1

2

3

4

5

6

7

8

9                                        Research Article

10        Perceptual mechanism underlying gaze guidance in chimpanzees and humans

11

12                                   Fumihiro Kano (1, 2)

13                                   Masaki Tomonaga (1)

14                        1.   Primate Research Institute, Kyoto University

15                            2.   Japan Society for Promotion of Science

16

17                                fkanou@pri.kyoto-u.ac.jp

18                                   +81-80-6902-5013

19

20    Abstract

21    Previous studies comparing eye movements between humans and their closest relatives,

22    chimpanzees, have revealed similarities and differences between the species in terms of

23    where individuals fixate their gaze during free viewing of a naturalistic scene, including

24    social stimuli (e.g. body and face). However, those results were somewhat confounded

25    by the fact that gaze behavior is influenced by low-level stimulus properties (e.g., color

26    and form), and by high-level processes such as social sensitivity and knowledge about

27    the scene. Given the known perceptual and cognitive similarities between chimpanzees

28    and humans, it is expected that such low-level effects do not play a critical role in

29    explaining the high-level similarities and differences between the species. However,

30    there is no quantitative evidence to support this assumption. To estimate the effect of

31    local stimulus saliency on such eye-movement patterns, this study used a

32    well-established bottom-up saliency model. In addition, to elucidate the cues that the

33    viewers use to guide their gaze, we presented scenes in which we had manipulated

34    various stimulus properties. As expected, the saliency model did not fully predict the

35    fixation patterns actually observed in chimpanzees and humans. In addition, both

36    species used multiple cues to fixate socially significant areas such as the face. There

37    was no evidence suggesting any differences between chimpanzees and humans in their

38    responses to low-level saliency. Therefore, this study found a substantial amount of

39    similarity in the perceptual mechanisms underlying gaze guidance in chimpanzees and

40    humans, and thereby offers a foundation for direct comparisons between them.

41    Key words: chimpanzees, eye-tracking, face, picture perception, saliency

42

43    Introduction

44    Eye-tracking methodology in human and nonhuman primates has been used for over 50

45    years (Fuchs, 1967; Yarbus, 1967). Eye-movement patterns of nonhuman primates show

46    a significant degree of similarity with those of humans under similar experimental

47    conditions (Kano and Tomonaga, 2009; Shepherd, Steckenfinger, Hasson, and

48    Ghazanfar, 2010). Comparative studies of human and nonhuman primates have directly

49    compared the species in order to clarify both similarities and differences in their

50    eye-movement characteristics (Dahl, Wallraven, Bulthoff, and Logothetis, 2009;

51    Gothard, Brooks, and Peterson, 2009; Guo, Robertson, Mahmoodi, Tadmor, and Young,

52    2003; Kano and Tomonaga, 2009; Keating and Keating, 1982; Nahm, Perret, Amaral,

53    and Albright, 1997; Shepherd et al., 2010). Those similarities and differences have been

54    an important source of information for the study of the evolution of visual behavior,

55    social perception, and high-level cognition (Kano and Tomonaga, 2009; Shepherd et al.,

56    2010). Although apes have been essential for this comparative approach, their

57    eye-movement characteristics are largely unknown compared to those of the

58    well-studied macaque species.

59        Recently, eye-tracking studies in chimpanzees (*Pan troglodytes*), the species

60    most closely related to humans, have been reported (Hattori, Kano, and Tomonaga,

61    2010; Hirata, Fuwa, Sugama, Kusunoki, and Fujita, 2010; Kano and Tomonaga, 2009,

62    2010). Those studies have presented naturalistic images of scenes (including faces,

63    bodies, etc.) to chimpanzees and humans under free-viewing conditions and compared

64    their fixation patterns under similar experimental conditions. There are several

65    advantages of comparing chimpanzees and humans for a free-viewing task. First,

66    chimpanzees are the species most closely related to humans and are known to have

67  similar perceptual mechanisms (Matsuno, Kawai, and Matsuzawa, 2004; Matsuzawa,

68  1985, 1990; Tomonaga and Matsuzawa, 1992). Second, the demands of a free-viewing

69  task are small; the participants of both species need only look at the stimuli

70  spontaneously and are not trained to solve particular problems using their eye

71  movements. Third, for the same reason, a free-viewing task is relatively independent of

72  the effect of reward or training. Therefore, we are able to efficiently and directly

73  compare the species, find both similarities and differences between them, and discuss

74  the extent to which chimpanzees and humans are similar and different in their

75  perception and cognition.

76      In the previous study comparing chimpanzees and humans in a free-viewing

77  task, it was found that the species were very similar in terms of where to fixate (i.e.

78  scanpath similarity). For example, when presented with a scene including an entire body

79  of a chimpanzee, a human, or another animal, both chimpanzees and humans

80  concentrated fixations on the body, especially the face, rather than on the background.

81  In addition, both species fixated on the face immediately after the image presentation

82  (within the first few fixations). However, those responses differed quantitatively

83  between the species; humans showed a higher proportion of face fixations than did

84  chimpanzees. There seem to be several functional reasons for those similarities and

85  differences between the species. First, faces are the most important source of social

86  information (such as individuality and emotions) for both chimpanzees and humans

87  (Chevalier-Skolnikoff, 1973; Parr, Dove, and Hopkins, 1998), and thus frequent

88  inspection and immediate detection of facial characteristics may benefit them by

89  enabling them to obtain such information efficiently. Second, humans have a specific

90  form of facial communication; humans often engage in lengthy face-to-face

91 communication, accompanied by intense eye contact (Argyle and Cook, 1976).

92 Therefore, more frequent inspection of faces may benefit humans more specifically than

93 chimpanzees in the context of their own form of facial communication.

94        There seem to be several factors that determine such similarities and

95 differences. These include, for example, the perception of low-level visual properties

96 (e.g. color, form), the perception of bodies and faces, and knowledge about the scenes

97 (which the viewers had obtained through daily lives or experimental instructions).

98 Previous studies using forced-choice discrimination paradigms have found that the

99 perceptions of low-level visual properties involving color (Matsuno et al., 2004;

100 Matsuzawa, 1985), form (Matsuzawa, 1990; Tomonaga and Matsuzawa, 1992) are

101 largely similar between chimpanzees and humans. In addition, the mechanisms

102 involving advanced social perceptions involving faces (Parr et al., 1998; Parr, Hecht,

103 Barks, Preuss, and Votaw, 2009; Tomonaga, 2007, 2010; Tomonaga and Imura, 2009)

104 and bodies (Tomonaga and Imura, 2008) are also similar between the species.

105        Because of these similarities between chimpanzees and humans, it is expected

106 that the influence of low-level stimulus properties on their eye-movement patterns

107 appears similarly in the two species and does not play a critical role in explaining for

108 the overall similarities and differences between the species. However, there are no

109 quantitative data to support this assumption. It is important to separate low-level from

110 higher-level influences on eye-movement patterns in order to provide a foundation for

111 direct comparison between the species. Therefore, this study aimed to elucidate the

112 influence of low-level stimulus properties on the eye-movements of chimpanzees and

113 humans.

114        We used two approaches in order to separate low-level from high-level

115    influences on eye movements. First, to simulate responses to local stimulus properties,

116    we used the well-established bottom-up saliency model (Itti and Koch, 2000; Walther

117    and Koch, 2006). This model estimates the local saliency of an image based on its

118    low-level components -- such as color, intensity, and component orientations -- and

119    predicts the locations of attention based on these local saliency values. The second

120    approach used global manipulation of stimulus properties (e.g., stimulus properties such

121    as color, configuration, frequency components, orientation, complexity, and location of

122    scene features) and observed how participants changed their patterns of scanning in

123    response to the manipulations.

124          In this study, we used similar stimulus sets to those used by Kano and

125    Tomonaga (2009) and analyzed the participants' responses to social stimuli, especially

126    to the face, as a main measure. In the previous study, chimpanzees and humans fixated

127    the face more frequently than any other part of the scene. The frequent fixation to the

128    face is most likely caused not only by the low-level saliency of the faces, but also by the

129    participants' sensitivity to the social stimuli. This study aimed to investigate the extent

130    to which such facial fixation patterns could be explained by the bottom-up saliency

131    model and could be influenced by the global manipulation of stimulus properties in the

132    scene.

133          Therefore, the topics we addressed in this study were as follows. (1) The

134    degree of similarity in fixation distribution patterns between chimpanzees, humans, and

135    those predicted by the bottom-up saliency model; we expected similar patterns of

136    fixation distribution between chimpanzees and humans even when we controlled for the

137    low-level saliency. (2) The extent to which the gaze of chimpanzees and humans is

138    attracted by low-level saliency. (3) The extent to which the two species' facial fixation

139   is influenced by the global manipulation of stimulus properties in the scene. For (2) and

140   (3), again we expected a significant degree of similarity between the species given their

141   perceptual similarities.

142

143   Method

144   Participants

145         Six chimpanzees (five females, one male; aged 9–31 years) and 16 humans (11

146   females, five males; aged 18–31 years; all Japanese) participated in this experiment.

147   The chimpanzees were members of a social group of 14 individuals living in enriched

148   outdoor compounds and attached indoor residences (Matsuzawa, Tomonaga, and Tanaka,

149   2006). They were highly experienced in observing pictorial representations appearing

150   on a computer screen (Matsuzawa et al., 2006). No food or water deprivation occurred

151   during the study period. Care and use of the chimpanzees adhered to the 2002 version of

152   the Guidelines for the Care and Use of Laboratory Primates published by the Primate

153   Research Institute, Kyoto University. The experimental protocol was approved by the

154   Animal Welfare and Care Committee of the Institute and by the Animal Research

155   Committee of Kyoto University. The human participants were graduate and

156   undergraduate students, who participated in the experiment voluntarily. Informed

157   consent was obtained from all human participants.

158   Apparatus

159         Both species used the same apparatus, in order to ensure the possibility of

160   direct comparison between the species. Participants sat still and unrestrained in an

161   experimental booth, with the eye-tracking apparatus and the experimenter separated by

162   transparent acrylic panels (see S1). A table-mounted eye tracker measured their eye

163  movements using infrared corneal reflection techniques (60 Hz; Tobii X120, Tobii

164  Technology AB). This eye-tracker has wide-angle lenses (±40 degrees in the semicircle

165  above the camera) and thus obviated the necessity to restrain the subjects. The

166  eye-tracker and the 17-inch LCD monitor (1280×1024) were mounted on a movable

167  platform, and the distance between the platform and the participants was adjusted to the

168  point at which the gaze was most accurately recorded (60 ±10 cm). This flexible

169  adjustment of the distance between the platform and the participants enabled us to

170  record the gaze movements of chimpanzees without any head restraint. The participant's

171  gaze was recorded as a relative coordinate with respect to the monitor size (i.e. not as

172  the gaze angle). One degree of gaze angle corresponded to approximately 1 cm on the

173  screen at a typical 60-cm viewing distance.

174      As a result of the training conducted during the study performed by Kano and

175  Tomonaga (2009), the chimpanzees were already skilled at sitting still in front of an

176  eye-tracker and looking upon request at a fixation point that appeared on the screen.

177  Five-point calibration was conducted for humans; for chimpanzees, the calibration

178  points were reduced to two in order to decrease the time required for each calibration

179  process. However, for chimpanzees, the calibration was repeated until the maximum

180  accuracy was obtained. The accuracy was checked by presenting to both species five

181  fixation points on the screen. Using these calibration procedures, six participants of both

182  species were tested for accuracy, and the errors were found to be small and comparable

183  between the species (mean errors of 0.62 ±0.06 and 0.52 ±0.05 cm ± s.e.m. on the

184  monitor for chimpanzees and humans, respectively). The drift (the calibration error due

185  to changes occurring in the eye surface) was checked occasionally by presenting the

186  fixation points to the participants again.

187    Stimuli

188          We prepared 20 color photographs of naturalistic scenes containing a human

189    figure (Figure 2). We used only human figures (all Japanese; no chimpanzees or other

190    animals) in this study because a previous study using an identical experimental

191    procedure (Kano and Tomonaga 2009) found similar fixation patterns in both species

192    for all animal figures. These 20 images served as the control condition. Eight

193    experimental conditions were additionally prepared (for the details of manipulation

194    procedure, see Table 1). In the monochrome, line drawing, and schematic drawing

195    conditions, we eliminated color, low-spatial frequency component, and complex lines,

196    respectively, from the original color scene and aimed to examine the influence of

197    realistic appearance of a scene on the participants' response to the faces. In the blurred

198    and silhouette conditions, we blurred and eliminated local features of face and body

199    from the scene and aimed to examine the influence of those features on the responses.

200    In the upside down and scramble conditions, we inverted and scrambled the scene,

201    respectively, and aimed to examine the influence of orientation and arrangements of

202    bodily parts on the response (i.e. we checked whether participants used only

203    information that the head is above the body). In the headless condition, we eliminated

204    the head from the body and aimed to examine whether participants used only bodily

205    information to fixate the location where the head ought to be. Overall, these conditions

206    aimed to observe whether participants used multiple cues to detect the location of faces

207    in the scene. Each experimental condition was represented by five examples created by

208    manipulating the control images. These five examples were pseudo-randomly selected

209    from the entire set of 20 control images so that each control image was used at least

210    once in the experimental conditions. In total, 60 stimuli were used (40 experimental and

211    20 control images). The images were converted into $1000 \times 800$ pixel images with

212    surrounding gray frames ($1280 \times 1024$ pixels in total; $37 \times 30$ degrees at a typical

213    60-cm viewing distance). We used Adobe Photoshop CS3 to process the images.

214    Procedure

215         Procedural differences for testing chimpanzees and humans were minimized to

216    allow for direct comparisons between species. In each trial, an image was presented

217    after participants focused on a fixation point that appeared at a random position on the

218    screen. Participants were then allowed to view images freely. The participants of both

219    species rarely kept gazing at the fixation point after the image presentation (i.e.

220    spontaneous scanning was almost always observed). Stimuli were presented for 3 sec

221    each. The presentation order of conditions and trials were randomized for each

222    participant so that the same conditions were not presented more than three times in

223    succession. 20 other stimuli depicting various interesting scenes (e.g. pictures of funny

224    faces) were presented occasionally during the sessions in order to keep the participants

225    interested. The entire session therefore consisted of 80 trials: 60 experimental stimuli

226    and 20 others. The entire session was conducted on a single day for humans, whereas

227    the session was divided among 15 days for the chimpanzees in order to decrease the

228    time required for each daily experiment (each day used six examples for the

229    chimpanzees). In a preliminary session, we confirmed that our human participants

230    showed similar scanning patterns of bodies/faces when tested on separate days

231    (comparing the results from this study with those from Kano and Tomonaga (2009)).

232    Daily experiments lasted 10–15 min for the chimpanzees and 20 min for the humans.

233    Human participants received book coupons as rewards after the session, and

234    chimpanzees received a small piece of apple after each trial. The reward was given for

235    chimpanzees in return for the initial fixation at the beginning of the trial, and thus was

236    given independently of their viewing behavior during the image presentation. Overall,

237    those procedural differences between the species were made in an effort to increase the

238    motivation of both species to participate in the daily experiment, and to keep their

239    interest during the presentation of each image (3 sec). Trials in which participants only

240    glanced at the monitor (one or two fixations) were repeated after the whole session and

241    were replaced by the new trials. As a result, we had no loss of trials for both species.

242

243    Data analysis

244    *Fixation definition*

245        A fixation was scored if the gaze remained stationary within a radius of 50

246    pixels for at least 75 ms (more than five measurement samples). Otherwise, the recorded

247    sample was defined as part of a saccade. The records during the first 200 ms were

248    eliminated from the analysis, thereby eliminating fixations that followed the offset of

249    the initial fixation point.

250    *Area of Interest (AOI)*

251        Each stimulus was divided into areas of interest (AOI) for the purpose of

252    quantitative comparison. Each scene was divided into background, face, torso, arms,

253    and legs (Figure 3, bottom). Each AOI was drawn 20 pixels larger than the precise

254    outline of the features to avoid errors in gaze estimations. The AOI of background, torso,

255    legs, arms, and face were laid above in this order (i.e. face is the topmost). If two or

256    more AOIs were duplicated, the samples were added to the upper AOI.

257    *Chance level*

258        The chance level was set on the assumption that participants would view

259    images randomly. However, the participants generally showed a central bias in fixation

260    distribution, while the model did not (evident by inspection of Figure 2). This needs to

261    be controlled to compare participants with the model, because such central bias is

262    known to be caused either by the participants' bias in scanning images or by the

263    experimenter's bias in the arrangement of main objects in the scene (Henderson,

264    Brockmole, Castelhano, and Mack, 2007; Tatler, 2007) (i.e., caused independently of

265    the low-level stimulus properties). Therefore, in this study, we modified the definition

266    of chance level by controlling for such particular bias shown by each participant.

267    Specifically, we compared the particular scanpath, which was obtained from a

268    participant (or the model) in a trial, with all the other scanpaths, which were obtained

269    from the same participants (or the model) in all the other trials of the experiment. All

270    data shown in this study were calculated as differences between the value obtained from

271    the particular scanpath and the mean value obtained from the other control scanpaths

272    (i.e., the chance level).

273    *Saliency model*

274         We used the well-established bottom-up model to estimate the low-level

275    saliency of the images (Itti and Koch, 2000; Walther and Koch, 2006). This model

276    processes the image with respect to several features -- such as color (red-green,

277    blue-yellow), intensity, and orientation (0, 45, 90, 135 degrees) -- then extracts the local

278    discontinuities in each feature, and finally combines them into a single 'saliency map'

279    (Figure 1). The model then predicts a scanpath based on the saliency map, selecting

280    salient locations in order of decreasing saliency. In this experiment, the saliency maps

281    and    the    model    scanpaths    were    generated    by    Saliency    Toolbox    2.2

282    (http://www.saliencytoolbox.net) in Matlab with all-default parameters. We used the

283  original resolution of images (1280×1024; including the surrounding gray frame) for the

284  simulation in the model. Because this model does not predict the duration of each

285  fixation, we arbitrarily set the scanpath length of the model as 9 fixations (about as long

286  as chimpanzee scanpaths in 3-s viewing) to compare the model with the participants.

287  There is no variance in the output when repeating the simulation.

288       To determine the saliency value at each fixated area, we employed the

289  following procedures. First, saliency value was normalized within each map to a range

290  of 0 (not salient) to 1 (highly salient). Second, to avoid errors in gaze estimation, the

291  saliency map was divided into a 12×9 grid, and all saliency values (i.e. $1280 \times 1024$

292  samples, in total) were summed within each grid (i.e. each grid had approx. $100 \times 100$

293  samples). The fixated area was defined as the grid where the fixation was observed, and

294  the saliency value of each grid was used for the saliency value at each fixated area.

295

296  Results

297       Figure 3a shows the distribution patterns of fixation over the scene in each

298  species/model (the data were sampled from 20 control images). Comparing between

299  species and between AOIs, we found a significant interaction ($F_{2.3, 47} = 9.52$, $p < 0.001$,

300  $\eta^2 = 0.32$)[1] because chimpanzees distributed their fixations over the scene more widely

301  than did humans. Comparing between AOIs respectively for each species, we found

302  significant main effects for both chimpanzees ($F_{1.8, 9.3} = 23.80$, $p < 0.001$, $\eta^2 = 0.82$) and

303  humans ($F_{2.4, 37} = 358.86$, $p < 0.001$, $\eta^2 = 0.96$) because both species showed higher

304  proportion of fixations on particular areas (the bodies, especially faces, rather than

305  backgrounds) than would be expected by chance (represented as zero in the figures).

306  This pattern of results emerged even when the model was subtracted from each species:

307     chimpanzees ($F_{1.8, 9.3} = 5.61$, $p = 0.003$, $\eta^2 = 0.52$) and humans ($F_{2.4, 37} = 159.74$, $p <$

308     $0.001$, $\eta^2 = 0.91$). This pattern emerged no later than the first two fixations, as shown in

309     Figures 3b and 3c, and is consistent with the previous reports in humans (Crouzet,

310     Kirchner, and Thorpe, 2010; Fletcher-Watson, Findley, Leekam, and Benson, 2008;

311     Honey, Kirchner, VanRullen, 2008). The global similarities in distribution patterns of

312     fixation among chimpanzees, humans, and the model suggest that the saliency model

313     partially (but not fully) explained those patterns for the two species. Although

314     chimpanzees were more similar to the model than were humans in that regard, it should

315     be noted that this does not mean that the low-level visual saliency influenced

316     chimpanzees more strongly than humans; this means that chimpanzees distributed their

317     fixations over the scene more widely than did humans, but less widely than did the

318     model.

319             Indeed, chimpanzees and humans did not significantly differ in their responses

320     to low-level visual saliency. There was no significant effect of species in the saliency

321     values at fixation (Figure 4); neither the main effect of species ($F_{1, 20} = 0.014$, $p = 0.90$,

322     $\eta^2 = 0.001$) nor the interaction between species and fixation order ($F_{5, 100} = 0.46$, $p =$

323     $0.80$, $\eta^2 = 0.023$) was significant. Overall, however, both species fixated on salient

324     regions in the scene more than would be expected by chance: the mean saliency values

325     for the first 6 fixations were significantly higher than zero in both chimpanzees ($t(5) =$

326     $9.83$, $p < 0.001$) and humans ($t(15) = 19.27$, $p < 0.001$). This pattern emerged more

327     strongly for the earlier than for the later fixations: saliency value decreased as a function

328     of increasing fixation order ($F_{5, 100} = 3.20$, $p = 0.010$, $\eta^2 = 0.13$). These results suggest

329     that the saliency model predicted the distribution patterns of fixation in both

330     chimpanzees and humans better than chance, especially for the early fixations. However,

331     it should be noted that this result does not necessarily mean that the low-level saliency

332     alone influenced the species' distribution patterns of fixation, because such frequently

333     fixated areas (e.g., bodies and faces) were in general more visually salient (because of

334     the complexity of lines, for example) as well as more informative than the other areas of

335     the scene (Figure 3; refer to (Henderson et al., 2007) for a similar discussion).

336             We then examined the effect of image manipulations on the fixation patterns of

337     chimpanzees and humans (Figure 5). Figure 5b shows the proportion of fixations on the

338     faces as a function of image manipulations. There was no interaction of species with

339     condition ($F_{3.5,\ 70} = 1.13$, $p = 0.34$, $\eta^2 = 0.05$). The main effect of species was

340     significant: humans showed a higher proportion of fixations on faces than did

341     chimpanzees ($F_{1,\ 20} = 5.51$, $p = 0.029$, $\eta^2 = 0.21$), which is consistent with the

342     aforementioned result. The main effect of condition ($F_{3.5,\ 70} = 5.33$, $p < 0.001$, $\eta^2 = 0.21$)

343     was significant: participants showed a lower proportion of face fixations in the headless

344     than the other conditions (as was revealed by the pair-wise comparisons with

345     Bonferroni's correction).

346             However, even in the headless condition, both species showed a higher

347     proportion of fixations on the face original locations than would be expected by the

348     model (as was revealed by the post-hoc $t$-tests). This means that even when a head was

349     actually absent from the scene, both species concentrated fixations on the area where

350     the face would have been (i.e. above the body).

351             Figure 5c shows the mean saliency values at the first 6 fixations as a function

352     of image manipulations. The main effect of condition was significant ($F_{8,\ 160} = 46.93$, $p$

353     $< 0.001$, $\eta^2 = 0.70$), probably modulated by saliency (or informativeness) in local

354     features of the scene, which was an outcome of image manipulations. Importantly, there

355   was no effect of species despite these image manipulations, either the main effect of

356   species ($F_{1,\ 20} = 0.017$, $p = 0.89$, $\eta^2 = 0.001$) or the interaction between species and

357   condition ($F_{8,\ 160} = 1.18$, $p = 0.31$, $\eta^2 = 0.05$).

358

359 Discussion

360 Chimpanzees and humans distributed fixations over the scene non-randomly, and

361 showed higher fixation proportions on particular areas of the scene, especially faces,

362 than would be expected by the saliency model. However, humans showed an even

363 higher proportion of fixation on the bodies and faces than did chimpanzees. These

364 results emerged even at the first two fixations, at the earliest moments of scene

365 inspection, suggesting that those fixation patterns reflect automatic rather than voluntary

366 control of gaze. Saliency values of chimpanzees and humans in the fixated region were

367 higher than would be expected by chance, suggesting that low-level saliency partially

368 (but not fully) predicted the species' distribution patterns of fixation. However,

369 chimpanzees and humans did not significantly differ in their responses to low-level

370 saliency. None of global manipulations of stimulus properties in the scene (color,

371 configuration, frequency components, orientation, complexity, and local features)

372 critically altered both species' strong tendency toward fixating faces, suggesting that

373 both species used multiple cues to fixate faces. In addition, although those

374 manipulations changed the extent to which low-level saliency influenced both species,

375 chimpanzees and humans did not differ in the degree of change in the response.

376 Therefore, chimpanzees and humans seem to be qualitatively similar in the

377 sense that both species have an enhanced perceptual mechanism to guide their fixation

378 location, one which is more complex than would be presumed on the basis of the

379 saliency model (i.e. color, intensity, and orientations), and have multiple strategies to

380 perceive the location of faces. Quantitatively, these two species did not differ

381 significantly in their responses to low-level saliency, suggesting that they have similar

382 perceptual mechanisms to guide the fixation locations.

383    Einhäuser et al. (Einhuäuser, Kruse, Hoffmann, and König, 2006) used the

384    standard saliency model to predict the fixation location of monkeys (rhesus macaques)

385    and humans when presented with the still images of naturalistic scene (without social

386    contents). They found that monkeys and humans did not differ significantly in their

387    responses to low-level saliency when viewing those images, which is consistent with the

388    present study comparing chimpanzees and humans. However, when the

389    luminance-contrast (or the saliency) was manipulated locally in the image, the monkeys

390    responded to those manipulated areas more strongly than did the humans. In the similar

391    analysis to that of Einhäuser et al. (2006) and this study, Berg et al. (Berg, Boehnke,

392    Marino, and Munoz, and Itti, 2009) found that, when presented with dynamic scenes

393    including various social, non-social, and narrative contents, humans responded to the

394    low-level visual saliency more strongly than did monkeys (perhaps because monkeys

395    tended to move their eyes independently of the stimuli (e.g. inattentiveness to the

396    stimuli) or show a large degree of individual differences in their fixation patterns),

397    which is somewhat inconsistent with Einhäuser et al. (2006) and this study. Therefore,

398    multiple factors seem to be involved in the species difference in the responses to the

399    low-level saliency. To clarify those factors, it is necessary to directly compare between

400    the three species for their fixation patterns when presented with various contents of still

401    and dynamic scenes.

402    Cerf et al. (Cerf, Harel, Einhäuser, and Koch, 2008) have shown that the

403    addition of a "face channel" into the standard saliency model better predicts the fixation

404    patterns of human participants viewing a naturalistic scene that includes faces. They

405    used an established face detector algorithm for that purpose, which predicts the location

406    of faces based on local facial features (e.g. local discontinuities in intensity around eye

407 and nose regions). The distribution patterns of fixation observed in this study suggest

408 that chimpanzees and humans have such a face perception channel in addition to the

409 low-level channels. However, the mechanism underlying such a face channel seems

410 more complex in chimpanzees and humans than would be assumed by the face detector

411 algorithm. This is because chimpanzees and humans concentrated fixations on the faces

412 even when local features of faces were significantly reduced (schematic and blurred) or

413 completely silhouetted out of the scene. They did so even when the faces were removed

414 completely (headless), suggesting that chimpanzees and humans can use the bodily

415 configuration alone to fixate where faces ought to be. On the other hand, chimpanzees

416 and humans also seem to be able to use local cues to fixate faces, because they

417 concentrated fixations on the face parts even when bodily configuration was disrupted

418 (scrambled). Therefore, chimpanzees and humans seem to have an enhanced perceptual

419 mechanism to guide their fixations to a face, a mechanism that is more complex than

420 would be assumed by the standard saliency model or the saliency model combined with

421 face detection.

422         Notwithstanding those similarities between the species, chimpanzees and

423 humans differ quantitatively in the distribution patterns of fixations. Humans showed a

424 higher proportion of fixations on bodies and faces than did chimpanzees. As clarified

425 above, it is unlikely that this species difference resulted from their differential responses

426 to the low-level visual properties (or in their differential tendencies for central bias). It

427 is also unlikely that this species difference resulted from the use of human, as opposed

428 to chimpanzee figures as stimuli, because a previous study (Kano and Tomonaga, 2009)

429 obtained the same patterns of results when using chimpanzees and other mammals as

430 the stimulus models. Therefore, we interpreted the results in the following two ways.

431  First, although the results suggested that both species have similar mechanism to guide

432  their gaze to the social stimuli (body and face), those mechanisms may operate

433  differently in each species. For example, humans may put more emphasis on the

434  body/face channels to create the saliency map, and so humans may perceive bodies and

435  faces as more salient than chimpanzees do. Second, humans, compared to chimpanzees,

436  may have a stronger tendency to process scenes in a top-down rather than a bottom-up

437  manner, and thus would be expected to show a higher proportion of fixations on the

438  semantically informative areas such as bodies and faces. Further studies are necessary to

439  test these two possibilities.

440        In summary, this study presented, to chimpanzees and humans, naturalistic

441  (unmanipulated) scenes including body, face, and their manipulated representations. We

442  then compared among the two species and the saliency model for the fixation patterns

443  on the images. We found that the saliency model did not fully predict the fixation

444  patterns actually observed in chimpanzees and humans. In addition, both species used

445  multiple cues to fixate the face. There was no evidence suggesting any differences

446  between chimpanzees and humans in the perception of low-level saliency (e.g. color,

447  intensity, or orientations). Therefore, we showed a substantial amount of similarities in

448  the perceptual mechanism underlying gaze guidance between chimpanzees and humans,

449  and thereby offer a foundation for the direct comparison between the species. Further

450  studies are necessary to elucidate the high-level similarities and differences between the

451  species (e.g. social sensitivity, knowledge-based attention).

466    References

467    Argyle, M., Cook, M. (1976). *Gaze and mutual gaze*. Cambridge: Cambridge

468          University Press

469    Berg, D. J., Boehnke, S. E., Marino, R. A., Munoz, D. P., & Itti, L. (2009). Free

470          viewing of dynamic stimuli by humans and monkeys. *J Vis, 9*(5), 1-15.

471    Cerf, M., Harel, J., Einhuäuser, W., Koch, C. (2008). Predicting human gaze using

472          low-level saliency combined with face detection. *Adv Neural Inform Process*

473          *Syst, 20*, 241-248.

474    Chevalier-Skolnikoff, S. (1973). Facial expression of emotion in nonhuman primates. In

475          P. Ekman (Ed.), *Darwin and facial expression: A century of research in review*

476          (pp. 11-89). New York: Academic Press.

477    Crouzet, S. M., Kirchner, H., & Thorpe, S. J. Fast saccades toward faces: Face detection

478          in just 100 ms. *J Vis, 10*(4), 1-17.

479    Einhäuser, W., Kruse, W., Hoffmann, K. P., & König, P. (2006). Differences of monkey

480          and human overt attention under natural conditions. *Vision Res, 46*(8-9),

481          1194-1209.

482    Dahl, C. D., Wallraven, C., Bulthoff, H. H., Logothetis, N. K. (2009). Humans and

483          macaques employ similar face-processing strategies. *Curr Biol, 19*(6), 509-513.

484    Fuchs, A. F. (1967). Saccadic and smooth pursuit eye movements in the monkey. *J*

485          *Physiol, 191*(3), 609-631.

486    Fletcher-Watson, S., Findlay, J. M., Leekam, S. R., & Benson, V. (2008). Rapid

487          detection of person information in a naturalistic scene. *Perception, 37*, 571-583.

488    Gothard, K. M., Brooks, K. N., Peterson, M. A. (2009). Multiple perceptual strategies

489          used by macaque monkeys for face recognition. *Anim Cogn, 12*(1), 155-167.

490    Guo, K., Robertson, R. G., Mahmoodi, S., Tadmor, Y., Young, M. P. (2003). How do

491          monkeys view faces?—a study of eye movements. *Exp Brain Res, 150*(3),

492          363-374.

493   Hattori, Y., Kano, F., Tomonaga, M. (2010). Differential sensitivity to conspecific and

494        allospecific cues in chimpanzees and humans: A comparative eye-tracking study.

495        *Biology Lett. 6*(5), 610-613

496   Henderson, J. M., Brockmole, J. R., Castelhano, M. S., Mack, M. (2007). Visual

497        saliency does not account for eye movements during visual search in real-world

498        scenes. In R. P. G. van Gompel, M. H. Fischer, W. S. Murray & R. L. Hill (Eds.),

499        *Eye movements: A window on mind and brain* (pp. 537-562). Neitherlands:

500        Elsevier.

501   Honey, C., Kirchner, H., & VanRullen, R. (2008). Faces in the cloud: Fourier power

502        spectrum biases ultrarapid face detection. *J Vis, 8*(12), 1-13.

503   Hirata, S., Fuwa, K., Sugama, K., Kusunoki, K., Fujita, S. (2010). Facial perception of

504        conspecifics: chimpanzees (*Pan troglodytes*) preferentially attend to proper

505        orientation and open eyes. *Anim Cogn, 13*(5), 679-688.

506   Itti, L., Koch, C. (2000). A saliency-based search mechanism for overt and covert shifts

507        of visual attention. *Vision Res, 40*(10-12), 1489-1506.

508   Kano, F., Tomonaga, M. (2009). How chimpanzees look at pictures: a comparative

509        eye-tracking study. *Proc Roy Soc B, 276*(1664), 1949-1955.

510   Kano, F., Tomonaga, M. (2010). Face scanning in chimpanzees and humans: continuity

511        and discontinuity. *Anim Behav, 79*, 227-235.

512   Keating, C. F., Keating, E. G. (1982). Visual scan patterns of rhesus monkeys viewing

513        faces. *Perception, 11*(2), 211-219.

514   Matsuno, T., Kawai, N., Matsuzawa, T. (2004). Color classification by chimpanzees

515        (*Pan troglodytes*) in a matching-to-sample task. *Behav Brain Res, 148*(1-2),

516        157-165.

517     Matsuzawa, T. (1985). Colour naming and classification in a chimpanzee (*Pan*

518          *troglodytes*). *J Hum Evol, 14*(3), 283-291.

519     Matsuzawa, T. (1990). Form perception and visual acuity in a chimpanzee. *Folia*

520          *Primatol, 55*(1), 24-32.

521     Matsuzawa, T., Tomonaga, M., Tanaka, M. (2006). *Cognitive development in*

522          *chimpanzees*. Tokyo: Springer.

523     Nahm, F. K. D., Perret, A., Amaral, D. G., Albright, T. D. (1997). How do monkeys look

524          at faces? *J Cogn Neurosci, 9*(5), 611-623.

525     Parr, L. A., Dove, T., Hopkins, W. D. (1998). Why faces may be special: evidence of the

526          inversion effect in chimpanzees. *J Cogn Neurosci, 10*(5), 615-622.

527     Parr, L. A., Hecht, E., Barks, S. K., Preuss, T. M., Votaw, J. R. (2009). Face processing

528          in the chimpanzee brain. *Curr Biol, 19*(1), 50-53.

529     Shepherd, S. V., Steckenfinger, S. A., Hasson, U., Ghazanfar, A. A. (2010).

530          Human-monkey gaze correlations reveal convergent and divergent patterns of

531          movie viewing. *Curr Biol, 20*(7), 649-656.

532     Tatler, B. W. (2007). The central fixation bias in scene viewing: Selecting an optimal

533          viewing position independently of motor biases and image feature distributions.

534          *J Vis, 7*(14).

535     Tomonaga, M. (2007). Visual search for orientation of faces by a chimpanzee (*Pan*

536          *troglodytes*): face-specific upright superiority and the role of facial configural

537          properties. *Primates, 48*(1), 1-12.

538     Tomonaga, M. (2010). Chimpanzee eyes have it? Social cognition on the basis of gaze

539          and attention from the comparative-cognitive-developmental perspective. In E.

540          Lornsdorf, S. Ross & T. Matsuzawa (Eds.), *The mind of the chimpanzee:*

541        *Ecological and empirical perspectives*. Chicago: University of Chicago Press.

542    Tomonaga, M., Imura, T. (2008). Chimp in the shadow: Efficient detection of

543        chimpanzee body by chimpanzees? *Primate Res, 24*(S), 14-15 (Japanese abstract

544        only).

545    Tomonaga, M., Imura, T. (2009). Faces capture the visuospatial attention of

546        chimpanzees(*Pan troglodytes*): evidence from a cueing experiment. *Front Zool,*

547        *6*(1), 14.

548    Tomonaga, M., Matsuzawa, T. (1992). Perception of complex geometric-figures in

549        chimpanzees (Pan troglodytes) and humans (Homo sapiens): Analysis of visual

550        similarity on the basis of choice reaction time. *J Comp Psychol, 106*(1), 43-52.

551    Walther, D., Koch, C. (2006). Modeling attention to salient proto-objects. *Neural*

552        *Networks, 19*(9), 1395-1407.

553    Yarbus, A. L. (1967). *Eye Movements and Vision*. New York: Plenum Press.

554

555

556

557    Content Note

558    1. In the ANOVAs, in cases in which the assumption of homogeneity of variance was

559    violated, the Greenhouse-Geisser correction was applied, and corrected $P$ values were

560    calculated.

561 Tables

Table 1. *Procedures Used for Image Manipulation*

| condition | n | procedure |
|---|---|---|
| control | 20 | |
| monochrome | 5 | The color was removed from the original photographs. |
| line drawing | 5 | Only edges were extracted from the monochrome photograph (with a Photoshop function), and binary image processing techniques simplified the image (emphasizing the fat lines and eliminating the thin lines and small dots). |
| schematic drawing | 5 | The edges were roughly traced with simple black circles and lines. |
| blurred | 5 | The edges were blurred to the extent that the facial features were not recognizable (a Gaussian blur 20 pixels in diameter). |
| silhouette | 5 | The figure was colored in black, and binary image processing techniques transformed the background into black and white patches. |
| upside down | 5 | The original photographs were turned upside down. |
| scrambled | 5 | The original scenes were superimposed into a $6 \times 5$ matrix, and each block of the matrix was randomly rearranged. A matrix was defined so that a block includes the whole face (i.e. the face was intact). |
| headless | 5 | The head was eliminated so that the background was visible through the regions in which the head was previously located. To this end, the headless figure was cropped in the first image and superimposed on the second image that contains only background. |

562

563

564

565    Figure legends

566    Figure 1

567    Scanpaths of a chimpanzee and a human, each superimposed on the naturalistic scene

568    (a) and fine art painting (b; Paul Klee, 1923, "*Puppet Theater*"; see Supporting material

569    for the quantitative data). Fixations and saccades are indicated by dots and lines,

570    respectively. The stimuli were presented for 3 sec. each. Also shown are a raw saliency

571    map and the scanpath predicted by the model. Bright areas indicate areas of high

572    saliency. By design, the model made 9 fixations on the images in the order of decreasing

573    saliency.

574    Figure 2

575    The locations of all fixations made by a chimpanzee, a human, and the model. While the

576    model showed a relatively even distribution of fixations over the scene, the chimpanzee

577    and the human showed a central bias in the distribution. Therefore, the chance level

578    (random gaze pattern) was adjusted to control for this observed bias (see text).

579    Figure 3

580    (a) Proportion of fixations on each area of interest (AOI; see the diagram for an

581    example) in each image by chimpanzees (n = 6) and humans (n = 16). (b) Proportion of

582    images (n = 20) in which a fixation was observed in each AOI at each fixation order.

583    The first 6 fixations are presented here. (c) The sum proportion of images at the first

584    two fixations, showing that the results from (a) are evident no later than the first two

585    fixations. The data are from the control condition. All data are shown as the difference

586    from the chance level. T-tests compared between chimpanzees and humans, and

587    between each species and the model (one-sample). * $p < 0.05$, ** $p < 0.01$, *** $p <$

588    0.001. Error bars indicate s.e.m.

589    Figure 4

590    The saliency values at the first 6 fixations. The saliency value was standardized, and

591    ranges from 0 (not salient) to 1 (highly salient). The data are taken from the control

592    condition. n.s. not significant. Error bars indicate s.e.m.

593    Figure 5

594    (a) Examples of stimuli presented in each experimental condition. Note that the original

595    stimuli were in color. (b) Proportion of fixations on the face in each image by

596    chimpanzees (n = 6), humans (n = 16), and the model. (c) The mean saliency values at

597    the first 6 fixations for chimpanzees and humans. All data are shown as the difference

598    from the chance level. T-tests compared between chimpanzees and humans, and

599    between each species and the model (one-sample). * $p < 0.05$, ** $p < 0.01$, *** $p <$

600    0.001, n.s. not significant (the P values for Figure 5c are 0.75, 0.20, 0.29, 0.78, 0.74,

601    0.19, 0.26, 0.86, 0.19, for each condition, from left to right). Error bars indicate s.e.m.

602

603    Figure



604

605    Figure 1

606



607

608    Figure 2

609

610    Figure 3

611

612    Figure 4

613

614    Figure 5

615

616