1  **Title: Are great apes able to reason from multi-item samples to**

2  **populations of food items?**

3  Short title: Are great apes able to reason from sample to population?

4  **Johanna Eckert[1,2], Hannes Rakoczy[2], Josep Call[1,3]**

5  [1] Department of Developmental and Comparative Psychology, Max Planck Institute for

6  Evolutionary Anthropology, Deutscher Platz 6, 04103 Leipzig, Germany

7  [2] Department of Developmental Psychology, University of Goettingen, Waldweg 26,

8  37073 Goettingen, Germany

9  [3] School of Psychology and Neuroscience, University of St Andrews, St Andrews KY16

10  9JP, UK

11  Conflict of interest: none

12  Corresponding Author: Johanna Eckert

13  Email: johanna_eckert@eva.mpg.de

14  Postal address: Deutscher Platz 6, 04103 Leipzig, Germany

15  Phone +49 (341) 3550 – 424

16

## **Abstract**

Inductive learning from limited observations is a cognitive capacity of fundamental importance. In humans, it is underwritten by our intuitive statistics, the ability to draw systematic inferences from populations to randomly drawn samples and vice versa. According to recent research in cognitive development, human intuitive statistics develops early in infancy. Recent work in comparative psychology has produced first evidence for analogous cognitive capacities in great apes who flexibly drew inferences from populations to samples. In the present study, we investigated whether great apes (*Pongo abelii*, *Pan troglodytes*, *Pan paniscus*, *Gorilla gorilla*) also draw inductive inferences in the opposite direction, from samples to populations. In two experiments, apes saw an experimenter randomly drawing one multi-item sample from each of two populations of food items. The populations differed in their proportion of preferred to neutral items (24:6 vs. 6:24) but apes saw only the distribution of food items in the samples that reflected the distribution of the respective populations (e.g. 4:1 vs. 1:4). Based on this observation they were then allowed to choose between the two populations. Results show that apes seemed to make inferences from samples to populations and thus chose the population from which the more favorable (4:1) sample was drawn in Experiment 1. In this experiment, the more attractive sample not only contained proportionally but also absolutely more preferred food items than the less attractive sample. Experiment 2, however, revealed that when absolute and relative frequencies were disentangled, apes performed at chance level. Whether these limitations in apes' performance reflect true limits of cognitive competence or merely performance limitations due to accessory task demands is still an open question.

40 **Keywords:** Intuitive statistics; Probabilistic reasoning; Comparative cognition; Non-

41 human primates; Numerical cognition

## Introduction

43       Making general inferences from limited data is one of the key components of

44 human inductive learning [see e.g. Skyrms, 1975; Holland, 1986; Tenenbaum et al. ,

45 2006; Denison and Xu, 2012]. Traditionally, statistical reasoning was deemed to be

46 difficult and error-prone [e.g. Tversky and Kahneman, 1974; Tversky and Kahneman,

47 1981; Cosmides and Tooby, 1996] and dependent on language and formal education

48 [e.g. Piaget and Inhelder, 1975]. However, recent studies suggest that even very young

49 human infants have an astonishingly broad understanding of statistical relations: They

50 are able to generalize from small samples to larger populations [Xu and Garcia, 2008;

51 Denison et al. , 2013], make predictions about single event probabilities [e.g. Teglas et

52 al. , 2007] and use these predictions to guide their search for desired objects [Feigenson

53 et al. , 2002; Denison and Xu, 2010b; Denison and Xu, 2014]. In one remarkable study,

54 for example, infants were confronted with two jars containing mixtures of preferred and

55 non-preferred types of candy [Denison and Xu, 2010b]. After they had watched the

56 experimenter randomly sampling one piece of candy from each jar and placing it in an

57 occluded cup, most infants searched in the cup that contained a sample from the jar with

58 a higher proportion of their preferred candy [Denison and Xu, 2010b]. Hence, infants

59 seem to have used the proportional information provided by the populations to reason

60 about the samples. Moreover, infants can integrate probabilistic information with

61 information from other domains such as intuitive physics or intuitive psychology [Xu

62 and Denison 2009; Teglas et al. , 2011; Denison et al. , 2014]. For example, infants

63  understand that a preference of the experimenter for a certain type of object can turn a

64  sampling process into a non-random event. If the same experimenter, however, is

65  blindfolded, infants expect the sampled objects to reflect the proportions within

66  populations [Xu and Denison 2009]. These findings imply that at least at the age of 6

67  months, humans already flexibly use intuitive statistics to predict the outcome of events.

68  Being apparently independent of language or formal education, this raises the question

69  whether these kinds of probabilistic reasoning represent an evolutionary ancient trait

70  that is shared with other species.

71      Many species are capable of numerical cognition: For example, great apes [e.g.

72  Boysen and Berntson, 1989; Call, 2000; Hanus and Call, 2007; Beran et al. , 2013], old-

73  and new-world monkeys [e.g. Beran et al. , 2008; Barnard et al. , 2013; Beran and

74  Parrish, 2016], elephants [Perdue et al. , 2012], bears [Vonk and Beran, 2012], raccoons

75  [Davis, 1984], dogs [Ward and Smuts, 2007], cats [Pisa and Agrillo, 2009], birds [e.g.

76  Rugani et al. , 2013], fish [e.g. Potrich et al. , 2015], and even insects [bees: Dacke and

77  Srinivasan, 2008; ants: Reznikova and Ryabko, 2011] are able to compare quantities,

78  suggesting that representing numerosity is an evolutionary ancient trait. The practical

79  advantages of such a capacity are obvious: in the context of foraging, for example,

80  comparing quantities is a highly useful tool to identify the most profitable feeding

81  location [see e.g. Farnsworth and Smolinski, 2006 and Hunt et al. , 2008 for field

82  experiments on quantity discrimination in a foraging context]. In the context of

83  competition comparing ones´ own group size with that of a rival group can help to

84  estimate the chances of winning a potential fight [e.g. McComb et al. , 1994; Wilson et

85  al. , 2002; Benson-Amram et al. , 2011]. Chimpanzees, for example, have been found to

86  attack an opponent group only if their own group outnumbers those of their conspecifics

87  by at least 1.5 [Wilson et al. , 2002].

88  Relatedly, one can imagine that in some situations it would make sense for an

89  animal to be able to make probability judgments instead of straightforward quantity

90  comparisons. Efficient foraging, for instance, requires an individual to search for food

91  in locations that most likely provide the best payoff in relation to foraging time [Geary

92  et al. , 2015; for a review about optimal foraging theory see e.g. Hamilton, 2010]. One

93  possibility to identify the best payoff per time unit is to use the relative frequency of

94  past successes in a feeding location. Imagine a group of chimpanzees that has to decide

95  in the morning in which direction to go: Either towards feeding ground A or towards

96  feeding ground B. The apes might want to compare the proportion of times they visited

97  each feeding ground and obtained a sufficient amount of food instead of simply

98  comparing the absolute number of times they were successful in each location. Hence,

99  non-human animals could clearly benefit from an ability exceeding a mere estimation of

100  absolute or relative numerosity, namely a sense for probabilistic relations, i.e. intuitive

101  statistics. Future research will need to investigate both when and due to which selection

102  pressures intuitive statistics evolved.

103  A recent comparative study investigated intuitive statistical abilities in non-

104  human great apes with the same kinds of methods used in infancy research [Rakoczy et

105  al. , 2014]. Individuals of four great ape species were presented with two populations of

106  food items. Both populations consisted of the same two types of food (one type clearly

107  preferred over the other) but with different relative frequency distributions. The

108  experimenter drew a one-object-sample from each population and gave the subject a

109  choice between the two hidden samples. Hence, subjects had to infer which population

110  was more likely to yield a preferred food item as a sample. Interestingly, individuals of

111  all tested great ape species were able to form correct expectations about the probability

112  of the sampling events, even when absolute and relative frequencies within the

113  populations were disentangled. Apes´ inferences were, therefore, not only based on

114  information about absolute frequency, but instead they were truly based on probabilistic

115  information. Most recently, another representative of the primate order was tested in the

116  same paradigm: Capuchin monkeys *Sapajus sp* [Tecwyn et al. , 2016]. In a series of

117  four experiments, the monkeys were allowed to choose between the randomly drawn

118  samples of two populations of food items with different proportions of preferred and

119  non-preferred food. Results revealed that a few individuals might have drawn

120  probabilistic inferences based on proportional information (control conditions excluded

121  the usage of simpler choice heuristics). However, monkeys´ performance in a baseline

122  control condition was unexpectedly low, questioning whether they truly fully

123  understood the procedure. It remains, therefore, an open question whether primates

124  other than great apes are capable of intuitive statistics.

125      Although the findings of Rakoczy and colleagues [2014] raised the possibility

126  that apes and human infants may operate with the same cognitive capacities for intuitive

127  statistics, it leaves many open questions for future research. One fundamentally

128  important question is whether apes´ intuitive statistics reveal the same kinds of

129  flexibility and generality as those found in human infants. In particular, does their

130  ability to form expectations about samples randomly drawn from populations (inference

131  population → sample; Rakoczy et al. 2014) extend to the ability to reason from a given

132  sample to the corresponding population (inference sample → population)?

133    In human infants, this question was addressed using the violation of expectation

134    (VOE) looking-time paradigm [Xu and Garcia, 2008; Denison et al. , 2013].   In one

135    study [Xu and Garcia 2008] eight-month-old infants were presented with boxes

136    containing populations of red and white Ping-Pong balls. The distribution of red to

137    white balls was either 9:1 or 1:9. During test-trials, the box containing one of the two

138    populations of Ping-Pong balls was covered and the infants watched the experimenter

139    drawing (apparently randomly) a sample of either 4 red and 1 white Ping-Pong balls or

140    1 red and 4 white Ping-Pong balls. Subsequently, the experimenter removed the cover

141    of the box she had drawn from and revealed the population. Infants looked longer at the

142    "mostly red"- sample when it was drawn from the "mostly white" population

143    (unexpected) than when it was drawn from the "mostly red" population (expected). The

144    analogue was true for the "mostly white"-sample. In a control condition it could be

145    ruled out that infants simply reacted to the perceptual mismatch between sample and

146    population: Instead of drawing the balls as samples from the box, the experimenter

147    pulled them out of her pocked and placed them next to the box, resulting in equal

148    looking-times at both matched and mismatched outcomes. This implies that, confronted

149    with a sample, infants were able to make inferences about the associated population.

150    Applying the same paradigm, a second study [Denison et al. 2013] showed that even 6-

151    month-old infants had intuitions about relationships between samples and populations,

152    suggesting that the ability to make inferences based on samples develops very early in

153    human ontogeny. The results of these two infant studies [Xu and Garcia, 2008; Denison

154    et al. , 2013] together with the findings of the first study on intuitive statistics in great

155    apes [Rakoczy et al. , 2014] may indicate that the capacity of making inductive

156    inferences is shared with other non-human species.

157        In the current study we tested great apes' ability to reason from samples to

158    populations. Combining the methodology of Rakoczy et al. [2014] and Xu and Garcia

159    [2008], we confronted apes with two covered containers holding populations of food

160    items that differed in their proportion of preferred to neutral food (24:6 vs. 6:24). In two

161    experiments, the apes witnessed the experimenter drawing one multi-item sample from

162    each population. The distribution within the multi-item samples reflected the

163    distribution of the respective populations (e.g. 4:1 vs. 1:4). Based on the observation of

164    these representative samples, the apes were allowed to choose between the two covered

165    populations. Hence, to receive the more favorable population, they were required to use

166    proportional information provided by the samples. In Experiment 1 we tested whether

167    apes were able to reason from samples to populations. In the crucial 4:1 vs. 1:4 Test

168    condition, apes were confronted with two populations (24:6 vs. 6:24) and watched the

169    samples 4:1 vs. 1:4 being drawn from them. Two control conditions were designed to

170    rule out alternative explanations. In Control condition 1 apes did not see the available

171    populations beforehand. This manipulation tested whether the information provided by

172    the samples alone was sufficient for the apes to infer about the distribution within the

173    populations. In Control condition 2, the samples were not re-inserted into the

174    populations after the drawing process. This manipulation tested whether apes' success

175    in previous conditions might have reflected a tendency to choose the population where

176    the more favorable sample was inserted, without necessarily having to reason about the

177    drawing process. In Experiment 2, apes were tested in two further Test conditions in

178    which absolute and relative frequencies of preferred food items within samples were

179    disentangled: In the 2:1 vs. 4:8 Test condition, the absolute number of preferred food

180    items was lower in the sample drawn from the more favorable population and therefore

181  misleading. In the 4:1 vs. 4:8 Test condition, the absolute number of preferred food

182  items was the same in both samples and therefore inconclusive. Hence, to receive the

183  more favorable population in Experiment 2, apes had to take into account proportions,

184  rather than absolute numbers.

185  **Experiment 1: Can apes reason from samples to populations?**

186  In this experiment we sought to investigate whether apes were able to reason from

187  multi-item samples to populations. In the Test condition, the experimenter presented the

188  apes with two covered containers holding populations of food items (24:6 vs. 6:24).

189  After watching representative samples being drawn from those populations (4:1 vs. 1:4),

190  subjects were allowed to choose between the two containers. Two control conditions

191  tested (1) whether apes inferred from the samples alone which distribution the

192  populations had and (2) ruled out that subjects used the simple heuristic of choosing the

193  container where the more attractive sample was inserted after the sampling process (see

194  Fig 1 for an illustration of the different Test conditions). Based on the results of

195  Rakoczy et al. [2014] we expected no inter-specific differences.

196  *Methods*

197  Subjects

198      Twenty-six individuals (female N = 20) of four great ape species participated:

199  Gorillas (*Gorilla gorilla*, N = 4), Bonobos (*Pan paniscus*, N = 6), Chimpanzees (*Pan*

200  *troglodytes*, N = 10) and Orangutans (*Pongo abelii*, N = 6). One further chimpanzee

201  was tested but excluded from data analysis since he did not complete all sessions due to

202  lack of motivation. Subjects were housed at the Wolfgang Koehler Primate Research

203  Center (WKPRC) in the Leipzig Zoo and were tested between November 2014 and

204  September 2015. Their age ranged between six and 48 years (Mean = 18 years) and

205  about 25% were hand-reared. The remaining 75% were mother-reared (see Table 1 for

206  more subject information). All apes were already experienced in participating in

207  cognitive tasks with food-rewards as reinforcement. To control for potential order

208  effects, 15 of the subjects underwent Experiment 1 first and then proceeded to

209  Experiment 2, the remaining 11 subjects experienced Experiment 2 first and were tested

210  in Experiment 1 afterwards (see Fig 2).

211      The study was ethically approved by an internal committee at the Max Planck

212  Institute for Evolutionary Anthropology. Research and animal husbandry comply with

213  the "EAZA Minimum Standards for the Accommodation and Care of Animals in Zoos

214  and Aquaria", the "EEP Bonobo Husbandry Manual", the "WAZA Ethical Guidelines

215  for the Conduct of Research on Animals by Zoos and Aquariums" and the "Guidelines

216  for the Treatment of Animals in Behavioral Research and Teaching" of the Association

217  for the Study of Animal Behavior (ASAB). This research adhered to the American

218  Society of Primatologists principles for the ethical treatment of primates.

219  Materials

220      Subjects were tested individually in their sleeping cages or in special test cages.

221  A Plexiglas panel fitted on the cage mesh separated ape and experimenter. The panel

222  had two small holes ($\varnothing$ 2 cm; distance between holes 59 cm) through which subjects

223  could insert a finger to indicate a choice. Perpendicular to the Plexiglas panel, a sliding

224  table (45 x 79 cm) was mounted on the cage and could be moved both towards the

225  subject and the experimenter. Underneath the table there were two small concealed

226    compartments in which food items could be hidden prior to each test trial without the

227    subject noticing it. To prevent subjects from watching, e.g. preparation of a trial, a

228    screen (27 x 79 cm) could be fixed via metal brackets at the end of the table closer to

229    the ape's side. During test trials, apes were presented with two transparent Plexiglas

230    containers (⌀ 8 cm), each containing a population of food items, namely pieces of fruit

231    pellets and pieces of carrots of roughly equal shape and size. The containers could be

232    covered with opaque occluders of the same diameter that prevented subjects from seeing

233    the content of the containers (see Fig 3 for an illustration of the setup).

234    Design and Procedure

235        Before the actual test started, subjects underwent a familiarization session.

236    Subsequently, we carried out one test and two control conditions to investigate whether

237    apes were able to reason from multi-item samples to populations. All conditions

238    consisted of 12 test trials, divided into three sessions. Each session started with two

239    preference trials with single pellet and carrot pieces (see below). Thus, each session

240    consisted of two preference trials and four test trials.

241    *Familiarization*

242        Each subject that had not experienced Experiment 2 before received one session

243    with six trials of familiarization. In this session, the ape was confronted with one

244    transparent container holding a population of carrot and pellet pieces (distribution

245    12:12). The experimenter presented the container to the ape, shook it several times to

246    give a good overview of the population and subsequently placed it in the center of the

247    sliding table. During the first three trials the subject watched the experimenter drawing a

248    random sample (three to five items) out of the population, presenting it on the palm of

249 the hand, and re-inserting it into the container. After that, the experimenter moved the

250 container to the edge of the table and pushed the sliding table forward, so that the ape

251 could point to the container. Subsequently, the subject received the content of the

252 container as reward. During the last three trials of familiarization, the procedure was the

253 same as explained above, but this time the container was placed in an opaque occluder

254 after the ape had seen the population. Thus, the subject did not see the population during

255 the sampling process and when pointing to it. The familiarization should ensure that

256 subjects were familiar with the material and that they understood that "inserting a hand

257 in an occluded container" meant that a sample was drawn from the contained

258 population.

259 *Preference trials*

260 The preference trials aimed at assuring the apes' constant preference for one of

261 the two single-item types and were conducted prior to each of the test sessions. In each

262 trial the experimenter placed one pellet piece and one carrot piece on the sliding table

263 close to the Plexiglas panel, directly in front of the holes. The side on which the pellet

264 piece was positioned was counterbalanced. Apes indicated their choice with their finger

265 and immediately received the selected food item as reinforcement. Subsequently, the

266 test trials of the respective condition were conducted. The criterion for an ape to be

267 included in the analysis was choosing the pellet piece in at least 75% of the trials.

268 *Test trials*

269 All apes participated in three conditions. To control for a possible effect of

270 order, 15 subjects were tested in the first order of conditions (Test condition - Control 1

271 - Control 2). The remaining eleven subjects were tested in the reverse order of

272 conditions (Control 2 - Control 1 – Test condition). When we decided to split up

273 subjects in the two groups of orders, all gorillas had already been tested in the first order

274 of conditions. Thus, order was counterbalanced across subjects for all species except for

275 the four gorillas, which were all tested in the original order of conditions (see Table 1

276 for information about the order of conditions each subject experienced). In all

277 conditions, the populations consisted of 30 items each: Population A was composed of

278 24 pellet pieces and 6 carrot pieces; population B was composed of 6 pellet pieces and

279 24 carrot pieces. (These ratios were chosen because Rakoczy et al. [2014] showed that

280 apes can reliably discriminate between multiples of the ratio 4:1 vs. 1:4. To not exceed

281 the upper limit of caloric intake recommended for the apes, we had to limit the absolute

282 number of food items to a certain extent. This limitation also reduced the risk of

283 satiation and thereby helped to keep up apes´ motivation over the course of trials).

284 Test 4:1 vs. 1:4: In this condition the samples reflected the distribution of the

285 populations one-to-one. More specifically, the sample apparently drawn from

286 population A (24 pellets : 6 carrots) consisted of 4 pellet and 1 carrot pieces, and the

287 sample apparently drawn from population B (6 pellets : 24 carrots) consisted of 1 pellet

288 and 4 carrot pieces. Before a trial started, the experimenter fixed the screen on the table

289 to prevent the subject from watching preparations. Subsequently, she positioned the pre-

290 prepared multi-item samples in the small compartments underneath the table. The two

291 containers holding the populations were placed next to each other in the center of the

292 table and the two opaque occluders were positioned over them.

293 The trial started when the screen was removed from the sliding table unblocking

294 the view over the table for the subject. The experimenter simultaneously removed the

295     two occluders from the containers, and subsequently showed each population to the ape

296     by lifting the container, tilting it forward and shaking it slightly. After the subject had

297     seen both populations, the experimenter repositioned the occluders over the containers

298     and put the screen back into the metal brackets. Then she shuffled both containers.

299     Hence, subjects knew the two available populations, but did not know which population

300     was which. Revealing the populations at the beginning of each trial ensured that apes

301     were aware of both containers holding a relatively high number of food items (higher

302     than the number of items subsequently drawn). During the shuffling process, the

303     experimenter reached into the two compartments underneath the table, retrieved the

304     hidden samples and put them into her fists to make sure that the subject did not see them

305     there. After removing the screen again, the experimenter pretended to draw

306     simultaneously out of each population by inserting her fists into the two covered

307     containers and moving them around while looking upwards (maintaining a pretence of

308     random drawing). While the subject was watching, she simultaneously removed both

309     hands out of the containers and presented the samples on the palms of her hands close to

310     the Plexiglas panel saying "look!". After the ape had seen both samples, the

311     experimenter let them fall back into the containers. Subsequently, the experimenter

312     closed her eyes to minimize unintended cueing and pushed the sliding table slightly

313     forward so that each container, covered by an occluder, was positioned directly in front

314     of one of the holes. By inserting a finger into one of the holes, the ape could indicate her

315     choice, which was coded live by the experimenter after she had opened her eyes again.

316     In cases where the subject pointed towards both containers, the sliding table was pulled

317     backwards with the words "just one", and then pushed forward again, giving the ape a

318     new choice between the populations. After the ape had made her decision, the occluder

319    of the chosen container was removed, revealing the selected population. Finally, the

320    subject received the chosen population (see Fig 1 for an illustration of the procedure).

321    <u>Control 1: Samples as only source of information</u>. To investigate whether apes were

322    able to infer from the samples alone which distribution the populations most likely had,

323    we carried out Control 1, in which the subjects did not see the available populations

324    prior to the sampling process. The procedure of Control 1 was the same as in the Test

325    condition, with the following exception: In the beginning of the trials, the experimenter

326    did not remove the two occluders from the containers, preventing the apes from seeing

327    the two available populations. Instead, she shook the containers with the occluders

328    consecutively, making sure that the apes were aware of something being in the

329    containers, but leaving them in uncertainty about the exact content (see Fig 1 for an

330    illustration of the procedure).

331    <u>Control 2: No replacement of samples.</u> One alternative explanation for subjects

332    succeeding in the Test condition as well as in Control 1 could be that apes did not make

333    inferences about the drawing process and the populations as a whole, but based their

334    choices on the side where the "more attractive" sample was inserted. More specifically,

335    apes could have tracked their preferred sample and chosen the population in which this

336    sample was dropped in. To rule that out, we conducted Control 2, in which the samples

337    were not re-inserted into the populations. The procedure was the same as in the Test

338    condition, but instead of letting the samples fall back into the containers, the

339    experimenter threw them away in a bucket next to the table. Thus, the apes were

340    prevented from basing their choice on the side where the "more attractive" sample was

341    inserted and could instead use the samples only as a hint for the composition of the

342    populations (see Fig 1 for an illustration of the procedure).

343    *Follow-up tests*

344        A pre-requisite for the correct interpretation of results was that apes recognized

345    and had a preference for the population containing a higher proportion of pellet pieces.

346    Therefore we conducted two follow-up tests. Each of them was tested within a single

347    session consisting of four trials. Note that the follow-up tests were the last conditions

348    subjects underwent in this study, i.e. individuals that underwent Experiment 1 first, were

349    tested in the follow-up tests after completion of Experiment 2. Subjects that were tested

350    in Experiment 2 first, received the follow-up tests after completion of Experiment 1 (see

351    Fig 2). This was to ensure that none of the subjects had any prior experience regarding

352    the populations before starting the test.

353     "Open population"-test: In the "open population" test, apes were presented with the

354    same populations as during test conditions (A 24:6; B 6:24). For each trial, populations

355    were placed in transparent containers standing next to each other in the center of the

356    sliding table. The experimenter shook both containers successively and tilted them

357    forward to give a full view of the available populations. Once the ape had seen both

358    populations, the experimenter positioned the containers on the edge of the sliding table,

359    each in front of one of the holes. Subsequently, she pushed the table forward and the ape

360    could indicate her choice by pointing through one of the holes and received the content

361    of the chosen container. The criterion for an ape to be included in the analysis was

362    choosing the population containing more pellets in at least 75% of trials.

363 "Covered population"-test: The procedure of the "covered population" test was the

364 same as in the "open population" test, except the fact that the experimenter pulled

365 opaque occluders on the containers after the subject had seen the content. Thus, when

366 making a choice, the ape was prevented from seeing the two populations; instead she

367 had to memorize the position of her preferred population for a few seconds. This second

368 follow-up test with covered containers was conducted to test for the possibility that

369 some apes might not have been able to choose the correct container throughout the test

370 trials due to the fact that it was not visible when the choice had to be made. Subjects

371 were considered successful when they chose the pellet-population in at least 75% of

372 trials. Based on previous studies that have shown that apes can solve quantity

373 discrimination tasks that require encoding and mental comparison of quantities [e.g.

374 Call, 2000; Beran et al. , 2005], we expected that apes would be able to cope with the

375 type of stimuli occlusion involved in this test.

376 Coding and Data Analysis

377 The apes´ choice was coded live by the experimenter. A second blind observer

378 coded 25% of the trials from video. Both raters were in excellent agreement (K = 0.95,

379 N = 168). Data of five subjects (one bonobo, two chimpanzees and two gorillas, see SI

380 Table 1 for individual data) had to be excluded because those individuals did not reach

381 criterion in the follow-up tests (see above). No ape had to be excluded on the basis of

382 the preference trials. Data of all conditions were analyzed separately using R [R Core

383 Team 2014]. Subjects' choices were the dependent measure and were defined as

384 "correct" if the chosen container contained the population with the more favorable ratio

385 of pellets to carrots (24:6). The apes' overall performance (percent correct across trials)

386  was tested against chance level using a two-tailed one-sample t-test (R function t.test).

387  The effect sizes were obtained applying the package "lsr" [Navarro 2015]. In addition,

388  we tested apes' first trial performance against chance level using an exact binomial test

389  (R function binom.test) to detect potential learning effects. In order to test whether

390  performance differed between species we used a one-way ANOVA (R function aov).

391  This was justified as residuals were normally distributed and homogenous as verified by

392  visual inspection of residuals plotted against fitted values and qqplot. For Tukey's post-

393  hoc test we used the R function TukeyHSD.

394  ***Results and discussion***

395  Test 4:1 vs. 1:4

396       Apes as a group chose the more favorable population on average on 72 % of

397  trials (see Fig 4 and supplementary material Table 1 for individual data), significantly

398  more often than predicted by chance ($t(20) = 6.12$, $P < 0.001$, 95% CI [0.64, 0.79], N =

399  21; Cohen's $d = 1.34$). This pattern was also visible in the first trial performance (Mean

400  = 71 %; Binomial test: $P = 0.04$, N = 21; Cohen's $g = 0.43$). Hence, the apes'

401  performance seems to reflect an intuitive capacity rather than a learning effect. We

402  detected no difference between species (ANOVA: $F(3, 17) = 0.2$, df = 3, $P = 0.895$).

403  These results suggest that all tested species of great apes were able to intuitively use the

404  information provided by the samples to receive the preferred population, therefore

405  giving a first hint towards apes being able to reason from samples to populations.

406  Control 1: Samples as only source of information

407          Apes as a group chose the more favorable population on average on 69 % of

408     trials (see Fig 4 and supplementary material Table 1 for individual data), which is

409     significantly above chance level ($t (20) = 5.20$, $P < 0.001$, 95% CI [0.62, 0.77], $N = 21$;

410     Cohen's $d = 1.13$). However, this pattern was not found considering only the

411     performance in the first trial (Mean = 52 %; Binomial test: $P = 0.5$, $N = 21$). This is

412     perhaps best explained by insecurity about the available populations. Control 1 was the

413     only condition in which subjects did not know the two possible answers (i.e. the two

414     available populations) before making their decision. Hence, in the very first trial they

415     could not be sure whether both populations were of the same size or whether, e.g. the

416     population associated with the "worse" sample contained four times more items than the

417     population from which the "better" sample was drawn. Potentially, apes had to

418     experience during the first trial that, even though they had not seen the containers'

419     content, there were two different populations of food items with the same absolute

420     quantity. This first trial data suggest that subjects did not necessarily expect the

421     populations to be the same as in other conditions, making it unlikely that subjects had

422     learned and remembered the composition of the populations during the previous

423     session(s). We detected no difference between species (ANOVA: $F (3, 17) = 0.99$, df =

424     3, $P = 0.421$). In sum, these results show that the information provided by the samples

425     was sufficient for the apes to infer about the distribution within the populations.

426     Control 2: No replacement of samples

427          Apes as a group chose the more favorable population on average on 66 % of

428     trials (see Fig 4 supplementary material Table 1 for individual data), which is

429     significantly more often than expected by chance ($t (20) = 4.97$, $P < 0.001$, 95% CI

430  [0.59, 0.73], N = 21; Cohen's $d$ = 1.08). This pattern was also reflected in the first trial

431  performance (Mean = 76 %; Binomial test: P = 0.01, N = 21; Cohen's $g$ = 0.52) and

432  thus cannot be due to learning. In this condition we detected differences between

433  species (ANOVA: $F_{(3, 17)}$ = 4.88, df = 3, P = 0.01, $R^2$ = 0.46). Tukey multiple

434  comparison of means revealed that bonobos performed significantly worse than gorillas

435  (Mean bonobos = 53 %, N = 5; Mean gorillas = 88 %, N = 2, P = 0.015). However,

436  considering the fact that we could only include the data of two gorillas (compared to

437  five bonobos) in the final analysis, it is questionable whether this result truly reflects

438  differences between species, or rather random variation or individual differences

439  between subjects. The findings of Control 2 rule out the possibility that the apes solved

440  the task by means of a simple heuristic: "choose the container where the more attractive

441  sample was inserted". Instead, apes seem to have considered the drawing process and

442  inferred about the population as a whole.

443       In sum, the results of Experiment 1 show that all tested species of great apes

444  were able to use information provided by multi-item samples to track their preferred

445  populations, and they did so even when they did not know the composition of the

446  populations beforehand (Control 1) and when samples were not replaced after drawing

447  (Control 2). These findings suggest that great apes might engage in intuitive statistical

448  inferences from samples to populations in a comparable way human infants do [Xu and

449  Garcia 2008; Denison et al. , 2013]. However, an alternative explanation for these

450  results could be that apes simply associated the preferable sample (i.e. the sample

451  containing absolutely more pellets), with the container that it was drawn from. To

452  address this alternative explanation, we tested subjects in Experiment 2 with samples in

453  which absolute and relative frequencies of pellets were disentangled.

454

## Experiment 2: Do apes take into account relative, rather than absolute frequencies?

Although results of Experiment 1 tentatively suggest that apes were able to reason from multi-item samples to populations, it is an open question to what extent the subjects relied on absolute quantities rather than on proportions to solve the task. More specifically, in all conditions of Experiment 1, absolute and relative frequencies were confounded within the samples, i.e. the sample which contained the higher proportion of preferred food items than the alternative (4:1 vs. 1:4), also contained the higher absolute quantity of preferred food items (4 vs. 1). Thus, Experiment 1 alone cannot tease apart whether apes truly compared the proportion of pellets to carrots in both samples (4:1 versus 1:4), or if they based their choice on the absolute amount of pellets (4 vs. 1) and used the heuristic: "choose the container where more pellets were drawn from". To address this question we tested apes in Experiment 2 in two further conditions. In both of them, absolute and relative frequencies within the samples were arranged in such a way that apes could not perform above chance level if they focused on absolute numbers only (see Fig 1 for an illustration of the Test conditions).

*Methods*

Subjects

The same 26 individuals as in Experiment 1 participated in this experiment. One additional chimpanzee was tested but excluded from data analysis as he did not complete all sessions due to a lack of motivation.

476    Materials

477        We used the same materials as in Experiment 1 (see Fig 3 for an illustration of

478    the experimental setup).

479    Design and Procedure

480        The general procedure was the same as in Experiment 1. To tease apart whether

481    apes truly compared the proportion of preferred to neutral food items in both samples,

482    or if they based their choice on the absolute amount of preferred food, we tested apes in

483    two conditions with varying sample composition. Again, each condition consisted of 12

484    test trials, divided into three sessions. Prior to the test trials, two preference trials with

485    single pellet and carrot pieces were carried out. Thus, each session consisted of two

486    preference trials and four test trials.

487    *Familiarization*

488        Each subject that had not experienced Experiment 1 before received one session

489    with six trials of familiarization. The procedure of the familiarization phase was exactly

490    as described for Experiment 1.

491    *Preference trials*

492        The procedure of the preference trials was the same as in Experiment 1.

493    *Test trials*

494        All apes participated in two Test conditions. To control for a possible effect of

495    order, 15 subjects were tested in the first order of conditions, starting with the 2:1 vs.

496    4:8 test, through to the 4:1 vs. 4:8 test. The remaining eleven subjects were tested in the

497    reverse order of conditions (see Table 1 for information about the order of conditions

498    each subject experienced). Again, in all conditions the populations consisted of 30 items

499    each: Population A was composed of 24 pellet pieces and 6 carrot pieces; population B

500    was composed of 6 pellet pieces and 24 carrot pieces.

501    Test 2:1 vs. 4:8: The procedure was the same as described for the Test condition of

502    Experiment 1. However, the composition of the samples was varied in such a way that

503    choosing the container from which the sample with the higher absolute number of

504    pellets was drawn, resulted in receiving the less attractive population. In particular, the

505    sample apparently drawn from population A (24 pellets : 6 carrots) consisted of 2 pellet

506    and 1 carrot pieces, and the sample apparently drawn from population B (6 pellets : 24

507    carrots) consisted of 4 pellet and 8 carrot pieces. Thus, even though sample B contained

508    double the amount of pellets compared to sample A, the proportion of pellets to carrots

509    was more favorable in sample A. If apes´ choice was based on absolute quantities, we

510    expected them to choose the "wrong" container more often than the "correct" one. If

511    they, however, took into account the proportion of pellets to carrots, we expected them

512    to choose the "correct" container more often than the foil (see Fig 1 for an illustration of

513    the procedure).

514    Test 4:1 vs. 4:8: Again, the procedure was the same as described for the Test condition

515    of Experiment 1.  However, here the composition of the samples was varied in a way

516    that both samples contained the same absolute number of pellets. More specifically, the

517    sample apparently drawn from population A (24 pellets : 6 carrots) consisted of 4 pellet

518    and 1 carrot pieces, and the sample apparently drawn from population B (6 pellets : 24

519   carrots) consisted of 4 pellet and 8 carrot pieces. Assuming that apes based their choice

520   on absolute quantities only, we expected them to choose both containers at similar rates,

521   as the absolute number of pellets did not provide any conclusive information. If they

522   instead reasoned about the proportion of pellets to carrots, we predicted that they chose

523   the correct container more often than expected by chance (see Fig 1 for an illustration of

524   the procedure).

525   *Follow-up tests*

526   Those individuals that underwent Experiment 2 after Experiment 1 received the

527   two follow-up tests. The procedure was exactly the same as described for Experiment 1.

528   Coding and Data Analysis

529   The apes´ choice was coded live by the experimenter. A second blind observer

530   coded 25% of the trials from video. Both raters were in excellent agreement (K= 0.95, N

531   = 120). Data of five subjects (one bonobo, two chimpanzees and two gorillas, see SI

532   Table 1 for individual data) had to be excluded because those individuals did not reach

533   criterion in the follow-up tests. No further ape had to be excluded on the basis of the

534   preference trials. Data analysis was the same as described for Experiment 1.

535   **Results and discussion**

536   Test 2:1 vs. 4:8

537   Apes as a group chose the more favorable population on average on 44 % of

538   trials (see Fig 4 and supplementary material Table 1 for individual data).  Though this

539   pattern is not different from what was expected by chance ($t (20) = -1.84$, $P = 0.08$, 95%

540    CI [0.36, 0.51], N = 21), it indicates a (non-significant) trend such that apes tended to

541    choose the less favorable population more often than the more favorable one. We

542    detected no differences between species (ANOVA: F (3, 17) = 1.66, df = 3,P = 0.213).

543    This pattern was also reflected in the first trial performance (Mean = 47 %; Binomial

544    test: P = 1, N = 21). Hence, all tested species of great apes were unable to extrapolate

545    from samples to populations, when the absolute number of preferred food-items was

546    misleading. Instead, they tended to choose the population where the sample with the

547    higher amount of preferred food-items was drawn from. This finding gives a first hint

548    that the strategy applied by the apes might have been a comparison of absolute numbers

549    between samples, rather than an extrapolation of proportions.

550    Test 4:1 vs. 4:8

551        Apes as a group chose the more favorable population on average on 51 % of

552    trials (see Fig 4 and supplementary material Table 1 for individual data), which is not

553    different from chance level (t (20) = 0.37, P = 0.715, 95% CI [0.44, 0.58], N = 21). We

554    detected no differences between species (ANOVA: F (3, 17) = 1.35, df = 3, P = 0.292).

555    The same pattern was found considering only the performance in the first trial (Mean =

556    43 %; Binomial test: P = 0.664, N = 21). This implies that apes failed to use the

557    information provided by the samples to reason about the populations and strengthens the

558    theory that apes might have relied on absolute, rather than relative frequencies.

559    **General discussion**

560        In Experiment 1, we investigated whether great apes are able to reason from

561    multi-item samples to populations of food items. Results showed that great apes did

562     extrapolate from samples to populations, irrespective of whether they knew the

563     composition of the available populations beforehand or not (Control 1) and if samples

564     were replaced after drawing or not (Control 2). The results of Control 2 are especially

565     revealing, as they rule out the possibility of a simple heuristic: "choose the container

566     where the more attractive sample was inserted". Instead, apes seem to have considered

567     the drawing process and inferred about the population as a whole from the first trial

568     onwards. This implies that apes seem to possess similar kinds of capacities as found in

569     human infants [Xu and Garcia, 2008; Denison et al. , 2013]. In fact, our findings even

570     go one step further than those of the two existing studies that tested infants' ability to

571     reason from sample to population: While the apes in our study drew inferences from

572     samples to populations in an active choice paradigm, the human infants in the above

573     mentioned studies were only tested using the VOE looking-time paradigm. There is

574     some evidence that findings of studies using the VOE looking time paradigm dissociate

575     from findings of studies using active choice measures [e.g. Ahmed and Ruffman, 1998;

576     Shinskey and Munakata, 2005; Charles and Rivera, 2009]. This is probably due to the

577     fact that a subject that is able to perceive something is not necessarily able to act

578     accordingly. As it is currently unknown whether human infants would succeed in an

579     active choice paradigm testing for their capacities to reason from sample to population,

580     we conclude that great apes' intuitive statistical abilities in this regard seem to be at

581     least at a comparable level as those of young human infants. However, based on

582     Experiment 1 alone it is impossible to rule out that apes used alternative strategies based

583     on the absolute number of preferred food items. The aim of Experiment 2, therefore,

584     was to investigate whether great apes can successfully reason from samples to

585     populations when prevented from relying on absolute quantities. Apes performed at

586 chance level both when the sample drawn from the more favorable population contained

587 less preferred food items than the sample drawn from the less favorable population, and

588 when both samples contained the same number of preferred food items. Thus, apes did

589 not rely on inferences from samples to populations in this experiment. There are at least

590 two interpretations for these findings.

591       One interpretation is that apes' failure in Experiment 2 reflects true limitations of

592 their cognitive competences. The most obvious difference between Experiment 1 and 2

593 is that only in the latter subjects could not rely on absolute numbers of preferred food

594 items. Hence, one could conclude that apes are able to reason and draw inferences about

595 absolute, but not relative frequencies. Assuming that apes simply compared the absolute

596 quantity of pellets in both samples and chose the population from which more pellets

597 were drawn, we expected the following pattern of results: When the number of pellets in

598 the samples was inconclusive (because it was the same in both samples), apes should

599 have chosen randomly between both populations. When the number of pellets was

600 misleading, i.e. higher in the sample drawn from the non-preferred population, apes

601 should have chosen the "wrong" population more often. While apes indeed chose

602 randomly between populations when the number of pellets was the same in both

603 samples, they also did so when the number of pellets was misleading. Yet, it should be

604 noted that even though there was no significant effect in this condition (misleading

605 number of pellets in both samples), apes nevertheless revealed a non-significant

606 tendency to choose the more favorable population less often than the more favorable

607 one. Consequently, it cannot be ruled out that apes mainly relied on absolute quantities

608 in this experiment.

609         This opens up an alternative explanation for the apes´ success in Experiment 1:

610 Subjects might have not drawn any inference from sample to population, but instead

611 simply associated the more favorable sample (i.e. the one containing absolutely more

612 preferred items than the other) with the container it was drawn from, since it was

613 temporally and spatially most closely associated with that container. In other words,

614 apes might have followed a heuristic like "chose the container where you saw

615 something good (i.e. more pellets) coming from". Future studies need to determine

616 whether subjects truly relied on associating containers with "better" and "worse", or if

617 they in fact perceived the samples as a representation of populations. One possible way

618 to disentangle the two explanations would entail presenting apes with two opaque

619 containers filled with two populations of food items (similar to the current study).

620 Crucially, the experimenter would already have the samples (i.e., pellets and carrots in

621 4:1 distribution in one hand, 1:4 in the other) in her hands. She would then show the

622 contents of her hands to the ape, insert her hands into the containers and remove them

623 again, showing the same items as before. Subsequently, she would discard the

624 "samples" and give the apes the choice between the two containers.  If apes merely

625 associated the two containers with "good" or "bad" according to the distribution they

626 had seen on each side, we would expect them to choose the side where the "sample"

627 with absolutely more pellets was shown. In contrast, if they recognized a randomly

628 drawn sample as representation of the population, they should pick both containers

629 equally often since no drawing took place, and therefore, no inference can be made.

630         Recall that Rakoczy et al. [2014] showed that great apes did take proportions

631 into account when reasoning the other way around, i.e. from populations to samples,

632 ruling out that subjects used a simple association mechanism to solve the task. If our

633    results reflected true limitations in apes' cognitive competences, they would, therefore,

634    suggest that nonhuman primates' statistical abilities could be unidirectional. This would

635    question whether apes have a true understanding of drawing processes and the relation

636    between populations and samples.

637        A different interpretation for the negative findings of Experiment 2 is that they

638    may merely reflect performance limitations imposed by the task's cognitive demands,

639    which may have masked apes' true competence. One of these task demands could be the

640    memory component required by our procedure. At the exact moment when apes were

641    asked to make a choice, the information necessary to do so (i.e. the samples) was not

642    available anymore. Instead, apes had to memorize this information for a few seconds

643    and recall it to choose between the two populations.  Note that this was not the case in

644    Rakoczy et al. [2014], where subjects were still able to see the populations during their

645    choice. Even though it may seem trivial to remember information for a few seconds,

646    results of the follow-up test with covered populations showed that this was indeed a

647    crucial factor for some of the subjects: Four of the 26 subjects were not able to choose

648    the more attractive population when it was covered while the decision was made, even

649    though they showed a clear preference for that population during the preference test

650    with open populations. Furthermore, other studies have shown the importance of

651    working memory in different problem solving tasks. For instance, in Seed et al. [2012]

652    four chimpanzees solved a tool-use task requiring causal inferences when the time-span

653    over which information had to be memorized was minimized. By contrast, in a related

654    previous study [Povinelli, 2000] that involved a higher working memory load, all

655    chimpanzees failed to do so. Although working memory demands, potentially in

656    combination with lack of attention, may have influenced the apes´ performance to a

657    certain extent, working memory alone cannot fully explain the fact that apes were not

658    able to use proportional information in this experiment. Recall that those subjects who

659    had difficulties remembering the populations' position were excluded from the analysis

660    and did therefore not bias the results in a negative way. Moreover, Experiment 1 also

661    required a memory component, and still subjects succeeded.

662          Another factor that could have made this task more difficult as compared to

663    Rakoczy et al. [2014] is the type of inferences required. Retrospective inferences seem

664    to be harder than prospective ones [Völter and Call, 2017]. This means that going from

665    samples back to populations (retrospective) may be more demanding than going from

666    populations forward to samples (prospective). The majority of knowledge that we have

667    about the origin and development of intuitive statistics derives from the extensive study

668    of pre-verbal infants. In the last decades, numerous such studies have tested infants both

669    for their abilities in reasoning from populations to samples as well as from samples to

670    populations. As mentioned above, to our knowledge there is no study testing pre-verbal

671    infants for their ability to reason from samples to populations in an active choice

672    measure. This type of methodology was, so far, only used in studies investigating

673    infants´ capacity to reason from population to sample [Feigenson et al. , 2002; Denison

674    and Xu, 2010b; Denison and Xu, 2014]. In these studies, infants were allowed to choose

675    between the covered samples of two populations of preferred and non-preferred items in

676    different ratios. Control conditions disentangled absolute and relative frequencies with

677    the result that infants indeed used proportional information, not a comparison of

678    absolute quantities, to retrieve their preferred item. The two existing studies

679    investigating the reverse ability, i.e. reasoning from samples to populations [Xu and

680    Garcia, 2008; Denison et al. , 2013], both used a VOE looking-time paradigm, a

681    methodology that is less comparable with the methodology applied for great apes.

682    Moreover, in both above-mentioned studies probability was confounded with quantity,

683    and no control condition tested for the fact that infants could have used the shortcut of

684    focusing on absolute quantities only. As a consequence, it remains unclear whether

685    reasoning from samples to populations represents a cognitively more challenging task

686    than the other way around. It would be of great interest to fill that gap of knowledge by

687    applying an active choice paradigm to investigate pre-verbal infants´ ability to reason

688    from samples to populations, including a control condition for absolute vs. relative

689    information.

690       A third task demand that may have masked apes' true competence in Experiment

691    2 is the poorer discriminability of the samples as compared to the samples used in

692    Experiment 1. As an index for discriminability we calculated the ratio of ratios

693    (hereafter: ROR) of the two samples for each of the conditions in the following way

694    [following Drucker et al. , 2016]:

695   
$$\frac{\text{Ratio of pellets to carrots in the sample drawn from the preferred population}}{\text{Ratio of pellets to carrots in the sample drawn from the non} - \text{preferred population}}$$

696       In all conditions of Experiment 1, the ROR was (4/1)/(1/4)=16 (in Rakoczy et al.

697    2014 the ROR was ≥ 16 in all conditions). In Experiment 2, the ROR was (2/1)/(4/8)=4

698    in the 2:1 vs. 4:8 test, and (4/1)/(4/8)=8 in the 4:1 vs. 4:8 test. Thus, in both conditions

699    of Experiment 2, the ROR was less than or equal to half the one used in Experiment 1.

700    This discrepancy was caused by our methodological constraints that prevented us from

701    using larger RORs. More specifically, a larger ROR would have required larger samples

702    and thus larger populations. As the apes received the *whole* chosen population as

703    reinforcement we had to minimize the number of food items within the populations for

704    the purpose of not exceeding their allowed daily caloric intake. Moreover, given that the

705    food items were kept in the experimenter´s fist, larger samples would have required a

706    different sampling method than the one applied here. As a consequence, in this study it

707    was not possible to disentangle absolute and relative information with the same ROR as

708    in Experiment 1. Recent research suggests that indeed the magnitude of difference

709    between two proportions is crucial for non-human primates to discriminate

710    probabilities. Hanus and Call [2014] presented chimpanzees with two trays, each of

711    them with a different ratio of hidden food items to potential hiding locations and

712    therefore a different likelihood of finding food. This study revealed that subjects´

713    performance was influenced by the relative difference between the two probabilities as

714    soon as a certain threshold thereof was reached. Moreover, the apes relied on the ratio

715    between probabilities, even in conditions where one tray depicted an absolute safe

716    option— a probability of finding food of 100%. This study emphasizes the importance

717    of the magnitude of difference between the two ratios to be discriminated, rather than

718    the magnitude of difference within the single ratios.

719        With regard to the present study this means the following: Although the

720    quantities within one sample were presumably easy to discriminate [for reviews about

721    quantity discrimination see e.g. Feigenson et al. , 2004; Nieder, 2005], it was probably

722    the ratio between the ratios of both samples that influenced the decision of the apes and

723    it could well be that the present RORs were simply below the threshold for

724    discriminating two ratios and thus failed to constitute notable differences. In a study

725    using a touch screen setup [Drucker et al. , 2016] rhesus macaques (*Macaca mulatta*)

726    were presented with arrays containing different ratios of positive to negative stimuli.

727    The monkeys learned to choose those arrays with the greater ratio of positive to

728    negative stimuli and were able to generalize to novel ratios. Similarly as in the

729    previously mentioned study with chimpanzees [Hanus and Call, 2014], the performance

730    was directly influenced by the magnitude of difference between the two ratios to be

731    discriminated. Interestingly, just as human infants [McCrink and Wynn, 2007], the two

732    macaques tested were able to discriminate a ROR of 2, which is much lower than those

733    used in our experiments. However, given the fact that those subjects received extensive

734    training in such discrimination tasks before the actual test, it remains unclear to which

735    extent those methods are comparable to the ones used here with apes.

736    *Conclusion*

737        The aim of the current study was to investigate whether apes can use samples of

738    items to infer the composition of the population from where the samples came from.

739    While apes performed competently when the samples from the more favorable

740    population were more attractive than the samples from the less favorable population not

741    only in terms of relative but also in terms of absolute frequencies of preferred over non-

742    preferred food items, they failed to do so when absolute and relative frequencies were

743    disentangled. The present study, therefore, cannot determine whether non-human

744    primates engage in intuitive statistical inferences from randomly drawn samples to

745    populations in a comparable way human infants have recently been found to do [Xu and

746    Garcia, 2008; Denison et al. , 2013]. It is an open question for future research whether

747    these limitations in apes' performance reflect true limits of cognitive competence or

748    merely performance limitations due to accessory task demands.

749 **Acknowledgements**

756 **References**

757

758 Ahmed A, Ruffman T. 1998. Why do infants make A not B errors in a search task, yet show

759   memory for the location of hidden objects in a nonsearch task? Developmental

760   Psychology 34(3):441-453.

761 Barnard AM, Hughes KD, Gerhardt RR et al. 2013. Inherently Analog Quantity Representations

762   in Olive Baboons (Papio anubis). Frontiers in Psychology 4:253.

763 Benson-Amram S, Heinen VK, Dryer SL, Holekamp KE. 2011. Numerical assessment and

764   individual call discrimination by wild spotted hyaenas, *Crocuta crotua*. Animal

765   Behaviour 82: 743-752.

766 Beran MJ, Beran MM, Harris EH, Washburn DA. 2005. Ordinal judgments and summation of

767   nonvisible sets of food items by two chimpanzees (*Pan troglodytes*) and a rhesus

768   macaque (*Macaca mulatta*). Journal of Experimental Psychology: Animal Behavior

769   Processes 31: 351-362.

770 Beran MJ, Evans TA, Leighty KA, Harris EH, Rice D. 2008. Summation and quantity

771   judgments of sequentially presented sets by capuchin monkeys (Cebus apella).

772   American Journal of Primatology 70(2):191-4.

773   Beran MJ, McIntyre JM, Garland A, Evans TA. 2013. What counts for 'counting'?
774        Chimpanzees, Pan troglodytes, respond appropriately to relevant and irrelevant
775        information in a quantity judgment task. Animal Behaviour 85(5):987-993.

776   Beran MJ, Parrish AE. 2016. Capuchin monkeys (Cebus apella) treat small and large numbers
777        of items similarly during a relative quantity judgment task. Psychonomic Bulletin and
778        Review 23(4):1206-13.

779   Boysen ST, Berntson GG. 1989. Numerical competence in a chimpanzee (Pan troglodytes).
780        Journal of Comparative Psychology 103(1):23-31.

781   Call J. 2000. Estimating and operating on discrete quantities in orangutans (Pongo pygmaeus).
782        Journal of Comparative Psychology 114(2):136-147.

783   Charles EP, Rivera SM. 2009. Object permanence and method of disappearance: looking
784        measures further contradict reaching measures. Developmental Science 12(6):991-1006.

785   Cosmides L, Tooby J. 1996. Are humans good intuitive statisticians after all? Rethinking some
786        conclusions from the literature on judgment under uncertainty. Cognition 58(1):1-73.

787   Dacke M, Srinivasan MV. 2008. Evidence for counting in insects. Animal Cognition 11(4):683-
788        689.

789   Davis H. 1984. Discrimination of the number three by a raccoon. Animal Learning & Behavior
790        12: 409-413.

791   Denison S, Reed C, Xu F. 2013. The emergence of probabilistic reasoning in very young
792        infants: evidence from 4.5- and 6-month-olds. Developmental Psychology 49(2):243-9.

793   Denison S, Trikutam P, Xu F. 2014. Probability versus representativeness in infancy: can
794        infants use naive physics to adjust population base rates in probabilistic inference?
795        Developmental Psychology 50(8):2009-19.

796   Denison S, Xu F. 2010a. Integrating physical constraints in statistical inference by 11-month-
797        old infants. Cognitive Science 34(5):885-908.

798   Denison S, Xu F. 2010b. Twelve- to 14-month-old infants can predict single-event probability
799        with large set sizes. Developmental Science 13(5):798-803.

800     Denison S, Xu F. 2012. Probabilistic inference in human infants. Adv Child Dev Behav 43:27-
801           58.

802     Denison S, Xu F. 2014. The origins of probabilistic inference in human infants. Cognition
803           130(3):335-47.

804     Drucker CB, Rossa MA, Brannon EM. 2016. Comparison of discrete ratios by rhesus macaques
805           (Macaca mulatta). Animal Cognition 19(1):75-89.

806     Farnsworth GL, Smolinski JL. 2006. Numerical discrimination by wild Northern Mockingbirds.
807           Condor 108(4):953-957.

808     Feigenson L, Carey S, Hauser M. 2002. The representations underlying infants' choice of more:
809           object files versus analog magnitudes. Psychological Science 13(2):150-6.

810     Feigenson L, Dehaene S, Spelke E. 2004. Core systems of number. Trends in Cognitive
811           Sciences 8(7):307-14.

812     Geary DC, Berch BB, Mann Koepke K (2015). The Evolution of Number Systems In: Geary
813           DC, Berch DC, Mann Koepke K, editors. Mathematical Cognition and Learning
814           Volume 1: Evolutionary Origins and Early Development of Number Processing.
815           London, San Diego, Waltham, Oxford: Elsevier. pp 335-352.

816     Hamilton IM (2010) Foraging theory. In: Westneat DF, Fox CW, editors. Evolutionary
817           behavioral ecology. New York: Oxford University Press. pp 177-193.

818     Hanus D, Call J. 2007. Discrete quantity judgments in the great apes (Pan paniscus, Pan
819           troglodytes, Gorilla gorilla, Pongo pygmaeus): the effect of presenting whole sets
820           versus item-by-item. Journal of Comparative Psychology 121(3):241-9.

821     Hanus D, Call J. 2014. When maths trumps logic: probabilistic judgements in chimpanzees.
822           Biological Letters 10:20140892 doi:10.1098/rsbl.2014.0892.

823     Holland PW. 1986. Statistics and Causal Inference. Journal of the American Statistical
824           Association 81(396):945-960.

825     Hunt S, Low J, Burns KC. 2008. Adaptive numerical competency in a food-hoarding songbird.
826           Proceedings of the Royal Society B: Biological Sciences 275(1649):2373-2379.

827    McComb K, Packer C, Pusey A. 1994. Roaring and Numerical Assessment in Contests Between

828         Groups of Female Lions, *Panthera leo*. Animal Behaviour 47(2):379-387.

829    McCrink K, Wynn K. 2007. Ratio abstraction by 6-month-old infants Psychological Sciences

830         18:740-745 doi:10.1111/j.1467-9280.2007.01969.x

831    Navarro DJ. 2015. Learning statistics with R: A tutorial for psychology students and other

832         beginners (Version 0.5). University of Adelaide. Adelaide, Australia.

833    Nieder A. 2005. Counting on neurons: The neurobiology of numerical competence. Nature

834         Reviews Neuroscience 6(3):177-190.

835    Perdue BM, Talbot CF, Stone AM, Beran MJ. 2012. Putting the elephant back in the herd:

836         relative quantity judgments match those of other species. Animal Cognition 15: 955-

837         961.

838    Piaget J, Inhelder Br. 1975. The origin of the idea of chance in children. New York,: Norton. xx,

839         251 p. p.

840    Pisa P, Agrillo C. 2009. Quantity discrimination in felines: a preliminary investigation of the

841         domestic cat (Felis silvestris catus). Journal of Ethology 27(2):289-293.

842    Potrich D, Sovrano VA, Stancher G, Vallortigara G. 2015. Quantity discrimination by zebrafish

843         (Danio rerio). Jorunal of Comparative Psychology 129(4):388-93.

844    Povinelli DJ. 2000. Folk physics for apes : the chimpanzee's theory of how the world works.

845         Oxford University Press, Oxford ; New York.

846    R Core Team. 2014. A language and environment for statistical computing. R Foundation for

847         Statistical Computing, Vienna, Austria.

848    Rakoczy H, Cluver A, Saucke L et al. 2014. Apes are intuitive statisticians. Cognition

849         131(1):60-8.

850    Reznikova Z, Ryabko B. 2011. Numerical competence in animals, with an insight from ants.

851         Behaviour 148(4):405-434.

852  Rugani R, Cavazzana A, Vallortigara G, Regolin L. 2013. One, two, three, four, or is there
853      something more? Numerical discrimination in day-old domestic chicks. Animal
854      Cognition 16(4):557-64.

855  Seed A, Seddon E, Greene B, Call J. 2012. Chimpanzee 'folk physics': bringing failures into
856      focus. Philosophical Transactions of the Royal Society B Biological Sciences 367:2743-
857      2752 doi:10.1098/rstb.2012.0222.

858  Shinskey JL, Munakata Y. 2005. Familiarity breeds searching - Infants reverse their novelty
859      preferences when reaching for hidden objects. Psychological Science 16(8):596-600.

860  Skyrms B. 1975. Choice and chance : an introduction to inductive logic. Encino, Calif.:
861      Dickenson Pub. Co. 220 p. p.

862  Tecwyn EC, Denison S, Messer EJ, Buchsbaum D. 2016. Intuitive probabilistic inference in
863      capuchin monkeys. Animal Cognition 130 (3): 335-347.

864  Teglas E, Girotto V, Gonzalez M, Bonatti LL. 2007. Intuitions of probabilities shape
865      expectations about the future at 12 months and beyond. Proceedings of the National
866      Academy of Sciences USA 104(48):19156-9.

867  Teglas E, Vul E, Girotto V et al. 2011. Pure reasoning in 12-month-old infants as probabilistic
868      inference. Science 332(6033):1054-9.

869  Tenenbaum JB, Griffiths TL, Kemp C. 2006. Theory-based Bayesian models of inductive
870      learning and reasoning. Trends in Cognitive Sciences 10(7):309-18.

871  Tversky A, Kahneman D. 1974. Judgment under Uncertainty - Heuristics and Biases. Science
872      185(4157):1124-1131.

873  Tversky A, Kahneman D. 1981. The Framing of Decisions and the Psychology of Choice.
874      Science 211(4481):453-458.

875  Vonk J, Beran MJ. 2012. Bears "count" too: quantity estimation and comparison in black bears,
876      *Ursus americanus*. Animal Behaviour 84: 231-238.

877  Völter CJ, Call J. (2017). Causal and inferential reasoning in animals. In J. Call (Ed. ), *APA*
878      Handbook of Comparative Psychology.

879    Ward C, Smuts BB. 2007. Quantity-based judgments in the domestic dog (Canis lupus

880        familiaris). Animal Cognition 10(1):71-80.

881    Wilson ML, Britton NF, Franks NR. 2002. Chimpanzees and the Mathematics of Battle.

882        Proceedings of the Royal Society of London B Biological Sciences 269(1496):1107-

883        1112.

884    Xu F, Denison S. 2009. Statistical inference and sensitivity to sampling in 11-month-old infants.

885        Cognition 112: 97-104.

886    Xu F, Garcia V. 2008. Intuitive statistics by 8-month-old infants. Proceedings of the National

887        Academy of Sciences USA 105(13):5012-5.

888

889