# Automatic recognition of macaque facial expressions for detection of affective states

**Automatic recognition of macaque facial expressions for detection of affective states** 1

2

Anna Morozov[1], Lisa Parr[2,3], Katalin Gothard[4*], Rony Paz[1*], Raviv Pryluk[1*] 3

[1] Department of Neurobiology, Weizmann Institute of Science, Israel 4

[2] Department of Psychiatry and Behavioral Science, Emory University School of Medicine 5

[3] Yerkes Primate National Research Center, Emory University 6

[4] Department of Physiology, College of Medicine, University of Arizona 7

* Equal contribution 8

Correspondence should be addressed to: R.Pa. (rony.paz@weizmann.ac.il), K.G. 12
(kgothard@email.arizona.edu), R.Pr. (ravivpryluk@gmail.com) 13

Number of Figures: 5 14

Number of Tables: 1 15

Number of Multimedia: 0 16

Number of words for Abstract: 150 17

Number of words for Significance Statement: 119 18

Number of words for Introduction: 500 19

Number of words for Discussion: 914 20

Acknowledgements 21

Conflict of Interest 24

Authors report no conflict of interest 25

Funding sources 26

29

30

**Automatic recognition of macaque facial expressions for detection of affective states** 31

32

**Abstract** 33

Internal affective states produce external manifestations such as facial expressions. In humans, the Facial 34
Action Coding System (FACS) is widely used to objectively quantify the elemental facial action-units 35
(AUs) that build complex facial expressions. A similar system has been developed for macaque monkeys 36
- the Macaque Facial Action Coding System (MaqFACS); yet unlike the human counterpart, which is 37
already partially replaced by automatic algorithms, this system still requires labor-intensive coding. Here, 38
we developed and implemented the first prototype for automatic MaqFACS coding. We applied the 39
approach to the analysis of behavioral and neural data recorded from freely interacting macaque monkeys. 40
The method achieved high performance in recognition of six dominant AUs, generalizing between 41
conspecific individuals (*Macaca mulatta*) and even between species (*Macaca fascicularis*). The study 42
lays the foundation for fully automated detection of facial expressions in animals, which is crucial for 43
investigating the neural substrates of social and affective states. 44

45

**Significance Statement** 46

MaqFACS is a comprehensive coding system designed to objectively classify facial expressions based on 47
elemental facial movements designated as Actions Units (AUs). It allows the comparison of facial 48
expressions across individuals of same or different species based on manual scoring of videos, a labor- 49
and time-consuming process. We implemented the first automatic prototype for AUs coding in macaques. 50
Using machine learning, we trained the algorithm on video-frames with AU labels, and showed that after 51
parameter tuning, it classified six AUs in new individuals. Our method demonstrates concurrent validity 52
with manual MaqFACS coding and supports the usage of automated MaqFACS. Such automatic coding is 53
useful not only for social- and affective- neuroscience research but also for monitoring animal health and 54
welfare. 55

56

**Introduction**                                                                                    57

Facial expressions are both a means of social communication and also a window to the internal states of   58
an individual. The expression of emotions in man and animals was discussed first by Darwin in his        59
eponymous treatise in which he attributed the shared features of emotional expression in multiple species 60
to a common ancestor (Darwin 1872). Further elaboration of these ideas came from detailed                61
understanding of the neuromuscular substrate of facial expressions, i.e., the role of each muscle in moving 62
facial features into configurations that have social communicative value. These studies brought to light 63
the homologies, but also the differences in how single facial muscles, or groups of muscles give rise to a 64
relatively stereotypical repertoire of facial expressions (Ekman 1989, Ekman and Keltner 1997, Burrows, 65
Waller et al. 2006, Vick, Waller et al. 2007, Parr, Waller et al. 2010).                                 66

The affective states that give rise to facial expressions are instantiated by distinct patterns of neural 67
activity (Panksepp 2004) in areas of the brain that have projections to the facial motor nucleus in the  68
pons. The axons of the motor neurons in the facial nucleus distribute to the facial musculature, including 69
the muscles that move the pinna (Jenny and Saper 1987, Welt and Abbs 1990). Of all possible facial       70
muscle movements, only a small set of coordinated movements give rise to unique facial configurations    71
that correspond, with some variations, to primary affective states. Human studies of facial expressions  72
proposed six primary affective states or "universal emotions" that were present in facial displays across 73
cultures (Ekman and Friesen 1986, Fridlund, Ekman et al. 1987, Ekman and Friesen 1988, reviewed by       74
Ekman, Friesen et al. 2013). The cross-cultural features of facial expressions allowed the development of 75
an anatomically based Facial Action Coding System (FACS) (Friesen and Ekman 1978, Ekman, Friesen         76
et al. 2002). In this system, a numerical code is assigned for each elemental facial action that is identified 77
as an Action Unit (AU). Considering the phylogenetic continuity in the facial musculature across primate 78
species (Burrows and Smith 2003, Burrows, Waller et al. 2006, Burrows, Waller et al. 2009, Parr, Waller  79
et al. 2010), a natural extension of human FACS was the homologous MaqFACS (Parr, Waller et al.          80
2010), developed for coding the facial action units in Rhesus macaques (for multi-species FACS review    81
see: Waller, Julle-Daniere et al. 2020).                                                                 82

The manual scoring of action units (AUs) requires lengthy training and a meticulous certification process 83
for FACS coders, that is a time-consuming process. Therefore, considerable effort has been made towards  84
the development of automatic measurement of human facial behaviour (Sariyanidi, Gunes et al. 2015,       85
reviewed by Barrett, Adolphs et al. 2019). These advances do not translate seamlessly to macaque        86
monkeys, and importantly, similar developments are desirable because macaques are commonly used to       87
investigate and understand the neural underpinnings of communication via facial expressions (Livneh,    88
Resnik et al. 2012, Pryluk, Shohat et al. 2020). We therefore aimed to develop and test an automatic    89
system to classify AUs in macaques, one that would allow comparison of elicited facial expressions and  90
neural responses at similar temporal resolutions.                                                       91

Like humans, macaque monkeys do not normally activate a full set of action units required for a classical 92
stereotypical expression, and partial sets of uncommon combination of action units are also probable and 93
give rise to mixed or ambiguous facial expressions (Chevalier-Skolnikoff 1973, Ekman and Friesen       94
1976). Therefore, we chose to classify not only the fully developed facial expressions (Blumrosen,      95
Hawellek et al. 2017) but also action units that were shown to play a role in exhibition of affective states 96
and social communication among macaque monkeys. We included even relatively rare facial expressions    97
as long as certain action unit were reliably involved in these expressions. We test the automatic       98

3

recognition of facial configurations and show that it generalizes to new situations, between conspecific 99
individuals, and even across macaque species. Taken together, this work demonstrates concurrent validity 100
with manual MaqFACS coding and supports the usage of automated MaqFACS in social- and affective- 101
neuroscience research, as well as in monitoring animal health and welfare. 102

103

**Materials and Methods** 104

*Video datasets* 105

We used videos from two different datasets. The first, *Rhesus dataset (RD)*, consists of 53 videos from 106
five Rhesus macaques (selected out of 10 Rhesus monkeys). Part of this dataset was used for training and 107
testing our system within and across Rhesus subjects. The second, *Fascicularis dataset (FD)*, includes 108
two videos from two Fascicularis macaques and was used only for testing our system across Fascicularis 109
subjects. 110

All the videos in both sets capture frontal (or near-frontal) views of head-fixed monkeys. The video- 111
frames were coded for the AUs present in each frame (none, one, or many). 112

The subjects and the videos for RD were selected with respect to the available data in FD, considering the 113
scale similarity, the filming angle and the AU frequencies occurring in the videos. 114

*The Rhesus Macaque Facial Action Coding System (MaqFACS)* 115

There are several stereotypical facial expressions that macaques produce (Fig. 1A), that represent, as in 116
humans, only a subset of the full repertoire of all the possible facial movements. For example, (Fig. 1B) 117
represents three common facial expressions from the Fascicularis monkey dataset (FD) (left, blue) and 118
two other facial configurations that, among others, occurred in our experiments (right, yellow). Therefore, 119
to allow the potential identification of all the possible facial movements (both the common and the less 120
common ones), we chose to work in the MaqFACS domain and to recognize AUs, rather than searching 121
for predefined stereotypical facial expressions. MaqFACS contains three main groups of AUs based on 122
facial sectors: upper face, lower face, and ears (Parr, Waller et al. 2010). Each facial expression is 123
instantiated by a select combination of AUs (Fig. 1C). 124

*AU selection* 125

The criteria for AU selection for the analysis in this work, were their frequencies (which should be 126
sufficient for training and testing purposes) and the importance of each AU for affective communication 127
(Fig. 1D, E) (Parr, Waller et al. 2010, Ballesta et al. 2016, Mosher et al. 2016). Frequent combinations of 128
lower face AUs together with upper face AUs (Fig. 1F outside the magenta and green frames) may hint at 129
the most recurring facial expressions in the test set. For example, UpperNone AU together with lower 130
face AU25, generate a near-neutral facial expression. Considering that our aim is to recognize single AUs 131
(as opposed to complete predefined facial expressions), lower face and upper face AUs were not merged 132
into single analysis units. This approach is also supported by the MaqFACS coding process, which is 133
performed separately for the lower and upper face. 134

The most frequent upper face AUs in FD were the none-action AU (defined here as "UpperNone"), the Brow Raiser AU1+2 and AU43_5, which is a union of Eye Closure AU43 and Blink AU45 (Fig. 1D). The two latter AUs differ only in the movement duration, and hence were joined.

There were five relatively frequent AUs in the lower face test set (Fig. 1E) that we merged into several AU groupings. All AUs that mostly co-occurred with other ones (within the same face region) were analyzed as a combination rather than single units (Fig. 1F inside the green frame). The upper face AUs however, rarely appeared as combination (Fig. 1F inside the magenta frame).

Overall, our system was trained to classify 6 units: AU1+2, AU43_5 and UpperNone in the upper face, and AU25+26, AU25+26+16 and AU25+26+18i in the lower face (Fig. 1G and 1H, left). Even though AU12 was one of the most prevalent AUs in the FD test set and often occurred in combination with other lower face AUs, it was eliminated from further analysis because it appeared too infrequently in the Rhesus monkey dataset (RD).

*Animals and procedures*

All surgical and experimental procedures were approved and conducted in accordance with the regulations of the Institute Animal Care and Use Committee (IACUC), following NIH regulations and with AAALAC accreditation.

Two male Fascicularis monkeys (*Macaca fascicularis*) and 10 Rhesus monkeys (*Macaca mulatta*) were videotaped while producing spontaneous facial movements. All monkeys were seated and head-fixed in a well-lit room during the experimental sessions.

The two monkeys produced facial behaviors in the context described in detail in (Pryluk, Shohat et al. 2020) (Fig. 2, Fig. 2-1, Fig. 2-2, Fig. 2-3). The facial movements obtained during neural recordings have not been previously analyzed in terms of action units. Earlier experiments showed that self-executed facial movements recruit cells in the amygdala (Livneh, Resnik et al. 2012, Mosher, Zimmerman et al. 2016) and the ACC (Livneh, Resnik et al. 2012) and that neural activity in these regions is temporally locked to different socially meaningful, communicative facial movements (Livneh, Resnik et al. 2012). The video data from these monkeys was captured using two Ximea_MQ013RG (Ximea GmbH, Munster, Germany) cameras (one camera for the whole face and one dedicated to the eyes), with Kowa (Kowa Optical Products Co. Ltd., Saitama, Japan) lenses mounted on them: 16mm LM16JC10M for the face- and 25mm LM25JC5M2 for the eye-camera. The frame rates of the face- and eye-videos are 34 frames per second (~29ms) and 17 frames per second (~59ms), respectively. The size parameters are 800x700 pixels for the facial videos and 700x300 pixels for the videos of eyes. Both video types have 8-bit precision for grayscale values. The lighting in the experiment room included white LED lamps and an infrared LED light bar (MetaBright Exolight ISO-14-IRN-24, Metaphase Technologies, Philadelphia, PA, USA) for face illumination.

The 10 Rhesus monkeys were filmed during baseline sessions as well as during provocation of facial movements by exposure to a mirror and to videos of other monkeys. Videos of facial expressions of the Rhesus macaques were recorded at 30 frames per second (~33ms) rate, with 1280x720 pixels size parameters and 24-bit precision for RGB values.

*Behavioural paradigms* 174

The intruder task is similar to the one described in (Pryluk, Shohat et al. 2020), including a monkey 175
intruder instead a of human (Fig. 2, Fig. 2-1, Fig. 2-2, Fig. 2-3). A single experimental block includes 6 176
interactions (trials) with a monkey intruder that is seated behind a fast LCD shutter (<1ms response time, 177
307mm x 407 mm) which is used to block the visual site. When the shutter opens, the monkeys are able 178
to see each other. Each trial is of ~9 sec and the shutter is closed for ~1 sec between the trials. Altogether, 179
the length of the interaction part (from the first shutter opening and until its last closure) is 60 sec. 180

We recorded the facial expressions of the subject monkey, along with monitoring the intruder monkey 181
behavior. When the intruder monkey was brought to or out from the room (the "enter-exit" stage), the 182
shutter was closed and the subject monkey could not see any part of the intruder unless the shutter was 183
open. The "enter" and the "exit" phases were of 30 sec long each. 184

*Data labeling* 185

Video-data annotation was carried out using Noldus software "The Observer XT" 186
(https://www.noldus.com/human-behavior-research/products/the-observer-xt). The recorded behavior 187
coding was exported in Excel (Microsoft Excel 2016) format for further processing. 188

RD videos were labeled by FACS- (Friesen and Ekman 1978, Ekman, Friesen et al. 2002) and 189
MaqFACS- (Parr, Waller et al. 2010) accredited coding expert. Another trained observer performed the 190
coding of all FD videos according to the MaqFACS manual based on (Parr, Waller et al. 2010). Facial 191
behavior definitions were discussed and agreed prior to the coding. To ensure consistency, we checked 192
the inter-rater reliability (IRR) for one of the two FD videos, against additional experienced coder. Our 193
target percentage of agreement between observers was set to 80% (Baesler and Burgoon 1987) and the 194
IRR test resulted with 88% agreement (Figure 5-1). 195

All the videos were coded for MaqFACS AUs along with their frequencies and intensities. Analyzed 196
frames with no labels were considered as frames with neutral expression. Upper- and lower-face AUs 197
were coded separately. This partition was inspired by observations indicating that facial actions in the 198
lower face have little influence on facial motion in the upper face and vice versa (Friesen and Ekman 199
1978). Moreover, neurological evidence suggests that lower and upper face are engaged differently by 200
facial expressions and their muscles are controlled by anatomically distinct motor areas (Morecraft, Louie 201
et al. 2001). 202

*Image preprocessing* 203

For each video from both datasets, seven landmark points (two corners of each eye, two corners of the 204
mouth and the mouth center) were manually located on the mean image of frames with neutral 205
expression. For image height $h$ and width $w$, the reference landmark points were defined by the following 206
coordinates: (0.42w, 0.3h) and (0.48w, 0.3h) for left eye corners, (0.52w, 0.3h) and (0.58w, 0.3h) for right 207
eye corners, (0.44w, 0.55h) for mouth left corner, (0.56w, 0.55h) for mouth right corner and (0.5w, 0.5h) 208
for the mouth center (Fig. 3-1). 209
Affine transformations (geometric transformations that preserve lines and parallelism, e.g., rotation) were 210
applied to all frames of all videos so that the landmark points were mapped to predefined reference 211
locations (Fig. 3A, Fig. 3-1). The alignment procedure was necessary to correct any movement, either 212
from the alignment of the camera (angle, distance, height) or movement of the monkey, that would shift 213

6

the facial landmarks between video frames. After the alignment procedure, total average image of all 214
mean neutral expression frames was calculated. Two rectangular ROIs (regions of interest), one for the 215
upper face and one for lower face, were marked manually on the total average image (Fig. 3B). Finally, 216
all the frames were cropped according to the ROI windows (Fig. 3C), resulting in 396x177 pixel upper 217
face images and 354x231 pixel lower face images. After this step, the originally RGB images were 218
converted to grayscale. For each video, one "optimal" neutral expression frame was selected out of all the 219
neutral expression images. Difference images (δ-images) were generated by subtraction of the optimal 220
neutral frame from all the frames of the video (Fig. 3D, Fig. 1G and 1H, right). The main idea behind this 221
operation was to eliminate variability due to texture differences in appearance (e.g. illumination changes), 222
and to analyze variability of facial distortions (e.g. action units) and individual differences in facial 223
distortion (Bartlett, Viola et al. 1996). In the last preprocessing step, upper face and lower face databases 224
(DBs) were created by converting the δ-images to single dimension vectors and storing them as a 2- 225
dimnesional matrix containing the pixel brightness values (one dimension is of size of the total image 226
pixels and the second represents the images quantity). The DBs were then used for construction of 227
training and test sets (Fig. 3E). 228

### *Eigenfaces: Dimensionality reduction and feature extraction* 229

Under controlled head-pose and imaging conditions, the statistical structure of facial expressions may be 230
efficiently captured by features extracted from Principal Component Analysis (PCA) (Calder, Burton et 231
al. 2001). This was demonstrated in the "EigenActions" technique (Donato, Bartlett et al. 1999), where 232
the facial actions were recognized separately for upper face and lower face images (the well-known 233
"Eigenfaces"). According to this technique, the PCA is used to compute a set of subspace basis vectors 234
(referred to as the ''eigenfaces'') for a dataset of facial images (the training set), which are then projected 235
into the compressed subspace. Typically, only the N eigenvectors associated with the largest eigenvalues 236
are used to define the subspace, where N is the desired subspace dimensionality (Draper, Baek et al. 237
2003). Each image in the training set may be represented and reconstructed by the mean image of the set 238
and a linear combination of its principal components (PCs). The PCs are the eigenfaces and the 239
coefficients of the PCs in the linear combination instance their weights. The test images are matched to 240
the training set by projecting them onto the basis vectors and finding the nearest compressed image in the 241
subspace (the eigenspace). 242

We applied the eigenfaces analysis on the training frames (the δ-images), which were first zero-meaned 243
(Fig. 3F). Once the eigenvectors were calculated, they were normalized to unit length, and the vectors 244
corresponding to the smallest eigenvalues (under $10^{-6}$) were eliminated. 245

### *Classification* 246

One of the benefits of the mean subtraction and the scaling to unit vectors is that this operation projects 247
the images into a subspace where Euclidean distance is inversely proportional to correlation between the 248
original images. Therefore, nearest neighbour matching in eigenspace establishes an efficient 249
approximation to image correlation (Draper, Baek et al. 2003). Consequently, we employed a K-Nearest 250
Neighbors (KNN) classifier in our system. Related to the choice of classifier, previous studies show that 251
when PCA is used, the choice of the subspace distance-measure depends on the nature of the 252
classification task (Draper, Baek et al. 2003). Based on this notion and other observations (Bartlett, 253
Donato et al. 2000), we chose the Euclidian distance and the cosine of the angle between feature vectors 254
to measure similarity. In addition, to increase the generality of our approach and to validate our results, 255

7

we also tested a Support Vector Machine (SVM) classifier. To evaluate the performance of the models we 256
define a classification trial as successful if the AU predicted by the classifier was the same as in the probe 257
image. To further justify the classification of AUs separately for upper face and lower face ROIs, it is 258
worth mentioning that evidence suggest that PCA-based techniques performed on full-face images lead to 259
poorer performance in emotion recognition compared to separate PCA for the upper and lower regions 260
(Padgett and Cottrell 1997, Bartlett 2001). 261

To train a classification model for AUs recognition, we used the weights of the principal components 262
(PCs) as predictors. To predict the AU of a new probe image, the probe should be projected onto the 263
eigenspace to estimate its weights (Fig. 3F). Once the weights are known, AU classification may be 264
applied. The output of the classifier of each facial ROI is the AU that is present in the frame (Fig. 3G). To 265
increase the generality of our approach and to validate our results, we used both K-Nearest Neighbors 266
(KNN) and Support Vector Machine (SVM) classifiers. 267

*Parameter selection* 268

In the KNN classification, we examined the variation of three main parameters: the number of the 269
eigenspace dimensions (PCs), the subspace distance metric and $k$ - the number of nearest neighbors in the 270
KNN classifier. 271

Multiple ranges of PCs were tested (the "pcExplVar" parameter), from PC quantity that cumulatively 272
explains 50% of the variance of each training set to 95%; $k$ was varied from 1 to 12 nearest neighbors and 273
the performance was also tested with Euclidian and cosine similarity measures. For each training set and 274
parameter set, the features were recomputed and the model performance was re-estimated. The process 275
was repeated across all the balanced training sets (see *Data under-sampling*). The parameters of the 276
models and the balanced training sets were selected according to the best classification performance in the 277
validation process. 278

*Data under-sampling* 279

The training sets in this study were composed of Rhesus Dataset (RD) frames from AU1+2, AU43_5 and 280
UpperNone categories in the upper face, and AU25+26, AU25+26+16 and AU25+26+18i in the lower 281
face (in a non-overlapping manner relatively to each ROI). For the training purposes, for both ROIs, the 282
RD frames were randomly under-sampled 3-10 times (depending on the data volume), producing the 283
"balanced training sets". The main reason for this procedure was to balance the frame quantity of the 284
different AUs in the training sets (He and Garcia 2009). For each dataset, the size of the balanced training 285
set was defined based on the smallest category size (Table 1). As a result, for the training processes in our 286
experiments we used upper face and lower face balanced training sets of size 3639 and 930 frames each, 287
correspondingly. 288

It should be noted that the under-sampling procedure influences only the training but not the test sets 289
composition (only the frames for training are selected from the balanced training sets). The test set 290
composition depends on the subjects and the videos selected for the testing, and considers all the available 291
frames that fit the task criteria (consequently, they are the same across all the balanced training sets). 292

*Validation and model evaluation* 293

We tested three types of generalization. For each type of generalization, the performance was evaluated 294
independently for upper face and lower face, using holdout validation for the Fascicularis data (Geisser 295

1975) and leave-one-out cross validation (CV) for the Rhesus data (Tukey 1958). The leave-one-out    296
technique is advantageous for small datasets because it maximizes the available information for training,    297
removing only a small amount of training data in each iteration. Applying the leave-one-out CV, data    298
from all subjects (or videos) but one, was used for the system training, and the testing was performed on    299
the one remaining subject (or video). We designed the CV partitions constraining equal number of frames    300
in each class of the training sets. In both the leave-one-out CV and the holdout validation, images of the    301
test sets were not part of the corresponding training sets, and only the training frames were retrieved from    302
the balanced training sets. To ensure the data sufficiency for training and testing, a subject (or video) was    303
included in the partition for CV only if it had enough frames of the three AU classes (separately for upper    304
face and lower face).    305

For each generalization type, the training and the testing sets were constructed as following:    306

1. *Within subject (Rhesus):* for each CV partition, frames from all videos but one, from the same    307
   Rhesus subject, were used for training. Frames of the remaining video were used for testing.    308
   Performed on RD, on three balanced training sets. To be included in a CV partition for testing,    309
   the training and the test sets for a video had to consist of at least 20 and 5 frames per class,    310
   correspondingly. Some subjects did not meet the condition, and this elimination process resulted    311
   with three subjects for upper face and four subjects for lower face CV.    312

2. *Across subjects (Rhesus):* for each CV partition, frames from all videos of all Rhesus monkeys    313
   but one, were used for training. Each test set was composed of frames from videos of the one    314
   remaining monkey. Performed on RD, on three balanced training sets. To be included for testing    315
   in the CV, the training and the test sets for a subject had to contain at least 150 and 50 frames of    316
   each class, correspondingly. In total, four subjects were included in the upper face and three in    317
   the lower face testing.    318

3. *Across Species:* frames from all videos of the five Rhesus monkeys were used for training.    319
   Frames from the two Fascicularis monkeys were used for validation and testing. In this case, a    320
   holdout model validation was carried out independently for each Fascicularis monkey (each    321
   subject had a different set of model parameters selected). For this matter, each Fascicularis    322
   monkey's dataset was randomly split 100 times in a stratified manner (so the sets will have    323
   roughly the same class proportions as in the original dataset) to create two sets: validation set    324
   with 80% of the data and test set with 20% of the data. Overall, the training sets were constructed    325
   from 10 balanced training sets of the Rhesus dataset. Validation and test sets (produced by 100    326
   splits in total) included 80% and 20% of the Fascicularis dataset, correspondingly. The best    327
   model parameters were selected according to the mean performance in validation set (over 100    328
   splits), and the final model evaluation was calculated based on the test set mean performance    329
   (over the 100 splits, as well).    330

***Performance measures***    331

Although the balanced training sets and the CV partitions were constructed to maintain the total number    332
of actions as even as possible, the subjects and their videos in these sets possessed different quantities of    333
actions. In addition, while we constrained the sizes of the classes within each training set to be equal, we    334
used the complete available data for the test sets. Since the overall classification correct rate (accuracy)    335
may be an unreliable performance measure due to its dependency on the targets to non-targets proportion    336
(Pantic and Bartlett 2007), we also applied a sensitivity measure (Benitez-Quiroz, Srinivasan et al. 2017)    337
for each AU (where the target is the particular AU and the non-targets are the two remaining AUs).    338

9

We used the average sensitivity measure (average true positive rate - $\overline{TPR}$) to select the best parameter set. To compare the performance of the classifiers, we present the generalization results on a *subject* (i.e., individual monkey) level (rather than *video*), for each classification type. Performance on Fascicularis dataset is reported as the mean performance of two parameter sets (one set per subject). 339 340 341 342

*Single-neuron activity analysis* 343

We analyzed a subset of neurons which were previously reported in (Pryluk, Shohat et al. 2020) and corresponded to the relevant blocks of monkey-monkey interactions. The neural analysis was performed with respect to facial AUs, focusing on 400-700 ms before and after the start of AU elicitation by the subject monkey. 344 345 346 347

Neural activity was normalized according to the baseline activity before the relevant block, using the same window length (300 ms) to calculate the mean and s.d. of the firing rate. 348 349

Therefore, the normalized (z-scored) firing rate (FR) was: 350

$$FR_{normalized} = \frac{FR - mean_{baseline}}{s.d._{baseline}}$$

*Software* 351

A custom code for automatic MaqFACS recognition and data analysis was written in Matlab R2017b (https://www.mathworks.com/). The code described in the paper is freely available online at [URL redacted for double-blind review]. The code is available as Extended Data. 352 353 354

**Results** 355

*Eigenfaces – unraveling the hidden space of facial expressions* 356

Intuitively, light and dark pixels in the *eigenfaces* (Fig. 4A, B) reveal the variation of facial features across the dataset. To further interpret their putative meaning, we varied the eigenface weights to demonstrate their range in the training set, producing an image sequence for each PC (Fig. 4C, D). This suggests that PC1 of this upper face set (Fig. 4C top, left-to-right) codes brows raising (AU1+2) and eyes opening (AU43_5). In contrast, PC2 resembles eyes closure (Fig. 4C bottom, bottom-up). Similarly, PC1 of the lower face set (Fig. 4D top, left-to-right) probably describes nose and jaw movement. Finally, PC2 for the lower face (Fig. 4D bottom, bottom-up) plausibly correspond to nose, jaw and lip movements, reminding the transition from pushed forward lips (AU25+26+18i) to depressed lower lip (AU25+26+16). 357 358 359 360 361 362 363 364 365

To illustrate the *eigenspace* concept, we present decision surfaces of two trained classifiers (Fig. 4E,F), along their first two dimensions (the weights of PC1 and PC2) which account for changes in facial appearance in (Fig. 4C,D). We show several training and test samples along with their locations following the projection onto the eigenspace. The projection of the samples is performed to estimate their weights, which are then used by the classifier as predictors. 366 367 368 369 370

371

10

*Parameter selection* 373

Example of parameter selection (Materials and Methods) for a Fascicularis subject is shown in Fig. 5A. 374
Interestingly, this upper face classification required much larger pcExplVar (93% versus 60% in the lower 375
face; the difference observed in both Fascicularis subjects). Specifically, this upper face classifier 376
achieved its best performance with 264 PCs, opposed to the lower face classifier succeeding with only 15 377
PCs (Fig. 5B). The most likely explanation is the large difference between the training-set sizes (3639 378
upper face versus 930 lower face images). Additionally, the eye-movement in the upper face images may 379
require many PCs to express its variance. 380

In contrast, the pcExplVar parameter behaved differently for generalizations *within* and *across* Rhesus 381
subjects: their best upper face classifiers required pcExplVar of 85%, and 83% in the lower face sets. The 382
notable difference between the parameters of these datasets suggests that one should tune a different 383
parameter set for each dataset. Generally, the Rhesus dataset required much larger pcExplVar to describe 384
the lower face than the Fascicularis dataset. 385

*Performance analysis* 386

Overall, the best parameter set for generalization to a new video *within subject (Rhesus)* using KNN 387
(Materials and Methods), performed with 81% accuracy and 74% $\overline{TPR}$ per subject for upper face, along 388
with 69% accuracy and 62% $\overline{TPR}$ for lower face, where the chance-level is 33% (Fig. 5C, left). Best 389
generalization *across subjects (Rhesus)* yielded $\overline{TPR}$ of 72% and 53% for upper and lower face 390
respectively, with corresponding accuracy of 75% and 43% (Fig. 5C, middle), compared to 33% chance- 391
level. The better performance in the upper face may be explained by its larger number of subjects in the 392
CV (four in the upper face, only three in the lower face) and by greater number of examples available for 393
training. Interestingly, applying the best parameter set of generalization *within subject* to classifiers 394
generalizing *across subjects*, produced close-to-best performance (upper face 71% and lower face 50% 395
$\overline{TPR}$). This finding suggests that tuning KNN parameters for generalization *within* Rhesus subjects, might 396
be enough also for *across-Rhesus-subjects* generalization. 397

The finest results, however, were achieved in generalization *between species* with 84% $\overline{TPR}$ for upper 398
face and 83% for lower face, with corresponding accuracy of 81% and 90%, concerning 33% chance- 399
level (Fig. 5C, right). To examine whether our findings depend on the particular classification algorithm, 400
we additionally tested this generalization with multiclass Support Vector Machine (SVM) approach. This 401
improved the $\overline{TPR}$ to 89% for both ROIs, indicating the advantage of using eigenfaces-based techniques 402
for MaqFACS AUs classification. 403

Finally, we have also compared the performance of the classifier to the human coders to determine 404
whether the algorithm is superior or inferior to the average, the slow and somewhat subjective human 405
decision. Due to the variability between raters, we found that that the algorithm was more accurate for 406
certain AUs whereas the human raters were more accurate for other AUs (data shown in Figure 5-1). 407
Specifically, for UpperNone AU, the classifier had average sensitivity of 84% vs. 81% in the human 408
coding, and for AU 1+2 its average sensitivity was 71% vs. raters' sensitivity of 92.3%. For AU 43_5, the 409
classifier performed with average sensitivity of 96%, which is similar to the sensitivity of the human 410
coders. For the lower face, the average sensitivity of the classifier for AU 25+26+16, AU 25+26+18i and 411

AU 25+26 was 70%, 88% and 91% as opposed to 63.6%, 100% and 87.5% sensitivity of the human   412
coders, respectively. Overall, our method generalized to Fascicularis monkeys with average accuracy of   413
81% for upper face and 90% for lower, as compared to the human inter-rater reliability (IRR) of 88%.   414

Altogether, the upper face KNN classifiers (Fig. 5D, top) separated well AU43_5, and had typical   415
confusions between UpperNone and AU1+2. Most lower face misclassifications (Fig. 5D, bottom) were   416
between AU25+26+16 versus AU25+26 and AU25+26+18i versus AU25+26. Characteristic outputs from   417
the system are shown in Fig. 5E.   418

*Behavioral analysis* 419

To demonstrate the potential applications of our method, we used it to analyze the facial expressions 420
produced by subject monkeys when exposed to a real life "intruder" (Fig. 2, Fig. Fig. 2-1, Fig. 2-2, Fig. 2- 421
3) (Pryluk, Shohat et al. 2020). The subject monkey was sitting behind a closed shutter, when the 422
"intruder" monkey was brought into the room ("enter" period). The shutter opened allowing the two 423
monkeys to see each other 18 times. After the last closure of the shutter, the intruder was taken out from 424
the room ("exit" period). 425

As the subject monkey was in head-immobilization, the facial expressions produced under these 426
conditions were a reduced version of the natural facial expressions that often include head and body 427
movements. To test the ethological validity of such reduced, or schematic facial expressions, we 428
determined whether they carry signal value, i.e., they are sufficient to elicit a situation-appropriate 429
reciprocation for a social partner. We found that when monkeys familiar with each other found 430
themselves in an unusual situation (open shutter) they reassured each other with reciprocal lip-smacking 431
facial expressions as shown in Fig. 2-1, Fig. 2-2 and Fig. 2-3. We verified, therefore, that multiple pairs of 432
monkeys can meaningfully communicate with each other when one of the social partners is in head 433
immobilization. 434

Statistical analysis of classification results for subject monkey B (Fig. 6A) revealed that in the presence of 435
intruder, he produced several facial expressions including UpperNone and AU25+26+18i, often 436
associated with cooing behavior. Cooing was more frequent during the "enter-exit" and open-shutter 437
periods, than during closed-shutter periods (Fig. 6B top, Fig. 6-1a left, $\chi2$, p<1e-3). Moreover, subject B 438
produced AU1+2 and AU25+26 combination more frequently during the "enter-exit" and closed-shutter 439
periods, than during the open-shutter periods (Fig. 6B bottom, Fig. 6-1a right, $\chi2$, p<1e-3). We interpret 440
this pattern as an expression of the monkey's alertness and interest in events that were signaled by 441
auditory but not visual inputs. Similarly, subject monkey D (Fig. 6C) produced action unit AU1+2 and 442
AU25+26+18i together most frequently when the intruder was visible, and on occasions when the shutter 443
was closed (intruder behind the shutter), but infrequently during the "enter-exit" periods (Fig. 6D, Fig. 6- 444
1b, $\chi2$ , p<1e-3). In a social context, this pattern is associated with the lip-smacking behavior (Parr, 445
Waller et al. 2010), representing an affiliative, appeasing social approach (Hinde and Rowell 1962). 446

*Neural analysis* 447

Finally, to validate the concept and strengthen the relevancy of automatic MaqFACS for neuroscience 448
applications, we used our method to determine whether neural activity recorded from brain regions 449
involved in facial communication (see *Materials and Methods*) is related to specific AUs (Fig. 2). Indeed, 450
neurons in the amygdala and anterior cingulate cortex (ACC) were previously shown to respond with 451
changes in firing rate during the production of facial expression (Livneh, Resnik et al. 2012). In the 452
monkeys' interaction block, responses were computed from the time when the subject monkey started 453
initiating AU25+26+18i (*Materials and Methods*). Re-analyzing the previously obtained data (Pryluk, 454
Shohat et al. 2020) showed that neurons responded before (Fig. 6E left) or after (Fig. 6E right) the 455
production of the socially meaningful AU25+26+18i. This finding supports the hypothesis that these 456
regions hold neural representations for the production of single AUs or socially meaningful AU 457
combinations. 458

13

**Discussion** 460

This work pioneers the development of an automatic system for the recognition of facial action units in 461
macaque monkeys. We based our approach on well-established methods that were successfully applied in 462
human studies of facial action units (Donato, Bartlett et al. 1999). Our system achieved high accuracy and 463
sensitivity and the results are easily interpretable in the framework of facial communication among 464
macaques. We tested our algorithm using different macaque-videos datasets in three different 465
configurations: within individual Rhesus monkeys, across individuals of Rhesus monkeys, and across 466
Rhesus and Fascicularis monkeys (generalizing across species). Performance (recognition rates) was 467
obtained for both upper and lower face and using several classification approaches, indicating that the 468
success of this method does not depend on a particular algorithm. 469

We aimed to build on commonly used and well-established tools, in order to enhance applicability and 470
ease-of-use. The pipeline of our system includes (A) alignment to predefined facial landmarks (B) 471
definition of upper and lower face ROIs (C) cropping the images to ROIs (D) generation of (difference) δ- 472
images (E) creation of lower and upper face δ-images databases (F) eigenfaces analysis, and (G) 473
classification. Our classification algorithm utilizes supervised learning, and its main challenge is the need 474
of a labeled dataset for training. Likewise, to generalize between species, a parameter fine-tuning should 475
be performed on the new species dataset. This requires a sample labeled set of the new species images. 476
The other manual operations are rather simple and not time consuming. They include a choice of neutral 477
frames and annotation of seven landmark points on a mean neutral image of a video. 478

Interestingly, unlike the *within-Rhesus* classifications, the generalization between species required a 479
larger number of components (explained variance) for classification of upper face AUs than for lower 480
face AUs. This might suggest that a separate set of parameters should be fine-tuned for each dataset and 481
ROI (lower and upper face). On the other hand, our findings show that tuning parameters for 482
generalization *within Rhesus* subjects, might suffice also for *across-Rhesus-subjects* generalization. 483
Further, and somewhat surprisingly, the *across-species* generalization performed better than *within-* and 484
*across- Rhesus-subjects* generalizations. One possible explanation is that unlike in the Rhesus dataset, the 485
Fascicularis dataset had better conditions for automatic coding, as its videos were well-controlled for 486
angle, scale, illumination, stabilization, and occlusion. This finding has an important implication, as it 487
shows that training on a large natural set of behaviors in less-controlled videos (Fig. 3-1), can be later 488
used for studying neural substrates of facial expressions in more controlled environments during 489
electrophysiology (Livneh, Resnik et al. 2012, Pryluk, Shohat et al. 2020). 490

A direct comparison to performance of human AUs-recognition systems is not straightforward. The 491
systems designed for humans are highly variable, due to differences in subjects, validation methods, the 492
number of test samples and the targeted AUs (Sariyanidi, Gunes et al. 2015). In addition, some human 493
datasets are posed, possibly exaggerating some AUs while our macaque datasets are the results of 494
spontaneous behavior. Automatic FACS achieve great accuracy (>90%) in well-controlled conditions, 495
where the facial view is strictly frontal and not occluded, the face is well illuminated, and AUs are posed 496
in a controlled manner (reviewed by Barrett, Adolphs et al. 2019). When the recordings are less 497
choreographed and the facial expressions are more spontaneous, the performance drops, (e.g. 83% in 498
Benitez-Quiroz, Srinivasan et al. 2017). Our MaqFACS recognition system performed comparably with 499

14

the human automated FACS systems despite the spontaneous nature of the macaque expressions and lack of controlled settings for the filming of Rhesus dataset. 500 501

We showed that our method can be used to add detail and depth to the analysis of neural data recorded during real-life social interactions between two macaques. This approach might pave the way toward experimental designs that capture spontaneous behaviors that may be variable across trials rather than rely on perfectly repeatable evoked responses (Krakauer, Ghazanfar et al. 2017). A departure from paradigms that dedicate less attention to the ongoing brain activity (Pryluk, Kfir et al. 2019) or internal state patterns (Mitz, Chacko et al. 2017) will increase our ability to translate experimental finding in macaques to similar finding in humans that target real-life human behavior in health and disease (Adolphs 2017). Specifically, this will allow internal emotional states and the associated neural activity that gives rise to observable behaviors to be modeled and studied across phylogeny (Anderson and Adolphs 2014). Indeed, a novel study in mice reported neural correlates of automatically-classified emotional facial expressions (Dolensek, Gehrlach et al. 2020). Finally, this system could become useful for animal-welfare assessment and monitoring (Descovich, Wathan et al. 2017, Carvalho, Gaspar et al. 2019, Descovich 2019, reviewed by McLennan, Miller et al. 2019) and in aiding the 3R framework for the refinement of experimental procedures involving all animals (Russell, Burch et al. 1959). 502 503 504 505 506 507 508 509 510 511 512 513 514 515

Given that macaques are the most commonly used non-human primate species in neuroscience, an automated system that is based on facial action units is highly desirable and will effectively complement the facial recognition systems (Loos and Ernst 2013, Freytag, Rodner et al. 2016, Crouse, Jacobs et al. 2017, Witham 2017) that address only the identity but not the behavioral state of the animal. Compared to the recently introduced method for facial expressions recognition in Rhesus macaques (Blumrosen, Hawellek et al. 2017), our system does not rely on complete stereotypical and frequent facial expressions, rather, it classifies even partial, incomplete, or ambiguous (mixed) and infrequent facial expressions, given by combination of action units. Although our system requires several manual operations, its main potential lies in automatic annotation of large datasets after tagging an example set and tuning the parameters for the relevant species or individuals. We prototyped our system on six action units in two facial regions (upper and lower face) but more advanced versions are expected to classify additional action unit combinations, spanning multiple regions of interest and tracking action units as temporal events. Further refinement of our work will likely include additional image-processing procedures, such as object tracking and segmentation, image stabilization, artifacts removal and more advanced feature extraction and classification methods. These efforts will be greatly aided by large, labeled datasets, are emerging (Murphy and Leopold 2019) to assist ongoing efforts of taking cross-species and translational neuroscience research to the next step. 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532

**References**

Adolphs, R. (2017). "How should neuroscience study emotions? By distinguishing emotion states, concepts, and experiences." Social Cognitive and Affective Neuroscience **12**(1): 24-31.

Altmann, S. A. (1962). "A field study of the sociobiology of rhesus monkeys, Macaca mulatta." Annals of the New York Academy of Sciences **102**(2): 338-435.

Anderson, D. J. and R. Adolphs (2014). "A framework for studying emotions across species." Cell **157**(1): 187-200.

Baesler, E. J. and J. K. Burgoon (1987). "Measurement and reliability of nonverbal behavior." Journal of Nonverbal Behavior **11**(4): 205-233.

Ballesta, S., et al. (2016). "Social determinants of eyeblinks in adult male macaques." Scientific reports **6**: 38686.

Barrett, L. F., et al. (2019). "Emotional expressions reconsidered: Challenges to inferring emotion from human facial movements." Psychological Science in the Public Interest **20**(1): 1-68.

Bartlett, M. S. (2001). Face Image Analysis by Unsupervised Learning, Kluwer Academic Publishers.

Bartlett, M. S., et al. (2000). Image representations for facial expression coding. Advances in neural information processing systems.

Bartlett, M. S., et al. (1996). Classifying facial action. Advances in neural information processing systems.

Benitez-Quiroz, C. F., et al. (2017). "Emotionet challenge: Recognition of facial expressions of emotion in the wild." arXiv preprint arXiv:1703.01210.

Blumrosen, G., et al. (2017). Towards Automated Recognition of Facial Expressions in Animal Models. Proceedings of the IEEE International Conference on Computer Vision.

Burrows, A. M. and T. D. Smith (2003). "Muscles of facial expression in Otolemur, with a comparison to Lemuroidea." The Anatomical Record Part A: Discoveries in Molecular, Cellular, and Evolutionary Biology: An Official Publication of the American Association of Anatomists **274**(1): 827-836.

Burrows, A. M., et al. (2009). "Facial musculature in the rhesus macaque (Macaca mulatta): evolutionary and functional contexts with comparisons to chimpanzees and humans." Journal of anatomy **215**(3): 320-334.

533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571

Burrows, A. M., et al. (2006). "Muscles of facial expression in the chimpanzee (Pan troglodytes): descriptive, comparative and phylogenetic contexts." Journal of anatomy **208**(2): 153-167.

Calder, A. J., et al. (2001). "A principal component analysis of facial expressions." Vision Research **41**(9): 1179-1208.

Carvalho, C., et al. (2019). "Ethical and Scientific Pitfalls Concerning Laboratory Research with Non-Human Primates, and Possible Solutions." Animals **9**(1): 12.

Chevalier-Skolnikoff, S. (1973). "Facial expression of emotion in nonhuman primates." Darwin and facial expression: A century of research in review: 11-89.

Crouse, D., et al. (2017). "LemurFaceID: a face recognition system to facilitate individual identification of lemurs." BMC Zoology **2**(1): 2.

Darwin, C. (1872). "The expression of emotions in men and animals."

Descovich, K. (2019). "Opportunities for refinement in neuroscience: Indicators of wellness and post-operative pain in laboratory macaques." ALTEX.

Descovich, K., et al. (2017). "Facial expression: An under-utilised tool for the assessment of welfare in mammals."

Dolensek, N., et al. (2020). "Facial expressions of emotion states and their neuronal correlates in mice." Science **368**(6486): 89-94.

Donato, G., et al. (1999). "Classifying facial actions." IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE **21**(10): 974-989.

Draper, B. A., et al. (2003). "Recognizing faces with PCA and ICA." Computer Vision and Image Understanding **91**(1-2): 115-137.

Ekman, P. (1989). The argument and evidence about universals in facial expres-sions. Handbook of social psychophysiology: 143-164.

Ekman, P. and W. V. Friesen (1976). "Measuring facial movement." Environmental psychology and nonverbal behavior **1**(1): 56-75.

572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610

17

Ekman, P. and W. V. Friesen (1986). "A new pan-cultural facial expression of emotion." <u>Motivation and emotion</u> **10**(2): 159-168.

Ekman, P. and W. V. Friesen (1988). "Who knows what about contempt: A reply to Izard and Haynes." <u>Motivation and Emotion</u> **12**(1): 17-22.

Ekman, P., et al. (2013). <u>Emotion in the human face: Guidelines for research and an integration of findings</u>, Elsevier.

Ekman, P., et al. (2002). "Facial action coding system: The manual on CD ROM." <u>A Human Face, Salt Lake City</u>: 77-254.

Ekman, P. and D. Keltner (1997). "Universal facial expressions of emotion." <u>Segerstrale U, P. Molnar P, eds. Nonverbal communication: Where nature meets culture</u>: 27-46.

Freytag, A., et al. (2016). <u>Chimpanzee faces in the wild: Log-euclidean cnns for predicting identities and attributes of primates</u>. German Conference on Pattern Recognition, Springer.

Fridlund, A. J., et al. (1987). Facial expressions of emotion. <u>Nonverbal behavior and communication, 2nd ed</u>. Hillsdale, NJ, US, Lawrence Erlbaum Associates, Inc**:** 143-223.

Friesen, E. and P. Ekman (1978). "Facial action coding system: a technique for the measurement of facial movement." <u>Palo Alto</u> **3**.

Geisser, S. (1975). "The predictive sample reuse method with applications." <u>Journal of the American statistical Association</u> **70**(350): 320-328.

He, H. and E. A. Garcia (2009). "Learning from imbalanced data." <u>IEEE Transactions on knowledge and data engineering</u> **21**(9): 1263-1284.

Hinde, R. A. and T. Rowell (1962). <u>Communication by postures and facial expressions in the rhesus monkey (Macaca mulatta)</u>. Proceedings of the Zoological Society of London, Wiley Online Library.

Jenny, A. B. and C. B. Saper (1987). "Organization of the facial nucleus and corticofacial projection in the monkey: a reconsideration of the upper motor neuron facial palsy." <u>Neurology</u> **37**(6): 930-930.

Krakauer, J. W., et al. (2017). "Neuroscience needs behavior: correcting a reductionist bias." <u>Neuron</u> **93**(3): 480-490.
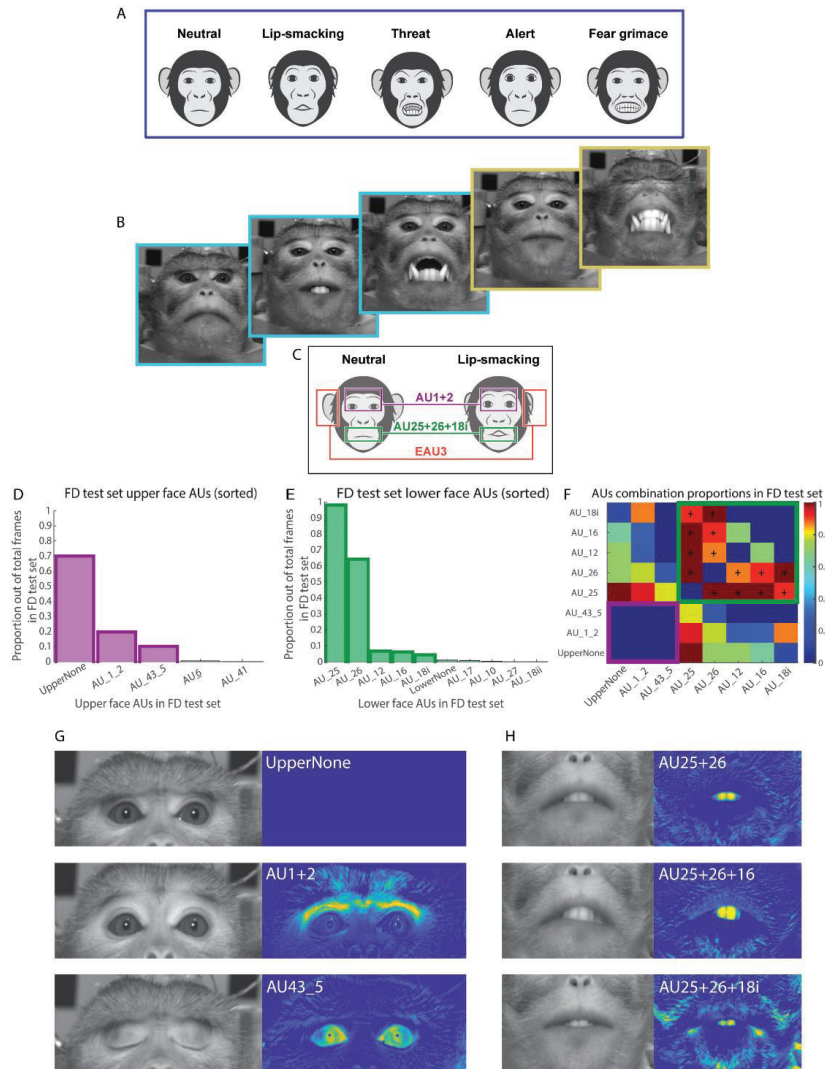
Livneh, U., et al. (2012). "Self-monitoring of social facial expressions in the primate amygdala and cingulate cortex." Proc Natl Acad Sci U S A.

Loos, A. and A. Ernst (2013). "An automated chimpanzee identification system using face detection and recognition." EURASIP Journal on Image and Video Processing **2013**(1): 49.

McLennan, K. M., et al. (2019). "Conceptual and methodological issues relating to pain assessment in mammals: The development and utilisation of pain facial expression scales." Applied Animal Behaviour Science.

Mitz, A. R., et al. (2017). "Using pupil size and heart rate to infer affective states during behavioral neurophysiology and neuropsychology experiments." Journal of Neuroscience Methods **279**: 1-12.

Morecraft, R. J., et al. (2001). "Cortical innervation of the facial nucleus in the non-human primate: a new interpretation of the effects of stroke and related subtotal brain trauma on the muscles of facial expression." Brain **124**(1): 176-208.

Mosher, C. P., et al. (2016). "Tactile Stimulation of the Face and the Production of Facial Expressions Activate Neurons in the Primate Amygdala." eNeuro **3**(5).

Murphy, A. P. and D. A. Leopold (2019). "A parameterized digital 3D model of the Rhesus macaque face for investigating the visual processing of social cues." Journal of Neuroscience Methods.

Padgett, C. and G. W. Cottrell (1997). Representing face images for emotion classification. Advances in neural information processing systems.

Panksepp, J. (2004). Affective Neuroscience: The Foundations of Human and Animal Emotions, Oxford University Press.

Pantic, M. and M. S. Bartlett (2007). Machine analysis of facial expressions. Face recognition, InTech.

Parr, L. A., et al. (2010). "Brief communication: MaqFACS: A muscle-based facial movement coding system for the rhesus macaque." Am J Phys Anthropol **143**(4): 625-630.

Pryluk, R., et al. (2019). "A Tradeoff in the Neural Code across Regions and Species." Cell **176**(3): 597-609.e518.

Pryluk, R., et al. (2020). "Shared yet dissociable neural codes across eye gaze, valence and expectation." Nature **586**(7827): 95-100.

650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689

Russell, W. M. S., et al. (1959). <u>The principles of humane experimental technique</u>, Methuen London. 690

691
Sariyanidi, E., et al. (2015). "Automatic Analysis of Facial Affect: A Survey of Registration, 692
Representation, and Recognition." <u>IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE</u> 693
<u>INTELLIGENCE</u> **37**(6): 1113-1133. 694

695
Tukey, J. (1958). "Bias and confidence in not quite large samples." <u>Ann. Math. Statist.</u> **29**: 614. 696

697
Vick, S.-J., et al. (2007). "A Cross-species Comparison of Facial Morphology and Movement in Humans 698
and Chimpanzees Using the Facial Action Coding System (FACS)." <u>Journal of Nonverbal Behavior</u> **31**(1): 1- 699
20. 700

701
Waller, B., et al. (2020). "Measuring the evolution of facial 'expression'using multi-species FACS." 702
<u>Neuroscience & Biobehavioral Reviews</u>. 703

704
Welt, C. and J. H. Abbs (1990). "Musculotopic organization of the facial motor nucleus in Macaca 705
fascicularis: a morphometric and retrograde tracing study with cholera toxin B-HRP." <u>Journal of</u> 706
<u>Comparative Neurology</u> **291**(4): 621-636. 707

708
Witham, C. L. (2017). "Automated face recognition of rhesus macaques." <u>Journal of Neuroscience</u> 709
<u>Methods</u>. 710

711

712

713

20

**Figure 1. Motivation for using automatic MaqFACS to analyze facial expressions**
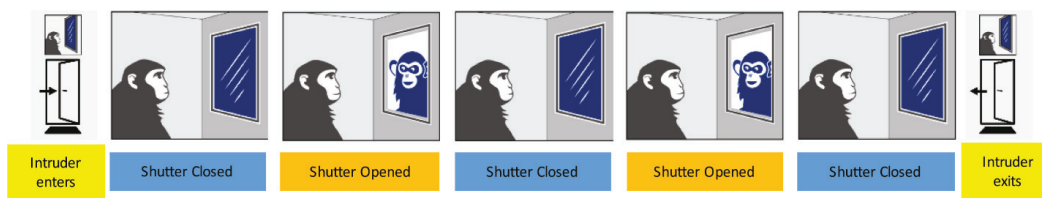
714

715



716

717

A. The stereotypical facial expressions in macaque monkeys include the 'neutral', 'lip-smacking', 'threat', 'alert' and 'fear grimace' expressions (Altmann 1962, Hinde and Rowell 1962).

718
719

B. Some of the facial expressions that monkeys produce during the experiments that require head immobilization match the stereotypical expressions produced during natural behaviors (for example, see the three images with blue frames on the left, correspond to the neutral, lip-smacking and threat expressions). We have also observed facial expressions that were less frequently described in the literature (two images with yellow frames on the right).

720
721
722
723
724

21

C. A comparison between the neutral and lip-smacking facial expression shows that the lip-smacking example contains AU1+2 (Brow Raiser) in the upper face, AU25+26+18i (Lips part, Jaw drop and True Pucker) in the lower face, and EAU3 (Ear Flattener) in the ear region. 725 726 727

D. The proportion of each upper face AU in the Fascicularis data (FD) test set. Bars with the solid outline (first three highest bars) represent the most frequent AUs, which were chosen for the analysis in this work. 728 729 730

E. Same as (D) but for lower face. First five most frequent AUs were chosen for the analysis. 731

F. Proportion matrix of AU combinations in the FD test set, for the most frequent AUs. Cells inside the magenta (bottom left) and green frames (top right) represent the combinations of upper face and lower face AUs, correspondingly. AUs that frequently occurred in combination with other AUs (in the upper face or the lower face, separately) are denoted by "+". Cell values were calculated as the ratio between the number of frames containing the combination of the two AUs and the total frames number containing the less frequent AU. 732 733 734 735 736 737

G. Left: images of upper face AUs from the FD test set. UpperNone: no coded action in the upper face. AU1+2: Brow raiser. AU43_5: Eye closure. Right: the difference of the images from the neutral face image. 738 739 740

H. Same as (G) but for lower face. AU25+26: Lips part and Jaw drop. AU25+26+16: Lips part, Jaw drop and Lower lip depressor. AU25+26+18i: Lips part, Jaw drop and True Pucker. 741 742

743

744

22

**Figure 2. Monkey-intruder behavioral paradigm**

| Intruder enters | Shutter Closed | Shutter Opened | Shutter Closed | Shutter Opened | Shutter Closed | Intruder exits |

Monkey Intruder Block: The subject monkey sitting behind a closed shutter. The intruder monkey is brought into the room and seated behind the shutter, which remains closed. The shutter opens and closes 18 times, and the monkeys are able to see each other while it is open. The subject monkey could not see any part of the intruder unless the shutter is open. At the end of the block, the shutter closes and the intruder monkey is taken out from the room.

For examples of monkey interactions, see extended Figures 2-1, 2-2 and 2-3.

757

758



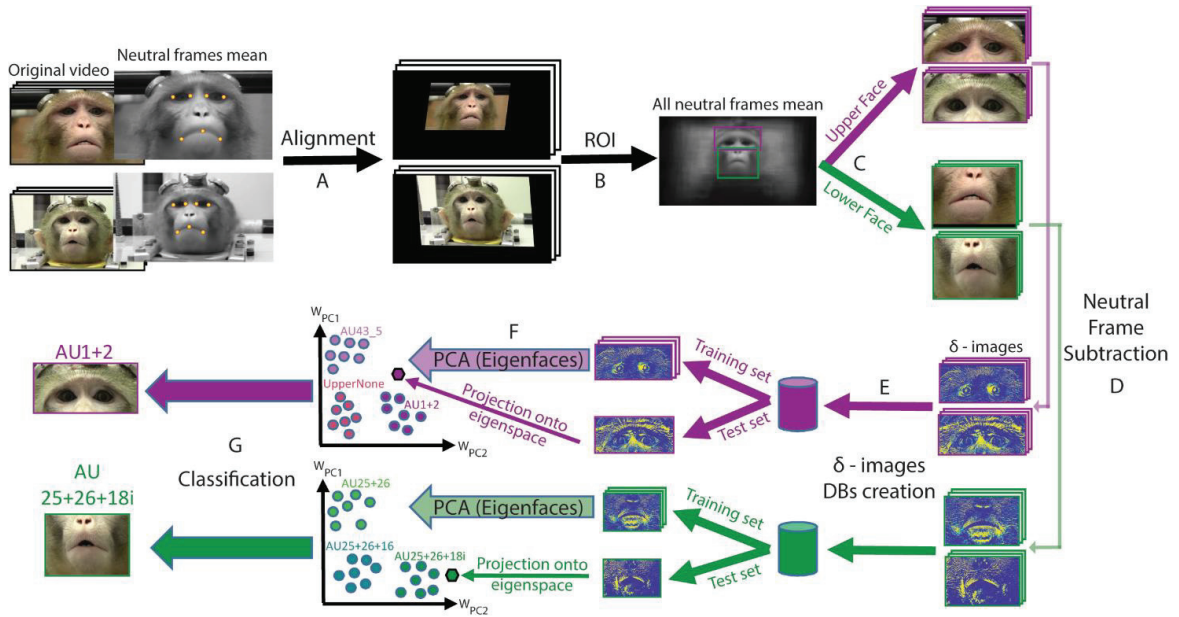Alignment of frames from the original video stream (example of two videos from two different Rhesus Dataset (RD) monkeys. Seven landmark points were manually selected on the mean of all neutral frames of each video. In the next step, these points were mapped to corresponding predefined positions (reference landmarks, common for all videos). The resulting affine transformation for each video was then applied to all its frames. For more examples, see extended Figure 3-1.
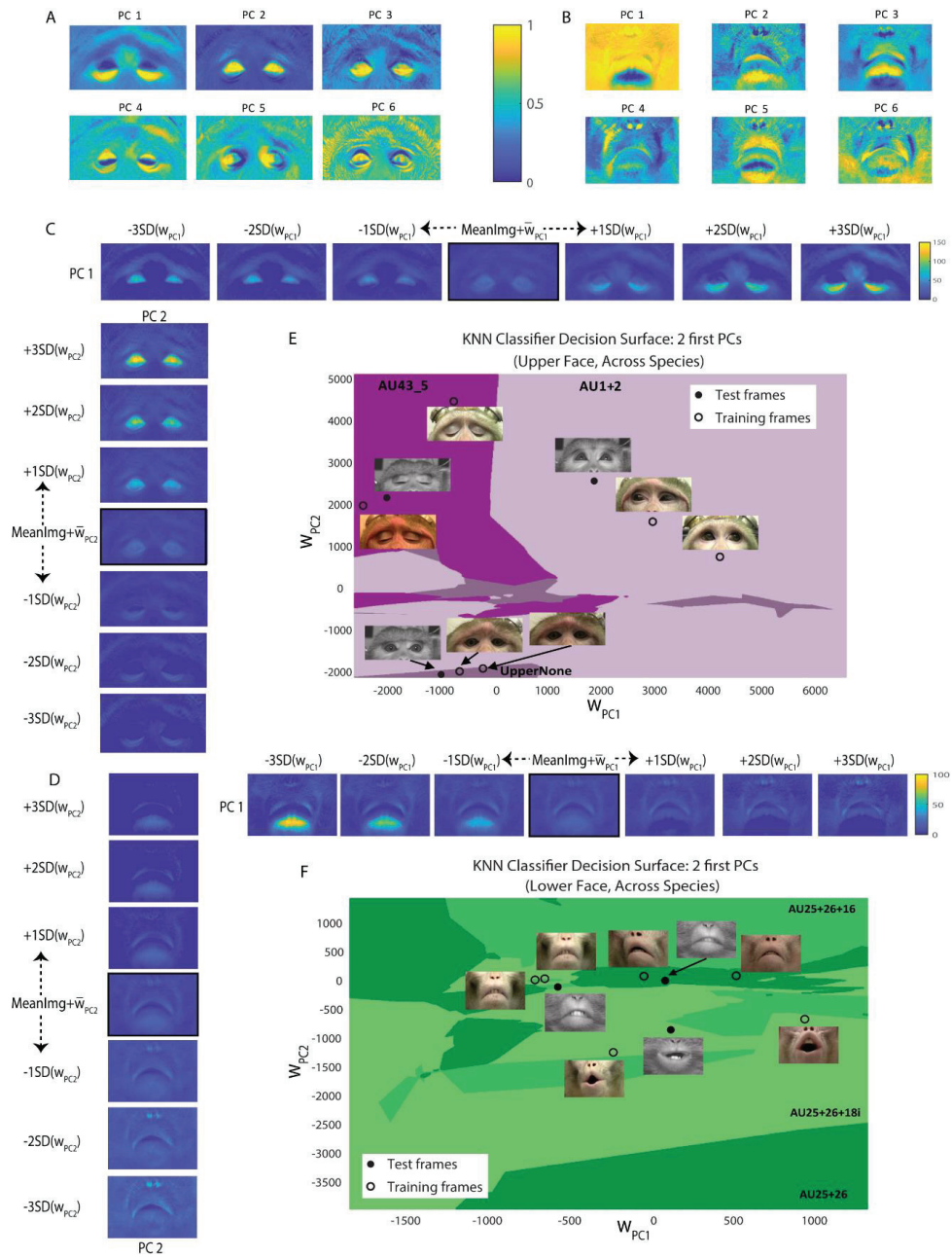
760
761
762
763
764

A.  Manual definition of upper face and lower face ROIs on the mean of all neutral frames. Magenta: upper face ROI, green: lower face ROI. The "All neutral frames mean" image in this scheme was calculated from all RD videos.

765
766
767

B.  Cropping of all the frames according to upper face and lower face ROIs.

768

C.  Generation of δ-images by subtracting the optimal neutral frame of each video from all its frames. The contrast and the color map of the gray scale images were adjusted for a better representation.

769
770

D.  Construction of lower face and upper face δ-images databases, consisting of 2-dimensional matrices where each row corresponds to one image.

771
772

E.  Eigenfaces extraction from the training images and projection of the training and test images onto the eigenspace (following the desired training and test sets construction). $W_{PC1}$ and $W_{PC2}$ denote the weights of PC1 and PC2, correspondingly.

773
774
775

F.  Classification of the testing images to upper face and lower face AUs. KNN (and SVM) classification was applied based on the distances between the testing and the training images in the eigenspace.
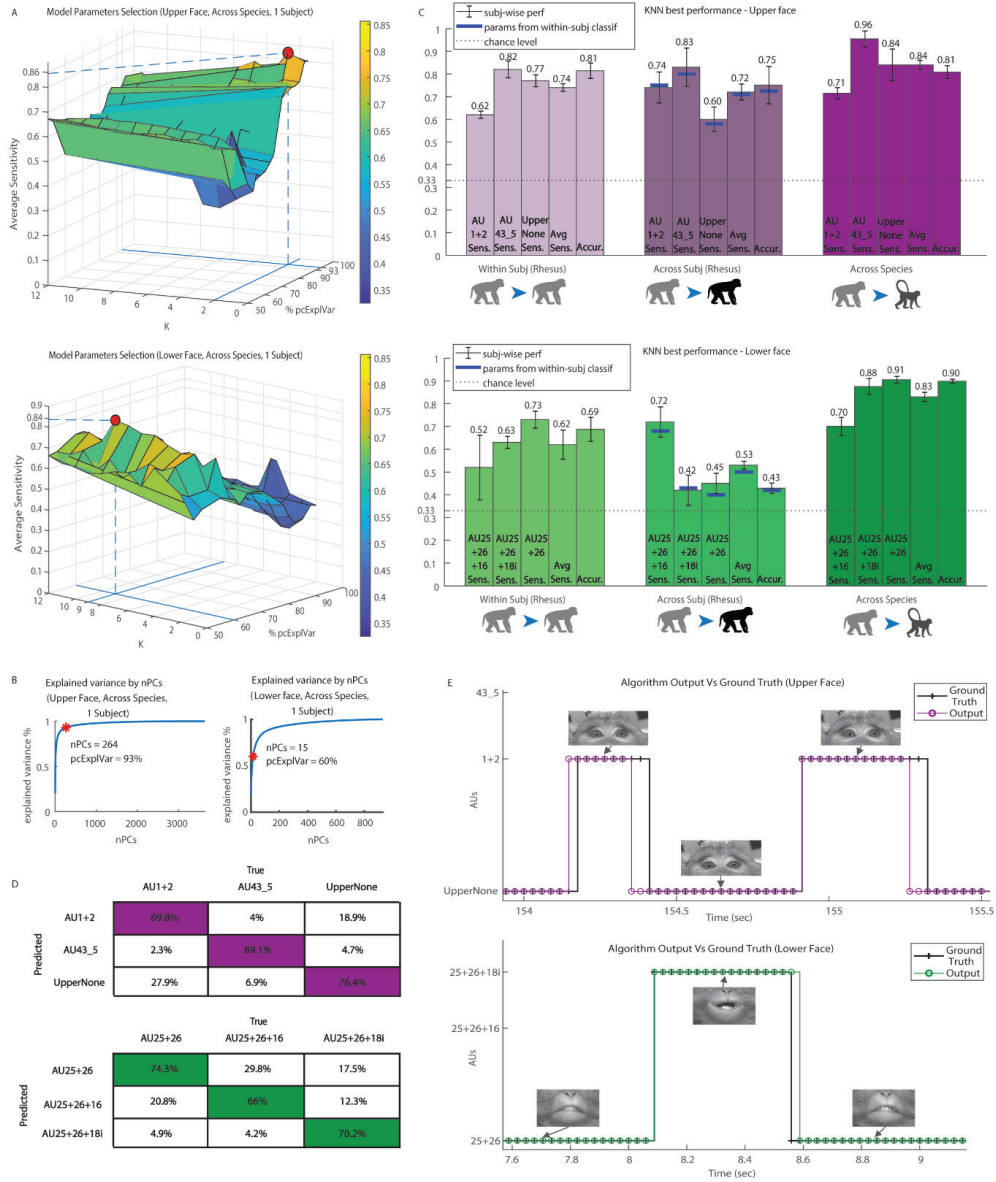
776
777

778

24

**Figure 4: Eigenfaces analysis**

A. Example of eigenfaces: six first eigenfaces (PCs) of one of the upper face training sets, containing all five Rhesus subjects from RD. The grayscale values were normalized to 0-1 range and the image contrast and color map were adjusted for a better representation. The color bar corresponds to pixel grayscale values.   783 784 785

B. Same as (A) but for lower face.   786

C. Example of the information coded by the first two eigenfaces.   787

Top: the image sequence demonstrates the first eigenface from (A), added to the mean image (*MeanImg*) and varied. The middle image is the mean image of the training set (described in A), with the first eigenface added after being weighted by its mean weight ($\overline{w}_{PC1}$). In each sequence, the weights were varied from -3SD to +3SD from the mean weight, and the weighted PC was then added to the mean image of the training set. This procedure resulted in a different facial image for each 1SD step. The images in the sequence are ordered from left to right: the first image contains the variation by -3SD (i.e. PC1 weighted by -3SD of its weights and added to the middle image), and the last one is the variation by +3SD.   788 789 790 791 792 793 794

Bottom: same as top but for the second eigenface (PC2). The image sequence is ordered from bottom to top.   795 796

The grayscale values were normalized to 0-150 range and the image contrast and color map were adjusted for a better representation. The color bar corresponds to pixel grayscale values, and is mutual for both top and bottom schemes.   797 798 799

D. Same as (C) but for lower face and with grayscale normalization to range 0-100.   800

E. Example of decision surface for upper face KNN classifier, trained for generalization *across species*. The training set is the one described in (A) and the test set is Fascicularis monkey D frames from FD. The decision surface is presented along the first two dimensions – weights of PC1 and PC2 ($w_{PC1}$ and $w_{PC2}$, correspondingly). Each colored region denotes one of the three upper face AU classes. The frames in color are training set images and the gray-scaled ones are from the test set. The classification decision is based on the test frames' proximity to samples of a certain class in this compressed subspace. For better illustration, the images shown here are frames after alignment, but before the neutral frame subtraction.   801 802 803 804 805 806 807

F. Same as (E) but for the lower face and Fascicularis monkey B from FD test set.   808

809

810

26

**Figure 5: Results – parameters selection and model performance**



A.  Top: example of parameter selection for upper face KNN classifier, trained for generalization *across* 814
*species*. The training set in the example is the one described in fig. 4(A), the test set is monkey D frames 815
from FD and the distance metric is set to be Euclidean. The surface represents the performance of KNN 816
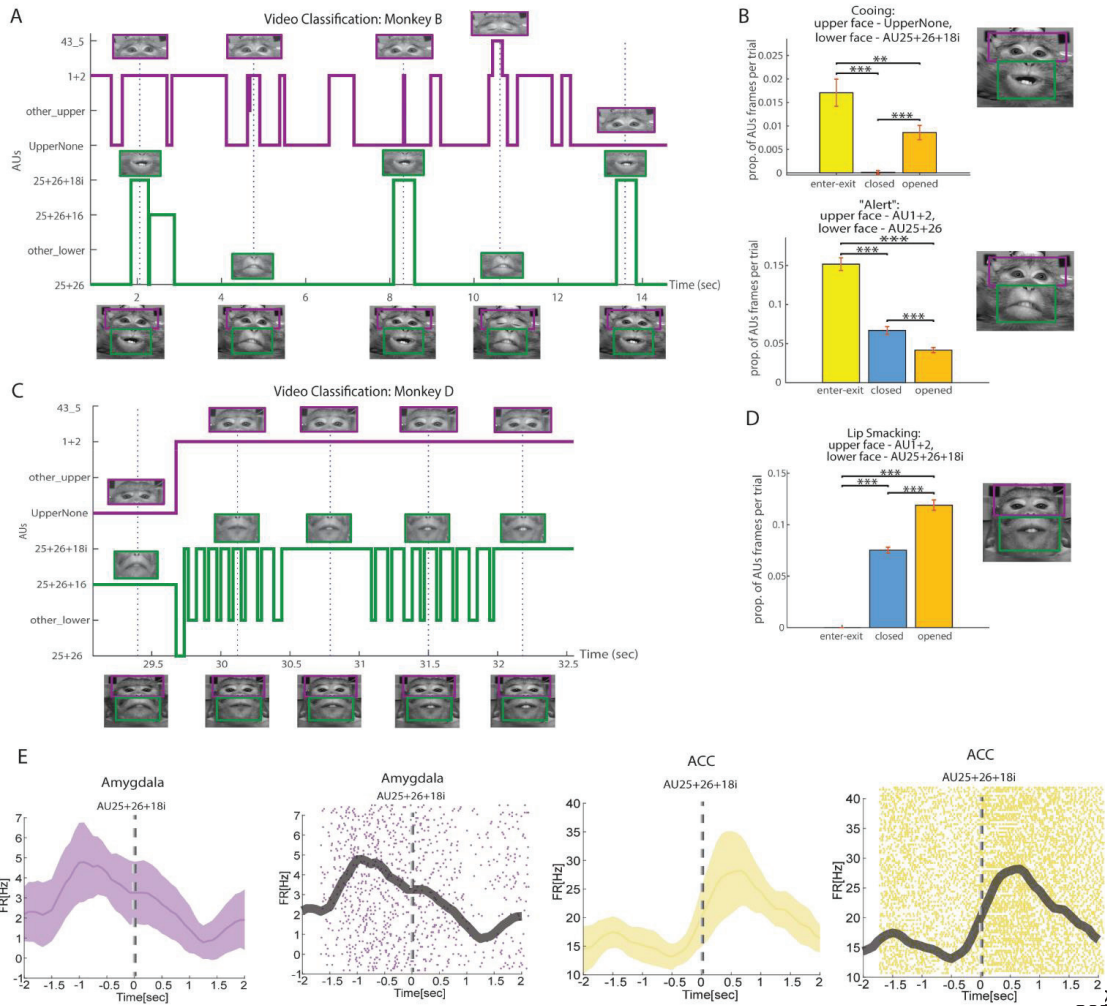
classifiers with two parameters varied: k (number of nearest neighbors, varied from 1 to 12), and the     817
percentage of the training set variance explained by the eigenfaces ("pcExplVar", varied from 50% to     818
95%). Z-axis is the average sensitivity value of each model (i.e. average of the sensitivity values for the     819
classification of three upper face AUs). The red dot denotes the highest point on the surface and hence the     820
parameters yielding the best performance. With the selected parameters k=2 and pcExplVar = 93% the     821
model average sensitivity value is 0.86.     822

Bottom: Same as on top but for the lower face. The training set is one of the lower face training sets,     823
containing all five Rhesus subjects from RD, and the test set is monkey D frames from FD. The distance     824
metric is set to be Euclidean. The selected model has the average sensitivity of 0.84 with the parameters:     825
k=9 and pcExplVar = 60%.     826

B. The curves demonstrate the number of the eigenfaces that should be used to cumulatively capture a given     827
percentage of the dataset variance. The red asterisk denotes the pcExplVar parameter value selected in (A).     828

Left: the curve corresponds to the dataset described in (A) top. To express 93% of the dataset variance, at     829
least 264 vectors (eigenfaces) should span the eigenspace. Right: same as left but regarding (A) bottom. To     830
express 60% of the dataset variance, at least 15 vectors (eigenfaces) should span the eigenspace.     831

C. Best performance of KNN classification for each generalization type. Each bar group contains five bars     832
(from left to right): three bars describing the classifier's sensitivity for single AUs; sensitivity averaged for     833
three classified AUs; and the total accuracy of the classifier. The mean and the error are calculated     834
regarding the recognition performance on a new subject. The horizontal dashed line denotes the chance     835
level.     836

The first bar group demonstrates the results for generalization of the classification within the same Rhesus     837
subject (*Within Subject (Rhesus)*: training on videos of a subject and testing on a new video of the same     838
subject).     839

The second group shows the generalization performance of a classifier to new Rhesus subjects (*Across     840
Subjects (Rhesus)*: training on videos from several subjects and testing on videos of a new subject). The     841
blue lines denote the performance of the classifier *across subjects* using the parameters selected in *Within     842
Subject (Rhesus)* case.     843

The third group displays the generalization performance to new Fascicularis subjects (*Across Species*:     844
training on videos from several Rhesus subjects and testing on videos of a new Fascicularis subject). In this     845
case, the parameters should be tuned for each Fascicularis subject, and the results are the mean     846
performance of two parameter sets (for the two Fascicularis subjects).     847

Top: performance for upper face. Bottom: performance for lower face.     848

D. Averaged confusion matrices of the KNN best performance results (of the three cases presented in (C)).     849
The columns in each matrix represent the true labels, and the rows stand for the predicted labels.     850

Top: upper face confusion matrices. Bottom: lower face confusion matrices.     851

For Confusion matrix of interrater variability, see extended Figure 5-1.     852

E. Example of the KNN classification performance demonstrating correctly recognized frames along with     853
some recognition errors. Each data point denotes a frame in a video. The classified AUs (magenta and     854
green lines) are shown in comparison to the ground truth labels (the black lines). Video time is displayed in     855
the X-axis. Sample frames of the original video stream (after alignment and ROI cropping) are shown     856
above the lines. The video for the example is taken from FD.     857

Top: output example for upper face video. Bottom: output example for lower face video.     858

    859

28

**Figure 6: Examples of the Method Applications**



A. Example of the final system output for monkey B from FD. Classification labels are presented on the Y-axis, while the frame time of the video-stream is on the X. "Other_upper" and "other_lower" labels are for video frames that were not part of the task of the classifier but exist in the original video and were labeled manually. Frames of the original video (with no preprocessing) are shown on the bottom and the dashed lines denote their corresponding timing. The magenta and green lines demonstrate the outputs from the upper face and lower face algorithms, respectively. Images above the output lines exhibit the frames as they were processed in the algorithm, after alignment and ROI cropping. The estimated locations of the ROIs, comprising the full facial expressions, are illustrated in frames on the bottom by magenta and green rectangles (the positions are not precise since the original images on the bottom are not aligned).

863
864
865
866
867
868
869
870
871
872

B. Facial expressions analysis following frames classification. Bars demonstrate the proportion of a specific facial configuration monkey B (from FD) elicited during one block of the experiment described in fig. 2. This value is calculated as the ratio between frames containing the AUs combination and the total frames, per trial. Yellow bars denote the block part when the intruder monkey enters and exists the room, the blue one is for phases with the closed shutter (after the first shutter opening and before its last closure), and the orange bars stand for periods of open shutter. An example image of the analyzed expression is shown on the right (taken from the examples in B).

Top: proportions of cooing facial expression events comprised of UpperNone AU for the upper face and AU25+26+18i for the lower face. Bottom: same as in top, but for "alert" facial expression – AU1+2 and AU25+25 in the upper face and lower face, correspondingly.

For analysis following classification by human coders, see extended Figure 6-1a.

** represents p<1e-2, *** represents p<1e-3

C. Same as (A) but for monkey D from FD.

D. Same as (B) but for monkey D from FD and lip-smacking facial expression with upper face AU1+2 and lower face AU25+26+18i.

For analysis following classification by human coders, see extended Figure 6-1b.

E. PSTHs and raster plots of one neuron in the amygdala and one in the ACC, temporally locked to the socially-associated AU25+26+18i, during monkey intruder block.

**Table 1: Data under-sampling (RD)**

Upper Face

|  | AU1+2 | AU43_5 | UpperNone | Undersampled per class | Total balanced training set |
|---|---|---|---|---|---|
| #frames | 1213 | ~19,500 | ~150,000 | 1213 | 3639 |

Lower Face

|  | AU25+26+16 | AU25+26+18i | AU25+26 | Undersampled per class | Total balanced training set |
|---|---|---|---|---|---|
| #frames | 310 | ~15,000 | ~15,000 | 310 | 930 |

In the upper face, the smallest category was AU1+2 with only 1213 frames (in total, from all RD subjects). On the contrary, AU43_5 category had around 19,500 frames (after eliminating RD AU45 frames due to time synchronization errors), and UpperNone class included over 150,000 images. Consequently, balanced training sets were generated each including all the AU1+2 frames, and randomly selected 1213 frames from AU43_5 along with 1213 randomly selected UpperNone frames. Therefore, the upper face balanced training sets were comprised of 3639 frames each. The same was done for the lower face, where the smallest category was AU25+26+16 with only 310 frames. Categories AU25+26+18i and AU25+26 contained over 15,000 images each. Accordingly, each lower face balanced training set included 930 frames.
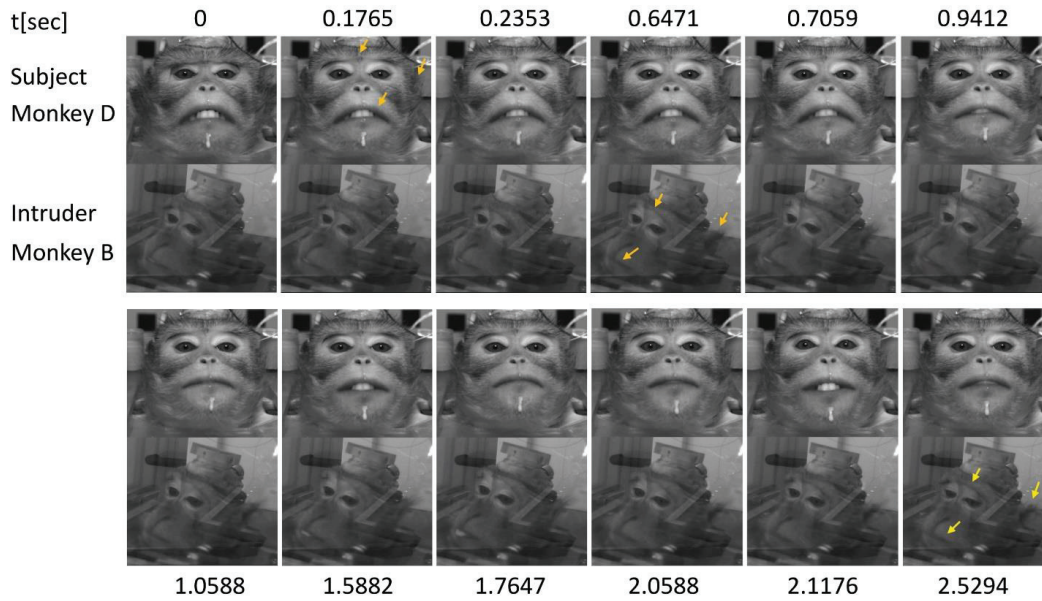
898
899
900
901
902
903
904
905

906

**Figure 2-1: Lip-smacking interactions**

| t[sec] | 0 | 0.1765 | 0.2353 | 0.6471 | 0.7059 | 0.9412 |

Subject Monkey D

Intruder Monkey B

| 1.0588 | 1.5882 | 1.7647 | 2.0588 | 2.1176 | 2.5294 |

Examples of dynamics and progression of lip-smacking interactions captured during the monkey-intruder experiment, where the subject monkey is the first to initiate the movement. Each sequence demonstrates sample frames of the Fascicularis subject D with his head fixed (first row), along with the corresponding frames of the intruder Fascicularis monkey (second row). The subject monkey D was filmed using the facial camera (Materials and Methods). The intruder monkey was filmed using another monitoring camera, from the direction of the subject monkey and through the opened shutter (hence the reflections on the screen). The time presented relative to the first frame in the sequence, which starts with a neutral expression of the subject monkey. Yellow arrows indicate the change in the movement of brows, ears and lips at the onset of the lip-smacking movement (for the subject and the intruder monkeys) and the offset of the movement (for the intruder monkey).
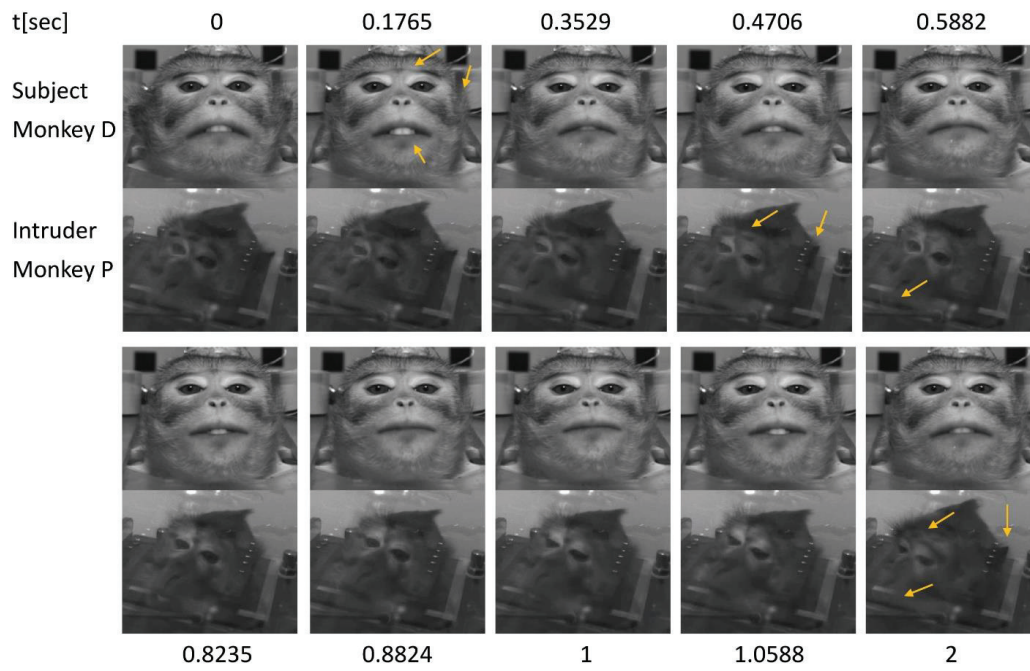
In the example: sequence with intruder monkey B.

32

**Figure 2-2: Lip-smacking interactions**



Same setup as in Fig. 2-1, but with intruder monkey P.

**Figure 2-3: Lip-smacking interactions**



Same setup as in Fig. 2-1, but with intruder monkey N.

930
931
932
933
934
935

**Figure 3-1: Motivation for alignment**

a



b



c

Seven reference landmark points (yellow, predefined and common for all videos) displayed on sample neutral frames of original video streams.

A. Sample neutral frames from five different videos of each of the five Rhesus monkeys (K, L, M, Q, R).
B. Sample neutral frames from one video of Rhesus monkey K.
C. Sample neutral frames of the two Fascicularis monkeys (D and B).

**Figure 5-1: Confusion matrix: interrater variability**

|  | AU43_5 | Upper None | AU1+2 | AU1+2 +43_5 | Other Upper | AU25 +26 | AU25+ 26+18i | Other Lower | AU25+26 +16 |
|---|---|---|---|---|---|---|---|---|---|
| AU43_5 | 96% | <1% | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Upper None | <1% | 81% | 6.2% | 0 | 6.1% | 0 | 0 | 0 | 0 |
| AU1+2 | 3% | 18.9% | 92.3% | 12.5% | 6.1% | 0 | 0 | 0 | 0 |
| AU1+2+43_5 | <1% | 0 | <1% | 87.5% | 0 | 0 | 0 | 0 | 0 |
| Other Upper | 0 | <1% | 1% | 0 | 87.8% | 0 | 0 | 0 | 0 |
| AU25+26 | 0 | 0 | 0 | 0 | 0 | 87.5% | 0 | 3.1% | 2.6% |
| AU25+26+18i | 0 | 0 | 0 | 0 | 0 | 1.4% | 100% | <1% | 0 |
| Other Lower | 0 | 0 | 0 | 0 | 0 | 10% | 0 | 95.5% | 33.8% |
| AU25+26+16 | 0 | 0 | 0 | 0 | 0 | 1.1% | 0 | 1.3% | 63.6% |

Confusion matrix for the interrater variability between two experienced human coders, for a video from FD. "Other Upper" and "Other Lower" represent all the upper-face and lower-face labels which were not part of the task of the automatic classifier.

a

Cooing:
upper face - neutral,
lower face - 25+26+18i



"Alert":
upper face - AU 1+2,
lower face - AU 25+26



b

Lip Smacking:
upper face - AU 1+2,
lower face - 25+26+18i

5C,E but deduced from ground-truth labels.

a. Monkey B from FD
b. Monkey D from FD

**Extended Data 1**                                                                              999

The archive "autoMaqFACS_code.zip" contains Matlab code for autoMaqFACS classification.        1000

**A**

Neutral    Lip-smacking    Threat    Alert    Fear grimace

**B**

**C**

Neutral      Lip-smacking

AU1+2

AU25+26+18i

EAU3

**D** FD test set upper face AUs (sorted)

Proportion out of total frames in FD test set

UpperNone   AU_1_2   AU_43_5   AU6   AU41

Upper face AUs in FD test set

**E** FD test set lower face AUs (sorted)

Proportion out of total frames in FD test set

AU_25   AU_26   AU_12   AU_16   AU_18i   LowerNone   AU_17   AU_10   AU_27   AU_18ii

Lower face AUs in FD test set

**F** AUs combination proportions in FD test set

AU_18i
AU_16
AU_12
AU_26
AU_25
AU_43_5
AU_1_2
UpperNone

UpperNone   AU_1_2   AU_43_5   AU_25   AU_26   AU_12   AU_16   AU_18i

**G**

UpperNone

AU1+2

AU43_5

**H**

AU25+26

AU25+26+16

AU25+26+18i

| Intruder enters | Shutter Closed | Shutter Opened | Shutter Closed | Shutter Opened | Shutter Closed | Intruder exits |

A  Model Parameters Selection (Upper Face, Across Species, 1 Subject)

Average Sensitivity

K
% pcExplVar

Model Parameters Selection (Lower Face, Across Species, 1 Subject)

Average Sensitivity

K
% pcExplVar

C  KNN best performance - Upper face

subj-wise perf
params from within-subj classif
chance level

0.62 | 0.82 | 0.77 | 0.74 | 0.81
0.74 | 0.83 | 0.60 | 0.72 | 0.75
0.71 | 0.96 | 0.84 | 0.84 | 0.81

AU 1+2 Sens. | AU 43_5 Sens. | Upper None Sens. | Avg Sens. | Accur.

Within Subj (Rhesus)   Across Subj (Rhesus)   Across Species

KNN best performance - Lower face

subj-wise perf
params from within-subj classif
chance level

0.52 | 0.63 | 0.73 | 0.62 | 0.69
0.72 | 0.42 | 0.45 | 0.53 | 0.43
0.70 | 0.88 | 0.91 | 0.83 | 0.90

AU25 +26 +16 Sens. | AU25 +26 +18i Sens. | AU25 +26 Sens. | Avg Sens. | Accur.

Within Subj (Rhesus)   Across Subj (Rhesus)   Across Species

B  Explained variance by nPCs (Upper Face, Across Species, 1 Subject)

explained variance %
nPCs = 264
pcExplVar = 93%
nPCs

Explained variance by nPCs (Lower face, Across Species, 1 Subject)

explained variance %
nPCs = 15
pcExplVar = 60%
nPCs

D

| | True | | |
| | AU1+2 | AU43_5 | UpperNone |
| AU1+2 | 69.8% | 4% | 18.9% |
| AU43_5 | 2.3% | 89.1% | 4.7% |
| UpperNone | 27.9% | 6.9% | 76.4% |

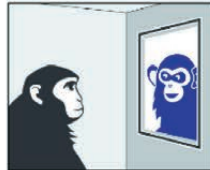Predicted

| | True | | |
| | AU25+26 | AU25+26+16 | AU25+26+18i |
| AU25+26 | 74.3% | 29.8% | 17.5% |
| AU25+26+16 | 20.8% | 66% | 12.3% |
| AU25+26+18i | 4.9% | 4.2% | 70.2% |

Predicted

E  Algorithm Output Vs Ground Truth (Upper Face)

Ground Truth
Output

43_5
1+2
AUs
UpperNone

154   154.5   155   155.5
Time (sec)

Algorithm Output Vs Ground Truth (Lower Face)

Ground Truth
Output

25+26+18i
25+26+16
AUs
25+26

7.6   7.8   8   8.2   8.4   8.6   8.8   9
Time (sec)

A — Video Classification: Monkey B

B — Cooing: upper face – UpperNone, lower face – AU25+26+18i

"Alert": upper face – AU1+2, lower face – AU25+26

C — Video Classification: Monkey D

D — Lip Smacking: upper face – AU1+2, lower face – AU25+26+18i

E — Amygdala    Amygdala    ACC    ACC