

1 **Empirical examination of the replicability of associations between brain structure and**
2 **psychological variables**

3 Shahrzad Kharabian Masouleh^{1,2}, Simon B. Eickhoff^{1,2}, Felix Hoffstaedter^{1,2} and Sarah
4 Genon^{1,2}

5 ¹Institute of Neuroscience and Medicine (INM-7: Brain and Behaviour), Research Centre
6 Jülich, Jülich, Germany

7 ²Institute of Systems Neuroscience, Heinrich Heine University Düsseldorf, Düsseldorf,
8 Germany;

9

10 **Author's email addresses:**

11 s.kharabian@fz-juelich.de

12 s.eickhoff@fz-juelich.de

13 f.hoffstaedter@fz-juelich.de

14 s.genon@fz-juelich.de

15 **Corresponding authors:**

16 Shahrzad Kharabian Masouleh.
17 Brain and Behaviour (INM-7)
18 Institute for Neuroscience and Medicine
19 Research Centre Jülich
20 52425 Jülich
21 E-mail: s.kharabian@fz-juelich.de

22

23 Dr. Sarah Genon.
24 Brain and Behaviour (INM-7)
25 Institute for Neuroscience and Medicine
26 Research Centre Jülich
27 52425 Jülich
28 E-mail: s.genon@fz-juelich.de

29

30

31 **Abstract**

32 Linking interindividual differences in psychological phenotype to variations in brain structure
33 is an old dream for psychology and a crucial question for cognitive neurosciences. Yet,
34 replicability of the previously-reported “structural brain behavior” (SBB)-associations has been
35 questioned, recently. Here, we conducted an empirical investigation, assessing replicability of
36 SBB among healthy adults. For a wide range of psychological measures, the replicability of
37 associations with gray matter volume was assessed. Our results revealed that among healthy
38 individuals 1) finding an association between performance at standard psychological tests and
39 brain morphology is relatively unlikely 2) significant associations, found using an exploratory
40 approach, have overestimated effect sizes and 3) can hardly be replicated in an independent
41 sample. After considering factors such as sample size and comparing our findings with more
42 replicable SBB-associations in a clinical cohort and replicable associations between brain
43 structure and non-psychological phenotype, we discuss the potential causes and consequences
44 of these findings.

45

46

47

48

49

50

51

52

53

54

55 **Introduction:**

56 The early observations of inter-individual variability in human psychological skills and traits
57 have triggered the search for defining their correlating brain characteristics. Studies using in-
58 vivo neuroimaging have provided compelling evidence of a relationship between human skills
59 and traits and brain morphometry that were further influenced by individuals' years of
60 experience, as well as level of expertise. More subtle changes were also shown following new
61 learning/training (Draganski et al., 2004; Taubert et al., 2011), hence further demonstrating
62 dynamic relationships between behavioral performance and brain structural features. Such
63 observations quickly generated a conceptual basis for growing number of studies aiming to
64 map subtle inter-individual differences in observed behavior such as personality traits (Nostro
65 et al., 2017), impulsivity traits (Matsuo et al., 2009) or political orientation (Kanai et al., 2011);
66 to normal variations in brain morphology (for review see (Genon et al., 2018; Kanai and Rees,
67 2011)). Altogether, these studies created an empirical background supporting the assumption
68 that the morphometry of the brain in humans is related to the wide spectrum of aspects observed
69 in human behavior. Such reports on structural brain behavior (SBB) associations may not only
70 have important implications in psychological sciences and clinical research (Ismaylova et al.,
71 2018; Kim et al., 2015; Luders et al., 2013, 2012; McEwen et al., 2016), but also possibly hold
72 an important key for our understanding of brain functions (Genon et al., 2018) and thus concern
73 basic research in cognitive neurosciences.

74 Yet, along with the general replication crisis affecting psychological sciences (Button et al.,
75 2013; De Boeck and Jeon, 2018; Open Science Collaboration, 2015), replicability of the
76 previously reported SBB-associations were also questioned recently. In particular, Boekel et al.
77 (2015) in a purely confirmatory replication study, picked on few specific previously reported

78 SBB-associations. Strikingly, for almost all the findings under scrutiny, they could not find
79 support for the original results in their replication attempt.

80 In another study we demonstrated a lack of robustness of correlations between cognitive
81 performance and measures of gray matter volume (GMV) in a-priori defined sub-regions of the
82 dorsal premotor cortex in two samples of healthy adults (Genon et al., 2017). Although our
83 study did not primarily aim to address the scientific qualities of SBB, it revealed, in line with
84 Boekel et al. (2015), that a replication issue in SBB-associations could seriously be considered.
85 However, ringing the warning bell of a replication crisis would be premature since these
86 previous studies have approached replicability questions within very specific contexts and
87 methods and using small sample sizes (Muhlert and Ridgway, 2016).

88 In particular, Boekel et al. and Genon et al.'s studies were performed by focusing on a-priori
89 defined regions-of-interest (ROIs). However, most SBB studies are commonly assessed in
90 groups of dozens of individuals, using an exploratory setting employing a mass-univariate
91 approach. Thus, the null findings of the two questioning studies could be related to the focus
92 and averaging of GMV within specific region-of-interests as suggested by (Kanai, 2016) and
93 discussed in (Genon et al., 2017).

94 In stark contrast with this argument, in whole-brain exploratory SBB studies, the multitude of
95 statistical tests that is performed (as the associations are tested for each voxel, separately) likely
96 yield many false positives. Directly addressing this limitation, several strategies for multiple
97 comparison correction have been proposed to control the rate of false positives (Eklund et al.,
98 2016). We could hence assume that the high number of multiple tests and general low power
99 of neuroimaging studies combined with the flexible analysis choices (Button et al., 2013;
100 Poldrack et al., 2017; Turner et al., 2018) represent critical factors likely to lead to the detection
101 of spurious and not replicable associations. Nevertheless, an empirical evaluation of the

102 replicability of the findings yielded by an exploratory approach is still crucially lacking to allow
103 questioning the replicability of exploratory SBB studies.

104 Thus in the following study, we empirically examined replicability rates of SBB-association
105 over broad range of psychological scores, among healthy adults. Similar to the commonly used
106 approach in the literature, we first identified “significant” findings with an exploratory
107 approach, searching for associations of GMV with psychometric variables across the whole
108 brain. Here a linear model is fit between inter-individual variability in the psychological score
109 and GMV at each voxel. Inference is then made at cluster level, using a threshold-free cluster
110 enhancement approach (Smith and Nichols, 2009). We then investigated the reproducibility of
111 these findings, across resampling, by conducting similar whole-brain voxel-wise exploratory
112 analysis within 100 randomly generated subsamples of individuals and comparing the spatial
113 overlap of the significant findings that survive multiple comparison correction, across all
114 samples (discovery samples). Additionally, for each of the 100 exploratory analyses, we
115 assessed replicability of SBB-associations using a confirmatory approach. Here, average GMV
116 within regions showing significant SBB-association in the initial exploratory analysis, i.e.
117 ROIs, are calculated among a demographically-matched independent sample and their
118 association with the same psychological score is compared between the matched-discovery and
119 -replication sub-samples (see Methods for more detail).

120 In line with the replication literature, we further examined the influence of sample size and
121 replication power on reproducibility of SBB-associations. We also investigated the relationship
122 between the effect size of exploratory and confirmatory analyses. In line with previous studies
123 and the reproducibility literature, we included the Bayes Factors as an indicator of evidence in
124 favor of the null or alternative hypotheses (Boekel et al., 2015). Finally, in order to promote
125 discussion on the underlying reality which is aimed to be captured by SBB in the framework of
126 the psychology of individual differences, we included as benchmarks non-psychological
127 phenotypical measures, i.e. age and body-mass-index (BMI), and extended our analysis to a

128 clinical sample, where SBB-associations are expected to enjoy higher biological validity. For
129 this purpose, a subsamples of patients drawn from Alzheimer's Disease Neuroimaging Initiative
130 (ADNI) database were selected, in which replicability of structural associations of immediate-
131 recall score from Rey auditory verbal learning task (RAVLT) (Schmidt, 1996) was assessed
132 (see Methods). Due to availability of the same score within the healthy cohort, this later analysis
133 is used as a “conceptual” benchmark.

134

135

136 **Results:**

137 A total of 10800 exploratory whole brain SBB associations (each with 1000 permutations) were
138 tested to empirically identify the replicability of the associations of 36 psychological scores
139 with GMV over 100 splits in independent matched subsamples, at three pre-defined sample
140 sizes, within the *healthy* cohort.

141 Altogether, in contrast to GMV-associations with age and BMI, significant SBB-associations
142 were highly unlikely. For the majority of the tested psychological variables no significant
143 association with GMV were found in beyond 90% of the whole brain analyses.

144 ***SBB-associations among the healthy population:***

145 *Replicability of “whole brain exploratory SBB-associations”:*

146 Age and BMI structural associations: Voxel-wise associations of age and BMI with GMV, as
147 suggested by previous studies (Fjell et al., 2014; Kharabian Masouleh et al., 2016; Salat et al.,
148 2004; Willette and Kapogiannis, 2014), were widespread and strong. In order to avoid large
149 clusters that simultaneously cover several cortical and subcortical regions, we focused on local
150 peaks of associations by increasing the voxel-level t-threshold of the statistical maps. The
151 modified voxel-level t-threshold was set to 8 and 3, for defining age- and BMI-associated
152 clusters, respectively. These *arbitrary* thresholds were chosen such that the very large clusters
153 would divide into smaller ones, while still retaining the general spatial pattern of the significant
154 regions.

155 Even with these adapted thresholds, for almost all subsamples, we found highly consistent
156 widespread negative associations of age with GMV (see figure 1A for aggregate maps of spatial
157 overlap of exploratory findings and density plots, summarizing distribution of “frequency of
158 significant findings” within each map).

159 When decreasing the sample size of the discovery cohort, the spatial overlap of significant
160 findings over 100 splits decreased. More specifically, for the discovery sample of 326 subjects,
161 more than half of the significant voxels were consistently found as being significant in beyond

162 90% of the whole-brain exploratory analyses (i.e. high level of spatial consistency of significant
163 findings). As the size of the subsamples decreased, the shape of the distribution also changed,
164 and the median of the density plots fell around 50% and even 10% for samples consisting of
165 232 and 138 individuals, respectively.

166 Similar results, though with much lower percent of consistently overlapping voxels, were seen
167 for negative associations of BMI with GMV. The density plots and the spatial maps of Figure
168 1B show that for the larger samples (consisting of 326 and 232 subjects) few voxels were
169 consistently found in “all” (100%) subsamples as having significant negative association with
170 BMI. For the smaller samples (with 138 participants) the maximum replicable association was
171 found in 93% of the splits and 4 out of 100 exploratory analyses did not result in any significant
172 clusters (Table 1). Additionally, as Figure 2B shows, the majority of significant voxels had a
173 replicability bellow 50%.

174 These results highlight the influence of sample size on the replicability (frequency of overlap)
175 of whole-brain significant associations, even for age and BMI, for which we expected more
176 stable associations with morphological properties of the brain.

177 Structural associations of the psychological scores: In contrast, for most of the psychological
178 scores, only few of the 100 discovery subsamples yielded significant clusters. Table 1 and
179 supplementary Table 2 show the number of splits for which the exploratory whole-brain SBB-
180 analysis resulted in *at least one* significant positively or negatively associated cluster for each
181 score. These results reveal that finding significant SBB-associations using the exploratory
182 approach in healthy individuals is highly *unlikely* for most of the psychological variables.
183 Furthermore, the significant findings were spatially very diverse, that is, spatially overlapping
184 findings were very rare.

185 We here retained for further analyses the three psychological scores for which the discovery
186 samples most frequently resulted in at least one significantly associated cluster. These three
187 scores were the Perceptual reasoning score of WASI (Wechsler, 1999), the number of correct

188 responses in word-context test and the interference time in the color-word interference task. For
189 example, for the discovery samples of 326 adults, in 83 out of 100 randomly generated
190 discovery samples, at least one cluster (not necessarily overlapping) showed a significant
191 positive association between perceptual reasoning and GMV (Table 1)). Of note, these more
192 frequently found associations were in the direction linking better task performance with higher
193 GMV.

194 Yet again, in line with our observations for BMI associations, the probability of finding at least
195 one significant cluster tend to decrease in smaller discovery samples (see Table 1). Likewise,
196 as the discovery sample size decreased, the maximum rate of spatial overlap, as denoted by the
197 height of the density plots, decreased (see Figure 1C-F). The width of these plots show that the
198 majority (> 50%) of the significant voxels spatially overlapped only in less than 10% of the
199 discovery samples. In the same line, the variability depicted by the spatial maps highlight that
200 many voxels are found as significant only in one out of 100 analyses.

201 These results highlight that finding a significant association between normal variations on
202 behavioral scores and voxel-wise measures of GMV among healthy individuals is highly
203 unlikely, for most of the tested domains. Furthermore, they underscore the extent of spatial
204 inconsistency and the *poor replicability* of the significant SBB-associations from *exploratory*
205 *analyses*.

206 -----**Table 1** -----

207 -----**figure1**-----

208 *Confirmatory ROI-based SBB-replicability:*

209 Age and BMI effects: Irrespective of the size of the test subsamples and definition used to
210 identify “successful” replication (see Methods), for all ROIs negative age-GMV associations
211 were “successfully” replicated in the matched test samples. Unlike the perfect replication of
212 age-associations, replication rate of BMI effects depended highly on the test sample size and
213 the criteria used to characterize “successful” replication. Over all three tested sample sizes, in

214 more than 90% of the a-priori defined ROIs, BMI associations were found to be in the same
215 “direction” in the discovery and test samples (i.e. replicated based on “sign” criteria). The
216 examination of replicated findings based on “statistical significance” revealed replicated effects
217 in more than 57% of ROIs. This rate of ROI-based replicability increased from ~57% to 75%,
218 as the test sample size increased from 140 to 328 individuals (see figure 2). Furthermore, as the
219 dark blue segments in the outer layers of figure 2 indicates, Bayesian hypothesis testing
220 revealed moderate-to-strong evidence for H1 in more than 30% of the ROIs.

221 -----figure2 -----

222 Psychological variables: Figure 2 also illustrates the replicability rates of structural associations
223 of the top three psychological measures from the whole brain analyses (the perceptual reasoning
224 score of WASI, the number of correct responses in word-context test and the interference time
225 in the color-word interference task).

226 Despite the structural associations of perceptual reasoning score being in the same direction
227 (positive SBB-association), for the majority of the ROIs (>85%), less than 31% of all ROIs
228 showed replicated effects based on “statistical significance” criterion. Finally, less than 4% of
229 the ROIs were identified as “successfully replicated” based on the Bayes factors. (Figure 2).

230 For the three tested samples sizes, associations of the word-context task were in the same
231 direction (positive SBB-association) in the discovery and test pairs in ~75% of ROIs.
232 Nevertheless, again, the rate of statistically “significantly”-replicated ROIs ranged between 17
233 to 26%. Furthermore, even less than 8% of all ROIs showed replicated effects based on the
234 Bayes factors (moderate-to-strong evidence for H1) (Figure 2).

235 Finally, negative correlations between interference time of the color-word interference task and
236 average GMV were depicted in ~70 % of the ROIs, but significant-replication was found in
237 only 11% to 17% of all ROIs, for the three test sample sizes. Along the same line, replication
238 based on the Bayes factors was below 5% (Figure 2E).

239 In general, these results show the span of replicability of structural associations from highly
240 replicable age-effects to very poorly replicable psychological associations. They also highlight
241 the influence of the sample size, as well as the criteria that is used to define successful
242 replication on the rate of replicability of SBB-effects in independent samples.

243 *Effect size in the discovery sample and its link with effect size of the test sample and actual*
244 *replication:*

245 Figure 3 plots discovery versus replication effect size for each ROI and for three test sample
246 sizes. Focusing on by-“sign” replicated ROIs (blue), for the three psychological scores
247 (perceptual reasoning, word-context and CWI) revealed that the discovery samples resulted in
248 overall larger effects (magnitude) compared to the test samples. Indeed, the marginal
249 distributions are centered around smaller effect sizes in the y-dimension (test sample) compared
250 to the x-axis (discovery samples). Furthermore, for these by-“sign” replicated ROIs, there was
251 no positive relationship between the effect sizes of the behavioral associations in the discovery
252 and test samples (blue lines in each subplot).

253 For BMI and age, however, the effect sizes of the discovery and test pairs were generally
254 positively correlated, suggesting that the ROIs with greater negative structural association with
255 BMI (or age) in the discovery sample, also tended to show stronger negative associations within
256 the matched test sample.

257 To investigate if the replication power, estimated using the effect size of the discovery samples,
258 was linked to a higher probability of *actual* replication in the test samples, the ROIs were
259 grouped into replicated and not-replicated, based on the “statistical significance” criterion.
260 While the estimations of statistical power were generally higher among the replicated compared
261 to not-replicated ROIs for BMI associations (p-value of the Mann-Whitney U tests $< 10^{-5}$), for
262 structural associations of the psychological scores, this was not the case. Strikingly, for the
263 structural associations of perceptual reasoning, over all sample sizes, the significantly
264 replicated ROIs tended to have **lower** estimated power compared to the ROIs that actually were

265 not-replicated (p-value of the Mann-Whitney U tests $< 10^{-5}$). These unexpected findings
266 highlight the unreliable aspect of effect size estimations of SBB-associations within the
267 discovery samples among healthy individuals. They also demonstrate that these inflated effect
268 sizes result in flawed and thus uninformative estimated statistical power.

269 -----figure3 -----

270

271 ***Structural associations of total immediate recall score in ADNI cohort:***

272 *Replicability of “whole brain exploratory associations”:*

273 Within the sample of patients from ADNI-cohort, 84 out of the 100 whole-brain exploratory
274 analyses resulted in *at least one* significant cluster showing a positive association between the
275 immediate-recall score and GMV. In the healthy population, however, the same score resulted
276 in a significant cluster in only less than 10% of exploratory analyses, for any of the three
277 discovery sample sizes (supplementary Table 2 and supplementary Figure 1).

278 As could be seen in the spatial maps of Figure 4, significant associations in the ADNI cohort
279 were found across several brain regions including the bilateral lateral and medial temporal lobe,
280 the lateral occipital cortex, the precuneus, the superior parietal lobule, the orbitofrontal cortex
281 and the thalamus. Although most of the significant voxels were found by less than 10% of the
282 splits, some voxels in the bilateral hippocampus were found to be significantly associated with
283 the recall score in more than 70% of the subsamples (peak of spatial overlap; see Figure 4A,
284 B).

285 *Confirmatory ROI-based SBB-replicability:*

286 Figure 4D shows the rates of “successful replication” of associations between the immediate-
287 recall score and GMV within each ROI in the independent, matched-samples. As the most inner
288 layer shows, in more than 94% of ROIs, GMV correlated positively with the recall score in the
289 test subsamples, corroborating the “sign” of the association in the paired-discovery samples.
290 These correlations were significant in 72% of all ROIs. Furthermore, in more than 50% of all

291 ROIs the correlations in the test sample supported, at least moderately, the link between higher
292 GMV and higher recall score (using the Bayes factors).

293 *Association between discovery and replication effect size:*

294 The marginal histograms in Figure 4C suggest that overall the size of effects in the discovery
295 samples are slightly larger than the effects sizes in the paired replication samples. When looking
296 at the ROIs that were successfully replicated (by-sign), there was a positive association between
297 the discovery and replication effect size (spearman's $\rho = 0.38$, $p\text{-value} < 10^{-11}$).

298 Finally, the median replication power was higher among “significantly replicated” ROIs,
299 compared to not replicated (defined using “statistical significance criterion”) ROIs ($p\text{-value}$ of
300 the mann-whitney U test $< 10^{-3}$). These results showed the superior, yet not perfect,
301 replicability of SBB-associations within the clinical population (see supplementary Figure 2 for
302 structural associations of immediate recall within healthy cohort). The observed somewhat
303 robustness of the findings in ADNI suggest that, when the population under study shows clear
304 variations in both brain structural markers and psychological measurements, such as the patient
305 group in ADNI cohort, the associations between brain structure and psychological performance
306 could be relatively reliably characterized. Nevertheless, again, the occurrence of not-replicated
307 results highlight the importance of confirmatory analyses for a robust characterization of brain-
308 behavior associations.

309

310 -----figure4 -----

311

312 **Discussion:**

313 Our empirical investigation of the replicability of SBB in healthy adults showed that significant
314 associations between psychological phenotype and GMV are not frequent when probing a range
315 of psychometric variables with an exploratory approach. Where significant associations were
316 found, these associations showed a poor replicability.

317 In the following, we first discussed implications of the very low rate of significant findings
318 revealed by the exploratory approach. We then discussed the possible causes of the observed
319 spatial variability of SBB-associations. Those pattern of findings are then compared with the
320 pattern observed in the clinical cohort. Finally, in line with the replication literature in
321 psychological sciences and neurosciences (Button et al., 2013; Poldrack et al., 2017; Turner et
322 al., 2018), we devoted our last section to sample size and power issues in SBB studies in healthy
323 adults and proposed some recommendations.

324 *Infrequent significant SBB associations in healthy individuals: Importance of reporting null* 325 *findings*

326 According to the scientific literature, associations between psychological phenotype (cognitive
327 performance and psychological trait) and local brain structure are not uncommon (Kanai and
328 Rees, 2011). However, in our exploratory analyses, when looking at a range of psychological
329 variables, significant associations with GMV were very rare. It is worth noting that here by
330 having a-priori fixed analysis design and inference routines, we aimed to avoid “fishing” for
331 significant findings (Gelman and Loken, 2014). Flexible designs and flexible analyses routines
332 (Simmons et al., 2011) as well as p-hacking (John et al., 2012) are considered as inappropriate
333 but frequent research practices (Poldrack et al., 2017). Based on our findings of infrequent
334 significant SBB-associations, we could assume that flexible analyses routines, p-hacking and
335 most importantly *publication bias* (Dwan et al., 2013) have contributed to the high number of
336 significant SBB-reports in the literature.

337 When considering potential impacts of biased SBB-reports on our confidence of psychological
338 measures, as well as our conception and apprehension of brain-behavior relationships and
339 psychological interindividual differences, we would strongly argue for null findings reports.
340 Such reports would contribute to a more accurate and balanced apprehension of associations
341 between differences in psychological phenotype and brain morphometric features, but it would
342 also help to progressively disentangle factors that mediate or modulate the relationship between
343 brain structure and behavioral outcomes.

344 *Poor spatial overlap of SBB across resampling: possible causes and recommendations*

345 In addition to the low likelihood of finding “any” significant SBB-association using the
346 exploratory approach, clusters that do survive the significance thresholding did not often
347 overlap in different subsamples. Furthermore, the probability of spatial overlap further dropped
348 as the number of participants in the subsamples decreased (Figure 1). Putting this finding in
349 light of the literature brings two main hypotheses.

350 First, from the conceptual level, we could hypothesize that the pattern of correlation between a
351 psychological measure is by nature spatially diffuse at the brain level. Psychological measures
352 aim to conceptually articulate *behavioral functions and processes*, thus, in most cases, they
353 have not been developed to identify specific localized *brain functions*. Following this
354 philosophical segregation between psychological sciences and neurosciences, it is now widely
355 acknowledged that there is no one-to-one mapping between behavioral functions and brain
356 regions (Anderson, 2015; Genon et al., 2018; Pessoa, 2014). Instead, mapping a psychological
357 concept to brain features usually result in a diffuse brain spatial pattern with small effect sizes
358 (Bressler, 1995; Poldrack, 2010; Tononi et al., 1998). From this axiom, we can expect that
359 several studies conducted in small samples (specifically after rigorous corrections for multiple
360 comparisons) are likely to each capture a partial and minor aspect of the whole true association
361 pattern, resulting in a poor replication rate for each individual study (i.e. high type II error).

362 Alternatively, a more parsimonious hypothesis is a methodological one questioning the truth or
363 validity of the found significant associations hence considering them as spurious (i.e. type I
364 error). Psychological and MRI measurements are both relatively indirect estimations of
365 respectively, behavioral features and brain structural features and thus are susceptible to noise.
366 Correlations in small samples in the presence of noise for both type of variables is likely to
367 produce spurious significant results (Loken and Gelman, 2017) by fitting a correlation or
368 regression between random noise in both variables.

369 Thus, the pattern of poor spatial consistency of SBB findings could result either from factors at
370 the object of study level, i.e. the relationship between brain and behavior, or, from factors at
371 the measurement and analysis level. While the latter hypothesis is more parsimonious, one
372 argument for the former hypothesis comes from the relatively substantial replications by-sign
373 observed in our confirmatory analyses, of three top behavioral scores (see figure 2). If the
374 significant SBB findings would be purely driven by noise in the data, we would expect them to
375 show purely random signs across resampling, which was not the case (but also see
376 Supplementary figure S2 for example of scores with lower replicability and higher inconsistent
377 associations across resampling). Therefore, it is actually likely that both hypotheses hold true
378 and that the spatial variability of significant SBB findings result from both, factors at the
379 analyses levels and factors at the object level, potentially interacting together.

380 It is worth noting that similar complexity and uncertainty have been described for task-based
381 functional associations studies (Cremers et al., 2017; Turner et al., 2018). In particular, Cremers
382 et al. (2017) using simulated and empirical data demonstrated that the task-based functional
383 activations have a generally weak and diffuse pattern. Therefore, Cremers et al. concluded that
384 most whole-brain analyses in small samples, specifically when combined with stringent
385 correction for multiple comparison, to control the false positive rates, would most likely
386 frequently overlook global meaningful effects and depict results with poor replicability (type II
387 error). On the other hand, in the present study, higher spatial extent and lower consistency of

388 significant findings in smaller samples in Figure 1, also suggests a higher number of spurious
389 associations (type I error) in smaller samples (due to winners curse (Button et al., 2013;
390 Forstmeier and Schielzeth, 2011)) than in the larger samples.

391 These factors, added to the complexity of human behavior, renders the objective of capturing
392 covariations with psychometric variables in brain structure *locally* particularly challenging. For
393 that reason, in exploratory studies whose aim is to identify brain structural features correlating
394 with a given (set of) psychological variable(s), a multivariate approach could be advised
395 (Habeck and Stern, 2010; McIntosh and Mišić, 2013). As all methods, multivariate analyses
396 have their own limitations: in particular, the ensuing difficulty of interpretability of the revealed
397 pattern. While some authors argue either for one or the other approach, the use of these
398 approaches are far from being mutually exclusive (Moeller and Habeck, 2006). Combining both
399 approaches in small datasets indeed revealed that the results of the univariate approach reflect
400 the “tip of the iceberg” of the behavior’s brain correlates, whose spatial extent are more
401 comprehensively captured with the multivariate analysis, but interpretability is facilitated by
402 the use of univariate analyses; e.g. (Genon et al., 2016, 2014).

403 Thus, to partially address the previously described concerns of small and spatially diffuse
404 effects at the brain level in exploratory whole-brain-behavior study, we here recommend to
405 combine a univariate and a multivariate approach. This solution may help to reduce the false
406 negatives, yet it does not provide any protection against the influence of noise that may affect
407 both approaches.

408 *Confirmatory replication of exploratory SBB findings: importance of out of sample replication*
409 ROI-based analysis further highlighted that significant associations, which have been
410 discovered when starting with a psychological measure and searching within the whole brain
411 for a significant association (i.e. “evidenced in an exploratory study”), show poor replicability
412 (using significance and Bayes factor criteria, but also using similar sign criterion for most
413 psychometric scores; For example, see Supplementary Figures S1 and S2.) in a confirmatory

414 ROI-based study (in line with what was previously shown by Boekel et al. (2015)). These
415 findings thus call for a general acknowledgment of the uncertainty and fragility of exploratory
416 findings and the need for *out of sample* confirmatory replications to provide evidence about the
417 robustness of the reported effects (Ioannidis, 2018; Tukey, 1980).

418 *Further factors influencing replicability of SBB-findings: power of replication and object of*
419 *study*

420 Another clear finding of our study is the overestimation of the effect size in the exploratory
421 approach (Kriegeskorte et al., 2010), specifically in smaller samples (see marginal distributions
422 of the x- and y-axis in Figure 3). For the majority of the psychological scores, in the ROI-based
423 approach, we failed to find a clear association between effect size in the discovery and
424 replication samples. Instead, we observed a rather high estimated statistical power for
425 replication (due to an inflated effect size estimation (Ioannidis, 2008)), despite very low actual
426 rate of replicated effects in the independent samples. These findings are particularly important
427 when considering the current research context, in which power analyses are encouraged to
428 justify the allocation of financial and human investment in specific future researches.
429 Prospective studies with power analyses are frequently proposed, where power is computed
430 based on the findings from previous exploratory analyses in a small sample (Albers and Lakens,
431 2018a). An inflated effect size estimation from the exploratory study results in an unreliable
432 high power, which in turn lead to confidence in prospective studies to find relevant findings
433 and hence in the allocation and possibly waste of (frequently public) resources (Albers and
434 Lakens, 2018b; Poldrack et al., 2017). Nevertheless, this provocative conclusion does not imply
435 that SBB studies should be banished to hell. Our conclusion here mainly concerns the study of
436 association between variations at *standard psychological measures* and variations in *measures*
437 *of gray matter* in “*small*” *samples of healthy individuals*. Our results further show that different
438 types of SBB exploratory studies should not be epistemologically all put in the same pot.

439 In support for this argument, in ADNI sample, despite the additional confounding effect of
440 different scanners and/or scanning parameters due to the multi-site nature of the cohort,
441 associations between immediate-recall score and GMV were relatively stable. Compared to
442 associations of the same measure of verbal learning performance within the healthy population
443 (see supplementary Figure 1), these results highlight the superior reliability of SBB-
444 associations that are defined in a clinical context. These findings have important conceptual
445 implications. From an epistemological and conceptual point of view, our comparative
446 investigation suggests that the object of study matters in the replicability of SBB. Searching for
447 correlation between variations in cognitive performance and GMV in healthy adults, on one
448 hand, and in neurodegenerative patients, on the other hand, appear as two different objects of
449 study, with different replicability rates. While several SBB results in healthy population are
450 likely to be spurious (see supplementary Table 2), it seems that SBB in clinical population are
451 more likely to capture true relationships.

452 Thus, maybe the conceptual objective itself should be questioned: should we expect the
453 association between normal psychological phenotype, in particular cognitive performance, in
454 healthy population to be substantially driven by local brain macrostructure morphology? Brain
455 structure can certainly not be questioned as the primary substrates of behavior and more than
456 a century of lesion studies recalls this primary principle to our attention (Broca, 1865; Scoville
457 and Milner, 1957), but this does not imply that “normal” variations at standard psychological
458 tests can be related to variations in markers of local brain macrostructure. Our results suggest
459 that reliable answer to this important question requires substantially big samples (bigger than
460 those used here) and independent replications.

461 *Further recommendation: Large sample sizes are important both for exploratory as well as*
462 *replication analyses*

463 The sample size and related power issues hold a central position in the current discussions of
464 the replication crisis in behavioral sciences, as well as in neuroimaging studies (Button et al.,

465 2013; Ioannidis, 2005; Lilienfeld, 2017; Munafò et al., 2017; Open Science Collaboration,
466 2015). Our ROI-based confirmatory analysis suggests that samples consisting of ~200-300
467 participants have in reality yet low power to identify reliable SBB-associations among healthy
468 participants. However, the sample size of SBB studies is usually substantially smaller. Figure
469 5 depicts the distribution of sample sizes (log-scale) of published studies examining GMV in
470 human participants with the standard voxel-based morphometry approach across previous years
471 (BrainMap data (Vanasse et al., 2018)). SBB studies in healthy adults also fall under this
472 general trend. Based on our current work, we would argue that *the probability of finding*
473 *spurious or inconclusive results and exaggerated effect size estimations in these studies is thus*
474 *quite high* (Albers and Lakens, 2018b; Schönbrodt and Perugini, 2013; Yarkoni, 2009).
475 In addition, to underscore the importance of the sample size, our analyses and results further
476 show that the size of the *replication sample* also matters when examining the replicability of a
477 previous SBB findings. This is an obvious factor that has been frequently neglected in the
478 discussions about replication crisis. Yet, while many replication studies straightforwardly
479 blame the sample size of the original studies, it is important to keep in mind that a replication
480 failure might also come from a too small sample size of the replication study (Muhler and
481 Ridgway, 2016).

482 -----figure5 -----

483 *Summary and conclusions*

484 Overall, our work and review of the recent and concomitant replication literature in related
485 fields demonstrate that several improvements could be recommended to get more accurate
486 insight on the relationship between psychological phenotype and brain structure and to
487 progressively answer open questions. Importantly, our recommendations and suggestions
488 concern different levels of SBB researches: the dataset level, the analyses level, as well as at
489 the post-publication and replication level.

490 *At the dataset level*, our study pointed out the need for big data samples to identify robust
491 associations between psychological variables and brain structure, with sample size of at least
492 several hundreds of participants. It should be acknowledged that this conclusion is easier to
493 achieve than to implement in research practice. Nevertheless, large scale cohort datasets from
494 healthy adult populations, such as eNKI used in the current study, human connectome project
495 (HCP) (Van Essen et al., 2013) and UK-biobank (Miller et al., 2016) are now openly available,
496 hence facilitating endeavor in that direction.

497 *At the analysis level*, we recommend the combined use of multivariate, for comprehensive
498 assessment of spatial extent of associations and univariate, to facilitate interpretability, analyses
499 to study brain structural covariates of psychological measures. Furthermore, we emphasize on
500 the importance of *independent* confirmatory replications to provide evidence about the
501 robustness of the effects.

502 Finally, *at the post-analysis level*, we concluded from our observations that publication of null
503 findings should be more encouraged. In addition to directly contributing in generation of a more
504 objective picture of SBB-associations, these null-reports could contribute to new quantitative
505 approaches. In particular, meta-analyses of published literature (Vanasse et al., 2018) would
506 clearly benefit from such unbiased reports of null findings.

507 Sharing raw data would undoubtedly improve the problem of low statistical power, but if not
508 possible, sharing the unthresholded statistical maps (e.g. through platforms such as Neurovault

509 (Gorgolewski et al., 2015)) could also be a significant scientific contribution. In addition to
510 directly contribute to our understanding of brain-behavior relationship, such efforts would open
511 up new possibilities for estimating the correct size and extent of effects by integrating
512 unthresholded statistical maps in the estimation of the effects sizes throughout the brain. Thus,
513 we could hope that sharing initiatives will also contribute indirectly to more valid and insightful
514 SBB studies in the remote future and hence to a better allocation of resources.

515

516

517 **Methods:**

518 ***Participants:***

519 Healthy adults' data were selected from the enhanced NKI (eNKI) Rockland cohort (Nooner et
520 al., 2012). Data collection received ethics approval through both the Nathan Klein Institute and
521 Montclair State University. Written informed consent was obtained from all participants.

522 We focused only on participants for which good quality T1-weighted scans was available along
523 with timewise-corresponding psychological data. Exclusion criteria consisted of alcohol or
524 substance dependence or abuse (current or past), psychiatric illnesses (eg. Schizophrenia) and
525 current depression (major or bipolar). Furthermore, we excluded participants with missing
526 information on important confounders (age, gender, education) or bad quality of structural
527 scans after pre-processing, resulting in a total sample of 466 healthy participants (age: 48 ± 19 ,
528 153 male).

529 Replicability of SBB-associations within the clinical sample was investigated using a
530 subsample drawn from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database,
531 which was launched in 2003 as a public-private partnership and led by Principal Investigator
532 Michael W. Weiner. The primary goal of ADNI has been to test whether serial magnetic
533 resonance imaging (MRI), positron emission tomography (PET), other biological markers, and
534 clinical and neuropsychological assessment can be combined to measure the progression of
535 mild cognitive impairment (MCI) and early Alzheimer's disease (AD). For up-to-date
536 information, see www.adni-info.org.

537 We used the baseline measurements from 371 patients (age : 71 ± 7 , 200 male ; 47 with
538 significant memory complaint, 177 early MCI, 85 late MCI and 62 AD patients (defined based
539 on ADNI diagnostic criteria, see [http://adni.loni.usc.edu/wp-content/themes/freshnews-dev-](http://adni.loni.usc.edu/wp-content/themes/freshnews-dev-v2/documents/clinical/ADNI-2_Protocol.pdf)
540 [v2/documents/clinical/ADNI-2_Protocol.pdf](http://adni.loni.usc.edu/wp-content/themes/freshnews-dev-v2/documents/clinical/ADNI-2_Protocol.pdf)), in whom anatomical brain scans had been
541 acquired in a 3Tesla scanner (from 39 different sites).

542

543 ***Phenotypical measurements:***

544 *Non-psychological measurements:*

545 Age and body mass index (BMI) are highly reliably assessed and their association with brain
546 morphology has been frequently examined in previous studies on healthy adults (Fjell et al.,
547 2014; Kharabian Masouleh et al., 2016; Salat et al., 2004; Willette and Kapogiannis, 2014).

548 Accordingly, they served here as the initial benchmarks among which SBB framework was
549 tested in healthy individuals.

550 *Psychological measurements:*

551 The psychological measurements consisted in standard psychometrics and neuropsychological
552 tests. The testing included: the attention network task (ANT) probing attention sub-functions
553 (Fan et al., 2002), the Delis-Kaplan testing battery assessing different aspects of executive
554 functions (Delis et al., 2001) (including trail-making test, color-word interference task, verbal
555 fluency, 20 questions, proverbs and word-context task) , the Rey auditory verbal learning task
556 (RAVLT) (Schmidt, 1996) assessing verbal memory performance, as well as the WASI-II
557 intelligence test (Wechsler, 1999). Psychological phenotyping also included anxiety (state and
558 trait) (Spielberger et al., 1970) and personality questionnaires (McCrae and Costa, 2004) in the
559 eNKI cohort. For each test, we used several commonly derived sub-scores to assess the
560 replicability of their structural associations. For each psychological measure, participants
561 whose performance deviated more than 3 standard deviation (SD) from mean of the whole
562 sample were considered as outliers and thus were excluded from further analysis (See
563 supplementary Table 1).

564 The list-learning task is a common measure of verbal learning performance and has been
565 implemented using the same standard tool (RAVLT) in both the eNKI and the ADNI cohort.
566 Previous studies have shown that the immediate-recall score (sum of recalled items over the
567 first 5 trials) could be reliably predicted from whole brain MRIs in AD patients (Moradi et al.,
568 2017). Since this score is a standard measure commonly used in healthy and clinical dataset

569 and its relations to brain structure in clinical data has been previously suggested, in the current
570 work we performed SBB with this score in the ADNI cohort as a “conceptual benchmark”.

571 ***MRI acquisition and preprocessing:***

572 The imaging data of the eNKI cohort were all acquired using a single scanner (Siemens
573 Magnetom TrioTim, 3.0 T). T1-weighted images were obtained using a MPRAGE sequence
574 (TR = 1900 ms; TE = 2.52 ms; voxel size = 1 mm isotropic).

575 ADNI, on the other hand, is a multisite dataset. Here we selected a subset of this data, which
576 has been acquired in a 3.0 T scanner (baseline measurements from ADNI2 and ADNI GO
577 cohort) from 39 different sites; see <http://adni.loni.usc.edu/methods/documents/> for more
578 information.

579 Both datasets were preprocessed using the CAT12 toolbox (Gaser and Dahnke, 2016). Briefly,
580 each participant’s T1-weighted scan was corrected for bias-field inhomogeneities, then
581 segmented into gray matter (GM), white matter (WM), and cerebrospinal fluid (CSF)
582 (Ashburner and Friston, 2005). The segmentation process was further extended for accounting
583 for partial volume effects (Tohka et al., 2004) by applying adaptive maximum a posteriori
584 estimations (Rajapakse et al., 1997). The gray matter segments were then spatially normalized
585 into standard (MNI) space using Dartel algorithm (Ashburner, 2007) and further modulated (for
586 non-linear transformations only) to preserve the total volume after spatial normalization.
587 Finally gray matter images were smoothed with an isotropic gaussian kernel of 8 mm (full-
588 width-half-maximum).

589

590 ***Statistical analysis:***

591 SBB-associations are commonly derived in an exploratory setting using a mass-univariate
592 approach, in which a linear model is used to fit interindividual variability in the psychological
593 score to GMV at each voxel. Inference is then usually made at cluster level, in which groups of
594 adjacent voxels that support the link between GMV and the tested score are clustered together.

595 Replicability of thus-defined associations could be assessed by conducting a similar whole-
596 brain voxel-wise exploratory analysis in another sample of individuals and comparing the
597 spatial location of the significant findings that survive multiple comparison correction, between
598 the two samples. Alternatively, replicability could be assessed, using a confirmatory approach,
599 in which only regions showing significant SBB-association in the initial exploratory analysis,
600 i.e. regions of interest (ROIs), are considered for testing the existence of the association between
601 brain structure and the same psychological score in an independent sample. The latter procedure
602 commonly focuses on a summary measure of GMV within each ROI and tests for existence of
603 the SBB-association in the direction suggested by the initial exploratory analysis. Thus this
604 approach circumvents the need for multiple comparison correction and therefore increases the
605 power of replication.

606 Here we assessed replicability of associations between each behavioral measure and gray mater
607 structure, using both approaches: the whole brain replication approach and the ROI replication
608 approach, which are explained in details in the following sections.

609

610 *Replicability of whole brain exploratory SBB-associations:*

611

612 Whole-brain GLM analyses: 100 random subsamples (of same size) were drawn from the main
613 cohort (eNKI or ADNI). Hereafter, each of these subsamples is called a “discovery sample”. In
614 each of these samples, SBB-associations were identified using the voxel-wise exploratory
615 approach after controlling for confounders. This was done by using the general linear model
616 (GLM) as implemented in the “randomise” tool
617 (<https://fsl.fmrib.ox.ac.uk/fsl/fslwiki/Randomise>), with 1000 permutations. Age, sex and
618 education were modeled as confounders in the eNKI data. As the ADNI dataset is a multi-site
619 study, we further added site and disease category as dummy-coded confounders to GLMs for
620 the analyses in that dataset. Inference was then made using threshold-free cluster enhancement

621 (TFCE) (Smith and Nichols, 2009), which unlike other cluster-based thresholding approaches
622 does not require an arbitrary a-priori cluster forming threshold. Significance was set at $P < 0.05$
623 (extent threshold of 100 voxels).

624 Spatial consistency maps and density plots: To quantify the spatial overlap of significant SBB
625 associations over 100 subsamples, spatial consistency maps were generated. To do so, the
626 binarized maps of all clusters that showed significant association in the same direction between
627 each psychological score and GMV were generated (i.e. voxels belonging to a significant
628 cluster get the value “1” and all other voxels were labeled “0”) and added over all 100
629 subsamples. These aggregate maps denote the frequency of finding a *significant* association
630 between the behavioral score and GMV, at each voxel. Accordingly, a voxel with value of 10
631 in the aggregate map has been found to be significantly associated with the phenotypical score
632 in 10 out of 100 subsamples. Density plots were also generated to represent the distribution of
633 values within each such map, i.e. the distribution of “frequency of significant finding”. Hence,
634 the spatial voxel-wise “significance overlap maps” as well as density plots of the distribution
635 of values within each map give indications of the replicability of “whole brain exploratory SBB-
636 associations” for each psychological score.

637

638 *Replicability of SBB-associations using confirmatory ROI-based approach:*

639 ROI-based confirmatory analyses: The replicability of the SBB associations was also evaluated
640 with the ROI-based confirmatory approach. For each of the 100 discovery subsamples, an age-
641 and sex-matched “test sample” was generated from the remaining participants of the main
642 cohort. In the clinical cohort the discovery and test pairs were additionally matched for “site”.
643 In this analysis, for each psychological variable, the significant clusters from the above-
644 mentioned exploratory approach from every “discovery sample” were used as a-priori ROIs.
645 Average GMV over all voxels within the ROI was then calculated for each participant in the
646 respective “discovery” and “test” pair subsamples. Within each subsample, association between

647 the average GMV and the psychological variable was assessed using ranked-partial correlation,
648 controlling for confounding factors. The correlation coefficient was then compared between
649 each discovery and test pair, providing means to assess “ROI-based SBB replicability” rates
650 for each psychological score. Accordingly, each ROI was examined only once, to identify if
651 associations between average GMV in this ROI and the psychological score from the discovery
652 subsample could be confirmed in the paired test sample. Replicability rates were quantified
653 according to different indexes (see below) over all ROIs from 100 discovery samples, yielding
654 a percentage of “successfully replicated” ROIs based on each index.

655 Indexes of replicability:

656 **Sign:** First, we used a lenient definition of replication, in which we compared only the sign of
657 correlation coefficients of associations within each ROI between the discovery and the
658 matched-test sample. Accordingly, any effect that was in the same direction in both samples
659 (even if very close to zero) was defined as a “successful” replication.

660 **Statistical Significance:** Another straightforward method for evaluating replication simply
661 defines statistically significant effects (e.g. p -value < 0.05) that are in the same direction as the
662 original effects (from the discovery sample) as “successful” replication. This criteria is
663 consistent with what is commonly used in the psychological sciences to decide whether a
664 replication attempt “worked” (Open Science Collaboration, 2015). Yet, a key weakness of this
665 approach is that it treats the threshold ($p < 0.05$) as a bright-line criterion between replication
666 success and failure. Furthermore, it does not quantify the decisiveness of the evidence that the
667 data provides for and against the presence of the correlation (Boekel et al., 2015; Wagenmakers
668 et al., 2015). However, such an estimation can be provided by using the “Bayes factors”.

669 **Bayes Factor:** To compare the evidence that the “test subsample” provided for or against the
670 presence of an association (H_1 and H_0 , respectively), we additionally quantified SBB-
671 replication within each ROI, using Bayes factors (Jeffreys, 1961). Similar to Boekel et al.
672 (2015), here we used the adjusted (one-sided) Jeffry’s test (Jeffreys, 1961) based on a uniform

673 prior distribution for the correlation coefficient. As we intended to confirm the SBB-
674 associations defined in the discovery subsamples, the alternative hypothesis (H1) in this study
675 was considered one-sided (in line with Boekel et al. (2015)). We used implementation of the
676 Bayes Factors for correlations from the R function available at
677 http://www.josineverhagen.com/?page_id=76.

678 To facilitate the interpretation, Bayes factors (BF) were summarized into four categories as
679 illustrated in the bar legend of Figure 2. A BF_{01} lower than 1/3 shows that the data is three times
680 or more likely to have happened under H1 than H0. Accordingly, this value defines the
681 “successful” replication.

682 *Investigation on factors influencing replicability of SBB-associations among healthy*
683 *individuals:*

684 Sample size: In order to study the influence of sample size on the replicability of SBB-
685 associations, for each psychological measure, the healthy sample (eNKI) was divided into
686 discovery and test pairs at three different ratios: 70% discovery and 30% test, 50% discovery
687 and 50% test and finally 30% discovery and 70% test. As mentioned earlier, in each case, the
688 discovery and test counterparts were randomly generated 100 times in order to quantify the
689 replication rates. For example, to assess the replicability of brain structural associations of age,
690 in the case of “70% discovery and 30% test”, the entire NKI sample ($n = 466$) was divided into
691 a discovery group of $n = 326$ participants and an age- and sex-matched test pair sample of $n =$
692 138 and this split procedure was repeated 100 times. Similarly, for generating equal-sized
693 discovery and test subsamples, 100 randomly generated age and sex matched split-half samples
694 were generated from the main NKI cohort.

695 Due to the multi-site structure of the ADNI cohort, when generating unequal sized discovery
696 and test samples, we did not achieve a good simultaneous matching of age, sex and site, while
697 trying to maintain samples sizes in each subgroup reasonably large. Thus, in this cohort, we did
698 not directly study the influence of the sample size and the replicability rates were only

699 quantified for equal sized discovery and test samples (187 participants matched for age, sex and
700 site between discovery and test pairs).

701 Effect size: Furthermore, to study the influence of the effect size on the replication rates, we
702 focused on the effect sizes within each a-priori ROI in the discovery samples. Here we tested
703 the following two assumptions:

704 1) ROIs with larger effect sizes in the discovery sample result in larger effect sizes in the test
705 sample pairs (i.e. positive association between effect size in the discovery and test samples).

706 2) ROIs with larger effect sizes in the discovery sample are more likely to result in a
707 “significant” replication in the independent sample.

708 To test the first assumption, in the “ROI-based SBB-replicability” the association between
709 effect size in the discovery and test pairs were calculated for each psychological measure. These
710 associations were calculated separately for the replicated (defined using “sign” criterion) and
711 not-replicated ROIs. We expected to find a positive association between discovery and
712 confirmatory effect sizes, for the “successfully replicated effects”.

713 To test the second assumption, for each ROI, we calculated its replication statistical power and
714 compared it between replicated and not-replicated ROIs (here replication was defined using
715 “Statistical Significance” criterion). The statistical power of a test is the probability that it will
716 correctly reject the null hypothesis when the null is false. In a bias-free case, the power of the
717 replication is a function of the replication sample size, real size of the effect and the nominal
718 type I error rate (α). In this study, the replication power was estimated based on the size of the
719 effects as they were defined in the discovery sample and a significant threshold of 0.05 (one-
720 sided) and was calculated using “pwr” library in R (<https://www.r-project.org>).

721 These analyses were performed for each discovery-test split size, separately (i.e. 70%-30%,
722 50%-50% and 30%-70% discovery-test sample sizes, respectively).

723

724 **Acknowledgement:**

725 This work was supported by the Deutsche Forschungsgemeinschaft (DFG, GE 2835/1-1, EI
726 816/4-1), the Helmholtz Portfolio Theme ‘Supercomputing and Modelling for the Human
727 Brain’ and the European Union’s Horizon 2020 Research and Innovation Programme under
728 Grant Agreement No. 720270 (HBP SGA1) and Grant Agreement No. 785907 (HBP SGA2).
729 Clinical data collection and sharing for this project was funded by the Alzheimer’s Disease
730 Neuroimaging Initiative (ADNI) (National Institutes of Health Grant U01 AG024904) and
731 DOD ADNI (Department of Defense award number W81XWH-12-2-0012). ADNI is funded
732 by the National Institute on Aging, the National Institute of Biomedical Imaging and
733 Bioengineering, and through generous contributions from the following: AbbVie, Alzheimer’s
734 Association; Alzheimer’s Drug Discovery Foundation; Araclon Biotech; BioClinica, Inc.;
735 Biogen; Bristol-Myers Squibb Company; CereSpir, Inc.; Cogstate; Eisai Inc.; Elan
736 Pharmaceuticals, Inc.; Eli Lilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and its
737 affiliated company Genentech, Inc.; Fujirebio; GE Healthcare; IXICO Ltd.; Janssen Alzheimer
738 Immunotherapy Research & Development, LLC.; Johnson & Johnson Pharmaceutical Research
739 & Development LLC.; Lumosity; Lundbeck; Merck & Co., Inc.; Meso Scale Diagnostics,
740 LLC.; NeuroRx Research; Neurotrack Technologies; Novartis Pharmaceuticals Corporation;
741 Pfizer Inc.; Piramal Imaging; Servier; Takeda Pharmaceutical Company; and Transition
742 Therapeutics. The Canadian Institutes of Health Research is providing funds to support ADNI
743 clinical sites in Canada. Private sector contributions are facilitated by the Foundation for the
744 National Institutes of Health (www.fnih.org). The grantee organization is the Northern
745 California Institute for Research and Education, and the study is coordinated by the Alzheimer’s
746 Therapeutic Research Institute at the University of Southern California. ADNI data are
747 disseminated by the Laboratory for Neuro Imaging at the University of Southern California.

748 **Competing interests:** The authors declare no competing interests.

749

750 **References:**

- 751 Albers C, Lakens D. 2018a. When power analyses based on pilot data are biased: Inaccurate
752 effect size estimators and follow-up bias. *J Exp Soc Psychol* **74**:187–195.
753 doi:10.1016/j.jesp.2017.09.004
- 754 Albers C, Lakens D. 2018b. When power analyses based on pilot data are biased: Inaccurate
755 effect size estimators and follow-up bias. *J Exp Soc Psychol* **74**:187–195.
756 doi:10.1016/j.jesp.2017.09.004
- 757 Anderson ML. 2015. Précis of after Phrenology: Neural Reuse and the Interactive Brain.
758 *Behav Brain Sci*. doi:10.1017/S0140525X15000631
- 759 Ashburner J. 2007. A fast diffeomorphic image registration algorithm. *Neuroimage* **38**:95–
760 113. doi:10.1016/j.neuroimage.2007.07.007
- 761 Ashburner J, Friston KJ. 2005. Unified segmentation. *Neuroimage* **26**:839–851.
762 doi:10.1016/j.neuroimage.2005.02.018
- 763 Boekel W, Wagenmakers EJ, Belay L, Verhagen J, Brown S, Forstmann BU. 2015. A purely
764 confirmatory replication study of structural brain-behavior correlations. *Cortex* **66**:115–
765 133. doi:10.1016/j.cortex.2014.11.019
- 766 Bressler SL. 1995. Large-scale cortical networks and cognition. *Brain Res Rev*.
767 doi:10.1016/0165-0173(94)00016-I
- 768 Broca P. 1865. Sur le siège de la faculté du langage articulé. *Bull la Société d'anthropologie*
769 *Paris* **6**:377–393. doi:10.3406/bmsap.1865.9495
- 770 Button KS, Ioannidis JPA, Mokrysz C, Nosek BA, Flint J, Robinson ESJ, Munafò MR. 2013.
771 Power failure: why small sample size undermines the reliability of neuroscience. *Nat Rev*
772 *Neurosci* **14**:365–76. doi:10.1038/nrn3475
- 773 Cremers HR, Wager TD, Yarkoni T. 2017. The relation between statistical power and
774 inference in fMRI. *PLoS One* **12**:1–20. doi:10.1371/journal.pone.0184923
- 775 De Boeck P, Jeon M. 2018. Perceived Crisis and Reforms: Issues, Explanations, and
776 Remedies. *Psychol Bull* **144**:757–777. doi:10.1037/bul0000154
- 777 Delis DC, Kaplan E, Kramer JH. 2001. Delis-Kaplan Executive Function System (D-KEFS)
778 examiner's manual. San Antonio, TX: The Psychological Corporation.
- 779 Draganski B, Gaser C, Busch V, Schuierer G, Bogdahn U, May A. 2004. Changes in grey
780 matter induced by training Newly honed juggling skills show up as a transient feature on
781 a brain-imaging scan. *Nature* **427**:311–312. doi:10.1038/427311a
- 782 Dwan K, Gamble C, Williamson PR, Kirkham JJ. 2013. Systematic Review of the Empirical
783 Evidence of Study Publication Bias and Outcome Reporting Bias - An Updated Review.
784 *PLoS One*. doi:10.1371/journal.pone.0066844
- 785 Eklund A, Nichols TE, Knutsson H. 2016. Cluster failure: Why fMRI inferences for spatial
786 extent have inflated false-positive rates. *Proc Natl Acad Sci U S A* **113**:7900–5.
787 doi:10.1073/pnas.1602413113
- 788 Fan J, McCandliss BD, Sommer T, Raz A, Posner MI. 2002. Testing the Efficiency and
789 Independence of Attentional Networks. *J Cogn Neurosci* **14**:340–347.
- 790 Fjell AM, Westlye LT, Grydeland H, Amlie I, Espeseth T, Reinvang I, Raz N, Dale AM,
791 Walhovd KB. 2014. Accelerating cortical thinning: unique to dementia or universal in
792 aging? *Cereb Cortex* **24**:919–34. doi:10.1093/cercor/bhs379
- 793 Forstmeier W, Schielzeth H. 2011. Cryptic multiple hypotheses testing in linear models:
794 Overestimated effect sizes and the winner's curse. *Behav Ecol Sociobiol* **65**:47–55.
795 doi:10.1007/s00265-010-1038-5
- 796 Gaser C, Dahnke R. 2016. CAT - A Computational Anatomy Toolbox for the Analysis of
797 Structural MRI Data. *HBM Conf 2016* **32**:7743.
- 798 Gelman A, Loken E. 2014. The garden of forking paths: Why multiple comparisons can be a
799 problem, even when there is no “fishing expedition” or “p-hacking” and the research

- 800 hypothesis was posited ahead of time. *Psychol Bull* **140**:1272–1280.
801 doi:[dx.doi.org/10.1037/a0037714](https://doi.org/10.1037/a0037714)
- 802 Genon S, Bastin C, Angel L, Collette F, Bahri MA, Salmon E. 2014. A partial least squares
803 analysis of the self reference effect in Alzheimer’s disease: A reply to Irish. *Cortex*.
804 doi:[10.1016/j.cortex.2014.02.003](https://doi.org/10.1016/j.cortex.2014.02.003)
- 805 Genon S, Reid A, Langner R, Amunts K, Eickhoff SB. 2018. How to Characterize the
806 Function of a Brain Region. *Trends Cogn Sci*. doi:[10.1016/j.tics.2018.01.010](https://doi.org/10.1016/j.tics.2018.01.010)
- 807 Genon S, Simon J, Bahri MA, Collette F, Souchay C, Jaspar M, Bastin C, Salmon E. 2016.
808 Relating pessimistic memory predictions to Alzheimer’s disease brain structure. *Cortex*
809 **85**:151–164. doi:[10.1016/j.cortex.2016.09.014](https://doi.org/10.1016/j.cortex.2016.09.014)
- 810 Genon S, Wensing T, Reid A, Hoffstaedter F, Caspers S, Grefkes C, Nickl-Jockschat T,
811 Eickhoff SB. 2017. Searching for behavior relating to grey matter volume in a-priori
812 defined right dorsal premotor regions: Lessons learned. *Neuroimage* **157**:144–156.
813 doi:[10.1016/j.neuroimage.2017.05.053](https://doi.org/10.1016/j.neuroimage.2017.05.053)
- 814 Gorgolewski KJ, Varoquaux G, Rivera G, Schwarz Y, Ghosh SS, Maumet C, Sochat V V.,
815 Nichols TE, Poldrack RA, Poline J-B, Yarkoni T, Margulies DS. 2015. NeuroVault.org:
816 a web-based repository for collecting and sharing unthresholded statistical maps of the
817 human brain. *Front Neuroinform* **9**:8. doi:[10.3389/fninf.2015.00008](https://doi.org/10.3389/fninf.2015.00008)
- 818 Habeck C, Stern Y. 2010. Multivariate data analysis for neuroimaging data: overview and
819 Application to Alzheimer’s disease. *Cell Biochem Biophys* **58**:53–67.
820 doi:[10.1007/s12013-010-9093-0](https://doi.org/10.1007/s12013-010-9093-0).Multivariate
- 821 Ioannidis JPA. 2018. Why replication has more scientific value than original discovery.
822 *Behav Brain Sci* **41**:e137. doi:[10.1017/S0140525X18000729](https://doi.org/10.1017/S0140525X18000729)
- 823 Ioannidis JPA. 2008. Why most discovered true associations are inflated. *Epidemiology*
824 **19**:640–648. doi:[10.1097/EDE.0b013e31818131e7](https://doi.org/10.1097/EDE.0b013e31818131e7)
- 825 Ioannidis JPA. 2005. Why most published research findings are false. *PLoS Med*.
826 doi:[10.1371/journal.pmed.0020124](https://doi.org/10.1371/journal.pmed.0020124)
- 827 Ismaylova E, Di Sante J, Gouin J-P, Pomares FB, Vitaro F, Tremblay RE, Booi L. 2018.
828 Associations Between Daily Mood States and Brain Gray Matter Volume, Resting-State
829 Functional Connectivity and Task-Based Activity in Healthy Adults. *Front Hum*
830 *Neurosci* **12**:168. doi:[10.3389/fnhum.2018.00168](https://doi.org/10.3389/fnhum.2018.00168)
- 831 Jeffreys H. 1961. Theory of probability. Oxford, Uk.: Oxford University Press.
- 832 John LK, Loewenstein G, Prelec D. 2012. Measuring the prevalence of questionable research
833 practices with incentives for truth telling. *Psychol Sci* **23**:524–32.
834 doi:[10.1177/0956797611430953](https://doi.org/10.1177/0956797611430953)
- 835 Kanai R. 2016. Open questions in conducting confirmatory replication studies: Commentary
836 on Boekel et al., 2015. *Cortex*. doi:[10.1016/j.cortex.2015.02.020](https://doi.org/10.1016/j.cortex.2015.02.020)
- 837 Kanai R, Feilden T, Firth C, Rees G. 2011. Political orientations are correlated with brain
838 structure in young adults. *Curr Biol* **21**:677–680. doi:[10.1016/j.cub.2011.03.017](https://doi.org/10.1016/j.cub.2011.03.017)
- 839 Kanai R, Rees G. 2011. The structural basis of inter-individual differences in human
840 behaviour and cognition. *Nat Rev Neurosci* **12**:231–242. doi:[10.1038/nrn3000](https://doi.org/10.1038/nrn3000)
- 841 Kharabian Masouleh S, Arélin K, Horstmann A, Lampe L, Kipping JA, Luck T, Riedel-Heller
842 SG, Schroeter ML, Stumvoll M, Villringer A, Witte AV. 2016. Higher body mass index
843 in older adults is associated with lower gray matter volume: Implications for memory
844 performance, Neurobiology of Aging. Elsevier Ltd.
845 doi:[10.1016/j.neurobiolaging.2015.12.020](https://doi.org/10.1016/j.neurobiolaging.2015.12.020)
- 846 Kim EJ, Pellman B, Kim JJ. 2015. Stress effects on the hippocampus: A critical review.
847 *Learn Mem*. doi:[10.1101/lm.037291.114](https://doi.org/10.1101/lm.037291.114)
- 848 Kriegeskorte N, Lindquist MA, Nichols TE, Poldrack RA, Vul E. 2010. Everything you never
849 wanted to know about circular analysis, but were afraid to ask. *J Cereb Blood Flow*
850 *Metab*. doi:[10.1038/jcbfm.2010.86](https://doi.org/10.1038/jcbfm.2010.86)

- 851 Lilienfeld SO. 2017. Psychology's Replication Crisis and the Grant Culture: Righting the
852 Ship. *Perspect Psychol Sci* **12**:660–664. doi:10.1177/1745691616687745
- 853 Loken E, Gelman A. 2017. Measurement error and the replication crisis. *Science* (80-).
854 doi:10.1126/science.aal3618
- 855 Luders E, Kurth F, Mayer EA, Toga AW, Narr KL, Gaser C. 2012. The Unique Brain
856 Anatomy of Meditation Practitioners: Alterations in Cortical Gyri-fication. *Front Hum*
857 *Neurosci* **6**:34. doi:10.3389/fnhum.2012.00034
- 858 Luders E, Kurth F, Toga AW, Narr KL, Gaser C. 2013. Meditation effects within the
859 hippocampal complex revealed by voxel-based morphometry and cytoarchitectonic
860 probabilistic mapping. *Front Psychol* **4**:398. doi:10.3389/fpsyg.2013.00398
- 861 Matsuo K, Nicoletti M, Nemoto K, Hatch JP, Peluso MAM, Nery FG, Soares JC. 2009. A
862 voxel-based morphometry study of frontal gray matter correlates of impulsivity. *Hum*
863 *Brain Mapp* **30**:1188–1195. doi:10.1002/hbm.20588
- 864 McCrae RR, Costa PT. 2004. A contemplated revision of the NEO Five-Factor Inventory.
865 *Pers Individ Dif* **36**:587–596. doi:10.1016/S0191-8869(03)00118-1
- 866 McEwen BS, Nasca C, Gray JD. 2016. Stress Effects on Neuronal Structure: Hippocampus,
867 Amygdala, and Prefrontal Cortex. *Neuropsychopharmacology*.
868 doi:10.1038/npp.2015.171
- 869 McIntosh AR, Mišić B. 2013. Multivariate Statistical Analyses for Neuroimaging Data. *Annu*
870 *Rev Psychol* **64**:499–525. doi:10.1146/annurev-psych-113011-143804
- 871 Miller KL, Alfaro-Almagro F, Bangerter NK, Thomas DL, Yacoub E, Xu J, Bartsch AJ,
872 Jbabdi S, Sotiropoulos SN, Andersson JLR, Griffanti L, Douaud G, Okell TW, Weale P,
873 Dragonu I, Garratt S, Hudson S, Collins R, Jenkinson M, Matthews PM, Smith SM.
874 2016. Multimodal population brain imaging in the UK Biobank prospective
875 epidemiological study. *Nat Neurosci* **19**:1523–1536. doi:10.1038/nn.4393
- 876 Moeller JR, Habeck CG. 2006. Reciprocal benefits of mass-univariate and multivariate
877 modeling in brain mapping: Applications to event-related functional MRI, H215O-, and
878 FDG-PET. *Int J Biomed Imaging* **2006**:1–13. doi:10.1155/IJBI/2006/79862
- 879 Moradi E, Hallikainen I, Hänninen T, Tohka J. 2017. Rey's Auditory Verbal Learning Test
880 scores can be predicted from whole brain MRI in Alzheimer's disease. *NeuroImage Clin*
881 **13**:415–427. doi:10.1016/j.nicl.2016.12.011
- 882 Muhlert N, Ridgway GR. 2016. Failed replications, contributing factors and careful
883 interpretations: Commentary on Boekel et al., 2015. *Cortex*.
884 doi:10.1016/j.cortex.2015.02.019
- 885 Munafò MR, Nosek BA, Bishop DVM, Button KS, Chambers CD, Percie Du Sert N,
886 Simonsohn U, Wagenmakers EJ, Ware JJ, Ioannidis JPA. 2017. A manifesto for
887 reproducible science. *Nat Hum Behav* **1**:1–9. doi:10.1038/s41562-016-0021
- 888 Nooner KB, Colcombe SJ, Tobe RH, Mennes M, Benedict MM, Moreno AL, Panek LJ,
889 Brown S, Zavitz Stephen TT, Li Q, Sikka S, Gutman D, Bangaru S, Schlachter RT,
890 Anwar SMK, Hinz CM, Kaplan MS, Rachlin AB, Adelsberg S, Cheung B, Khanuja R,
891 Yan C, Courtney CCC, King M, Wood D, Cox CL, Kelly AMC, Petkova E, Reiss PT,
892 Duan N, Thomsen D, Biswal B, Coffey B, Hoptman MJ, Javitt DC, Pomara N, Sidtis JJ,
893 Koplewicz HS, Castellanos FX, Leventhal BL, Milham MP. 2012. The NKI-Rockland
894 sample: A model for accelerating the pace of discovery science in psychiatry. *Front*
895 *Neurosci*. doi:10.3389/fnins.2012.00152
- 896 Nostro AD, Müller VI, Reid AT, Eickhoff SB. 2017. Correlations between Personality and
897 Brain Structure: A Crucial Role of Gender. *Cereb Cortex* **27**:3698–3712.
898 doi:10.1093/cercor/bhw191
- 899 Open Science Collaboration OS. 2015. Estimating the reproducibility of psychological
900 science. *Science* **349**:aac4716. doi:10.1126/science.aac4716
- 901 Pessoa L. 2014. Understanding brain networks and brain organization. *Phys Life Rev.*

- 902 doi:10.1016/j.plrev.2014.03.005
903 Poldrack RA. 2010. Mapping mental function to brain Structure: How can cognitive
904 Neuroimaging Succeed? *Perspect Psychol Sci* **5**:753–761.
905 doi:10.1177/1745691610388777
906 Poldrack RA, Baker CI, Durnez J, Gorgolewski KJ, Matthews PM, Munafò MR, Nichols TE,
907 Poline J-B, Vul E, Yarkoni T. 2017. Scanning the horizon: towards transparent and
908 reproducible neuroimaging research. *Nat Rev Neurosci* **18**:115–126.
909 doi:10.1038/nrn.2016.167
910 Rajapakse JC, Giedd JN, Rapoport JL. 1997. Statistical approach to segmentation of single-
911 channel cerebral mr images. *IEEE Trans Med Imaging* **16**:176–186.
912 doi:10.1109/42.563663
913 Salat DH, Buckner RL, Snyder AZ, Greve DN, Desikan RSR, Busa E, Morris JC, Dale AM,
914 Fischl B. 2004. Thinning of the cerebral cortex in aging. *Cereb Cortex* **14**:721–30.
915 doi:10.1093/cercor/bhh032
916 Schmidt M. 1996. RAVLT (Rey Auditory Verbal Learning Test: A Handbook).
917 Schönbrodt FD, Perugini M. 2013. At what sample size do correlations stabilize? *J Res Pers*
918 **47**:609–612. doi:10.1016/j.jrp.2013.05.009
919 Scoville WB, Milner B. 1957. Loss of recent memory after bilateral hippocampal lesions. *J*
920 *Neurol Neurosurg Psychiatry* **20**:11–21. doi:10.1136/jnnp-2015-311092
921 Simmons JP, Nelson LD, Simonsohn U. 2011. False-positive psychology: Undisclosed
922 flexibility in data collection and analysis allows presenting anything as significant.
923 *Psychol Sci* **22**:1359–1366. doi:10.1177/0956797611417632
924 Smith SM, Nichols TE. 2009. Threshold-free cluster enhancement: Addressing problems of
925 smoothing, threshold dependence and localisation in cluster inference. *Neuroimage*
926 **44**:83–98. doi:10.1016/j.neuroimage.2008.03.061
927 Spielberger CD, Gorsuch RL, Lushene RE. 1970. Manual for the State- Trait Anxiety
928 Inventory. Palo Alto, CA: Consulting Psychologists Press.
929 Taubert M, Lohmann G, Margulies DS, Villringer A, Ragert P. 2011. Long-term effects of
930 motor training on resting-state networks and underlying brain structure. *Neuroimage*
931 **57**:1492–1498. doi:10.1016/j.neuroimage.2011.05.078
932 Tohka J, Zijdenbos A, Evans A. 2004. Fast and robust parameter estimation for statistical
933 partial volume models in brain MRI. *Neuroimage* **23**:84–97.
934 doi:10.1016/j.neuroimage.2004.05.007
935 Tononi G, Edelman GM, Sporns O. 1998. Complexity and coherency: Integrating information
936 in the brain. *Trends Cogn Sci*. doi:10.1016/S1364-6613(98)01259-5
937 Tukey JW. 1980. We need both Eploratory and Confirmatory We Need Both Exploratory. *Am*
938 *Stat* **34**:23–25. doi:10.2307/2682991
939 Turner BO, Paul EJ, Miller MB, Barbey AK. 2018. Small sample sizes reduce the
940 replicability of task-based fMRI studies. *Commun Biol* **1**:62. doi:10.1038/s42003-018-
941 0073-z
942 Van Essen DC, Smith SM, Barch DM, Behrens TEJ, Yacoub E, Ugurbil K. 2013. The WU-
943 Minn Human Connectome Project: An overview. *Neuroimage* **80**:62–79.
944 doi:10.1016/j.neuroimage.2013.05.041
945 Vanasse TJ, Fox PM, Barron DS, Robertson M, Eickhoff SB, Lancaster JL, Fox PT. 2018.
946 BrainMap VBM: An environment for structural meta-analysis. *Hum Brain Mapp* 1–18.
947 doi:10.1002/hbm.24078
948 Wagenmakers E-J, Verhagen J, Ly A, Bakker M, Lee MD, Matzke D, Rouder JN, Morey RD.
949 2015. A power fallacy. *Behav Res Methods* **47**:913–917. doi:10.3758/s13428-014-0517-
950 4
951 Wechsler D. 1999. Wechsler Abbreviated Scale of Intelligence. San Antonio, TX: The
952 Psychologica Corporation.

- 953 Willette A a, Kapogiannis D. 2014. Does the brain shrink as the waist expands? *Ageing Res*
954 *Rev* 1–12. doi:10.1016/j.arr.2014.03.007
955 Yarkoni T. 2009. Big Correlations in Little Studies. *Perspect Psychol Sci* 4:294–298.
956 doi:10.1111/j.1745-6924.2009.01127.x
957

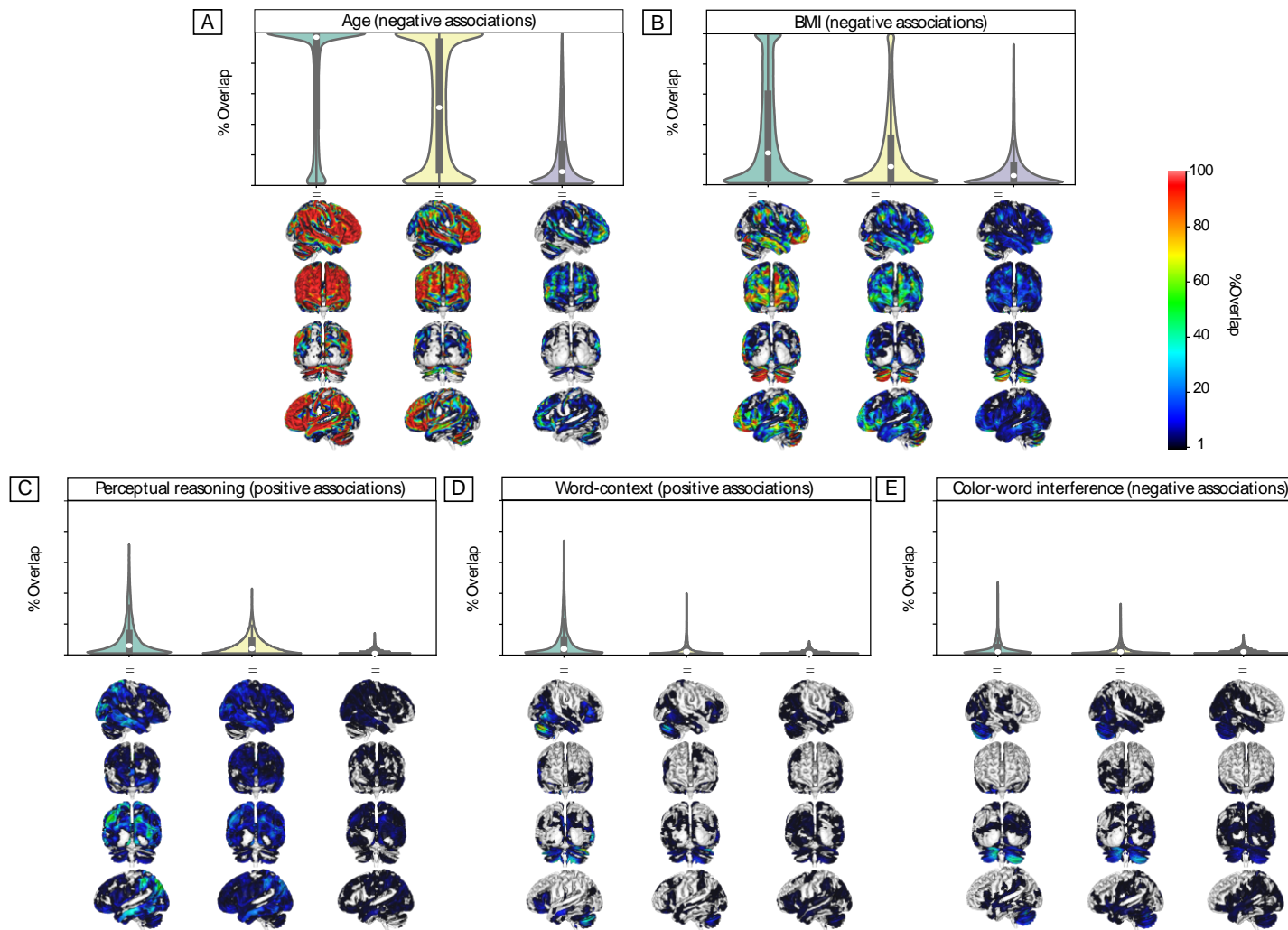


Figure 1. Replicability of exploratory results within healthy cohort. Frequency of spatial overlap (density plots and aggregate maps) of significant findings from exploratory analysis over 100 random subsamples, calculated for three different sample sizes (x-axis). Here in addition to age and BMI (A,B), which are used as benchmarks, the top three behavioral scores with the highest frequency of overlapping findings are depicted (C-E). Warmer colors on spatial maps denote higher number of samples with a significant association at the respective voxel. BMI : body mass index; CWI : color-word interference.

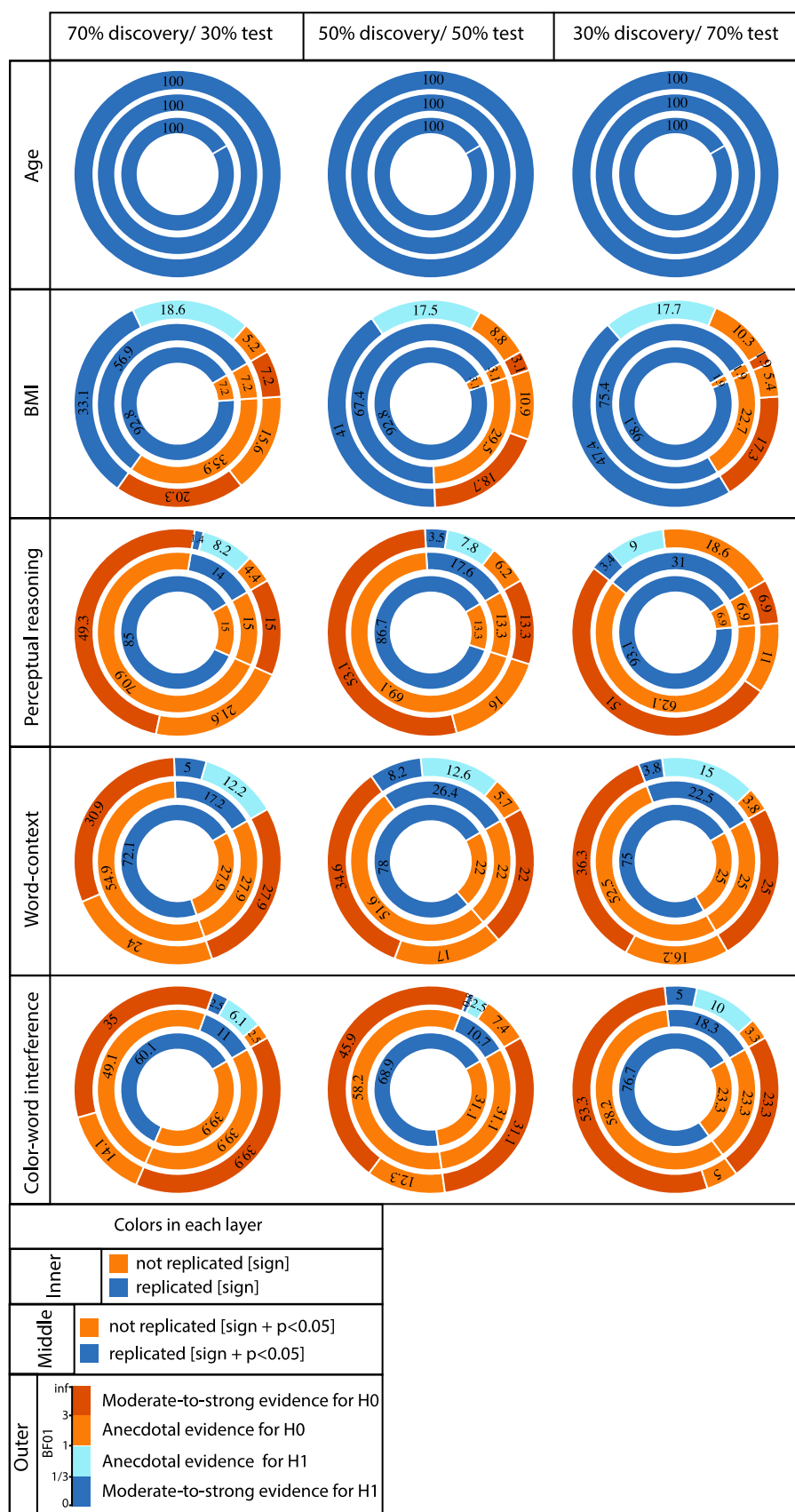


Figure 2. ROI-based confirmatory replication results within healthy cohort. Donut plots summarising ROI-based replication rates (% of ROI) using three different criteria for three different sample sizes among healthy participants. The most inner layers depict replication using “sign” only (blue: replicated, orange: not replicated). The middle layers define replication based on similar “sign” as well as “statistical significance” (i.e. $p < 0.05$) (blue: replicated, orange: not replicated). The most outer layers define replication using “bayes factor” (blue: “moderate-to-string evidence for H1, light blue: anecdotal evidence for H1; light orange: anecdotal evidence for H0, orange: “moderate-to-string evidence for H0”);

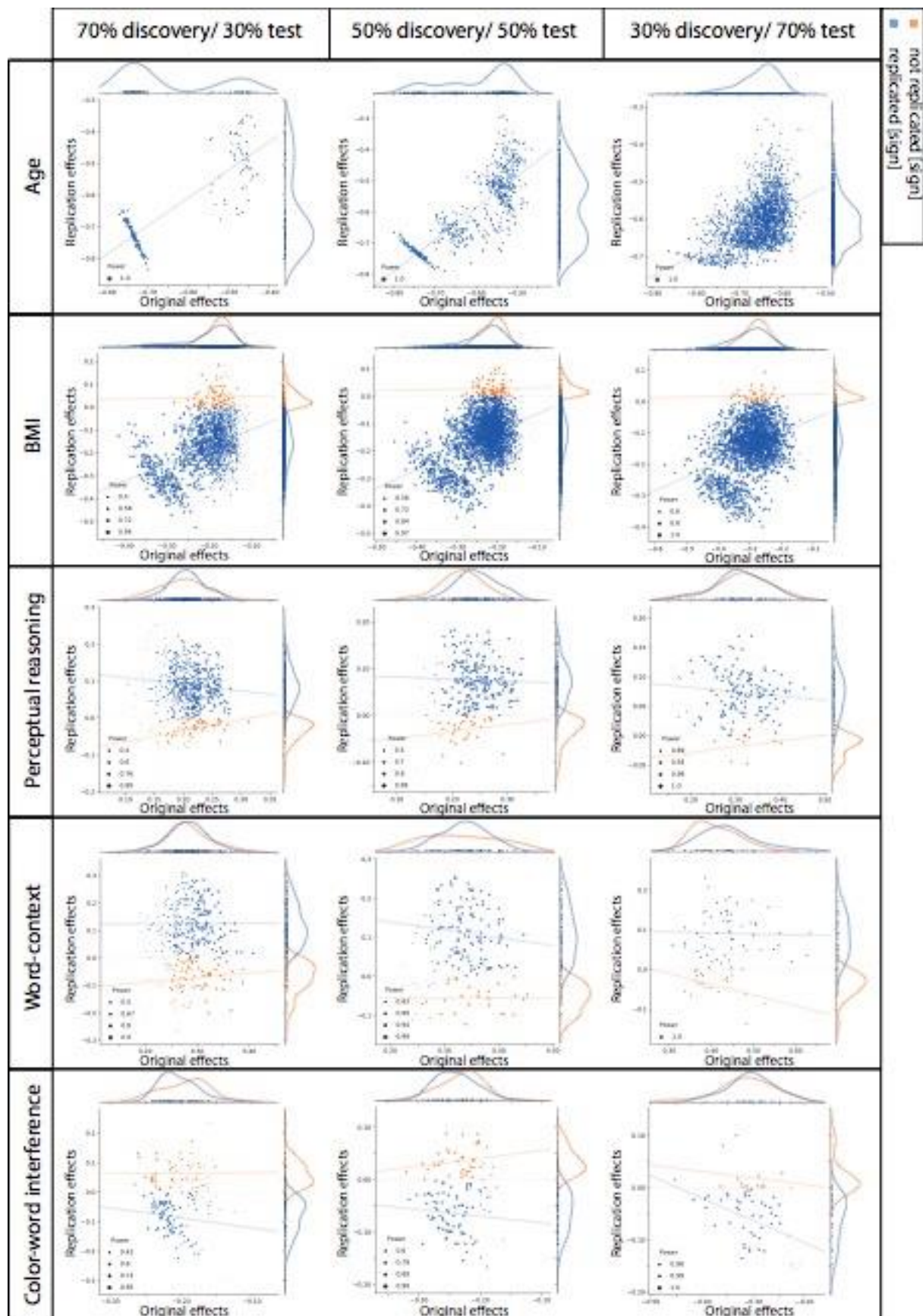


Figure 3. Discovery versus replication effects sizes: Scatter plots of effect sizes in the discovery versus replication sample for all ROIs from 100 splits within healthy cohort; each point denote one ROI, which is color-coded based on its replciation status (by-“sign”). Size of each point is proportional to its estimated statistical power of replication. Regression lines are drawn for the replciated and unreplicated ROIs, separately.

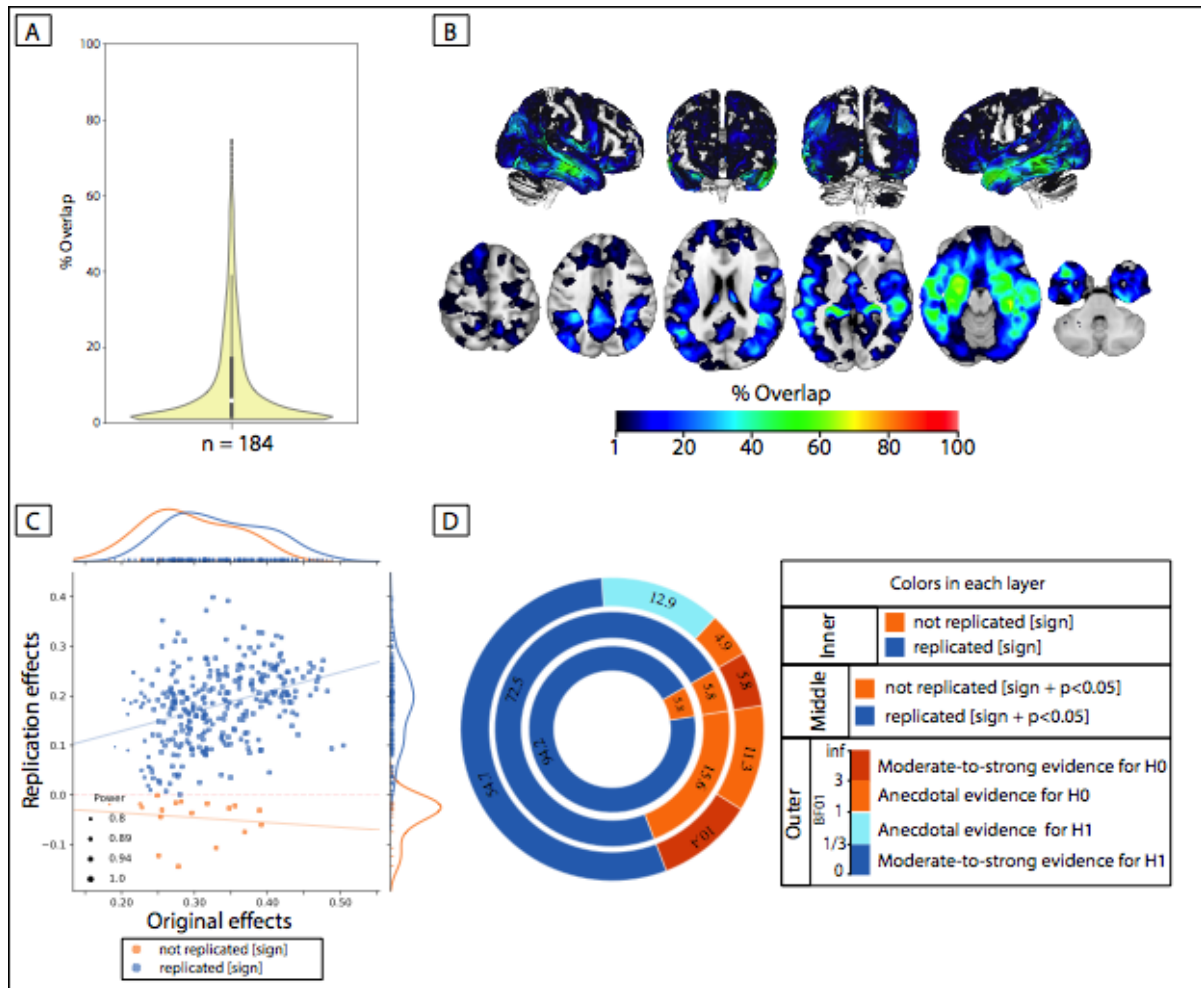


Figure 4. Replicability of positive association between immediate-recall and GMV within ADNI cohort. A, B: Replicability of exploratory results: Frequency of spatial overlaps (density plot and aggregate maps) over 100 random subsamples. C, D: ROI-based confirmatory replication results: C: Original versus replication effects sizes for all ROIs from 100 splits; points are color-coded based on their replication status (by-“sign”) and size of each point is proportional to the estimated statistical power of replication. Regression lines are drawn for the replicated and unreplicated ROIs, separately. D: Donut plots summarising ROI-based replicability rates using three different criteria. The most inner layer depicts replicability using “sign” only (blue: replicated, orange: not replicated). The middle layer, defines replication based on similar “sign” as well as “statistical significance” (i.e. $p < 0.05$) (blue: replicated, orange: not replicated). The most outer layer defines replicability using bayes factor” (blue: “moderate-to-string evidence for H1, light blue: anecdotal evidence for H1; light orange: anecdotal evidence for H0, orange: “moderate-to-string evidence for H0); Discovery and replication samples have equal size ($n = 184$) and are matched for age, sex and site.

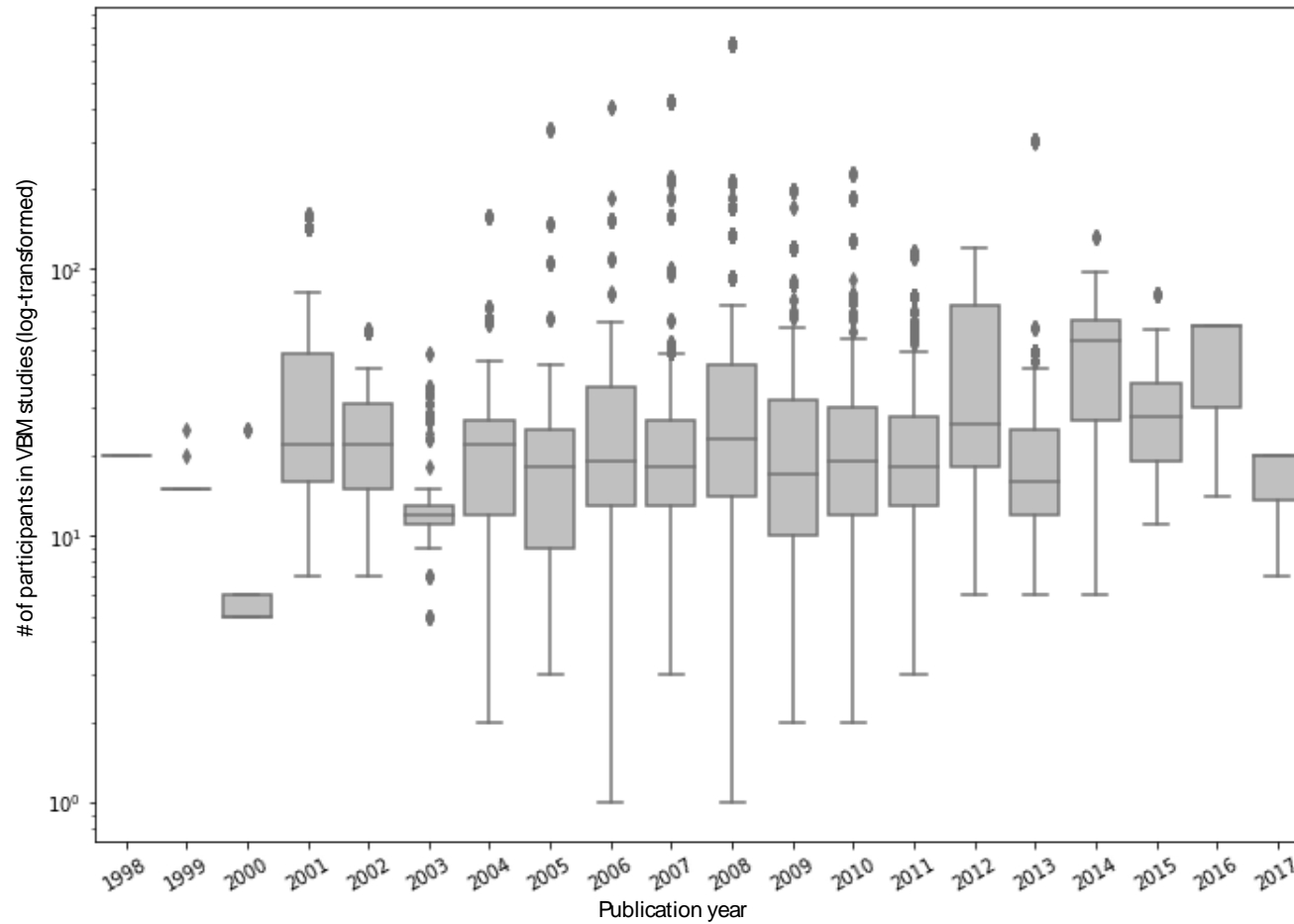


Figure 5. box-plots showing distribution of sample sizes (log-scale) of VBM studies by their publication year (data from the BrainMap database; see (Vanasse et al., 2018))

Table 1. Summary of exploratory findings. For each discovery sample size, the number of clusters in which grey matter volume is positively or negatively associated with the tested phenotypic or psychological score is reported. The number of splits (out of 100) in which the clusters were detected are noted in parentheses (i.e. % of splits with at least one significant cluster [in the respective direction])

	n_discovery = 70% n_total		n_discovery = 50% n_total		n_discovery = 30% n_total	
	# positively associated clusters (split%)	# negatively associated clusters (split%)	# positively associated clusters (split%)	# negatively associated clusters (split%)	# positively associated clusters (split%)	# negatively associated clusters (split%)
Healthy cohort						
Age (years) n-total = 466	77 (54%)	154 (100%)	5 (4%)	522 (100%)	1 (1%)	1781 (100%)
BMI (kg/m ²) n-total = 466	0	1741 (100%)	0	2276 (100%)	0	1937 (96%)
Perceptual IQ (sum of t-scores) n-total = 466	499 (83%)	0	256 (58%)	0	145 (33%)	0
Word-context (# of consecutively correct) n-total = 262	337 (80%)	0	159 (47%)	0	80 (21%)	0
CWI (interference) (sec) n-total = 449	0	163 (53%)	1 (1%)	122 (39%)	6 (1%)	60 (26%)
Clinical cohort	-		n_discovery = 50% n_total		-	
RAVLT (# total immediate recall)	-	-	309 (84%)	0	-	-

Abbreviations: BMI : body mass index; IQ : intelligence quotient, CWI: color-word interference task; RAVLT : Rey auditory verbal learning task;

Supplementary material:

Supplementary Table legends:

Table S1. Distribution of the raw phenotypical and psychological scores in the whole sample.

Table S2. Summary of the exploratory findings. For each discovery sample size, the number of clusters in which grey matter volume is positively or negatively associated with the tested psychological score is reported. Number of splits (out of 100) in which the clusters were detected are noted in parentheses.

Supplementary Figure legends:

Figure S1. Summary of replication of positive associations between immediate-recall and GMV within healthy cohort. A: Frequency of spatial overlap (density plots and aggregate maps) of significant findings from exploratory analysis over 100 random subsamples, calculated for three different sample sizes (x-axis). B: ROI-based confirmatory replication results: Top row : Donut plots summerising ROI-based replicability rates (% of ROI) using three different criteria for three different sample sizes. The most inner layers depict replicability using “sign” only (blue: replicated, orange: not replciated). The middle layers define replication based on similar “sign” as well as “statistical significance” (i.e. $p < 0.05$) (blue: replicated, orange: not replciate). The most outer layers define replicability using bayes factor ” (blue: “moderate-to-string evidence for H1, light blue: anecdotal evidence for H1; light orange: anecdotal evidence for H0, orange: “moderate-to-string evidence for H0); Bottom row: Scatter plots of effect sizes in the discovery versus replication sample for all ROIs from 100 splits within healthy cohort; Points are color-coded based on their replciation status (by-“sign”) and size of each point is proportional to the estimated statistical power of replication. Regression lines are drawn for the replciated and unreplicated ROIs, separately.

Figure S2. ROI-based confirmatory replication results for five personality subscores within healthy cohort. Donut plots summerising ROI-based replication rates (% of ROI) using three different criteria for three different sample sizes among heathy participants. The most inner

layers depict replication using “sign” only (blue: replicated, orange: not replicated). The middle layers define replication based on similar “sign” as well as “statistical significance” (i.e. $p < 0.05$) (blue: replicated, orange: not replicated). The most outer layers define replication using “bayes factor” (blue: “moderate-to-string evidence for H1, light blue: anecdotal evidence for H1; light orange: anecdotal evidence for H0, orange: “moderate-to-string evidence for H0”);