![JNeurosci — THE JOURNAL OF NEUROSCIENCE]

*Research Articles: Behavioral/Cognitive*

# Power-up: a reanalysis of 'power failure' in neuroscience using mixture modelling

**Camilla L Nord[1], Vincent Valton[1], John Wood[2] and Jonathan P Roiser[1]**

[1]*Institute of Cognitive Neuroscience, University College London, 17 Queen Square, London, UK, WC1N 3AZ*

[2]*Research Department of Primary Care and Population Health, University College London Medical School, Rowland Hill Street, London, UK, NW3 2PF*

**Correspondence should be addressed to** Please address correspondence to: Camilla L Nord, 17 Queen Square, London WC1N 3AZ, camilla.nord.11@ucl.ac.uk, +442076791138

**Alerts:** Sign up at www.jneurosci.org/cgi/alerts to receive customized email alerts when the fully formatted version of this article is published.

# Power-up: a reanalysis of 'power failure' in neuroscience using mixture modelling

3  Abbreviated title: Power-up

4  Nord, Camilla L[1*], Valton, Vincent[1*], Wood, John[2], Roiser, Jonathan P[1]

5  [1]Institute of Cognitive Neuroscience, University College London, 17 Queen Square,
6  London, UK, WC1N 3AZ

7  [2] Research Department of Primary Care and Population Health, University College
8  London Medical School, Rowland Hill Street, London, UK, NW3 2PF

9  *These authors contributed equally to this work.

10

11  Please address correspondence to:

12  Camilla L Nord

13  17 Queen Square, London WC1N 3AZ

14  camilla.nord.11@ucl.ac.uk

15  +442076791138

16

17  **No. pages: 31**

18  **No. figures: 5;  no. tables: 2**

19  **No. words abstract: 169**

28

29

30

31

32

33

## Abstract

34

35  Evidence for endemically low statistical power has recently cast neuroscience findings

36  into doubt. If low statistical power plagues neuroscience, this reduces confidence in

37  reported effects. However, if statistical power is not uniformly low, such blanket mistrust

38  might not be warranted. Here, we provide a different perspective on this issue, analysing

39  data from an influential paper reporting a median power of 21% across 49 meta-

40  analyses (Button et al., 2013). We demonstrate, using Gaussian mixture modelling, that

41  the sample of 730 studies included in that analysis comprises several subcomponents;

42  therefore the use of a single summary statistic is insufficient to characterise the nature of

43  the distribution. We find that statistical power is extremely low for studies included in

44  meta-analyses that reported a null result; and that it varies substantially across subfields

45  of neuroscience, with particularly low power in candidate gene association studies.

46  Thus, while power in neuroscience remains a critical issue, the notion that studies are

47  systematically underpowered is not the full story: low power is far from a universal

48  problem.

49

50

51

52

53

54

55

**Significance statement**

Recently, researchers across the biomedical and psychological sciences have become

concerned with the reliability of results. One marker for reliability is statistical power: the

probability of finding a statistically significant result, given that the effect exists. Previous

evidence suggests that statistical power is low across the field of neuroscience. Our

results present a more comprehensive picture of statistical power in neuroscience: on

average, studies are indeed underpowered—some very seriously so—but many studies

show acceptable or even exemplary statistical power. We show that this heterogeneity in

statistical power is common across most subfields in neuroscience (psychology,

neuroimaging, etc.). This new, more nuanced picture of statistical power in neuroscience

could affect not only scientific understanding, but potentially policy and funding decisions

for neuroscience research.

## Introduction

Trust in empirical findings is of vital importance to scientific advancement, but publishing biases and questionable research practices can cause unreliable results (Nosek et al., 2012; Button et al., 2013). In recent years, scientists and funders across the biomedical and psychological sciences have become concerned with what has been termed a crisis of replication and reliability (Barch and Yarkoni, 2013).

One putative marker for the reliability of results is statistical power: the probability that a statistically significant result will be declared, given that the null hypothesis is false (i.e., a real effect exists). It can be shown that, in the context of field-wide underpowered studies, a smaller proportion of significant findings will reflect true positives than if power is universally high (Ioannidis, 2005). A recent influential paper by Button and colleagues (Button et al., 2013) calculated statistical power across all meta-analyses published in 2011 that were labelled as "neuroscience" by Thomson Reuters Web of Science. It concluded that neuroscience studies were systematically underpowered, with a median statistical power of 21%, and that the proportion of statistically significant results that reflect true positives is therefore likely to be low. The prevalence of very low power has serious implications for the field. If the majority of studies are indeed underpowered, statistically significant findings are untrustworthy, and scientific inference will often be misinformed. This analysis provoked considerable debate in the field about whether neuroscience does indeed suffer from endemic low statistical power (Bacchetti, 2013; Quinlan, 2013). We sought to add nuance to this debate by re-analysing the original dataset using a more fine-grained approach, and provide a different perspective on statistical power in neuroscience.

We extended the analyses of Button and colleagues (Button et al., 2013), using data from all 730 individual studies, which provided initial results that were consistent with the

102    original report (which used only the median-sized study in 49 meta-analyses). To

103    quantify the heterogeneity of the dataset we made use of Gaussian mixture modelling

104    (GMM) (Corduneanu and Bishop, 2001), which assumes that the data may be described

105    as being composed of multiple Gaussian components. We then used model comparison

106    to find the most parsimonious model for the data. We also categorised each study based

107    on its methodology to examine whether low power is common to all fields of

108    neuroscience.

109    We find strong evidence that the distribution of power across studies is multi-modal, with

110    the most parsimonious model tested including four components. Moreover, we show that

111    candidate gene association studies and studies from meta-analyses with null results

112    make up the majority of extremely low powered studies in the analysis of Button and

113    colleagues. Although median power in neuroscience is low, the distribution of power is

114    heterogeneous, and there are clusters of adequately and even well-powered studies in

115    the field. Thus, our in-depth analysis reveals that the crisis of power is not uniform:

116    instead, statistical power is extremely diverse across neuroscience.

## 117    **Methods**

118    Experimental design and analysis

119    *Re-analysing 'power failures'*

120    Our initial analysis took a similar approach to that of Button and colleagues, but contrary

121    to their protocol (which reported power only for the median-sized study in each meta-

122    analysis: N=49), we report power for each of the 730 individual studies (see Figure 3a

123    and Table 1). As in the original analysis, we defined power as the probability that a given

124    study would declare a significant result, assuming that the population effect size was

125    equal to the weighted mean effect size derived from the corresponding meta-analysis

126  (note that this differs from 'post-hoc' power, in which the effect size would be assumed to

127  be equal to the reported effect size from each individual study (O'Keefe, 2007)).

128  For experiments with a binary outcome, power was calculated by assuming that the

129  expected incidence or response rate for the control group (i.e. the base rate) was equal

130  to that reported in the corresponding meta-analysis and, similarly, used an assumed

131  "treatment effect" (odds or risk ratio) equal to that given by each meta-analysis. The test

132  statistic used for the calculation was the log odds-ratio divided by its standard error. The

133  latter was derived from a first order approximation, and estimated by the square root of

134  the sum of the reciprocals of the expected values of the counts in the 2-by-2 summary

135  table. The test statistic itself was then referenced to the standard normal distribution for

136  the purposes of the power calculation. For studies reporting Cohen's *d*, the assumed

137  treatment effect was again taken directly from the corresponding meta-analysis, and all

138  power calculations were based on the standard noncentral *t*-distribution. For

139  comparability with the original study we calculated the median power across all 730

140  individual studies which was equal to 23%, close to the 21% reported by Button and

141  colleagues (2013).

142  Figure 1 shows an overview of our analytical process. We additionally classified each

143  study according to methodology: candidate gene association studies (N=234);

144  psychology (N=198); neuroimaging (N=65); treatment trials (N=145); neurochemistry

145  (N=50); and a miscellaneous category (N=38 studies from N=2 meta-analyses). Two

146  independent raters categorized the 49 meta-analyses into these six subfields, with 47/49

147  classified consistently; the remaining two were resolved following discussion. Before

148  continuing our analysis in more depth, we present the reader with results that are directly

149  comparable with the analysis of Button and colleagues (with the addition of the

150  subfields; Table 2). These results are intended for comparison with our more nuanced

151  characterisation of the distributions using GMMs presented below; given the results of

152    those GMMs (which suggest the these distributions are multi-modal and therefore not

153    well characterised by a single measure of central tendency) they should not be used to

154    draw strong inferences.

155

156    **Figure 1. Classification of studies for analysis**

157    Description of study methodology. GMM=Gaussian mixture model.

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

| First author of study | *k* | Cohen's *d* | Odds ratio | CI | Significance | Classification |
|---|---|---|---|---|---|---|
| Babbage (Babbage et al., 2011) | 13 | -1.11 | | -0.97 to -1.25 | * | Psychology |
| Bai (Bai, 2011) | 18 | | 1.47 | 1.22 to 1.77 | * | Genetic |
| Bjorkhelm-Bergman (Björkhem-Bergman et al., 2011) | 6 | -1.20 | | 1.6 to 8.0 | * | Treatment |
| Bucossi (Bucossi et al., 2011) | 21 | .41 | | .17 to .65 | * | Neurochemistry |
| Chamberlain (Chamberlain et al., 2011) | 11 | -.51 | | .825 to 1.08 | * | Psychology |
| Chang (Chang et al., 2011a) | 56 | -.19 | | -.29 to -.1 | * | Psychology |
| Chang (Chang et al., 2011b) | 6 | | .98 | .86 to 1.12 | - | Genetic |
| Chen (Chen et al., 2011) | 12 | | .6 | .52 to .69 | * | Miscellaneous |
| Chung (Chung and Chua, 2011) | 11 | | .67 | .43 to 1.04 | - | Treatment |
| Domellof (Domellöf et al., 2011) | 14 | | 2.12 | 1.59 to 2.78 | * | Psychology |
| Etminan (Etminan et al., 2011) | 14 | | 0.8 | .7 to .92 | * | Treatment |
| Feng (Feng et al., 2011) | 4 | | 1.20 | 1.04 to 1.4 | * | Genetic |
| Green (Green et al., 2011) | 17 | -.59 | | -.93 to -.257 | * | Neurochemistry |
| Han (Han et al., 2011) | 14 | | 1.35 | 1.06 to 1.72 | * | Genetic |
| Hannestad (Hannestad et al., 2011) | 13 | -.13 | | -.55 to .29 | - | Treatment |
| Hua (Hua et al., 2011) | 27 | | 1.13 | 1.05 to 1.21 | * | Genetic |
| Lindson (Lindson and Aveyard, 2011) | 8 | | 1.05 | .92 to 1.19 | - | Treatment |
| Liu (Liu et al., 2011a) | 12 | | 1.04 | .88 to 1.22 | - | Genetic |
| Liu (Liu et al., 2011b) | 6 | | .89 | .82 to .96 | * | Genetic |
| MacKillop (MacKillop et al., 2011) | 57 | .58 | | .509 to .641 | * | Psychology |
| Maneeton (Maneeton et al., 2011) | 5 | | 1.67[†] | 1.23 to 2.26 | * | Treatment |
| Ohi (Ohi et al., 2011) | 6 | | 1.12 | 1.00 to 1.26 | * | Genetic |
| Olabi (Olabi et al., 2011) | 14 | -.4 | | -.62 to -.19 | * | Brain imaging |
| Oldershaw (Oldershaw et al., 2011) | 10 | -.51 | | -.73 to -.28 | * | Psychology |
| Oliver (Oliver et al., 2011) | 7 | | .86 | 0.79 to .95 | * | Treatment |
| Peerbooms (Peerbooms et al., 2011) | 36 | | 1.26 | 1.09 to 1.46 | * | Genetic |
| Pizzagalli (Pizzagalli, 2011) | 22 | .92 | | .442 to 1.393 | * | Treatment |
| Rist (Rist et al., 2011) | 5 | | 2.06 | 1.33 to 3.19 | * | Miscellaneous |
| Sexton (Sexton et al., 2011) | 8 | .43 | | .063 to .799 | * | Brain imaging |
| Shum (Shum et al., 2011) | 11 | .89 | | .75 to 1.02 | * | Psychology |
| Sim (Sim et al., 2011) | 2 | | 1.23[†] | 1.08 to 1.52 | * | Treatment |
| Song (Song et al., 2011) | 12 | .15 | | .043 to .264 | * | Neurochemistry |
| Sun (Sun et al., 2011) | 6 | | 1.93 | 1.55 to 2.41 | * | Genetic |
| Tian (Tian et al., 2011) | 4 | 1.26 | | .947 to 1.568 | * | Treatment |
| Trzesniak (Trzesniak et al., 2011) | 11 | | 1.98 | 1.33 to 2.94 | * | Brain imaging |
| Veehof (Veehof et al., 2011) | 8 | .37 | | .20 to .53 | * | Treatment |
| Vergouwen (Vergouwen et al., 2011) | 24 | | .83 | .74 to .93 | * | Treatment |
| Vieta (Vieta et al., 2011) | 10 | | .68[†] | .60 to .77 | * | Treatment |
| Wisdom (Wisdom et al., 2011) | 53 | -.14 | | -.21 to -.07 | * | Genetic |
| Witteman (Witteman et al., 2011) | 26 | -1.41 | | -1.76 to -1.05 | * | Psychology |
| Woon (Woon and Hedges, 2011) | 24 | -.60 | | -.83 to -.37 | * | Brain imaging |
| Xuan (Xuan et al., 2011) | 20 | | 1.00 | .861 to 1.156 | - | Genetic |
| Yang (cohort) (Yang et al., 2011a) | 14 | | 1.38[†] | 1.18 to 1.61 | * | Miscellaneous |
| Yang (case control) (Yang et al., 2011a) | 7 | | 2.48 | 1.93 to 3.19 | * | Miscellaneous |
| Yang (Yang et al., 2011b) | 3 | 0.67 | | .43 to .92 | * | Treatment |
| Yuan (Yuan et al., 2011) | 14 | | 4.98 | 3.97 to 6.23 | * | Genetic |
| Zafar (Zafar et al., 2011) | 8 | | 1.07[†] | .91 to 1.27 | - | Treatment |
| Zhang (Zhang et al., 2011) | 12 | | 1.27 | 1.01 to 1.59 | * | Genetic |
| Zhu (Zhu et al., 2011) | 8 | 0.84 | | .18 to 1.49 | * | Brain imaging |

183 **Table 1. Characteristics and classification of included meta-analyses**

184 Classification performed by two independent raters. *k*: number of studies; [†] indicates relative risk; CI:
185 confidence interval; * indicates *p*<0.05.

| Group of studies | Median power (%) | Min. power (%) | Max. power (%) | 2.5[th] and 97.5[th] percentile (based on raw data) | 95% HDI (based on GMMs) | Total *k* |
|---|---|---|---|---|---|---|
| All studies | 23 | 0.05 | 1 | [0.05 to 1.00] | [0.00 to 0.72], [0.8 to 1.00] | 730 |
| All studies excluding null | 30 | 0.05 | 1 | [0.05 to 1.00] | [0.01 to 0.73], [0.79 to 1.00] | 638 |
| Genetic | 11 | 0.05 | 1 | [0.05 to 0.94] | [0.00 to 0.44], [0.63 to 0.93] | 234 |
| Treatment | 20 | 0.05 | 1 | [0.05 to 1.00] | [0.00 to 0.65], [0.91 to 1.00] | 145 |
| Psychology | 50 | 0.07 | 1 | [0.07 to 1.00] | [0.02 to 0.24], [0.28 to 1.00] | 198 |
| Imaging | 32 | 0.11 | 1 | [0.11 to 1.00] | [0.03 to 0.54], [0.71 to 1.00] | 65 |
| Neurochemistry | 47 | 0.07 | 1 | [0.07 to 1.00] | [0.02 to 0.79], [0.92 to 1.00] | 50 |
| Miscellaneous | 57 | 0.11 | 1 | [0.11 to 1.00] | [0.09 to 1.00] | 38 |

186    **Table 2. Median power by study type**

187    Median, maximum, and minimum power subdivided by study type. We also provide the 2.5[th] and
188    97.5[th] percentile of the frequency distribution of power estimates of individual studies for the raw data
189    and 95% highest-density intervals (95% HDI) for the GMMs. We used highest density intervals (HDI)
190    to summarise the intervals of the most probable values from the distribution. HDIs differ from CIs in
191    that they represent the most probable values of the distribution rather than symmetric credible
192    intervals in a central tendency. As a result, HDIs are more suitable for summarising skewed and
193    multimodal distributions than CIs. HDIs were computed using the HDRCDE R toolbox, which finds the
194    shortest intervals such that these intervals encompass the 95% most probable values of the
195    distribution. Multiple intervals may be identified if a region between modes of the distribution is
196    unrepresentative of the distribution (i.e. below the 5% threshold) (Wand et al., 1991; Hyndman, 1996;
197    Samworth and Wand, 2010), which occurs for multimodal data.

198

199    *One or many populations?*

200    The common measures of central tendency (mean, median, and mode) may not

201    always characterise populations accurately, because distributions can be complex,

202    and made up of multiple 'hidden' subpopulations. Consider the distribution of height

203    in the United States: the mean is 168.8±13.04 cm (Fryar et al., 2012). This statistic is

204    rarely reported because the distribution comprises two distinct populations: male

205    (175.9 ±15.03 cm) and female (162.1 cm ±10.8 cm). The mean of the male

206   population is greater than the 95[th] percentile of the female population. Thus, a single

207   measure of central tendency fails to describe this distribution adequately.

208   In an analogous fashion, the original paper of Button and colleagues reported a

209   median of 21% power, which could be interpreted as implying a degree of statistical

210   homogeneity across neuroscience. The use of the median as a summary statistic,

211   while having the straightforward interpretation of 'half above and half below', also

212   implies that the power statistics are drawn from a distribution with a single central

213   tendency. As we show below, this assumption is contradicted by our analyses, which

214   makes the median statistic difficult to interpret. It should be noted that Button and

215   colleagues themselves described their results as demonstrating a 'clear bimodal

216   distribution'. Therefore we next explored the possibility that the power data originated

217   from a combination of multiple distributions, using GMM.

218   GMM (similar to latent class analysis and factor models (Lubke and Muthén, 2005))

219   can be used to represent complex density functions where the central limit theorem

220   does not apply, such as in the case of bimodal or multi-modal distributions. We fit

221   GMMs with varying numbers of 'K' unknown components to the data and performed

222   model selection using the Bayesian Information Criteria (BIC) scores to compare

223   models with different fit and complexity (the higher the number of 'K' unknown

224   components the more complex the model). This allowed us to take a data-driven

225   approach, as opposed to direct mixture models using a set number of components:

226   thus, we were agnostic as to the number of components that emerged from the

227   model. The GMM with the lowest BIC identifies the most parsimonious model,

228   trading model fit against model complexity. A difference in BIC between models of 10

229   or above on a natural logarithm scale is indicative of strong evidence in support of

230   the model with the lower score (Kass and Raftery, 1995). To ensure that we used the

231    most suitable GMM for this dataset, we ran different GMM models: standard GMMs,

232    regularized GMMs, and Dirichlet Process GMMs (see below for full methods, and

233    Figure 2 for model comparison, and model selection). The results were similar using

234    each of these techniques (see Figure 2).

235    *Finite Gaussian mixture model*

236    For a finite GMM, the corresponding likelihood function is given by (Corduneanu and

237    Bishop, 2001):

$$P(D|\pi,\theta) = \prod_{n=1}^{N}\left[\sum_{i=1}^{K}\pi_i \, \mathcal{N}(x_n|\theta_i)\right]$$

238    where $\pi_i$ denotes the mixing coefficient (proportions of the *i*–th component),

239    $\mathcal{N}(x_n|\theta_i)$   denotes the conditional probability of the observation $x_n$ given by a

240    Gaussian distribution with parameters $\theta_i$ and *D* denotes the whole dataset of

241    observations, $x_n$. Generally speaking, this means that we believe that there is an

242    underlying generative structure to the observed data, and that a mixture of Gaussian

243    components would a reasonable description/approximation of the true generative

244    process of this data. That is, we assume that the data *D* has been generated from a

245    mixture of Gaussians distributions with varying means, variances, and weights

246    (model parameters), which we want to uncover. To do so, we perform model

247    inversion and find the point estimates of the model parameters that maximize the

248    likelihood (see eq. 1 above) of the observed data (maximum likelihood estimation).

249

250    Model inversion is performed using the iterative EM (expectation-maximisation)

251    algorithm, which finds a local maximum of the likelihood function given initial starting

252    parameters. We performed 50 restarts with kmeans++ initialization (Arthur and

253   Vassilvitskii, 2007). Multiple restarts were performed in order to find the global

254   maximum of the likelihood (i.e., the best GMM for the data; that is, the parameters

255   that maximize the chance of observing the data), as opposed to a local maximum.

256   This allowed us to ensure that convergence was achieved for all GMMs, on all

257   datasets.

258   Traditionally, finite mixture modelling approaches require the number of components

259   to be specified in advance of analysing the data. That is, for each finite Gaussian

260   mixture model fitted to the data, one is required to input the number of components $K$

261   present in the mixture (model inversion only estimates the parameters for each

262   component). Finding the number of components present in the data is a model

263   selection problem, and requires fitting multiple GMMs with varying numbers of

264   components to the data, then comparing the model evidence for each fit, and

265   selecting the most parsimonious model for the data in question (Bishop, 2006;

266   Gershman and Blei, 2012; Murphy, 2012).

267   It is worth noting, however, that GMMs can be subject to instabilities, such as

268   singularities of the likelihood function. Specifically, it is possible for one component to

269   'collapse' all of its variance onto a single data point, leading to an infinite likelihood

270   (Bishop, 2006; Murphy, 2012) and to incorrect parameter estimation for the model.

271   Multiple techniques have been developed in order to address this problem. The

272   simplest and most commonly used technique is to introduce a regularization

273   parameter. Another is to adopt a fully Bayesian approach and apply soft constraints

274   on the possible range of likely parameter values, therefore preventing problematic

275   and unrealistic parameter values. Both methodologies were used in this study, and

276   we report on the resulting analysis for both implementations in the model selection

277   section (below).

278    *Finite Gaussian mixture model with regularization*

279    In typical finite mixture models, a regularization parameter can be added in order to

280    avoid likelihood singularities. To do so, a very small value is added to the diagonal of

281    the covariance matrix, enforcing positive-definite covariance and preventing infinitely

282    small precision parameters for individual components. This model specification

283    enables one to address the issue of 'collapsing' components but also enforces

284    simpler explanations of the data, favouring models with fewer components. The

285    larger the regularization parameter, the simpler the models will be, as single

286    components will tend to encompass a larger subspace of the data partition. In this

287    study we introduced a regularization parameter of 0.001, which represents a

288    reasonable trade-off between preventing over-fitting components to noise in the

289    dataset, while capturing the most salient features from the data (the separate peaks);

290    therefore providing a better generative model of the data than using non-regularized

291    GMMs. We used this approach for our primary inferences.

292    *Dirichlet Process Gaussian mixture model (DPGMM)*

293    Dirichlet Process (DP) Gaussian mixture models (DPGMMs) are a class of Bayesian

294    non-parametric methods that avoid the issue of model selection when identifying the

295    optimal number of components in a mixture model (Gershman and Blei, 2012;

296    Murphy, 2012). With DPGMM, we expand the original GMM model to incorporate a

297    prior over the mixing distribution, and a prior over the component parameters (mean

298    and variance of components). Common choices for DPGMM priors are conjugate

299    priors such as the normal-inverse-Wishart distribution over the mean and covariance

300    matrix of components, and a non-parametric prior over mixing proportions based on

301    the DP.

302    The DP, often referred to as the Chinese restaurant process or the stick-breaking

303    process, is a distribution over infinite partitions of integers (Gershman and Blei,

304    2012; Murphy, 2012). As a result, the DPGMM theoretically allows for an infinite

305    number of components as it lets the number of components grow as the amount of

306    data increases. The DP assigns each observation to a cluster with a probability that

307    is proportional to the number of observations already assigned to that cluster. That

308    is, the process will tend to cluster data points together, dependent on the population

309    of the existing cluster and a concentration parameter $\alpha$. The smaller the $\alpha$

310    parameter, the more likely it is that an observation will be assigned to an existing

311    cluster with probability proportional to the number of elements already assigned to

312    this cluster. This phenomenon is often referred to as the 'rich get richer'. This

313    hyperparameter $\alpha$ indirectly controls how many clusters one expects to see from the

314    data (another approach is to treat $\alpha$ as unknown, using a gamma hyperprior over $\alpha$,

315    and letting the Bayesian machinery infer the value (Blei and Jordan, 2006)).

316    Implementation and analysis for the non-regularized finite GMMs, regularized finite

317    GMMs, and DPGMMs was performed using Matlab R2015b (Mathworks Inc.), using

318    the Statistics and Machine Learning toolbox, the Lightspeed toolbox and the vdpgm

319    toolbox (Kurihara et al., 2007).

320    *Model selection*

321    The traditional mixture modelling approach requires the number of clusters or

322    components to be specified in advance of analysing the data. However, in many

323    settings, including here, one does not know the number of underlying components

324    and would like to estimate this directly from the data. One approach typically used

325    with finite mixture models is to fit the data with varying number of components and

326   then to select the model that provides the best trade-off between model fit (how well

327   the model explains the data) and model complexity (how many component

328   parameters are used in the model). A metric commonly used in this setting is the

329   Bayesian Information Criterion (BIC), which allows one to compute an approximation

330   to the Bayes factor (relative evidence) for a model. The BIC typically has two terms,

331   the likelihood (how well the model fits the data) and a complexity term that penalizes

332   more complex models with more free parameters (e.g. the number of components).

333   The model with the lowest BIC metric is usually preferred as it provides the most

334   parsimonious and generalizable model of the data.

335   For each one of the following datasets model fits were performed using non-

336   regularized and regularized finite mixtures with up to 15 components (up to 10

337   components for the subfield categories – Figure 2): the original dataset; the original

338   dataset excluding null studies; each methodological subfield within the original

339   dataset (Genetics, Psychology, Neurochemistry, Treatment, Imaging, and

340   Miscellaneous studies); and the original dataset excluding each methodological

341   subfield. Model selection was then performed using the BIC in order to select the

342   most parsimonious model for each dataset. Figure 2 presents (for each dataset) the

343   corresponding BIC metric for increasing levels of model complexity. Plain blue lines

344   denote the BIC metric using non-regularized GMMs, while plain red lines denote the

345   BIC using regularized GMMs. The BIC metric curve for non-regularized GMMs (blue

346   line) exhibits wide jumps (Figure 2), while the function should remain relatively

347   smooth as seen with regularized-GMMs (red line). This suggests that non-

348   regularized GMMs results were prone to overfitting and were inadequate for some of

349   our datasets.

350    Finally, we compared different modelling methodologies, in order to select and report

351    the most robust findings in terms of the estimation of the number of components. We

352    compared non-regularized GMMs, regularized GMMs and DPGMMs on the same

353    datasets (Figure 2), and found that regularized GMMs provided the most

354    conservative estimation of the number of components. We therefore opted to report

355    these results as the main findings.

356

357

358    **Figure 2. Model comparison and model selection analysis for Gaussian mixture models**
359    **(GMM), regularized GMMs and Dirichlet process GMMs (DPGMMs).** The blue and red lines
360    display Bayesian Information Criterion (BIC) scores (natural log scale) for non-regularized GMMs and
361    regularized GMMs, respectively, for different levels of model complexity (number of mixture
362    components). The lowest BIC score indicates the model that provides the best compromise between
363    model fit (likelihood) and model complexity for the given dataset. Winning models for GMMs (purple
364    dotted-dash vertical line), regularized GMMs (yellow dashed vertical line), and DPGMMs (green
365    dotted vertical line) are clearly present for each dataset, enabling direct comparison of the output for
366    each methodology. The regularized GMM approach provided the most parsimonious interpretation of
367    the data on the two main datasets: all studies (a), excluding null studies (b) as well as 5 out of 6
368    subfield datasets – (c) to (h).

369    ## Results

370    We analysed the original sample of 730 powers (see histogram in Figure 3a). If the

371    median were the most appropriate metric to describe the distribution of powers across

372    studies, we would expect the GMM to produce a solution containing only a single

373    component. Instead, the most parsimonious GMM solution included four components,

374    with strong evidence in favour of this model versus either of the next best models (i.e.

375    GMMs with 3 or 5 components - see Figure 2). Importantly, this model revealed that the

376    overall distribution of power appears to be composed of sub-groups of lower and higher

377    powered studies (overlay in Figure 3a). We next explored possible sources of this

378    variability, considering the influence of both null effects and specific subfields of

379    neuroscience.

380

381

382

383

384    **Figure 3. Power of studies**

385    **Figure 3a-b: Histograms depicting the distribution of study powers across all 730 studies (a)**
386    **and across studies excluding null meta-analyses (b).** However, we note that excluding power
387    statistics from studies included in null meta-analyses may provide an overestimation of power,
388    because in many instances there remains uncertainty as to whether or not a true effect exists. Pale
389    overlay: results of the regularised Gaussian mixture model (GMM), identifying four components (C1,
390    C2, C3, C4) and their relative weights within the dataset. Below the histogram, pie charts depict
391    methodological subfields, as well as null meta-analyses, contributing to each component. The null
392    studies (white pie-chart sections) comprise 52 genetic studies and 40 treatment studies. The dark
393    blue line shows the sum of the components (overall GMM prediction). c-h: histograms depicting the
394    distribution of study powers across all meta-analyses, separated by subfield: candidate gene
395    association studies (c); psychology studies (d); neurochemistry studies (e); treatment studies (f);
396    imaging studies (g); miscellaneous studies (h). Pale overlays show the results of the regularised GMM
397    for each subfield; the dark lines show the sum of the components (overall GMM prediction).

398

399    *When is an effect not an effect?*

400    The first important source of variability we considered relates to the concept of power

401    itself. The calculation of power depends not just on the precision of the experiment

402    (heavily influenced by the sample size), but also on the true population effect size.

403    Logically, power analysis requires that an effect (the difference between population

404    distributions) actually exists. Conducting a power analysis when no effect exists violates

405    this predicate, and will therefore yield an uninterpretable result. Indeed, when no effect

406    exists the power statistic becomes independent of the sample size and is simply equal to

407    the Type I error rate; which by definition is the probability of declaring a significant result

408    under the null hypothesis.

409    To illustrate this point, consider the meta-analysis titled 'No association between APOE

410    epsilon 4 allele and multiple sclerosis susceptibility' (Xuan et al., 2011), which included a

411   total of 5,472 cases and 4,727 controls. The median effect size (odds ratio) reported was

412   precisely 1.00, with a 95% confidence interval from 0.861-1.156. Button and colleagues

413   calculated the median power to be 5%, which is equal to the Type I error rate. However,

414   as is evident from the paper's title, this meta-analysis was clearly interpreted by its

415   authors as indicating a null effect, which is consistent with the observed result. Indeed,

416   in this case the power is 5% for both the largest (N>3000) and the smallest (N<150)

417   study in the meta-analysis. In such cases the estimate of 5% power is not easily

418   interpretable.

419   On the other hand, it is problematic to assume that a non-significant meta-analytic

420   finding can be taken as evidence there is no true effect; in the Frequentist statistical

421   framework, failure to reject the null hypothesis cannot be interpreted as unambiguous

422   evidence that no effect exists (due to the potential for false negative results). For

423   example, reference 16 ('Effects on prolongation of Bazett's corrected QT interval of

424   seven second-generation antipsychotics in the treatment of schizophrenia: a meta-

425   analysis') reported a median effect size (odds ratio) of 0.67, with a 95% confidence

426   interval from 0.43-1.04. While this result was non-significant, the point estimate of the

427   effect size is greater than those from several meta-analyses that did achieve statistical

428   significance, and in our view it would be premature to conclude that this effect does not

429   exist.

430   These examples illustrate the difficulty in deciding whether conducting a power analysis

431   is appropriate. Even tiny effect sizes could hypothetically still exist: in any biological

432   system the probability that an effect is precisely null is itself zero – therefore all effects

433   "exist" by this definition (with certain exceptions, e.g. in the context of randomization),

434   even if to detect them we might need to test more individuals than are currently alive.

435   However, the notion of "falsely rejecting the null hypothesis" then loses its meaning

436   (Jacob Cohen, 1994). One approach would be to assume that an effect does not exist

437  until the observed evidence suggests that the null hypothesis can be rejected, consistent

438  with the logical basis of classical statistical inference. This would avoid any potential bias

439  towards very low power estimates due to non-existent effects. On the other hand, this

440  approach raises the potential problem of excluding effects that are genuinely very small,

441  which may cause a bias in the other direction. Within the constraints of the null

442  hypothesis significance testing framework, it is impossible to be confident that an effect

443  does not exist at all. Therefore, we cannot simply assume an effect does not exist after

444  failing to reject the null hypothesis, since a small effect could go undetected.

445  Motivated by this logic, we initially included studies from 'null meta-analyses' (i.e. where

446  the estimated effect size from the meta-analysis was not significantly different from the

447  null at the conventional alpha=0.05) in our GMMs (Figure 3a). However, we note that

448  excluding power statistics from studies included in null meta-analyses may provide an

449  overestimation of power, because in many instances there remains uncertainty as to

450  whether or not a true effect exists. Nonetheless, with the above caveats in mind, we also

451  wished to assess the degree to which null meta-analyses may have impacted the

452  results. Null results occurred in 7 of the 49 meta-analyses (92 of the 730 individual

453  studies), contributing a substantial proportion of the extremely low powered studies

454  (<10% power; Figure 3a, white pie chart segment of C1). When we restricted our

455  analysis only to studies within meta-analyses that reported statistically significant results

456  ('non-null' meta-analyses), the median study power (unsurprisingly) increased, but only

457  slightly, to 30%, and the nature of the resulting GMM distribution did not change

458  substantially (see Figure 3b). Thus, excluding null meta-analyses does not provide a

459  radically different picture. Therefore, we also examined another potential contributor to

460  power variability in neuroscience: the influence of specific subfields of neuroscience.

461  *Power in neuroscience subfields*

462    As described above, we categorised each meta-analysis into one of six methodological

463    subfields. Interestingly, statistical power varied significantly according to subfield

464    (permutation test of equivalence: $p$<0.001), with genetic association studies lower (11%

465    median power) than any other subfield examined (all Mann-Whitney U tests $p$<0.001).

466    This is consistent with the original report by Button and colleagues, which reported the

467    median power of animal studies (18% and 31% for two meta-analyses) and structural

468    brain imaging studies (8% across 41 meta-analyses). However, even within specific

469    subfields, the distribution of power is multimodal (see Figure 3c-h). This could represent

470    variability in statistical practices across studies, but another possible explanation is that

471    the size of the effect being studied varies substantially between meta-analyses, even

472    within the same subfield. This alternative explanation may, at least in part, account for

473    the variability between (and within) subfields of neuroscience.

474    The large number of extremely low powered candidate gene association studies

475    warrants additional comment. These were included in the original analysis because the

476    Web of Science classifies such studies as "neuroscience" if the phenotypes in question

477    are neurological or psychiatric disorders. However, modern genome-wide association

478    studies have revealed that the overwhelming majority of candidate gene association

479    studies have been underpowered, because the reliable associations that have been

480    identified are extremely small (Flint and Munafò, 2013); thus, very low power is expected

481    within this subgroup, which our analysis confirms (see Figure 3c). This subgroup of

482    studies can offer important lessons to the rest of neuroscience: without large genetic

483    consortia, the field of neuropsychiatric genetics might still be labouring under the

484    misapprehension that individual common variants make substantial contributions to the

485    risk for developing disorders. Providing that sampling and measurement are

486    standardised, pooling data across multiple sites has the potential to improve dramatically

487    not only statistical power, but also the precision on estimates of effect size.

488   Since numerous studies report that candidate gene association studies are severely

489   underpowered (Klerk et al., 2002; Colhoun et al., 2003; Duncan and Keller, 2011), and

490   given that candidate gene association studies comprised over one-third of our total

491   sample of studies, we suspected that they might contribute heavily to the lowest-power

492   peak in our distribution. We confirmed this: in the absence of genetic studies, many

493   studies remained underpowered, but the distribution contained proportionally fewer

494   studies in the lowest-power peak (around 10% power) (Figure 4a). Although low power

495   is clearly not limited to candidate gene association studies, they nonetheless seem to

496   have a greater influence on the overall power distribution than any other subfield,

497   skewing the distribution towards the lowest-power peak (Figure 4b-f).

498

499   **Figure 4. Gaussian Mixture Models (GMMs) excluding each subfield.**

500   GMMs for the whole population of studies excluding genetic studies (a), excluding psychology studies
501   (b), excluding neurochemistry studies (c), excluding treatment studies (d), excluding imaging studies
502   (e), and excluding the remaining miscellaneous studies (f). Compare with the distribution including all
503   studies (Figure 3a).

504

505   *Estimations of effect size*

506   An important factor contributing to the estimation of power is whether the effect size was

507   estimated accurately *a priori*. If researchers initially overestimated the effect size, even

508   the sample size specified by a power calculation would be insufficient to detect a real,

509   but smaller effect. Interestingly, our analysis also shows the existence of very high

510   powered studies within neuroscience, in which far more subjects have been included

511   than would technically be warranted by a power analysis. In this case, an *a priori*

512   underestimate of effect size could yield a very high powered study, if an effect proves to

513   be larger than initially expected (which has occasionally been reported (Open Science

514   Collaboration, 2015)). Another important consideration is that an over-estimation of

515  effect size might occur due to publication bias, which will skew effect size estimates from

516  meta-analyses upwards, resulting in an optimistic power estimate. This is an important

517  caveat to the results we report here: a bias toward publishing significant results means

518  that the power estimates we report will represent upper bounds on the true power

519  statistics. Unfortunately, we could not adequately address this potential confound

520  directly, since tests of publication bias themselves have very low power, particularly if

521  the number of studies in a meta-analysis is low. However, publication bias has long

522  been reported in psychology (Francis, 2012) and neuroscience (Sena et al., 2010), so it

523  is reasonable to assume that it has inflated estimates of statistical power in these

524  analyses.

525  *Simulating power in hypothetical fields*

526  One clear conclusion of our analyses is that the interplay between the proportion of true

527  effects and the power to detect those effects is crucial in determining the power

528  distribution of a field. We simulated four power graphs for hypothetical fields to illustrate

529  this point: one with low power (~50%), but where all effects exist (Figure 5a); one with

530  high power (~90%), where all effects exist (Figure 5b); one with low power (~50%),

531  where only a minority (25%) of effects exist (Figure 5c); and high power (~90%), but

532  where only a minority (25%) of effects exist (Figure 5d). We found that the 'low power'

533  field did not resemble the distribution of power in neuroscience we observed (Figure 3a).

534  Instead, our findings were closest to a mixture of two distributions: Figure 5c, with low

535  (~50%) power, and where only 25% of findings are true effects; and Figure 5d, with high

536  (~90%) power, but where only 25% of findings are true effects. This would be consistent

537  with the notion that the absence of true effects may contribute to the distribution of

538  statistical power in neuroscience.

539  **Figure 5. Simulated power distributions for four hypothetical fields.** (a) 'Easy field' with low
540  power (~0.5) and all effects exist; (b) 'Easy field' with high power (~0.9) and all effects exist; (c) 'Hard

541   field' with low power (~0.5) (for those effects that exist), but where effects exist in only 25% of cases;
542   (d) 'Hard field' with high power (~0.9) (for those effects that exist), but where effects exist in only exist
543   in 25% of cases. Power distributions were simulated by generating 50,000 samples with fixed sample-
544   size (N=45) while varying effect-size. For each panel, the effect-size was sampled from a truncated
545   (effect-size>0) Gaussian distribution with mean 0.3 (a & c) or 0.49 (b & d), so as to represent low or
546   high power respectively. For the 'hard' fields (c & d), 75% of the effect-size sample was generated
547   from a half-Gaussian distribution with mean=0. SD was set to 0.07 for all effect size distributions.
548   Similar results can be obtained by fixing the effect size and varying the sample size.

549   **Discussion**

550   *Implications for neuroscience*

551   We argue that a very influential analysis (cited over 1500 times at the time of writing)

552   does not adequately describe the full variety of statistical power in neuroscience. Our

553   analyses show that the dataset is insufficiently characterized by a single distribution.

554   Instead, power varies considerably, including between subfields of neuroscience, and is

555   particularly low for candidate gene association studies. Conducting power analyses for

556   null effects may also contribute to low estimates in some cases, though determining

557   when this has occurred is challenging. Importantly, however, power is far from adequate

558   in every subfield.

559   Our analyses do not negate the importance of the original work in highlighting poor

560   statistical practice in the field, but they do reveal a more nuanced picture. In such a

561   diverse field as neuroscience, it is not surprising that statistical practices differ. While

562   Button and colleagues were careful to point out that they identified a range of powers in

563   neuroscience, their reporting of a median result could be interpreted as implying that the

564   results were drawn from a single distribution, which our analyses suggest is not the

565   case. We confirm that low power is clearly present in many studies, and agree that

566   focusing on power is a critical step in improving the replicability and reliability of findings

567   in neuroscience. However, we also argue that low statistical power in neuroscience is

568   neither consistent nor universal.

569    Ethical issues accompany both under- and over-powered studies. Animal sacrifices,

570    drugs taken to human trials, and government funding are all wasted if power is too low.

571    However, blindly increasing sample size across the board, simply to satisfy concerns

572    about field-wide power failures, is also not the best use of resources. Instead, each

573    study design needs to be considered on its own merits. In this vein, one response to the

574    original article pointed out that any measure of a study's projected value suffers from

575    diminishing marginal returns: every additional animal or human participant adds less

576    statistical value than the previous one (Bacchetti et al., 2005, 2008; 2013).

577    Studies with extremely large sample sizes can also fall prey to statistically significant

578    findings for trivial effects that are unlikely to be either theoretically or clinical important

579    (Lenth, 2001; Ioannidis, 2005; Friston, 2012; Quinlan, 2013). In other words, the

580    assessment of power is determined by what we consider to be an interesting (i.e.

581    nontrivial) effect size (Cohen, 1988). This dependency means that power considerations

582    are meaningless in the absence of assumptions about how large effect sizes need to be

583    in order to be considered theoretically or clinically important; and this may vary

584    dramatically across different fields. This is particularly relevant in fields where multiple

585    comparisons are performed routinely, such as genetics and neuroimaging (Friston,

586    2012). Conversely, smaller studies can only detect large effect sizes, and may suffer

587    from imprecise estimates of effect size and interpretive difficulties. Crucially, there is no

588    single study design that will optimise power for every genetic locus or brain area. In fact,

589    power estimates for individual studies are themselves extremely noisy and may say little

590    about the actual power in any given study. However, a move away from presenting only

591    p-values and towards reporting point estimates and confidence intervals (as long

592    advocated by statisticians), and towards sharing data to improve such estimates, would

593    allow researchers to make better informed decisions about whether an effect is likely to

594    be clinically or theoretically useful.

595

596

*Conclusion*

We have demonstrated the great diversity of statistical power in neuroscience. Do our

findings lessen concerns about statistical power in neuroscience? Unfortunately not. In

fact, the finding that the distribution of power is highly heterogeneous demonstrates an

undesirable inconsistency, both within and between methodological subfields. Yet within

this variability are several appropriately, and even very high powered studies. Therefore,

we should not tar all studies with the same brush, but instead encourage investigators to

engage in the best research practices, including preregistration of study protocols

(ensuring the study will have sufficient power), routine publication of null results, and

avoiding practices such as p-hacking that lead to biases in the published literature.

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

## References

Babbage DR, Yim J, Zupan B, Neumann D, Tomita MR, Willer B (2011) Meta-analysis of facial affect recognition difficulties after traumatic brain injury. Neuropsychology 25:277.

Bacchetti P (2013) Small sample size is not the real problem. Nat Rev Neurosci 14:585–585.

Bacchetti P, McCulloch CE, Segal MR (2008) Simple, defensible sample sizes based on cost efficiency. Biometrics 64:577–585.

Bacchetti P, Wolf LE, Segal MR, McCulloch CE (2005) Ethics and sample size. Am J Epidemiol 161:105–110.

Bai H (2011) Meta-analysis of 5, 10-methylenetetrahydrofolate reductase gene poymorphism as a risk factor for ischemic cerebrovascular disease in a Chinese Han population. Neural Regen Res 6:277–285.

Barch DM, Yarkoni T (2013) Introduction to the special issue on reliability and replication in cognitive and affective neuroscience research. Cogn Affect Behav Neurosci 13:687–689.

Bishop CM (2006) Pattern recognition and machine learning. springer New York.

Björkhem-Bergman L, Asplund AB, Lindh JD (2011) Metformin for weight reduction in non-diabetic patients on antipsychotic drugs: a systematic review and meta-analysis. J Psychopharmacol (Oxf) 25:299–305.

Blei DM, Jordan MI (2006) Variational inference for Dirichlet process mixtures. Bayesian Anal 1:121–143.

Bucossi S, Ventriglia M, Panetta V, Salustri C, Pasqualetti P, Mariani S, Siotto M, Rossini PM, Squitti R (2011) Copper in Alzheimer's disease: a meta-analysis of serum, plasma, and cerebrospinal fluid studies. J Alzheimers Dis 24:175–185.

Button KS, Ioannidis JP, Mokrysz C, Nosek BA, Flint J, Robinson ES, Munafò MR (2013) Power failure: why small sample size undermines the reliability of neuroscience. Nat Rev Neurosci 14:365–376.

Chamberlain SR, Robbins TW, Winder-Rhodes S, Müller U, Sahakian BJ, Blackwell AD, Barnett JH (2011) Translational approaches to frontostriatal dysfunction in attention-deficit/hyperactivity disorder using a computerized neuropsychological battery. Biol Psychiatry 69:1192–1203.

Chang W, Arfken CL, Sangal MP, Boutros NN (2011a) Probing the relative contribution of the first and second responses to sensory gating indices: A meta-analysis. Psychophysiology 48:980–992.

656    Chang X-L, Mao X-Y, Li H-H, Zhang J-H, Li N-N, Burgunder J-M, Peng R, Tan E-K (2011b) Functional
657            parkin promoter polymorphism in Parkinson's disease: new data and meta-analysis. J Neurol
658            Sci 302:68–71.

659    Chen C, Xu T, Chen J, Zhou J, Yan Y, Lu Y, Wu S (2011) Allergy and risk of glioma: a meta-analysis. Eur
660            J Neurol 18:387–395.

661    Chung AK, Chua S (2011) Effects on prolongation of Bazett's corrected QT interval of seven second-
662            generation antipsychotics in the treatment of schizophrenia: a meta-analysis. J
663            Psychopharmacol (Oxf) 25:646–666.

664    Cohen J (1988) Statistical power analysis for the behavioral sciences. Vol. 2. Lawrence Earlbaum
665            Assoc Hillsdale NJ.

666    Corduneanu A, Bishop CM (2001) Variational Bayesian model selection for mixture distributions. In,
667            pp 27–34. Morgan Kaufmann Waltham, MA.

668    Domellöf E, Johansson A-M, Rönnqvist L (2011) Handedness in preterm born children: a systematic
669            review and a meta-analysis. Neuropsychologia 49:2299–2310.

670    Etminan N, Vergouwen MD, Ilodigwe D, Macdonald RL (2011) Effect of pharmaceutical treatment on
671            vasospasm, delayed cerebral ischemia, and clinical outcome in patients with aneurysmal
672            subarachnoid hemorrhage: a systematic review and meta-analysis. J Cereb Blood Flow
673            Metab 31:1443–1451.

674    Feng X, Wang F, Zou Y, Li W, Tian Y, Pan F, Huang F (2011) Association of FK506 binding protein 5
675            (FKBP5) gene rs4713916 polymorphism with mood disorders: a meta-analysis. Acta
676            Neuropsychiatr 23:12–19.

677    Flint J, Munafò MR (2013) Candidate and non-candidate genes in behavior genetics. Curr Opin
678            Neurobiol 23:57–61.

679    Francis G (2012) Too good to be true: Publication bias in two prominent studies from experimental
680            psychology. Psychon Bull Rev 19:151–156.

681    Friston K (2012) Ten ironic rules for non-statistical reviewers. Neuroimage 61:1300–1310.

682    Fryar C, Gu Q, Ogden (2012) Anthropometric Reference Data for Children and Adults: United States,
683            2007-2010. U.S. Department of Health and Human Services.

684    Gershman SJ, Blei DM (2012) A tutorial on Bayesian nonparametric models. J Math Psychol 56:1–12.

685    Green M, Matheson S, Shepherd A, Weickert C, Carr V (2011) Brain-derived neurotrophic factor
686            levels in schizophrenia: a systematic review with meta-analysis. Mol Psychiatry 16:960–972.

687    Han X-M, Wang C-H, Sima X, Liu S-Y (2011) Interleukin-6− 174G/C polymorphism and the risk of
688            Alzheimer's disease in Caucasians: A meta-analysis. Neurosci Lett 504:4–8.

689    Hannestad J, DellaGioia N, Bloch M (2011) The effect of antidepressant medication treatment on
690            serum levels of inflammatory cytokines: a meta-analysis. Neuropsychopharmacology
691            36:2452–2459.

692   Hua Y, Zhao H, Kong Y, Ye M (2011) Association between the MTHFR gene and Alzheimer's disease: a
693       meta-analysis. Int J Neurosci 121:462–471.

694   Hyndman RJ (1996) Computing and graphing highest density regions. Am Stat 50:120–126.

695   Ioannidis JP (2005) Why most published research findings are false. PLoS Med 2:e124.

696   Jacob Cohen (1994) The earth is round (p<0.05). Am Psychol 49:997–1003.

697   Kass RE, Raftery AE (1995) Bayes factors. J Am Stat Assoc 90:773–795.

698   Kurihara K, Welling M, Teh YW (2007) Collapsed Variational Dirichlet Process Mixture Models. In, pp
699       2796–2801.

700   Lenth RV (2001) Some practical guidelines for effective sample size determination. Am Stat 55:187–
701       193.

702   Lindson N, Aveyard P (2011) An updated meta-analysis of nicotine preloading for smoking cessation:
703       investigating mediators of the effect. Psychopharmacology (Berl) 214:579–592.

704   Liu H, Liu M, Wang Y, Wang X-M, Qiu Y, Long J-F, Zhang S-P (2011a) Association of 5-HTT gene
705       polymorphisms with migraine: a systematic review and meta-analysis. J Neurol Sci 305:57–
706       66.

707   Liu J, Sun Q, Tang B, Hu L, Yu R, Wang L, Shi C, Yan X, Pan Q, Xia K (2011b) PITX3 gene polymorphism
708       is associated with Parkinson's disease in Chinese population. Brain Res 1392:116–120.

709   Lubke GH, Muthén B (2005) Investigating population heterogeneity with factor mixture models.
710       Psychol Methods 10:21.

711   MacKillop J, Amlung MT, Few LR, Ray LA, Sweet LH, Munafò MR (2011) Delayed reward discounting
712       and addictive behavior: a meta-analysis. Psychopharmacology (Berl) 216:305–321.

713   Maneeton N, Maneeton B, Srisurapanont M, Martin SD (2011) Bupropion for adults with attention-
714       deficit hyperactivity disorder: Meta-analysis of randomized, placebo-controlled trials.
715       Psychiatry Clin Neurosci 65:611–617.

716   Murphy KP (2012) Machine learning: a probabilistic perspective. MIT press.

717   Nosek BA, Spies JR, Motyl M (2012) Scientific utopia II. Restructuring incentives and practices to
718       promote truth over publishability. Perspect Psychol Sci 7:615–631.

719   Ohi K, Hashimoto R, Yasuda Y, Fukumoto M, Yamamori H, Umeda-Yano S, Kamino K, Ikezawa K,
720       Azechi M, Iwase M (2011) The SIGMAR1 gene is associated with a risk of schizophrenia and
721       activation of the prefrontal cortex. Prog Neuropsychopharmacol Biol Psychiatry 35:1309–
722       1315.

723   O'Keefe DJ (2007) Brief report: post hoc power, observed power, a priori power, retrospective
724       power, prospective power, achieved power: sorting out appropriate uses of statistical power
725       analyses. Commun Methods Meas 1:291–299.

726    Olabi B, Ellison-Wright I, McIntosh AM, Wood SJ, Bullmore E, Lawrie SM (2011) Are there progressive
727          brain changes in schizophrenia? A meta-analysis of structural magnetic resonance imaging
728          studies. Biol Psychiatry 70:88–96.

729    Oldershaw A, Hambrook D, Stahl D, Tchanturia K, Treasure J, Schmidt U (2011) The socio-emotional
730          processing stream in anorexia nervosa. Neurosci Biobehav Rev 35:970–988.

731    Oliver BJ, Kohli E, Kasper LH (2011) Interferon therapy in relapsing-remitting multiple sclerosis: a
732          systematic review and meta-analysis of the comparative trials. J Neurol Sci 302:96–105.

733    Open Science Collaboration (2015) Estimating the reproducibility of psychological science. Science
734          349:aac4716.

735    Peerbooms OL, van Os J, Drukker M, Kenis G, Hoogveld L, De Hert M, Delespaul P, van Winkel R,
736          Rutten BP (2011) Meta-analysis of MTHFR gene variants in schizophrenia, bipolar disorder
737          and unipolar depressive disorder: evidence for a common genetic vulnerability? Brain Behav
738          Immun 25:1530–1543.

739    Pizzagalli DA (2011) Frontocingulate dysfunction in depression: toward biomarkers of treatment
740          response. Neuropsychopharmacology 36:183–206.

741    Quinlan PT (2013) Misuse of power: in defence of small-scale science. Nat Rev Neurosci 14:585–585.

742    Rist PM, Diener H-C, Kurth T, Schürks M (2011) Migraine, migraine aura, and cervical artery
743          dissection: a systematic review and meta-analysis. Cephalalgia 31:886–896.

744    Samworth R, Wand M (2010) Asymptotics and optimal bandwidth selection for highest density
745          region estimation. Ann Stat 38:1767–1792.

746    Sena ES, Van Der Worp HB, Bath PM, Howells DW, Macleod MR (2010) Publication bias in reports of
747          animal stroke studies leads to major overstatement of efficacy. PLoS Biol 8:e1000344.

748    Sexton CE, Kalu UG, Filippini N, Mackay CE, Ebmeier KP (2011) A meta-analysis of diffusion tensor
749          imaging in mild cognitive impairment and Alzheimer's disease. Neurobiol Aging 32:2322-e5.

750    Shum D, Levin H, Chan RC (2011) Prospective memory in patients with closed head injury: a review.
751          Neuropsychologia 49:2156–2165.

752    Sim H, Shin B-C, Lee MS, Jung A, Lee H, Ernst E (2011) Acupuncture for carpal tunnel syndrome: a
753          systematic review of randomized controlled trials. J Pain 12:307–314.

754    Song F, Poljak A, Valenzuela M, Mayeux R, Smythe GA, Sachdev PS (2011) Meta-analysis of plasma
755          amyloid-β levels in Alzheimer's disease. J Alzheimers Dis 26:365–375.

756    Sun Q, Fu Y, Sun A, Shou Y, Zheng M, Li X, Fan D (2011) Correlation of E-selectin gene polymorphisms
757          with risk of ischemic stroke A meta-analysis. Neural Regen Res 6.

758    Tian Y, Kang L, Wang H, Liu Z (2011) Meta-analysis of transcranial magnetic stimulation to treat post-
759          stroke dysfunction. Neural Regen Res 6.

760    Trzesniak C, Kempton MJ, Busatto GF, de Oliveira IR, Galvao-de Almeida A, Kambeitz J, Ferrari MCF,
761          Santos Filho A, Chagas MH, Zuardi AW (2011) Adhesio interthalamica alterations in

762     schizophrenia spectrum disorders: A systematic review and meta-analysis. Prog
763     Neuropsychopharmacol Biol Psychiatry 35:877–886.

764  Veehof MM, Oskam M-J, Schreurs KM, Bohlmeijer ET (2011) Acceptance-based interventions for the
765     treatment of chronic pain: a systematic review and meta-analysis. PAIN® 152:533–542.

766  Vergouwen MD, Etminan N, Ilodigwe D, Macdonald RL (2011) Lower incidence of cerebral infarction
767     correlates with improved functional outcome after aneurysmal subarachnoid hemorrhage. J
768     Cereb Blood Flow Metab 31:1545–1553.

769  Vieta E, Günther O, Locklear J, Ekman M, Miltenburger C, Chatterton ML, Åström M, Paulsson B
770     (2011) Effectiveness of psychotropic medications in the maintenance phase of bipolar
771     disorder: a meta-analysis of randomized controlled trials. Int J Neuropsychopharmacol
772     14:1029–1049.

773  Wand MP, Marron JS, Ruppert D (1991) Transformations in density estimation. J Am Stat Assoc
774     86:343–353.

775  Wisdom NM, Callahan JL, Hawkins KA (2011) The effects of apolipoprotein E on non-impaired
776     cognitive functioning: a meta-analysis. Neurobiol Aging 32:63–74.

777  Witteman J, van IJzendoorn MH, van de Velde D, van Heuven VJ, Schiller NO (2011) The nature of
778     hemispheric specialization for linguistic and emotional prosodic perception: a meta-analysis
779     of the lesion literature. Neuropsychologia 49:3722–3738.

780  Woon F, Hedges DW (2011) Gender does not moderate hippocampal volume deficits in adults with
781     posttraumatic stress disorder: A meta-analysis. Hippocampus 21:243–252.

782  Xuan C, Zhang B-B, Li M, Deng K-F, Yang T, Zhang X-E (2011) No association between APOE epsilon 4
783     allele and multiple sclerosis susceptibility: a meta-analysis from 5472 cases and 4727
784     controls. J Neurol Sci 308:110–116.

785  Yang W, Kong F, Liu M, Hao Z (2011a) Systematic review of risk factors for progressive ischemic
786     stroke. Neural Regen Res 6:346–352.

787  Yang Z, Li W, Huang T, Chen J, Zhang X (2011b) Meta-analysis of Ginkgo biloba extract for the
788     treatment of Alzheimer's disease. Neural Regen Res 6:1125–1129.

789  Yuan H, Yang X, Kang H, Cheng Y, Ren H, Wang X (2011) Meta-analysis of tau genetic polymorphism
790     and sporadic progressive supranuclear palsy susceptibility. Neural Regen Res 6:353–359.

791  Zafar SN, Iqbal A, Farez MF, Kamatkar S, de Moya MA (2011) Intensive insulin therapy in brain injury:
792     a meta-analysis. J Neurotrauma 28:1307–1317.

793  Zhang Y, Zhang J, Tian C, Xiao Y, Li X, He C, Huang J, Fan H (2011) The− 1082G/A polymorphism in IL-
794     10 gene is associated with risk of Alzheimer's disease: a meta-analysis. J Neurol Sci 303:133–
795     138.

796  Zhu Y, He Z-Y, Liu H-N (2011) Meta-analysis of the relationship between homocysteine, vitamin B 12,
797     folate, and multiple sclerosis. J Clin Neurosci 18:933–938.

798

799

**All studies
(N = 730)**

**GMM**

Studies excluding null
(N=638)

Genetic (N = 234)

Treatments (N = 145)

Psychology (N = 198)

Brain imaging (N = 65)

Neurochemistry (N = 50)

Miscellaneous (N = 38)

Studies excluding genetic
(N = 496)

Studies excluding treatments
(N = 595)

Studies excluding psychology
(N = 532)

Studies excluding brain imaging
(N = 665)

Studies excluding
neurochemistry (N = 680)

Studies excluding miscellaneous
(N = 692)

**a** All data

**b** All data excluding Null

- BIC metric without regularization
- BIC metric with regularization
- Best model without regulazitation
- Best model using regularization
- Best model using Variational Bayes Dirichlet Process GMM

**c** Genetic Category

**d** Psychology Category

**e** Neurochemistry Category

**f** Treatment Category

**g** Imaging Category

**h** Miscellaneous Category

**a** Data: Original Study



**b** Data: Data excluding Null



C1
44.40%

C2
27.28%

C3
15.81%

C4
12.51%

C1
38.29%

C2
30.39%

C3
17.29%

C4
14.03%

C1    C2    C3    C4

2% 4%
5%
25%    17%
37%    10%

2% 6%
8%
17%    19%
25%    23%

11% 10%
12%    16%
46%    6%

13% 1% 11%
5%
12%    7%
51%

C1    C2    C3    C4

2% 6%
6%
49%    22%
15%

6% 8%
16%    20%
24%    26%

10% 10%
12%    16%
45%    6%

13% 11%
5%
12%    7%
51%

Null Studies    Treatment Studies    Brain Imaging Studies
Genetic Studies    Psychology Studies    Neurochemistry Studies    Miscellaneous Studies

Genetic Studies    Psychology Studies    Neurochemistry Studies
Treatment Studies    Brain Imaging Studies    Miscellaneous Studies

**c** Data: Genetic Studies



**d** Data: Psychology Studies



**e** Data: Neurochemistry Studies



**f** Data: Treatment Studies



**g** Data: Imaging Studies



**h** Data: Miscellaneous Studies

**a** Data: All data minus Genetics subfield
C1 30.91%
C2 39.01%
C3 16.93%
C4 13.15%

**b** Data: All data minus Psychology subfield
C1 49.97%
C2 29.36%
C3 12.99%
C4 7.67%

**c** Data: All data minus Neurochemistry subfield
C1 45.72%
C2 26.62%
C3 14.37%
C4 13.29%

**d** Data: All data minus Treatment subfield
C1 44.08%
C2 20.66%
C3 20.18%
C4 15.08%

**e** Data: All data minus Imaging subfield
C1 47.35%
C2 18.86%
C3 20.68%
C4 13.11%

**f** Data: All data minus Misc. subfield
C1 46.21%
C2 26.68%
C3 15.57%
C4 11.54%

**a**

BENCHMARK POWER DISTRIBUTION

Easy field, low power
~50% power, 100% true findings



**b**

BENCHMARK POWER DISTRIBUTION

Easy field, excellent power
~90% power, 100% true findings



**c**

BENCHMARK POWER DISTRIBUTION

Hard field, low power
~50% power, only 25% true findings



**d**

BENCHMARK POWER DISTRIBUTION

Hard field, excellent power
~90% power, only 25% true findings