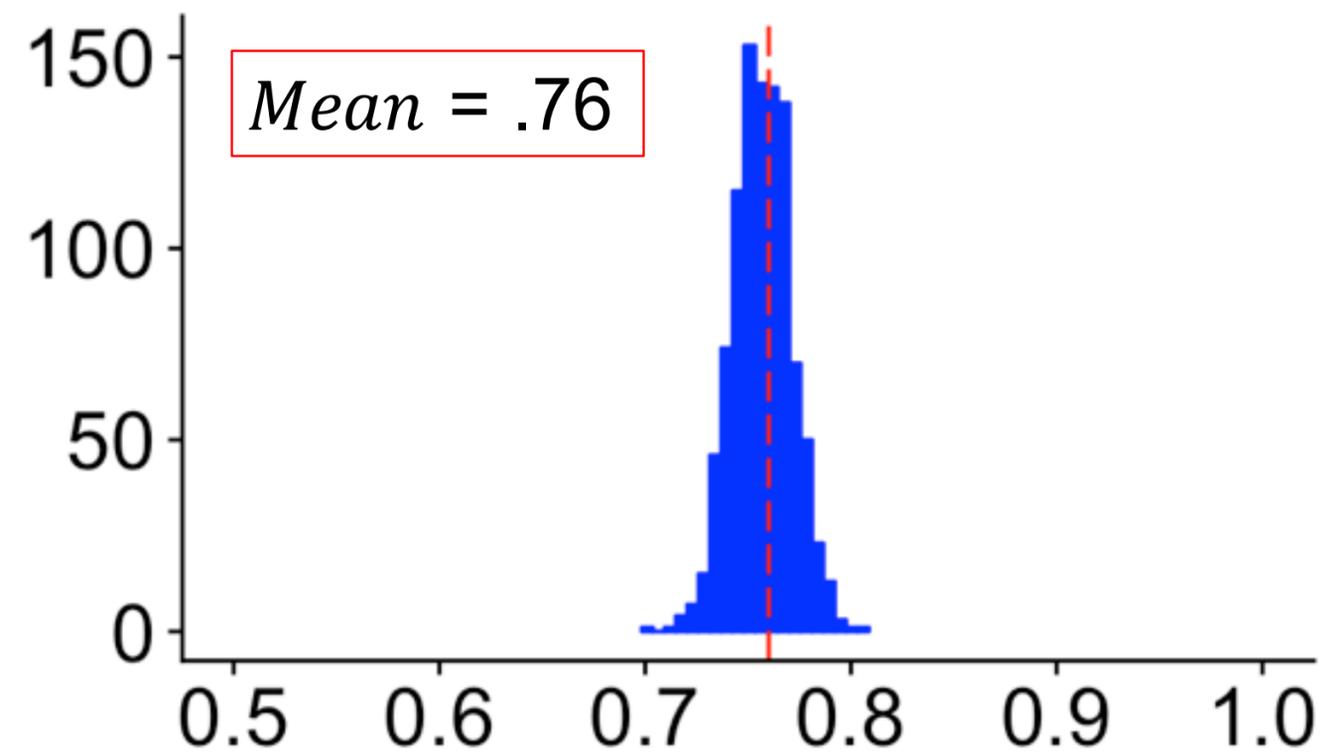
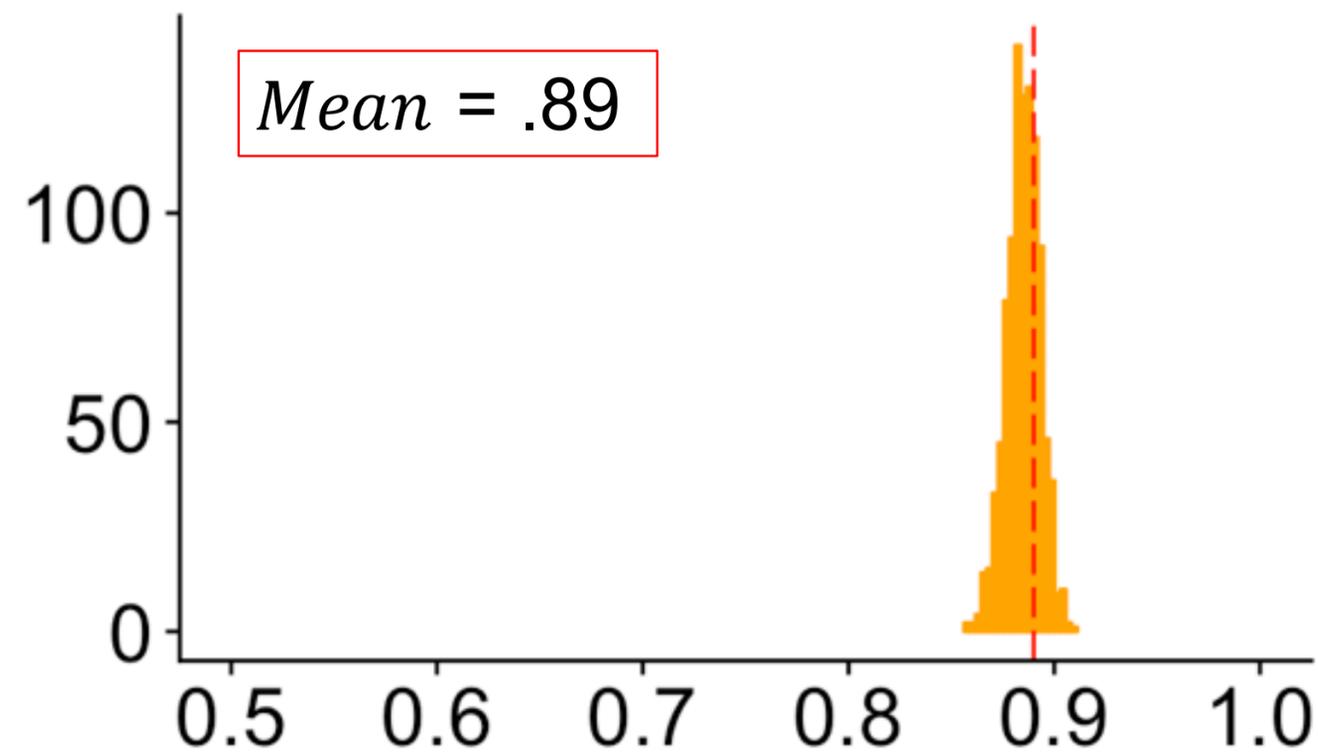
Negative

Frequency

Training



Test



Across-participant Correlation

Positive

Negative

Frequency

Training

Test

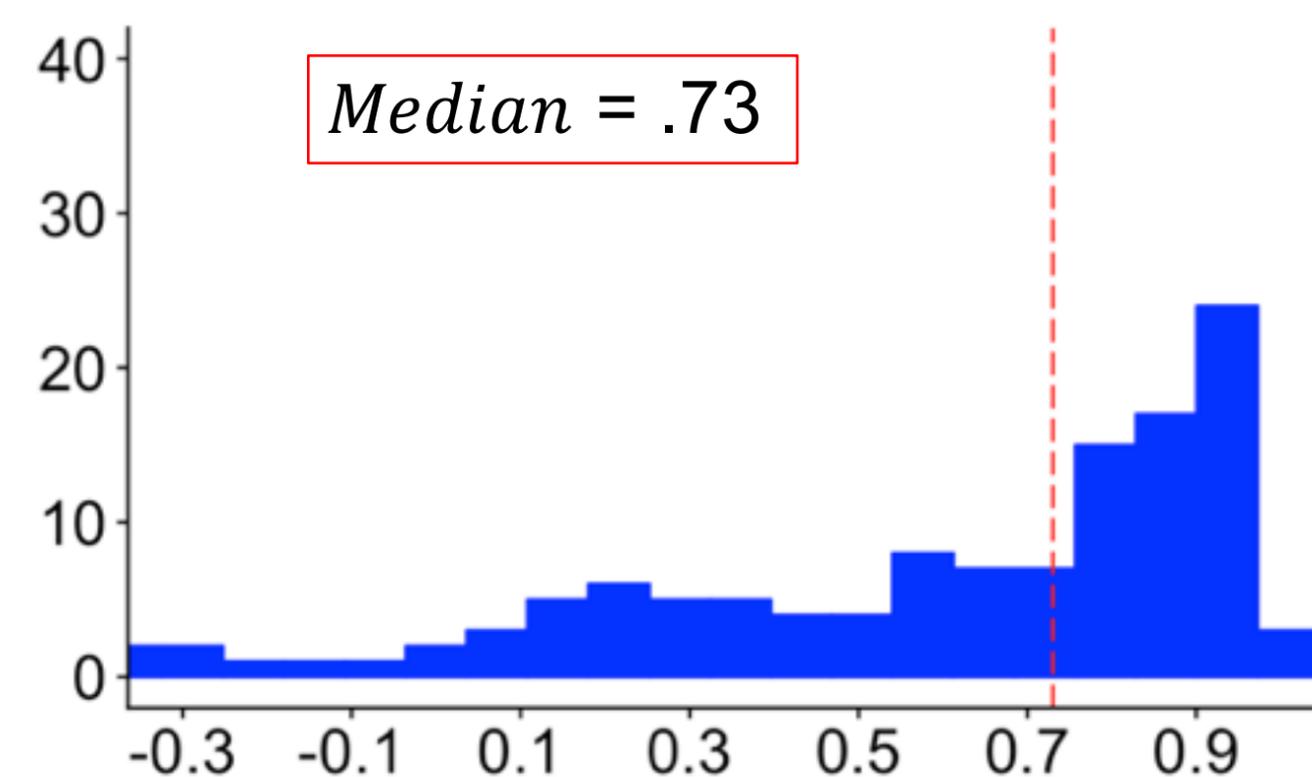
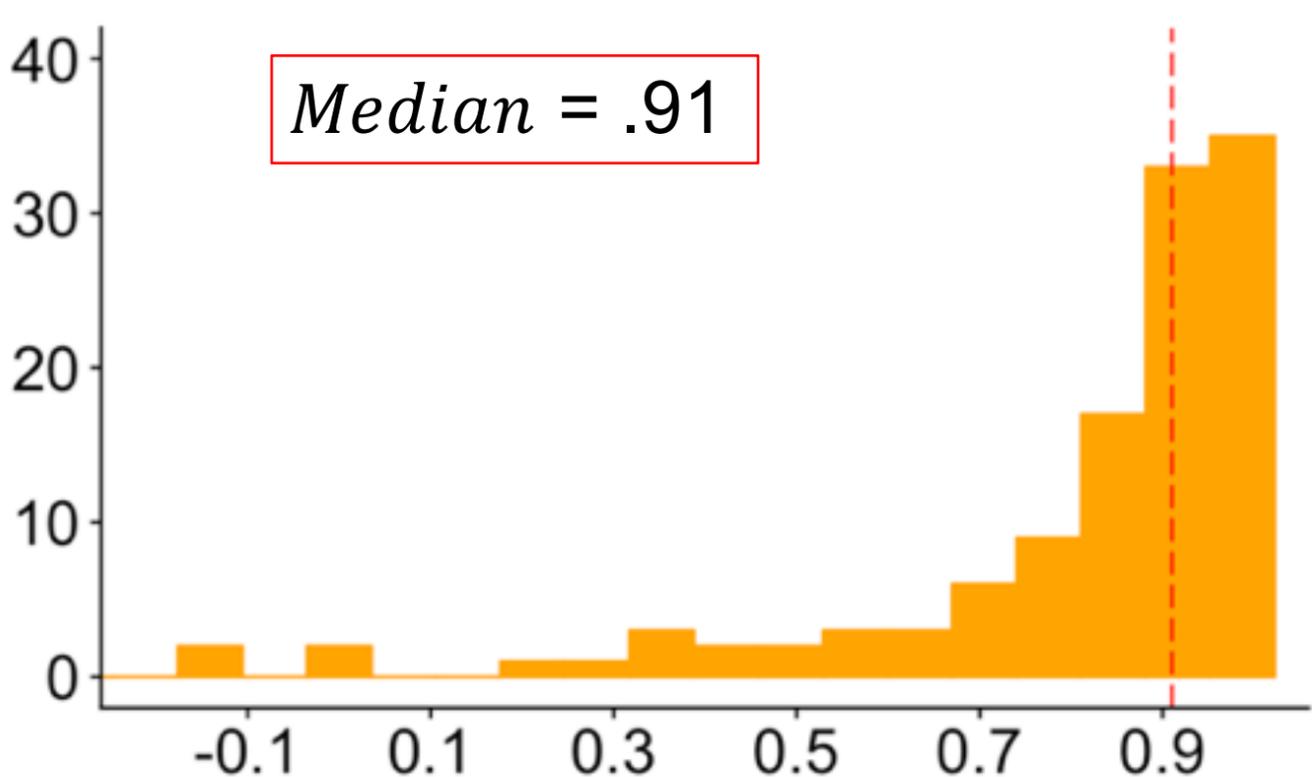
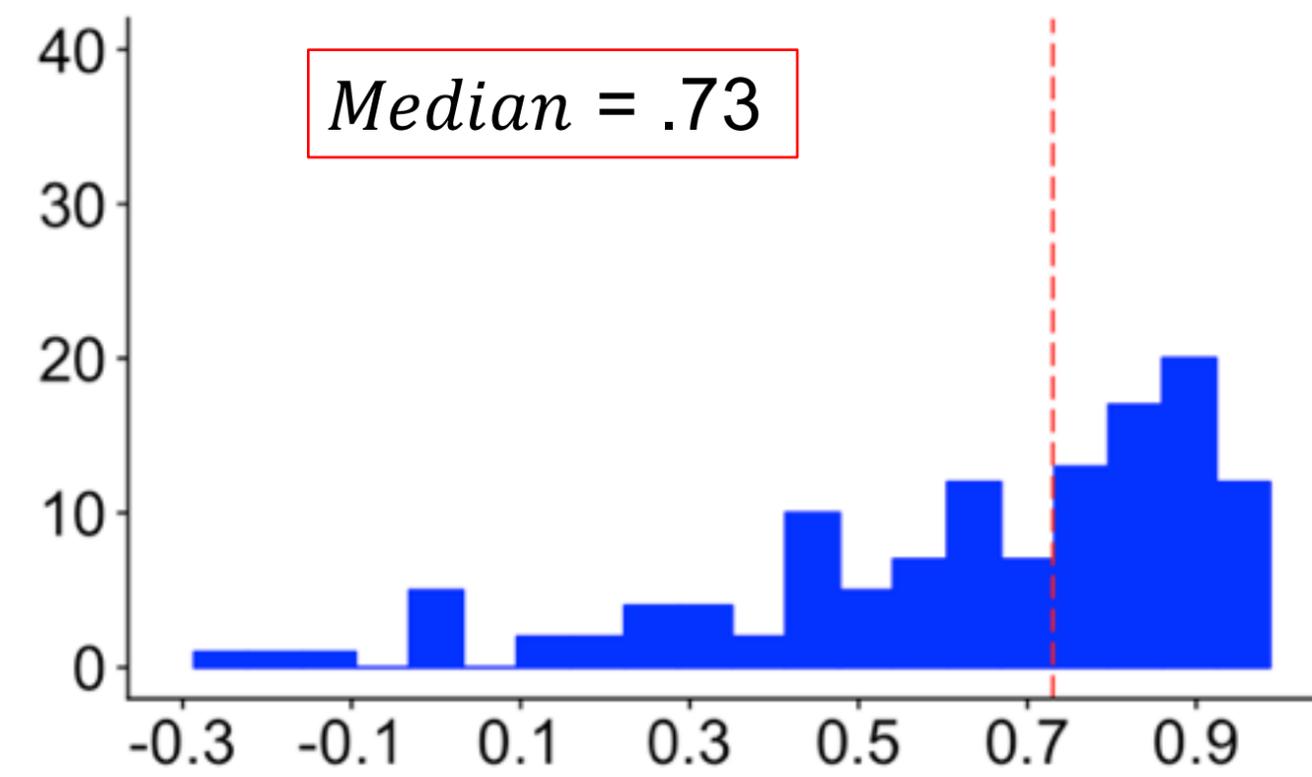
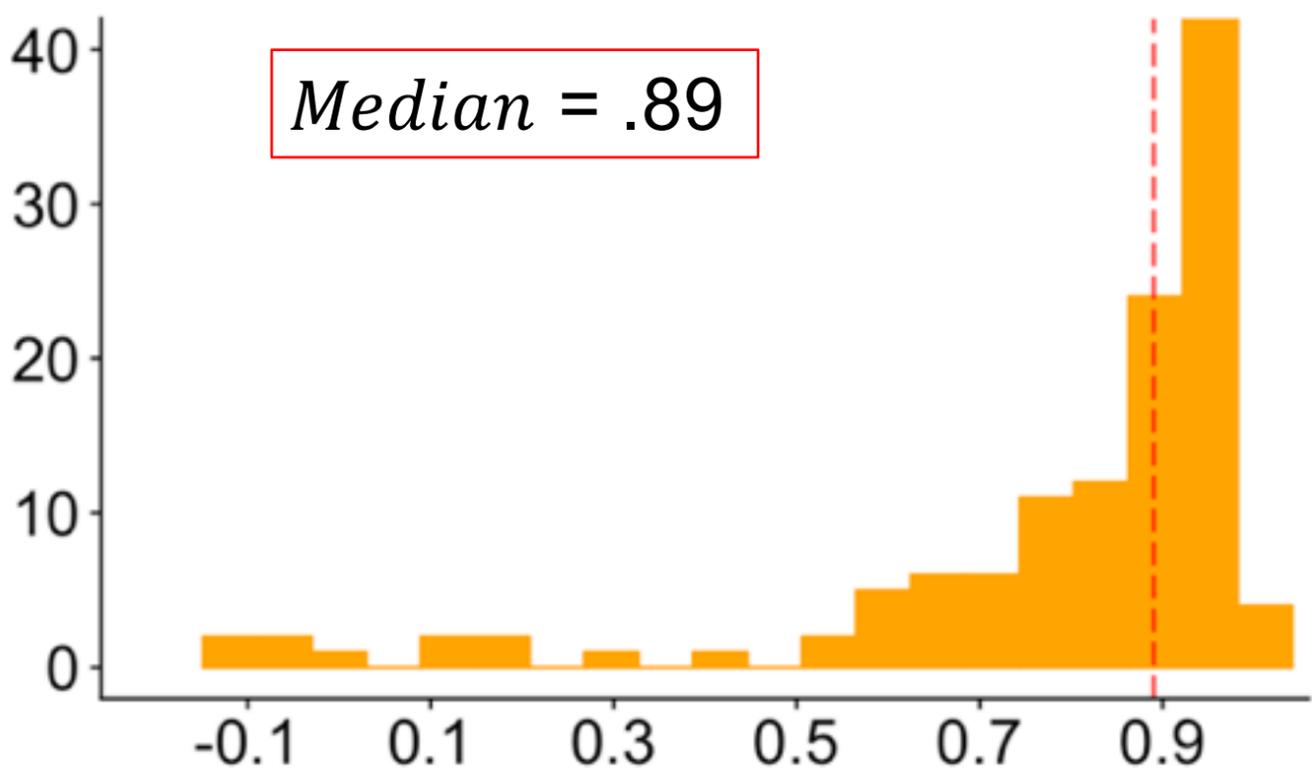
Within-participant Correlation

Median = .89

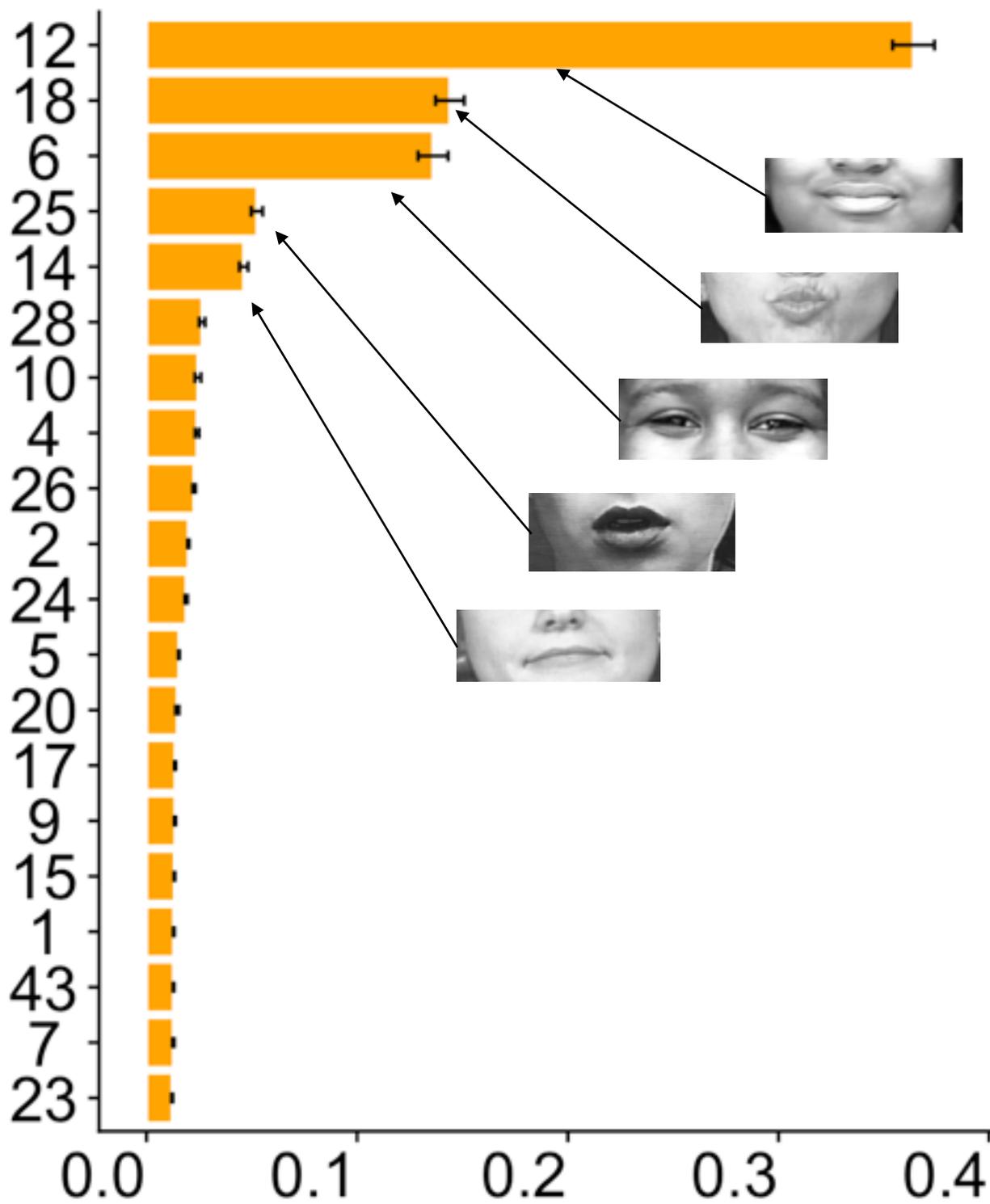
Median = .73

Median = .91

Median = .73



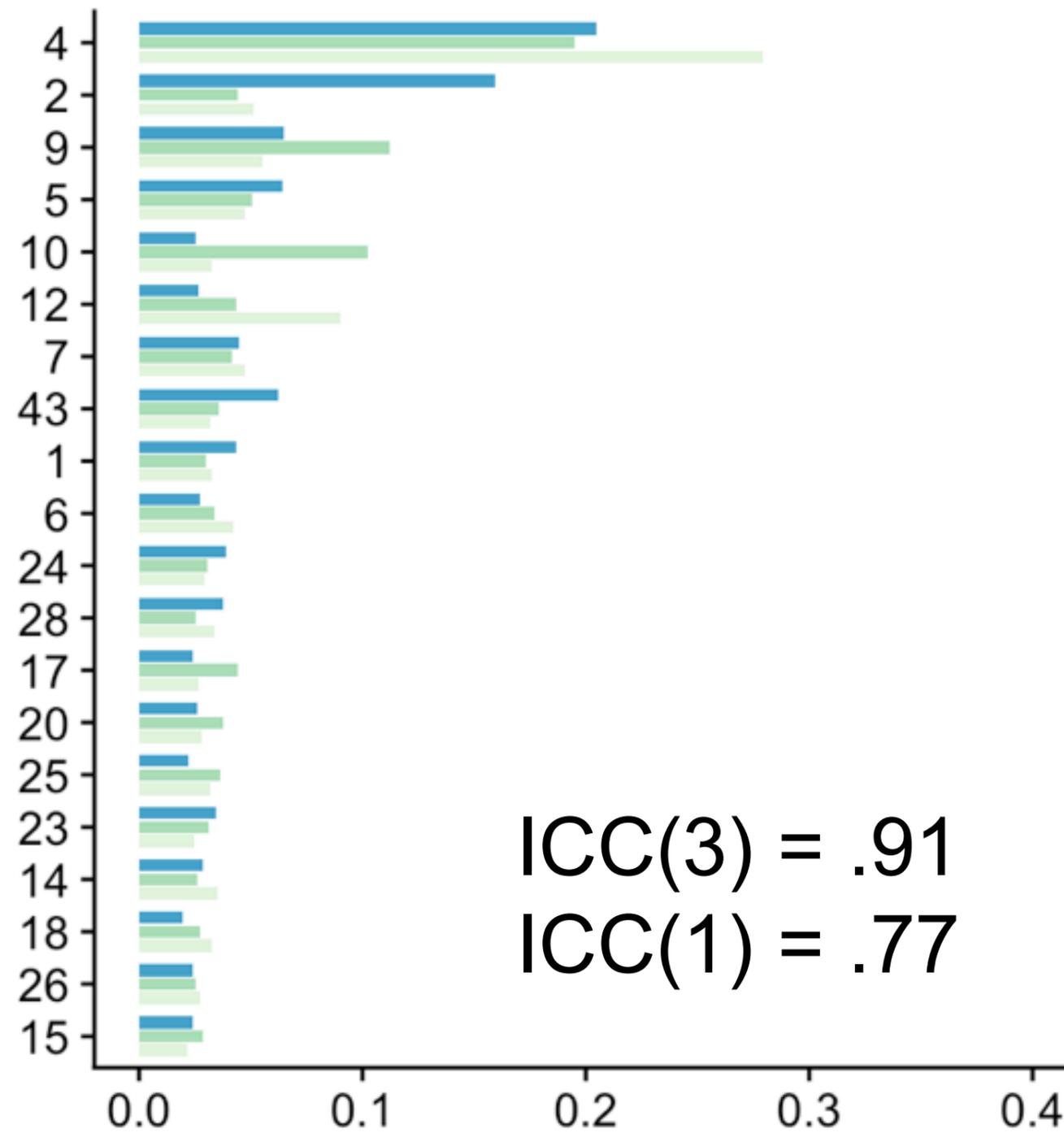
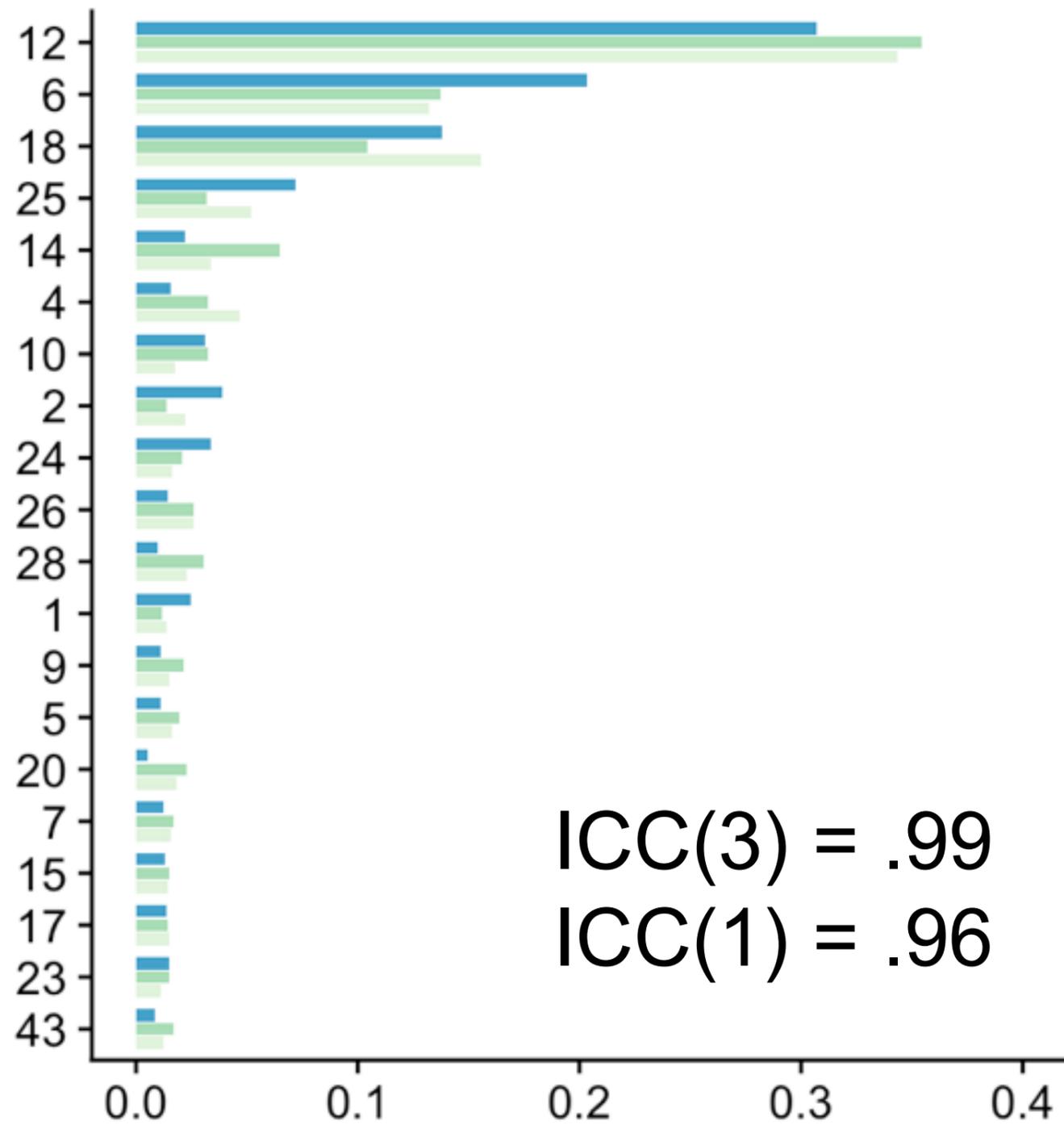
Action Unit



Positive

Negative

Action Unit



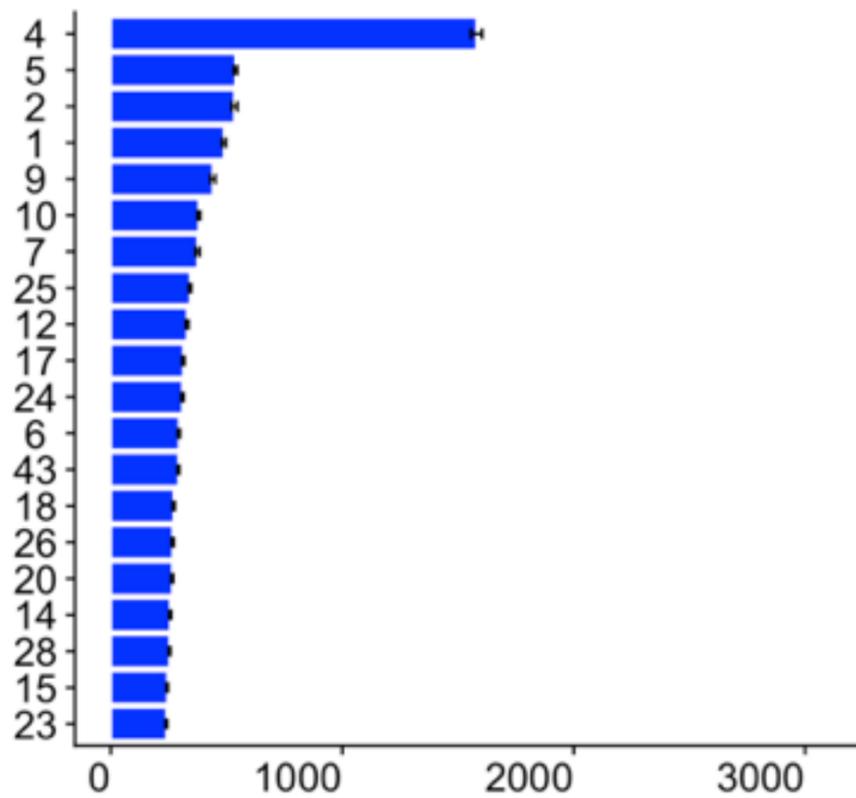
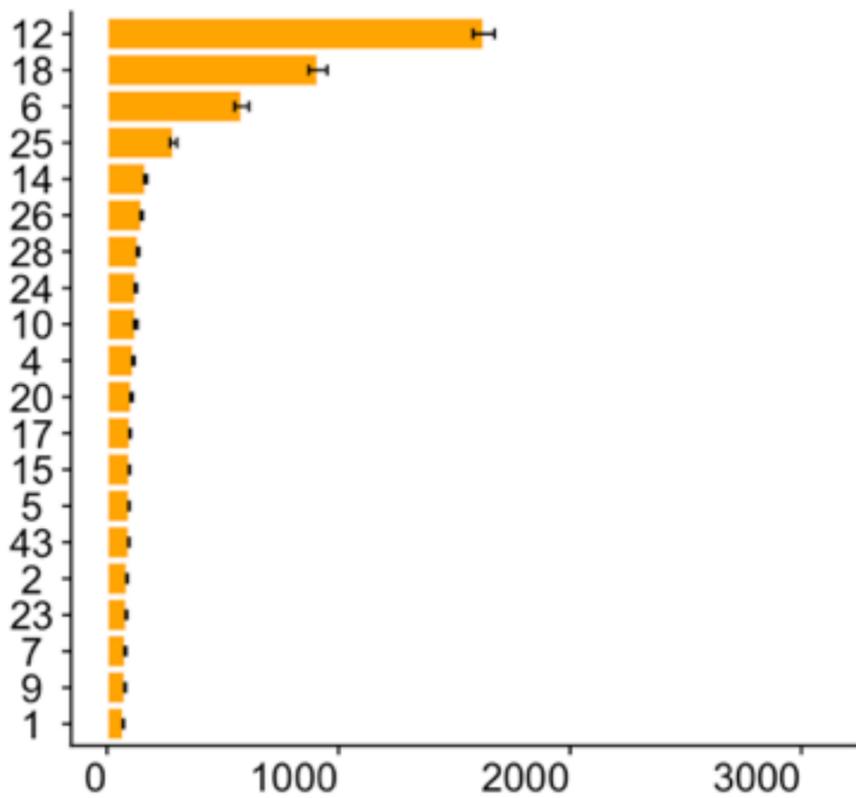
Condition
Express
Normal
Suppress

Relative Importance

Positive

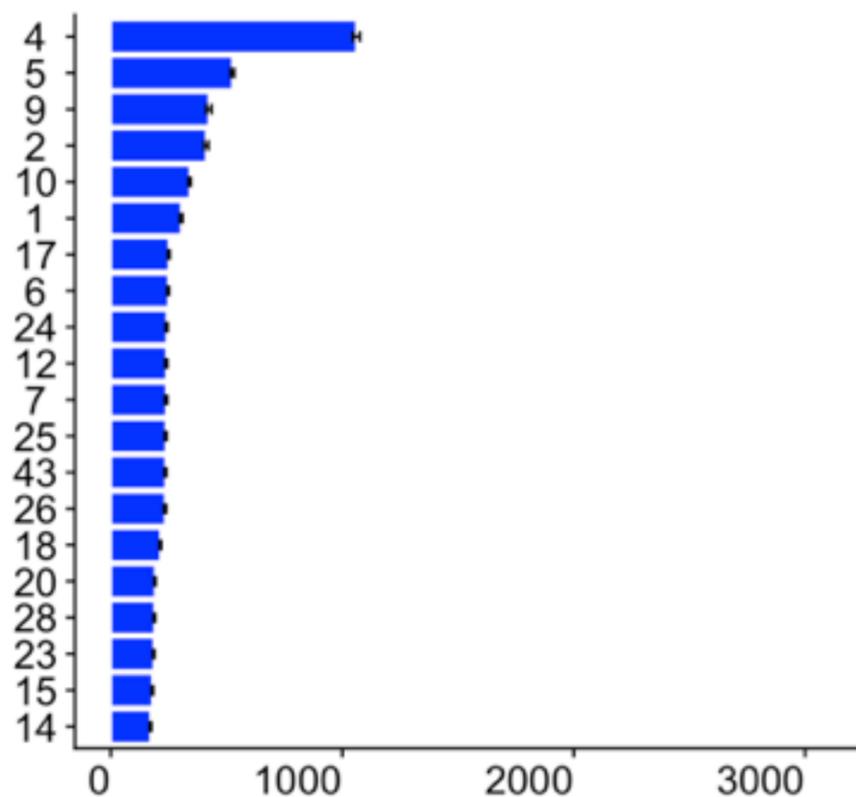
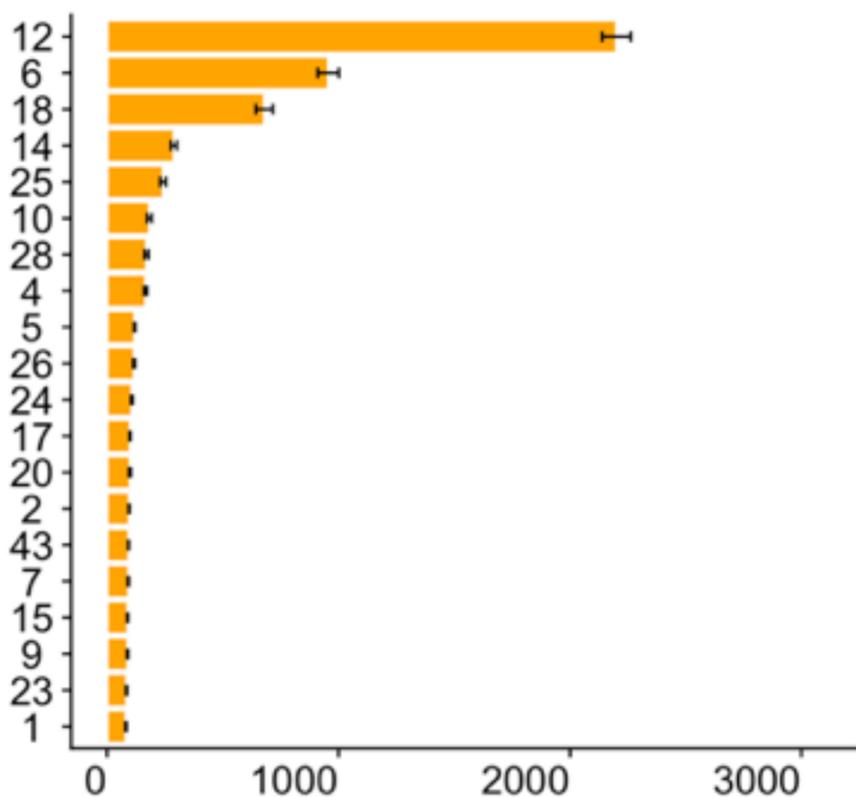
Negative

Coder 1

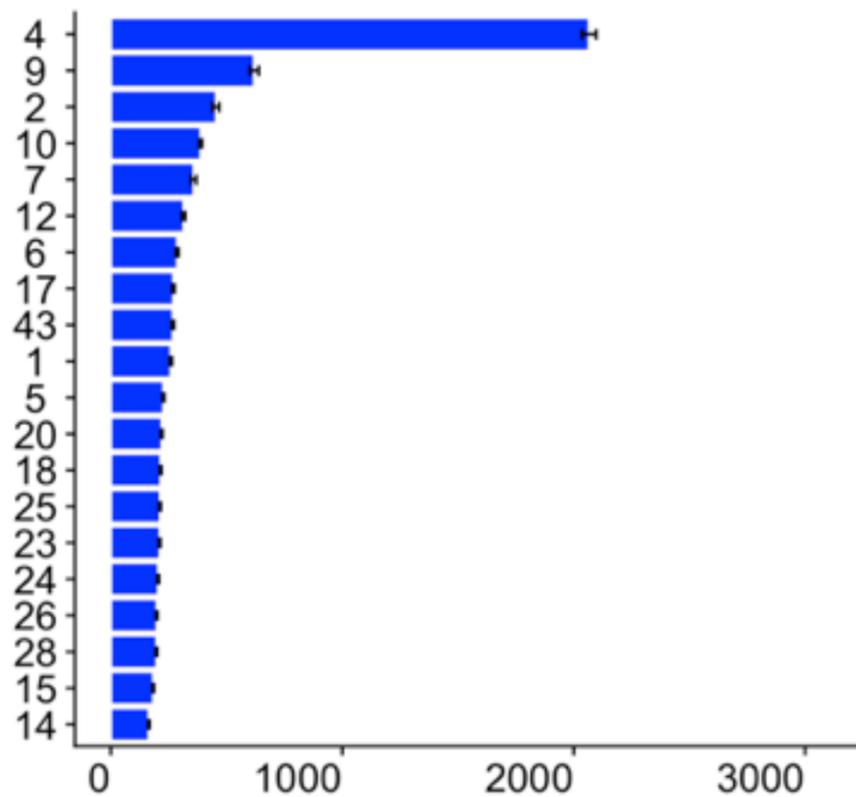
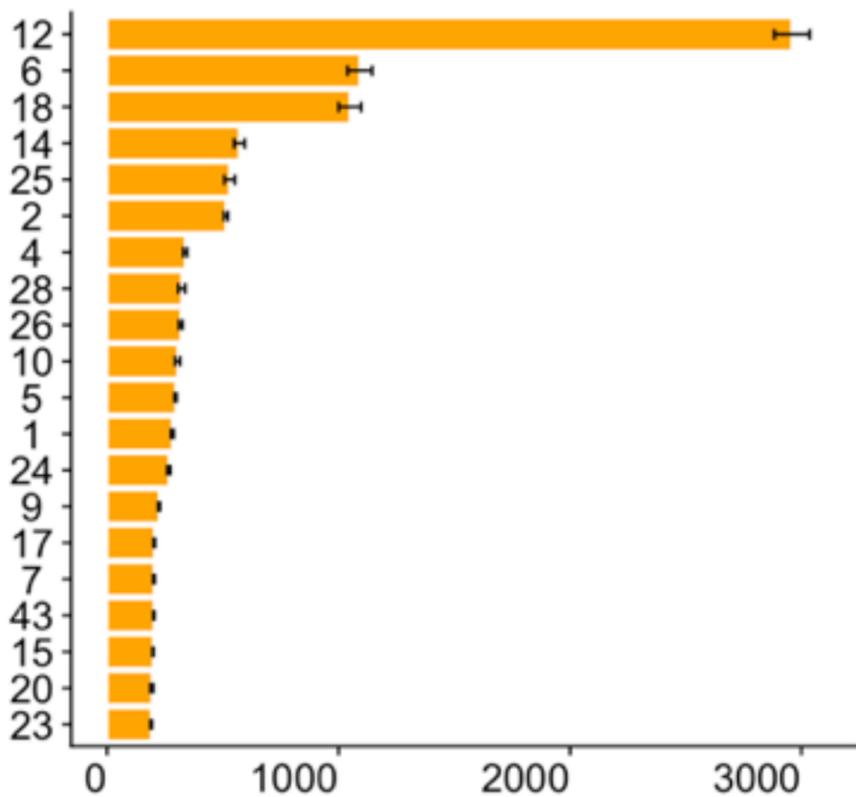


Action Unit

Coder 2



Coder 3



Increase in Node Purity

Positive

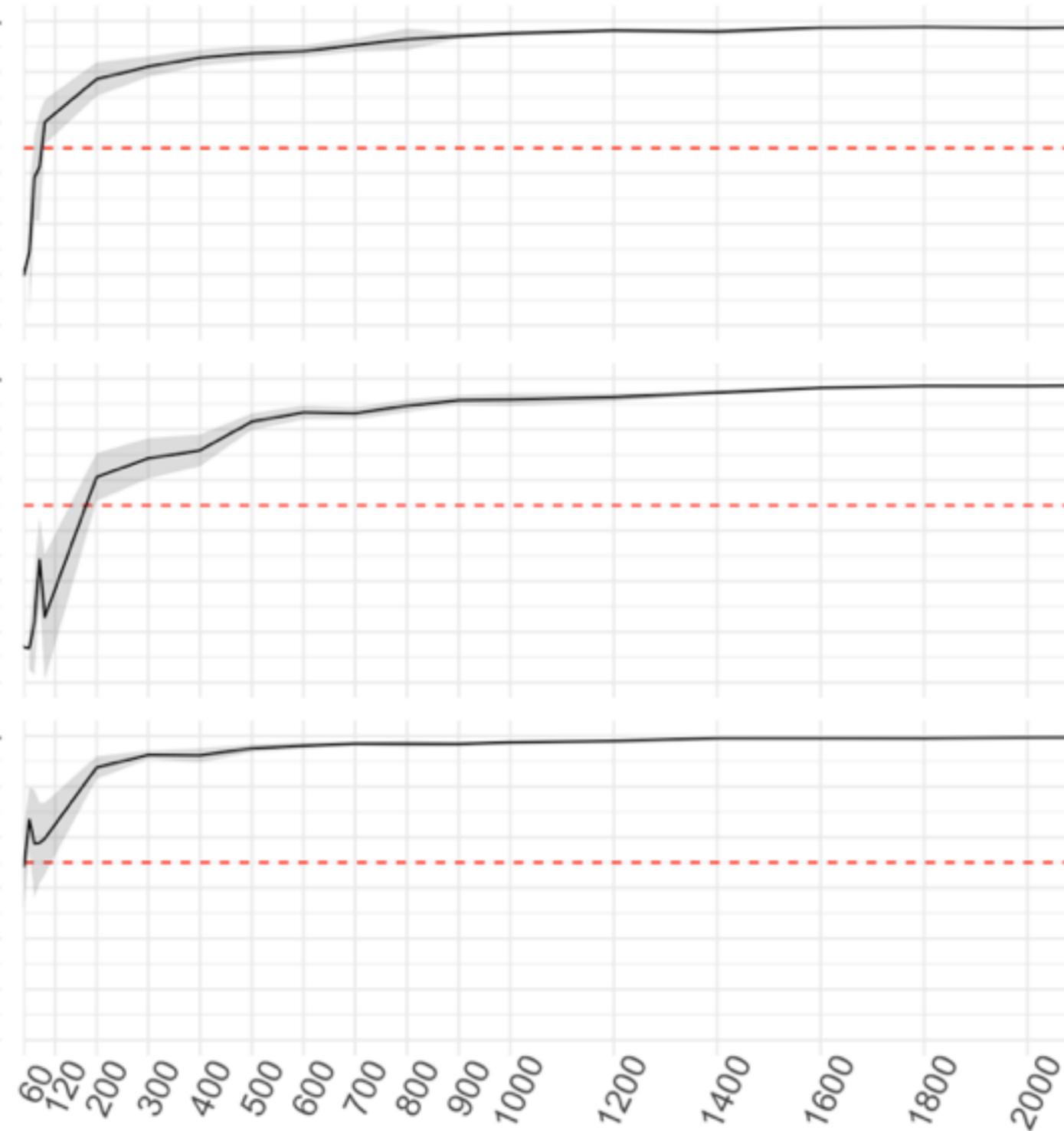
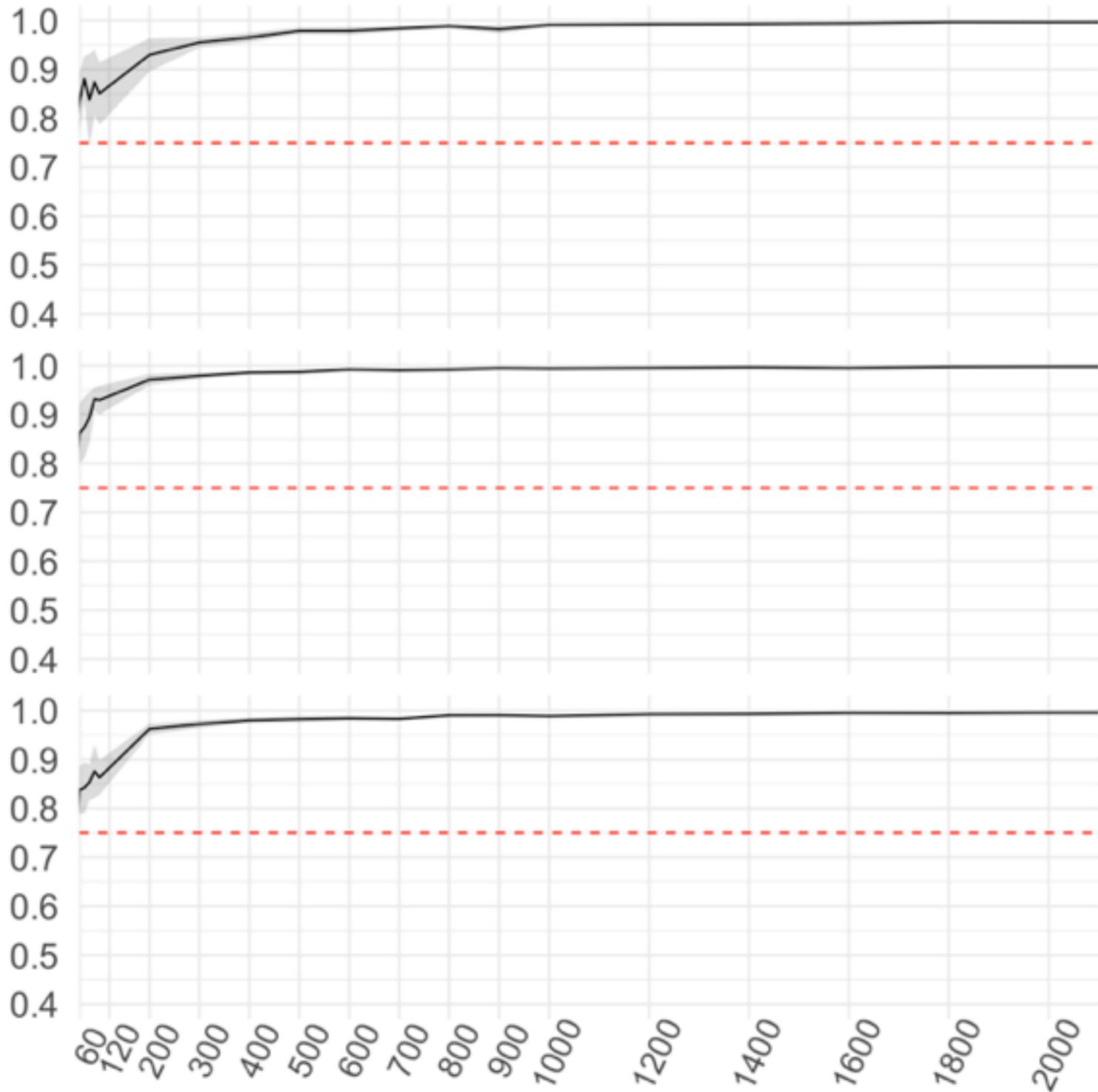
Negative

Mean ICC

Coder 1

Coder 2

Coder 3



Number of Recordings