

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32

Main Manuscript for

Boundaries of executive functions in large language models: GPT-4o selectively replicates human performance

Tongyi Zhang¹, Jiang Qiu², and Xin Zhao^{1*}

¹School of Psychology, Northwest Normal University, Lanzhou, China

²Key Laboratory of Cognition and Personality of the Ministry of Education, Faculty of Psychology, Southwest University

*Corresponding author: Xin Zhao

Email: psyzhaoxin@nwnu.edu.cn

Phone: +86-931-7975316

Authors' Contributions: T.Z. conceived and designed the experiments, conducted the experiments, analyzed the data, and drafted and revised the manuscript. J.Q. provided validated samples with fMRI data resources and contributed to manuscript revision. X.Z. conceived and designed the experiments, supervised the study, contributed to participant recruitment and data collection, and drafted and revised the manuscript. All authors reviewed and approved the final manuscript.

Conflicts of Interest: The authors declare no conflicts of interest.

Funding Statement: This work was supported by the National Natural Science Foundation of China (Grant No. 32260207 to X.Z.).

Data, Materials, and Software Availability. Transcripts and numerical coding of transcript data have been deposited in OSF (<https://osf.io/hfnt6>).

Classification: Biological Sciences / Psychological and Cognitive Sciences

Keywords: Large language models; Executive functions; GPT; Log probability; Prefrontal cortex

This PDF file includes:

Main Text
Figures 1 to 6

33 **Abstract**

34 Large language models (LLMs) have demonstrated considerable capabilities in complex
35 reasoning and problem-solving tasks that, in humans, depend on executive functions (EFs).
36 However, the extent to which these models reproduce human EF patterns remains unclear. We
37 systematically evaluated the performance of GPT-4o across three core EF dimensions—inhibitory
38 control, working memory, and cognitive flexibility—using established behavioral paradigms.
39 Additionally, we examined whether model-internal log probability parameters could serve as
40 quantitative indicators of cognitive processing analogous to human neural activity. Using two
41 independent datasets ($N_1 = 1,970$; $N_2 = 39$), we simulated trial-by-trial responses in GPT-4o while
42 recording log probability metrics. Bayesian analyses revealed a selective replication of human EF
43 patterns. GPT-4o reproduced human-like performance in interference inhibition (*Stroop* task),
44 working memory capacity (*digit span* task), and working memory updating (*n-back* task).
45 Conversely, the model showed divergent patterns in response inhibition (*Go/No-Go* task), time-
46 sensitive working memory updating (*running memory* task, 1,750 and 750 ms presentation times),
47 and cognitive flexibility (*number-switching* task). Log probability parameters demonstrated task-
48 specific associations with behavioral measures and corresponded with activation patterns in EF-
49 related brain regions during working memory and task-switching paradigms. These findings
50 implied that GPT-4o captures specific aspects of human EFs—particularly those involving
51 symbolic representation and static information maintenance—while showing limitations in
52 dynamic control and temporal processing. This selective replication pattern offers insights into
53 both the computational basis of EF and the cognitive boundaries of current LLM architectures.
54 Notably, log probability parameters may provide a valuable methodological framework for
55 evaluating artificial cognitive mechanisms.

56 **Significance statement**

57 Executive functions (EFs)—cognitive control processes central to human intelligence—are
58 traditionally linked to prefrontal architecture. We systematically evaluated GPT-4o using
59 established EF paradigms to determine which capacities emerge from statistical language
60 learning. GPT-4o selectively replicated human patterns: successfully simulating interference
61 resolution and working memory maintenance, but showing divergent patterns in inhibitory control
62 and cognitive flexibility. This suggests certain EFs arise from domain-general computational
63 principles, while others require specialized mechanisms absent in current LLMs. Model-internal
64 log probabilities correlated with behavioral performance and brain activation, establishing a
65 quantitative framework for comparing artificial and biological cognition. These findings advance
66 understanding of EF computational foundations and provide methodologies for evaluating
67 cognitive capabilities in artificial systems, with implications for cognitive science and AI
68 development.

69

70

71 **Introduction**

72

73 Large language models (LLMs) now achieve remarkable performance on tasks traditionally used
74 to assess human reasoning and problem-solving, including multi-step reasoning (1), strategic
75 planning (2), and abstract concept manipulation (3). This convergence between machine and
76 human performance raises fundamental questions regarding the computational requirements of
77 cognition: Can statistical learning systems develop processing mechanisms analogous to those
78 underlying human cognitive functions? A positive answer to this question would profoundly impact
79 our understanding of cognitive computation and provide new methodological tools for cognitive
80 science (4–8). By systematically comparing LLMs and humans behaviorally, we can begin to
81 identify which cognitive capacities arise from large-scale language learning and which depend on
82 biological or embodied constraints unique to human cognition (9, 10).

83 Executive functions (EFs) provide a particularly informative test case, as they constitute the
84 cognitive control system coordinating goal-directed behavior and enabling complex reasoning and
85 decision-making (11–13). Initial investigations into LLM performance on EF tasks have yielded
86 mixed results (details in SI Appendix, Table S1). Regarding working memory, LLMs display
87 capacity constraints similar to those in humans (14, 15), although performance varies
88 considerably across models and appears more dependent on training scale than architecture
89 (16). The scenario is more complex for cognitive flexibility and inhibitory control, where LLMs
90 consistently struggle with dynamic rule adaptation, conflict resolution, and prospective planning,
91 remaining well below human benchmarks (2, 17, 18). Moreover, while LLMs can reproduce
92 specific psychological phenomena, they exhibit distinct patterns of cognitive biases and
93 capabilities that do not straightforwardly improve with model scale (17, 19, 20). These findings
94 collectively suggest that the correspondence between LLM and human EF is selective rather than
95 general.

96
97 The multidimensional nature of EF presents both challenges and opportunities for this
98 investigation. According to the influential framework proposed by Miyake et al. (2000), EF is
99 composed of three core, yet separable, components: inhibitory control, working memory, and
100 cognitive flexibility (21). This tripartite structure is supported by converging neuroimaging
101 evidence linking working memory with the dorsolateral prefrontal cortex (22), conflict monitoring
102 with the anterior cingulate cortex (23), and response inhibition with the inferior frontal gyrus (24).
103 Despite this well-established framework, previous studies have typically examined LLM
104 performance on individual executive tasks, rather than systematically evaluating the full tripartite
105 structure. This piecemeal approach leaves unanswered whether LLMs possess selective
106 strengths and weaknesses across executive domains or demonstrate more uniform capabilities.
107 Mapping these patterns would clarify which EFs can emerge from statistical learning alone *versus*
108 those requiring specialized cognitive architecture (25).

109
110 A key challenge for current research is how to relate the underlying computational processes to
111 behavioral outcomes. Accordingly, establishing meaningful parallels between LLM and human
112 cognition requires moving beyond surface-level similarities to understanding the mechanisms.
113 Whereas in cognitive neuroscience processing dynamics are inferred from neural activity,
114 comparable methods for examining LLM internal states remain underdeveloped. The log
115 probability (logprobs) parameters available in OpenAI’s GPT models offer one promising
116 approach (26). These metrics—perplexity, probability variance, average log probability, and
117 minimum log probability—serve to quantify the model’s uncertainty and computational complexity
118 during response generation (27–29). Systematic patterns in these parameters may reveal
119 processing demands analogous to cognitive load in humans, potentially bridging behavioral and
120 mechanistic levels of analysis. However, the relationship between these internal metrics and
121 cognitive processes remains largely unexplored.

122
123 In this study, to address these knowledge gaps, we conducted a comprehensive evaluation of
124 GPT-4o’s performance across the three-factor EF framework, employing standard cognitive tasks
125 mirroring human studies. We tested whether current LLMs exhibit patterns consistent with
126 established cognitive phenomena in inhibitory control, working memory, and cognitive flexibility
127 (Fig. 1). Beyond behavioral assessment, we also analyzed log probability parameters to
128 characterize the model’s internal processing dynamics and examine the correspondence of these
129 parameters with human neural activity patterns derived from task-based functional magnetic
130 resonance imaging (fMRI) studies.

131 132 133 **Results**

134 135 **Selective replication of human EF effects in GPT-4o**

136 **EF dimension I: Inhibitory Control**

137 In the *Stroop* task, human participants exhibited the typical interference effect. One-way analysis
138 of variance (ANOVA) revealed that significant differences existed in reaction times across the
139 three conditions ($p = 0.001$) (Fig. 2A). Post-hoc analyses indicated that reaction times did not
140 significantly differ between the congruent and neutral trial types ($p = 0.363$), but were significantly
141 longer in the incongruent condition than in both the congruent ($p < 0.001$) and neutral ($p = 0.009$)
142 conditions. GPT-4o exhibited a comparable interference pattern on the same task. One-way
143 ANOVA showed that reaction times differed significantly across conditions ($p < 0.001$) (Fig. 2A).
144 Specifically, reaction times in the incongruent trial were significantly longer than those in both the
145 congruent ($p < 0.001$) and neutral ($p = 0.002$) states, while no significant difference was observed
146 between congruent and neutral conditions ($p = 0.384$). Notably, Bayesian analysis further
147 indicated that the mean interference effect of GPT-4o (89.210 ms) exceeded that of humans
148 (51.570 ms), with a 95% Highest Density Interval (HDI) of [26.160, 52.620] ms, and a Bayes
149 Factor (BF_{10}) of >100 . Moreover, analysis of GPT-4o's output probability distributions revealed
150 that there were significant differences in perplexity across conditions ($p < 0.001$), with incongruent
151 trials showing higher perplexity than either congruent ($p < 0.001$) or neutral ($p < 0.001$) ones (Fig.
152 2A). Similar patterns were observed for log probability variance ($p < 0.001$), average log
153 probability ($p < 0.001$), and minimum log probability ($p < 0.001$). Furthermore, correlation
154 analyses showed that GPT-4o's reaction times were positively correlated with perplexity ($r =$
155 0.306 , $p_{\text{fdr}} < 0.001$) and log probability variance ($r = 0.280$, $p_{\text{fdr}} < 0.001$), but negatively correlated
156 with average log probability ($r = -0.314$, $p_{\text{fdr}} < 0.001$) and minimum log probability ($r = -0.272$, p_{fdr}
157 < 0.001) (Fig. 2C). Condition-specific analyses revealed that in congruent trials, reaction times
158 correlated positively with perplexity ($r = 0.339$, $p < 0.001$) and log probability variance ($r = 0.404$,
159 $p_{\text{fdr}} < 0.001$). In incongruent trials, reaction times showed a positive correlation with perplexity ($r =$
160 0.358 , $p_{\text{fdr}} < 0.001$) and a negative correlation with average log probability ($r = -0.383$, $p_{\text{fdr}} <$
161 0.001). No significant correlations were observed in the neutral condition (all $p_{\text{fdr}} > 0.050$).
162

163 Similarly, in the *Go/No-Go* task, human participants exhibited typical response inhibition
164 difficulties, showing significantly lower accuracy on No-Go trials than on Go trials ($p < 0.001$) (Fig.
165 2B). In contrast, GPT-4o achieved 100% accuracy on both trial types. Consequently, Bayesian
166 analysis yielded a posterior mean difference in inhibitory control effects of 0.110 (95% HDI:
167 [0.086, 0.133], $BF_{10} > 100$). Despite this perfect behavioral accuracy, GPT-4o's internal processing
168 states revealed differential patterns, with No-Go trials showing significantly higher perplexity than
169 Go trials ($p < 0.001$). This difference was also evident across log probability variance ($p < 0.001$),
170 average log probability ($p < 0.001$), and minimum log probability ($p < 0.001$) (Fig. 2B).
171

172 **EF dimension II: working memory**

173 In the *digit span* task, human participants demonstrated classic working memory capacity
174 characteristics, with forward digit span significantly exceeding backward digit span ($p < 0.001$)
175 (Fig. 3A). GPT-4o replicated this effect pattern ($p < 0.001$), though its digit span capacity
176 exceeded that of human participants (-2.830 , 95% HDI [-3.660 , -1.960], BF_{10} approaching 0).
177 Furthermore, analysis of GPT-4o's internal processing states revealed significant differences in
178 perplexity between conditions ($p < 0.001$), with backward tasks showing higher perplexity than
179 forward tasks. Log probability variance also differed significantly ($p < 0.001$), with backward tasks
180 showing greater variance. Both average log probability ($p < 0.001$) and minimum log probability (p
181 < 0.001) were significantly lower in backward tasks, reflecting reduced predictive certainty (Fig.
182 3A). Moreover, correlation analyses showed that in forward tasks, sequence length correlated
183 positively with perplexity ($r = 0.306$, $p_{\text{fdr}} < 0.001$), and negatively with both average log probability
184 ($r = -0.307$, $p_{\text{fdr}} < 0.001$) and minimum log probability ($r = -0.313$, $p_{\text{fdr}} < 0.001$). In backward
185 tasks, these correlations were markedly stronger: $r = 0.601$, $p_{\text{fdr}} < 0.001$ for perplexity; $r = -0.594$,
186 $p_{\text{fdr}} < 0.001$ for average log probability; and $r = -0.597$, $p_{\text{fdr}} < 0.001$ for minimum log probability
187 (Fig. 3C). In contrast, in the *running memory* task, human participants showed a refresh time
188 effect, achieving higher average accuracy under the 1750-ms stimulus presentation than under
189 the 750-ms stimulus ($p = 0.001$) (Fig. 3B). However, GPT-4o showed no significant difference

190 between conditions ($p = 0.833$). Accordingly, Bayesian analysis revealed an effect difference of
191 -0.012 , with a 95% HDI of $[-0.049, 0.020]$, and a $BF_{10} = 3.090$. Nevertheless, analysis of GPT-
192 4o's internal processing states showed that there were significant differences in perplexity
193 between conditions ($p = 0.002$), with lower perplexity observed in the 750 ms condition than in the
194 1750 ms condition. Similarly, both average log probability ($p = 0.002$) and minimum log probability
195 ($p = 0.002$) values were higher in the 750 ms condition (Fig. 3B).
196

197 **EF dimension III: cognitive flexibility**

198 In the *number-switching* task, human participants demonstrated classic task-switching costs,
199 achieving significantly longer reaction times in mixed-task switch trials than in repeat trials (switch
200 cost, $p < 0.001$), and markedly longer reaction times in mixed-task repeat trials than in single-task
201 trials (mixing cost, $p < 0.001$). While GPT-4o also showed a significant switch cost ($p = 0.008$),
202 this cost was smaller than that observed in humans (Fig. 4A), an effect that was further confirmed
203 through Bayesian analysis (-145.090 ms, 95% HDI $[-202.210, -90.730]$ ms, BF_{10} approaching
204 infinity). Additionally, GPT-4o's mixing cost was 58.190 ms. Between-condition comparisons
205 showed that there was no significant difference between mixed-task switch trials and repeat trials
206 ($p = 0.958$); however, significant differences were found between mixed-task conditions and
207 single-task conditions (switch vs. single: $p = 0.008$; repeat vs. single: $p = 0.009$). Notably, analysis
208 of GPT-4o's internal processing states did not identify significant between-condition differences in
209 perplexity (all $p > 0.050$), log probability variance, or average log probability (Fig. 4A).
210

211 **Replication in an independent sample across three EF dimensions**

212 To validate our findings, we analyzed three EF tasks from Dataset 2. In the *n*-back task (Fig. 5A),
213 LLMs demonstrated working memory load effects. Specifically, accuracy decreased with
214 increasing *n*-back levels (0-back: 100%; 1-back: 93.2%; 2-back: 82.1%), resulting in significant
215 main effects ($p < 0.001$). Post-hoc comparisons revealed significant differences across all
216 conditions: 0-back versus 1-back ($d = 2.850$), 0-back versus 2-back ($d = 2.340$), and 1-back
217 versus 2-back ($d = 1.380$). While reaction times increased with working memory load, the
218 increase was not significant ($p = 0.228$). Overall, *n*-back effects led to a 17.910% accuracy
219 decrease ($p < 0.001$). Additionally, internal state parameters varied with working memory load.
220 Perplexity and log probability variance decreased with increasing load, with *n*-back effects of
221 -0.011 and -0.005 , respectively (both $p < 0.001$). Conversely, average and minimum log
222 probabilities became more positive as the load increased (both $p < 0.001$). Furthermore,
223 correlation analyses (Fig. 5D) showed that model accuracy was positively correlated with
224 perplexity ($r = 0.529$, $p_{\text{fdr}} < 0.001$) and log probability variance ($r = 0.534$, $p_{\text{fdr}} < 0.001$), and
225 negatively correlated with average log probability ($r = -0.533$, $p_{\text{fdr}} < 0.001$) and minimum log
226 probability ($r = -0.587$, $p_{\text{fdr}} < 0.001$). Reaction times showed only weak correlations with log
227 probability parameters (all $p_{\text{fdr}} > 0.050$).
228

229 Similarly, in the *stop-signal* task (Fig. 5B), LLMs demonstrated response inhibition capabilities.
230 Successful stop trials showed increased perplexity compared to successful Go trials ($p < 0.001$, d
231 $= 0.880$) and higher log probability variance ($p = 0.002$, $d = 0.690$). Correspondingly, average log
232 probability decreased ($p < 0.001$, $d = -0.870$), and minimum log probability showed more
233 negative values ($p < 0.001$, $d = -1.030$). Finally, in the *letter-number-switching* task (Fig. 5C),
234 LLMs did not show typical human switch cost patterns. The reaction time switch cost was -0.630
235 ms ($p = 0.228$, $d = 0.022$). However, internal state parameters revealed switch cost effects.
236 Specifically, perplexity was higher in switch trials than in repeat trials, yielding a switch cost of
237 0.006 ($p < 0.001$, $d = 1.070$). Similarly, log probability variance, average log probability, and
238 minimum log probability all showed significant switch costs (all $p < 0.001$, $d = 1.030-1.110$).
239

240 **Potential correspondence between LLM log probability parameters and human brain 241 activation**

242 After examining LLM behavioral performance, we explored whether internal computational
243 parameters might correspond to human neural activity patterns. In the *stop-signal* task

244 "successful stop > successful Go" contrast (Fig. 6A), we found no significant correlations between
245 brain region activation and LLM internal computational parameters. However, in the *n*-back
246 working memory task "2-back > 0-back" contrast (Fig. 6B), we identified brain-model
247 correspondences. Specifically, left gyrus rectus activation was positively correlated with model
248 perplexity ($r = 0.340$, $p_{\text{fdr}} = 0.037$) and negatively correlated with average log probability ($r =$
249 -0.338 , $p_{\text{fdr}} = 0.038$). The right medial superior frontal gyrus showed comparable patterns, with its
250 activation positively correlating with perplexity ($r = 0.328$, $p_{\text{fdr}} = 0.044$) and negatively correlating
251 with average log probability ($r = -0.325$, $p_{\text{fdr}} = 0.046$). Furthermore, in the switch cost analysis of
252 the *task-switching* paradigm (Fig. 6C), left thalamus activation showed a positive correlation with
253 model minimum log probability ($r = 0.323$, $p_{\text{fdr}} = 0.048$) and a negative correlation with perplexity
254 ($r = -0.322$, $p_{\text{fdr}} = 0.048$).

255

256 **Sensitivity analysis in subsampling and different temperature settings**

257 **Subsampling Analysis:** To assess the robustness of our findings and ensure valid comparisons
258 between the large human sample ($N_1 = 1,970$) and the GPT-4o responses ($N = 120$), we
259 conducted bootstrap resampling analyses (Fig. 4B). We randomly sampled 120 human
260 participants with replacement across 3,000 iterations and compared each subsample to the GPT-
261 4o data. For the *Stroop* task, 99.8% of human subsamples showed significant interference effects
262 (mean $d = 0.372$, 95% CI [0.371, 0.374]). GPT-4o demonstrated a larger effect ($d = 0.516$), with
263 all iterations confirming this difference ($\text{BF}_{10} \geq 10$). For the *Go/No-Go* task, human participants
264 consistently showed inhibitory control difficulties (mean $d = -0.866$, 95% CI [-0.870, -0.861]),
265 whereas GPT-4o achieved near-perfect accuracy ($d < 0.001$; $\text{BF}_{10} > 100$ across all iterations). In
266 the *digit span* task, both groups showed forward-backward span differences (human: $d = -0.701$,
267 95% CI [-0.705, -0.697]; GPT-4o: $d = -2.238$), with 97% of iterations revealing qualitatively
268 similar patterns despite magnitude differences ($\text{BF}_{10} < 0.1$). Regarding the *working memory*
269 *updating* task, while 37.8% of human subsamples showed time-dependent effects (mean $d =$
270 -0.226 , 95% CI [-0.230, -0.223]), GPT-4o exhibited no such effect ($d = -0.027$, $p = 0.833$), with
271 98.5% of iterations supporting the null hypothesis (median $\text{BF}_{10} = 3.07$). In the *task-switching*
272 paradigm, all human subsamples showed robust switch costs (mean $d = 0.662$, 95% CI [0.659,
273 0.664]). GPT-4o exhibited smaller but significant costs ($d = 0.347$; $\text{BF}_{10} > 100$ across all
274 iterations). These analyses confirmed the stability of our findings and demonstrated the selectivity
275 in GPT-4o's replication of human EF patterns.

276

277 **Temperature parameter analysis:** The temperature parameter in LLMs controls response
278 variability, with lower values producing more deterministic outputs and higher values increasing
279 stochasticity—potentially paralleling individual differences in human cognitive processing (30). To
280 determine whether our findings reflected stable patterns rather than artifacts of specific parameter
281 settings, we repeated all analyses at temperature = 1. The temperature manipulation did not alter
282 the main effect patterns across EF tasks, consistent with recent evidence showing that
283 temperature parameters exert only a minimal impact on LLM performance in standardized
284 psychological assessments (20, 31, 32). Specifically, in the *Stroop*, *Go/No-Go*, *running memory*,
285 and *task-switching* paradigms, the direction of effects, statistical significance, and effect sizes
286 remained consistent between temperature = 0 and temperature = 1 conditions (see SI Appendix,
287 Fig. S2–4).

288

289

290 **Discussion**

291

292 In this study, we systematically evaluated the performance of GPT-4o when instructed to simulate
293 human-like responses across behavioral tasks measuring three core dimensions of EF, while
294 simultaneously exploring the potential of model log probability parameters as indicators of
295 cognitive processing. Our results revealed that GPT-4o successfully replicated typical human
296 cognitive patterns in tasks assessing interference inhibition, working memory maintenance, and
297 updating. However, the model exhibited patterns that diverged from human performance in

298 response inhibition, time-constrained working memory updating, and cognitive flexibility tasks.
299 This selective replication pattern delineates specific boundaries in current LLM architectures for
300 simulating human executive control mechanisms (SI Appendix, Fig. S5). Notably, log probability
301 parameters showed promise as cognitive processing assessment tools in certain task contexts,
302 thereby offering new empirical approaches for understanding the computational foundations of
303 EFs.

304 Our findings demonstrate that GPT-4o, when prompted to simulate human-like responses,
305 exhibited patterns consistent with human EF in tasks involving information representation
306 conflicts and the maintenance of static information. In the *Stroop* task, the model displayed
307 significantly prolonged reaction times in incongruent conditions, with effect sizes marginally
308 exceeding those typical of human performance. This difference may reflect the fact that LLMs are
309 computational systems and are thus unaffected by fatigue or attentional fluctuations. These
310 results align with the observation by Cui et al. (2025) that LLMs reliably replicate main effects
311 while struggling with those related to interaction (20). In working memory tasks, GPT-4o
312 demonstrated structural constraints paralleling those of humans. Specifically, the model
313 replicated the backward difficulty effect in *digit span* tasks and exhibited performance degradation
314 with increasing load in *n*-back tasks, suggesting that even with adequate task comprehension, the
315 model displays structural working memory limitations (33). Meanwhile, GPT-4o failed to
316 demonstrate typical human patterns in tasks demanding dynamic control and temporal sensitivity.
317 The model achieved perfect accuracy in response inhibition tasks (*Go/No-Go* and *stop-signal*),
318 showed no effects of time pressure in working memory updating tasks, and exhibited minimal
319 switch costs in cognitive flexibility tasks. These findings converge with recent Wisconsin Card
320 Sorting Test (WCST) studies, where Langis et al. (2025) and Goto et al. (2025) documented
321 similar failures in cognitive flexibility. The latter study reported that most LLMs achieved less than
322 30% accuracy with systematic rule-response mismatches (17, 34), suggesting that they have
323 fundamental limitations in dynamic cognitive control across multiple paradigms. The “zero
324 variance” phenomenon reported by Cui et al. (2025) further suggests that the deterministic nature
325 of LLM response generation may be incompatible with the inherent variability of executive control
326 mechanisms. Notably, this selective replication pattern provides new computational insights into
327 EF models. Tasks that ostensibly measure the same EF dimension may, in fact, engage distinct
328 computational processes. While both the *Stroop* and *Go/No-Go* tasks purportedly assess
329 inhibitory control (12), the former can be replicated by LLMs while the latter cannot, indicative of
330 computational-level differences in their underlying inhibitory mechanisms. LLMs appear capable
331 of replicating processes based on symbolic representation and static maintenance but exhibit
332 divergent characteristics in processes requiring dynamic control, temporal sensitivity, or intrinsic
333 variability. This computational heterogeneity aligns with Xu et al.’s (2025) finding that LLM-human
334 representational similarity decreases from non-sensorimotor to sensory domains (35), and is
335 further implied by Goto et al.’s (2025) demonstration that models struggle to integrate bottom-up
336 card information with top-down rule-based reasoning (34)—a coordination failure that may
337 fundamentally constrain dynamic cognitive control. Such computational heterogeneity may help
338 explain the variations in EF factor structures reported across different studies (36).

339 Extending these behavioral findings, we introduced log probability parameters as a novel
340 approach to cognitive assessment in GPT-4o, revealing task-specific associations with the
341 model’s behavioral performance. In the *Stroop* task, reaction times correlated positively with
342 perplexity and negatively with average log probability. The associations between internal
343 uncertainty and behavioral delays became more pronounced during semantic conflict (*Stroop*
344 incongruent condition) and with increased cognitive load (*backward digit span*). Log probabilities
345 may capture “internal certainty” rather than “actual accuracy”—elevated perplexity might indicate
346 a state of “surface fluency but internal confusion” (26). Critically, log probability parameters
347 revealed differences in cognitive processing even in tasks where behavioral performance reached
348 ceiling levels. For example, in the *Go/No-Go* task, despite GPT-4o achieving perfect accuracy,
349 No-Go trials exhibited significantly higher perplexity than Go trials, suggesting that the model
350 experienced inhibitory conflict processing at the internal representation level. Similarly, in task-
351 switching paradigms, internal state parameters showed significant between-condition differences,

352 even though behavioral switch costs were minimal. This dissociation between behavioral and
353 internal states—where statistical uncertainty fails to manifest as overt behavioral costs—mirrors
354 the rule-response mismatches observed by Goto et al. (2025) in the *WCST* (34). Collectively,
355 these findings suggest that a fundamental disconnect exists between internal representation and
356 behavioral execution within the information processing hierarchies in LLMs.

357
358 To further establish the biological plausibility of log probability parameters as cognitive indicators,
359 we compared them to human neural activity (37). In the *n*-back task, activation in the left gyrus
360 rectus and right medial superior frontal gyrus correlated significantly with model perplexity. This
361 correspondence between these key working memory network nodes and the model's internal
362 states may indicate functional equivalence. Additionally, in *task-switching* paradigms, left
363 thalamus activation correlated with log probability parameters, suggesting that these parameters
364 might capture comparable control processes (38). However, no significant brain-model parameter
365 correlations appeared in the stop-signal task, which contrasts with the correlations observed in
366 the working memory and *task-switching* tasks. This selective pattern suggests that when the
367 model's computational strategies fundamentally diverge from human cognitive processes, internal
368 parameters may not map onto human neural activity patterns—a finding consistent with recent
369 research showing that the presence or absence of model-brain correspondences can reveal
370 computational commonalities and differences between human brains and artificial intelligence
371 (39).

372
373 Collectively, log probability parameters offer distinct advantages as internal indicators of cognitive
374 processing. Unlike hidden layer activations in deep neural networks (40), these parameters are
375 clearly statistically interpretable, while remaining accessible through API interfaces, substantially
376 reducing technical barriers for cognitive research using LLMs. Moreover, they provide diagnostic
377 tools for identifying potential “surface fluency but internal confusion” states, where models
378 generate seemingly coherent responses despite underlying computational uncertainty.

379 **Implications**

381 This study advances the understanding of the computational nature of EFs by employing LLMs as
382 computational tools. The selective replication pattern observed suggests that there may be
383 potential computational heterogeneity within EF components (41), providing empirical evidence
384 for understanding the relationships between different tasks in Miyake's three-factor model.
385 Whereas traditional research examines covariance relationships between tasks using factor
386 analysis, this study explored potential task differences from the perspective of computational
387 mechanisms, comparing which tasks LLMs can replicate and which they cannot. The value of this
388 approach lies in its potential to move beyond correlational analysis and provide insights into
389 causal computational mechanisms (5, 8). Selective replication by LLMs helps identify which EF
390 components might be based on domain-general computational principles and which might require
391 specific neurobiological foundations (42). These findings are not directly obtainable from
392 traditional human participant studies, as human participants simultaneously employ multiple
393 cognitive mechanisms, making it difficult to isolate their individual contributions.

394
395 Furthermore, in this study, we established a novel methodological framework for assessing
396 cognitive abilities in LLMs. While traditional assessments primarily focus on task accuracy, we
397 introduced three levels of analysis: behavioral performance, measured through accuracy and
398 reaction times; internal-state level log probability parameters; and brain-model correspondence at
399 the neural correlation level. This multilevel analytical framework enables a more comprehensive
400 characterization of the models' cognitive features, helping avoid incomplete judgments based
401 solely on behavioral performance. For instance, in the *Go/No-Go* task, examining only behavioral
402 performance might lead to the conclusion that “GPT-4o possesses perfect response inhibition
403 ability.” However, internal parameter analysis indicated that the model still experienced inhibitory
404 conflict, while brain-model correlation analysis showed that this internal state lacks
405 correspondence with human response inhibition networks. These three levels of evidence

406 collectively suggest that the model's "perfect performance" may not stem from response inhibition
407 ability in the human sense, but from different computational mechanisms (35). Building on the
408 work of Bauer et al. (2025), who proposed using log probability parameters to assess LLM
409 uncertainty, in this study, we extended this method to cognitive psychology, exploring its potential
410 applications in EF assessment (27). Future research could further investigate relationships
411 between other internal parameters (e.g., attention weights, hidden layer activations) (43) and
412 cognitive processing, thereby establishing a more comprehensive cognitive assessment toolkit.
413

414 **Limitations and future considerations**

415 Despite the meaningful findings, this study has several limitations that warrant consideration.
416 First, we focused exclusively on GPT-4o as a model, similar to previous studies (44, 45). Different
417 model architectures, scales, and training paradigms may yield distinct EF performance patterns
418 (46). For instance, models incorporating more human feedback reinforcement learning might
419 perform better at following "inhibition instructions," although whether this would constitute genuine
420 "response inhibition ability" is unclear (47). Because the observed selective replication pattern
421 may partially reflect specific characteristics of GPT series models, future research should assess
422 multiple mainstream models to establish benchmark datasets for LLM EFs. Second, while this
423 study covered three major EF dimensions based on the Miyake framework, the range of tasks
424 within each component was limited. Comprehensive EF assessment necessitates a broader task
425 repertoire, including Simon tasks, complex span tasks, the WCST, and others (48). Future
426 research should expand the task range to achieve a more comprehensive characterization of
427 LLM EF features. Third, in this study, while we undertook a preliminary exploration of log
428 probability parameters as cognitive indicators, many unresolved questions remain. Regarding
429 parameter selection and combination, Application Programming Interfaces (APIs) offer additional
430 information such as top-*k* token probability distributions, and determining which parameter
431 combinations best predict cognitive performance requires further investigation (49). Moreover,
432 whether the psychological significance of log probability parameters remains consistent across
433 different models also needs verification. Furthermore, this study primarily relied on correlational
434 analysis; future work could attempt to manipulate log probability parameters and observe
435 corresponding behavioral changes to explore potential causal relationships. Finally, the
436 theoretical mechanisms underlying selective replication require deeper examination. Why can
437 certain EF components be acquired through text learning while others cannot? Does this
438 difference reflect evolutionary hierarchies of cognitive functions? Addressing these questions will
439 enhance our understanding of the uniqueness of human cognition and inform future directions for
440 artificial intelligence development.
441

442 **Summary**

443 In this study, we systematically evaluated GPT-4o's performance on EF tasks using standard
444 cognitive paradigms employed in human studies, and found that the model selectively replicated
445 human cognitive patterns. The model demonstrated performance comparable to humans in tasks
446 based on symbolic representation and maintenance of static information, but exhibited divergent
447 performance in tasks involving dynamic control and temporal sensitivity, as well as those
448 requiring intrinsic variability. Log probability parameters, serving as quantitative indicators of
449 model internal states, showed systematic associations with behavioral performance and human
450 brain activation patterns in certain tasks, providing new analytical tools for assessing cognitive
451 processing. This selective replication pattern not only clarifies the cognitive boundaries of current
452 LLMs but also provides valuable empirical methods for exploring the computational foundations of
453 EFs.
454

456 **Materials and Methods**

457 **Participants** 458

459 Two independent datasets were analyzed to investigate EF performance (SI Appendix, Table
460 S2). Dataset 1 consisted of behavioral data collected for this study from 1,970 native Mandarin-
461 speaking adults (915 males [46.45%], 1,055 females [53.55%]; mean age = 19.78 years, SD =
462 1.72) who completed five standardized EF tasks. Participants were recruited from Northwest
463 Normal University and surrounding communities. The inclusion criteria comprised native
464 Mandarin proficiency, normal or corrected-to-normal vision, and normal color vision. Individuals
465 with neurological or psychiatric disorders or those taking medications affecting cognitive function
466 were excluded. Dataset 2 consisted of fMRI data from He et al. (2021), involving 39 adults (13
467 males, 26 females; mean age = 19.21 ± 0.52 years) who performed three EF tasks during
468 scanning (50). All human participants provided written informed consent, and the study was
469 approved by the Ethics Committee of Northwest Normal University (Protocol #NUNU-20250105).

470

471 **Testing human EF performance**

472 Participants from Dataset 1 completed five EF tasks using E-Prime software (SI Appendix, Fig.
473 S1; for details, see SI Appendix, Methods). The task set evaluated core EF components, and
474 consisted of the following: (1) The *Stroop color-word* task, which assessed interference control by
475 requiring participants to identify ink colors while ignoring word meanings across congruent,
476 incongruent, and neutral conditions (18 practice trials, 108 test trials); (2) the *Go/No-Go* task,
477 which measured response inhibition, with participants responding to target letters while
478 withholding responses to non-targets in a counterbalanced design (4 blocks × 100 trials); (3) the
479 *number-switching* task, which evaluated cognitive flexibility through alternating magnitude
480 judgments for red digits and parity judgments for blue digits across single-task blocks (10 blocks
481 × 8 trials) and mixed-task blocks (10 blocks × 17 trials); (4) the *running memory* task, which
482 assessed working memory updating by requiring continuous updating of the last three digits
483 under simple (1750 ms presentation) and difficult (750 ms) conditions across sequences of
484 varying lengths (5–11 items); and (5) the *digit span* task, which measured working memory
485 maintenance through forward and backward recall using adaptive procedures beginning with
486 three-digit sequences.

487

488 Participants from Dataset 2 performed three tasks during fMRI scanning: (1) The *n*-back task,
489 which assessed working memory across 0-back, 1-back, and 2-back conditions using blocked
490 designs (4 blocks per condition, 29 s each with 15 stimuli); (2) the *stop-signal* task, which
491 evaluated response inhibition with 96 Go trials and 32 stop trials, employing dynamically adjusted
492 stop-signal delays; and (3) the *number-letter switching* task, which measured cognitive flexibility
493 with 48 repeat and 48 switch trials, requiring alternation between parity and consonant/vowel
494 judgments based on color cues.

495

496 **Testing LLM EF performance and log parameters**

497 GPT-4o (version 2025-03-26) was evaluated through OpenAI's official API with the temperature
498 set to 0 to ensure consistency with prior LLM cognitive research (19, 20, 31, 44, 45), and thereby
499 facilitate direct comparisons of results across studies. To verify the robustness of our findings, all
500 analyses were repeated at temperature = 1 as a sensitivity check (20, 31). Default settings were
501 used for all other parameters. The model completed EF tasks through prompt-driven multi-turn
502 conversations, with each simulated participant maintaining consistent context within a single
503 dialog session to ensure ecological validity and independence between participants. The task
504 prompts and procedures closely matched those used with human participants (35), and were
505 delivered to GPT-4o through natural language inputs (details in SI Appendix, Tables S3 and S4).
506 To quantify the model's internal computational processes (SI Appendix, Table S5), four log
507 probability parameters were extracted during response generation (see
508 https://cookbook.openai.com/examples/using_logprobs):

509

510 (1) Average log probability quantified overall model confidence:

$$\text{avg_logprob} = \frac{1}{n} \sum_{i=1}^n \log p(t_i)$$

511 (2) Perplexity measured model uncertainty:
512

$$\text{PPL} = \exp(-\text{avg_logprob})$$

513 (3) Log probability variance indicated confidence variability:
514
515

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (\log p(t_i) - \text{avg_logprob})^2$$

516 (4) Minimum log probability identified weakest predictions:
517
518

$$\text{min_logprob} = \min_{i=1, \dots, n} \log p(t_i)$$

519 All computations were implemented using *NumPy* (51).
520
521

522 MRI data acquisition and analysis

523 fMRI data were acquired using a 3T Siemens scanner with echo-planar imaging sequences (TR =
524 2000 ms, TE = 30 ms, voxel size = $3.4 \times 3.4 \times 3$ mm³). High-resolution T1-weighted structural
525 images were obtained using Magnetization Prepared Rapid Gradient Echo (MPRAGE)
526 sequences. Preprocessing followed standard procedures using Statistical Parametric Mapping 12
527 (SPM12) and Data Processing & Analysis for Brain Imaging (DPABI) toolboxes (52). These steps
528 included slice-timing correction, realignment, coregistration, Diffeomorphic Anatomical
529 Registration Through Exponentiated Lie Algebra (DARTEL) segmentation, normalization to
530 Montreal Neurological Institute (MNI) space, and spatial smoothing with an 8-mm Full Width at
531 Half Maximum (FWHM) Gaussian kernel. Task-based activation was modeled using general
532 linear models with the following contrasts: “2-back > 0-back” for working memory load,
533 “successful stop > successful Go” for inhibitory control, and “switch > repeat” for cognitive
534 flexibility costs. Mean activation values were extracted from 116 anatomically-defined regions of
535 interest using the Automated Anatomical Labeling (AAL) atlas (53) (see SI Appendix, Methods).
536

537 Sample size justification and considerations

538 Sample size determination was predicated on fundamental differences between human and LLM
539 response distributions. While human samples inherently reflect genuine individual differences in
540 EF, LLMs, when using fixed temperature parameters, generate statistical variation around a
541 model-determined central tendency (20). Given the more homogeneous response patterns
542 displayed by LLMs, excessive sampling introduces the risk of pseudo-precision and violates the
543 value of information principle (54, 55). A total of 120 simulated participants were generated for
544 Dataset 1, providing >90% power to detect medium effects ($d = 0.5$) and >99% power for large
545 effects ($d = 0.8$) at $\alpha = 0.05$, consistent with typical EF effect sizes ($d = 0.4$ – 0.8) (36). To ensure
546 fair comparisons and mitigate statistical biases from unequal sample sizes, 120 participants from
547 the original 1,970 human participants were randomly subsampled for matched analyses. For
548 Dataset 2, 39 simulated participants were generated to match the fMRI sample size, enabling a
549 direct comparison of internal computational parameters with neural activation patterns.
550

551 Data analysis

552 Bayesian ANOVA was used to compare EF performance (SI Appendix, Table S6) between AI
553 and humans (56), with evidence strength quantified through BFs (57). Analyses were
554 implemented in Python using Python Monte Carlo (PyMC3) with weakly informative priors (58).
555 Posterior distributions for effect size differences were computed, and the resulting posterior
556 means and medians were reported as point estimates along with 95% HDIs to characterize

557 uncertainty. BF served to quantify support for the hypothesis that AI effects exceed human
558 effects:

$$BF_{10} = \frac{P(\Delta_{\text{effect}} > 0 \mid \text{data})}{P(\Delta_{\text{effect}} \leq 0 \mid \text{data})}$$

559

560 Following established conventions (59), $BF_{10} > 10$ indicated strong evidence, $3 < BF_{10} < 10$
561 indicated moderate evidence, $1 < BF_{10} < 3$ indicated weak evidence, and $BF_{10} < 0.33$ supported the
562 null hypothesis.

563 To investigate the relationships between internal computational mechanisms and behavioral
564 performance, Pearson correlation analyses were performed between task performance metrics
565 (reaction times, accuracy) and log probability parameters. All correlation coefficients were
566 corrected for multiple comparisons using the Benjamini-Hochberg False Discovery Rate (FDR)
567 procedure at $\alpha = 0.05$ (60). For brain-model correspondence analyses, we identified significantly
568 activated ROIs through one-sample t -tests ($p_{\text{Bonferroni}} < 0.05$), then performed Pearson correlations
569 between ROI activation values and model log probability parameters. These analyses explored
570 statistical associations between neural activity and model output uncertainty metrics, though such
571 correlations should be interpreted as statistical relationships rather than evidence of shared
572 computational mechanisms. All correlation coefficients were also corrected for multiple
573 comparisons using FDR procedure at $\alpha = 0.05$.

574 **Data, materials, and software availability**

575 The transcripts and the numerical coding of the transcripts have been deposited in the Open
576 Science Framework (OSF) (<https://osf.io/hfnt6>).

577 **Acknowledgments**

578 We thank Peng Lu, Shuyan Gao, and Qiya Huang from the School of Psychology, Northwest
579 Normal University, for their assistance with data collection; Yajie Gong and Jia Li from the School
580 of Psychology, Northwest Normal University, for organizing the experimental materials; and all
581 participants for their participation in this study. We also thank Nai Ding from the College of
582 Biomedical Engineering and Instrument Sciences, Zhejiang University, for valuable insights on
583 the research approach during the early stages of this study.

584

585

586

References

- 587 1. T. Hagendorff, S. Fabi, M. Kosinski, Human-like intuitive behavior and reasoning biases
588 emerged in large language models but disappeared in ChatGPT. *Nat Comput Sci* **3**, 833–838
589 (2023).
- 590 2. K. Valmeekam, A. Olmo, S. Sreedharan, S. Kambhampati, Large Language Models Still Can't
591 Plan (A Benchmark for LLMs on Planning and Reasoning about Change) in (2022).
- 592 3. T. W. Webb, S. S. Mondal, C. Wang, B. Krabach, I. Momennejad, A Prefrontal Cortex-inspired
593 Architecture for Planning in Large Language Models. *CoRR* (2023).
- 594 4. A. Chemero, LLMs differ from human cognition because they are not embodied. *Nat Hum*
595 *Behav* **7**, 1828–1829 (2023).
- 596 5. M. C. Frank, Openly accessible LLMs can help us to understand human cognition. *Nat Hum*
597 *Behav* **7**, 1825–1827 (2023).
- 598 6. L. Ke, S. Tong, P. Cheng, K. Peng, Exploring the frontiers of LLMs in psychological
599 applications: a comprehensive review. *Artif Intell Rev* **58**, 305 (2025).
- 600 7. J. Meng, AI emerges as the frontier in behavioral science. *Proceedings of the National*
601 *Academy of Sciences* **121**, e2401336121 (2024).
- 602 8. M. Hardy, I. Sucholutsky, B. Thompson, T. Griffiths, Large language models meet cognitive
603 science: LLMs as tools, models, and participants. *Proceedings of the Annual Meeting of the*
604 *Cognitive Science Society* **45** (2023).
- 605 9. K. Mahowald, *et al.*, Dissociating language and thought in large language models. *Trends in*
606 *Cognitive Sciences* **28**, 517–540 (2024).

- 607 10. I. Yildirim, L. A. Paul, From task structures to world models: what do LLMs know? *Trends in*
608 *Cognitive Sciences* **28**, 404–415 (2024).
- 609 11. A. Diamond, Executive Functions. *Annu. Rev. Psychol.* **64**, 135–168 (2013).
- 610 12. N. P. Friedman, T. W. Robbins, The role of prefrontal cortex in cognitive control and executive
611 function. *Neuropsychopharmacol.* **47**, 72–89 (2022).
- 612 13. M. Carlén, What constitutes the prefrontal cortex? *Science* **358**, 478–482 (2017).
- 613 14. D. Gong, X. Wan, D. Wang, Working Memory Capacity of ChatGPT: An Empirical Study.
614 *AAAI* **38**, 10048–10056 (2024).
- 615 15. C. Zhang, Y. Jian, Z. Ouyang, S. Vosoughi, Working Memory Identifies Reasoning Limits in
616 Language Models in *Proceedings of the 2024 Conference on Empirical Methods in Natural*
617 *Language Processing*, Y. Al-Onaizan, M. Bansal, Y.-N. Chen, Eds. (Association for
618 Computational Linguistics, 2024), pp. 16896–16922.
- 619 16. X. Hu, R. L. Lewis, Do Language Models Understand the Cognitive Tasks Given to Them?
620 Investigations with the N-Back Paradigm. [Preprint] (2024). Available at:
621 <http://arxiv.org/abs/2412.18120> [Accessed 29 December 2024].
- 622 17. K. de Langis, *et al.*, A Framework for Robust Cognitive Evaluation of LLMs. [Preprint] (2025).
623 Available at: <http://arxiv.org/abs/2504.02789> [Accessed 26 April 2025].
- 624 18. R. Loconte, G. Orrù, M. Tribastone, P. Pietrini, G. Sartori, Challenging ChatGPT ' *Intelligence*'
625 with Human Tools: A Neuropsychological Investigation on Prefrontal Functioning of a Large
626 Language Model. [Preprint] (2023). Available at: <https://papers.ssrn.com/abstract=4377371>
627 [Accessed 8 July 2025].
- 628 19. J. Coda-Forno, M. Binz, J. X. Wang, E. Schulz, CogBench: a large language model walks into
629 a psychology lab. [Preprint] (2024). Available at: <http://arxiv.org/abs/2402.18225> [Accessed 2
630 December 2024].
- 631 20. Z. Cui, N. Li, H. Zhou, A large-scale replication of scenario-based experiments in psychology
632 and management using large language models. *Nat Comput Sci* **5**, 627–634 (2025).
- 633 21. A. Miyake, N. P. Friedman, The nature and organization of individual differences in executive
634 functions: Four general conclusions. *Current directions in psychological science* **21**, 8–14 (2012).
- 635 22. M. D'Esposito, B. R. Postle, The cognitive neuroscience of working memory. *Annual review of*
636 *psychology* **66**, 115 (2015).
- 637 23. M. M. Botvinick, T. S. Braver, D. M. Barch, C. S. Carter, J. D. Cohen, Conflict monitoring and
638 cognitive control. *Psychological review* **108**, 624 (2001).
- 639 24. A. R. Aron, T. W. Robbins, R. A. Poldrack, Inhibition and the right inferior frontal cortex: one
640 decade on. *Trends in Cognitive Sciences* **18**, 177–185 (2014).
- 641 25. M. Binz, *et al.*, A foundation model to predict and capture human cognition. *Nature* 1–8
642 (2025). <https://doi.org/10.1038/s41586-025-09215-4>.
- 643 26. H. Lin, Y. Zhang, Navigating the Risks of Using Large Language Models for Text Annotation
644 in Social Science Research. *Social Science Computer Review* 08944393251366243 (2025).
645 <https://doi.org/10.1177/08944393251366243>.
- 646 27. C. Bauer, L. T. Duc Dang, T. van den Beucken, J. Schuchhardt, R. Herwig, Systematic
647 analysis of hepatotoxicity: combining literature mining and AI language models. *Front. Artif. Intell.*
648 **8** (2025).
- 649 28. S. Chen, M. v. Maddali, C. Bluethgen, C. p. Langlotz, R. Raj, Leveraging Generative Pre-
650 trained Transformer (GPT) Large Language Models (LLMs) For Interstitial Lung Diseases (ILD)
651 Clinical Research. *Am J Respir Crit Care Med* **211**, A2086–A2086 (2025).
- 652 29. H. Gonen, S. Iyer, T. Blevins, N. A. Smith, L. Zettlemoyer, Demystifying prompts in language
653 models via perplexity estimation. *arXiv preprint arXiv:2212.04037* (2022).
- 654 30. M. Peeperkorn, T. Kouwenhoven, D. Brown, A. Jordanous, Is Temperature the Creativity
655 Parameter of Large Language Models? [Preprint] (2024). Available at:
656 <http://arxiv.org/abs/2405.00492> [Accessed 30 October 2025].
- 657 31. M. Kosinski, Evaluating large language models in theory of mind tasks. *Proceedings of the*
658 *National Academy of Sciences* **121**, e2405460121 (2024).
- 659 32. C. Li, Y. Qi, Toward accurate psychological simulations: Investigating LLMs' responses to
660 personality and cultural variables. *Computers in Human Behavior* **170**, 108687 (2025).

661 33. J. Huang, K. Sun, W. Wang, M. Dredze, Language Models Do Not Have Human-Like Working
662 Memory. [Preprint] (2025). Available at: <http://arxiv.org/abs/2505.10571> [Accessed 31 October
663 2025].

664 34. D. Goto, H. Idei, Y. Shiozuka, T. Ogata, Performance of Large Language Models and
665 Analysis of Responses in the Wisconsin Card Sorting Task in *2025 IEEE International
666 Conference on Development and Learning (ICDL)*, (2025), pp. 1–7.

667 35. Q. Xu, *et al.*, Large language models without grounding recover non-sensorimotor but not
668 sensorimotor features of human concepts. *Nat Hum Behav* **9**, 1871–1886 (2025).

669 36. J. E. Karr, *et al.*, The unity and diversity of executive functions: A systematic review and re-
670 analysis of latent variable studies. *Psychological Bulletin* **144**, 1147–1185 (2018).

671 37. C. Gao, *et al.*, Increasing alignment of large language models with language processing in the
672 human brain. *Nature computational science* 1–11 (2025).

673 38. M. M. Halassa, S. Kastner, Thalamic functions in distributed cognitive control. *Nat Neurosci*
674 **20**, 1669–1679 (2017).

675 39. A. Doerig, *et al.*, High-level visual representations in the human brain are aligned with large
676 language models. *Nature Machine Intelligence* 1–15 (2025).

677 40. A. Azaria, T. Mitchell, The Internal State of an LLM Knows When It's Lying. [Preprint] (2023).
678 Available at: <http://arxiv.org/abs/2304.13734> [Accessed 31 October 2025].

679 41. D. Haywood, F. D. Baughman, Multidimensionality in Executive Function Profiles in
680 Schizophrenia: a Computational Approach Using the Wisconsin Card Sorting Task. *Comput Brain
681 Behav* **4**, 381–394 (2021).

682 42. B. Alkhamissi, *et al.*, From Language to Cognition: How LLMs Outgrow the Human Language
683 Network. [Preprint] (2025). Available at: <http://arxiv.org/abs/2503.01830> [Accessed 31 October
684 2025].

685 43. A. Vaswani, *et al.*, Attention is all you need. *Advances in neural information processing
686 systems* **30** (2017).

687 44. G. Suri, L. R. Slater, A. Ziaee, M. Nguyen, Do large language models show decision
688 heuristics similar to humans? A case study using GPT-3.5. *Journal of Experimental Psychology:
689 General* **153**, 1066–1075 (2024).

690 45. S. A. Lehr, K. S. Saichandran, E. Harmon-Jones, N. Vitali, M. R. Banaji, Kernels of selfhood:
691 GPT-4o shows humanlike patterns of cognitive dissonance moderated by free choice. *Proc. Natl.
692 Acad. Sci. U.S.A.* **122**, e2501823122 (2025).

693 46. X. Wang, *et al.*, CogLM: Tracking Cognitive Development of Large Language Models.
694 [Preprint] (2025). Available at: <http://arxiv.org/abs/2408.09150> [Accessed 31 October 2025].

695 47. S. Chaudhari, *et al.*, RLHF Deciphered: A Critical Analysis of Reinforcement Learning from
696 Human Feedback for LLMs. *ACM Comput. Surv.* **58**, 53:1-53:37 (2025).

697 48. M. Yangüez, B. Bediou, J. Chanal, D. Bavelier, In search of better practice in executive
698 functions assessment: Methodological issues and potential solutions. *Psychological Review* **131**,
699 402–430 (2024).

700 49. L. Shan, S. Luo, Z. Zhu, Y. Yuan, Y. Wu, Cognitive Memory in Large Language Models.
701 [Preprint] (2025). Available at: <http://arxiv.org/abs/2504.02441> [Accessed 30 October 2025].

702 50. L. He, *et al.*, Unity and diversity of neural representation in executive functions. *Journal of
703 Experimental Psychology: General* **150**, 2193–2207 (2021).

704 51. C. R. Harris, *et al.*, Array programming with NumPy. *Nature* **585**, 357–362 (2020).

705 52. C.-G. Yan, X.-D. Wang, X.-N. Zuo, Y.-F. Zang, DPABI: data processing & analysis for
706 (resting-state) brain imaging. *Neuroinformatics* **14**, 339–351 (2016).

707 53. E. T. Rolls, C.-C. Huang, C.-P. Lin, J. Feng, M. Joliot, Automated anatomical labelling atlas 3.
708 *Neuroimage* **206**, 116189 (2020).

709 54. D. Lakens, Sample Size Justification. *Collabra: Psychology* **8**, 33267 (2022).

710 55. F. Anvari, D. Lakens, Using anchor-based methods to determine the smallest effect size of
711 interest. *Journal of Experimental Social Psychology* **96**, 104159 (2021).

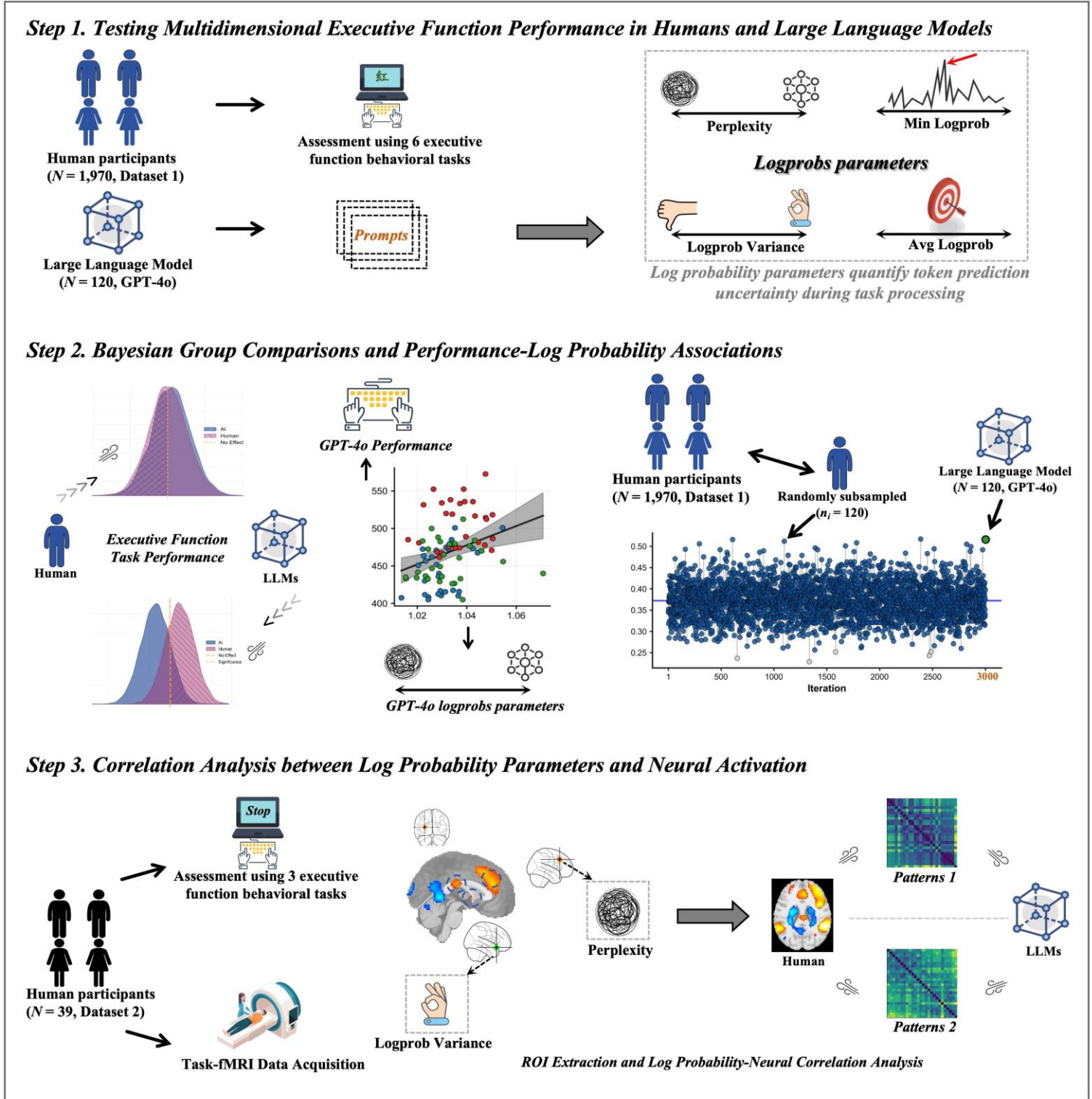
712 56. E. J. Miller, *et al.*, AI hyperrealism: Why AI faces are perceived as more real than human
713 ones. *Psychological Science* **34**, 1390–1403 (2023).

- 714 57. J. K. Kruschke, T. M. Liddell, The Bayesian New Statistics: Hypothesis testing, estimation,
715 meta-analysis, and power analysis from a Bayesian perspective. *Psychon Bull Rev* **25**, 178–206
716 (2018).
- 717 58. A. Vehtari, A. Gelman, J. Gabry, Practical Bayesian model evaluation using leave-one-out
718 cross-validation and WAIC. *Stat Comput* **27**, 1413–1432 (2017).
- 719 59. R. E. Kass, A. E. Raftery, Bayes factors. *Journal of the american statistical association* **90**,
720 773–795 (1995).
- 721 60. Y. Benjamini, D. Drai, G. Elmer, N. Kafkafi, I. Golani, Controlling the false discovery rate in
722 behavior genetics research. *Behavioural Brain Research* **125**, 279–284 (2001).

723

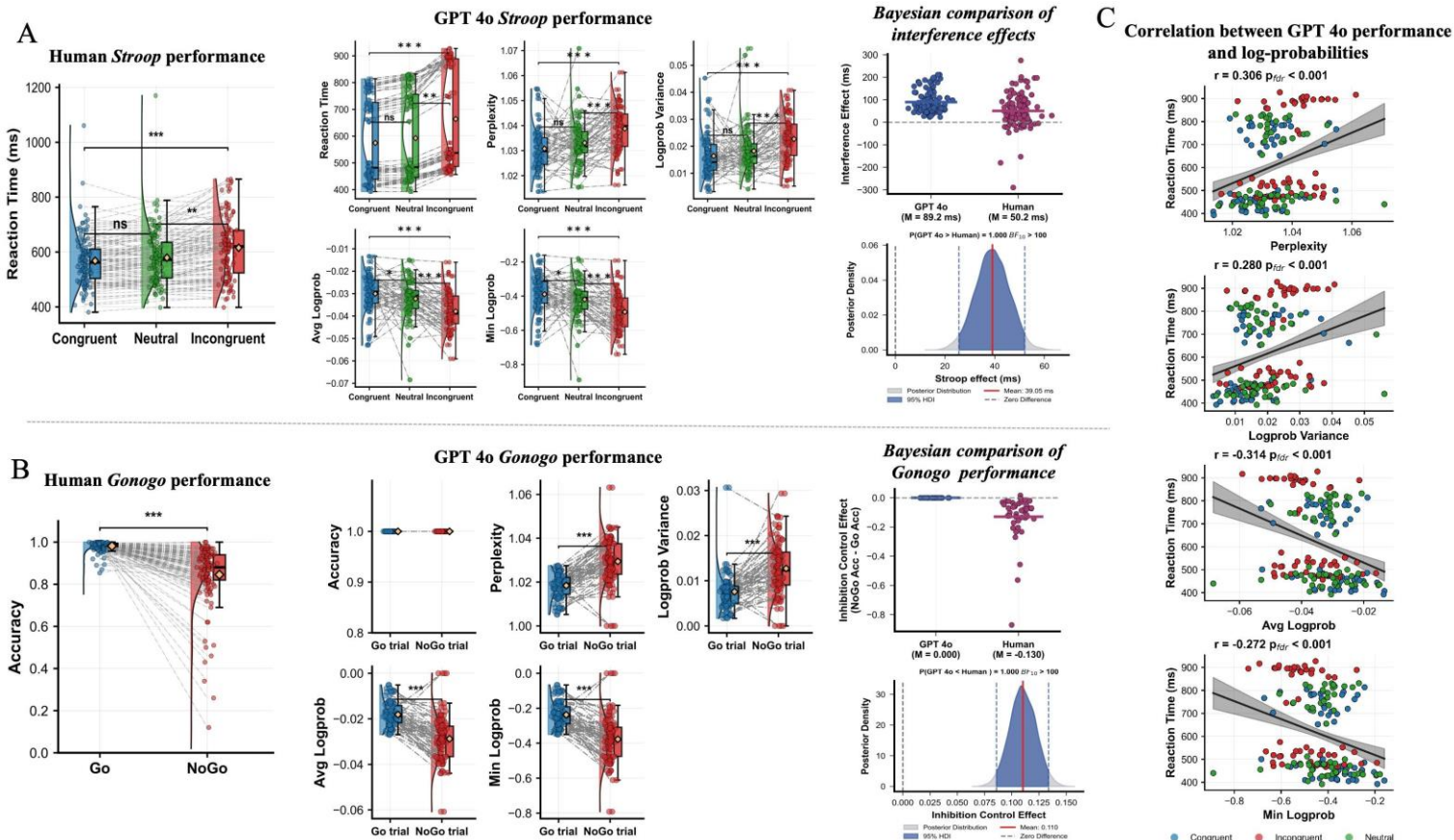
724 **Figures**

725



727 **Figure 1.** Study design and analytical framework. This study comprised three main steps. Step 1:
 728 Data collection. We collected behavioral data from 1,970 human participants who completed five
 729 standardized executive function (EF) tasks (Dataset 1). Concurrently, we tested GPT-4o (N =

730 120) on the same tasks through prompt-driven multi-turn dialogs, extracting log probability
731 parameters (perplexity, log probability variance, average log probability, and minimum log
732 probability) as indicators of the model's internal computational states. Step 2: Behavioral analysis.
733 We compared EF task performance between GPT-4o and humans using Bayesian hierarchical
734 ANOVA to examine canonical cognitive effects (e.g., Stroop interference, task-switching costs).
735 To ensure statistical robustness, we randomly subsampled 120 human participants to match the
736 GPT-4o sample size, repeating this procedure 3,000 times. Effect size differences were
737 quantified using Bayes factors. We additionally examined correlations between model log
738 probability parameters and behavioral performance metrics (FDR-corrected). Step 3: Neural
739 correspondence analysis. We analyzed fMRI data from 39 participants (Dataset 2) who
740 performed three EF tasks (*n*-back, *stop-signal*, and *number-letter switching*). Task-related
741 activation was modeled using general linear models with predefined contrasts. Mean activation
742 values were extracted from 116 brain regions defined by the Automated Anatomical Labeling
743 (AAL) atlas, and correlations with GPT-4o log probability parameters were assessed for regions
744 showing significant task-related activation.



745

746

Figure 2. Behavioral performance comparison between GPT-4o and human participants in inhibitory control tasks. (A) Reaction times across three Stroop task conditions for human participants ($N = 120$) and the corresponding reaction times with log probability parameters for GPT-4o. The rightmost panel shows Bayesian analysis comparing interference effects (incongruent minus congruent condition) between humans and GPT-4o. (B) Accuracy in Go versus No-Go trials for human participants ($N = 120$) and the corresponding accuracy with log probability parameters for GPT-4o in the Go/No-Go task. The rightmost panel presents Bayesian analysis comparing inhibitory control effects (No-Go accuracy minus Go accuracy) between humans and GPT-4o. (C) Correlation analysis between reaction times and log probability parameters across the three Stroop task conditions for GPT-4o.

748

749

750

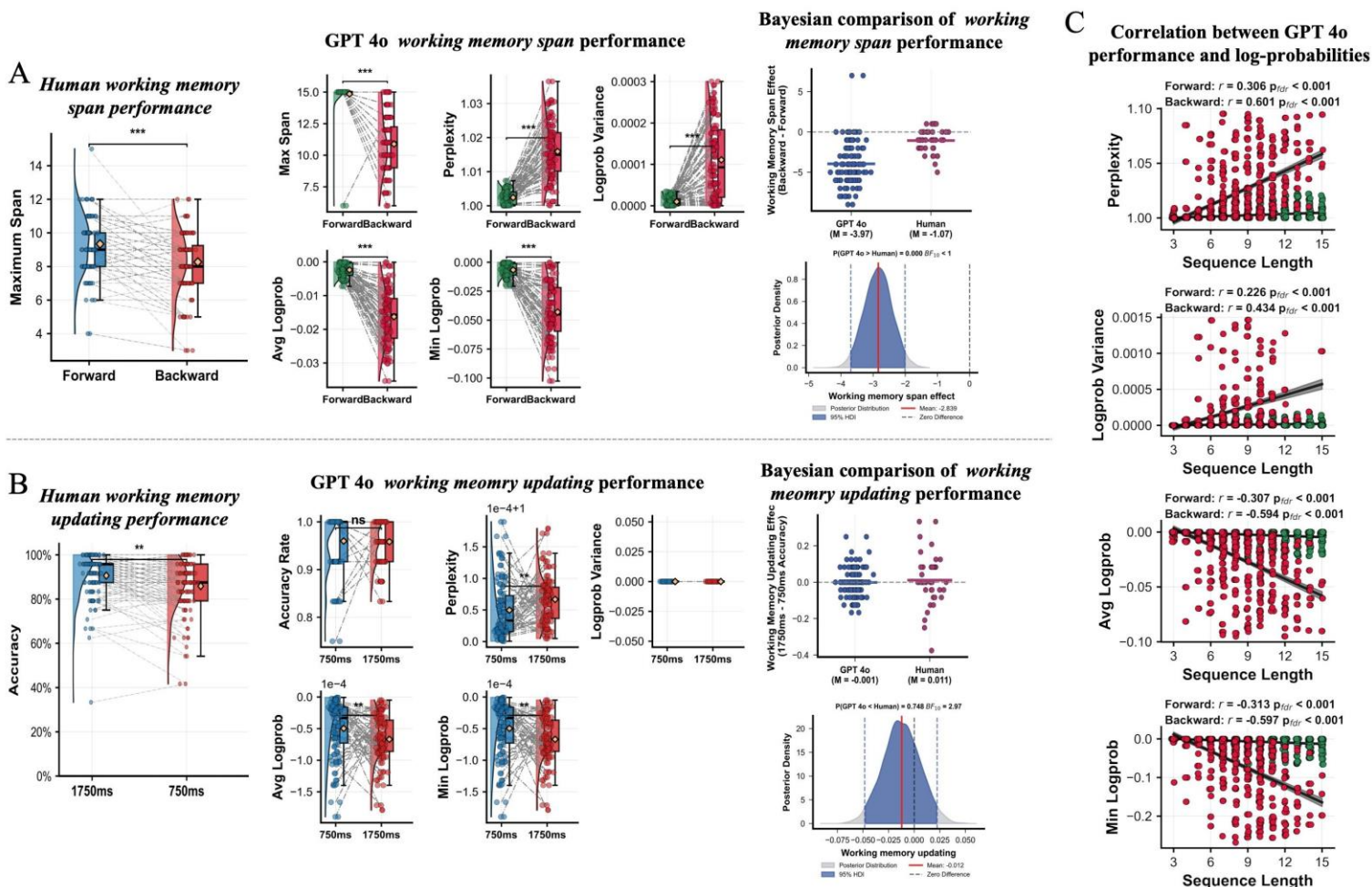
751

752

753

754

755



756

757

Figure 3. Performance comparison between GPT-4o and human participants in working memory tasks. (A) Working memory span in forward *versus* backward *digit span* tasks for human participants ($N = 120$) and the corresponding working memory span with log probability parameters for GPT-4o. (B) Accuracy at different stimulus presentation times (1750 and 750 ms) in the running digit memory task for human participants ($N = 120$) and the corresponding accuracy with log probability parameters for GPT-4o. The rightmost panel shows Bayesian analysis comparing refresh time effects (accuracy at 750 ms minus accuracy at 1750 ms) between humans and GPT-4o. (C) Bayesian analysis comparing working memory load effects (forward span minus backward span) between humans and GPT-4o. (D) Correlation analysis between sequence length and log probability parameters in forward and backward digit span tasks for GPT-4o.

758

759

760

761

762

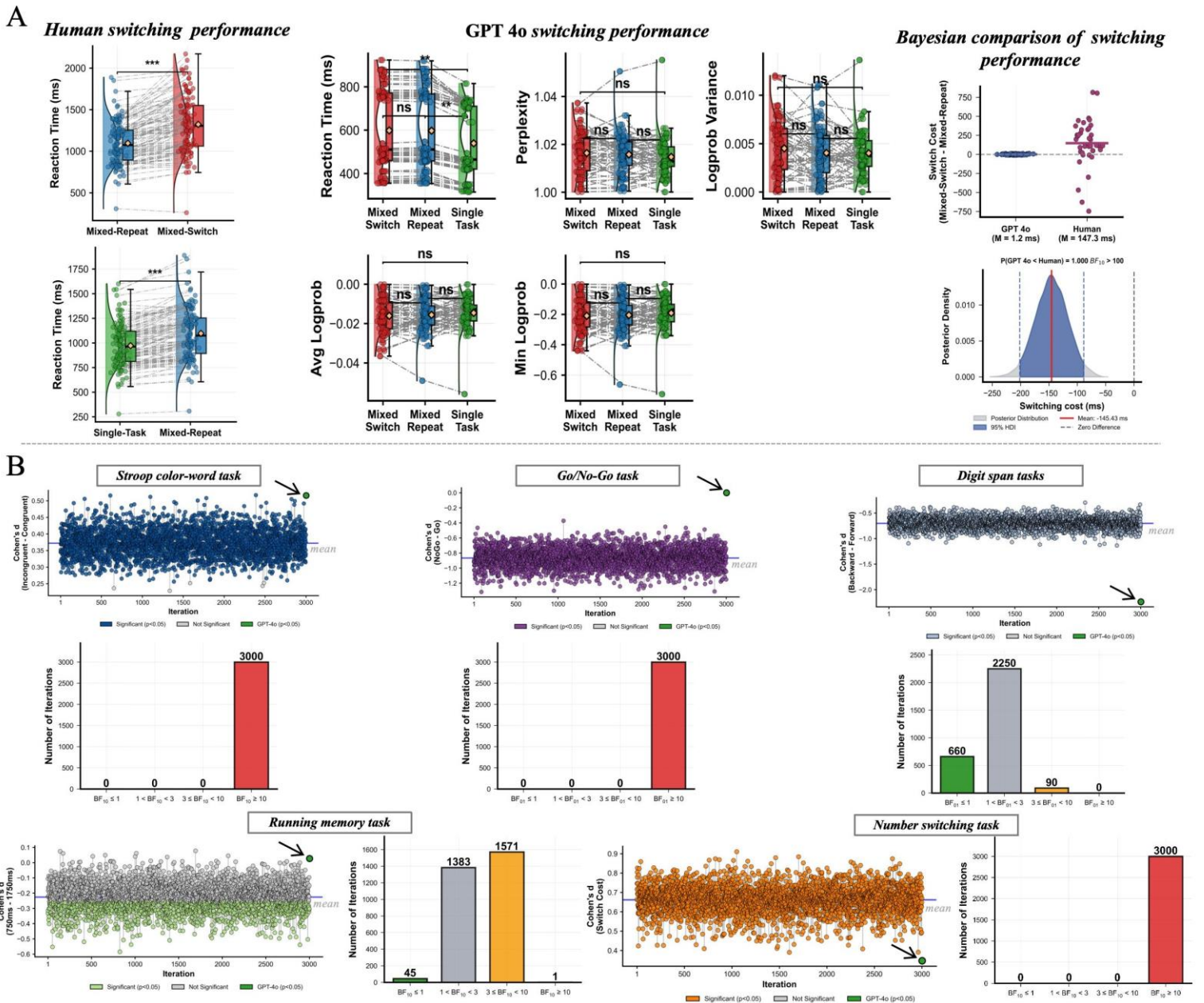
763

764

765

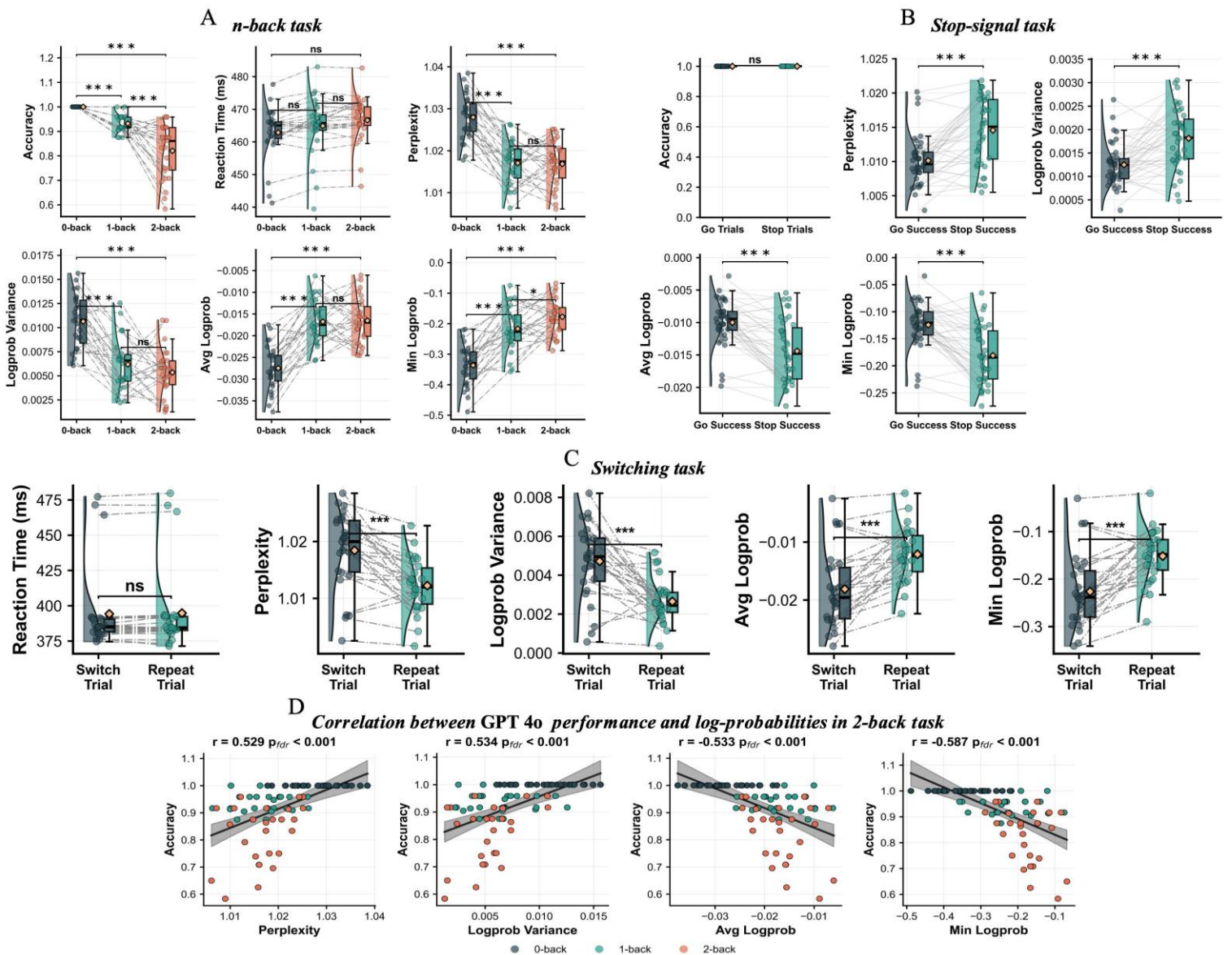
766

767



768

769 **Figure 4.** Comparison between GPT-4o and human participants in cognitive flexibility tasks and
 770 subsampling analysis. (A) Switching costs and mixing costs for human participants ($N = 120$) in
 771 the *number-switching* task, with corresponding reaction times and log probability parameters
 772 across three task conditions for GPT-4o. The rightmost panel presents a Bayesian analysis
 773 comparing switching costs between humans and GPT-4o. (B) Distribution of results from 3,000
 774 repeated random samples ($N = 120$ each) drawn from the full human dataset ($N = 1,970$),
 775 demonstrating the statistical robustness of the primary findings.



776

777 **Figure 5.** Behavioral performance in three executive function tasks from independent samples.

778 (A) Performance comparison between human participants and GPT-4o in the *n*-back task. (B)

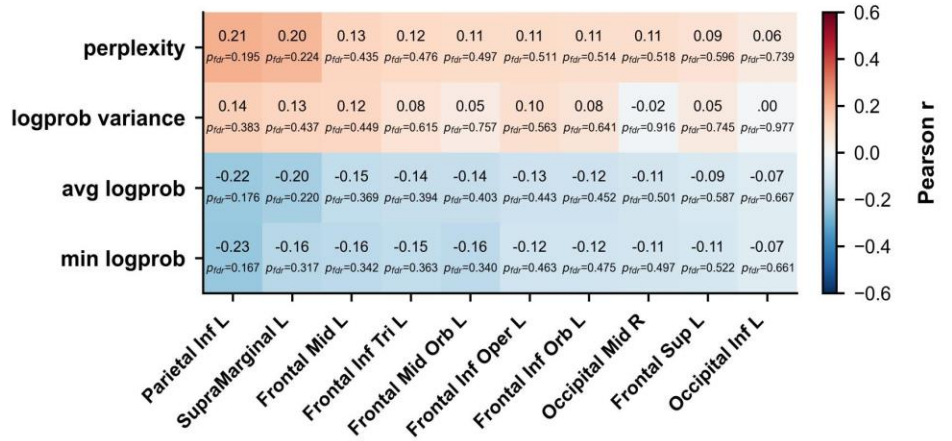
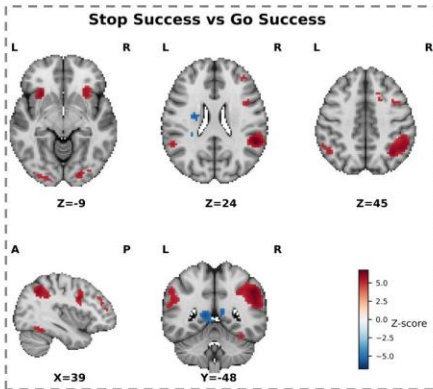
779 Performance comparison between human participants and GPT-4o in the *stop-signal* task. (C)

780 Performance comparison between human participants and GPT-4o in the *number-letter category-*

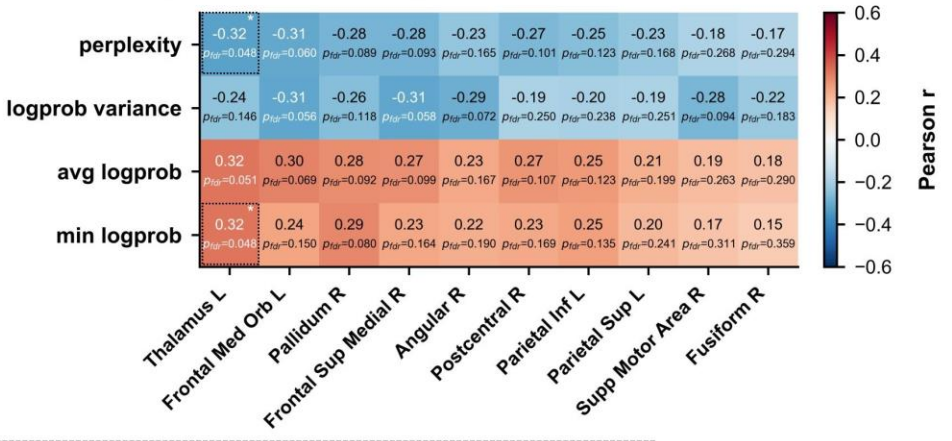
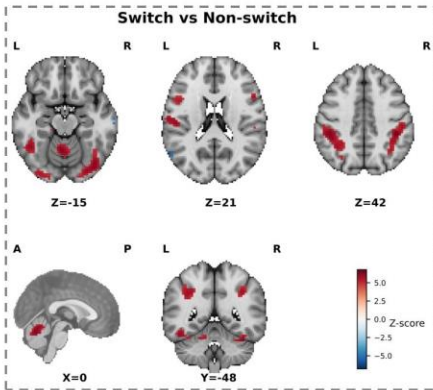
781 *switching* task. (D) Correlation analysis between behavioral performance metrics and log

782 probability parameters in the 2-back task for GPT-4o.

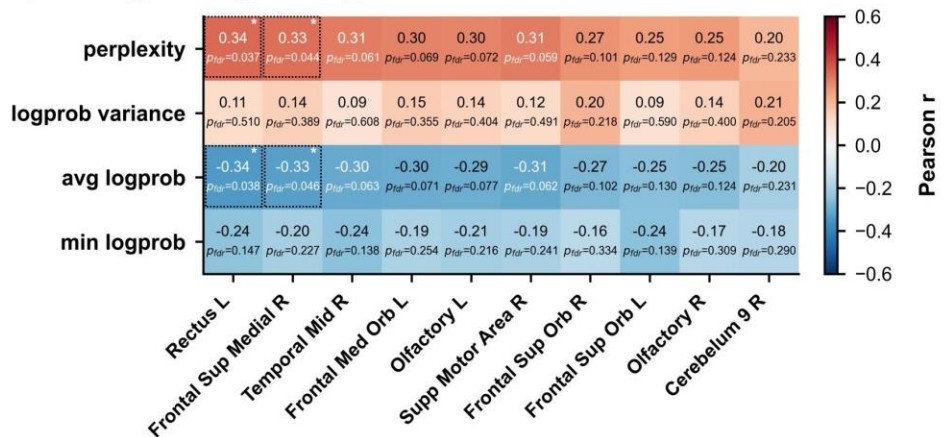
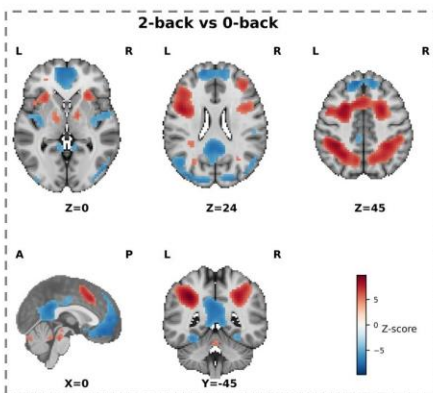
A *Stop-signal task (assessing inhibitory control)*



B *Switching task (assessing cognitive flexibility)*



C *n-back task (assessing working memory)*



783

784

785

786

787

Figure 6. Relationship between GPT-4o internal log probability parameters and human brain neural activity. (A) *Stop-signal task* (measuring response inhibition). Left panel displays group-level brain activation patterns for the “successful stop > successful Go” contrast (Bonferroni corrected, $p < 0.05$); right panel presents correlation analyses between mean activation values

788 from significantly activated brain regions (116 ROIs defined by the Automated Anatomical
789 Labeling [AAL] atlas) and GPT-4o log probability parameters. (B) *Number-letter category-*
790 *switching* task (measuring cognitive flexibility). The left panel shows brain activation patterns for
791 the “successful switch > successful repeat” contrast; the right panel displays correlation analyses
792 between significantly activated regions and model parameters. (C) *n*-Back task (measuring
793 working memory). The left panel illustrates brain activation patterns for the “2-back > 0-back”
794 contrast; the right panel shows correlation analyses between significantly activated regions and
795 model parameters.