

¹ Individual differences in great ape cognition across time and domains: stability, structure,
² and predictability

³ Manuel Bohn^{1,2}, Christoph Völter², Daniel Hanus², Nico Eisbrenner², Johanna Eckert², Jana
⁴ Holtmann³, & Daniel Haun²

⁵ ¹ Institute of Psychology in Education, Leuphana University Lüneburg

⁶ ² Department of Comparative Cultural Psychology, Max Planck Institute for Evolutionary
⁷ Anthropology, Leipzig, Germany

⁸ ³ Wilhelm Wundt Institute of Psychology, Leipzig University, Leipzig, Germany

¹⁰ Manuel Bohn was supported by a Jacobs Foundation Research Fellowship
¹¹ (2022-1484-00). We are grateful to thank all children and caregivers for participating in the
¹² study. We thank the Max Planck Society for the Advancement of Science.

¹³ The authors made the following contributions. Manuel Bohn: Conceptualization,
¹⁴ Formal Analysis, Writing - Original Draft Preparation, Writing - Review & Editing;
¹⁵ Christoph Völter: Conceptualization, Writing - Original Draft Preparation, Writing - Review
¹⁶ & Editing; Daniel Hanus: Conceptualization, Writing - Original Draft Preparation, Writing -
¹⁷ Review & Editing; Nico Eisbrenner: Formal Analysis, Writing - Original Draft Preparation,
¹⁸ Writing - Review & Editing; Johanna Eckert: Conceptualization, Writing - Original Draft
¹⁹ Preparation, Writing - Review & Editing; Jana Holtmann: Formal Analysis, Writing -
²⁰ Original Draft Preparation, Writing - Review & Editing; Daniel Haun: Conceptualization,
²¹ Writing - Review & Editing.

²² Correspondence concerning this article should be addressed to Manuel Bohn,
²³ Universitätsallee 1, 21335 Lüneburg, Germany. E-mail: manuel.bohn@leuphana.de

²⁴

Abstract

²⁵ 200 words

²⁶ *Keywords:* keywords

27 Individual differences in great ape cognition across time and domains: stability, structure,
28 and predictability

29 **Introduction**

30 Variation is the fodder of evolution. Individual differences in cognitive abilities are key
31 for understanding what evolves (Shaw & Schmelz, 2017; Völter, Tinklenberg, Call, & Seed,
32 2018). They inform us which aspects of cognition are invariant and which are more
33 malleable. They also inform us about the broader structure of the cognitive architecture
34 when studying relations between individual differences in different aspects of cognition.
35 Finally, they hold the key to understanding which socio-ecological factors shape cognition
36 during ontogeny and phylogeny.

37 structure:

38 Völter, Reindl, et al. (2022) Herrmann, Hernández-Lloreda, Call, Hare, and Tomasello
39 (2010)

40 Fichtel, Dinter, and Kappeler (2020)

41 Wobber, Herrmann, Hare, Wrangham, and Tomasello (2014); Beran and Hopkins
42 (2018); Hopkins, Russell, and Schaeffer (2014); MacLean et al. (2014); Kaufman, Reynolds,
43 and Kaufman (2019)

44 old vs newwold monkeys: schmitt2012old

45 factors influencing:

46 Most theorizing happens on a species level Dunbar and Shultz (2017); Rosati (2017)

47 big picture studies using aggregated data: Deaner, Van Schaik, and Johnson (2006)
48 Piantadosi and Kidd (2016) . Comparability of the data for different species is problematic

49 because they were often collected using very different methods (Schubiger, Fichtel, &
50 Burkart, 2020). – but see ManyPrimates et al. (2022)

51 small scale comparisons between e.g. more or less tolerant species from the same genus
52 Rosati and Santos (2017) Joly et al. (2017) or species showing different social dynamics
53 Amici, Aureli, and Call (2008)

54 much less on an individual level:

55 effects of enculturation - see Berio and Moore (2023) for a recent summary

56 Hopper et al. (2014) personality traits influence problem solving; Brosnan et al. (2015)
57 personality correlated with responses to inequity (see also Carter, Marshall, Heinsohn, and
58 Cowlishaw (2014) for monkeys)

59 Bard, Bakeman, Boysen, and Leavens (2014) human reared chimps are better at social
60 cognition - but groups differ in many respects.

61 Watson et al. (2018) hand reared individuals are more likely to use social information

62 Rosati, DiNicola, and Buckholtz (2018) no link between self control and cooperation

63 measurement: few studies on individual level becasue studying individual differnices is
64 hard. Matzel and Sauce (2017) Boogert, Madden, Morand-Ferron, and Thornton (2018)
65 reliability: Griffin, Guillette, and Healy (2015)

66 Bohn et al. (2023) studied

67 The current study extended previous work in two important aspects. First, we study a
68 broader range of cognitive domains including social cognition, reasoning about quantities,
69 executive functions and inferential reasoning. This allows us to assess whether the results
70 obtained by Bohn et al. (2023) replicate within domains and generalize to others. Second,

71 we explored the structure of great ape cognition in more depth: we pooled the data collected
72 here with the data from Bohn et al. (2023) to study the correlations between cognitive traits
73 within and across domains.

74

Methods

75 **Participants**

76 A total of 48 great apes participated at least in one tasks at one time point. This
77 included 12 Bonobos (*pan paniscus*, 4 females, age 3.60 to 40.70 years), 24 Chimpanzees (*pan*
78 *troglodytes*, 17 females, age 3.80 to 57.80 years), 6 Gorillas (*gorilla gorilla*, 4 females, age 4.40
79 to 24.40 years), and 6 Orangutans (*pongo abelii*, 5 females, age 4.70 to 43.10 years). The
80 sample size at the different time points ranged from 34 to 45 for the different species (see
81 supplementary material for details). All apes participated in cognitive research on a regular
82 basis. Apes were housed at the Wolfgang Köhler Primate Research Center located in Zoo
83 Leipzig, Germany. They lived in groups, with one group per species and two chimpanzee
84 groups (group A and B). Research was noninvasive and strictly adhered to the legal
85 requirements in Germany. Animal husbandry and research complied with the European
86 Association of Zoos and Aquaria Minimum Standards for the Accommodation and Care of
87 Animals in Zoos and Aquaria as well as the World Association of Zoos and Aquariums
88 Ethical Guidelines for the Conduct of Research on Animals by Zoos and Aquariums.
89 Participation was voluntary, all food was given in addition to the daily diet, and water was
90 available ad libitum throughout the study. The study was approved by an internal ethics
91 committee at the Max Planck Institute for Evolutionary Anthropology.

92 **Procedure**

93 Apes were tested in familiar sleeping or observation rooms by a single experimenter.
94 The basic setup comprised a sliding table positioned in front of a mesh or a clear plexiglas
95 panel. The experimenter sat on a small stool and used an occluder to cover the table (see

⁹⁶ Figure 1).

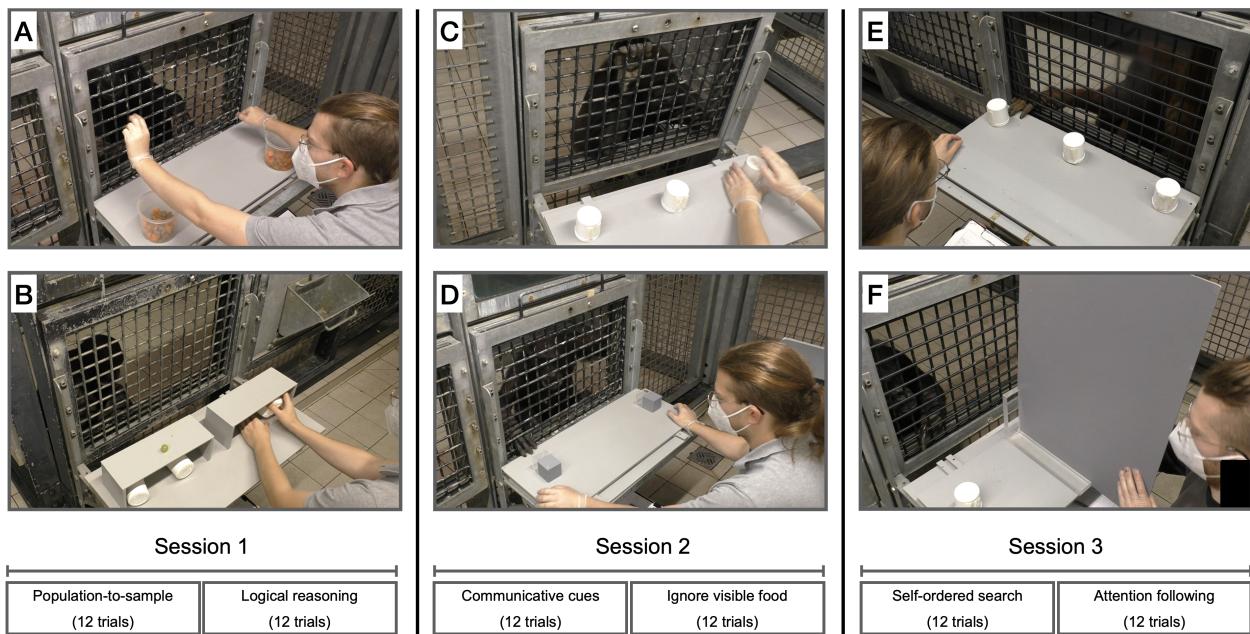


Figure 1. Setup used for the six tasks. A) population-to-sample, B) logical-reasoning, C) communicative-cues, D) ignore-visible-food, E) self-ordered-search and F) attention-following. Text at the bottom shows order of task presentation and trial numbers

⁹⁷ The study involved a total of six cognitive tasks. These were based on published
⁹⁸ procedures in the field of comparative psychology. The original publications often include
⁹⁹ control conditions to rule out alternative, cognitively less demanding ways to solve the tasks.
¹⁰⁰ We did not include such controls here and only ran the experimental conditions. For each
¹⁰¹ task, we refer to these papers if they want to know more about control conditions and/or a
¹⁰² detailed discussion of the nature of the underlying cognitive mechanisms. Example videos for
¹⁰³ each task can be found in the associated online repository. In the following, we give a brief
¹⁰⁴ description of each task. Additional details can be found in the supplementary material.

¹⁰⁵ **Attention-following.** The Attention-following task was loosely modeled after
¹⁰⁶ Kaminski, Call, and Tomasello (2004). The setup consisted of two identical cups placed on
¹⁰⁷ the sliding table and a large opaque screen that was longer than the width of the sliding
¹⁰⁸ table (Supplementary Figure 1F). The experimenter placed both cups on the table and

109 showed the ape that they were empty. Then, the experimenter baited both cups in view of
110 the ape and placed the opaque screen in the center between the two cups, perpendicular to
111 the mesh. Next, the experimenter moved to one side and looked at the cup in front of them.
112 Then, the experimenter pushed the sliding table forward and the ape was allowed to choose
113 one of the cups by pointing at it. If the ape chose the cup that the experimenter was looking
114 at, they received the food item. If they choose the other cup, they did not. We coded
115 whether the ape chose the side the experimenter was looking at (correct choice) or not. Apes
116 received twelve trials. The side at which the experimenter looked was counterbalanced with
117 same number of looks to each side and looks to the same side not more than two times in a
118 row. We assumed that apes follow the experimenters focus of attention to determine whether
119 or not their request could be seen and thus be successful.

120 **Communicative-cues.** This task was modeled after Schmid, Karg, Perner, and
121 Tomasello (2017). Three identical cups were placed equidistantly on a sliding table directly
122 in front of the ape (Figure 1C). In the beginning of a trial, the experimenter showed the ape
123 that all cups are empty. After placing an occluder between the subject and the cups, the
124 experimenter held up one food item and moved it behind the occluder, visiting all three cups
125 but baiting only one. Next, the occluder was lifted and E looked at the ape (ostensive cue),
126 called the name, and looked at one of the cups, while holding on to it with one hand and
127 tapping it with the other (continuous looking, 3 times tapping). Finally, the experimenter
128 pushed the sliding table forward for the ape to make a choice. If the ape chose the baited
129 cup, they received the reward – if not, not. We coded as correct choice if the ape chose the
130 indicated cup. Apes received twelve trials. The location of the indicated cup was
131 counterbalanced, with each cup being the target equally often and the same target not more
132 than two times in a row. We assumed that apes use the experimenter's communicative cues
133 to determine where the food is hidden.

134 **Ignore-visible-food.** The task was modeled after Völter, Tinklenberg, Call, and
135 Seed (2022). The task involved two opaque cups with an additional, sealed but transparent,

136 compartment attached to the front of each cup (facing the ape). For one cup, the
137 compartment contained a preferred food item that was clearly visible, for the other cup, the
138 compartment was empty (Figure 1D). In the beginning of the trial, the two cups were placed
139 upside down on the sliding table so that the ape could see that the opaque compartments of
140 both cups were empty. Next, the experimenter baited one of the cups in full view of the
141 subject. In non-conflict trials, the baited cup was the cup with the food item in the
142 transparent compartment. In conflict trials, the baited cup was the cup with the empty
143 compartment. After baiting the experimenter pushed the sliding table forwards and the ape
144 could chose by pointing. If the baited cup was chosen, the ape received the food. Apes
145 received 14 trials, twelve conflict trials and two non-conflict trials (1st and 8th trial). Only
146 conflict trials were analyzed. The location of the cup with the baited compartment was
147 counterbalanced, with the cup not being in the same location more than twice in a row. We
148 assumed that apes need to inhibit selecting the visible food item and instead use their
149 short-term memory to remember where the food was hidden.

150 **Logical-reasoning.** The task was modeled after Hanus and Call (2014). Three
151 identical cups were presented side-by-side on a sliding table, with the cup in the middle
152 sometimes positioned closer to the left cup and sometimes closer to the right.
153 (Supplementary Figure 1B). Two half-open boxes served as occluders to block the ape's view
154 when shuffling the cups. Each trial started by showing the ape that all three cups (one on
155 one side of the table, two on the other) were empty. After placing the occluders over both
156 sides of the table, the experimenter put one piece of food on top of each occluder. Next, the
157 experimenter hid each piece of food under the cup(s) behind the occluders. In case of the
158 occluder with the two cups, the food was randomly placed under one of the two cups while
159 both cups were visited and even shuffled. Finally, both occluders were lifted and the table
160 pushed forwards, allowing the ape to choose one of the three cups, from which they then
161 received the content. We coded whether the ape chose the certain cup (i.e. the cup from the
162 side of the table with only one cup). Apes received 12 trials. The side with one cup was

¹⁶³ counterbalanced, with the same constellation appearing not more than two times in a row on
¹⁶⁴ the same side. We assumed that apes would infer that the cup from the tray with only one
¹⁶⁵ cup certainly contains food while the other cups contain food only in 50% of cases.

¹⁶⁶ **Population-to-sample.** The task was modeled after Eckert, Call, Hermes,
¹⁶⁷ Herrmann, and Rakoczy (2018). During the test, apes saw two transparent buckets filled
¹⁶⁸ with pellets and carrot pieces (the carrot pieces had roughly the same size and shape as the
¹⁶⁹ pellets). Each bucket contained 80 food items. The distribution of pellets to carrot pieces
¹⁷⁰ was 4:1 in bucket A, and 1:4 in bucket B. Pellets are preferred food items in comparison to
¹⁷¹ carrots. The experimenter placed both buckets on a table, one left, one right (Figure 1A). In
¹⁷² the beginning of a trial, the experimenter picked up the bucket on the right side, tilted it
¹⁷³ forward so the ape could see inside, placed it back on the table and turned it around 360°.
¹⁷⁴ The same procedure was repeated with the other bucket. Next, the experimenter looked at
¹⁷⁵ the ceiling, inserted each hand in the bucket in front of it and drew one item from the bucket
¹⁷⁶ without the ape seeing which type (E picked always of the majority type). The food items
¹⁷⁷ remained hidden in the experimenter's fists. Next, the experimenter extended the arms (in
¹⁷⁸ parallel) towards the ape who was then allowed to make a choice by pointing to one of the
¹⁷⁹ fists. The ape received the chosen sample. In half of the trials, the experimenter crossed
¹⁸⁰ arms when moving the fists towards the ape to ensure that the apes made a choice between
¹⁸¹ samples and not just chose the side where the favorable population was still visible. In
¹⁸² between trials, the buckets were refilled to restore the original distributions. Apes received
¹⁸³ twelve trials. We coded whether the ape chose the sample from the population with the
¹⁸⁴ higher number of high quality food items. The location of the buckets (left and right) was
¹⁸⁵ counterbalanced, with the buckets in the same location no more than two times in a row.
¹⁸⁶ The crossing of the hands was also counterbalanced with no more than two crossings in a
¹⁸⁷ row. We assumed that apes reasoned about the probability of the sample being a high
¹⁸⁸ quality item based on observing the ratio in the population.

189 **Self-ordered-search.** The task was modeled after Völter, Mundry, Call, and Seed

190 (2019; Diamond, Prevor, Callender, and Druin, 1997; see also Petrides, 1995). Three

191 identical cups were placed equidistantly on a sliding table directly in front of the ape

192 (Supplementary Figure 1E). The experimenter baited all three cups in full view of the ape.

193 Next, the experimenter pushed the sliding table forwards for the ape to choose one of the

194 cups by pointing. After the choice, the table was pulled back and the ape received the food.

195 After a 3s pause, the table was pushed forward again for a second choice. This procedure

196 was repeated for a third choice. If the ape chose a baited cup, they received the food, if not,

197 not. We coded the number of times the ape chose an empty cup (i.e. chose a cup they

198 already chose before). Please note that this outcome variable differed from the other tasks in

199 two ways: first, possible values were 0, 1, and 2 (instead of just 0 and 1) and second, a lower

200 score indicated better performance. Apes received twelve trials. No counterbalancing was

201 needed. We assumed that apes use their working memory abilities to remember where they

202 had already searched and which cups still contained food.

203 **Predictor variables.** In addition to the data from the cognitive tasks, we collected

204 data for a range of predictor variables to predict individual differences in performance in the

205 cognitive tasks. Predictors could either vary with the individual (stable individual

206 characteristics: group, age sex, rearing history, and time spent in research), vary with

207 individual and time point (variable individual characteristics: rank, sickness, and sociality),

208 vary with group membership (group life: time spent outdoors, disturbances, and life events),

209 or vary with the testing arrangements and thus with individual, time point and session

210 (testing arrangements: presence of an observer, participation in other studies on the same

211 day or since the last time point). Predictors were collected from the zoo handbook with

212 demographic information about the apes, via a diary that the animal caretakers filled out on

213 a daily basis, or via proximity scans of the whole group. We provide a detailed description of

214 these variables in the supplementary material.

215 Data collection

216 Data collection started on April 28th, 2022, lasted until October 7th, 2023 and included
217 10 time points. One time point meant running all tasks with all participants. Within each
218 time point, the tasks were organized in three sessions (see Fig. 1). Session 1 included the
219 population-to-sample and logical-reasoning tasks, session 2 the communicative-cues and
220 ignore-visible-food tasks and session 3 the self-ordered-search and attention-following tasks.

221 The interval between two time points was planned to be eight weeks. However, it was
222 not always possible to follow this schedule so that some intervals were longer or shorter (see
223 supplementary material for details). The order of tasks was the same for all subjects. So was
224 the counterbalancing within each task. This exact procedure was repeated at each time point
225 so that the results would be comparable across participants and time points.

226 Analysis, results and discussion

227 To get an overview of the results, we first visualized the data (Fig. 2). Performance
228 was consistently above chance in the communicative-cues, ignore-visible-food and
229 population-to-sample tasks. For attention-following, this was the case only from time point 7
230 onward and for logical-reasoning, performance was, if anything, below chance. For the
231 self-ordered-search task, performance was below chance but here lower values reflect better
232 performance (i.e. systematic avoidance of the visible food item). For attention-following,
233 ignore-visible-food, communicative-cues and self-ordered-search there was a steady
234 improvement in performance over time.

235 In the following, we link performance in the tasks across time points to latent variables
236 representing cognitive abilities. We first ask how stable these abilities are over time and how
237 reliably they are measured. Next, we study the correlations between different abilities to
238 explore the internal structure of great ape cognition. Finally, we link performance in the
239 tasks to external predictors to shed light on the sources of individual differences in abilities.

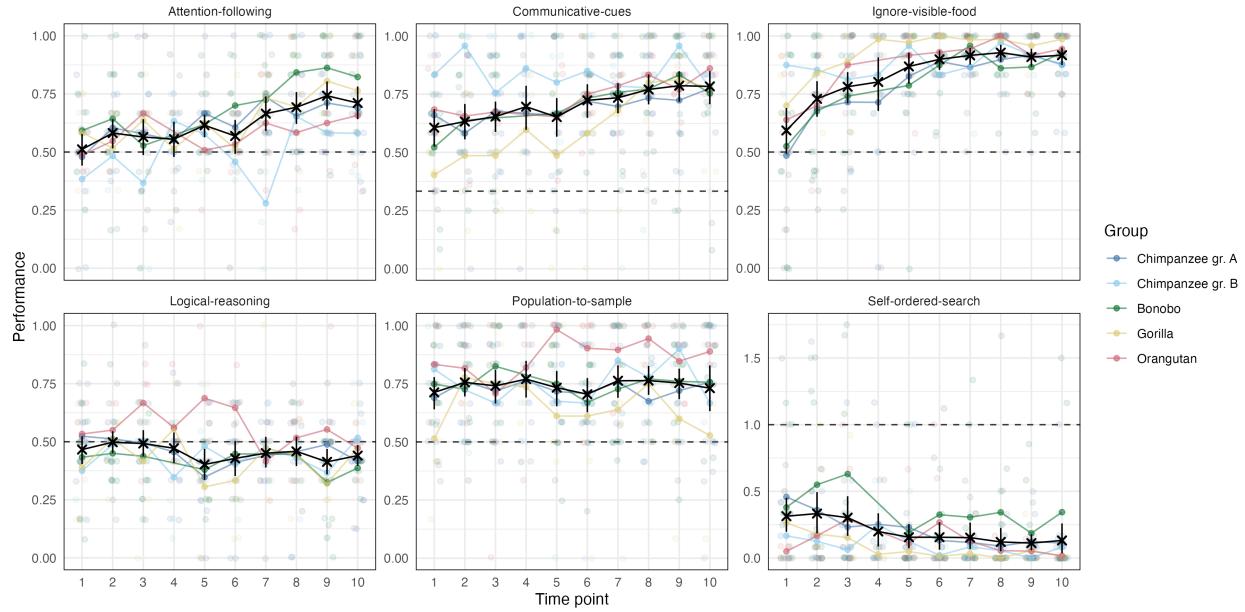


Figure 2. Results from the six cognitive tasks across time points. Black crosses show mean performance at each time point across species (with 95% CI). The sample size varied between time points and can be found in Supplementary Figure 1. Colored dots show mean performance by species. Dashed line shows chance level performance.

240 Each section uses different statistical techniques which we describe in the respective section.

241 Stability and reliability

242 We first asked how robust performance was on a task-level, how stable individual
 243 differences were and how reliable the measures were. We used *Structural Equation Modeling*
 244 (SEM) (Bollen, 1989; Hoyle, 2012) to address these questions¹. For each task we fit three
 245 types of models that addressed different questions. We provide a detailed, mathematical
 246 description of the models in the supplementary material.

¹ SEMs usually use larger sample sizes than available in the present study. Bohn et al. (2023) reported a simulation study showing that parameters could be accurately estimated using Bayesian estimation techniques and reasonable model restrictions with sample sizes comparable to one we have here. We lay out the restrictive assumptions we imposed on the parameters in the supplementary material.

We started with a latent state (LS) model. The goal of this model is to estimate a measurement-error free latent state, representing an individual's cognitive ability, for each time point. Measurement error is captured by dividing the trials from one time point into two test-halves. Roughly speaking, the correlation between these two test-halves is an indicator of measurement precision and used to estimate measurement error (and reliability). Robustness of task-level performance can be assessed by comparing the means of latent states across subjects for the different time points. Stability of individual differences can be assessed by correlating latent states across different time points.

The temporal robustness of latent state means varied across tasks (Fig. 3A). In attention-following, means increased over time and were significantly different from zero at later time points (9 and 10). Communicative-cues and ignore-visible-food exhibited steady increases, though ignore-visible-food saw a late-stage decline, with the latent mean at time point 10 still significantly different from 0. Self-ordered-search showed a decrease (reduction in errors) from time point 6 onward, while latent means for logical-reasoning and population-to-sample remained stable throughout the study.

Correlations between latent states illustrated varying degrees of stability of individual differences across tasks (Fig. 3B). Attention-following displayed low-to-moderate correlations at early time points (before time point 7), increasing substantially thereafter. Communicative-cues, ignore-visible-food, and self-ordered-search generally showed high correlations between latent states (with time point 1 of ignore-visible-food being an exception). Population-to-sample correlations were consistently high, while logical-reasoning showed generally low, sometimes even negative, correlations, suggesting no stability across time points.

Next, we fit two types of latent state-trait (LST) models. In comparison to the LS models, these models assume that there is a single latent trait, representing an individual's stable cognitive ability, that is the same across time points. This way we can partition

273 variation in performance on a given time point into variance due to the trait (consistency),
274 variance due to the occasion (occasion specificity; 1 - consistency), and measurement error
275 (used to estimate reliability). Like the latent states in the LS model, the trait in the LST
276 model is assumed to be measurement error free (Geiser, 2020; Steyer, Ferring, & Schmitt,
277 1992; Steyer, Mayer, Geiser, & Cole, 2015). The first LST model we fit assumed that neither
278 the absolute trait values nor the ranking of individuals changes over time (fixed means).
279 This is the classic version of an LSTM. The second model allowed the absolute trait values
280 to change over time while the ranking of individuals was fixed (varying means). Change over
281 time according to this model is thus seen as change that is the same for all individuals. In
282 both cases, stability of individual differences can be assessed by the proportion of variance
283 explained by the trait (consistency).

284 Consistency estimates varied across tasks (Fig. 3C). In attention-following, the
285 consistency coefficient was estimated to be 0.92 [Jana: possible to add 95%CrI?] for the fixed
286 means model and 0.95 for the varying means model, indicating that more than 90% of true
287 inter-individual differences were attributable to stable traits. However, given the low
288 reliability of measurement (see below), this result should be interpreted with caution. [Jana:
289 expand what that means]. For communicative-cues, consistency estimates differed the most
290 between models and were higher in the varying means model (0.76) compared to the fixed
291 means model (0.64). The reasons for this discrepancy is most likely the substantial change in
292 mean performance over time in the task (see Fig. 3A). Ignore-visible-food showed similar
293 consistency across models, with values of 0.59 (fixed means) and 0.64 (varying means).
294 Logical-reasoning showed a similiar pattern to attention-following: Cosistency was estimated
295 to be high (fixed means: 0.81; varying means: 0.80) but reliability was low so that the same
296 restrictions for interpretation apply. Self-ordered-search and population-to-sample had high
297 consistency estimates according to both models: 0.79 (fixed means) and 0.84 (varying means)
298 for self-ordered-search and 0.83 (fixed means) and 0.84 (varying means) for
299 population-to-sample.

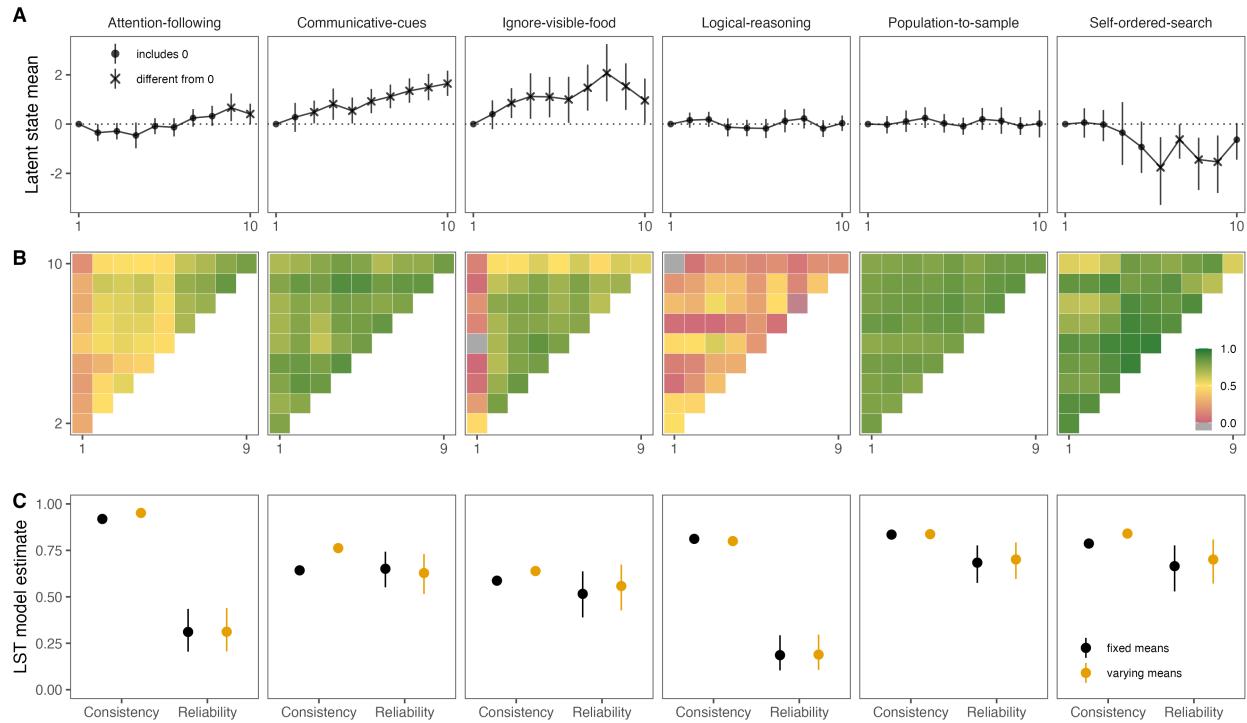


Figure 3. A) Latent mean estimates for each time point by task based on latent state model. Means at time point 1 are set to zero. Shape denotes whether the 95% CrI included zero (dashed line). The sample size varied between time points and can be found in Supplementary Fig. 1. B) Correlations between subject-level latent state estimates for the different time points by task. C) Mean estimates from latent state-trait models with fixed and varying means (color coded) with 95% CrI. Consistency refers to the proportion of (measurement-error-free) variance in performance explained by stable trait differences. Reliability refers to the proportion of true score variance to variance in raw scores.

300 Reliability of measurement also varied significantly across tasks, based on the LST
301 models (Fig. 3C). For attention-following, reliability was initially low (fixed means: 0.31;
302 varying means: 0.31), but was substantially higher when only considering time points 7 and
303 onward (see supplementary material). Communicative-cues showed moderate reliability
304 (fixed means: 0.65; varying means: 0.63). Ignore-visible-food also had moderate reliability
305 across time points (fixed means: 0.52; varying means: 0.56). As mentioned above,
306 logical-reasoning exhibited very low reliability (fixed means: 0.19; varying means: 0.19).
307 Population-to-sample showed acceptable reliability (fixed means: 0.68; varying means: 0.70).
308 Self-ordered-search also exhibited acceptable reliability levels (fixed means: 0.66; varying
309 means: 0.70).

310 To summarize the SEM results, we saw that the six tasks differed substantially in what
311 they revealed about group- and individual-level variation. What stands out is the
312 widespread change in performance over time. For all tasks except population-to-sample and
313 logical-reasoning we observed an improvement in performance over time. This group-level
314 change, however, has different individual-level interpretations for the different tasks. For
315 communicative-cues, ignore-visible-food and self-ordered-search, individual differences
316 remained relatively stable despite the group-level change suggesting stable individual
317 differences combined with a systematic learning effect across individuals. In contrast, for
318 attention-following, there was little stability in individual differences at earlier time points
319 and only towards the end emerged a more stable ordering of individuals. In combination
320 with the low reliability at earlier time points, this suggests that at least some individuals
321 changed their response strategy in the course of the study. The combination of low reliability,
322 chance-level performance and low correlation of latent states for logical-reasoning suggests
323 that this task is not suited to assess individual differences in logical reasoning abilities in
324 great apes. It is also noteworthy that the reliability estimates are on average lower compared
325 to a previous study testing the same individuals on different tasks (Bohn et al., 2023). One
326 explanation might be the increase in performance over time. At the beginning of the study,

327 more individuals might have chosen randomly instead of using the available information
328 provided in the task setup and the demonstrations. By definition, random variation is not
329 reliable. With time, more and more individuals started using the available information so
330 that inter-individual differences in how good they are in using it could be detected.

331 **Structure**

332 To explore the structure of great ape cognition we correlated latent trait estimates for
333 each task. In contrast to raw performance scores, these estimates take into account the
334 reliability of measurement and are considered to be measurement-error free. Bohn et al.
335 (2023) tested the same individuals and we therefore also include the data from tasks reported
336 there (data from phase 2). Even though the data in the two studies was collected at different
337 time points, we think it is justifiable to analyse them jointly because the trait estimates
338 represent stable, time-invariant individual differences in cognitive abilities. The estimates
339 were computed ...

340 Figure 4 shows the correlations between trait estimates for the different tasks. Overall,
341 most correlations were not significantly different from zero (i.e. the 95% CI did include zero).
342 Because of this low average level of correlations, we decided not to explore models with
343 higher-order factors and will only interpret the qualitative patterns.

344 Conceptually, the tasks can be clustered in the following broader domains: *social*
345 *cognition* (attention-following, gaze-following, communicative-cues), *reasoning about*
346 *quantities* (quantity-discrimination, population-to-sample), *executive functions*
347 (delay-of-gratification, self-ordered-search, ignore-visible-food) and *inferential reasoning*
348 (logical-reasoning, causal-inference, inference-by-exclusion). As a first step, we will evaluate
349 whether we find evidence for such a clustering in the data.

350 There was no significant correlation between any of the social cognition tasks.
351 Furthermore, attention-following and gaze-following did not correlate significantly with any

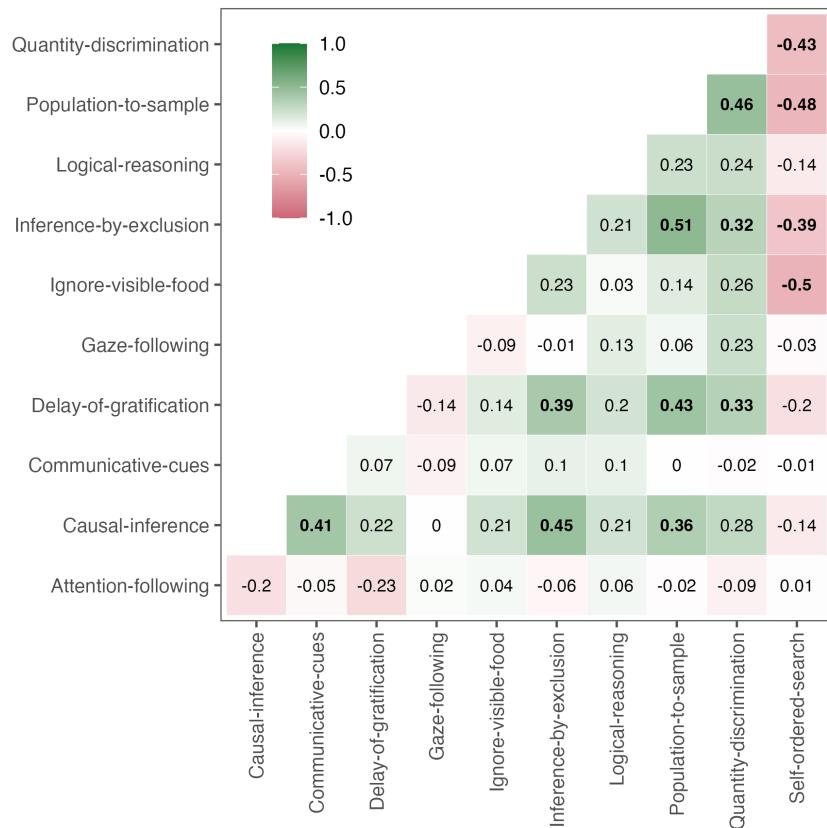


Figure 4. Correlations between ... trait estimates. Bold correlations have 95% CrI not overlapping with zero.

352 of the other tasks and communicative-cues correlated only with causal-inference – a result we
 353 will discuss below. Thus, and in line with previous work (Herrmann et al., 2010), we found
 354 no evidence for shared cognitive processes in tasks measuring different aspects of social
 355 cognition.

356 The two tasks measuring reasoning about quantities did correlate significantly. Both
 357 tasks require discriminating between different quantities, directly in the case of
 358 quantity-discrimination and as part of the decision making process in the case of
 359 population-to-sample. Deciding between the samples from the two populations requires
 360 discriminating between the relative quantities within each bucket from which the samples
 361 were drawn.

362 Within the executive functions measures, self-ordered-search and inhibit-visible-food
363 were significantly correlated but none of the two correlated with delay-of-gratification. The
364 significant correlation can be explained by the need to inhibit a premature response
365 (selecting visible food or a cup that was previously rewarded) in both tasks. It has been
366 argued that delay-of-gratification requires self control (tolerating a longer waiting time to
367 gain a more valuable reward) over and above behavioral inhibition (Beran, 2015). From this
368 point of view, individual differences in the delay-of-gratification task might be due to
369 differences in self control and less due to differences in inhibition.

370 Finally, for the three inferential reasoning measures we found a correlation between
371 inference-by-exclusion and causal-inference. Logical-reasoning did not correlate with either
372 (neither did it with any other task). This is not surprising given the results reported above:
373 the observed variation in the logical-reasoning task was largely noise and did not reflect
374 systematic individual differences. The correlation between causal-inference and
375 inference-by-exclusion is most likely due to the fact that both tasks involve making
376 inferences about the location of food based on reasoning about its physical properties.

377 Next we turn to the correlations across domains. Perhaps the most surprising finding is
378 the correlation between causal-inference and communicative-cues. On a closer look, the
379 origin might be task impurity in that there are two ways to solve the causal-inference task:
380 first, as hypothesized, by using the rattling sound to infer the location of the food. Second,
381 by interpreting the experimenter's shaking of the cup as a communicative cue, which is very
382 similar to the communicative-cues task. Thus, we suspect that at least some individuals
383 solved the task via the second route.

384 Finally, there was a cluster of significant correlations between delay-of-gratification,
385 self-ordered-search, inference-by-exclusion, causal-inference, population-to-sample and
386 quantity discrimination. Of the 15 possible correlations, only four were non-significant. One
387 commonality between these tasks that might – in part – explain this pattern is that they all

388 benefit from sustained attention to the task. Sustained attention facilitates the processing of
389 the experimenter's demonstrations (population-to-sample, inference-by-exclusion,
390 causal-inference, delay-of-gratification), ones one actions on the setup (self-ordered-search) or
391 visually complex stimuli (quantity discrimination). Tentative support for this idea comes
392 from the analysis of relevant predictors (see Bohn et al., 2023 and below) in which `time`
393 `spent in research` was selected as a relevant predictor of performance for all of these tasks
394 except causal-inference. This predictor reflects individual's experience with experimental
395 studies, which often involve sustained attention to distributions of food items, actions of
396 conspecifics and/or demonstrations by experimenters. Next, we turn to the sources of the
397 individual differences analysed here.

398 Predictability

399 In this section, we analysed which external variables accounted for for inter- and
400 intra-individual differences in task performance. That is, we asked which of the predictor
401 variables described above predicted performance in the different tasks. Given the large
402 number of predictor variables (14), this question translates to a variable selection problem:
403 selecting a subset of variables from a larger pool. We used the projection predictive inference
404 (Piironen, Paasiniemi, & Vehtari, 2020) approach because it is a state-of-the-art procedure
405 that provides an excellent trade-off between model complexity and accuracy (Pavone,
406 Piironen, Bürkner, & Vehtari, 2020; Piironen & Vehtari, 2017). The projection prediction
407 approach is a two-step process: The first step consists of building the best predictive model
408 possible, called the reference model. In our case, the reference model is a Bayesian multilevel
409 regression model – fit via `brms` (Bürkner, 2017) – including all available predictors (Catalina,
410 Bürkner, & Vehtari, 2020). In the second step, the goal is to replace the posterior
411 distribution of the reference model with a simpler distribution containing fewer predictors
412 compared to the reference model. The importance of a predictor is assessed by inspecting
413 the mean log-predictive density (`elpd`) and root-mean-squared error (`rmse`) of models

414 containing the predictor vs. not.

415 The output of the procedure is a ranking of the different predictors. That is, for each
416 task, we get a ranking of how important a predictor is for constructing the simpler
417 replacement distribution. In addition, we can make a qualitatively assessment of whether or
418 not a predictor is relevant or not. In addition to the global assessment, we also inspected the
419 projected posterior distribution of the predictors classified as relevant to see how they
420 influenced performance. In the supplementary material we provide a detailed description of
421 the procedure including how the different variables were handled and how the importance of
422 each predictor was assessed.

423 In addition to the external predictors, the models also included a random intercept
424 term for subject ((1 | subject) in `brms` notation). This predictor was handled in a special
425 way in that it was always considered last because it would otherwise have soaked up most of
426 the variance before the other predictors would have had a chance to explain any of it.

427 Fig. 5A summarizes the selected predictors across tasks. For all tasks, the random
428 intercept term improved model fit the most (not shown in Fig. 5A). In line with results
429 reported by Bohn et al. (2023), this suggests that idiosyncratic developmental processes or
430 genetic pre-dispositions, which operate on time-scales longer than what we captured in our
431 study, accounted for a substantial portion of the variance in cognitive abilities between
432 individuals.

433 However, for two tasks, other predictors had an comparable explanatory power –
434 something that was not observe in Bohn et al. (2023). For population-to-sample, `time`
435 `spent in research` improved the model fit even more than adding the random intercept at
436 the end did. This could be interpreted that performance in this task strongly depends on
437 having learned to pay attention to stimuli and the human experimenter. For
438 ignore-visible-food, `time point` had an influence exceeding that of the random intercept

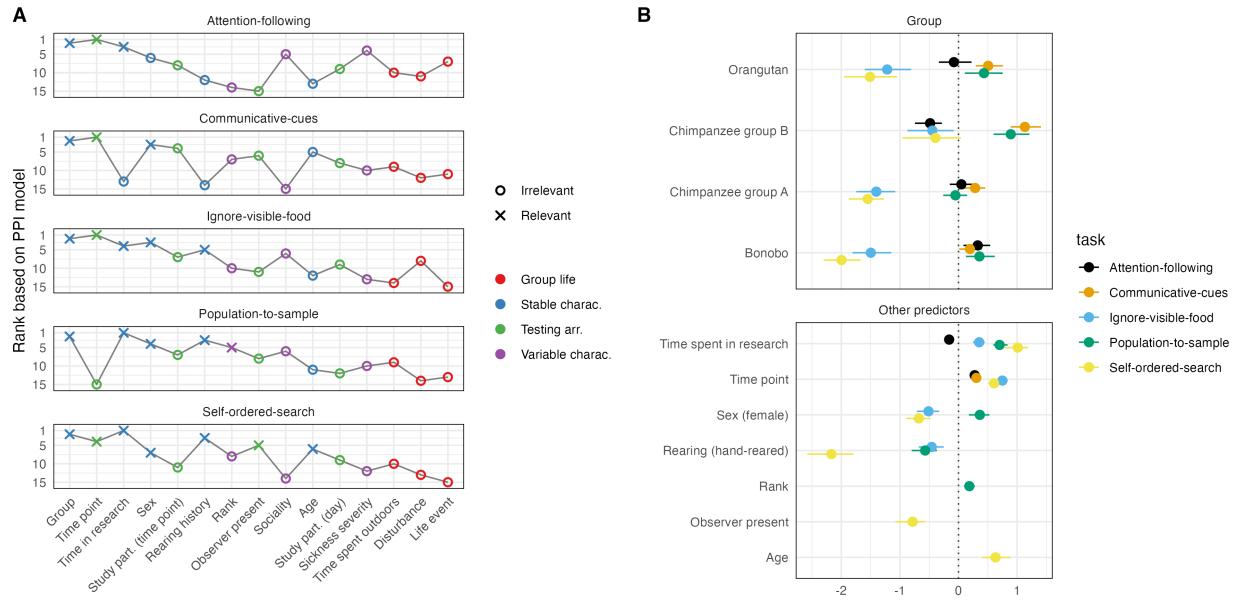


Figure 5. A) Predictor ranking and selection based on PPI models. Crosses mark predictors that were selected to be relevant based on the PPI models. Color shows the broader category each predictor belongs to. The x-axis is sorted by the average rank across tasks. B) Posterior model estimates for the selected predictors for each task based on data. Points show means with 95% Credible Interval. Color denotes task. For categorical predictors, the estimate gives the difference compared to the reference level (Gorilla for group).

439 term. We think this result reflects the strong within-task learning effect across subjects.
 440 Because performance increased substantially with time, most of the variation captured by
 441 **time point** exceeded the variation between individuals.

442 For the remaining predictors, the most highly-ranked and frequently selected ones
 443 came from the group of stable individual characteristics. The big exception being **time**
 444 **point**, which was ranked second across tasks. This pattern aligns with the SEM results, in
 445 which we saw that most of the variance in performance could be traced back to stable trait
 446 differences between individuals. The remaining occasion specific variation was largely due to
 447 improvement over time, most likely reflecting task-specific learning processes. The remaining
 448 time-varying predictors did not account for much variation.

449 The predictor selected most often was **group**. It was the only predictor that was
450 selected as relevant for all tasks. However, differences between groups were variable in that
451 the ranking of the groups changed from task to task (Fig. 5B). For example, Gorillas
452 performed best in ignore-visible-food and self-ordered-search, the Chimpanzee group B
453 performed best in communicative-cues and population-to-sample and the Bonobos performed
454 best in attention-following. This speaks against clear species or group differences in general
455 cognitive performance. Again, the most likely explanation for group differences is an
456 interaction between species specific dispositions and individual- / task-level developmental
457 processes.

458 The predictors that were selected more than once influenced performance in variable
459 ways (Fig. 5B). As mentioned above, **time point** always had a positive effect because
460 performance increased with time. Whenever **rearing** was selected to be relevant,
461 mother-reared individuals outperformed others. **Time spent in research** had a positive
462 effect, suggesting that more experience with research leads to better performance, except for
463 attention-following. The effect of **sex** was variable in that females outperformed males in
464 population-to-sample but males outperformed females in self-ordered-search and
465 ignore-visible-food.

466 General Discussion

467 no evidence for something like g contra eg (Banerjee et al., 2009)
468 kaufman2019structure. they did not take into account the measurement properties of the
469 different tasks but used sum scores

470 on an individual level, in long-lived species like primates sociality etc play littel roles. if
471 anythin rearing history.

472 Conclusion

References

- 473
- 474 Amici, F., Aureli, F., & Call, J. (2008). Fission-fusion dynamics, behavioral flexibility, and
475 inhibitory control in primates. *Current Biology*, 18(18), 1415–1419.
- 476 Banerjee, K., Chabris, C. F., Johnson, V. E., Lee, J. J., Tsao, F., & Hauser, M. D. (2009).
477 General intelligence in another primate: Individual differences across cognitive task
478 performance in a new world monkey (*saguinus oedipus*). *PLoS One*, 4(6), e5883.
- 479 Bard, K. A., Bakeman, R., Boysen, S. T., & Leavens, D. A. (2014). Emotional engagements
480 predict and enhance social cognition in young chimpanzees. *Developmental Science*,
481 17(5), 682–696.
- 482 Beran, M. J. (2015). The comparative science of “self-control”: What are we talking about?
483 *Frontiers in Psychology*, 6, 51.
- 484 Beran, M. J., & Hopkins, W. D. (2018). Self-control in chimpanzees relates to general
485 intelligence. *Current Biology*, 28(4), 574–579.
- 486 Berio, L., & Moore, R. (2023). Great ape enculturation studies: A neglected resource in
487 cognitive development research. *Biology & Philosophy*, 38(2), 17.
- 488 Bohn, M., Eckert, J., Hanus, D., Lugauer, B., Holtmann, J., & Haun, D. B. (2023). Great
489 ape cognition is structured by stable cognitive abilities and predicted by developmental
490 conditions. *Nature Ecology & Evolution*, 7(6), 927–938.
- 491 Bollen, K. A. (1989). *Structural equations with latent variables* (Vol. 210). John Wiley &
492 Sons.
- 493 Boogert, N. J., Madden, J. R., Morand-Ferron, J., & Thornton, A. (2018). Measuring and
494 understanding individual differences in cognition. *Philosophical Transactions of the Royal
495 Society B: Biological Sciences*, 373(1756), 20170280.
- 496 Brosnan, S. F., Hopper, L. M., Richey, S., Freeman, H. D., Talbot, C. F., Gosling, S. D., ...
497 Schapiro, S. J. (2015). Personality influences responses to inequity and contrast in
498 chimpanzees. *Animal Behaviour*, 101, 75–87.
- 499 Bürkner, P.-C. (2017). brms: An R package for Bayesian multilevel models using Stan.

- 500 *Journal of Statistical Software*, 80(1), 1–28.
- 501 Carter, A. J., Marshall, H. H., Heinsohn, R., & Cowlishaw, G. (2014). Personality predicts
502 the propensity for social learning in a wild primate. *PeerJ*, 2, e283.
- 503 Catalina, A., Bürkner, P.-C., & Vehtari, A. (2020). Projection predictive inference for
504 generalized linear and additive multilevel models. *arXiv Preprint arXiv:2010.06994*.
- 505 Deaner, R. O., Van Schaik, C. P., & Johnson, V. (2006). Do some taxa have better
506 domain-general cognition than others? A meta-analysis of nonhuman primate studies.
507 *Evolutionary Psychology*, 4(1), 147470490600400114.
- 508 Diamond, A., Prevor, M. B., Callender, G., & Druin, D. P. (1997). Prefrontal cortex
509 cognitive deficits in children treated early and continuously for PKU. *Monographs of the*
510 *Society for Research in Child Development*, i–206.
- 511 Dunbar, R., & Shultz, S. (2017). Why are there so many explanations for primate brain
512 evolution? *Philosophical Transactions of the Royal Society B: Biological Sciences*,
513 372(1727), 20160244.
- 514 Eckert, J., Call, J., Hermes, J., Herrmann, E., & Rakoczy, H. (2018). Intuitive statistical
515 inferences in chimpanzees and humans follow weber’s law. *Cognition*, 180, 99–107.
- 516 Fichtel, C., Dinter, K., & Kappeler, P. M. (2020). The lemur baseline: How lemurs compare
517 to monkeys and apes in the primate cognition test battery. *PeerJ*, 8, e10025.
- 518 Geiser, C. (2020). *Longitudinal structural equation modeling with mplus: A latent state-trait*
519 *perspective*. Guilford Publications.
- 520 Griffin, A. S., Guillette, L. M., & Healy, S. D. (2015). Cognition and personality: An
521 analysis of an emerging field. *Trends in Ecology & Evolution*, 30(4), 207–214.
- 522 Hanus, D., & Call, J. (2014). When maths trumps logic: Probabilistic judgements in
523 chimpanzees. *Biology Letters*, 10(12), 20140892.
- 524 Herrmann, E., Hernández-Lloreda, M. V., Call, J., Hare, B., & Tomasello, M. (2010). The
525 structure of individual differences in the cognitive abilities of children and chimpanzees.
526 *Psychological Science*, 21(1), 102–110.

- 527 Hopkins, W. D., Russell, J. L., & Schaeffer, J. (2014). Chimpanzee intelligence is heritable.
528 *Current Biology*, 24(14), 1649–1652.
- 529 Hopper, L. M., Price, S. A., Freeman, H. D., Lambeth, S. P., Schapiro, S. J., & Kendal, R. L.
530 (2014). Influence of personality, age, sex, and estrous state on chimpanzee
531 problem-solving success. *Animal Cognition*, 17, 835–847.
- 532 Hoyle, R. H. (2012). *Handbook of structural equation modeling*. Guilford press.
- 533 Joly, M., Micheletta, J., De Marco, A., Langermans, J. A., Sterck, E. H., & Waller, B. M.
534 (2017). Comparing physical and social cognitive skills in macaque species with different
535 degrees of social tolerance. *Proceedings of the Royal Society B: Biological Sciences*,
536 284(1862), 20162738.
- 537 Kaminski, J., Call, J., & Tomasello, M. (2004). Body orientation and face orientation: Two
538 factors controlling apes' begging behavior from humans. *Animal Cognition*, 7, 216–223.
- 539 Kaufman, A. B., Reynolds, M. R., & Kaufman, A. S. (2019). The structure of ape
540 (hominoidea) intelligence. *Journal of Comparative Psychology*, 133(1), 92.
- 541 MacLean, E. L., Hare, B., Nunn, C. L., Addessi, E., Amici, F., Anderson, R. C., et al.others.
542 (2014). The evolution of self-control. *Proceedings of the National Academy of Sciences*,
543 111(20), E2140–E2148.
- 544 ManyPrimates, Aguenounon, G., Allritz, M., Altschul, D. M., Ballesta, S., Beaud, A., et
545 al.others. (2022). The evolution of primate short-term memory. *Animal Behavior and*
546 *Cognition*, 9(4), 428–516.
- 547 Matzel, L. D., & Sauce, B. (2017). Individual differences: Case studies of rodent and primate
548 intelligence. *Journal of Experimental Psychology: Animal Learning and Cognition*, 43(4),
549 325.
- 550 Pavone, F., Piironen, J., Bürkner, P.-C., & Vehtari, A. (2020). *Using reference models in*
551 *variable selection*. Retrieved from <https://arxiv.org/abs/2004.13118>
- 552 Petrides, M. (1995). Impairments on nonspatial self-ordered and externally ordered working
553 memory tasks after lesions of the mid-dorsal part of the lateral frontal cortex in the

- monkey. *Journal of Neuroscience*, 15(1), 359–375.
- Piantadosi, S. T., & Kidd, C. (2016). Extraordinary intelligence and the care of infants. *Proceedings of the National Academy of Sciences*, 113(25), 6874–6879.
- Piironen, J., Paasiniemi, M., & Vehtari, A. (2020). Projective inference in high-dimensional problems: Prediction and feature selection. *Electronic Journal of Statistics*, 14(1), 2155–2197. <https://doi.org/10.1214/20-EJS1711>
- Piironen, J., & Vehtari, A. (2017). Comparison of bayesian predictive methods for model selection. *Statistics and Computing*, 27, 711–735.
<https://doi.org/10.1007/s11222-016-9649-y>
- Rosati, A. G. (2017). Foraging cognition: Reviving the ecological intelligence hypothesis. *Trends in Cognitive Sciences*, 21(9), 691–702.
- Rosati, A. G., DiNicola, L. M., & Buckholtz, J. W. (2018). Chimpanzee cooperation is fast and independent from self-control. *Psychological Science*, 29(11), 1832–1845.
- Rosati, A. G., & Santos, L. R. (2017). Tolerant barbary macaques maintain juvenile levels of social attention in old age, but despotic rhesus macaques do not. *Animal Behaviour*, 130, 199–207.
- Schmid, B., Karg, K., Perner, J., & Tomasello, M. (2017). Great apes are sensitive to prior reliability of an informant in a gaze following task. *PLoS One*, 12(11), e0187451.
- Schubiger, M. N., Fichtel, C., & Burkart, J. M. (2020). Validity of cognitive tests for non-human animals: Pitfalls and prospects. *Frontiers in Psychology*, 11, 1835.
- Shaw, R. C., & Schmelz, M. (2017). Cognitive test batteries in animal cognition research: Evaluating the past, present and future of comparative psychometrics. *Animal Cognition*, 20(6), 1003–1018.
- Steyer, R., Ferring, D., & Schmitt, M. J. (1992). States and traits in psychological assessment. *European Journal of Psychological Assessment*, 8, 79–98.
- Steyer, R., Mayer, A., Geiser, C., & Cole, D. A. (2015). A theory of states and traits—revised. *Annual Review of Clinical Psychology*, 11, 71–98.

- 581 Völter, C. J., Mundry, R., Call, J., & Seed, A. M. (2019). Chimpanzees flexibly update
582 working memory contents and show susceptibility to distraction in the self-ordered search
583 task. *Proceedings of the Royal Society B*, 286(1907), 20190715.
- 584 Völter, C. J., Reindl, E., Felsche, E., Civelek, Z., Whalen, A., Lugosi, Z., ... Seed, A. M.
585 (2022). The structure of executive functions in preschool children and chimpanzees.
586 *Scientific Reports*, 12(1), 1–16.
- 587 Völter, C. J., Tinklenberg, B., Call, J., & Seed, A. M. (2018). Comparative psychometrics:
588 Establishing what differs is central to understanding what evolves. *Philosophical
589 Transactions of the Royal Society B: Biological Sciences*, 373(1756), 20170283.
- 590 Völter, C. J., Tinklenberg, B., Call, J., & Seed, A. M. (2022). Inhibitory control and cue
591 relevance modulate chimpanzees'(*pan troglodytes*) performance in a spatial foraging task.
592 *Journal of Comparative Psychology*, 136(2), 105.
- 593 Watson, S. K., Vale, G. L., Hopper, L. M., Dean, L. G., Kendal, R. L., Price, E. E., et
594 al.others. (2018). Chimpanzees demonstrate individual differences in social information
595 use. *Animal Cognition*, 21, 639–650.
- 596 Wobber, V., Herrmann, E., Hare, B., Wrangham, R., & Tomasello, M. (2014). Differences in
597 the early cognitive development of children and great apes. *Developmental Psychobiology*,
598 56(3), 547–573.