

1 Great ape cognition is structured by stable cognitive abilities and predicted by  
2 developmental conditions

3 Manuel Bohn<sup>1</sup>, Johanna Eckert<sup>1</sup>, Daniel Hanus<sup>1</sup>, Benedikt Lugauer<sup>2</sup>, Jana Holtmann<sup>2</sup>, &  
4 Daniel Haun<sup>1, 3</sup>

5 <sup>1</sup> Department of Comparative Cultural Psychology, Max Planck Institute for Evolutionary  
6 Anthropology, Leipzig, Germany

7 <sup>2</sup> Wilhelm Wundt Institute of Psychology, Leipzig University, Leipzig, Germany

8 <sup>3</sup> Leipzig Research Centre for Early Child Development, Leipzig University, Leipzig,  
9 Germany

<sup>11</sup> The authors made the following contributions. Manuel Bohn: Conceptualization,  
<sup>12</sup> Formal Analysis, Writing - Original Draft Preparation, Writing - Review & Editing;  
<sup>13</sup> Johanna Eckert: Conceptualization, Writing - Original Draft Preparation, Writing -  
<sup>14</sup> Review & Editing; Daniel Hanus: Conceptualization, Writing - Original Draft Preparation,  
<sup>15</sup> Writing - Review & Editing; Benedikt Lugauer: Formal Analysis, Writing - Original Draft  
<sup>16</sup> Preparation, Writing - Review & Editing; Jana Holtmann: Formal Analysis, Writing -  
<sup>17</sup> Original Draft Preparation, Writing - Review & Editing; Daniel Haun: Conceptualization,  
<sup>18</sup> Writing - Review & Editing.

<sup>19</sup> Correspondence concerning this article should be addressed to Manuel Bohn, Max  
<sup>20</sup> Planck Institute for Evolutionary Anthropology, Deutscher Platz 6, 04103 Leipzig,  
<sup>21</sup> Germany. E-mail: manuel\_bohn@eva.mpg.de

22

## Abstract

23 Great ape cognition is used as a reference point to specify the evolutionary origins of  
24 complex cognitive abilities, including in humans. This research often assumes that great  
25 ape cognition consists of cognitive abilities (traits) that account for stable differences  
26 between individuals which change and develop in response to experience. Here, we test the  
27 validity of these assumptions by assessing repeatability of cognitive performance among  
28 captive great apes (gorilla gorilla, pongo abelii, pan paniscus, pan troglodytes) in five tasks  
29 covering a range of cognitive domains. We examine whether individual characteristics (age,  
30 group, test experience) or transient situational factors (life events, testing arrangements, or  
31 sociality) influence cognitive performance. Our results show that task-level performance is  
32 generally stable over time; four of the five tasks were reliable measurement tools.

33 Performance in the tasks was best explained by stable differences in cognitive abilities  
34 (traits) between individuals. Cognitive abilities were further correlated, suggesting shared  
35 cognitive processes. Finally, when predicting cognitive performance, we found stable  
36 individual characteristics to be more important than variables capturing transient  
37 experience. Taken together, this study shows that great ape cognition is structured by  
38 stable cognitive abilities that respond to different developmental conditions.

39       *Keywords:* cognition, evolution, comparative psychology, great apes, individual  
40 differences

41       Word count: X

42 Great ape cognition is structured by stable cognitive abilities and predicted by  
43 developmental conditions

44 **Introduction**

45 In their quest to understand the evolution of cognition, anthropologists,  
46 psychologists, and cognitive scientists face a major obstacle: cognition does not fossilize.  
47 Instead of studying the cognitive abilities of, e.g., extinct early hominins directly, we have  
48 to rely on inferences. We can, for example, study fossilized skulls and crania to  
49 approximate brain size and structure and use this information to infer cognitive abilities<sup>1,2</sup>.  
50 We can also study the material culture left behind by extinct species and try to infer its  
51 cognitive complexity<sup>3–5</sup>. Yet, the archaeological record is sparse and only goes back so far.  
52 Thus, additionally, we rely on backward inference about a last common ancestor based on  
53 the phylogenetically informed comparison of extant species. The so-called comparative  
54 method is one of the most fruitful approaches to investigating cognitive evolution. If  
55 species A and B both show cognitive ability X, the last common ancestor of A and B most  
56 likely also had ability X<sup>6–9</sup>. In this way, similarities and differences between species are  
57 used to make inferences about points of divergence in the evolutionary tree as well as about  
58 external drivers of this divergence. Following this approach, comparing humans to  
59 nonhuman great apes has been highly productive and provides the empirical basis for  
60 numerous theories about human cognitive evolution<sup>10–15</sup>.

61 The use of cross-species comparisons to make backward inferences about human  
62 cognitive evolution relies on a particular view of the nature and structure of great ape  
63 cognition. Cognition is seen as structured in the form of cognitive abilities that account for  
64 stable differences between individuals and which evolve and develop in response to  
65 enduring social and environmental conditions. Such differences in cognitive abilities are  
66 involved in generating variation in the behavior on which selection can act<sup>16,17</sup>. Without a  
67 stable cognitive basis that is systematically linked to behavior, cognitive evolution is not

68 possible – at least not in the way it is commonly theorized about<sup>18</sup>. In this study, we seek  
69 to provide empirical answers to a series of questions asking whether this view on great ape  
70 cognition holds. Alternatively, performance in cognitive tasks could be largely determined  
71 by transient situational factors and not capture stable abilities of individuals. Thus,  
72 because cognitive abilities cannot directly be observed, asking questions questions about  
73 the structure of great ape cognition inevitably comes with asking questions about the tools  
74 – experimental tasks – that are used to measure it.

75 The first question is whether studies on great ape cognition produce robust results:  
76 inferences about the cognitive abilities of great apes – as a clade, species, group or  
77 individual – should remain the same across repeated studies with different individuals or  
78 follow predictable patterns in studies with the same individuals. This is a critical  
79 requirement to build theories around the results of cross-species comparisons. In practice,  
80 the robustness of aggregated results is implicitly assumed but rarely tested<sup>19–21,21,22</sup>.

81 The second question is whether there are stable differences between individuals and  
82 whether tasks commonly used in great ape cognition research are able to reliably measure  
83 them. This is a prerequisite to investigate the extent to which differences between  
84 individuals in one ability co-vary with differences in other abilities to map out the internal  
85 structure of great ape cognition<sup>18,23–25</sup>. Once again, in practice, this is simply assumed to  
86 be the case but rarely tested empirically.

87 Finally, we ask which social and environmental conditions influence cognition. That  
88 is, we look for individual characteristics or everyday experiences that predict performance  
89 in our measures of cognitive ability. On the one hand, such predictive relationships inform  
90 us about the nature of cognitive performance: is it heavily influenced by transient and  
91 situational factors or malleable to long-term experiences? On the other hand, they inform  
92 us about the contexts in which cognitive abilities emerge and are the cornerstone for  
93 theorising about the ontogeny and phylogeny of cognitive abilities<sup>26,27</sup>. To summarise, to

94 date we know too little about the structure of great ape cognition to judge the validity of  
95 the comparative method as a way to study the origins of of human cognition.

96 There are several studies that provide a more comprehensive picture of one or more  
97 aspects of the nature and structure of great ape cognition<sup>24,28–32</sup>. Herrmann and  
98 colleagues<sup>33</sup> tested more than one hundred great apes (chimpanzees and orangutans) and  
99 human children in various tasks covering numerical, spatial, and social cognition. The  
100 results indicated pronounced group-level differences between great apes and humans in the  
101 social but not the spatial or numerical domain. Furthermore, relationships between the  
102 tasks pointed to a different internal structure of cognition, with a distinct social cognition  
103 factor for humans but not great apes<sup>34,35</sup>. Völter and colleagues<sup>36</sup> focused on the structure  
104 of executive functions. Using a multi-trait multi-method approach<sup>37</sup>, they developed a new  
105 test battery to assess memory updating, inhibition, and attention shifting in chimpanzees  
106 and human children. Overall, they found low correlations between tasks and, thus, no clear  
107 support for structures put forward by theoretical models built around adult human data.

108 Beyond great-apes, there have been numerous attempts to investigate the structure of  
109 cognition in other animals<sup>24</sup>. In many cases, test batteries have been used in order to find  
110 evidence for a “general cognitive ability”, i.e., a correlation of individual performance  
111 across tasks<sup>38–42</sup>. Such studies found consistent individual differences across two or more  
112 tasks in various species (e.g., insects<sup>43,44</sup>, rodents<sup>45–47</sup>, birds<sup>48,49</sup>). Some even correlated  
113 these differences with individual characteristics such as sex or relatedness<sup>43,44,47</sup>.

114 Despite their contributions to understanding the nature and structure of animal and  
115 great ape cognition, these studies suffer from one or more of the shortcomings outlined  
116 above: It is unclear if the results are robust. If the same individuals were tested again,  
117 would the results license the same conclusions about absolute differences between species?  
118 Furthermore, the psychometric properties of the tasks are unknown and it is thus unclear  
119 if, for example, low correlations between tasks reflect a genuine lack of shared cognitive

120 processes or simply measurement imprecision. Most importantly, which characteristics and  
121 experiences predict cognitive performance remains unclear. Establishing such a link is  
122 essential if we want to understand cognitive abilities and the driving forces behind their  
123 emergence and development.

124 The studies reported here address the shortcomings outlined above and seek to  
125 solidify the empirical grounds on which the use of the comparative method for investigating  
126 the evolution of human cognition rests. For one-and-a-half years, every two weeks, we  
127 administered a set of five cognitive tasks (see Figure 1) to the same population of great  
128 apes ( $N = 43$ ). The tasks spanned cognitive domains and were based on published  
129 procedures widely used in comparative psychology. As a test of social cognition, we  
130 included a gaze following task<sup>50</sup>. To assess causal reasoning abilities, we had a direct causal  
131 inference and an inference by exclusion task<sup>51</sup>. Numerical cognition was tested using a  
132 quantity discrimination task<sup>52</sup>. Finally, as a test of executive functions, we included a delay  
133 of gratification task<sup>53</sup> (Second half of the study only – see below. In the first half, we used  
134 a different measure of executive functions<sup>54</sup> . This task, however, failed to produce  
135 meaningful results. See section *Tasks - Switching* in the Supplementary Material and  
136 Supplementary Figures 8 and 9 for details).

137 In addition to the cognitive data, we continuously collected 14 variables that capture  
138 stable and variable aspects of our participants and their lives and used this to predict inter-  
139 and intra-individual variation in cognitive performance. These predictors included a) stable  
140 differences between individuals (group, age, sex, rearing history, experience with research),  
141 b) differences that varied within and between individuals (rank, sickness, sociality), c)  
142 differences that varied with group membership (time spent outdoors, disturbances, life  
143 events), and d) differences in testing arrangements (presence of observers, study  
144 participation on the same day and since the last time point).

145 Data collection was split into two phases; after Phase 1 (14 data collection time

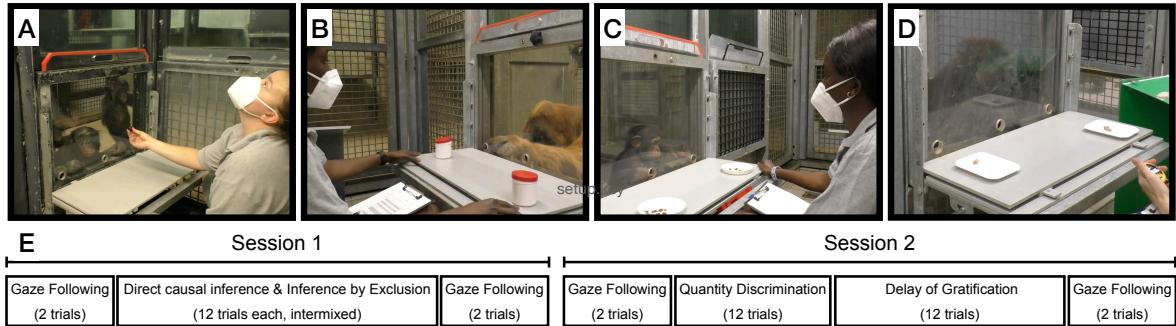
146 points), we analyzed the data and registered the results (<https://osf.io/7qyd8>). Phase 2  
147 lasted for another 14 time points and served to replicate and extend Phase 1. This  
148 approach allowed us to test a) how robust task-level results are, b) how reliable individual  
149 differences are measured and how stable they are over time, c) how individual differences  
150 are structured and d) what predicts cognitive performance.

151 **Results**

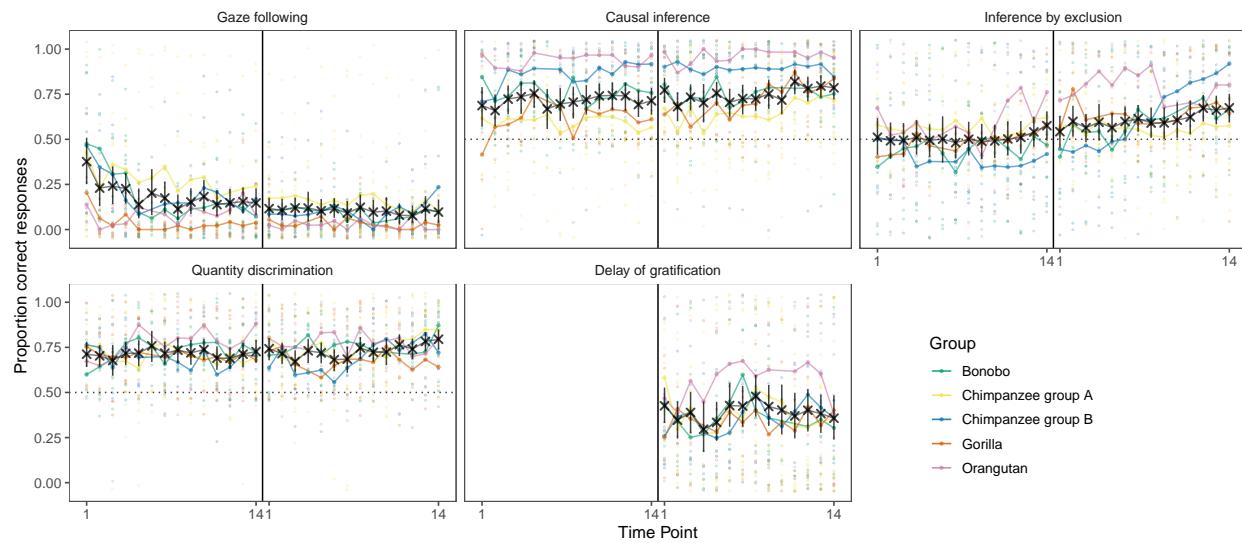
152 **Robustness of task-level performance**

153 As a first step, we asked whether the average performance of a given sample at a time  
154 is robust, that is, whether we could assume to find a similar average performance for a  
155 given sample of individuals if we repeated the task assessment. We assessed robustness in  
156 two ways: First, whenever there was a level of performance expected by chance (i.e. 50%  
157 correct), we checked if the 95% Confidence Interval (CI) for the mean proportion correct  
158 overlapped with chance. Second, we assessed temporal robustness using structural equation  
159 modeling, in particular, Latent State models (see Methods and Supplementary Figure 6 for  
160 details). These models partition the observed performance variable at a given time point  
161 into a latent state variable (time-specific true score variable) and a measurement error  
162 variable. The mean of the latent state variable for the first time point of each phase was  
163 fixed at zero and we assessed average change across time by asking whether the 95%  
164 Credible Intervals (CrI) for the latent state means of subsequent time points overlapped  
165 with zero (i.e. the mean of the first time point).

166 Task-level performance was largely robust or followed clear temporal patterns. Figure  
167 2 visualizes the proportion of correct responses for each task; Figure 3A shows the latent  
168 state means for each task and phase. The direct causal inference and quantity  
169 discrimination tasks were the most robust: in both cases, performance was different from  
170 chance across both phases with no apparent change over time. The rate of gaze following



*Figure 1.* Setup used for the five tasks. A) Gaze following: the experimenter looked to the ceiling. We coded if the ape followed gaze. B) Direct causal inference: food was hidden in one of two cups, the baited cup was shaken (food produced a sound) and apes had to choose the shaken cup to get food. Inference by exclusion: food was hidden in one of two cups. The empty cup was shaken (no sound), so apes had to choose the non-shaken cup to get food. C) Quantity discrimination: Small pieces of food were presented on two plates (5 vs. 7 items); we coded if subjects chose the larger amount. D) Delay of gratification (only Phase 2): to receive a larger reward, the subject had to wait and forgo a smaller, immediately accessible reward. E) Order of task presentation, trial numbers and organisation of tasks into sessions. In both phases, we ran the two sessions on two separate days.



*Figure 2.* Results from the five cognitive tasks across time points. Black crosses show mean performance at each time point across all individuals in the sample at that time point (with 95% CI). The sample size varied between time points and can be found in Supplementary Figure 1. Colored dots show mean performance by species. Light dots show individual means per time point. Dashed lines show chance level whenever applicable. The vertical black line marks the transition between phases 1 and 2.

171 declined at the beginning of Phase 1 but then settled on a low but stable level until the end  
 172 of Phase 2. This pattern was expected given that following the experimenter's gaze was  
 173 never rewarded – neither explicitly with food nor by bringing something interesting to the  
 174 participant's attention. The inference by exclusion task showed an inverse pattern with  
 175 task-level performance being at chance-level for most of Phase 1, followed by a small but  
 176 steady increase throughout Phase 2 so that from time point 6 in Phase 2 onwards,  
 177 performance was significantly different from the first time point of that Phase. These  
 178 temporal patterns most likely reflect training (or habituation) effects that are a  
 179 consequence of repeated testing. Performance in the delay of gratification task (Phase 2  
 180 only) was more variable but within the same general range for the whole testing period. In  
 181 sum, despite these exceptions, performance was very robust in that time points generally

182 supported the same task-level conclusions. For example, Figure 2 shows that performance  
 183 in the direct causal inference task was clearly above chance at all time points and, on a  
 184 descriptive level, consistently higher compared to the inference by exclusion task. Thus,  
 185 the tasks appeared well suited to study task-level performance.

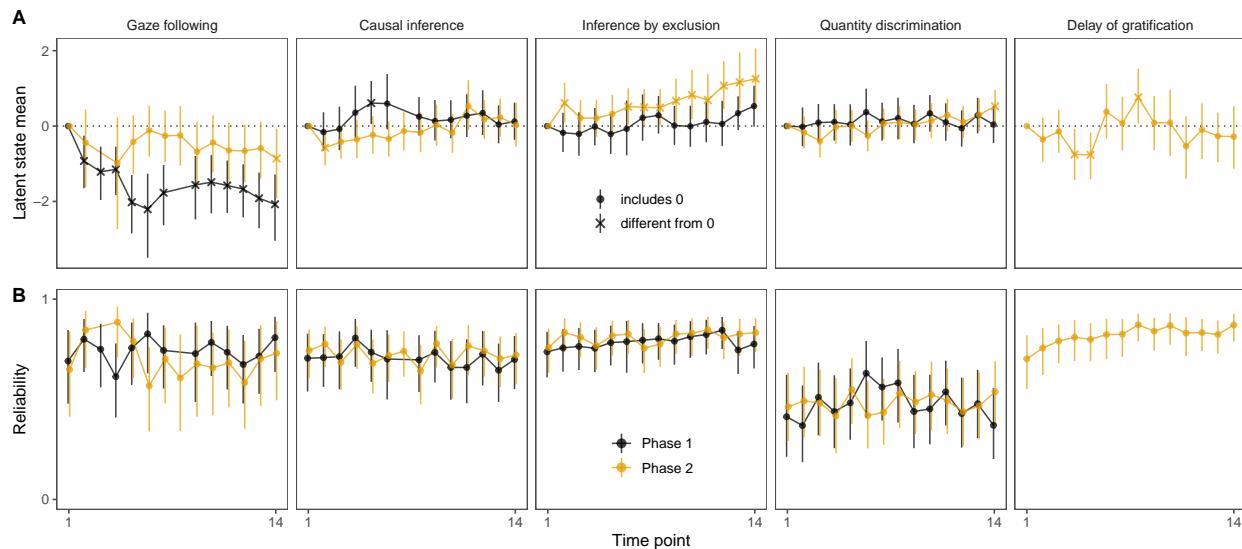


Figure 3. A) Latent state means for each time point by task and phase estimated via latent state models. The sample size varied between time points and can be found in Supplementary Figure 1. Shape denotes whether the 95% CrI included zero. B) Corresponding reliability estimates. Points show mean of the posterior distribution with 95%CrI

#### 186 Reliability of individual-level measurements

187 The reliability of a measure is defined as the proportion of true score variance to its  
 188 observed total variance. That is, a reliable measure captures inter-individual differences  
 189 with precision (i.e., perfect reliability corresponds to measurement without measurement  
 190 error) and is expected to produce similar results if repeated under identical conditions. In  
 191 practice, however, there may be a trade-off between aggregate and individual level  
 192 measurement goals – an observation that has been coined the “reliability paradox”<sup>55</sup>.

193 As a first step towards investigating individual differences, we inspected re-test

194 correlations of our five tasks. For that, we correlated the performance at the different time  
195 points in each task (see Figure 4). Correlations were generally high – some even  
196 exceptionally high – for animal cognition standards<sup>22</sup>. As expected, values were higher for  
197 more proximate time points<sup>56</sup>. The quantity discrimination task had lower correlations  
198 compared to the other tasks.

199 However, based on re-test correlations alone, we cannot say whether lower  
200 correlations reflect higher measurement error (low reliability) or inter-individual differences  
201 in (true) change of performance across time (low stability). To tease these components  
202 apart, we turned again to the LS models mentioned above. For each time point, we  
203 estimated a latent state variable (time-specific true score variable) using two test halves as  
204 indicators. These test halves were constructed by splitting the trials of each task per time  
205 point into two parallel subgroups. Thereby, the models allow us to estimate the reliability  
206 of the respective test halves (see Methods and section *Latent state models* in the  
207 Supplementary Material for details). We interpreted reliability estimates in the following  
208 way: acceptable = .7, good = .8 and high = .9. Please note that these estimates are for  
209 test-halves; the reliability of the full test would be higher.

210 Figure 3B shows that reliability was generally good (~.75) for all tasks at all time  
211 points, except for the quantity discrimination task which had reliability estimates  
212 fluctuating around .5. Thus, the lower re-test correlations for quantity discrimination most  
213 likely reflect low reliability instead of individual changes in cognitive performance across  
214 time. We will return to this point again in the next section. Taken together, these results  
215 suggest that the majority of tasks reliably measured differences between individuals.

216 As a final note, the results show that *task-level robustness does not imply*  
217 *individual-level stability* – and vice versa. The quantity discrimination task showed robust  
218 task-level performance above chance (Figure 2) but relatively poor reliability (Figure 3B).  
219 In other words, even though task-level performance was similar at all time points,

220 differences between individuals were measured with low precision. In contrast, task-level  
 221 performance in the inference by exclusion and gaze following tasks changed over time, with  
 222 satisfactory measurement precision and moderate to high stability of true inter-individual  
 223 differences (see next section).

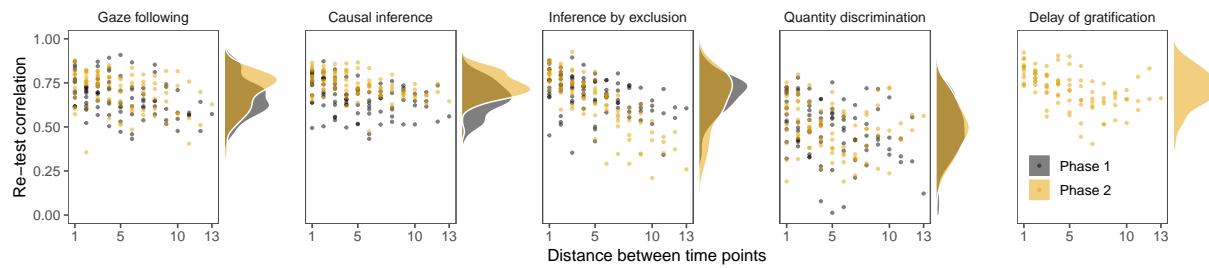


Figure 4. Re-test correlation coefficients are plotted against the temporal distance between the testing time points. Color shows the phase. Side: Distribution of re-test Pearson correlation coefficients.

#### 224 Structure and stability of inter-individual differences

225 Next, we investigated the structure of individual differences. Importantly – and in  
 226 contrast to earlier work<sup>34</sup> – with ‘structure’ we do not exclusively mean the relationship  
 227 between different cognitive tasks. As mentioned in the introduction, we start with a more  
 228 basic question: do individual differences in a given task reflect differences in cognitive  
 229 ability (e.g. ability to make causal inferences) that persist over time or rather differences in  
 230 transient factors (e.g., motivation or attentiveness) that vary from time point to time  
 231 point. The former would imply that individuals (true scores) are ranked similarly across  
 232 time points, while the latter would predict fluctuations in ranks.

233 To quantify to what extent stable or variable differences between individuals explain  
 234 performance, we used latent state-trait (LST) models which partitioned the observed  
 235 performance score into a latent trait variable, a latent state residual variable, and  
 236 measurement error<sup>57–59</sup>. We assume stable latent traits, such that one can think of a latent

237 trait as a stable cognitive ability (e.g., the ability to make causal inferences) and latent  
238 state residuals as variables capturing the effect of occasion-specific psychological conditions  
239 (e.g., being more or less attentive or motivated). The sum of the latent trait and the latent  
240 state residual variable corresponds to the true score of cognitive performance at a specific  
241 time point (latent state variable). We report additional models that account for the  
242 temporal structure of the data in the Supplementary Material (see Supplementary Note).

243 True individual differences were largely stable across time. Across tasks, more than  
244 75% of the reliable variance (true inter-individual differences) was accounted for by latent  
245 *trait* differences and less than 25% by *occasion-specific* variation between individuals  
246 (Figure 5A). The good reliability estimates ( $> .75$  for most tasks; Figure 5A) show that  
247 these latent variables accounted for most of the variance in raw test scores – with the  
248 quantity discrimination task being an exception (reliability = .47). Reflecting back on the  
249 results reported above, we can now say that the – relatively speaking – lower correlations  
250 between time points in the quantity discrimination task indicate a higher degree of  
251 measurement error rather than variable individual differences. In fact, once measurement  
252 error is accounted for, consistency estimates for the quantity discrimination task were close  
253 to 1, reflecting highly stable true differences between individuals.

254 Next, we compared the estimates for the two phases of data collection. We found  
255 estimates for consistency (proportion of true score variance due to latent trait variance)  
256 and occasion specificity (proportion of true score variance due to state residual variance) to  
257 be remarkably similar for the two phases. For inference by exclusion, the LST model did  
258 not fit the data from Phase 2 well (see section *Phase 2 - Inference by exclusion* in the  
259 Supplementary Material for details). Therefore, we divided Phase 2 into two parts (time  
260 points 1-8 and 9-14) and estimated a separate trait for each part. All estimates were  
261 similar for both parts (Figure 5A), and the two traits were highly correlated ( $r = .82$ ).  
262 Together with the LS model results reported in the robustness section, this suggests that  
263 the increase in group-level performance in Phase 2 was probably driven by a relatively

<sup>264</sup> sudden improvement of a few individuals, mostly from the chimpanzee B group (see Figure  
<sup>265</sup> 2). These individuals quickly improved in performance halfway through Phase 2 and  
<sup>266</sup> retained this level for the rest of the study.

<sup>267</sup> Finally, we investigated the relationship between latent traits. We asked whether  
<sup>268</sup> individuals with high abilities in one domain also have higher abilities in another. We fit  
<sup>269</sup> pairwise LST models that modeled the correlation between latent traits for two tasks (two  
<sup>270</sup> models for inference by exclusion in Phase 2). In Phase 1, the only substantial correlation  
<sup>271</sup> (i.e. coefficients indicated medium to large effects<sup>60</sup> and their 95% CrI did not include zero)  
<sup>272</sup> was between quantity discrimination and inference by exclusion. In Phase 2, this finding  
<sup>273</sup> was replicated, and, in addition, four more correlations turned out to be substantial (see  
<sup>274</sup> Figure 5B). One reason for this increase was the inclusion of the delay of gratification task.  
<sup>275</sup> Across phases, correlations involving the gaze following task were the closest to zero, with  
<sup>276</sup> quantity discrimination in Phase 2 being an exception. Taken together, the overall pattern  
<sup>277</sup> of results suggests substantial shared variance between tasks – except for gaze following.

## <sup>278</sup> Predictability of individual differences

<sup>279</sup> The results thus far suggest that individual differences originate from stable  
<sup>280</sup> differences between individuals in cognitive abilities that persist across time points.  
<sup>281</sup> Differences in ability outweigh fluctuations due to transient, occasion-specific factors such  
<sup>282</sup> as attentiveness or motivation. An alternative pattern would arise when time point-specific  
<sup>283</sup> variation in e.g., attentiveness or motivation would be responsible for differences in  
<sup>284</sup> performance between individuals. Of course, there can be stable differences between  
<sup>285</sup> individuals in attentiveness and motivation, in which case they would influence  
<sup>286</sup> performance in a consistent way over time and presumably also across tasks<sup>61–63</sup>. The  
<sup>287</sup> distinction we want to make here is between transient and stable factors influencing  
<sup>288</sup> cognitive performance.

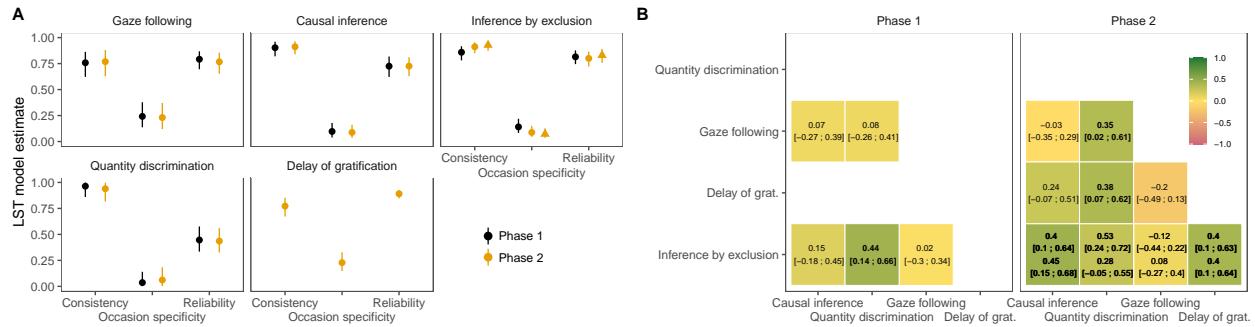


Figure 5. A) Mean estimates from latent state-trait models for Phase 1 and 2 with 95% CrI based on data from N = 43 participants. Consistency: proportion of (measurement-error-free) variance in performance explained by stable trait differences. Occasion specificity: proportion of true variance explained by variable state residuals. Reliability: proportion of true score variance to variance in raw scores. For inference by exclusion: different shapes show estimates for different parts of Phase 2 (see main text for details). B) Correlations between latent traits based on pairwise LST models between tasks with 95% CrI. Bold correlations have CrI not overlapping with zero. Inference by exclusion has one value per part in Phase 2. The models for quantity discrimination and direct causal inference showed a poor fit and are not reported here (see section \*Relations between tasks – Phase 2\* in Supplementary Material for details).

289 In the last set of analyses, we sought to explain the origins of individual differences.  
 290 That is, we analyzed whether inter- and intra-individual variation in cognitive performance  
 291 in the tasks could be predicted by non-cognitive variables that captured a) stable  
 292 differences between individuals (group, age, sex, rearing history, experience with research),  
 293 b) differences that varied within and between individuals (rank, sickness, sociality), c)  
 294 differences that varied with group membership (time spent outdoors, disturbances, life  
 295 events), and d) differences in testing arrangements (presence of observers, study  
 296 participation on the same day and since the last time point). We collected these predictor  
 297 variables using a combination of directed observations and caretaker questionnaires.

298 This large set of potentially relevant predictors poses a variable selection problem.

299 Thus, in our analysis, we sought to find the smallest number of predictors (main effects  
300 only) that allowed us to accurately predict performance in the cognitive tasks. We chose  
301 the projection predictive inference approach because it provides an excellent trade-off  
302 between model complexity and accuracy<sup>64-66</sup>. The outcome of this analysis is a ranking of  
303 the different predictors in terms of how important they are to predicting performance in a  
304 given task. Furthermore, for each predictor, we get a qualitative assessment of whether it  
305 makes a substantial contribution to predicting performance in the task or not.

306 Predictors capturing stable individual characteristics were ranked highest and selected

307 as relevant most often (Figure 6A). The three highest-ranked predictors belonged to this  
308 category. This result fits well with the LST model results reported above, in which we saw  
309 that most of the variance in performance could be traced back to stable trait differences  
310 between individuals. Here we saw that performance was best predicted by variables that  
311 reflect stable characteristics of individuals. This suggests that stable characteristics  
312 partially cause selective development that leads to differences in cognitive abilities. The  
313 tasks with the highest occasion-specific variance (gaze following and delay of gratification,  
314 see Figure 5A) were also those for which the most time point-specific predictors were  
315 selected. The quantity discrimination task did not fit this pattern in Phase 2; even though  
316 the LST model suggested that only a very small portion of the variance in performance was  
317 occasion-specific, four time-point-specific variables were selected to be relevant.

318 The most important predictor was group. Interestingly, differences between groups

319 were not systematic in that one group consistently outperformed the others across tasks.

320 Furthermore, group differences could not be collapsed into species differences as the two  
321 chimpanzee groups varied largely independently of one another (Figure 6B). Predictors  
322 that were selected more than once influenced performance in variable ways. The presence  
323 of observers always had a negative effect on performance. The more time an individual had  
324 been involved in research during their lifetime, the better performance was. On the other

325 hand, while the rate of gaze following increased with age in Phase 1, performance in the  
326 inference by exclusion task decreased. Females were more likely to follow gaze than males,  
327 but males were more likely to wait for the larger reward in the delay of gratification task.  
328 Finally, time spent outdoors had a positive effect on gaze following but a negative effect on  
329 direct causal inference (Figure 6B).

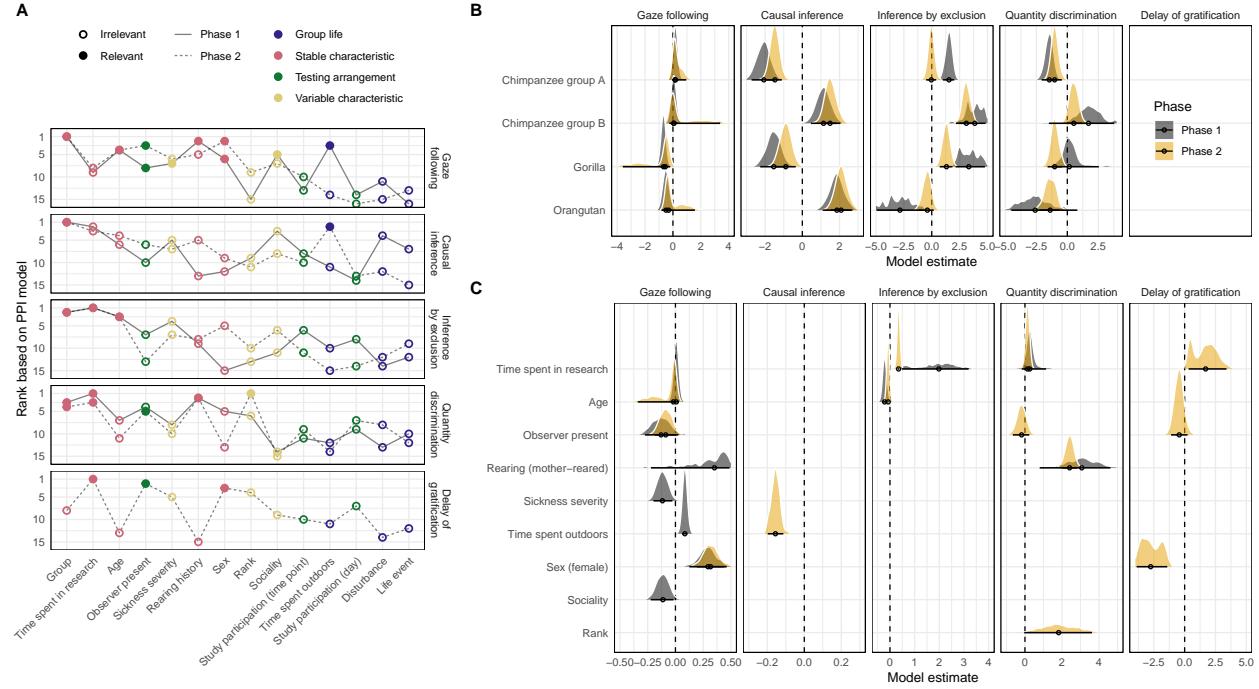
330 In sum, individual characteristics were most predictive of cognitive performance. In  
331 most cases, the corresponding predictors were selected as relevant in both phases. The  
332 influence of time-point-specific predictors was less consistent: except for the presence of an  
333 observer in the gaze following task, none of the variable predictors was selected as relevant  
334 in both phases. To avoid misinterpretation, this suggests that cognitive performance was  
335 influenced by temporal variation in group life, testing arrangements, and variable  
336 characteristics; however, the way this influence exerts itself was either less consistent or less  
337 pronounced (or both) compared to the influence of stable characteristics.

338 It is important to note, however, that in terms of absolute variance explained, the  
339 largest portion was accounted for by a random intercept term in the model (not shown in  
340 Figure 5) that simply captured the identity of the individual (see Supplementary Figures  
341 21 to 29). This suggests that idiosyncratic developmental processes and/or genetic  
342 pre-dispositions, which operate on a much longer time scale than what we captured in the  
343 present study, were responsible for most of the variation in cognitive performance.

344

## Discussion

345 This study aimed to test the assumptions of robustness, stability, reliability, and  
346 predictability that underlie much of comparative research and theorizing about cognitive  
347 evolution. We repeatedly tested a large sample of great apes in five tasks covering a range  
348 of different cognitive domains. We found task-level performance to be robust for most tasks  
349 so that conclusions drawn based on one testing occasion mirrored those on other occasions.



*Figure 6.* A. Ranking of predictors based on the projection predictive inference model for the five tasks in the two phases. Order (left to right) is based on average rank across phases. Solid points indicate predictors selected as relevant. Color of the points shows the category of the predictor. Line type denotes the phase. B. Posterior model estimates for the selected predictors for each task based on data from  $N = 43$  participants. Points show means with 95% Credible Interval. Color denotes phase. For categorical predictors, the estimate gives the difference compared to the reference level (Bonobo for group, no observer for observer, hand-reared for rearing, male for sex).

350 Most of the tasks measured differences between individuals in a reliable and stable way –  
351 making them suitable to study individual differences. Using structural equation models, we  
352 found that individual differences in performance were largely explained by traits – that is,  
353 stable differences in cognitive abilities between individuals. Furthermore, we found  
354 systematic relationships between cognitive abilities. When predicting variation in cognitive  
355 performance, we found stable individual characteristics (e.g., group or time spent in  
356 research) to be the most important. Variable predictors were also found to be influential at  
357 times but less systematically.

358 At first glance, the results send a reassuring message: most of the tasks we used  
359 produced robust task-level results and captured individual differences in a reliable and  
360 stable way. However, this did not apply to all tasks. As noted above, in the Supplementary  
361 Material, we report on a rule-switching task<sup>54</sup> that produced neither stable nor reliable  
362 results (Supplementary Figures 8 and 9). The quantity discrimination task was robust on a  
363 task level but did not measure individual differences reliably. We draw two conclusions  
364 based on this pattern. First, replicating studies – even if it is with the same animals –  
365 should be an integral part of primate cognition research<sup>67–69</sup>. Second, for individual  
366 differences research, it is crucial to assess the psychometric properties (e.g., reliability) of  
367 the measures involved<sup>70</sup>. If this step is omitted, it is difficult to interpret studies, especially  
368 when they produce null results. It is important to note that the sample size in the current  
369 study was large compared to other comparative studies (median sample size across studies  
370 = 7)<sup>67</sup>. With smaller sample sizes, task-level estimates are likely more variable and thus  
371 more likely to produce false-positive or false-negative conclusions<sup>71,72</sup>. Small samples in  
372 comparative research usually reflect the resource limitations of individual labs. Pooling  
373 resources in large-scale collaborative projects like *ManyPrimates*<sup>73,74</sup> will thus be vital to  
374 corroborate findings.

375 Continuing on this theme, the data reported here would be exciting to explore for  
376 species differences. For example, the descriptive results shown in Figure 2 suggest that

377 orangutans performed best in the nonsocial tasks but worse in the social task. However, we  
378 are hesitant to interpret such findings because of the small sample sizes per species and the  
379 substantial differences in sample size between species. Consequently, it is impossible to  
380 distinguish individual-level from species-level variation.

381 Given their good psychometric properties, our tasks offer insights into the structure  
382 of great ape cognition. We used structural equation modeling to partition reliable variance  
383 in performance into stable (trait) and variable (state residual) differences between  
384 individuals. We found traits to explain more than 75% of the reliable variance across tasks.  
385 This suggests that the patterns in performance we observed mainly originate from stable  
386 differences in cognitive abilities . This finding does not mean there cannot be  
387 developmental change over longer time periods. In fact, for the inference by exclusion task,  
388 we saw a relatively abrupt change in performance for some individuals, which stabilized on  
389 an elevated level, suggesting a sustained change in cognitive ability.

390 We found systematic relationships between traits estimated via LST models for the  
391 different tasks. Correlations tended to be higher among the non-social tasks compared to  
392 when the gaze-following task was involved, which could be taken to indicate shared  
393 cognitive processes. However, we feel such a conclusion would be premature and require  
394 additional evidence from more tasks and larger sample sizes<sup>34</sup>. One possibility is that  
395 stable, domain-general psychological processes – such as attentiveness or motivation – are  
396 responsible for the shared variance. Cognitive modeling could be used to explicate the  
397 processes involved in each task. Shared processes could be probed by comparing models  
398 that make different assumptions<sup>75,76</sup>.

399 The finding that stable differences in cognitive abilities explained most of the  
400 variation between individuals was also corroborated by the analyses focused on the  
401 predictability of performance. We found that predictors that captured stable individual  
402 characteristics (e.g., group, time spent in research, age, rearing history) were more likely to

403 be selected as relevant predictors. Aspects of everyday experience or testing arrangements  
404 that would influence performance on particular time points and thus increase the  
405 proportion of occasion-specific variation (e.g., life events, disturbances, participating in  
406 other tests) were ranked as less important. Despite this general pattern, there was  
407 variation across tasks in which individual characteristics were selected to be relevant. For  
408 example, rearing history was an important predictor for quantity discrimination and gaze  
409 following but less so for the other three tasks (Figure 6A). Group – the overall most  
410 important predictor – exerted its influence differently across tasks. Orangutans, for  
411 example, outperformed the other groups in direct causal inference but were the least likely  
412 to follow gaze. Together with the finding that the random intercept term explained the  
413 largest proportion of variance in performance across tasks, this pattern suggests that the  
414 cognitive abilities underlying performance in the different tasks respond to different –  
415 though sometimes overlapping – external conditions that together shape the individual’s  
416 developmental environment.

417 Our results also address a very general matter. Comparative psychologists often  
418 worry – or are told they should worry – that their results can be explained by  
419 mechanistically simpler associative learning processes<sup>77</sup>. Oftentimes such explanations are  
420 theoretically plausible and rarely disproved empirically<sup>78</sup>. The present study speaks to this  
421 issue in so far as we created the conditions for such associative learning processes to  
422 potentially unfold. Great apes were tested by the same experimenter in the same tasks,  
423 using differential reinforcement and the same counterbalancing for hundreds of trials.  
424 However, a steady increase in performance – uniform over individuals – did not show. This  
425 does not take away the theoretical possibility that associative learning accounts for  
426 improved performance over time on isolated tasks. In fact, we are agnostic as to whether or  
427 not a particular learning account might explain our results (or parts of them) and invite  
428 others to further analyze the data provided here.

**429 Conclusion**

430 The present study put the implicit assumptions underlying much of comparative  
431 research on cognitive evolution involving great apes to an empirical test. While we found  
432 reassuring results in terms of task-level stability and reliability of the measurement of  
433 individual differences, we also pointed out the importance of explicitly questioning and  
434 testing these assumptions, ideally in large-scale collaborative projects. Our results paint a  
435 picture of great ape cognition in which variation between individuals is predicted and  
436 explained by stable individual characteristics that respond to different – though sometimes  
437 overlapping – developmental conditions. Hence, an ontogenetic perspective is not auxiliary  
438 but fundamental to studying cognitive diversity across species. We hope these results  
439 contribute to a more solid and comprehensive understanding of the nature and origins of  
440 great ape and human cognition as well as provide useful methodological guidance for future  
441 comparative research.

**442 Methods****443 Participants**

444 A total of 43 great apes participated at least once in one of the tasks. This included 8  
445 bonobos (*pan paniscus*, 3 females, age 7.30 to 39), 24 chimpanzees (*pan troglodytes*, 18  
446 females, age 2.60 to 55.90), 6 gorillas (*gorilla gorilla*, 4 females, age 2.70 to 22.60), and 5  
447 orangutans (*pongo abelii*, 4 females, age 17 to 41.20). The overall sample size at the  
448 different time points ranged from 22 to 43 for the different species.

449 Apes were housed at the Wolfgang Köhler Primate Research Center located in Zoo  
450 Leipzig, Germany. They lived in groups, with one group per species and two chimpanzee  
451 groups (A and B). Studies were noninvasive and adhered to the legal requirements in  
452 Germany. Animal husbandry and research complied with the European Association of Zoos  
453 and Aquaria Minimum Standards for the Accommodation and Care of Animals in Zoos

454 and Aquaria as well as the World Association of Zoos and Aquariums Ethical Guidelines  
455 for the Conduct of Research on Animals by Zoos and Aquariums. Participation was  
456 voluntary, all food was given in addition to the daily diet, and water was available ad  
457 libitum throughout the study. The study was approved by an internal ethics committee at  
458 the Max Planck Institute for Evolutionary Anthropology.

## 459 Procedure

460 Apes were tested in familiar sleeping or test rooms by a single experimenter.  
461 Whenever possible, they were tested individually. The basic setup comprised a sliding table  
462 positioned in front of a clear Plexiglas panel with three holes in it. The experimenter sat  
463 on a small stool and used an occluder to cover the sliding table (see Figure 1).

464 The tasks we selected are based on published procedures and are commonly used in  
465 the field of comparative psychology. Example videos for each task can be found in the  
466 associated online repository (<https://github.com/ccp-eva/laac/tree/master/videos>).

467 **Gaze Following.** The gaze following task was modeled after a study by Bräuer and  
468 colleagues<sup>50</sup>. The experimenter sat opposite the ape and handed over food at a constant  
469 pace. That is, the experimenter picked up a piece of food, briefly held it out in front of her  
470 face and then handed it over to the participant. After a predetermined (but varying)  
471 number of food items had been handed over, the experimenter again picked up a food item,  
472 held it in front of her face and then looked up (i.e., moving her head up – see Figure 1A).  
473 The experimenter looked to the ceiling; no object of particular interest was placed there.  
474 After 10s, the experimenter looked down again, handed over the food and the trial ended.  
475 We coded whether the participant looked up during the 10s interval. Apes received eight  
476 gaze-following trials. We assume that participants look up because they assume that the  
477 experimenter's attention is focused on a potentially noteworthy object.

478       **Direct causal inference.** The direct causal inference task was modeled after a  
479 study by Call<sup>51</sup>. Two identical cups, each with a lid, were placed left and right on the table  
480 (Figure 1B). The experimenter covered the table with the occluder, retrieved a piece of  
481 food, showed it to the ape, and hid it in one of the cups outside the participant's view.  
482 Next, the experimenter removed the occluder, picked up the baited cup and shook it three  
483 times, which produced a rattling sound. Next, the cup was put back in place, the sliding  
484 table pushed forwards, and the participant made a choice by pointing to one of the cups. If  
485 they picked the baited cup, their choice was coded as correct, and they received the  
486 reward. If they chose the empty cup, they did not. Participants received 12 trials. The  
487 location of the food was counterbalanced; six times in the right cup and six times in the  
488 left. Direct causal inference trials were intermixed with inference by exclusion trials (see  
489 below). We assume that apes locate the food by reasoning that the food – a solid object –  
490 causes the rattling sound and, therefore, must be in the shaken cup.

491       **Inference by exclusion.** Inference by exclusion trials were also modeled after the  
492 study by Call<sup>51</sup> and followed a very similar procedure compared to direct causal inference  
493 trials. After covering the two cups with the occluder, the experimenter placed the food in  
494 one of the cups and covered both with the lid. Next, they removed the occluder, picked up  
495 the empty cup and shook it three times. In contrast to the direct causal inference trials,  
496 this did not produce any sound. The experimenter then pushed the sliding table forward  
497 and the participant made a choice by pointing to one of the cups. Correct choice was coded  
498 when the baited (non-shaken) cup was chosen. If correct, the food was given to the ape.  
499 There were 12 inference by exclusion trials intermixed with direct causal inference trials.  
500 The order was counterbalanced: six times the left cup was baited, six times the right. We  
501 assume that apes reason that the absence of a sound suggests that the shaken cup is  
502 empty. Because they saw a piece of food being hidden, they exclude the empty cup and  
503 infer that the food is more likely to be in the non-shaken cup.

504       **Quantity discrimination.** For this task, we followed the general procedure of  
505 Hanus and colleagues<sup>52</sup>. Two small plates were presented left and right on the table (see  
506 Figure 1C). The experimenter covered the plates with the occluder and placed five small  
507 food pieces on one plate and seven on the other. Then they pushed the sliding table  
508 forwards, and the participant made a choice. We coded as correct when the subject chose  
509 the plate with the larger quantity. Participants always received the food from the plate  
510 they chose. There were 12 trials, six with the larger quantity on the right and six on the  
511 left (order counterbalanced). We assume that apes identify the larger of the two food  
512 amounts based on discrete quantity estimation.

513       **Delay of gratification.** This task replaced the switching task in Phase 2. The  
514 procedure was adapted from Rosati and colleagues<sup>53</sup>. Two small plates, including one and  
515 two pieces of pellet, were presented left and right on the table. The experimenter moved  
516 the plate with the smaller reward forward, allowing the subject to choose immediately,  
517 while the plate with the larger reward was moved forward after a delay of 20 seconds. We  
518 coded whether the subject selected the larger delayed reward (correct choice) or the smaller  
519 immediate reward (incorrect choice) as well as the waiting time in cases where the  
520 immediate reward was chosen. Subjects received 12 trials, with the side on which the  
521 immediate reward was presented counterbalanced. We assume that, in order to choose the  
522 larger reward, apes inhibit choosing the immediate smaller reward.

523       **Interrater reliability.** A second coder unfamiliar to the purpose of the study  
524 coded 15% of all time points (four out of 28) for all tasks. Reliability was good to excellent.  
525 Gaze following: 92% agreement ( $\kappa = .64$ ), direct causal inference 99% agreement ( $\kappa = .98$ ),  
526 inference by exclusion: 99% agreement ( $\kappa = .99$ ), quantity discrimination: 99% agreement  
527 ( $\kappa = .97$ ), delay of gratification: 98% agreement ( $\kappa = .97$ ).

528 **Data collection**

529 We collected data in two phases. Phase 1 started on August 1st, 2020, lasted until  
530 March 5th, 2021, and included 14 time points. Phase 2 started on May 26th, 2021, and  
531 lasted until December 4th, 2021, and also had 14 time points. Phase 1 also included a  
532 strategy switching task. However, because it did not produce meaningful results, we  
533 replaced it with the delay of gratification task. Details and results can be found in the  
534 Supplementary Material available (Supplementary Figure 8 and 9).

535 One time point meant running all tasks with all participants. Within each time  
536 point, the tasks were organized in two sessions (see Figure 1E). Session 1 started with two  
537 gaze following trials. Next was a pseudo-randomized mix of direct causal inference and  
538 inference by exclusion trials with 12 trials per task but no more than two trials of the same  
539 task in a row. At the end of Session 1, there were again two gaze following trials. Session 2  
540 also started with two gaze following trials, followed by quantity discrimination and strategy  
541 switching (Phase 1) or delay of gratification (Phase 2). Finally, there were again two gaze  
542 following trials. The order of tasks was the same for all subjects. So was the positioning of  
543 food items within each task. The two sessions were usually spread out across two adjacent  
544 days. The interval between two time points was planned to be two weeks. However, it was  
545 not always possible to follow this schedule, so some intervals were longer or shorter.  
546 Supplementary Figure 1 shows the timing and spacing of the time points.

547 In addition to the data from the cognitive tasks, we collected data for a range of  
548 predictor variables. Predictors could either vary with the individual (stable individual  
549 characteristics: group, age, sex, rearing history, time spent in research), vary with  
550 individual and time point (variable individual characteristics: rank, sickness, sociality),  
551 vary with group membership (group life: e.g., time spent outdoors, disturbances, life  
552 events) or vary with the testing arrangements and thus with individual, time point and  
553 session (testing arrangements: presence of observers, study participation on the same day

554 and since the last time point). Most predictors were collected via a diary that the animal  
555 caretakers filled out on a daily basis. Here, the caretakers were asked a range of questions  
556 about the presence of a predictor and its severity. Other predictors were based on direct  
557 observations. A detailed description of the predictors and how they were collected can be  
558 found in the section *Predictors* in the Supplementary Material.

559 **Analysis**

560 In the following, we provide an overview of the analytical procedures we used. We  
561 encourage the reader to consult the section *Analytical Framework* in the Supplementary  
562 Material available online for additional details. We had two overarching questions. On the  
563 one hand, we were interested in the cognitive measures and the relationships between  
564 them. That is, we asked how robust performance was on a task-level, how stable individual  
565 differences were, and how reliable the measures were. We also investigated relationships  
566 between the different tasks. We used structural equation modeling (SEM)<sup>79,80</sup> to address  
567 these questions.

568 Our second question was, which predictors explained variability in cognitive  
569 performance. Here we wanted to see which of the predictors we recorded were most  
570 important to predict performance over time. This is a variable selection problem (selecting  
571 a subset of variables from a larger pool) and we used projection predictive inference\* for  
572 this<sup>66</sup>.

573 **Structural equation modeling.** We used SEM<sup>79,80</sup> to address the reliability and  
574 stability of each task, as well as relationships between tasks. SEMs allowed us to partition  
575 the variance in performance into latent variable (true-score) variance and measurement  
576 error variance. Latent variables are estimated using multiple observed indicators (here: two  
577 test halves, see below). Longitudinal data for each task was modeled with a latent state  
578 (LS) and a latent state-trait (LST) model<sup>57–59</sup>. All of the models were estimated as  
579 normal-ogive graded response models<sup>81,82</sup> due to the ordinal nature of the indicators. For

580 each task and time point we split the trials in two test halves, which served as indicators  
 581 for a common latent construct. Due to only few different observed values and skewed  
 582 distributions of the sum score for each test half, indicators were modeled as ordered  
 583 categorical variables, using a probit link function. That is, the models assume a continuous  
 584 latent ability underlying the discrete responses, with an increasing probability of more  
 585 correctly solved trials with increasing ability.

586 Formally speaking, the observed categorical variables  $Y_{it}$  for test half  $i$  at time point  $t$   
 587 result from a categorization of unobserved continuous latent variables  $Y_{it}^*$  which underlie  
 588 the observed categorical variables. In the LS models,  $Y_{it}^*$  is decomposed into a latent  
 589 state variable  $S_t$  and a measurement error variable  $\epsilon_{it}$ <sup>83</sup>. At each time point  $t$ , the two  
 590 latent variables  $Y_{1t}^*$  and  $Y_{2t}^*$  are assumed to capture a common latent state variable  $S_t$ . To  
 591 test for possible mean changes of ability across time, the means of the latent state variables  
 592 were freely estimated (assuming invariance of the threshold parameters  $\kappa_{sit}$  across time).

593 As an estimate of reliability, we computed the proportion of true score variance  
 594 relative to the total variance of the continuous latent variables  $Y_{it}^*$ :

$$Rel(Y_{it}^*) = \frac{Var(S_t)}{Var(S_t) + Var(\epsilon_{it})} = \frac{Var(S_t)}{Var(S_t) + 1} \quad (1)$$

595 For the LST model, the continuous latent variable  $Y_{it}^*$  was decomposed into a latent  
 596 trait variable  $T_{it}$ , a latent state residual variable  $\zeta_{it}$ , and a measurement error variable. The  
 597 latent trait variables  $T_{it}$  are time-specific dispositions, that is, they capture the expected  
 598 value of the latent state (i.e., true score) variable for an individual at time  $t$  across all  
 599 possible situations the individual might experience at time  $t$ <sup>83,84</sup>. The state residual  
 600 variables  $\zeta_{it}$  capture the deviation of a momentary state from the time-specific disposition  
 601  $T_{it}$ . We assumed that latent traits were stable across time. In addition, we assumed  
 602 common latent trait and state residual variables across the two test halves, which leads to  
 603 the following measurement equation for parcel  $i$  at time point  $t$ :

$$Y_{it}^* = T + \zeta_t + \epsilon_{it} \quad (2)$$

604 Here,  $T$  is a stable (time-invariant) latent trait variable, capturing stable  
 605 inter-individual differences. The state residual variable  $\zeta_t$  captures time-specific deviations  
 606 of the respective true score from the trait variable at time  $t$ , and thereby captures  
 607 deviations from the trait due to situation or person-situation interaction effects.  $\epsilon_{it}$  denotes  
 608 a measurement error variable, with  $\epsilon_{it} \sim N(0, 1) \forall i, t$ . This allowed us to compute the  
 609 following variance components.

610 Consistency: Proportion of true variance (i.e., measurement-error-free variance) that  
 611 is due to true inter-individual stable trait differences.

$$Con(Y_{it}^*) = \frac{Var(T)}{Var(T) + Var(\zeta_t)} \quad (3)$$

612 Occasion specificity: Proportion of true variance (i.e., measurement-error-free  
 613 variance) that is due to true inter-individual differences in the state residual variables (i.e.,  
 614 occasion-specific variation not explained by the trait).

$$OS(Y_{it}^*) = 1 - Con(Y_{it}^*) = \frac{Var(\zeta_t)}{Var(T) + Var(\zeta_t)} \quad (4)$$

615 As state residual variances  $Var(\zeta_t)$  were set equal across time,  $OS(Y_{it}^*)$  is constant  
 616 across time (as well as across item parcels  $i$ ).

617 To investigate associations between cognitive performance in different tasks, the LST  
 618 models were extended to multi-trait models. Due to the small sample size, we could not  
 619 combine all tasks in a single, structured model. Instead, we assessed relationships between  
 620 tasks in pairs.

621 We used Bayesian estimation techniques to estimate the models. In the section  
 622 *Estimation* in the Supplementary Material, we report the prior settings used for estimation

623 as well as the restrictions we imposed on the model parameters. We justify these settings  
624 via simulation studies described in the Supplementary Note.

625 **Projection predictive inference.** The selection of relevant predictor variables  
626 constitutes a variable selection problem, for which a range of different methods are available  
627 e.g., shrinkage priors<sup>85</sup>. We chose to use *Projection Predictive Inference* (PPI) because it  
628 provides an excellent trade-off between model complexity and accuracy<sup>64,66</sup>, especially when  
629 the goal is to identify a minimal subset of predictors that yield a good predictive model<sup>65</sup>.

630 The PPI approach can be viewed as a two-step process: The first step consists of  
631 building the best predictive model possible, called the reference model. In the context of  
632 this work, the reference model is a Bayesian multilevel regression model with repeated  
633 measurements nested in apes, fit using the package `brms`<sup>86</sup>, including all 14 predictors and  
634 a random intercept term for the individual (R notation: `DV ~ predictors + (1 |`  
635 `subject)`). Note that this reference model only included main effects and no interactions  
636 between predictors. Including interactions would have increased the number of predictors  
637 to consider exponentially.

638 In the second step, the goal is to replace the posterior distribution of the reference  
639 model with a simpler distribution. This is achieved via a forward step-wise addition of  
640 predictors that decrease the Kullback-Leibler (KL) divergence from the reference model to  
641 the projected model.

642 The result of the projection is a list containing the best model for each number of  
643 predictors from which the final model is selected by inspecting the mean log-predictive  
644 density (`elpd`) and root-mean-squared error (`rmse`). The projected model with the smallest  
645 number of predictors is chosen, which shows similar predictive performance as the reference  
646 model.

647 We built separate reference models for each phase and task and ran them through the  
648 above-described projection predictive inference approach. The dependent variable for each

task was the cognitive performance of the apes, that is, the number of correctly solved trials per time point and task. The model for the delay of gratification task was only estimated once (Phase 2).

We used the R package `projpred`<sup>87</sup>, which implements the aforementioned projection predictive inference technique. The predictor relevance ranking is measured by the Leave-One-Out (LOO) cross-validated mean log-predictive density and root-mean-squared error. To find the optimal submodel size, we inspected summaries and the plotted trajectories of the calculated `elpd` and `rmse`.

The order of relevance for the predictors and the random intercept (together called terms) is created by performing forward search. The term that decreases the KL divergence between the reference model's predictions and the projection's predictions the most goes into the ranking first. Forward search is then repeated  $N$  times to get a more robust selection. We chose the final model by inspecting the predictive utility of each projection. To be precise, we chose the model with  $p$  terms where  $p$  depicts the number of terms at the cutoff between the term that increases the `elpd` and the term that does not increase the `elpd` by any significant amount. In order to get a useful predictor ranking, we manually delayed the random intercept (and random slope for time point for gaze following) term to the last position in the predictor selection process. The random intercept delay is needed because if the random intercept were not delayed, it would soak up almost all of the variance of the dependent variable before the predictors are allowed to explain some amount of the variance themselves.

## 670 Data Availability

All data can be found in the following public repository:  
<https://github.com/ccp-eva/laac>. The same repository also contains example videos for the different tasks.

**674 Code Availability**

675 All analysis code needed to reproduce the results and figures reported in the paper  
676 and the Supplementary Material can be found in the following public repository:  
677 <https://github.com/ccp-eva/laac>.

**678 Acknowledgements**

679 We thank Damilola Olaoba, Anna Wolff, Nico Eisenbrenner for the data collection.  
680 We are very grateful to Matthias Allritz for his helpful comments on an earlier version of  
681 the paper. Furthermore, we thank all keepers at the WKPRC for their help conducting  
682 this study. The authors received no specific funding for this work.

**683 Competing interest**

684 The authors declare that no competing interests exist.

685

## References

- 686 1. Coqueugniot, H., Hublin, J.-J., Veillon, F., Houët, F. & Jacob, T. Early brain growth  
687 in homo erectus and implications for cognitive ability. *Nature* **431**, 299–302 (2004).
- 688 2. Gunz, P. *et al.* Australopithecus afarensis endocasts suggest ape-like brain organiza-  
689 tion and prolonged brain growth. *Science Advances* **6**, eaaz4729 (2020).
- 690 3. Coolidge, F. L. & Wynn, T. An introduction to cognitive archaeology. *Current Di-  
691 rections in Psychological Science* **25**, 386–392 (2016).
- 692 4. Currie, A. & Killin, A. From things to thinking: Cognitive archaeology. *Mind &  
693 Language* **34**, 263–279 (2019).
- 694 5. Haslam, M. *et al.* Primate archaeology evolves. *Nature Ecology & Evolution* **1**, 1431–  
695 1437 (2017).
- 696 6. Martins, E. P. & Martins, E. P. *Phylogenies and the comparative method in animal  
697 behavior*. (Oxford University Press, 1996).
- 698 7. MacLean, E. L. *et al.* How does cognition evolve? Phylogenetic comparative psychol-  
699 ogy. *Animal Cognition* **15**, 223–238 (2012).
- 700 8. Burkart, J. M., Schubiger, M. N. & Schaik, C. P. van. The evolution of general  
701 intelligence. *Behavioral and Brain Sciences* **40**, (2017).
- 702 9. Shettleworth, S. J. *Cognition, evolution, and behavior*. (Oxford university press,  
703 2009).
- 704 10. Laland, K. & Seed, A. Understanding human cognitive uniqueness. *Annual Review  
705 of Psychology* **72**, 689–716 (2021).
- 706 11. Heyes, C. *Cognitive gadgets*. (Harvard University Press, 2018).
- 707
- 708 12. Tomasello, M. *Becoming human*. (Harvard University Press, 2019).
- 709

- 710 13. Penn, D. C., Holyoak, K. J. & Povinelli, D. J. Darwin's mistake: Explaining the  
discontinuity between human and non-human minds. *Behavioral and Brain Sciences*  
**31**, 109–130 (2008).
- 711
- 712 14. Dunbar, R. & Shultz, S. Why are there so many explanations for primate brain  
evolution? *Philosophical Transactions of the Royal Society B: Biological Sciences*  
**372**, 20160244 (2017).
- 713
- 714 15. Dean, L. G., Kendal, R. L., Schapiro, S. J., Thierry, B. & Laland, K. N. Identification  
of the social and cognitive processes underlying human cumulative culture. *Science*  
**335**, 1114–1118 (2012).
- 715
- 716 16. Call, J. E., Burghardt, G. M., Pepperberg, I. M., Snowdon, C. T. & Zentall, T. E.  
*APA handbook of comparative psychology: Basic concepts, methods, neural substrate,*  
*and behavior, vol. 1.* (American Psychological Association, 2017).
- 717
- 718 17. Darwin, C. *On the origin of species.* (Routledge, 1859).
- 719
- 720 18. Thornton, A. & Lukas, D. Individual variation in cognitive performance: Develop-  
mental and evolutionary perspectives. *Philosophical Transactions of the Royal Society*  
*B: Biological Sciences* **367**, 2773–2783 (2012).
- 721
- 722 19. Uher, J. Three methodological core issues of comparative personality research. *Euro-  
pean Journal of Personality* **22**, 475–496 (2008).
- 723
- 724 20. Griffin, A. S., Guillette, L. M. & Healy, S. D. Cognition and personality: An analysis  
of an emerging field. *Trends in Ecology & Evolution* **30**, 207–214 (2015).
- 725
- 726 21. Soha, J. A., Peters, S., Anderson, R. C., Searcy, W. A. & Nowicki, S. Performance  
on tests of cognitive ability is not repeatable across years in a songbird. *Animal*  
*Behaviour* **158**, 281–288 (2019).
- 727

- 728 22. Cauchoix, M. *et al.* The repeatability of cognitive performance: A meta-analysis.  
*Philosophical Transactions of the Royal Society B: Biological Sciences* **373**, 20170281  
(2018).
- 729 23. Völter, C. J., Tinklenberg, B., Call, J. & Seed, A. M. Comparative psychometrics:  
Establishing what differs is central to understanding what evolves. *Philosophical  
Transactions of the Royal Society B: Biological Sciences* **373**, 20170283 (2018).
- 730 24. Shaw, R. C. & Schmelz, M. Cognitive test batteries in animal cognition research:  
Evaluating the past, present and future of comparative psychometrics. *Animal Cog-  
nition* **20**, 1003–1018 (2017).
- 731 25. Matzel, L. D. & Sauce, B. Individual differences: Case studies of rodent and primate  
intelligence. *Journal of Experimental Psychology: Animal Learning and Cognition*  
**43**, 325 (2017).
- 732 26. Horn, L., Cimarelli, G., Boucherie, P. H., Šlipogor, V. & Bugnyar, T. Beyond the  
dichotomy between field and lab—the importance of studying cognition in context.  
*Current Opinion in Behavioral Sciences* **46**, 101172 (2022).
- 733 27. Damerius, L. A. *et al.* Orientation toward humans predicts cognitive performance in  
orang-utans. *Scientific Reports* **7**, 1–12 (2017).
- 734 28. Wobber, V., Herrmann, E., Hare, B., Wrangham, R. & Tomasello, M. Differences in  
the early cognitive development of children and great apes. *Developmental Psychobi-  
ology* **56**, 547–573 (2014).
- 735 29. Beran, M. J. & Hopkins, W. D. Self-control in chimpanzees relates to general intelli-  
gence. *Current Biology* **28**, 574–579 (2018).
- 736 30. Hopkins, W. D., Russell, J. L. & Schaeffer, J. Chimpanzee intelligence is heritable.  
*Current Biology* **24**, 1649–1652 (2014).
- 737 31. MacLean, E. L. *et al.* The evolution of self-control. *Proceedings of the National  
Academy of Sciences* **111**, E2140–E2148 (2014).

- 747
- 748 32. Kaufman, A. B., Reynolds, M. R. & Kaufman, A. S. The structure of ape (homi-  
749 noidea) intelligence. *Journal of Comparative Psychology* **133**, 92 (2019).
- 750 33. Herrmann, E., Call, J., Hernández-Lloreda, M. V., Hare, B. & Tomasello, M. Humans  
751 have evolved specialized skills of social cognition: The cultural intelligence hypothesis.  
*Science* **317**, 1360–1366 (2007).
- 752 34. Herrmann, E., Hernández-Lloreda, M. V., Call, J., Hare, B. & Tomasello, M. The  
753 structure of individual differences in the cognitive abilities of children and chim-  
panzees. *Psychological Science* **21**, 102–110 (2010).
- 754 35. Schmitt, V., Pankau, B. & Fischer, J. Old world monkeys compare to apes in the  
755 primate cognition test battery. *PloS one* **7**, e32024 (2012).
- 756 36. Völter, C. J. *et al.* The structure of executive functions in preschool children and  
757 chimpanzees. *Scientific Reports* **12**, 1–16 (2022).
- 758 37. Campbell, D. T. & Fiske, D. W. Convergent and discriminant validation by the  
759 multitrait-multimethod matrix. *Psychological bulletin* **56**, 81 (1959).
- 760 38. Anderson, B. Evidence from the rat for a general factor that underlies cognitive  
761 performance and that relates to brain size: intelligence? *Neuroscience Letters* **153**,  
98–102 (1993).
- 762 39. Matzel, L. D. *et al.* Individual differences in the expression of a ‘general’ learning  
763 ability in mice. *Journal of Neuroscience* **23**, 6423–6433 (2003).
- 764 40. Light, K. R. *et al.* Working memory training promotes general cognitive abilities in  
765 genetically heterogeneous mice. *Current biology* **20**, 777–782 (2010).
- 766 41. Keagy, J., Savard, J.-F. & Borgia, G. Complex relationship between multiple measures  
767 of cognitive ability and male mating success in satin bowerbirds, *ptilonorhynchus*  
*violaceus*. *Animal Behaviour* **81**, 1063–1070 (2011).

- 768 42. Isden, J., Panayi, C., Dingle, C. & Madden, J. Performance in cognitive and problem-solving tasks in male spotted bowerbirds does not correlate with mating success.  
769 *Animal Behaviour* **86**, 829–838 (2013).
- 770 43. Chandra, S. B., Hosler, J. S. & Smith, B. H. Heritable variation for latent inhibition  
771 and its correlation with reversal learning in honeybees (*apis mellifera*). *Journal of  
Comparative Psychology* **114**, 86 (2000).
- 772 44. Raine, N. E. & Chittka, L. No trade-off between learning speed and associative flex-  
773 ibility in bumblebees: A reversal learning test with multiple colonies. (2012).
- 774 45. Kolata, S. *et al.* Variations in working memory capacity predict individual differences  
775 in general learning abilities among genetically diverse mice. *Neurobiology of Learning  
and Memory* **84**, 241–246 (2005).
- 776 46. Wass, C. *et al.* Covariation of learning and ‘reasoning’ abilities in mice: Evolutionary  
777 conservation of the operations of intelligence. *Journal of Experimental Psychology:  
Animal Behavior Processes* **38**, 109 (2012).
- 778 47. Galsworthy, M. J. *et al.* Assessing reliability, heritability and general cognitive ability  
779 in a battery of cognitive tasks for laboratory mice. *Behavior Genetics* **35**, 675–692  
(2005).
- 780 48. Boogert, N. J., Anderson, R. C., Peters, S., Searcy, W. A. & Nowicki, S. Song reper-  
781toire size in male song sparrows correlates with detour reaching, but not with other  
cognitive measures. *Animal Behaviour* **81**, 1209–1216 (2011).
- 782 49. Bouchard, J., Goodyer, W. & Lefebvre, L. Social learning and innovation are posi-  
783 tively correlated in pigeons (*columba livia*). *Animal Cognition* **10**, 259–266 (2007).
- 784 50. Bräuer, J., Call, J. & Tomasello, M. All great ape species follow gaze to distant  
785 locations and around barriers. *Journal of Comparative Psychology* **119**, 145 (2005).

- 786 51. Call, J. Inferences about the location of food in the great apes (*pan paniscus*, *pan*  
787 *troglodytes*, gorilla *gorilla*, and *pongo pygmaeus*). *Journal of Comparative Psychology*  
788 **118**, 232 (2004).
- 789 52. Hanus, D. & Call, J. Discrete quantity judgments in the great apes (*pan paniscus*,  
790 *pan troglodytes*, gorilla *gorilla*, *pongo pygmaeus*): The effect of presenting whole sets  
versus item-by-item. *Journal of Comparative Psychology* **121**, 241 (2007).
- 791 53. Rosati, A. G., Stevens, J. R., Hare, B. & Hauser, M. D. The evolutionary origins of  
human patience: Temporal preferences in chimpanzees, bonobos, and human adults.  
*Current Biology* **17**, 1663–1668 (2007).
- 792 54. Haun, D. B., Call, J., Janzen, G. & Levinson, S. C. Evolutionary psychology of spatial  
793 representations in the hominidae. *Current Biology* **16**, 1736–1740 (2006).
- 794 55. Hedge, C., Powell, G. & Sumner, P. The reliability paradox: Why robust cognitive  
tasks do not produce reliable individual differences. *Behavior Research Methods* **50**,  
795 1166–1186 (2018).
- 796 56. Uher, J. Individual behavioral phenotypes: An integrative meta-theoretical frame-  
797 work. Why ‘behavioral syndromes’ are not analogs of ‘personality’. *Developmental  
Psychobiology* **53**, 521–548 (2011).
- 798 57. Steyer, R., Ferring, D. & Schmitt, M. J. States and traits in psychological assessment.  
799 *European Journal of Psychological Assessment* **8**, 79–98 (1992).
- 800 58. Steyer, R., Mayer, A., Geiser, C. & Cole, D. A. A theory of states and traits—revised.  
801 *Annual Review of Clinical Psychology* **11**, 71–98 (2015).
- 802 59. Geiser, C. *Longitudinal structural equation modeling with mplus: A latent state-trait  
803 perspective*. (Guilford Publications, 2020).
- 804 60. Cohen, J. A power primer. *Psychological Bulletin* **112**, 155–159 (1992).

- 806 61. Altschul, D. M., Wallace, E. K., Sonnweber, R., Tomonaga, M. & Weiss, A. Chimpanzee intellect: Personality, performance and motivation with touchscreen tasks.  
807 *Royal Society Open Science* **4**, 170169 (2017).
- 808 62. Morton, F. B., Lee, P. C. & Buchanan-Smith, H. M. Taking personality selection bias  
809 seriously in animal cognition research: A case study in capuchin monkeys (*sapajus  
apella*). *Animal Cognition* **16**, 677–684 (2013).
- 810 63. Altschul, D. M., Terrace, H. S. & Weiss, A. Serial cognition and personality in  
811 macaques. *Animal Behavior and Cognition* **3**, 46 (2016).
- 812 64. Piironen, J. & Vehtari, A. Comparison of bayesian predictive methods for model  
813 selection. *Statistics and Computing* **27**, 711–735 (2017).
- 814 65. Pavone, F., Piironen, J., Bürkner, P.-C. & Vehtari, A. Using reference models in  
815 variable selection. (2020).
- 816 66. Piironen, J., Paasiniemi, M. & Vehtari, A. Projective inference in high-dimensional  
817 problems: Prediction and feature selection. *Electronic Journal of Statistics* **14**, 2155–  
2197 (2020).
- 818 67. ManyPrimates *et al.* Collaborative open science as a way to reproducibility and new  
819 insights in primate cognition research. *Japanese Psychological Review* **62**, 205–220  
(2019).
- 820 68. Stevens, J. R. Replicability and reproducibility in comparative psychology. *Frontiers  
821 in Psychology* **8**, 862 (2017).
- 822 69. Farrar, B., Boeckle, M. & Clayton, N. Replications in comparative cognition: What  
823 should we expect and how can we improve? *Animal Behavior and Cognition* **7**, 1  
(2020).
- 824 70. Fried, E. I. & Flake, J. K. Measurement matters. *APS Observer* **31**, (2018).

- 826 71. Oakes, L. M. Sample size, statistical power, and false conclusions in infant looking-  
827 time research. *Infancy* **22**, 436–469 (2017).
- 828 72. Forstmeier, W., Wagenmakers, E.-J. & Parker, T. H. Detecting and avoiding likely  
829 false-positive findings—a practical guide. *Biological Reviews* **92**, 1941–1968 (2017).
- 830 73. ManyPrimates *et al.* Establishing an infrastructure for collaboration in primate cog-  
831 nition research. *PLoS One* **14**, e0223675 (2019).
- 832 74. ManyPrimates *et al.* The evolution of primate short-term memory. *Animal Behavior  
and Cognition* **9**, 428–516 (2022).
- 833 75. Bohn, M., Liebal, K. & Tessler, M. H. Great ape communication as contextual social  
834 inference: A computational modeling perspective. *Philosophical Transactions of the  
Royal Society B: Biological Sciences* **377**, 20210096 (2022).
- 835 76. Devaine, M. *et al.* Reading wild minds: A computational assay of theory of mind so-  
836 phistication across seven primate species. *PLoS Computational Biology* **13**, e1005833  
837 (2017).
- 838 77. Hanus, D. Causal reasoning versus associative learning: A useful dichotomy or a  
839 strawman battle in comparative psychology? *Journal of Comparative Psychology*  
**130**, 241 (2016).
- 840 78. Heyes, C. Simple minds: A qualified defence of associative learning. *Philosophical  
Transactions of the Royal Society B: Biological Sciences* **367**, 2695–2703 (2012).
- 841 79. Bollen, K. A. *Structural equations with latent variables*. **210**, (John Wiley & Sons,  
842 1989).
- 843 80. Hoyle, R. H. *Handbook of structural equation modeling*. (Guilford press, 2012).
- 844 81. Samejima, F. Estimation of latent ability using a response pattern of graded scores.  
845 *Psychometrika* 1–97 (1969).
- 846 847

- 848 82. Samejima, F. The graded response model. in *Handbook of modern item response*  
849 *theory* (eds. Linden, W. van der & Hambleton, R.) 85–100 (Springer, 1996).
- 850 83. Eid, M. & Kutscher, T. Statistical models for analyzing stability and change in hap-  
851 piness. in *Stability of happiness: Theories and evidence on whether happiness can*  
*change* (eds. Sheldon, K. & Lucas, R.) 261–297 (Elsevier, 2014).
- 852 84. Eid, M., Holtmann, J., Santangelo, P. & Ebner-Priemer, U. On the definition of latent-  
853 state-trait models with autoregressive effects: Insights from LST-r theory. *European*  
*Journal of Psychological Assessment* **33**, 285 (2017).
- 854 85. Van Erp, S., Oberski, D. L. & Mulder, J. Shrinkage priors for bayesian penalized  
855 regression. *Journal of Mathematical Psychology* **89**, 31–50 (2019).
- 856 86. Bürkner, P.-C. brms: An R package for Bayesian multilevel models using Stan. *Jour-*  
857 *nal of Statistical Software* **80**, 1–28 (2017).
- 858 87. Piironen, J., Paasiniemi, M., Catalina, A., Weber, F. & Vehtari, A. projpred: Projec-  
859 tion predictive feature selection. (2022).