

# Evidence for the null hypothesis in functional magnetic resonance imaging using group-level Bayesian inference

Ruslan Masharipov, Yaroslav Nikolaev, Alexander Korotkov, Michael Didur, Denis Cherednichenko and Maxim Kireev\*

1 - N.P. Bechtereva Institute of the Human Brain, Russian Academy of Sciences, Akademika Pavlova Street 9, St. Petersburg, 197376, Russia

\*Corresponding author: Maxim Kireev;

Address: Akademika Pavlova Street 9, St. Petersburg, 197376, Russia

E-mail: [kireev@ihb.spb.ru](mailto:kireev@ihb.spb.ru)

Phone number: +7 812 670-09-51

Fax number: +7 812 2343247

## Abstract

Classical null hypothesis significance testing is limited to the rejection of the point-null hypothesis; it does not allow the interpretation of non-significant results. Moreover, studies with a sufficiently large sample size will find statistically significant results even when the effect is negligible and may be considered practically equivalent to the ‘null effect’. This leads to a publication bias against the null hypothesis. There are two main approaches to assess ‘null effects’: shifting from the point-null to the interval-null hypothesis and considering the practical significance in the frequentist approach; using the Bayesian parameter inference based on posterior probabilities, or the Bayesian model inference based on Bayes factors. Herein, we discuss these statistical methods with particular focus on the application of the Bayesian parameter inference, as it is conceptually connected to both frequentist and Bayesian model inferences. Although Bayesian methods have been theoretically elaborated and implemented in commonly used neuroimaging software, they are not widely used for ‘null effect’ assessment. To demonstrate the advantages of using the Bayesian parameter inference, we compared it with classical null hypothesis significance testing for fMRI data group analysis. We also consider the problem of choosing a threshold for a practically significant effect and discuss possible applications of Bayesian parameter inference in fMRI studies. We argue that Bayesian inference, which directly provides evidence for both the null and alternative hypotheses, may be more intuitive and convenient for practical use than frequentist inference, which only provides evidence against the null hypothesis. Moreover, it may indicate that the obtained data are not sufficient to make a confident inference. Because interim analysis is easy to perform using Bayesian inference, one can evaluate the data as the sample size increases and decide to terminate the experiment if the obtained data are sufficient to make a confident inference. To facilitate the application of the Bayesian parameter inference to ‘null effect’ assessment, scripts with a simple GUI were developed.

**Keywords:** Null effects, Practical significance, Practical equivalence, Bayesian inference, fMRI

## 1. Introduction

In the neuroimaging field, it is a common practice to identify statistically significant differences in local brain activity using the general linear model approach for mass-univariate null hypothesis significance testing (NHST) (Friston et al., 1994). NHST considers the probability of obtaining the observed data, or more extreme data, given that the null hypothesis of no difference is true. This probability, or p-value, of 0.01, means that, on average, in one out of 100 ‘hypothetical’ replications of the experiment, we find a difference no less than the one found under the null hypothesis. We conventionally suppose that this is unlikely, therefore, we ‘reject the null’; that is, NHST employs ‘proof by contradiction’ (Cohen, 1994). Conversely, when the p-value is large, it is tempting to ‘accept the null’. However, the absence of evidence is not evidence of absence (Altman and Bland, 1995). Using NHST, we can only state that we have ‘failed to reject the null’. Therefore, in the classical NHST framework, the question of interpreting non-significant results remains.

The most pervasive misinterpretation of non-significant results is that they provide evidence for the null hypothesis that there is no difference, or ‘no effect’ (Nickerson, 2000; Greenland et al., 2016; Wasserstein

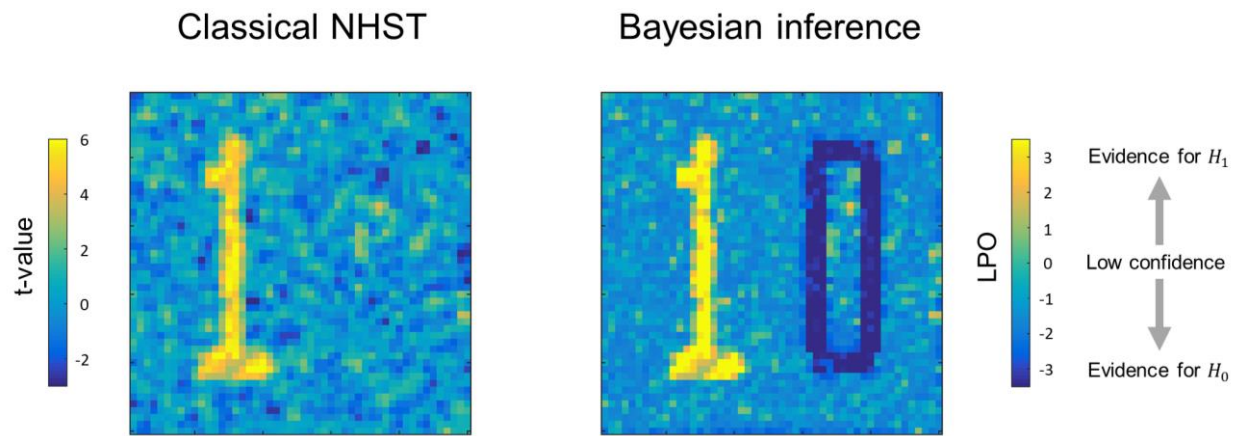
52 and Lazar, 2016). In fact, non-significant results can be obtained in two cases (Dienes, 2014): 1) the data  
53 are insufficient to distinguish the alternative from the null hypothesis, or 2) an effect is indeed null or trivial.  
54 To date, the extent to which the problem of making ‘no effect’ conclusions from non-significant results  
55 have affected the field of neuroimaging remains unclear, particularly in functional magnetic resonance  
56 imaging (fMRI) studies<sup>1</sup>. Regarding other fields of science such as psychology, neuropsychology, and  
57 biology, it was found that in 38–72% of surveyed articles, the null hypothesis was accepted based on non-  
58 significant results only (Finch et al., 2001; Schatz et al., 2005; Fidler et al., 2006; Hoekstra et al., 2006;  
59 Aczel et al., 2018).

60 Not mentioning non-significant results at all is another problem. Firstly, some authors may consider non-  
61 significant results disappointing or not worth publishing. Secondly, papers with non-significant results are  
62 less likely to be published. At the same time, NHST is usually based on the point-null hypothesis, that is,  
63 the hypothesis that the effect is *exactly* zero. However, the probability thereof is zero (Meehl, 1967; Friston  
64 et al., 2002a). This means that studies with a sufficiently large sample size will find statistically significant  
65 differences even when the effect is trivial or has no *practical* significance (Cohen, 1965, 1994; Serlin and  
66 Lapsley, 1985; Kirk, 1996). Therefore, ignoring non-significant results systematically biases our  
67 knowledge of true effects (Greenwald, 1975). This publishing bias is also known as the ‘file-drawer  
68 problem’ (Rosenthal, 1979; Ioannidis et al., 2014; De Winter and Dodou, 2015; for evidence in fMRI  
69 studies, see Jennings and Van Horn, 2012; Acar et al., 2018; David et al., 2018; Samartsidis et al., 2020).

70 Having the means to assess non-significant results would mitigate these problems. To this end, two main  
71 alternatives are available: Firstly, there are frequentist approaches that shift from point-null to interval-null  
72 hypothesis testing, for example, equivalence testing based on the two one-sided tests (TOST) procedure  
73 (Shuirmann, 1987; Wellek, 2010). Secondly, Bayesian approaches that are based on posterior parameter  
74 distributions (Lindley, 1965; Greenwald, 1975; Kruschke, 2010) and Bayes factors (Jeffreys, 1939/1948;  
75 Kass and Raftery, 1995; Rouder et al., 2009). The advantage of frequentist approaches is that they do not  
76 require a substantial paradigm shift (Campbell and Gustafson, 2018; Lakens, 2017). However, it has been  
77 argued that Bayesian approaches may be more natural and straightforward than frequentist approaches  
78 (Edwards et al., 1963; Lindley, 1975; Friston et al., 2002a; Wagenmakers, 2007; Rouder et al., 2009;  
79 Denies, 2014; Kruschke and Liddell, 2017b). It has long been noted that we tend to perceive lower p-values  
80 as stronger evidence for the alternative hypothesis, and higher p-values as evidence for the null, i.e., the  
81 ‘inverse probability’ fallacy as it is referred to by Cohen (1994). This is what we obtain in Bayesian  
82 approaches by calculating posterior probabilities. Instead of considering infinite ‘hypothetical’ replications  
83 and employing probabilistic ‘proof by contradiction’, Bayesian approaches directly provide evidence for  
84 the null and alternative hypotheses given the data, updating our prior beliefs in light of new relevant  
85 information. Bayesian inference allows us to ‘reject and accept’ the null hypothesis on an equal footing.  
86 Moreover, it allows us to talk about ‘low confidence’, indicating the need to either accumulate more data  
87 or revise the study design (see Fig. 1).

---

<sup>1</sup> Here are some examples of ‘no effect’ conclusions that can be found in the fMRI literature: a) brain area was not activated, b) brain area was not involved in the function, c) no effect was found in the brain area ( $p > 0.05$ ), d) both groups showed no differences, which can be interpreted as evidence against the alternative hypothesis; e) patients have similar responses to both conditions ( $p > 0.05$ ), that is, they have difficulties in differentiating these conditions; f) lack of significant correlation during treatment suggest a protective impact of the therapy on brain areas.



88

89 **Figure 1.** Possible results for the same data, obtained using classical NHST and Bayesian parameter  
90 inference. Classical NHST detects only areas with a statistically significant difference ('number one').  
91 Bayesian parameter inference based on the logarithm of posterior probability odds (*LPO*) provides us with  
92 additional information that is not available in classical NHST: a) it provides relative evidence for the null  
93 ( $H_0$ ) and alternative ( $H_1$ ) hypotheses, b) it detects areas with a trivial effect size ('number zero'), c) it  
94 indicates 'low confidence' areas surrounding the 'number one' and 'number zero'.

95 Despite the importance of this issue, and the high level of theoretical elaboration and implementation of  
96 Bayesian methods in common neuroimaging software programs, for example, Statistical Parametric  
97 Mapping 12 (SPM12) and FMRIB's Software Library (FSL), to date, only a few fMRI studies implemented  
98 the Bayesian inference to assess 'null effects' (for example, see subject-level analysis in Magerkurth et al.,  
99 2015, group-level analysis in Dandolo and Schwabe, 2019; Feng et al., 2019). Therefore, this study is  
100 intended to introduce fMRI practitioners to the methods for assessing 'null effects'. In particular, we focus  
101 on Bayesian parameter inference (Friston and Penny, 2003; Penny and Ridgway, 2013), as implemented in  
102 SPM12. Although Bayesian methods have been described elsewhere, the distinguishing feature of this study  
103 is that we aim to demonstrate the practical implementation of Bayesian inference to the assessment of 'null  
104 effects', and reemphasize its contributions over and above those of classical NHST. We deliberately aim to  
105 avoid mathematical details, which can be found elsewhere (Friston et al., 2002a, 2002b; Friston and Penny,  
106 2003; Penny et al., 2003, 2005, 2007, 2013). Firstly, we briefly review the frequentist and Bayesian  
107 approaches for the assessment of the 'null effect'. Next, we compare the classical NHST and Bayesian  
108 parameter inference on the Human Connectome Project (HCP) and the UCLA Consortium for  
109 Neuropsychiatric Phenomics datasets, focusing on group-level analysis. We then consider the choice of the  
110 threshold of the effect size for Bayesian parameter inference and estimate the typical effect sizes in different  
111 fMRI task designs. To demonstrate how the common sources of variability influence NHST and Bayesian  
112 parameter inference, we examine their behaviour for different sample sizes and spatial smoothing. Finally,  
113 we discuss practical research and clinical applications of Bayesian inference.

## 114 2. Theory

115 In this section, we briefly describe the classical NHST framework and review statistical methods which can  
116 be used to assess the 'null effect'. We also considered two historical trends in statistical analysis: the shift  
117 from point-null hypothesis testing to interval estimation and interval-null hypothesis testing (Murphy and  
118 Myers, 2004; Wellek, 2010; Cumming, 2013), and the shift from frequentist to Bayesian approaches  
119 (Kruschke and Liddell, 2017b).

### 120 2.1. Classical NHST framework

121 Most task-based fMRI studies rely on the general linear model approach (Friston et al., 1994; Poline and  
122 Brett, 2012). It provides a simple way to separate blood-oxygenated-level dependent (BOLD) signals  
123 associated with particular task conditions from nuisance signals and residual noise when analysing single-  
124 subject data (subject-level analysis). At the same time, it allows us to analyse BOLD signals within one  
125 group of subjects or between different groups (group-level analysis). Firstly, we must specify a general  
126 linear model and estimate its parameters ( $\beta$  values). Some of these parameters reflect the mean amplitudes

127 of BOLD responses evoked in particular task conditions or groups of subjects. The linear contrast of these  
128 parameters,  $\theta = c\beta$ , represents the experimental effect of interest (hereinafter ‘*the effect*’), expressed as the  
129 difference between two conditions or groups of subjects. Next, we test the effect against the point-null  
130 hypothesis,  $H_0: \theta = \gamma$  (usually,  $\theta = 0$ ). To do this, we use test statistics that summarise the data in a single  
131 value, for example, the t-value. For the one-sample case, the t-value is the ratio of the discrepancy of the  
132 estimated effect from the hypothetical null value to its standard error. Thus, it represents the contrast-to-  
133 noise ratio, which is similar to the signal-to-noise ratio (Welvaert and Rosseel, 2013). Finally, we calculate  
134 the probability of obtaining the observed t-value or a more extreme value, given that the null hypothesis is  
135 true (p-value). This is also commonly formulated as the probability of obtaining the observed data or more  
136 extreme data, given that the null hypothesis is true (Cohen, 1994). It can be simply written as a conditional  
137 probability  $P(D + | H_0)$ , where ‘ $D +$ ’ denotes the observed data or more extreme data which can be obtained  
138 in infinite ‘hypothetical’ replications under the null (Schneider, 2014, 2018). If this probability is lower  
139 than some conventional threshold, or alpha level (for example,  $\alpha = 0.05$ ), then we can ‘reject the null  
140 hypothesis’ and state that we found a statistically significant effect. When this procedure is repeated for a  
141 massive number of voxels, it is referred to as ‘mass-univariate analysis’. However, if we consider  $m = 100$   
142 000 voxels with no true effect and repeat significance testing for each voxel at  $\alpha = 0.05$ , we would expect  
143 to obtain 5000 false rejections of the null hypothesis (false positives). To control the number of false  
144 positives, we must reduce the alpha level for each significance test by applying the multiple comparison  
145 correction (Worsley et al., 1992; Genovese et al., 2002; Nichols and Hayasaka, 2003; Nichols, 2012).

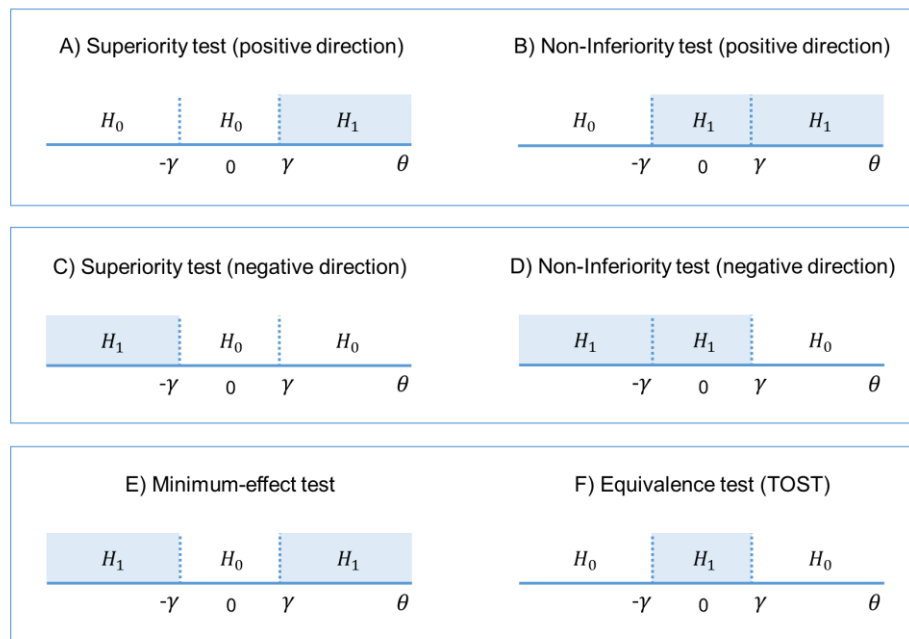
146 To date, the classical NHST has been the most widely used statistical inference method in neuroscience,  
147 psychology, and biomedicine (Szucs and Ioannidis, 2017, 2020; Ioannidis, 2019). It is often criticised for  
148 the use of the point-null hypothesis (Meehl, 1967), also known as the ‘nil null’ (Cohen, 1994) or ‘sharp  
149 null’ hypothesis (Edwards et al., 1963). It was argued that the point-null hypothesis could be appropriate  
150 only in hard sciences such as physics, but it is always false in soft sciences; this problem is sometimes  
151 known as the Meehl’s paradox (Meehl, 1967, 1978; Serlin and Lapsley, 1985, 1993; Cohen, 1994; Kirk,  
152 1996). In the case of fMRI research, we face complex brain activity which is influenced by numerous  
153 psychophysiological factors. This means that with a large amount of data, we find a statistically significant  
154 effect in all voxels for any linear contrast (Friston, 2002a). For example, Gonzalez-Castillo et al. (2012)  
155 showed a statistically significant effect in over 95% of the brain for simple visual stimulation when  
156 averaging single-subject data from 100 runs (approximately four hours of scanning). Approximately half  
157 of the brain areas showed activation, whereas the other half showed deactivation. Whole-brain  
158 (de)activations can also be found when analysing large datasets such as the HCP ( $N \sim 1000$ ) or UK Biobank  
159 ( $N \sim 10\,000$ ) datasets. When we increase the sample size, the effect estimate does not change much. Still,  
160 the standard error in the denominator of the t-value becomes increasingly smaller, resulting in negligible  
161 effects becoming statistically significant. Thus, the classical NHST ignores the magnitude of the effect.  
162 Attempts to overcome this problem led to the proposal of making a distinction between ‘statistical  
163 significance’ and ‘material significance’ (Hodges and Lehmann, 1954) or ‘practical significance’ (Cohen,  
164 1965; Kirk, 1996). That is, we can test whether the effect size is larger or smaller than some practically  
165 meaningful value using interval-null hypothesis testing (Friston et al., 2002a, 2002b, 2013).

## 166 **2.2. Frequentist approach to interval-null hypothesis testing**

167 Interval-null hypothesis testing is widely used in medicine and biology (Meyners, 2012). Consider, for  
168 example, a pharmacological study designed to compare a new treatment with an old treatment that has  
169 already shown its effectiveness. Let  $\beta_{new}$  be the mean effect on brain activity of the new treatment and  $\beta_{old}$   
170 the mean effect of the old treatment. Then,  $\theta = (\beta_{new} - \beta_{old})$  is the relative effect of the new treatment. The  
171 practical significance is defined by the effect size (ES) threshold  $\gamma$ . If a larger effect on brain activity is  
172 preferable, then we can test whether there is a practically meaningful difference in a positive direction ( $H_1:$   
173  $\theta > \gamma$  vs.  $H_0: \theta \leq \gamma$ ). This procedure is known as the *superiority test* (see Fig. 2A). We can also test whether  
174 the effect of the new treatment is no worse (practically smaller) than the effect of the old treatment ( $H_1:$   
175  $\theta > -\gamma$  vs.  $H_0: \theta \leq -\gamma$ ). This procedure is sometimes known as the *non-inferiority test* (see Fig. 2B). If a smaller  
176 effect on brain activity is preferable, we can use the superiority or non-inferiority test in the opposite  
177 direction (see Fig. 2C–D). The combination of these two superiority tests allows us to find a practically  
178 meaningful difference in both directions ( $H_1: \theta > \gamma$  and  $\theta < -\gamma$  vs.  $H_0: -\gamma \leq \theta \leq \gamma$ ), that is, the *minimum-effect*  
179 *test* (see Fig. 2E). The combination of the two non-inferiority tests allows us to reject the hypothesis of  
180 practically meaningful differences in any direction ( $H_1: -\gamma \leq \theta \leq \gamma$  vs.  $H_0: \theta > \gamma$  and  $\theta < -\gamma$ ). This is the most



181 widely used approach to *equivalence testing*, known as the *two one-sided tests* (TOST) procedure (see Fig.  
 182 2F). For more details on the superiority and minimum-effect tests, see (Serlin and Lapsey, 1985, 1993;  
 183 Murphy and Myers, 1999, 2004). For more details on the non-inferiority test and TOST procedure see  
 184 (Schuirmann, 1987; Rogers et al., 1993; Wellek, 2010; Meyners, 2012; Lakens, 2017).



185

186 **Figure 2.** The alternative ( $H_1$ ) and null ( $H_0$ ) hypotheses for different types of interval-null hypotheses tests.  
 187 A–B) One-sided tests in the positive direction (‘the larger is better’). C–D) One-sided tests in the negative  
 188 direction (‘the smaller is better’). E) Combination of both superiority tests. F) Combination of both non-  
 189 inferiority tests. Scheme modified from Aisbett et al. (2020)

190 The interval  $[-\gamma; \gamma]$  defines trivially small effect sizes that we consider to be equivalent to the ‘null effect’  
 191 for practical purposes. This interval is also known as the ‘equivalence interval’ (Schuirmann, 1987) or  
 192 ‘region of practical equivalence (ROPE)’ (Kruschke, 2011). The TOST procedure, in contrast to classical  
 193 NHST, allows us to assess the ‘null effects’. If we reject the null hypothesis of a practically meaningful  
 194 difference, we can conclude that the effect is trivially small. The TOST procedure can also be intuitively  
 195 related to frequentist interval estimates, known as confidence intervals (‘confidence interval approach’,  
 196 Westlake, 1972; Schuirmann, 1987). Confidence intervals reflect the uncertainty in the point estimation of  
 197 the parameters defined by its standard error. The confidence level of  $(1 - \alpha)$  means that among infinite  
 198 ‘hypothetical’ replications,  $(1 - \alpha)\%$  of the confidence intervals will contain the true effect under the null.  
 199 Therefore, the TOST procedure is operationally identical to considering whether the  $(1 - 2\alpha)\%$  confidence  
 200 interval falls entirely into the ROPE, as it uses two one-sided tests with an alpha level of  $\alpha$ .

201 Interval-null hypothesis testing can be used in fMRI studies not only to compare the effects of different  
 202 treatments. For example, we can apply superiority tests in the positive and negative directions to detect  
 203 ‘activated’ and ‘deactivated’ voxels and additionally apply the TOST procedure to detect ‘not activated’  
 204 voxels. However, even though we can solve the Meehl’s paradox and assess the ‘null effects’ by switching  
 205 from point-null to interval-null hypothesis testing within the frequentist approach, this approach still has  
 206 fundamental philosophical and practical difficulties which can be effectively addressed using Bayesian  
 207 statistics.

### 208 2.3. Pitfalls of the frequentist approach

209 The pitfalls of the frequentist approach have been actively discussed by statisticians and researchers for  
 210 decades. We believe that the Bayesian approach is clearer and more coherent (Lindley, 1990). Unlike the  
 211 frequentist approach, it does not require special modifications or adjustments to achieve similar practical  
 212 goals. Here, we briefly mention a few of the main problems associated with the frequency approach.

213 (1) NHST is a hybrid of Fisher’s approach that focuses on the p-value (thought to be a measure of evidence  
214 against the null hypothesis), and Neyman-Pearson’s approach that focuses on controlling false positives  
215 with the alpha level while maximising true positives in long-run replications. These two approaches are  
216 argued to be incompatible and have given rise to several misinterpretations among researchers, for example,  
217 confusing the meaning of p-values and alpha levels (Edwards et al., 1963; Gigerenzer, 1993; Goodman,  
218 1993; Finch et al., 2001; Berger, 2003; Hubbard and Bayarri, 2003; Royall, 1997; Turkheimer et al., 2004,  
219 Perezgonzalez, 2015; Schneider, 2014; Szucs and Ioannidis, 2017; Greenland, 2019).

220 (2) The logical structure of NHST is the same as that of ‘proof by contradiction’ or ‘indirect proof’, which  
221 becomes formally invalid when applied to probabilistic statements (Pollard and Richardson, 1987; Cohen,  
222 1994; Falk and Greenbaum, 1995; Nickerson, 2000; Sober, 2008; Schneider, 2014, 2018; Wagenmakers et  
223 al., 2017; but see Hagen, 1997). Valid ‘proof by contradiction’ can be expressed in syllogistic form as: 1)  
224 ‘If A, then B’ (Premise №1), 2) ‘Not B’ (Premise №2), 3) ‘Therefore not A’ (Conclusion). Probabilistic  
225 ‘proof by contradiction’ in relation to NHST can be formulated as: 1) ‘If  $H_0$  is true, then  $D +$  are highly  
226 unlikely, 2) ‘ $D +$  was obtained’, 3) ‘Therefore  $H_0$  is highly unlikely’. This problem is also referred to as the  
227 ‘illusion of probabilistic proof by contradiction’ (Falk and Greenbaum, 1995). To illustrate the fallacy of  
228 such logic, consider the following example from Pollard and Richardson (1987): 1) ‘If a person is an  
229 American ( $H_0$ ), then he is most probably not a member of Congress’, 2) ‘The person is a member of  
230 Congress’, 3) ‘Therefore the person is most probably not an American’. Based on this, one ‘rejects the null’  
231 and makes an obviously wrong inference, as only American citizens can be a member of Congress. At the  
232 same time, using Bayesian statistics, we can show that the null hypothesis (‘the person is an American’) is  
233 true (see the Bayesian solution of the ‘Congress example’ in the Supplementary Materials). The ‘illusion  
234 of probabilistic proof by contradiction’ leads to widespread confusion between the probability of obtaining  
235 the data, or more extreme data, under the null  $P(D + |H_0)$  and the probability of the null under the data  
236  $P(H_0|D)$  (Pollard and Richardson, 1987; Gigerenzer, 1993; Cohen, 1994; Falk and Greenbaum, 1995;  
237 Nickerson, 2000; Finch et al., 2001; Hoekstra et al., 2006; Goodman, 2008; Wasserstein and Lazar, 2016;  
238 Greenland, 2016; Amrhein and Roth, 2017). The latter is a posterior probability calculated based on Bayes’  
239 rule. The fact that researchers usually treat the p-value as a continuous measure of evidence (the Fisherian  
240 interpretation) only exacerbates this problem. ‘The lower the p-value, the stronger the evidence against the  
241 null’ statement can be erroneously transformed to statements such as ‘the lower the p-value, the stronger  
242 the evidence for the alternative’ or ‘the higher the p-value, the stronger the evidence for the null’. NHST  
243 can only provide evidence *against*, but never *for*, a hypothesis. In contrast, posterior probability provides  
244 direct evidence for a hypothesis; hence, it has a simple intuitive interpretation.

245 (3) The p-value is not a plausible measure of evidence (Cornfield, 1966; Royall, 1986, 1997; Berger and  
246 Sellke, 1987; Berger and Berry, 1988; Goodman, 1993; Wagenmakers et al., 2007, 2008, 2017; Hubbard  
247 and Lindsay, 2008; Johansson, 2011; Wasserstein and Lazar, 2016; but see Greenland, 2019). The  
248 frequentist approach considers infinite ‘hypothetical’ replications of the experiment (sampling  
249 distribution); that is, the p-value depends on unobserved data. One of the most prominent theorists of  
250 Bayesian statistics, Harold Jeffreys, put it as follows: ‘*What the use of P implies, therefore, is that a*  
251 *hypothesis that may be true may be rejected because it has not predicted observable results that have not*  
252 *occurred*’ (Jeffreys, 1948, p. 357). In turn, the sampling distribution depends on the researcher’s intentions.  
253 These intentions may include different kinds of *multiplicities*, such as multiple comparisons, double-sided  
254 comparisons, secondary analyses, subgroup analyses, exploratory analyses, preliminary analyses, and  
255 interim analyses of sequentially obtained data with optional stopping (Gopalan and Berry, 1998). Two  
256 researchers with different intentions may obtain different p-values based on the same dataset. The problem  
257 is that these intentions are usually unknown. When null findings are considered disappointing, it is tempting  
258 to increase the sample size until one obtains a statistically significant result. However, a statistically  
259 significant result may arise when the null is, in fact, true, which can be shown by Bayesian statistics. That  
260 is, the p-value usually exaggerates evidence against the null hypothesis. The discrepancy that may arise  
261 between frequentist and Bayesian inference is also known as the Jeffreys–Lindley paradox (Jeffreys,  
262 1939/1948; Lindley, 1957). In addition, it is argued that a consistent measure of evidence should not depend  
263 on the sample size (Cornfield, 1966). However, identical p-values provide different evidence against the  
264 null hypothesis for small and large sample sizes (Wagenmakers, 2007). In contrast, evidence provided by  
265 posterior probabilities and Bayes factors does not depend on the testing or stopping intentions or the sample  
266 size (Kruschke and Liddell, 2017b; Wagenmakers, 2007).

267 (4) Although frequentist interval estimates (Cohen, 1990, 1994; Cumming, 2013) and interval-based  
268 hypothesis testing (Murphy and Myers, 2004; Wellek, 2010; Meyners, 2012; Lakens, 2017) greatly  
269 facilitate the mitigation of the abovementioned pitfalls in data interpretation, they are still subject to some  
270 of the same types of problems as the p-values and classic NHST (Cortina and Dunlap, 1997; Nickerson,  
271 2000; Belia et al., 2005; Wagenmakers et al., 2008; Kruschke, 2013; Kruschke and Liddell, 2017a; Hoekstra  
272 et al., 2014; Morey et al., 2015; Greenland et al., 2016). Confidence intervals also depend on unobserved  
273 data and the intentions of the researcher. Moreover, the meaning of confidence intervals seems  
274 counterintuitive to many researchers. For example, one of the most common misinterpretations of the  $(1 -$   
275  $\alpha)\%$  confidence interval is that the probability of finding an effect within the confidence interval is  $(1 -$   
276  $\alpha)\%$ . In fact, it is a Bayesian interval estimate known as a *credible* interval.

277 Nevertheless, we would like to emphasise that we do not advocate abandoning the frequency approach.  
278 Correctly interpreted frequentist interval-based hypothesis testing with a priori power analysis defining the  
279 sample size and proper multiplicity adjustments often lead to conclusions similar to those of Bayesian  
280 inference (Lakens et al., 2018). However, even honest researchers with transparent intentions may find it  
281 logically and practically difficult to carry out an appropriate power analysis and make multiplicity  
282 adjustments (Berry and Hochberg, 1999; Cramer et al., 2015; Schönbrodt et al., 2015; Streiner, 2015;  
283 Sjölander and Vansteelandt, 2019). These procedures may be even more complicated in fMRI research than  
284 in psychological or social studies (see discussion on power analysis in Mumford and Nichols, 2008; Joyce  
285 and Hayasaka, 2012; Mumford, 2012; Cremes et al., 2017; Poldrack et al., 2017; multiple comparisons in  
286 Nichols and Hayasaka, 2003; Nichols, 2012; Eklund et al., 2016; and other types of multiplicities in  
287 Turkheimer et al., 2004; Chen et al., 2018, 2019, 2020; Alberton et al., 2020). For example, at the beginning  
288 of a long-term study, one may want to check whether stimulus onset timings are precisely synchronised  
289 with fMRI data collection and perform preliminary analysis on the first five subjects. The question of  
290 whether the researcher must make an adjustment for this technical check when reporting the results for the  
291 final sample become important in the frequentist approach. Such preliminary analyses (or other forms of  
292 interim analyses) are not a source of concern in Bayesian inference. Or, for example, one may want to find  
293 both ‘(de)activated’ and ‘not activated’ brain areas and use two superiority tests in combination with the  
294 TOST procedure. It is not trivial to make appropriate multiplicity adjustments in this case. In contrast,  
295 Bayesian inference suggests a single decision rule without the need for additional adjustments. Moreover,  
296 to our knowledge, practical implementations of superiority tests and the TOST procedure in common  
297 software for fMRI data analysis do not yet exist. At the same time, Bayesian analysis has already been  
298 implemented in SPM12 (<https://www.fil.ion.ucl.ac.uk/spm/software/spm12>) and is easily accessible to end-  
299 users. It consists of two steps: Bayesian parameter estimation and Bayesian inference. In general, it is not  
300 necessary to use Bayesian analysis at the subject level of analysis to apply it at the group level. One can  
301 combine computationally less demanding frequentist parameter estimation for single subjects with  
302 Bayesian estimation and inference at the group level. In the next sections, we consider the group-level  
303 Bayesian analysis implemented in SPM12.

#### 304 **2.4. Bayesian parameter estimation**

305 Bayesian statistics is based on Bayes’ rule:

$$306 \quad P(H|D) = \frac{P(D|H)P(H)}{P(D)} \quad (1)$$

307 where  $P(H|D)$  is the probability of the hypothesis given the obtained data or posterior probability.  $P(D|H)$   
308 is the probability of obtaining the *exact* data given the hypothesis or the likelihood (notice the difference  
309 from  $P(D + |H)$ , which includes *more extreme* data).  $P(H)$  is the prior probability of the hypothesis (our  
310 knowledge of the hypothesis before we obtain the data).  $P(D)$  is a normalising constant ensuring that the  
311 sum of posterior probabilities over all possible hypotheses equals one. For example, if we consider two  
312 mutually exclusive hypotheses  $H_0$  and  $H_1$ , then  $P(D) = P(D|H_0)P(H_0) + P(D|H_1)P(H_1)$  and  $P(H_0|D) +$   
313  $P(H_1|D) = 1$ .

314 In verbal form, Bayes’ rule can be expressed as:

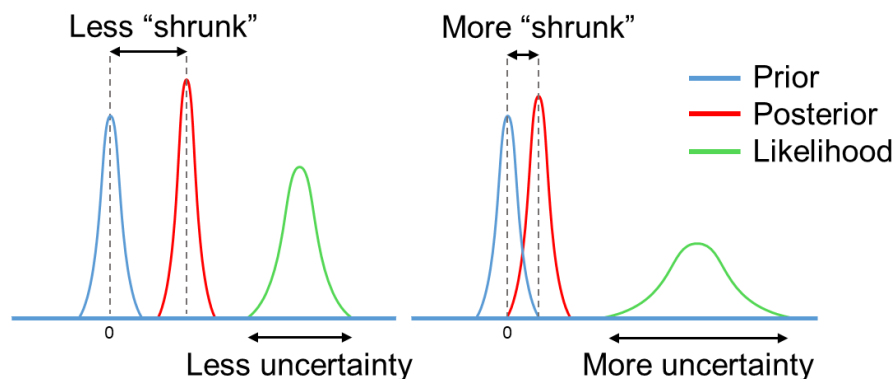
$$315 \quad \text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

316 This means that we can update our prior beliefs about the hypothesis based on the obtained data. The main  
 317 difficulty in using Bayesian statistics lies in choosing the appropriate prior assumptions. If the data are  
 318 organised hierarchically, which is the case for neuroimaging data, priors can be specified based on the data  
 319 using an empirical Bayesian approach. The lower level of the hierarchy corresponds to the experimental  
 320 effects at any given voxel, and the higher level of the hierarchy comprises the effect over all voxels. Thus,  
 321 the variance of the experimental effect over all voxels can be used as the prior variance of the effect at any  
 322 given voxel. This approach is known as the parametric empirical Bayes (PEB) with the ‘global shrinkage’  
 323 prior (Friston and Penny, 2003). The prior variance is estimated from the data under the assumption that  
 324 the prior probability density corresponds to a Gaussian distribution with zero mean. In other words, a global  
 325 experimental effect is assumed to be absent. An increase in local activity can be detected in some brain  
 326 areas; a decrease can be found in others, but the total change in neural metabolism in the whole brain is  
 327 approximately zero. This is a reasonable physiological assumption because studies of brain energy  
 328 metabolism have shown that the global metabolism is ‘remarkably constant despite widely varying mental  
 329 and motoric activity’ (Raichle and Gusnard, 2002), and ‘the changes in the global measurements of blood  
 330 flow and metabolism’ are ‘too small to be measured’ by functional imaging techniques such as PET and  
 331 fMRI (Gusnard and Raichle, 2001).

332 Now, we can rewrite Bayes’ rule (eq. 1) for the effect  $\theta = c\beta$ :

333 
$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)} \quad (2)$$

334 In the process of Bayesian updating with the ‘global shrinkage’ prior, the effect estimate ‘shrinks’ toward  
 335 zero. The greater the uncertainty of the effect estimate (variability) in a particular voxel, the less confidence  
 336 in this estimate, and the more it shrinks (see Fig. 3).



337

338 **Figure 3.** Schematic of Bayesian updating with the ‘global shrinkage’ prior. Scheme modified from  
 339 Stephan (2016).

340 After Bayesian parameter estimation, we can apply one of the two main types of Bayesian inference (Penny  
 341 and Ridgway, 2013): *Bayesian parameter inference (BPI)* or *Bayesian model inference (BMI)*. BPI is also  
 342 known as Bayesian parameter estimation (Kruschke and Liddell, 2017b). However, we deliberately  
 343 separate these two terms, as they correspond to two different steps of data analysis in SPM12. BMI is also  
 344 known as Bayesian model comparison, Bayesian model selection, or Bayesian hypothesis testing (Kruschke  
 345 and Liddell, 2017b). We chose the term BMI as it is consonant with the term BPI.

### 346 2.5. Bayesian parameter inference

347 The BPI is based on the posterior probability of finding the effect within or outside the ROPE. Let effects  
 348 larger than the ES threshold  $\gamma$  be ‘activations’, those smaller than  $-\gamma$  be ‘deactivations’, and those falling  
 349 within the ROPE  $[-\gamma; \gamma]$  be ‘no activations’. Then, we can classify voxels as ‘activated’, ‘deactivated’, or  
 350 ‘not activated’ if:

351 
$$P_{act} = P(\theta > \gamma|D) \geq P_{thr} \quad (3.1)$$



352 
$$P_{deact} = P(\theta < -\gamma | D) \geq P_{thr} \quad (3.2)$$

353 
$$P_{null} = P(-\gamma \leq \theta \leq \gamma | D) \geq P_{thr} \quad (3.3)$$

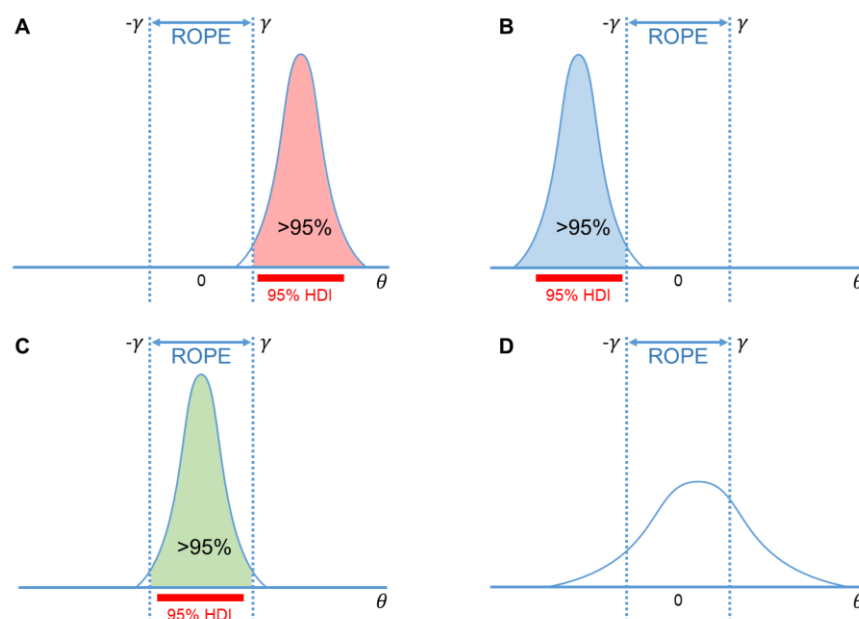
354 where  $P_{thr}$  is the posterior probability threshold (usually  $P_{thr} = 95\%$ ). Note that  $P_{act} + P_{deact} + P_{null} = 1$ .

355 If none of the above criteria are satisfied, the data in a particular voxel are insufficient to distinguish voxels  
 356 that are ‘(de)activated’ from those that are ‘not activated’. Hereinafter, we refer to them as ‘low confidence’  
 357 voxels (Magerkurth et al., 2015). This decision rule is also known as the ‘ROPE-only’ rule (Kruschke and  
 358 Liddell, 2017a, see also Greenwald (1975), Wellek (2010), Liao et al. (2019). To the best of our knowledge,  
 359 the application of this decision rule to neuroimaging data was pioneered by Friston et al. (2002a, 2002b,  
 360 2003). For convenience and visualisation purposes, we can use the natural logarithm of the posterior  
 361 probability odds (LPO), for example:

362 
$$LPO_{null} = \ln\left(\frac{P_{null}}{P_{act} + P_{deact}}\right) = \ln\left(\frac{P_{null}}{1 - P_{null}}\right) \quad (4)$$

363 This allows us to more effectively discriminate voxels with a posterior probability close to unity (Penny  
 364 and Ridgway, 2013).  $LPO_{null} > 3$  corresponds to  $P_{null} > 95\%$ . In addition,  $LPO$  also allows us to identify the  
 365 connection between BPI and BMI. The maps of the  $LPO$  are termed posterior probability maps (PPMs) in  
 366 SPM12.

367 Another possible decision rule considers the overlap between ROPE and the 95% highest density interval  
 368 (HDI). HDI is a type of credible interval (Bayesian analogue of the confidence interval), which contains  
 369 only the effects with the highest posterior probability density. If the HDI falls entirely inside the ROPE, we  
 370 can classify voxels as ‘not activated’. In contrast, if the HDI lies completely outside the ROPE, we can  
 371 classify voxels as either ‘activated’ or ‘deactivated’. If the HDI overlaps with the ROPE, we cannot make  
 372 a confident decision (we can consider them to be ‘low confidence’ voxels). This decision rule is known as  
 373 the ‘HDI+ROPE’ rule (Kruschke and Liddell, 2017a). It is more conservative than the ‘ROPE-only’ rule,  
 374 as it does not consider the effects from the tails of a posterior probability distribution. When the posterior  
 375 distributions are skewed, the ‘HDI+ROPE’ decision rule should be used. However, in SPM12, the  
 376 ‘HDI+ROPE’ and ‘ROPE-only’ decision rules should produce similar results because SPM12 utilises  
 377 normal distributions. These decision rules are illustrated in Fig. 4.



378

379 **Figure 4.** Possible variants of the posterior probability distributions of the effect  $\theta = c\beta$  in A) ‘activated’  
 380 voxels, B) ‘deactivated’ voxels, C) ‘not activated’ voxels, D) ‘low confidence’ voxels. The ‘ROPE only’

381 rule considers only the coloured parts of the distributions. The ‘HDI+ROPE’ rule considers overlap between  
 382 the ROPE and 95% HDI. Scheme modified from Magerkurth et al. (2015)

### 383 **2.6. Bayesian model inference**

384 With BPI, we consider the posterior probabilities of the linear contrast of parameters  $\theta = c\beta$ . Instead, we  
 385 can consider models using BMI.

386 Let  $H_{alt}$  and  $H_{null}$  be two non-overlapping hypotheses represented by models  $M_{alt}$  and  $M_{null}$ . These models  
 387 are defined by two parameter spaces: 1)  $M_{alt}$ :  $\theta > \gamma$  and  $\theta < -\gamma$ , and 2)  $M_{null}$ :  $-\gamma \leq \theta \leq \gamma$ .

388 Now, we can rewrite Bayes’ rule (eq. 1) for  $M_{alt}$  and  $M_{null}$

$$389 \quad P(M_{alt}|D) = \frac{P(D|M_{alt})P(M_{alt})}{P(D)} \quad (5.1)$$

$$390 \quad P(M_{null}|D) = \frac{P(D|M_{null})P(M_{null})}{P(D)} \quad (5.2)$$

391 If we divide equation (5.1) by (5.2),  $P(D)$  is cancelled out, and we obtain:

$$392 \quad \frac{P(M_{alt}|D)}{P(M_{null}|D)} = \frac{P(D|M_{alt})}{P(D|M_{null})} \frac{P(M_{alt})}{P(M_{null})} \quad (6)$$

393 In verbal form equation (6) can be expressed as:

$$394 \quad \text{Posterior Odds} = \text{Bayes Factor} \times \text{Prior Odds}$$

395 The Bayes factor ( $BF$ ) is a multiplier that converts prior model probability odds to posterior model  
 396 probability odds. It indicates the relative evidence for one model against another. For example, if  $BF_{null} =$   
 397  $\frac{p(D|M_{null})}{p(D|M_{alt})} = 2$ , then the observed data are twice as likely under the null model than under the alternative.

398 A connection exists between the BPI (eq. 2–4), and BMI (eq. 6) (see Morey and Rouder, 2011; Liao et al.,  
 399 2019):

$$400 \quad BF_{null} = \left( \frac{P(-\gamma \leq \theta \leq \gamma|D)}{1 - P(-\gamma \leq \theta \leq \gamma|D)} \right) \left( \frac{1 - P(-\gamma \leq \theta \leq \gamma)}{P(-\gamma \leq \theta \leq \gamma)} \right) \quad (7)$$

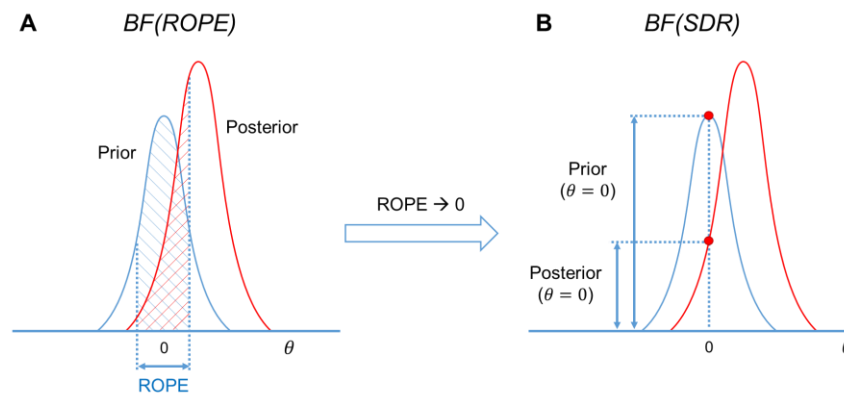
401 or, in verbal form:

$$402 \quad BF(ROPE)_{null} = \frac{\text{Posterior}(\theta \text{ in ROPE})}{\text{Posterior}(\theta \text{ outside ROPE})} \frac{\text{Prior}(\theta \text{ outside ROPE})}{\text{Prior}(\theta \text{ in ROPE})}$$

403 For convenience,  $BF$  may also be expressed in the form of a natural logarithm:

$$404 \quad \text{Log}BF(ROPE)_{null} \propto LPO_{null}$$

405 See Schematic illustration of  $BF(ROPE)_{null}$  in Fig. 5A.



406

407 **Figure 5.** Schematic of  $BF$ s used in BMI. A)  $BF(ROPE)$  is related to the areas under the functions of the  
 408 posterior and prior probability densities inside and outside the ROPE. B)  $BF(SDR)$  is the relation between  
 409 the posterior and prior probability at  $\theta = 0$ .

410 The calculation of  $BF$  may be computationally challenging, as it requires integration over the parameter  
 411 spaces. However, if the ROPE has zero width ( $\gamma = 0$ ), then the  $BF$  has an analytical solution known as the  
 412 Savage–Dickey ratio (SDR) (Wagenmakers et al., 2010; Friston and Penny, 2011; Rosa et al., 2012; Penny  
 413 and Ridgway, 2013). The interpretation of the SDR is simple: if the effect size is less likely to equal zero  
 414 after obtaining the data (posterior probability at  $\theta = 0$ ) than before (prior probability at  $\theta = 0$ ), then  
 415  $BF(SDR)_{null} < 1$ ; that is, we have more evidence for  $M_{alt}$  (see Fig. 5B).

## 416 2.7. Relations between frequentist and Bayesian approaches

417 Now we can point out the conceptual links between the frequentist and Bayesian approaches.

418 (1) **Parameter estimation.** When we have no prior information, that is, all parameter values are a priori  
 419 equally probable ('flat' prior), Bayesian parameter estimation reduces to frequentist parameter estimation  
 420 (maximum likelihood estimation; Friston et al., 2002a).

421 (2) **Multiplicity adjustments.** One of the major concerns in frequentist inference is the multiplicity  
 422 problem. In general, after the Bayesian parameter estimation, it is not necessary to classify any voxel as  
 423 '(de)activated' or 'not activated'. If we consider *unthresholded* maps of posterior probabilities, *LPOs*, or  
 424 *LogBFs*, the multiple comparisons problem does not arise (Friston and Penny, 2003). However, if we apply  
 425 a decision rule to classify voxels, we should control for wrong decisions across multiple comparisons  
 426 (Woolrich et al., 2009, see also possible loss functions in Muller et al., 2006; Kruschke and Liddell, 2017a).  
 427 The advantage of hierarchical PEB with the 'global shrinkage' prior is that it automatically accounts for  
 428 multiple comparisons without the need for ad hoc multiplicity adjustments (Berry, 1988; Friston and Penny,  
 429 2003; Scott and Berger, 2010; Gelman et al., 2012). The frequentist approach processes every voxel  
 430 independently, whereas the PEB algorithm considers joint information from all voxels. Frequentist  
 431 inference uncorrected for multiple independent comparisons is prone to label noise-driven, random  
 432 extremes as 'statistically significant'. Bayesian analysis specifies that extreme values are unlikely a priori,  
 433 and thus they shrink toward a common mean (Lindley, 1990; Westfall et al., 1997; Berry and Hochberg,  
 434 1999; Friston et al., 2002a, 2002b; Gelman et al., 2012; Kruschke and Liddell, 2017b). If we consider  
 435 *thresholded* maps of posterior probabilities, for example,  $P_{act} > 95\%$ , then as many as 5% of 'activated'  
 436 voxels could be falsely labelled so. This is conceptually similar to the false discovery rate (FDR) correction  
 437 (Berry and Hochberg, 1999; Frison et al., 2002b; Friston and Penny, 2003; Storey, 2003; Muller et al.,  
 438 2006; Schwartzman et al., 2009). In practice, BPI with  $\gamma = 0$  should produce similar results (in terms of the  
 439 number of '(de)activated' voxels) as classical NHST with FDR correction. If we increase the ES threshold,  
 440 fewer voxels will be classified as '(de)activated', and at some  $\gamma$  value, BPI will produce results similar to  
 441 the more conservative Family Wise Error (FWE) correction<sup>2</sup>.

<sup>2</sup> FDR correction controls the rate of false discoveries (false positives in frequentist terminology) among all significant voxels. FWE correction controls the rate of any false positives in the whole brain.

442 (3) **Interval-based hypothesis testing.** Frequentist interval-based hypothesis testing is conceptually  
443 connected with BPI, particularly, the ‘HDI+ROPE’ decision rule. The former considers the intersection  
444 between ROPE and the confidence intervals. The latter considers the intersection between ROPE and the  
445 HDI (credible intervals).

446 (4) **BPI and BMI.** BMI based on  $BF(ROPE)$  is conceptually linked to BPI based on the ‘ROPE-only’  
447 decision rule. The interval-based Bayes factor  $BF(ROPE)$  is proportional to the posterior probability odds.  
448 When ROPE is infinitesimally narrow,  $BF$  can be approximated using the  $SDR$ . Note that even though  
449  $BF(SDR)$  is based on the point-null hypothesis, it can still provide evidence for the null hypothesis, in  
450 contrast to BPI with  $\gamma = 0$ . However,  $BF(SDR)$  in PEB settings has not yet been tested using empirical fMRI  
451 data. Because the point-null hypothesis is always false (Meehl, 1967), BPI and  $BF(ROPE)$  are preferred  
452 over  $BF(SDR)$ .

### 453 **2.8. Definition of the effect size threshold**

454 The main difficulty in applying interval-based methods is the choice of the ES threshold  $\gamma$ . To date, only a  
455 few studies have been devoted to determining the minimal relevant effect size. One of them suggested a  
456 method to objectively define  $\gamma$  at the subject level of analysis which was calibrated by clinical experts and  
457 may be implemented for pre-surgical planning (Magerkurth et al., 2015). At the same time, the problem of  
458 choosing the ES threshold  $\gamma$  for group-level Bayesian analysis remains unresolved.

459 Several ways in which to define the ES threshold are available. Firstly, we can conduct a pilot study to  
460 determine the expected effect sizes. Secondly, we can use data from the literature to determine the typical  
461 effect sizes for the condition of interest. Thirdly, we can use the default ES thresholds that are commonly  
462 accepted in the field. One of the first ES thresholds proposed in the neuroimaging literature was  $\gamma = 0.1\%$   
463 (Friston et al., 2002b). This is the default ES threshold for the subject-level BPI in SPM12. For the group-  
464 level BPI, the default ES threshold is one prior standard deviation of the effect  $\gamma = 1 \text{ prior } SD_{\theta}$  (Friston and  
465 Penny, 2003). Fourthly,  $\gamma$  can be selected in such a way as to ensure maximum similarity of the activation  
466 patterns revealed by classical NHST and Bayesian inference. This would allow us to reanalyse the data  
467 using Bayesian inference, reveal similar activation patterns as was previously the case for classic inference,  
468 and detect the ‘not activated’ and ‘low confidence’ voxels. Lastly, we can consider the posterior  
469 probabilities at multiple ES thresholds or compute the ROPE maps (see below).

470 It should also be noted that the ES threshold can be expressed as unstandardised (raw  $\beta$  values or percent  
471 signal change) and standardised values (for example, Cohen’s  $d$ ). Raw  $\beta$  values calculated by SPM12 at the  
472 first level of analysis represent the BOLD signal in arbitrary units. However, they can be scaled to a more  
473 meaningful unit, the BOLD percent signal change (PSC) (Poldrack et al., 2011; Chen et al., 2017).  
474 Unstandardised and standardised values have disadvantages and advantages. Different ways exist in which  
475 to scale  $\beta$  values to PSC (Pernet, 2014; Chen et al., 2017), which is problematic when comparing the results  
476 of different studies. Standardised values represent the effect size in terms of the standard deviation units,  
477 which supposedly facilitate the comparison of results between different experiments. However,  
478 standardised values are relatively more unstable between measurements and less interpretable (Baguley,  
479 2009; Chen et al., 2017). Moreover, Cohen’s  $d$  is closely related to the t-value (for one sample case,  $d =$   
480  $t/\sqrt{N}$ ) and may share some drawbacks with t-values. Reimold et al. (2005) showed that spatial smoothing  
481 has a nonlinear effect on voxel variance. Using t-values or Cohen’s  $d$  for inference in neuroimaging may  
482 lead to spatially inaccurate results (spatial shift of local maxima in t-maps or Cohen’s  $d$  maps compared to  
483 PSC-maps). In this study, we focused on PSCs.

## 484 **3. Methods**

485 We used the HCP and UCLA datasets to: 1) compare classical NHST and BPI, 2) consider different  
486 approaches to ES thresholding and estimate typical effect sizes, 3) demonstrate the behaviour of classical  
487 NHST and BPI depending on the sample size and spatial smoothing, and 4) illustrate a possible practical  
488 application of BPI.

### 489 **3.1. Datasets**



490 Seven block-design tasks were considered from the HCP dataset, including working memory, gambling,  
491 motor, language, social cognition, relation processing, and emotion processing tasks (Barch et al., 2013).  
492 Two event-related tasks, including the stop-signal and task-switching tasks were considered from the  
493 UCLA dataset (Poldrack et al., 2016). The length, conditions, and number of scans of the tasks are provided  
494 in the Supplementary Materials (Table S1). A subset of 100 unrelated subjects (S1200 release) was selected  
495 from the HCP dataset (54 females, 46 males, mean age =  $29.1 \pm 3.7$  years) for assessment. A total of 115  
496 subjects from the UCLA dataset were included in the analysis (55 females, 60 males, mean age =  $31.7 \pm$   
497  $8.9$  years) after removing subjects with no data for the stop-signal task, a high level ( $>15\%$ ) of errors in the  
498 Go-trials, and those of which the raw data were reported to be problematic (Gorgolewski et al., 2017). See  
499 the fMRI acquisition parameters in the Supplementary Materials, Par. 1.

### 500 **3.2. Preprocessing**

501 The minimal preprocessing pipelines for the HCP and UCLA datasets were described by Glasser et al.  
502 (2013) and Gorgolewski et al. (2017), respectively. Spatial smoothing was applied to the preprocessed  
503 images with a 4 mm full width at half maximum (FWHM) Gaussian smoothing kernel. Additionally, to  
504 compare the extent to which the performance of classical NHST and BPI depended on the smoothing, we  
505 applied 8 mm FWHM smoothing to the emotion processing task. Spatial smoothing was performed using  
506 SPM12. The results are reported for the 4 mm FWHM smoothing filter, unless otherwise specified.

### 507 **3.3. Parameter estimation**

508 Frequentist parameter estimation was applied at the subject level of analysis. A detailed description of the  
509 general linear models for each task design is available in the Supplementary Materials, Par. 2. Fixation  
510 blocks and null events were not modelled explicitly in any of the tasks. Twenty-four head motion regressors  
511 were included in each subject-level model (six head motion parameters, six head motion parameters one  
512 time point before, and 12 corresponding squared items) to minimise head motion artefacts (Friston et al.,  
513 1996). Raw  $\beta$  values were converted to PSC relative to the mean whole-brain ‘baseline’ signal  
514 (Supplementary Materials, Par. 3). The linear contrasts of the  $\beta$  values were calculated to describe the  
515 effects of interest  $\theta = c\beta$  in different tasks. The sum of positive terms in the contrast vector,  $c$ , is equal to  
516 one. The list of contrasts calculated in the current study to explore typical effect sizes is presented in Table  
517 S1. At the group level of analysis, the Bayesian parameter estimation with the ‘global shrinkage’ prior was  
518 applied using SPM12 (v6906). We performed a one-sample test on the linear contrasts created at the subject  
519 level of analysis.

### 520 **3.4. Classical NHST and Bayesian parameter inference**

521 Classical inference was performed using voxel-wise FWE correction with  $\alpha = 0.05$ . This is the default SPM  
522 threshold and is known to be conservative and to guarantee protection from false positives (Eklund et al.,  
523 2016). Although voxel-wise FWE correction may be too conservative for small sample sizes, it is  
524 recommended when large sample sizes are available (Woo et al., 2014).

525 BPI, accessible via the SPM12 GUI, allows the user to declare only whether the voxels are ‘activated’ or  
526 ‘deactivated’. The classification of voxels as being either ‘not activated’ or ‘low confidence’ requires the  
527 posterior mean and variance. At the group level of analysis, SPM12 does not save the posterior variance  
528 image. However, the posterior variance can be reconstructed from the image of the noise hyperparameter  
529 using a first-order Taylor series approximation (Penny and Ridgway, 2013). Therefore, in the current study,  
530 BPI was performed using in-house scripts available at  
531 ([https://github.com/Masharipov/Bayesian\\_inference](https://github.com/Masharipov/Bayesian_inference)). For the ‘ROPE-only’ rule, the posterior probability  
532 threshold was  $P_{thr} = 95\%$  ( $LPO > 3$ ). For the ‘HDI+ROPE’ rule, we used 95% HDI.

533 We compared the number of ‘activated’ voxels (as a percentage of the total number of voxels) detected by  
534 Bayesian and classical inference. We also compared the number of ‘activated,’ ‘deactivated,’ and ‘not  
535 activated’ voxels detected using BPI with the ‘ROPE-only’ and ‘HDI+ROPE’ decision rules and different  
536 ES thresholds. To estimate the influence of the sample size on the results, all the above-mentioned analyses  
537 were performed with samples of different sizes: 5 to 100 subjects from the HCP dataset and 5 to 115 subjects  
538 from the UCLA dataset, in steps of 5 subjects. Ten random groups were sampled for each step.

### 539 3.5. Effect size thresholds

540 We considered three ES thresholds: firstly, the default ES threshold for the subject-level  $\gamma = 0.1\%$  (BOLD  
541 PSC); secondly, the default ES threshold for the group-level  $\gamma = 1 \text{ prior } SD_{\theta}$ ; thirdly, the  $\gamma(Dice_{max})$   
542 threshold, which ensures maximum similarity of the activation patterns revealed by classical NHST and  
543 BPI. The similarity was assessed using the Dice coefficient:

$$544 \quad Dice(\gamma) = \frac{2 * V_{overlap}(\gamma)}{V_{classic} + V_{bayesian}(\gamma)} \quad (8)$$

545 where  $V_{classic}$  is the number of ‘activated’ voxels detected using classical NHST,  $V_{bayesian}(\gamma)$  is the  
546 number of ‘activated’ voxels detected using BMI with the ES threshold  $\gamma$ , and  $V_{overlap}$  is the number of  
547 ‘activated’ voxels detected by both methods. A Dice coefficient of 0 indicates no overlap between the  
548 patterns, and 1 indicates complete overlap. Dice coefficients were calculated for  $\gamma$  ranging from 0% to 0.4%  
549 in steps of 0.001%.

### 550 3.6. Estimation of typical effect sizes

551 In the current study, we aimed to provide a reference set of typical effect sizes for different task designs  
552 (block and event-related) and different contrasts (‘task-condition > control-condition,’ ‘task-condition >  
553 baseline,’) in a set of a priori defined regions of interest (ROI). Effect sizes were expressed in PSC and  
554 Cohen’s d. ROI masks were defined using anatomical and a priori functional masks. For more details, see  
555 Supplementary Materials, Par. 4.

### 556 3.7. Evaluating BPI on contrasts with no expected practically significant difference

557 BPI should be able to detect the ‘null effect’ in the majority of voxels when comparing samples with no  
558 expected practically significant difference. For example, there may be two groups of healthy adult subjects  
559 performing the same task or two sessions with the same task instructions. To test this, we used fMRI data  
560 from the emotion processing task. To emulate two ‘similar’ *independent* samples, 100 healthy adult  
561 subjects’ contrasts (‘Emotion > Shape’) were randomly divided into two groups of 50 subjects. A two-  
562 sample test comparing the ‘Group #1’ and ‘Group #2’ was performed with the assumption of unequal  
563 variances between the groups (SPM12 default option). To emulate ‘similar’ *dependent* samples, we  
564 randomized ‘Emotion > Shape’ contrasts from right-to-left (RL) and left-to-right (LR) phase encoding  
565 sessions in the ‘Session #1’ and ‘Session #2’ samples. Each sample consisted of 50 contrasts from the RL  
566 session and 50 from the LR session. A paired test designed to compare ‘Session #1’ and ‘Session #2’ was  
567 equivalent to the one-sample test on 50 ‘RL > LR session’ and 50 ‘LR > RL session’ contrasts.

## 568 4. Results

### 569 4.1. Results for contrasts with no expected practically significant difference

570 Classical NHST did not show a significant difference between ‘Group #1’ and ‘Group #2’ (see  
571 Supplementary Materials, Fig. S1). BPI with the ‘ROPE-only’ decision rule and default ES threshold  $\gamma = 1$   
572  $\text{prior } SD_{\theta} = 0.190\%$  classified 83.4% of voxels as having ‘no difference’ in which the null hypothesis was  
573 accepted (see Fig. S1). The ‘HDI+ROPE’ rule classified 76.2% of voxels as having ‘no difference’.

574 Classical NHST did not reveal a significant difference between ‘Session #1’ and ‘Session #2’ (see Fig. S2).  
575 The  $\text{prior } SD_{\theta}$  was 0.005%. In this case, using the default ES threshold  $\gamma = 1 \text{ prior } SD_{\theta}$  did not allow the  
576 detection of any ‘no difference’ voxels, because the ROPE was unreasonably narrow. The ‘null effect’ was  
577 detected in all voxels beginning with a  $\gamma = 0.013\%$  threshold using the ‘ROPE-only’ and ‘HDI+ROPE’  
578 decision rules (see Fig. S2).

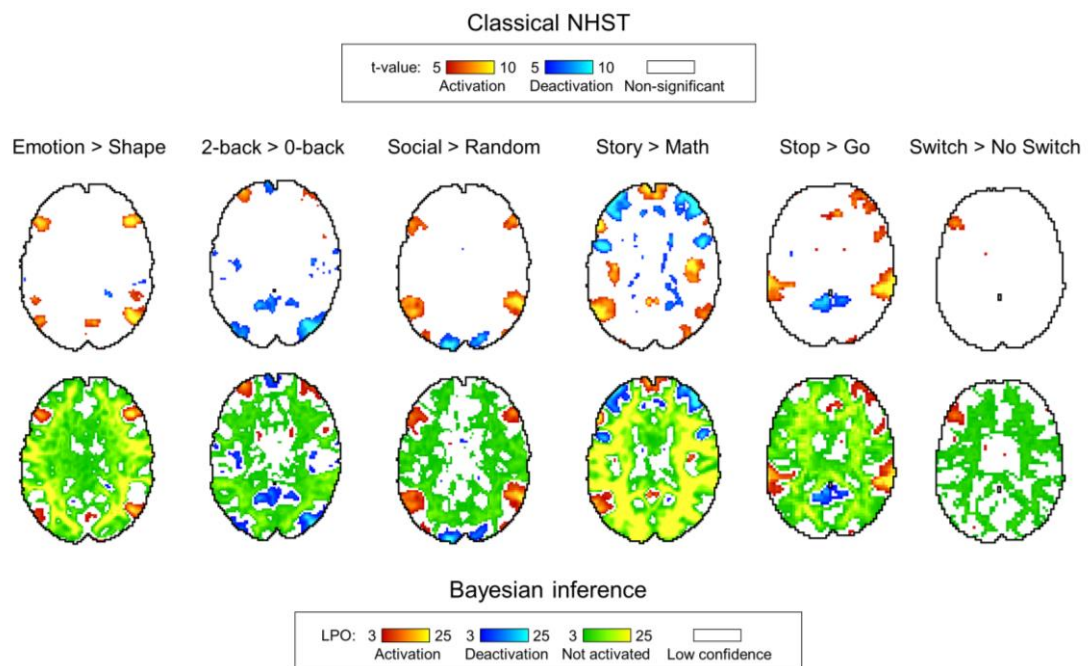
579 In this way, when comparing two ‘similar’ *independent* samples (two groups of healthy subjects performing  
580 the same task), BPI with the default group-level threshold ( $1 \text{ prior } SD_{\theta}$ ) allowed us to correctly label  
581 voxels as having ‘no difference’ for the majority of the voxels of the brain. However, when comparing two

582 ‘similar’ *dependent* samples (two sessions from the same task), the *one prior*  $SD_{\theta}$  threshold became  
 583 inadequately small.

584 Therefore, the default *one prior*  $SD_{\theta}$  threshold is not suitable when the difference between *dependent*  
 585 conditions is very small (paired sample test or one-sample test). In such cases, one can use an a priori  
 586 defined ES threshold based on previously reported effect sizes or provide an ES threshold at which most of  
 587 the voxels can be labelled as having ‘no difference’, allowing the critical reader to decide whether this  
 588 speaks in favour of the absence of differences.

#### 589 4.2. Comparison of classical NHST and BPI results

590 Generally, classical NHST with voxel-wise FWE correction and BPI with the ‘ROPE-only’ decision rule  
 591 and default group-level ES threshold  $\gamma = 1$  *prior*  $SD_{\theta}$  revealed similar (de)activation patterns in all  
 592 considered contrasts (see Fig. 6, Table 1, Supplementary materials, Tables S2–S10). The median ES  
 593 threshold based on  $Dice_{max}$  and median default group-level ES threshold across all considered contrasts were  
 594 close in magnitude to the default subject-level ES threshold  $\gamma = 0.1\%$ :  $\gamma(Dice_{max}) = 0.118\%$  and  $\gamma = 1$  *prior*  
 595  $SD_{\theta} = 0.142\%$ . The median  $Dice_{max}$  across all the considered contrasts reached 0.904. At the same time, BPI  
 596 allowed us to classify ‘non-significant’ voxels as ‘not activated’ or ‘low confidence’. As it can be clearly  
 597 seen from Figure 6, areas with ‘non-activated’ voxels surround clusters with ‘(de)activated’ voxels. Both  
 598 are separated by areas comprising ‘low confidence’ voxels.



599

600 **Figure 6.** Examples of results obtained with classical NHST and BPI. Six contrasts were chosen for the  
 601 illustration purposes (four event-related and two block-design tasks). Classical NHST was implemented  
 602 using voxel-wise FWE correction ( $\alpha = 0.05$ ). BPI was implemented using the ‘ROPE-only’ decision rule,  
 603  $P_{thr} = 95\%$  ( $LPO > 3$ ) and  $\gamma = 1$  *prior*  $SD_{\theta}$ . Axial slice  $z = 18$  mm (MNI152 standard space).

604 **Table 1. Maximum Dice coefficient and corresponding effect size thresholds for each task.**

Contrast, $\theta$	1 <i>prior</i> $SD_{\theta}$ , %	‘ROPE-only’ decision rule		‘HDI+ROPE’ decision rule	
		$\gamma(Dice_{max})$ , %	$Dice_{max}$	$\gamma(Dice_{max})$ , %	Dice
<b>Emotion processing</b>					
Emotion > Shape	0.135	0.116	0.904	0.104	0.912
<b>Working memory</b>					
2-back > baseline	0.325	0.136	0.925	0.125	0.932
2-back > 0-back	0.089	0.095	0.891	0.089	0.903
<b>Language</b>					

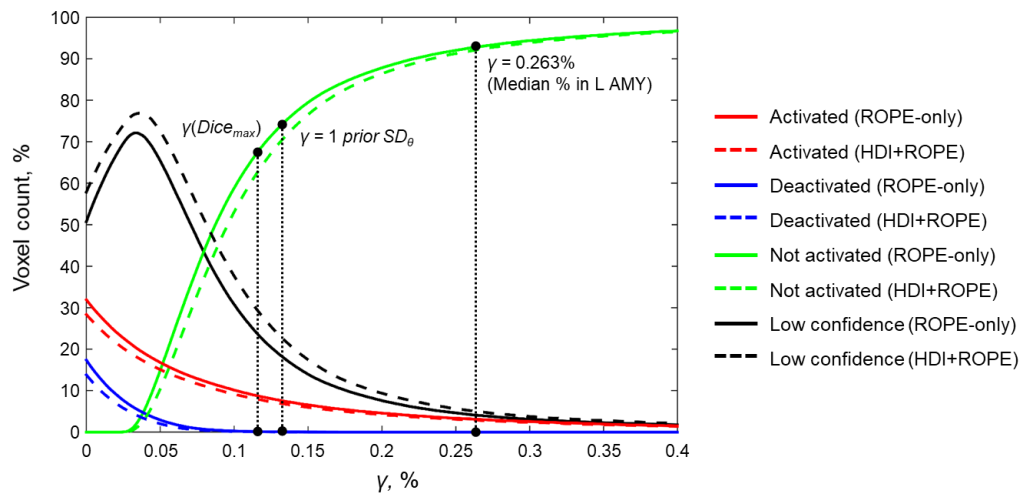
Story > Math	0.255	0.119	0.896	0.108	0.904
<b>Motor</b>					
LF > baseline	0.149	0.148	0.897	0.135	0.907
RF > baseline	0.171	0.160	0.886	0.144	0.897
T > baseline	0.268	0.205	0.904	0.181	0.913
<b>Gambling</b>					
Reward > baseline	0.254	0.132	0.917	0.122	0.924
Loss > baseline	0.249	0.134	0.918	0.118	0.925
Reward > Loss	0.032	0.044	0.894	0.037	0.886
<b>Social cognition</b>					
Social > baseline	0.325	0.139	0.939	0.124	0.944
Social > Random	0.104	0.114	0.896	0.104	0.907
<b>Relational processing</b>					
Relational > baseline	0.390	0.154	0.935	0.143	0.940
Relational > Match	0.051	0.073	0.892	0.066	0.894
<b>Stop-signal task</b>					
Correct Stop > baseline	0.069	0.066	0.895	0.061	0.906
Correct Stop > Go	0.064	0.052	0.906	0.047	0.917
<b>Task-switching</b>					
Switch > baseline	0.133	0.075	0.907	0.067	0.916
Switch > No switch	0.030	0.037	0.924	0.033	0.925
<b>Summary</b>					
Median	0.142	0.118	0.904	0.106	0.913

605 As expected, compared with the ‘HDI+ROPE’ rule, using the ‘ROPE-only’ rule slightly increases the  
 606 number of ‘(de)activated’ and ‘not activated’ voxels (see Table 1 and Tables S2-10). The ‘HDI+ROPE’  
 607 rule labelled more voxels as ‘low confidence’. In principle, the difference between these rules would only  
 608 be expected to increase with skewed distributions, which is not the case in SPM12. Thus, we recommend  
 609 using the ‘ROPE-only’ decision rule (as well as in Friston et al., 2002a, 2002b, 2003; Wellek, 2010; Liao  
 610 et al., 2019). In the following sections, we focus primarily on describing the results for the ‘ROPE-only’  
 611 rule.

#### 612 **4.3. Comparison of BPI results with different ES thresholds**

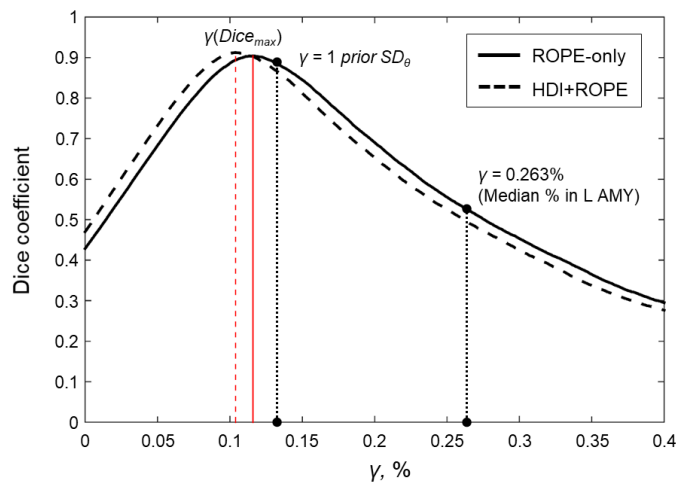
613 Here, we first consider the results for the emotional processing task and then consider other tasks. Using  
 614 the default single-subject ES threshold  $\gamma = 0.1\%$  for the emotional processing task (‘Emotion > Shape’  
 615 contrast), 58.8% of all voxels can be classified as ‘not activated,’ 30.8% as ‘low confidence,’ and 10.1% as  
 616 ‘activated’ (see Fig. 7, Table S2). The default group-level ES threshold  $\gamma = 1$  prior  $SD_{\theta} = 0.135\%$  allowed  
 617 us to classify 75.0% of all voxels as ‘non-activated,’ 17.5% as ‘low confidence,’ and 7.4% as ‘activated’  
 618 (see Fig. 7, Table S2). Both types of thresholds were comparable to those of classical NHST for the  
 619 detection of ‘activated’ voxels. The maximum overlap between ‘activations’ patterns revealed by classical  
 620 NHST and BPI was observed at  $\gamma(Dice_{max}) = 0.116\%$  (see Fig. 8, Table 1).





621

622 **Figure 7.** Number of voxels classified into the four categories depending on the ES threshold  $\gamma$ . The results  
 623 for the emotion processing task ('Emotion>Shape' contrast) are presented for illustration. Abbreviations:  
 624 L AMY – left amygdala.



625

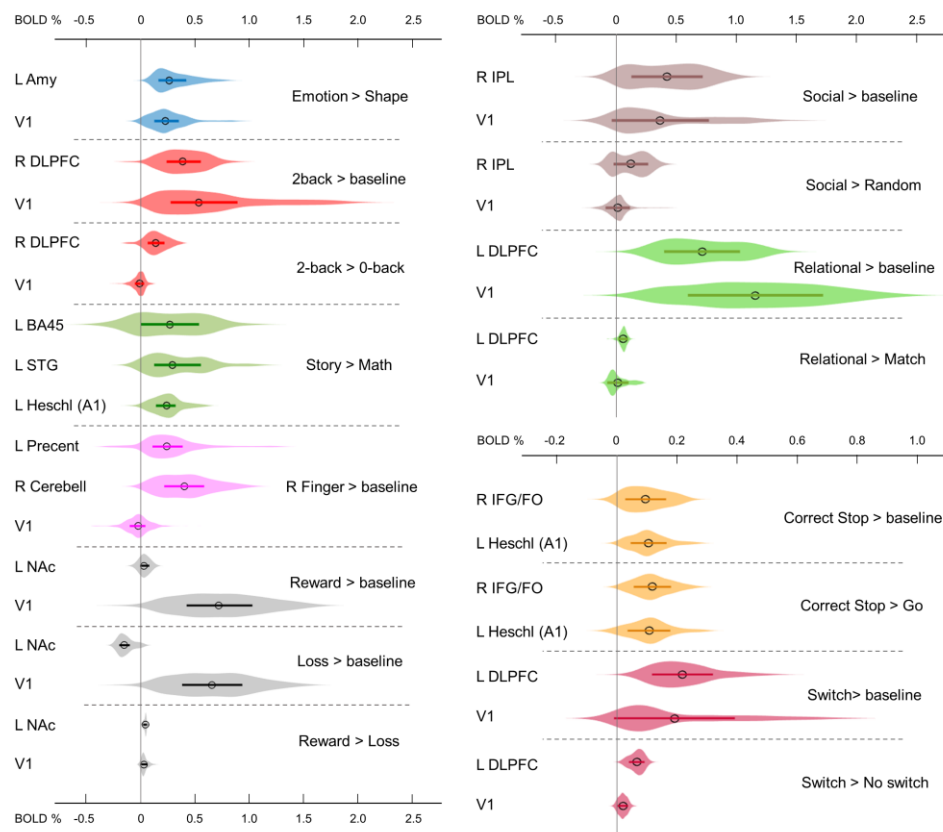
626 **Figure 8.** Dependence of the Dice coefficient on the ES threshold  $\gamma$ . Results for the emotion processing  
 627 task ('Emotion>Shape' contrast). The red lines denote  $\gamma(Dice_{max})$ . Abbreviations: L AMY – left amygdala.

628 For the '2-back > 0-back,' 'Left Finger > baseline,' 'Right Finger > baseline,' and 'Social > Random'  
 629 contrasts, the three ES thresholds that were considered—0.1%, *one prior SD<sub>θ</sub>*,  $\gamma(Dice_{max})$ —produced similar  
 630 results (see Table 1 and Tables S3, S5, S7). For the event-related stop-signal task ('Correct Stop > baseline'  
 631 and 'Correct Stop > Go' contrasts), *one prior SD<sub>θ</sub>* and  $\gamma(Dice_{max})$  were close in terms of their values but  
 632 smaller than 0.1% (see Table 1). Block designs tend to evoke higher BOLD PSC than event-related designs;  
 633 therefore, a lower *prior SD<sub>θ</sub>* should be expected for event-related designs and higher *prior SD<sub>θ</sub>* for block  
 634 designs. Within a single design, in contrasts such as 'task-condition > baseline', higher BOLD PSC and  
 635 *prior SD<sub>θ</sub>* would be expected than in contrasts in which the experimental conditions are compared directly.  
 636 For example, the contrast '2-back > baseline' has *prior SD<sub>θ</sub>* = 0.325% and contrast '2-back > 0-back' has  
 637 *prior SD<sub>θ</sub>* = 0.089%.

638 As previously noted, some contrasts did not elicit robust activations: 'Reward > Loss', 'Relational >  
 639 Match', (Barch et al., 2013) and 'Switch > No switch' (Gorgolewski et al., 2017). The corresponding  
 640  $\gamma(Dice_{max})$  thresholds were 0.044%, 0.073%, and 0.037% (see Tables 1, S6, S8, and S10). The *prior SD<sub>θ</sub>*  
 641 were 0.032%, 0.051%, and 0.030%. Correspondingly, BPI with the  $\gamma = 1$  *prior SD<sub>θ</sub>* threshold classified 0%,  
 642 18.4%, and 42.2% of voxels as 'not activated'. This demonstrates that when we compare conditions with  
 643 similar neural activity and minor differences, it becomes more difficult to separate '(de)activations' from  
 644 the 'null effects' using the  $\gamma = 1$  *prior SD<sub>θ</sub>* threshold.

#### 645 4.4. Typical effect sizes in fMRI studies

646 A complete list of effect sizes (BOLD PSC and Cohen's *d*) estimated for different tasks and a priori defined  
 647 ROIs is presented in the Supplementary Materials (Tables S11–19). Here, we focus only on the BOLD  
 648 PSC. The violin plots for some of these are shown in Figure 9.



649

650 **Figure 9.** Typical BOLD PSC in fMRI studies. The box plots inside the violins represent the first and third  
 651 quartile, and the black circles represent median values. Contrasts from the same task are indicated in one  
 652 colour. Abbreviations: L/R – left/right, AMY – amygdala, V1 – primary visual cortex, DLPFC –  
 653 dorsolateral prefrontal cortex, BA – Brodmann area, STG – superior temporal gyrus, A1 – primary auditory  
 654 cortex, NAc – nucleus accumbens, IPL – inferior parietal lobule, IFG/FO – inferior frontal gyrus/frontal  
 655 operculum.

656 For example, the median BOLD PSC in the left amygdala ROI, one of the key brain areas for emotional  
 657 processing, was 0.263%, which is approximately twice as large as *one prior*  $SD_{\theta}$  (see Fig. 7). Thus, using  
 658 this PSC as the ES threshold in future studies may cause the ROPE to become too wide compared to the  
 659 effect sizes typical for tasks with such designs. Therefore, such a threshold can be used to detect large and  
 660 highly localised effects. However, it may fail to detect small but widely distributed effects previously  
 661 described for HCP data (Cremers et al., 2018).

662 In general, median PSCs within ROIs were up to 1% for block designs and 0.5% for event-related designs.  
 663 The maximum PSCs reached 2.5% and were usually observed in the primary visual cortex (V1) for visual  
 664 tasks comparing experimental conditions with baseline activity. For ‘moderate’ physiological effects, PSC  
 665 varied in the range 0.1–0.2%, for example, for the ‘2-back > 0-back’ contrast, the median PSC in the right  
 666 dorsolateral prefrontal cortex (R DLPFC in Fig. 9) was 0.137%. Likewise, for the ‘Social > Random’  
 667 contrast, the right inferior parietal lobule (R IPL) median PSC was 0.137%, for the ‘Correct Stop > Go’,  
 668 the right inferior frontal gyrus/frontal operculum (R IFG/FO) median PSC was 0.120%. For more ‘strong’  
 669 physiological effects, the PSC was in the range 0.2–0.3%, for example, for the ‘Emotion > Shape’ contrast,  
 670 the median PSC in the left amygdala was 0.263%, and for the ‘Story > Math’ contrast, the median PSC in  
 671 the left Brodmann area 45 (Broca’s area) was 0.269%. For the motor activity, for example the ‘Right Finger

672 > baseline' contrast, the median PSC in the left precentral gyrus was 0.239%, in the left postcentral gyrus  
673 was 0.362%, in the left putamen was 0.290%, and in the right cerebellum was 0.401%. For the contrasts  
674 that did not elicit robust activations (Barch et al., 2013), the PSC was approximately 0.05%; for example,  
675 for the 'Reward > Loss' contrast, the median PSC in the left nucleus accumbens was 0.043%, and for the  
676 'Relational > Match' contrast, the median PSC in the left dorsolateral prefrontal cortex was 0.062%.

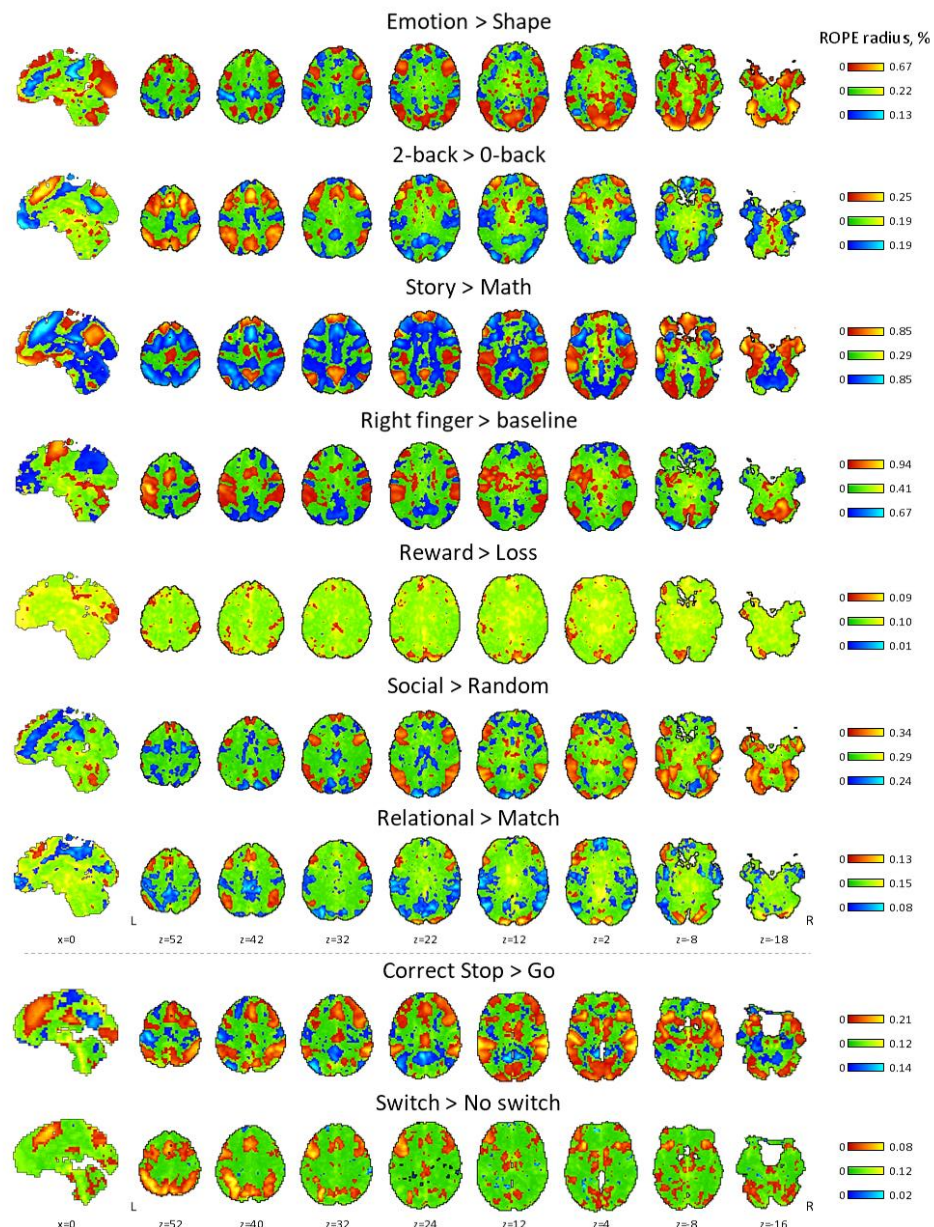
#### 677 **4.5. ROPE maps**

678 We considered BPI with two consecutive thresholding steps: 1) calculate the *LPOs* (or PPMs) with a  
679 selected ES threshold  $\gamma$ , 2) apply the posterior probability threshold  $p_{th} = 95\%$  or consider the overlap  
680 between 95% HDI and ROPE. We can *reverse the thresholding sequence* and calculate *the ROPE maps*.

681 For the '(de)activated' voxels, the ROPE map contains the maximum ES thresholds that allow voxels to be  
682 classified as '(de)activated' based on the 'ROPE-only' or 'HDI+ROPE' decision rules. For the 'not  
683 activated' voxels, the map contains the minimum effect size thresholds that allow voxels to be classified as  
684 'not activated'.

685 The procedure for calculating the ROPE map can be performed as follows. Let us consider a gradual  
686 increase in the ROPE radius (i.e., the half-width of ROPE or the ES threshold  $\gamma$ ) from zero to the maximum  
687 effect size in observed volume. (1) For voxels in which PSC is close to zero, at a certain ROPE radius, the  
688 posterior probability of finding the effect within the ROPE becomes higher than 95%. This width is  
689 indicated on the ROPE map for 'not activated' voxels. (2) For voxels in which the PSC deviates from zero,  
690 at a certain ROPE radius, the posterior probability of finding the effect outside the ROPE becomes lower  
691 than 95%. This width is indicated on the ROPE map for '(de)activated' voxels. The same maps can be  
692 calculated for the 'HDI+ROPE' decision rule.

693 Examples of the ROPE maps are shown in Figure 10. From our point of view, ROPE maps, as well as  
694 unstandardised effect size (PSC) maps, may facilitate an intuitive understanding of the spatial distribution  
695 of a physiological effect under investigation (Chen et al., 2017). They can also be a useful addition to  
696 standard PPMs.



697

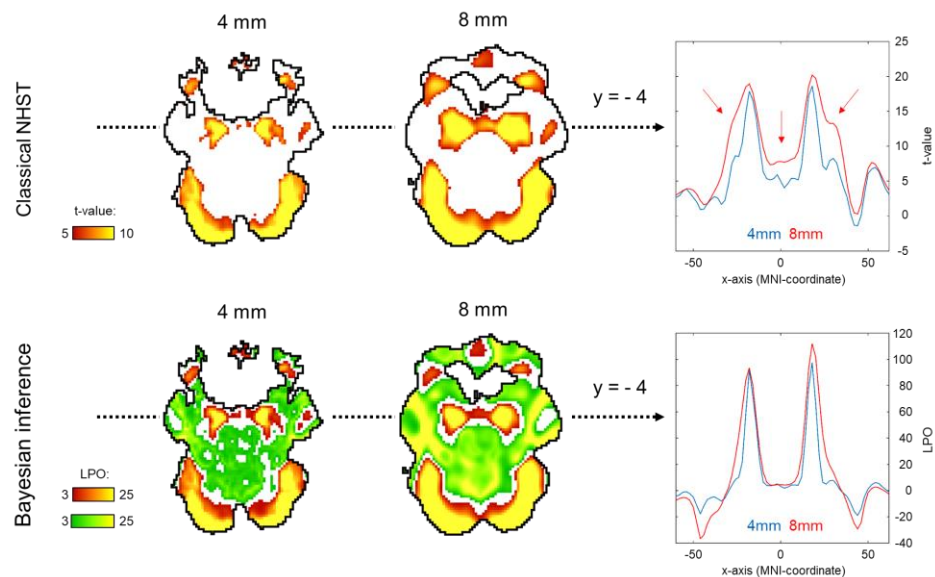
698 **Figure 10.** ROPE maps for all tasks considered in the present study. The ROPE maps are presented using  
 699 different colours for the ‘activated’, ‘deactivated’, and ‘not activated’ voxels. The green bars represent the  
 700 minimum ROPE radii at which voxels with a PSC close to zero can be classified as ‘not activated’ based  
 701 on the ‘ROPE-only’ decision rule. The red and blue bars represent the maximum ROPE radii at which  
 702 voxels of which the PSC deviates from zero can be classified as ‘activated’ and ‘deactivated’, respectively.

#### 703 **4.6. Effects of spatial smoothing on classical NHST and BPI**

704 Two main differences were identified with respect to the influence of the spatial smoothing between  
 705 classical NHST and BPI. Firstly, some voxels with a negligible BOLD PSC could be classified as ‘not  
 706 significant’ at lower smoothing and as ‘activated’ at higher smoothing using classical NHST. At the same  
 707 time, BPI classified these voxels as ‘low confidence’ at lower smoothing and as ‘not activated’ at higher  
 708 smoothing. Higher spatial smoothing increased the number of both ‘(de)activated’ and ‘not activated’  
 709 voxels classified by BPI, and decrease the number of ‘low confidence’ voxels.

710 Secondly, higher smoothing blurred the spatial localisation of local maxima of t-maps and PPMs (*LPO*-  
 711 maps) to a different extent. Consider, for example, the emotion processing task (‘Emotion > Shape’  
 712 contrast). The broadening of two peaks in the left and right amygdala was more noticeable on the t-map  
 713 than on the PPM (see Fig. 11).





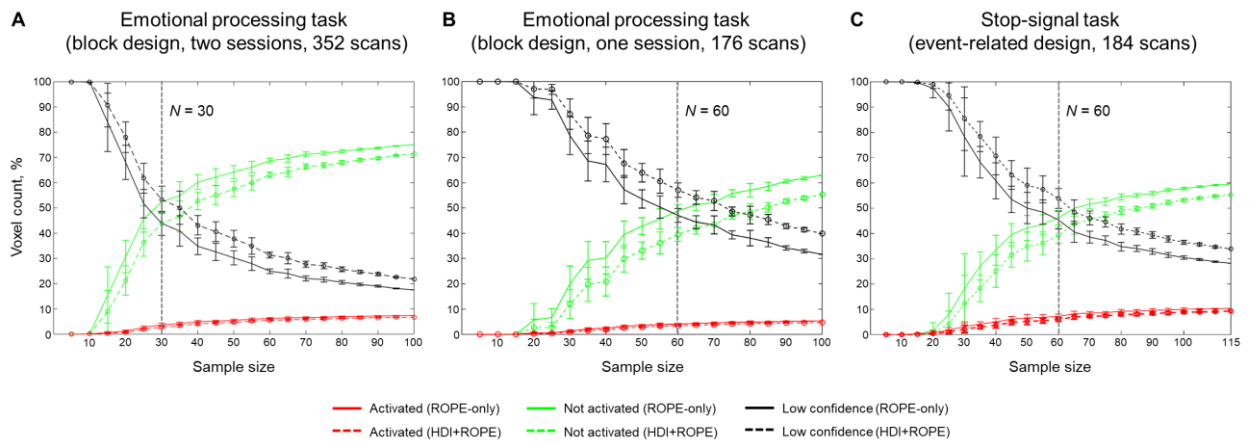
714

715 **Figure 11.** Influence of spatial smoothing on classical NHST and BPI: results for the emotion processing  
 716 task ('Emotion > Shape' contrast). Classical NHST was implemented using voxel-wise FWE correction ( $\alpha$   
 717 = 0.05). BPI was implemented using the 'ROPE-only' decision rule,  $P_{thr} = 95\%$  ( $LPO > 3$ ) and  $\gamma = 1$  prior  
 718  $SD_{\theta} = 0.135\%$ . Axial slice  $z = -14$  mm (MNI152 standard space). In the panels on the right, 1-D images  
 719 are presented for t-values and  $LPOs$  along the x-axis for  $y = -4$  mm. The red arrows indicate a noticeable  
 720 broadening of two peaks of local maxima (left and right amygdala) at higher smoothing.

721 Smoothing was previously shown to have a nonlinear effect on the voxel variances and thus to affect more  
 722 t-maps than  $\beta$  value maps, sometimes leading to counterintuitive artefacts (Reimold et al., 2005). This is  
 723 especially noticeable at the border between two different tissues or between the two narrow peaks of the  
 724 local maxima. If the peak is localised close to white matter voxels with low variability, then smoothing can  
 725 shift the peak to the white matter. If low-variance white matter voxels separate two close peaks, then after  
 726 smoothing, they may serve as a 'bridge' between the two peaks. To avoid this problem, Reimold et al.  
 727 (2005) recommended using masked  $\beta$  value maps. In the present study, we suggest that PPMs based on  
 728 BOLD PSC thresholding can mitigate this problem. Importantly, smoothing artefacts can also arise on  
 729 Cohen's d maps. Therefore, PPMs based on PSC thresholding may be preferable to PPMs based on Cohen's  
 730 d thresholding.

#### 731 4.7. Sample size dependencies for classical NHST and BPI

732 An enlargement of the sample size led to an increase in the number of 'activated' and 'not activated' voxels,  
 733 and a decrease in the number of 'low confidence' voxels. This is due to a decrease in the posterior variance.  
 734 The curve of the 'activated' voxels rose much slower than that of the 'not activated' voxels. For the emotion  
 735 processing task ('Emotion > Shape' contrast, block-design, two sessions, 352 scans), the largest gain in the  
 736 number of 'activated' and 'not activated' voxels can be noted from 20 to 30 subjects (see Fig. 12A). With  
 737 a sample size of  $N > 30$ , the number of 'activated' and 'not activated' voxels increased less steeply. The  
 738 'not activated' and 'low confidence' voxels curves intersected at  $N = 30$  subjects. After the intersection  
 739 point, the graphs reached a plateau.

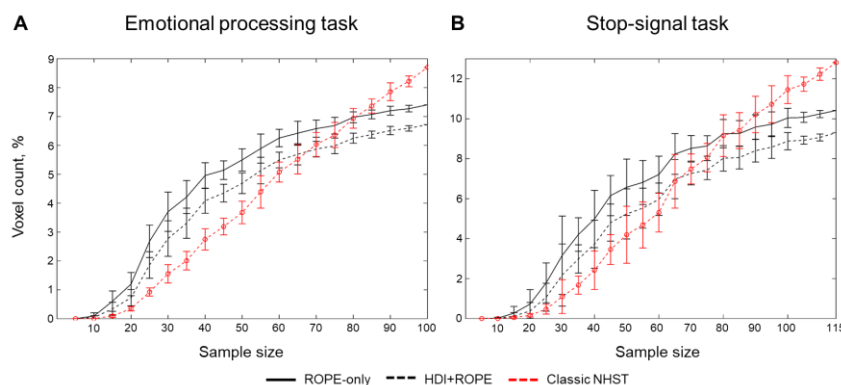


740

741 **Figure 12.** Dependencies of the number of ‘activated’, ‘not activated’, and ‘low confidence’ voxels on the  
 742 sample size. BPI was implemented using  $\gamma = 1$  prior  $SD_{\theta}$ . A) The emotional processing task (‘Emotion >  
 743 Shape’ contrast, two sessions). B) The emotional processing task (‘Emotion > Shape’ contrast, one session).  
 744 C) The stop-signal task (‘Correct Stop > Go’ contrast). The error bars represent the mean and standard  
 745 deviation across ten random groups.

746 Considering only half of the emotional processing task data (one session, 176 scans), the intersection point  
 747 shifted from  $N = 30$  to  $N = 60$  (see Fig. 12B). For the event-related task (‘Correct Stop > Go’ contrast, the  
 748 stop-signal task, 184 scans), all considered dependencies had the same features as for the block-design task,  
 749 and the point of intersection was at  $N = 60$  subjects (see Fig. 12C). Therefore, the moment at which the  
 750 graphs reach a plateau depends mainly on the amount of data at the subject level, that is, on the number of  
 751 scans, blocks, and events. The task designs from the HCP and UCLA datasets have relatively short durations  
 752 (for example, the stop-signal task has approximately 15 ‘Correct Stop’ trials per subject). Studies with a  
 753 longer scanning time generally require a smaller sample size to enable inferences to be made with  
 754 confidence.

755 Classical NHST with the voxel-wise FWE correction showed a steady linear increase in the number of  
 756 ‘activated’ voxels with increasing sample size (see Fig. 13). With a further increase in the sample size, the  
 757 number of statistically significant voxels revealed by classical NHST is expected to approach 100% (see,  
 758 for example, Gonzalez-Castillo et al., 2012). In contrast, the BPI with the  $\gamma = 1$  prior  $SD_{\theta}$  threshold  
 759 demonstrated exponential dependencies. We observed a steeper increase at small and moderate sample  
 760 sizes ( $N = 15-60$ ). The curve of the ‘activated’ voxels flattened at large sample sizes ( $N > 80$ ). BPI offers  
 761 protection against the detection of ‘trivial’ effects that can appear as a result of an increased sample size if  
 762 classical NHST with the point-null hypothesis is used (Friston et al., 2002a; Friston, 2012; Chen et al.,  
 763 2017). This is achieved by the ES threshold  $\gamma$ , which eliminates physiologically (practically) negligible  
 764 effects. Figure 13 presents an illustration of the Jeffreys-Lindley paradox, that is, the discrepancy between  
 765 results obtained using classical and Bayesian inference, which is usually manifested at higher sample sizes  
 766 (Jeffreys, 1939/1948; Lindley, 1957; Friston, 2012).



767

768 **Figure 13.** Dependencies of the number of ‘activated’ voxels on the sample size. Classical NHST was  
 769 implemented using FWE correction ( $\alpha = 0.05$ ). BPI was implemented using  $\gamma = 1$  prior  $SD_{\theta}$ . A) The

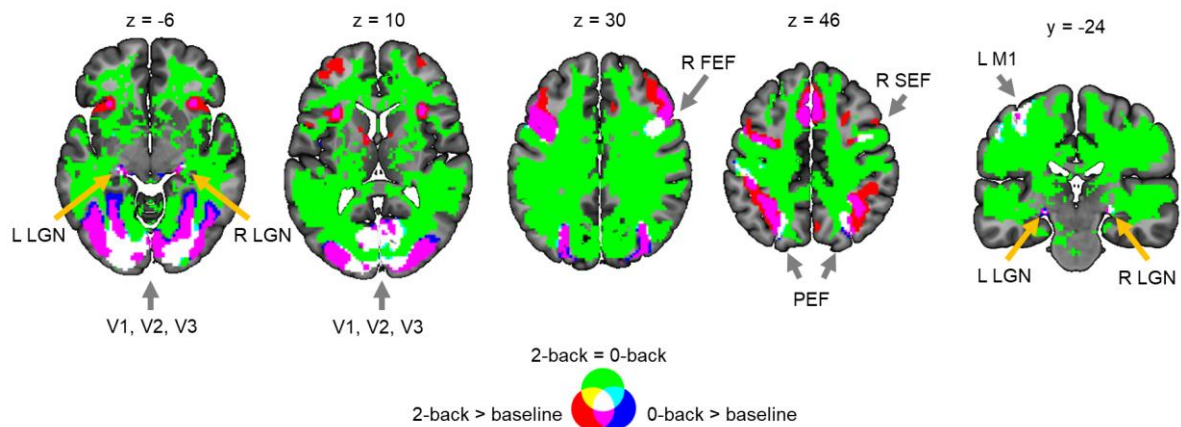
770 emotional processing task (block design, ‘Emotion > Shape’ contrast). B) The stop-signal task (event-  
771 related design, ‘Correct Stop > Go’ contrast). The error bars represent the mean and standard deviation  
772 across ten random groups.

#### 773 4.8. Example of practical application of BPI

774 In contrast to classical NHST, Bayesian inference allows us to:

- 775 (1) Provide evidence that there is no practically meaningful BOLD signal change in the brain area  
776 when comparing the two task conditions.
- 777 (2) Establish double dissociations (Friston et al., 2002a).
- 778 (3) Provide evidence for practically equivalent engagement of one area under different experimental  
779 conditions in terms of local brain activity.
- 780 (4) Provide evidence for the absence of a practically meaningful difference in BOLD signals between  
781 groups of subjects or repeated measures.

782 To illustrate a possible application of Bayesian inference in research practice, we used a working memory  
783 task. Let us consider an overlap between the ‘2-back > baseline’ and ‘0-back > baseline’ contrasts (see Fig  
784 14, purple areas). We cannot claim that brain areas revealed by this conjunction analysis were equally  
785 engaged in the ‘2-back’ and ‘0-back’ conditions. To provide evidence for this notion, we can use BPI and  
786 attempt to identify voxels with a practically equivalent BOLD signal in the ‘2-back’ and ‘0-back’ conditions  
787 (see Fig 14, green). Overlap between the ‘2-back > baseline’ and ‘0-back > baseline’ and the ‘2-back = 0-  
788 back’ effects was found in several brain areas: visual cortex (V1, V2, V3), frontal eye field (FEF), superior  
789 eye field (SEF), parietal eye field (PEF, or posterior parietal cortex), lateral geniculate nucleus (LGN) and  
790 left primary motor cortex (M1) (see Fig. 14, white). This result can be easily explained by the fact that both  
791 experimental conditions require the subject to analyse perceptually similar visual stimuli and push response  
792 buttons with the right hand, which should not depend much on the working memory load. At the same time,  
793 it does not follow directly from simple conjunction analysis.



794

795 **Figure 14.** Example of possible application of BPI based on the working memory task. Abbreviations: L/R  
796 – left/right, V1, V2, V3 – primary, secondary, and third visual cortex, FEF – frontal eye field, SEF –  
797 superior eye field, PEF – parietal eye field, LGN – lateral geniculate nucleus, M1 – primary motor cortex  
798 (M1).

#### 799 5. Discussion

800 BPI allows us to simultaneously find ‘(de)activated’, ‘not activated’, and ‘low confidence’ voxels using a  
801 single decision rule. The ‘not activated’ decision means that the effect is negligible to an extent that it can  
802 be considered equivalent to the null for practical purposes. The ‘low confidence’ decision means we need  
803 more data to make a confident inference, that is, we need to increase the scanning time, sample size, or  
804 revisit the task design. The use of hierarchical PEB with the ‘global shrinkage’ prior enables us to check

805 the results as the sample size increases and allows us to decide whether to optionally terminate the  
806 experiment if the obtained data are sufficient to make a confident inference. All the above features are  
807 absent from the classical NHST framework.

808 An important advantage of Bayesian inference is that we can use graphs such as those shown in Figure 12  
809 to determine when the obtained data are sufficient to make a confident inference. We can plot such graphs  
810 for the whole brain or for a priori defined ROIs. When the curves reach a plateau, the data collection can  
811 be stopped. If the brain area can be labelled as either ‘(de)activated’ or ‘not activated’ at a relatively small  
812 sample size, it will be still so at larger sample sizes. If the brain area can be labelled as ‘low confidence’,  
813 we must increase the sample size to make a confident inference. At a certain sample size, it could possibly  
814 be labelled as either ‘(de)activated’ or ‘not activated’. In the worst case, we can reach the plateau and still  
815 label the brain area as ‘low confidence’. However, even in this case, we can make a definite conclusion:  
816 the task design is not sensitive to the effect and should be revised. Bayesian inference allows us to monitor  
817 the evidence for the alternative or null hypotheses after each participant without special adjustment for  
818 multiplicity (Edwards et al., 1963; Berger and Berry, 1988; Wagenmakers, 2007, Schönbrodt et al., 2015;  
819 Kruschke and Liddell, 2017b). The optional stopping of the experiment not only allows more freedom in  
820 terms of the experimental design, but also saves limited resources and is even more ethically justified in  
821 certain cases<sup>3</sup> (Edwards et al., 1963; Wagenmakers, 2007).

822 In contrast, frequentist inference depends on the researcher’s intention to stop data collection and thus  
823 requires a definition of the stopping rule based on a priori power analysis. The sequential analysis and  
824 optional stopping in frequentist inference inflate the number of false positives. Moreover, even if the a  
825 priori defined sample size is reached, the researcher can still obtain a non-significant result. In this case,  
826 the researcher can follow two controversial paths within the classical NHST framework. Firstly, the sample  
827 size could be further increased to force an indecisive result to a decisive conclusion. The problem is that  
828 this conclusion would always be against the null hypothesis. Thus, an unbounded increase in the sample  
829 size introduces a discrepancy between classical NHST and Bayesian inference, also known as the Jeffreys-  
830 Lindley paradox. Secondly, one may argue that high a priori power and non-significant results provide  
831 evidence for the null hypothesis (see, for example, Cohen, 1990). However, even high *a priori* power and  
832 non-significant results do not provide direct evidence for the null hypothesis. In fact, a high-powered non-  
833 significant result may arise when the obtained data provide no evidence for the null over the alternative  
834 hypothesis, according to Bayesian inference (Denies and Mclatchie, 2017). This does not mean that power  
835 analysis is irrelevant from a Bayesian perspective. Although power analysis is not necessary for Bayesian  
836 inference, it can still be used within the Bayesian framework for study planning (Kruschke and Liddell,  
837 2017b). At the same time, power analysis is a critical part of frequentist inference, as it depends on  
838 researcher intentions, such as the stopping intention.

839 The main difficulty with the application of BPI is the need to define the ES threshold. However, the problem  
840 of choosing a practically meaningful effect size is not unique to fMRI studies, as it arises in every mature  
841 field of science. It should not discourage us from using BPI, as the point-null hypothesis is never true in the  
842 soft sciences. From our perspective, there are several ways to address this problem. Firstly, the ES threshold  
843 can be chosen based on previously reported effect sizes in studies with a similar design or perform a pilot  
844 study to estimate the expected effect size.

845 Based on the fMRI literature, the largest BOLD responses are evoked by sensory stimulation and vary  
846 within 1–5% of the overall mean whole-brain activity. In contrast, BOLD responses induced by cognitive  
847 tasks vary within 0.1–0.5% (Friston et al., 2002b; Poldrack et al., 2011; Chen et al., 2017). The results  
848 obtained in this study support this notion. Primary sensory effects were >1%, and motor effects were >0.3%.  
849 Cognitive effects can be classified into three categories.

- 850 (1) ‘Strong’ effects of 0.2–0.3% (for example, emotion processing in the amygdala, language  
851 processing in Broca’s area),
- 852 (2) ‘Moderate’ effects of 0.1–0.2% (for example, working memory load in DLPFC, social cognition  
853 in IPL, response inhibition in IFG/FO),

---

<sup>3</sup> This is especially true for PET studies. The BPI method described in this work can also be applied to PET data to reduce the sample size and thus exposure to radioactivity (Svensson et al., 2020).



854 (3) ‘Weak’ effects of 0.05% in contrasts without robust activations (for example, reward processing in  
855 the nucleus accumbens, relational processing in DLPFC).

856 However, choosing the ES threshold based on relatively older studies can be challenging because fMRI  
857 designs become increasingly complex over time, and it can be difficult to find previous experiments  
858 reporting unbiased effect size with a similar design. In this case, one can use the ES threshold equal to *one*  
859 *prior SD* of the effect (Friston and Penny, 2003), which can be thought as a neuronal ‘background noise  
860 level’ or a level of activity that is generic to the whole brain (Eickhoff et al., 2008). BPI with this ES  
861 threshold generally works well for both ‘(de)activated’ and ‘not activated’ voxel detection. However, it  
862 may not be suitable in cases with very small differences between the dependent samples. In addition,  
863 researchers who rely more on the frequentist inference may use the  $\gamma(Dice_{max})$  threshold to replicate the  
864 results obtained previously with classical NHST and additionally search for ‘not activated’ and ‘low  
865 confidence’ voxels. Finally, the degree to which the posterior probability is contained within the ROPEs of  
866 different widths could be specified or the ROPE maps in which the thresholding sequence is inverted could  
867 be calculated. The ROPE maps can be shared in public repositories, such as Neurovault, along with PPMs,  
868 and subsequently thresholded by any reasonable ES threshold.

869 Classical NHST is limited to the detection of ‘(de)activated’ voxels. However, in a properly designed and  
870 conducted study, the ‘null effect’ can be just as valuable and exciting as rejecting the null hypothesis.  
871 Although the interval-based frequentist methods can be used to assess the ‘null effects’, they are less  
872 intuitive and more complicated in practice. At the same time, the results obtained with BPI have intuitively  
873 simple interpretations, as they provide direct evidence *for* the null or the alternative hypothesis, given the  
874 obtained data. Moreover, BPI has already been implemented in SPM12, which greatly facilitates its use by  
875 fMRI practitioners.

876 The ability to provide evidence for the null hypothesis may be especially beneficial for clinical  
877 neuroimaging. Possible issues that can be resolved using this approach are:

- 878 (1) Let the brain activity in certain ROIs due to a neurodegenerative process decrease by more than  $\gamma$   
879 per year on average without any treatment. To prove that a new treatment *effectively protects*  
880 *against neurodegenerative processes*, we can provide evidence that, within one year of treatment,  
881 brain activity was reduced by less than X%.
- 882 (2) Assume that an effective treatment should change the brain activity in certain ROIs by at least X%.  
883 Then, we can prove that a new treatment is *practically ineffective* if the activity has changed by  
884 less than X%.
- 885 (3) Consider two groups of subjects taking a new treatment and a placebo, respectively. Using BPI, we  
886 can provide evidence that the result of the new treatment is *does not differ from that of the placebo*.
- 887 (4) Consider two groups of subjects taking an old effective treatment and a new treatment. Using BPI,  
888 we can provide evidence that the new treatment is *no worse than the old effective treatment*.
- 889 (5) Consider a new treatment for a disease that *is not related to brain function*. Using BPI, we can  
890 provide evidence that the new treatment *does not have side effects* on brain activity.

## 891 **6. Conclusions**

892 Herein, a discussion of the use of the Bayesian and frequentist approaches to assess the ‘null effects’ was  
893 presented. We demonstrated that Bayesian inference may be more intuitive and convenient in practice than  
894 frequentist inference. Crucially, Bayesian inference can detect ‘(de)activated’, ‘not activated,’ and ‘low  
895 confidence’ voxels using a single decision rule. Moreover, it allows for interim analysis and optional  
896 stopping when the obtained sample size is sufficient to make a confident inference. We considered the  
897 problem of defining a threshold for the effect size and provided a reference set of typical effect sizes in  
898 different fMRI designs. Bayesian inference and assessment of the ‘null effects’ may be especially beneficial  
899 for basic and applied clinical neuroimaging. The developed SPM12-based scripts with a simple GUI is  
900 expected to be useful for the assessment of ‘null effects’ using BPI.

## 901 **7. Limitations and future work**

902 Firstly, we did not consider BMI, which is currently mainly used for the analysis of effective connectivity.  
903 A promising area of future research would be to compare the advantages of BMI and BPI when analysing  
904 local brain activity. Secondly, the ‘global shrinkage’ prior must be compared with other possible priors.  
905 Thirdly, we used Bayesian statistics only at the group level. Future studies could consider the advantages  
906 of using the Bayesian approach at both the subject and group levels.

## 907 **8. Acknowledgments**

908 Data were provided [in part] by the Human Connectome Project, WU-Minn Consortium (Principal  
909 Investigators: David Van Essen and Kamil Ugurbil; 1U54MH091657) funded by the 16 NIH Institutes and  
910 Centers that support the NIH Blueprint for Neuroscience Research; and by the McDonnell Center for  
911 Systems Neuroscience at Washington University. Another part of the data were provided by UCLA dataset  
912 which was obtained from the OpenfMRI database (its accession number is ds000030) and data collection  
913 was funded by the Consortium for Neuropsychiatric Phenomics (NIH Roadmap for Medical Research  
914 grants UL1-DE019580, RL1MH083268, RL1MH083269, RL1DA024853, RL1MH083270,  
915 RL1LM009833, PL1MH083271, and PL1NS062410). We thank Dr. Irina Knyazeva for her helpful  
916 discussions and valuable suggestions for this manuscript. We also thank Andrey Ogai for the valuable help  
917 with script for visualisation of statistical maps. RM, AK and MV were supported by the Russian Science  
918 Foundation grant #19-18-00454. YN, MD and DC were supported by the state assignment of the Ministry  
919 of Education and Science of Russian Federation (theme number AAAA-A19-119101890066-2).

920

## 921 **References**

- 922 Acar, F., Seurinck, R., Eickhoff, S. B., & Moerkerke, B. (2018). Assessing robustness against potential  
923 publication bias in Activation Likelihood Estimation (ALE) meta-analyses for fMRI. *PLOS ONE*, *13*(11),  
924 e0208177. <https://doi.org/10.1371/journal.pone.0208177>
- 925 Aczel, B., Palfi, B., Szollosi, A., Kovacs, M., Szaszi, B., Szecsi, P., Zrubka, M., Gronau, Q. F., van den  
926 Bergh, D., & Wagenmakers, E. J. (2018). Quantifying Support for the Null Hypothesis in Psychology: An  
927 Empirical Investigation. *Advances in Methods and Practices in Psychological Science*, *1*(3), 357–366.  
928 <https://doi.org/10.1177/2515245918773742>
- 929 Aisbett J., Lakens D., Sainani K. (2020). Magnitude based inference in relation to one-sided hypotheses  
930 testing procedures. SportRxiv. <https://doi.org/10.31236/osf.io/pn9s3>
- 931 Alberton, B. A., Nichols, T. E., Gamba, H. R., & Winkler, A. M. (2020). Multiple testing correction over  
932 contrasts for brain imaging. *NeuroImage*, *216*, 116760. <https://doi.org/10.1016/j.neuroimage.2020.116760>
- 933 Altman, D. G., & Bland, J. M. (1995). Statistics notes: Absence of evidence is not evidence of absence.  
934 *BMJ*, *311*(7003), 485. <https://doi.org/10.1136/bmj.311.7003.485>
- 935 Amrhein, V., Korner-Nievergelt, F., & Roth, T. (2017). The earth is flat ( $p > 0.05$ ): significance thresholds  
936 and the crisis of unreplicable research. *PeerJ*, *5*, e3544. <https://doi.org/10.7717/peerj.3544>
- 937 Baguley, T. (2009). Standardized or simple effect size: What should be reported? *British Journal of*  
938 *Psychology*, *100*(3), 603–617. <https://doi.org/10.1348/000712608x377117>
- 939 Barch, D. M., Burgess, G. C., Harms, M. P., Petersen, S. E., Schlaggar, B. L., Corbetta, M., Glasser, M. F.,  
940 Curtiss, S., Dixit, S., Feldt, C., Nolan, D., Bryant, E., Hartley, T., Footer, O., Bjork, J. M., Poldrack, R.,  
941 Smith, S., Johansen-Berg, H., Snyder, A. Z., & Van Essen, D. C. (2013). Function in the human  
942 connectome: Task-fMRI and individual differences in behavior. *NeuroImage*, *80*, 169–189.  
943 <https://doi.org/10.1016/j.neuroimage.2013.05.033>
- 944 Belia, S., Fidler, F., Williams, J., & Cumming, G. (2005). Researchers Misunderstand Confidence Intervals  
945 and Standard Error Bars. *Psychological Methods*, *10*(4), 389–396. <https://doi.org/10.1037/1082-989x.10.4.389>
- 947 Berger, J. O. (2003). Could Fisher, Jeffreys and Neyman Have Agreed on Testing? *Statistical Science*,  
948 *18*(1). <https://doi.org/10.1214/ss/1056397485>
- 949 Berger, J. O., & Berry, D. A. (1988). Statistical analysis and the illusion of objectivity. *American Scientist*,  
950 *76*(2), 159–165.
- 951 Berger, J. O., & Sellke, T. (1987). Testing a Point Null Hypothesis: The Irreconcilability of P Values and  
952 Evidence: Rejoinder. *Journal of the American Statistical Association*, *82*(397), 135.  
953 <https://doi.org/10.2307/2289139>
- 954 Berry, D. (1988). Multiple comparisons, multiple tests, and data dredging: A Bayesian perspective. In  
955 Bernardo, J., DeGroot, M., Lindley, D., Smith, A. (Ed.), *Bayesian Statistics* (pp. 79–94). Oxford University  
956 Press.
- 957 Berry, D. A., & Hochberg, Y. (1999). Bayesian perspectives on multiple comparisons. *Journal of Statistical*  
958 *Planning and Inference*, *82*(1–2), 215–227. [https://doi.org/10.1016/s0378-3758\(99\)00044-0](https://doi.org/10.1016/s0378-3758(99)00044-0)
- 959 Campbell, H., & Gustafson, P. (2018). Conditional equivalence testing: An alternative remedy for  
960 publication bias. *PLOS ONE*, *13*(4), e0195145. <https://doi.org/10.1371/journal.pone.0195145>

- 961 Chen, G., Cox, R. W., Glen, D. R., Rajendra, J. K., Reynolds, R. C., & Taylor, P. A. (2018). A tail of two  
962 sides: Artificially doubled false positive rates in neuroimaging due to the sidedness choice with t -tests.  
963 *Human Brain Mapping*, 40(3), 1037–1043. <https://doi.org/10.1002/hbm.24399>
- 964 Chen, G., Taylor, P. A., & Cox, R. W. (2017). Is the statistic value all we should care about in  
965 neuroimaging? *NeuroImage*, 147, 952–959. <https://doi.org/10.1016/j.neuroimage.2016.09.066>
- 966 Chen, G., Taylor, P. A., Cox, R. W., & Pessoa, L. (2020). Fighting or embracing multiplicity in  
967 neuroimaging? neighborhood leverage versus global calibration. *NeuroImage*, 206, 116320.  
968 <https://doi.org/10.1016/j.neuroimage.2019.116320>
- 969 Chen, G., Xiao, Y., Taylor, P. A., Rajendra, J. K., Riggins, T., Geng, F., Redcay, E., & Cox, R. W. (2019).  
970 Handling Multiplicity in Neuroimaging Through Bayesian Lenses with Multilevel Modeling.  
971 *Neuroinformatics*, 17(4), 515–545. <https://doi.org/10.1007/s12021-018-9409-6>
- 972 Cohen, J. (1965). Some statistical issues in psychological research. In B. B. Wolman (Ed.), *Handbook of*  
973 *clinical psychology* (pp. 95–121). New York: McGraw-Hill.
- 974 Cohen, J. (1990). Things I have learned (so far). *American Psychologist*, 45(12), 1304–1312.  
975 <https://doi.org/10.1037/0003-066x.45.12.1304>
- 976 Cohen, J. (1994). The earth is round ( $p < .05$ ). *American Psychologist*, 49(12), 997–1003.  
977 <https://doi.org/10.1037/0003-066x.49.12.997>
- 978 Cornfield, J. (1966). Sequential Trials, Sequential Analysis and the Likelihood Principle. *The American*  
979 *Statistician*, 20(2), 18–23. <https://doi.org/10.1080/00031305.1966.10479786>
- 980 Cortina, J. M., & Dunlap, W. P. (1997). On the logic and purpose of significance testing. *Psychological*  
981 *Methods*, 2(2), 161–172. <https://doi.org/10.1037/1082-989x.2.2.161>
- 982 Cramer, A. O. J., van Ravenzwaaij, D., Matzke, D., Steingroever, H., Wetzels, R., Grasman, R. P. P. P.,  
983 Waldorp, L. J., & Wagenmakers, E. J. (2015). Hidden multiplicity in exploratory multiway ANOVA:  
984 Prevalence and remedies. *Psychonomic Bulletin & Review*, 23(2), 640–647.  
985 <https://doi.org/10.3758/s13423-015-0913-5>
- 986 Cremers, H. R., Wager, T. D., & Yarkoni, T. (2017). The relation between statistical power and inference  
987 in fMRI. *PLOS ONE*, 12(11), e0184923. <https://doi.org/10.1371/journal.pone.0184923>
- 988 Cumming, G. (2013). The New Statistics: Why and How. *Psychological Science*, 25(1), 7–29.  
989 <https://doi.org/10.1177/0956797613504966>
- 990 Dandolo, L. C., & Schwabe, L. (2019). Time-dependent motor memory representations in prefrontal cortex.  
991 *NeuroImage*, 197, 143–155. <https://doi.org/10.1016/j.neuroimage.2019.04.051>
- 992 David, S. P., Naudet, F., Laude, J., Radua, J., Fusar-Poli, P., Chu, I., Stefanick, M. L., & Ioannidis, J. P. A.  
993 (2018). Potential Reporting Bias in Neuroimaging Studies of Sex Differences. *Scientific Reports*, 8(1).  
994 <https://doi.org/10.1038/s41598-018-23976-1>
- 995 de Winter, J. C., & Dodou, D. (2015). A surge of p-values between 0.041 and 0.049 in recent decades (but  
996 negative results are increasing rapidly too). *PeerJ*, 3, e733. <https://doi.org/10.7717/peerj.733>
- 997 Dienes, Z. (2014). Using Bayes to get the most out of non-significant results. *Frontiers in Psychology*, 5.  
998 <https://doi.org/10.3389/fpsyg.2014.00781>
- 999 Dienes, Z., & Mclatchie, N. (2017). Four reasons to prefer Bayesian analyses over significance testing.  
1000 *Psychonomic Bulletin & Review*, 25(1), 207–218. <https://doi.org/10.3758/s13423-017-1266-z>



- 1001 Edwards, W., Lindman, H., & Savage, L. J. (1963). Bayesian statistical inference for psychological  
1002 research. *Psychological Review*, 70(3), 193–242. <https://doi.org/10.1037/h0044139>
- 1003 Eickhoff, S. B., Grefkes, C., Fink, G. R., & Zilles, K. (2008). Functional Lateralization of Face, Hand, and  
1004 Trunk Representation in Anatomically Defined Human Somatosensory Areas. *Cerebral Cortex*, 18(12),  
1005 2820–2830. <https://doi.org/10.1093/cercor/bhn039>
- 1006 Eklund, A., Nichols, T. E., & Knutsson, H. (2016). Cluster failure: Why fMRI inferences for spatial extent  
1007 have inflated false-positive rates. *Proceedings of the National Academy of Sciences*, 113(28), 7900–7905.  
1008 <https://doi.org/10.1073/pnas.1602413113>
- 1009 Falk, R., & Greenbaum, C. W. (1995). Significance Tests Die Hard. *Theory & Psychology*, 5(1), 75–98.  
1010 <https://doi.org/10.1177/0959354395051004>
- 1011 Feng, C., Forthman, K. L., Kuplicki, R., Yeh, H. W., Stewart, J. L., & Paulus, M. P. (2019). Neighborhood  
1012 affluence is not associated with positive and negative valence processing in adults with mood and anxiety  
1013 disorders: A Bayesian inference approach. *NeuroImage: Clinical*, 22, 101738.  
1014 <https://doi.org/10.1016/j.nicl.2019.101738>
- 1015 Fidler, F., Burgman, M. A., Cumming, G., Buttrose, R., & Thomason, N. (2006). Impact of Criticism of  
1016 Null-Hypothesis Significance Testing on Statistical Reporting Practices in Conservation Biology.  
1017 *Conservation Biology*, 20(5), 1539–1544. <https://doi.org/10.1111/j.1523-1739.2006.00525.x>
- 1018 Finch, S., Cumming, G., & Thomason, N. (2001). Colloquium on Effect Sizes: the Roles of Editors,  
1019 Textbook Authors, and the Publication Manual. *Educational and Psychological Measurement*, 61(2), 181–  
1020 210. <https://doi.org/10.1177/0013164401612001>
- 1021 Friston, K. (2012). Ten ironic rules for non-statistical reviewers. *NeuroImage*, 61(4), 1300–1310.  
1022 <https://doi.org/10.1016/j.neuroimage.2012.04.018>
- 1023 Friston, K. (2013). Sample size and the fallacies of classical inference. *NeuroImage*, 81, 503–504.  
1024 <https://doi.org/10.1016/j.neuroimage.2013.02.057>
- 1025 Friston, K., Glaser, D., Henson, R., Kiebel, S., Phillips, C., & Ashburner, J. (2002b). Classical and Bayesian  
1026 Inference in Neuroimaging: Applications. *NeuroImage*, 16(2), 484–512.  
1027 <https://doi.org/10.1006/nimg.2002.1091>
- 1028 Friston, K. J., Holmes, A. P., Worsley, K. J., Poline, J. P., Frith, C. D., & Frackowiak, R. S. J. (1994).  
1029 Statistical parametric maps in functional imaging: A general linear approach. *Human Brain Mapping*, 2(4),  
1030 189–210. <https://doi.org/10.1002/hbm.460020402>
- 1031 Friston, K. J., Williams, S., Howard, R., Frackowiak, R. S. J., & Turner, R. (1996). Movement-Related  
1032 effects in fMRI time-series. *Magnetic Resonance in Medicine*, 35(3), 346–355.  
1033 <https://doi.org/10.1002/mrm.1910350312>
- 1034 Friston, K., & Penny, W. (2003). Posterior probability maps and SPMs. *NeuroImage*, 19(3), 1240–1249.  
1035 [https://doi.org/10.1016/s1053-8119\(03\)00144-7](https://doi.org/10.1016/s1053-8119(03)00144-7)
- 1036 Friston, K., & Penny, W. (2011). Post hoc Bayesian model selection. *NeuroImage*, 56(4), 2089–2099.  
1037 <https://doi.org/10.1016/j.neuroimage.2011.03.062>
- 1038 Friston, K., Penny, W., Phillips, C., Kiebel, S., Hinton, G., & Ashburner, J. (2002a). Classical and Bayesian  
1039 Inference in Neuroimaging: Theory. *NeuroImage*, 16(2), 465–483. <https://doi.org/10.1006/nimg.2002.1090>
- 1040 Gelman, A., Hill, J., & Yajima, M. (2012). Why We (Usually) Don't Have to Worry About Multiple  
1041 Comparisons. *Journal of Research on Educational Effectiveness*, 5(2), 189–211.  
1042 <https://doi.org/10.1080/19345747.2011.618213>

- 1043 Genovese, C. R., Lazar, N. A., & Nichols, T. (2002). Thresholding of Statistical Maps in Functional  
1044 Neuroimaging Using the False Discovery Rate. *NeuroImage*, 15(4), 870–878.  
1045 <https://doi.org/10.1006/nimg.2001.1037>
- 1046 Gigerenzer, G. (1993). The superego, the ego, and the id in statistical reasoning. In G. Keren & C. Lewis  
1047 (Eds.), *A handbook for data analysis in the behavioral sciences: Methodological issues* (pp. 311–339).  
1048 Lawrence Erlbaum Associates, Inc.
- 1049 Glasser, M. F., Sotiropoulos, S. N., Wilson, J. A., Coalson, T. S., Fischl, B., Andersson, J. L., Xu, J., Jbabdi,  
1050 S., Webster, M., Polimeni, J. R., Van Essen, D. C., & Jenkinson, M. (2013). The minimal preprocessing  
1051 pipelines for the Human Connectome Project. *NeuroImage*, 80, 105–124.  
1052 <https://doi.org/10.1016/j.neuroimage.2013.04.127>
- 1053 Gonzalez-Castillo, J., Saad, Z. S., Handwerker, D. A., Inati, S. J., Brenowitz, N., & Bandettini, P. A. (2012).  
1054 Whole-brain, time-locked activation with simple tasks revealed using massive averaging and model-free  
1055 analysis. *Proceedings of the National Academy of Sciences*, 109(14), 5487–5492.  
1056 <https://doi.org/10.1073/pnas.1121049109>
- 1057 Goodman, S. (2008). A Dirty Dozen: Twelve P-Value Misconceptions. *Seminars in Hematology*, 45(3),  
1058 135–140. <https://doi.org/10.1053/j.seminhematol.2008.04.003>
- 1059 Goodman, S. N. (1993). p Values, Hypothesis Tests, and Likelihood: Implications for Epidemiology of a  
1060 Neglected Historical Debate. *American Journal of Epidemiology*, 137(5), 485–496.  
1061 <https://doi.org/10.1093/oxfordjournals.aje.a116700>
- 1062 Gopalan, R., & Berry, D. A. (1998). Bayesian Multiple Comparisons Using Dirichlet Process Priors.  
1063 *Journal of the American Statistical Association*, 93(443), 1130–1139.  
1064 <https://doi.org/10.1080/01621459.1998.10473774>
- 1065 Gorgolewski, K. J., Durnez, J., & Poldrack, R. A. (2017). Preprocessed Consortium for Neuropsychiatric  
1066 Phenomics dataset. *F1000Research*, 6, 1262. <https://doi.org/10.12688/f1000research.11964.2>
- 1067 Greenland, S. (2019). Valid P-Values Behave Exactly as They Should: Some Misleading Criticisms of P-  
1068 Values and Their Resolution With S-Values. *The American Statistician*, 73(sup1), 106–114.  
1069 <https://doi.org/10.1080/00031305.2018.1529625>
- 1070 Greenland, S., Senn, S. J., Rothman, K. J., Carlin, J. B., Poole, C., Goodman, S. N., & Altman, D. G. (2016).  
1071 Statistical tests, P values, confidence intervals, and power: a guide to misinterpretations. *European Journal*  
1072 *of Epidemiology*, 31(4), 337–350. <https://doi.org/10.1007/s10654-016-0149-3>
- 1073 Greenwald, A. G. (1975). Consequences of prejudice against the null hypothesis. *Psychological Bulletin*,  
1074 82(1), 1–20. <https://doi.org/10.1037/h0076157>
- 1075 Gusnard, D. A., & Raichle, M. E. (2001). Searching for a baseline: Functional imaging and the resting  
1076 human brain. *Nature Reviews Neuroscience*, 2(10), 685–694. <https://doi.org/10.1038/35094500>
- 1077 Hagen, R. L. (1997). In praise of the null hypothesis statistical test. *American Psychologist*, 52(1), 15–24.  
1078 <https://doi.org/10.1037/0003-066x.52.1.15>
- 1079 Hodges, J. L., & Lehmann, E. L. (1954). Testing the Approximate Validity of Statistical Hypotheses.  
1080 *Journal of the Royal Statistical Society: Series B (Methodological)*, 16(2), 261–268.  
1081 <https://doi.org/10.1111/j.2517-6161.1954.tb00169.x>
- 1082 Hoekstra, R., Finch, S., Kiers, H. A. L., & Johnson, A. (2006). Probability as certainty: Dichotomous  
1083 thinking and the misuse of p values. *Psychonomic Bulletin & Review*, 13(6), 1033–1037.  
1084 <https://doi.org/10.3758/bf03213921>

- 1085 Hoekstra, R., Morey, R. D., Rouder, J. N., & Wagenmakers, E. J. (2014). Robust misinterpretation of  
1086 confidence intervals. *Psychonomic Bulletin & Review*, *21*(5), 1157–1164. [https://doi.org/10.3758/s13423-](https://doi.org/10.3758/s13423-013-0572-3)  
1087 013-0572-3
- 1088 Hubbard, R., & Bayarri, M. J. (2003). Confusion Over Measures of Evidence ( $p$ 's) Versus Errors ( $\alpha$ 's) in  
1089 Classical Statistical Testing. *The American Statistician*, *57*(3), 171–178.  
1090 <https://doi.org/10.1198/0003130031856>
- 1091 Hubbard, R., & Lindsay, R. M. (2008). Why P Values Are Not a Useful Measure of Evidence in Statistical  
1092 Significance Testing. *Theory & Psychology*, *18*(1), 69–88. <https://doi.org/10.1177/0959354307086923>
- 1093 Ioannidis, J. P. A. (2019). What Have We (Not) Learnt from Millions of Scientific Papers with P Values?  
1094 *The American Statistician*, *73*(sup1), 20–25. <https://doi.org/10.1080/00031305.2018.1447512>
- 1095 Ioannidis, J. P., Munafò, M. R., Fusar-Poli, P., Nosek, B. A., & David, S. P. (2014). Publication and other  
1096 reporting biases in cognitive sciences: detection, prevalence, and prevention. *Trends in Cognitive Sciences*,  
1097 *18*(5), 235–241. <https://doi.org/10.1016/j.tics.2014.02.010>
- 1098 Jeffreys, H. (1948). *Theory of Probability*, 2nd ed. The Clarendon Press, Oxford.
- 1099 Jennings, R. G., & Van Horn, J. D. (2011). Publication Bias in Neuroimaging Research: Implications for  
1100 Meta-Analyses. *Neuroinformatics*, *10*(1), 67–80. <https://doi.org/10.1007/s12021-011-9125-y>
- 1101 Johansson, T. (2010). Hail the impossible: p-values, evidence, and likelihood. *Scandinavian Journal of*  
1102 *Psychology*, *52*(2), 113–125. <https://doi.org/10.1111/j.1467-9450.2010.00852.x>
- 1103 Joyce, K. E., & Hayasaka, S. (2012). Development of PowerMap: a Software Package for Statistical Power  
1104 Calculation in Neuroimaging Studies. *Neuroinformatics*, *10*(4), 351–365. [https://doi.org/10.1007/s12021-](https://doi.org/10.1007/s12021-012-9152-3)  
1105 012-9152-3
- 1106 Kass, R. E., & Raftery, A. E. (1995). Bayes Factors. *Journal of the American Statistical Association*,  
1107 *90*(430), 773–795. <https://doi.org/10.1080/01621459.1995.10476572>
- 1108 Kirk, R. E. (1996). Practical Significance: A Concept Whose Time Has Come. *Educational and*  
1109 *Psychological Measurement*, *56*(5), 746–759. <https://doi.org/10.1177/0013164496056005002>
- 1110 Kruschke, J. K. (2010). What to believe: Bayesian methods for data analysis. *Trends in Cognitive Sciences*,  
1111 *14*(7), 293–300. <https://doi.org/10.1016/j.tics.2010.05.001>
- 1112 Kruschke, J. K. (2011). Bayesian Assessment of Null Values Via Parameter Estimation and Model  
1113 Comparison. *Perspectives on Psychological Science*, *6*(3), 299–312.  
1114 <https://doi.org/10.1177/1745691611406925>
- 1115 Kruschke, J. K., & Liddell, T. M. (2017a). Bayesian data analysis for newcomers. *Psychonomic Bulletin &*  
1116 *Review*, *25*(1), 155–177. <https://doi.org/10.3758/s13423-017-1272-1>
- 1117 Kruschke, J. K., & Liddell, T. M. (2017b). The Bayesian New Statistics: Hypothesis testing, estimation,  
1118 meta-analysis, and power analysis from a Bayesian perspective. *Psychonomic Bulletin & Review*, *25*(1),  
1119 178–206. <https://doi.org/10.3758/s13423-016-1221-4>
- 1120 Lakens, D. (2017). Equivalence Tests. *Social Psychological and Personality Science*, *8*(4), 355–362.  
1121 <https://doi.org/10.1177/1948550617697177>
- 1122 Lakens, D., McLatchie, N., Isager, P. M., Scheel, A. M., & Dienes, Z. (2018). Improving Inferences About  
1123 Null Effects With Bayes Factors and Equivalence Tests. *The Journals of Gerontology: Series B*, *75*(1), 45–  
1124 57. <https://doi.org/10.1093/geronb/gby065>

- 1125 Liao J.G. , Midya V., Berg A. (2019). Connecting Bayes factor and the region of practical equivalence  
1126 (ROPE) procedure for testing interval null hypothesis. arXiv:1903.03153
- 1127 Lindley, D. (1965). *Introduction to Probability and Statistics from a Bayesian Viewpoint* (1st ed.).  
1128 Cambridge University Press.
- 1129 Lindley, D. V. (1957). A Statistical Paradox. *Biometrika*, 44(1/2), 187. <https://doi.org/10.2307/2333251>
- 1130 Lindley, D. V. (1975). The Future of Statistics: A Bayesian 21st Century. *Advances in Applied Probability*,  
1131 7, 106. <https://doi.org/10.2307/1426315>
- 1132 Lindley, D. V. (1990). The 1988 Wald Memorial Lectures: The Present Position in Bayesian Statistics.  
1133 *Statistical Science*, 5(1). <https://doi.org/10.1214/ss/1177012253>
- 1134 Magerkurth, J., Mancini, L., Penny, W., Flandin, G., Ashburner, J., Micallef, C., De Vita, E., Daga, P.,  
1135 White, M. J., Buckley, C., Yamamoto, A. K., Ourselin, S., Yousry, T., Thornton, J. S., & Weiskopf, N.  
1136 (2015). Objective Bayesian fMRI analysis - a pilot study in different clinical environments. *Frontiers in*  
1137 *Neuroscience*, 9. <https://doi.org/10.3389/fnins.2015.00168>
- 1138 Meehl, P. E. (1967). Theory-Testing in Psychology and Physics: A Methodological Paradox. *Philosophy*  
1139 *of Science*, 34(2), 103–115. <https://doi.org/10.1086/288135>
- 1140 Meehl, P. E. (2004). Theoretical risks and tabular asterisks: Sir Karl, Sir Ronald, and the slow progress of  
1141 soft psychology. *Applied and Preventive Psychology*, 11(1), 1. <https://doi.org/10.1016/j.appsy.2004.02.001>
- 1142 Meyners, M. (2012). Equivalence tests – A review. *Food Quality and Preference*, 26(2), 231–245.  
1143 <https://doi.org/10.1016/j.foodqual.2012.05.003>
- 1144 Morey, R. D., Hoekstra, R., Rouder, J. N., & Wagenmakers, E. J. (2015). Continued misinterpretation of  
1145 confidence intervals: response to Miller and Ulrich. *Psychonomic Bulletin & Review*, 23(1), 131–140.  
1146 <https://doi.org/10.3758/s13423-015-0955-8>
- 1147 Morey, R. D., & Rouder, J. N. (2011). Bayes factor approaches for testing interval null hypotheses.  
1148 *Psychological Methods*, 16(4), 406–419. <https://doi.org/10.1037/a0024377>
- 1149 Muller, P., Parmigiani, G., Rice, K., (2006). FDR and Bayesian multiple comparisons rules. In: Bernardo,  
1150 J.M., Bayarri, M.J., Berger, J.O., Dawid, A.P., Heckerman, D., Smith, A.F.M., West, M. (Ed.), *Bayesian*  
1151 *Statistics 8: Proceedings of the Eighth Valencia International Meeting* (pp. 366–368). Oxford University  
1152 Press.
- 1153 Mumford, J. A. (2012). A power calculation guide for fMRI studies. *Social Cognitive and Affective*  
1154 *Neuroscience*, 7(6), 738–742. <https://doi.org/10.1093/scan/nss059>
- 1155 Mumford, J. A., & Nichols, T. E. (2008). Power calculation for group fMRI studies accounting for arbitrary  
1156 design and temporal autocorrelation. *NeuroImage*, 39(1), 261–268.  
1157 <https://doi.org/10.1016/j.neuroimage.2007.07.061>
- 1158 Murphy, K. R., & Myers, B. (1999). Testing the hypothesis that treatments have negligible effects:  
1159 Minimum-effect tests in the general linear model. *Journal of Applied Psychology*, 84(2), 234–248.  
1160 <https://doi.org/10.1037/0021-9010.84.2.234>
- 1161 Murphy, K. R., & Myers, B. (2004). *Statistical Power Analysis: A Simple and General Model for*  
1162 *Traditional and Modern Hypothesis Tests* (2nd ed.). Lawrence Erlbaum Associates.
- 1163 Nichols, T. E. (2012). Multiple testing corrections, nonparametric methods, and random field theory.  
1164 *NeuroImage*, 62(2), 811–815. <https://doi.org/10.1016/j.neuroimage.2012.04.014>



- 1165 Nichols, T., & Hayasaka, S. (2003). Controlling the familywise error rate in functional neuroimaging: a  
1166 comparative review. *Statistical Methods in Medical Research*, 12(5), 419–446.  
1167 <https://doi.org/10.1191/0962280203sm341ra>
- 1168 Nickerson, R. S. (2000). Null hypothesis significance testing: A review of an old and continuing  
1169 controversy. *Psychological Methods*, 5(2), 241–301. <https://doi.org/10.1037/1082-989x.5.2.241>
- 1170 Penny, W. D., & Ridgway, G. R. (2013). Efficient Posterior Probability Mapping Using Savage-Dickey  
1171 Ratios. *PLoS ONE*, 8(3), e59655. <https://doi.org/10.1371/journal.pone.0059655>
- 1172 Penny, W. D., Trujillo-Barreto, N. J., & Friston, K. J. (2005). Bayesian fMRI time series analysis with  
1173 spatial priors. *NeuroImage*, 24(2), 350–362. <https://doi.org/10.1016/j.neuroimage.2004.08.034>
- 1174 Penny, W., Flandin, G., & Trujillo-Barreto, N. (2007). Bayesian comparison of spatially regularised general  
1175 linear models. *Human Brain Mapping*, 28(4), 275–293. <https://doi.org/10.1002/hbm.20327>
- 1176 Penny, W., Kiebel, S., & Friston, K. (2003). Variational Bayesian inference for fMRI time series.  
1177 *NeuroImage*, 19(3), 727–741. [https://doi.org/10.1016/s1053-8119\(03\)00071-5](https://doi.org/10.1016/s1053-8119(03)00071-5)
- 1178 Perezgonzalez, J. D. (2015). Fisher, Neyman-Pearson or NHST? A tutorial for teaching data testing.  
1179 *Frontiers in Psychology*, 6. <https://doi.org/10.3389/fpsyg.2015.00223>
- 1180 Pernet, C. R. (2014). Misconceptions in the use of the General Linear Model applied to functional MRI: a  
1181 tutorial for junior neuro-imagers. *Frontiers in Neuroscience*, 8. <https://doi.org/10.3389/fnins.2014.00001>
- 1182 Poldrack, R. A., Baker, C. I., Durnez, J., Gorgolewski, K. J., Matthews, P. M., Munafò, M. R., Nichols, T.  
1183 E., Poline, J. B., Vul, E., & Yarkoni, T. (2017). Scanning the horizon: towards transparent and reproducible  
1184 neuroimaging research. *Nature Reviews Neuroscience*, 18(2), 115–126.  
1185 <https://doi.org/10.1038/nrn.2016.167>
- 1186 Poldrack, R. A., Mumford, J. A., & Nichols, T. E. (2011). *Handbook of Functional MRI Data Analysis*.  
1187 Cambridge University Press.
- 1188 Poldrack, R., Congdon, E., Triplett, W., Gorgolewski, K., Karlsgodt, K., Mumford, J., Sabb, F., Freimer,  
1189 N., London, E., Cannon, T., & Bilder, R. (2016). A phenome-wide examination of neural and cognitive  
1190 function. *Scientific Data*, 3(1). <https://doi.org/10.1038/sdata.2016.110>
- 1191 Poline, J. B., & Brett, M. (2012). The general linear model and fMRI: Does love last forever? *NeuroImage*,  
1192 62(2), 871–880. <https://doi.org/10.1016/j.neuroimage.2012.01.133>
- 1193 Pollard, P., & Richardson, J. T. (1987). On the probability of making Type I errors. *Psychological Bulletin*,  
1194 102(1), 159–163. <https://doi.org/10.1037/0033-2909.102.1.159>
- 1195 Raichle, M. E., & Gusnard, D. A. (2002). Appraising the brain's energy budget. *Proceedings of the National  
1196 Academy of Sciences*, 99(16), 10237–10239. <https://doi.org/10.1073/pnas.172399499>
- 1197 Reimold, M., Slifstein, M., Heinz, A., Mueller-Schauenburg, W., & Bares, R. (2005). Effect of Spatial  
1198 Smoothing on t-Maps: Arguments for Going Back from t-Maps to Masked Contrast Images. *Journal of  
1199 Cerebral Blood Flow & Metabolism*, 26(6), 751–759. <https://doi.org/10.1038/sj.jcbfm.9600231>
- 1200 Rogers, J. L., Howard, K. I., & Vessey, J. T. (1993). Using significance tests to evaluate equivalence  
1201 between two experimental groups. *Psychological Bulletin*, 113(3), 553–565. <https://doi.org/10.1037/0033-2909.113.3.553>
- 1202  
1203 Rosa, M., Friston, K., & Penny, W. (2012). Post-hoc selection of dynamic causal models. *Journal of  
1204 Neuroscience Methods*, 208(1), 66–78. <https://doi.org/10.1016/j.jneumeth.2012.04.013>

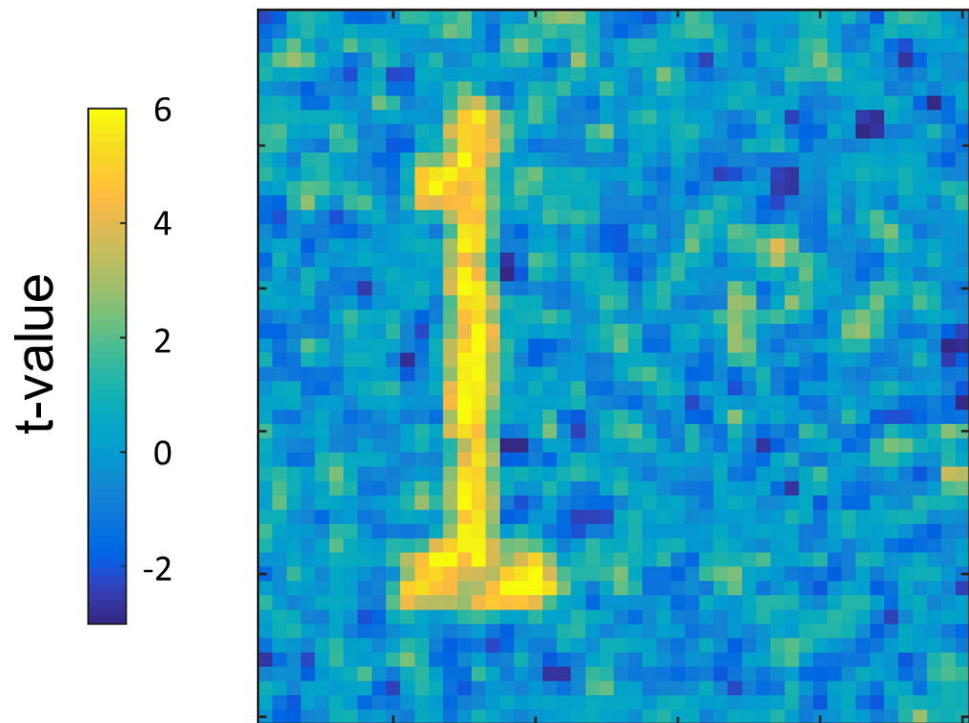
- 1205 Rosenthal, R. (1979). The file drawer problem and tolerance for null results. *Psychological Bulletin*, 86(3),  
1206 638–641. <https://doi.org/10.1037/0033-2909.86.3.638>
- 1207 Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t tests for accepting  
1208 and rejecting the null hypothesis. *Psychonomic Bulletin & Review*, 16(2), 225–237.  
1209 <https://doi.org/10.3758/pbr.16.2.225>
- 1210 Royall, R. M. (1986). The Effect of Sample Size on the Meaning of Significance Tests. *The American*  
1211 *Statistician*, 40(4), 313. <https://doi.org/10.2307/2684616>
- 1212 Royall, R. M. (1997). *Statistical Evidence: A Likelihood Paradigm*. CRC Press.
- 1213 Samartsidis, P., Montagna, S., Laird, A. R., Fox, P. T., Johnson, T. D., & Nichols, T. E. (2020). Estimating  
1214 the prevalence of missing experiments in a neuroimaging meta-analysis. *Research Synthesis Methods*,  
1215 11(6), 866–883. <https://doi.org/10.1002/jrsm.1448>
- 1216 Schatz, P., Jay, K., McComb, J., & McLaughlin, J. (2005). Misuse of statistical tests in publications.  
1217 *Archives of Clinical Neuropsychology*, 20(8), 1053–1059. <https://doi.org/10.1016/j.acn.2005.06.006>
- 1218 Schneider, J. W. (2014). Null hypothesis significance tests. A mix-up of two different theories: the basis  
1219 for widespread confusion and numerous misinterpretations. *Scientometrics*, 102(1), 411–432.  
1220 <https://doi.org/10.1007/s11192-014-1251-5>
- 1221 Schneider, J. W. (2018). NHST is still logically flawed. *Scientometrics*, 115(1), 627–635.  
1222 <https://doi.org/10.1007/s11192-018-2655-4>
- 1223 Schönbrodt, F. D., Wagenmakers, E. J., Zehetleitner, M., & Perugini, M. (2015). Sequential Hypothesis  
1224 Testing with Bayes Factors: Efficiently Testing Mean Differences. *SSRN Electronic Journal*. Published.  
1225 <https://doi.org/10.2139/ssrn.2604513>
- 1226 Schuirmann, D. J. (1987). A comparison of the Two One-Sided Tests Procedure and the Power Approach  
1227 for assessing the equivalence of average bioavailability. *Journal of Pharmacokinetics and*  
1228 *Biopharmaceutics*, 15(6), 657–680. <https://doi.org/10.1007/bf01068419>
- 1229 Schwartzman, A., Dougherty, R., Lee, J., Ghahremani, D., & Taylor, J. (2009). Empirical null and false  
1230 discovery rate analysis in neuroimaging. *NeuroImage*, 44(1), 71–82.  
1231 <https://doi.org/10.1016/j.neuroimage.2008.04.182>
- 1232 Scott, J. G., & Berger, J. O. (2010). Bayes and empirical-Bayes multiplicity adjustment in the variable-  
1233 selection problem. *The Annals of Statistics*, 38(5). <https://doi.org/10.1214/10-aos792>
- 1234 Serlin, R. C., & Lapsley, D. K. (1985). Rationality in psychological research: The good-enough principle.  
1235 *American Psychologist*, 40(1), 73–83. <https://doi.org/10.1037/0003-066x.40.1.73>
- 1236 Serlin, R. C., & Lapsley, D. K. (1993). Rational appraisal of psychological research and the good-enough  
1237 principle. In G. Keren & C. Lewis (Eds.), *A handbook for data analysis in the behavioral sciences: Methodological issues* (pp. 199–228). Lawrence Erlbaum Associates, Inc.
- 1239 Sjölander, A., & Vansteelandt, S. (2019). Frequentist versus Bayesian approaches to multiple testing.  
1240 *European Journal of Epidemiology*, 34(9), 809–821. <https://doi.org/10.1007/s10654-019-00517-2>
- 1241 Sober, E. (2008). *Evidence and Evolution: The Logic Behind the Science*. Cambridge University Press.
- 1242 Stephan, K. E. (2016). *Bayesian inference and generative models* [Slides]. Translational Neuromodeling  
1243 Unit.  
1244 [https://www.tnu.ethz.ch/fileadmin/user\\_upload/teaching/Methods\\_Models2016/10\\_BayesianInference\\_H](https://www.tnu.ethz.ch/fileadmin/user_upload/teaching/Methods_Models2016/10_BayesianInference_H)  
1245 [S2016\\_Handout.pdf](https://www.tnu.ethz.ch/fileadmin/user_upload/teaching/Methods_Models2016/10_BayesianInference_H_S2016_Handout.pdf)

- 1246 Storey, J. D. (2003). The positive false discovery rate: a Bayesian interpretation and the q-value. *The Annals*  
1247 *of Statistics*, *31*(6). <https://doi.org/10.1214/aos/1074290335>
- 1248 Streiner, D. L. (2015). Best (but oft-forgotten) practices: the multiple problems of multiplicity—whether  
1249 and how to correct for many statistical tests. *The American Journal of Clinical Nutrition*, *102*(4), 721–728.  
1250 <https://doi.org/10.3945/ajcn.115.113548>
- 1251 Svensson J., Schain M., Knudsen G.M., Ogden T., Plavén-Sigraý P. (2020). Early stopping in clinical PET  
1252 studies: how to reduce expense and exposure. MedRxiv. <https://doi.org/10.1101/2020.09.13.20192856>
- 1253 Szucs, D., & Ioannidis, J. P. (2020). Sample size evolution in neuroimaging research: An evaluation of  
1254 highly-cited studies (1990–2012) and of latest practices (2017–2018) in high-impact journals. *NeuroImage*,  
1255 *221*, 117164. <https://doi.org/10.1016/j.neuroimage.2020.117164>
- 1256 Szucs, D., & Ioannidis, J. P. A. (2017). When Null Hypothesis Significance Testing Is Unsuitable for  
1257 Research: A Reassessment. *Frontiers in Human Neuroscience*, *11*.  
1258 <https://doi.org/10.3389/fnhum.2017.00390>
- 1259 Turkheimer, F. E., Aston, J. A. D., & Cunningham, V. J. (2004). On the logic of hypothesis testing in  
1260 functional imaging. *European Journal of Nuclear Medicine and Molecular Imaging*, *31*(5), 725–732.  
1261 <https://doi.org/10.1007/s00259-003-1387-7>
- 1262 Wagenmakers, E. J. (2007). A practical solution to the pervasive problems of p values. *Psychonomic*  
1263 *Bulletin & Review*, *14*(5), 779–804. <https://doi.org/10.3758/bf03194105>
- 1264 Wagenmakers, E. J., Lee, M., Lodewyckx, T., & Iverson, G. J. (2008). Bayesian Versus Frequentist  
1265 Inference. In H. Hoijtink, I. Klugkist, & P. A. Boelen (Eds.), *Bayesian Evaluation of Informative*  
1266 *Hypotheses. Statistics for Social and Behavioral Sciences*. (pp. 181–207). Springer, New York, NY.
- 1267 Wagenmakers, E. J., Lodewyckx, T., Kuriyal, H., & Grasman, R. (2010). Bayesian hypothesis testing for  
1268 psychologists: A tutorial on the Savage–Dickey method. *Cognitive Psychology*, *60*(3), 158–189.  
1269 <https://doi.org/10.1016/j.cogpsych.2009.12.001>
- 1270 Wagenmakers, E. J., Verhagen, J., Ly, A., Matzke, D., Steingroever, H., Rouder, J. N., & Morey, R. D.  
1271 (2017). The Need for Bayesian Hypothesis Testing in Psychological Science. In S. O. Lilienfeld & I. D.  
1272 Waldman (Eds.), *Psychological science under scrutiny: Recent challenges and proposed solutions* (pp.  
1273 123–138). Wiley Blackwell.
- 1274 Wasserstein, R. L., & Lazar, N. A. (2016). The ASA Statement on p-Values: Context, Process, and Purpose.  
1275 *The American Statistician*, *70*(2), 129–133. <https://doi.org/10.1080/00031305.2016.1154108>
- 1276 Wellek, S. (2010). *Testing Statistical Hypotheses of Equivalence and Noninferiority, Second Edition*.  
1277 Taylor & Francis.
- 1278 Welvaert, M., & Rosseel, Y. (2013). On the Definition of Signal-To-Noise Ratio and Contrast-To-Noise  
1279 Ratio for fMRI Data. *PLoS ONE*, *8*(11), e77089. <https://doi.org/10.1371/journal.pone.0077089>
- 1280 Westfall, P., Johnson, W. O., & Utts, J. M. (1997). A Bayesian perspective on the Bonferroni adjustment.  
1281 *Biometrika*, *84*(2), 419–427. <https://doi.org/10.1093/biomet/84.2.419>
- 1282 Westlake, W. J. (1972). Use of Confidence Intervals in Analysis of Comparative Bioavailability Trials.  
1283 *Journal of Pharmaceutical Sciences*, *61*(8), 1340–1341. <https://doi.org/10.1002/jps.2600610845>
- 1284 Woo, C. W., Krishnan, A., & Wager, T. D. (2014). Cluster-extent based thresholding in fMRI analyses:  
1285 Pitfalls and recommendations. *NeuroImage*, *91*, 412–419.  
1286 <https://doi.org/10.1016/j.neuroimage.2013.12.058>

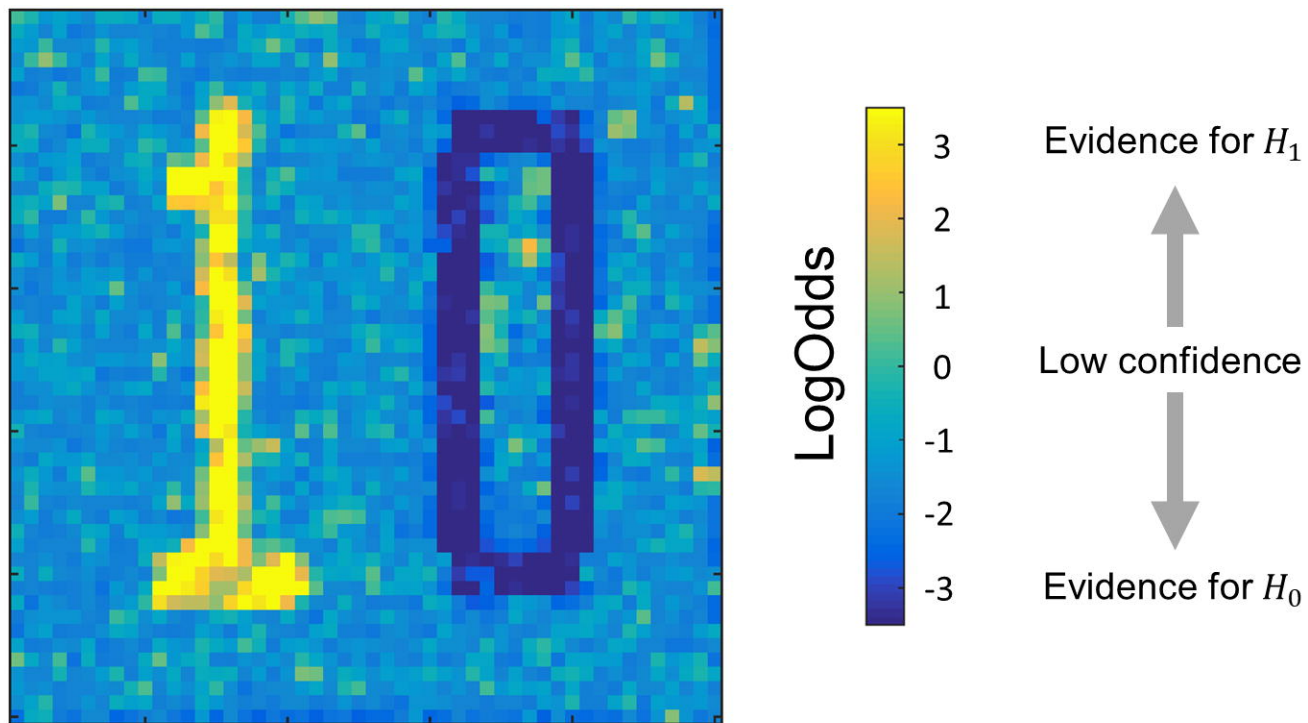
- 1287 Woolrich, M. W., Jbabdi, S., Patenaude, B., Chappell, M., Makni, S., Behrens, T., Beckmann, C.,  
1288 Jenkinson, M., & Smith, S. M. (2009). Bayesian analysis of neuroimaging data in FSL. *NeuroImage*, *45*(1),  
1289 S173–S186. <https://doi.org/10.1016/j.neuroimage.2008.10.055>
- 1290 Worsley, K., & Friston, K. (1995). Analysis of fMRI Time-Series Revisited—Again. *NeuroImage*, *2*(3),  
1291 173–181. <https://doi.org/10.1006/nimg.1995.1023>



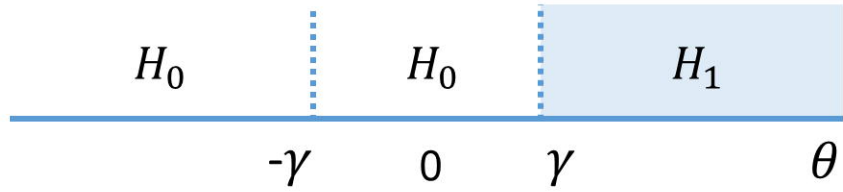
# Classical NHST



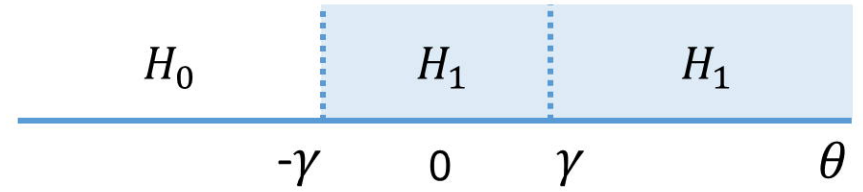
# Bayesian inference



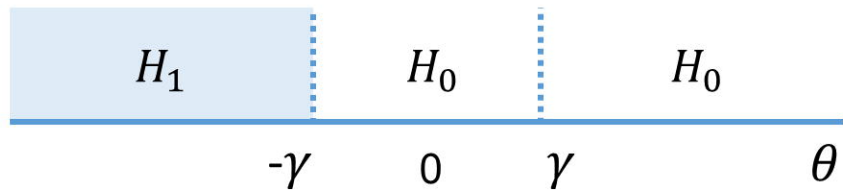
A) Superiority test (positive direction)



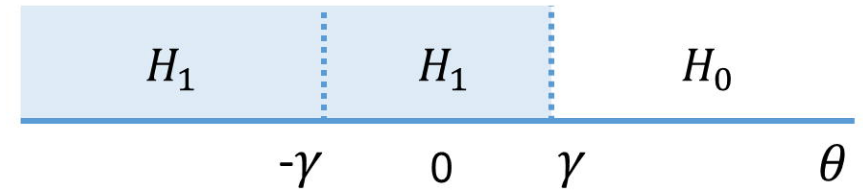
B) Non-Inferiority test (positive direction)



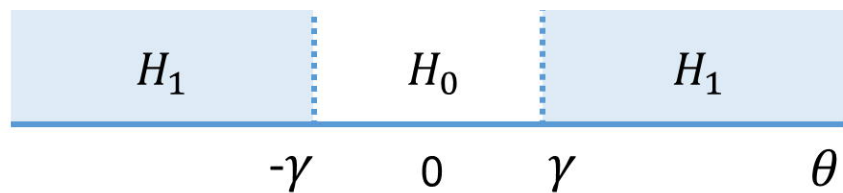
C) Superiority test (negative direction)



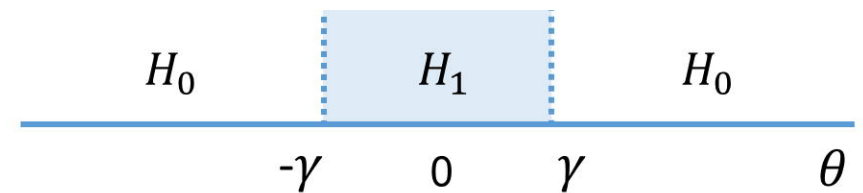
D) Non-Inferiority test (negative direction)

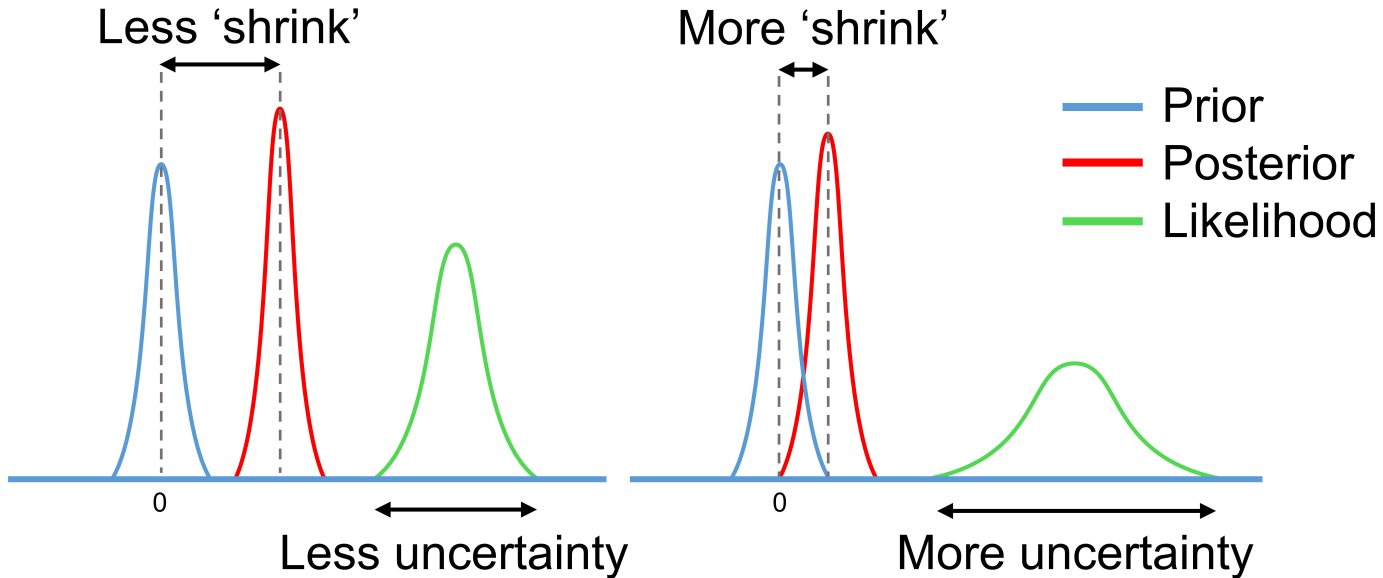


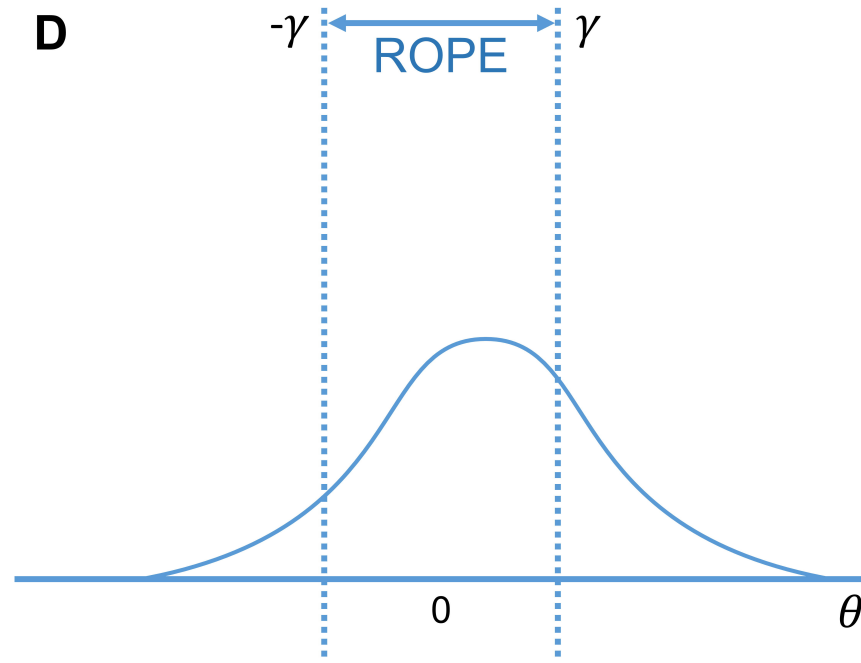
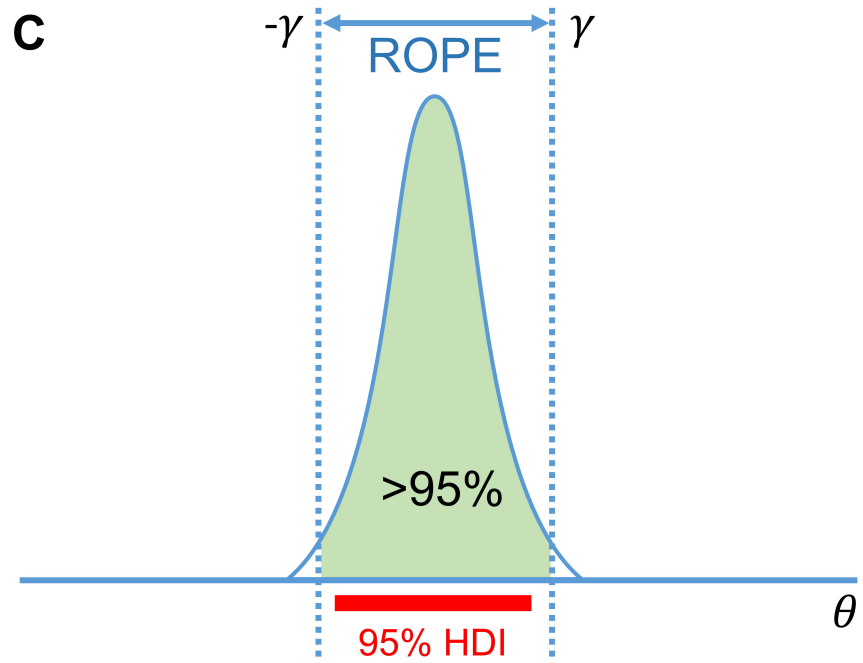
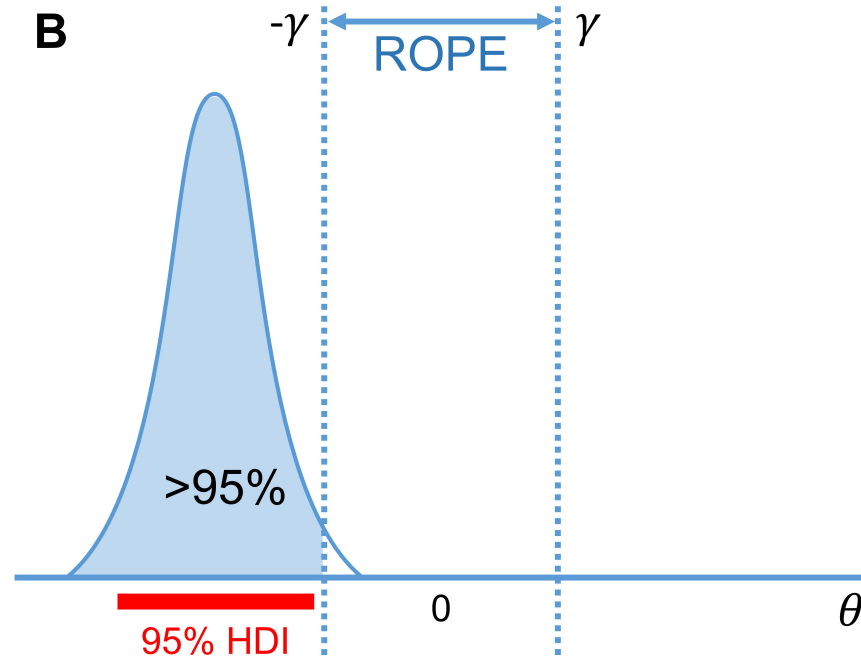
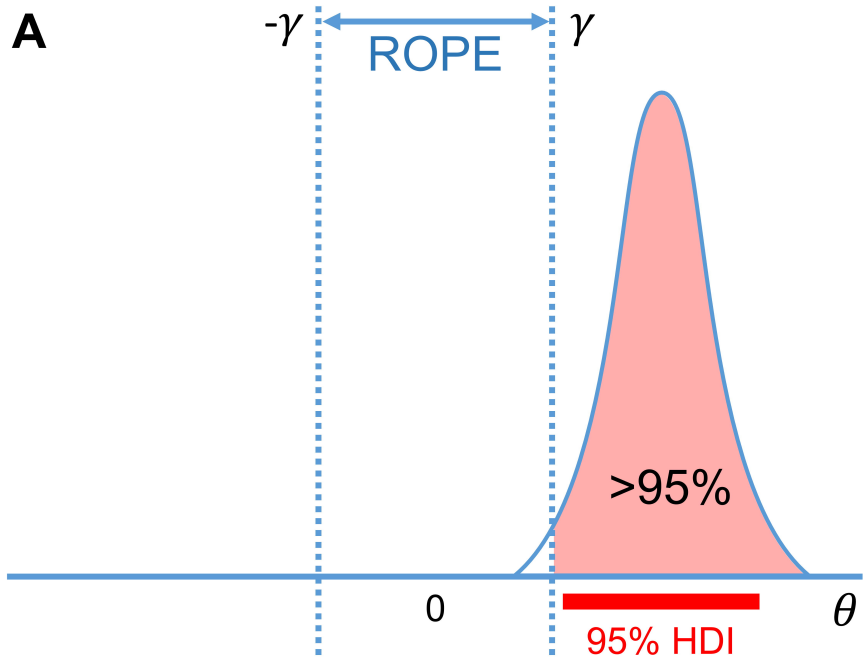
E) Minimum-effect test



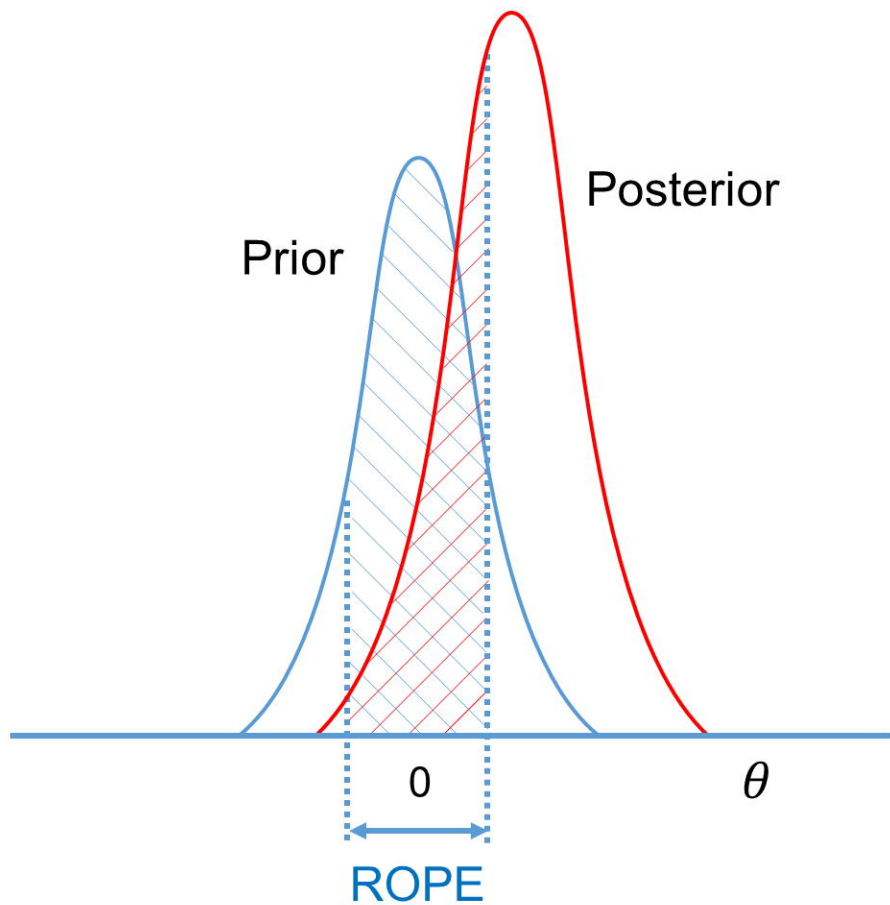
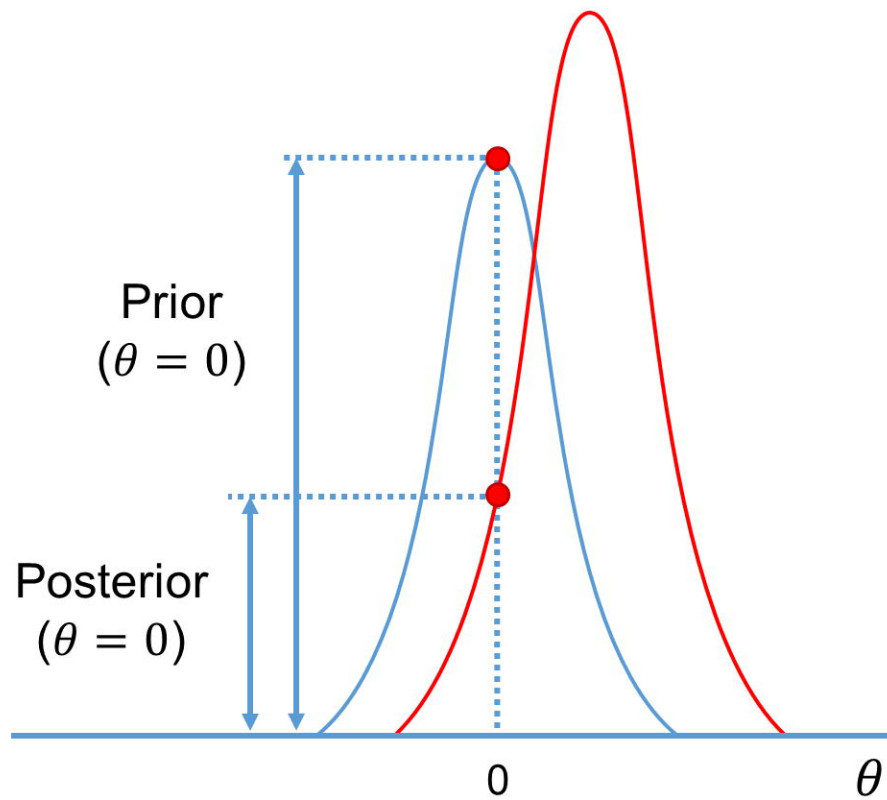
F) Equivalence test (TOST)





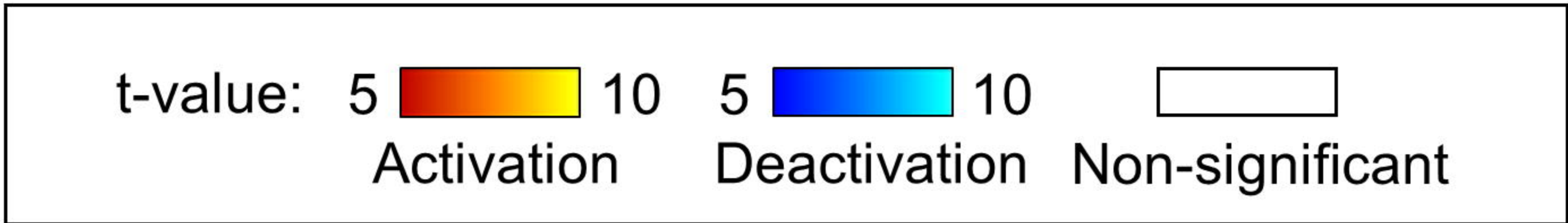




**A***BF(ROPE)*ROPE  $\rightarrow$  0**B***BF(SDR)*



# Classical NHST



Emotion > Shape

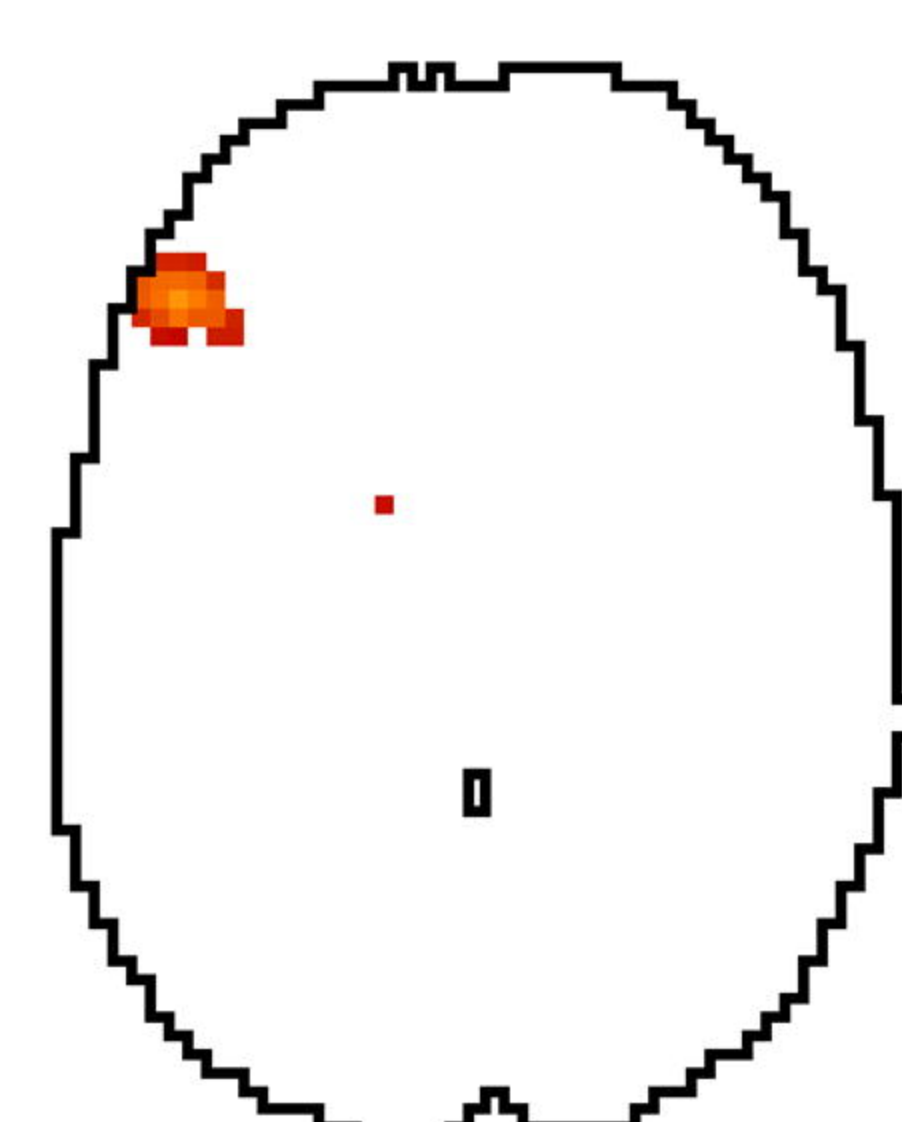
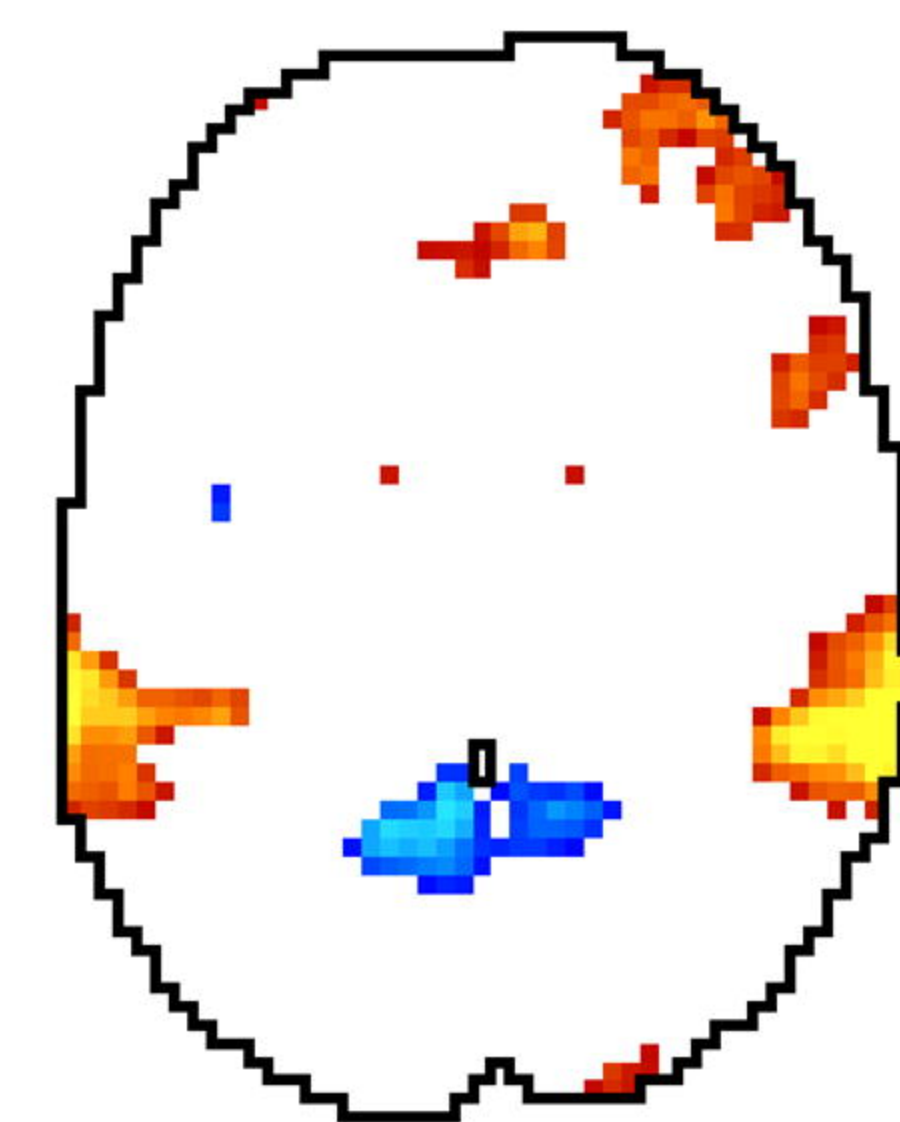
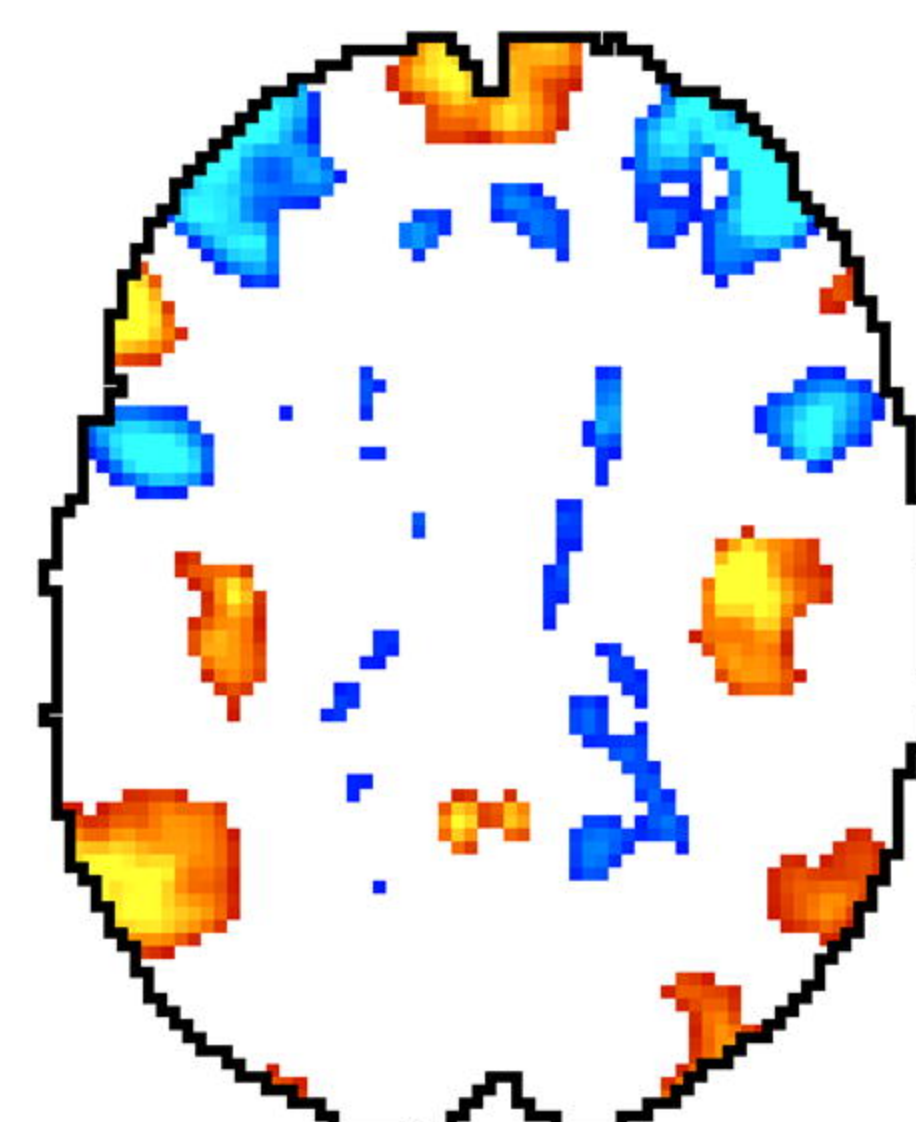
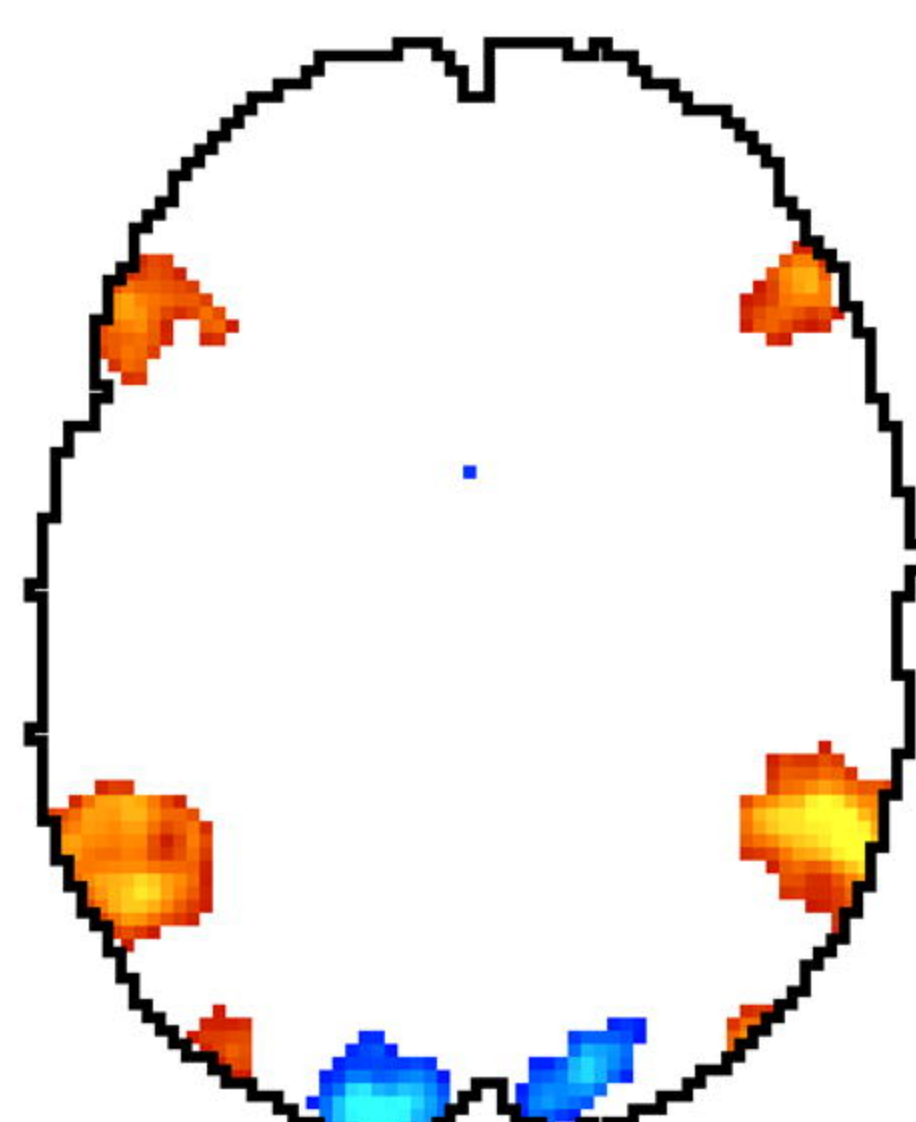
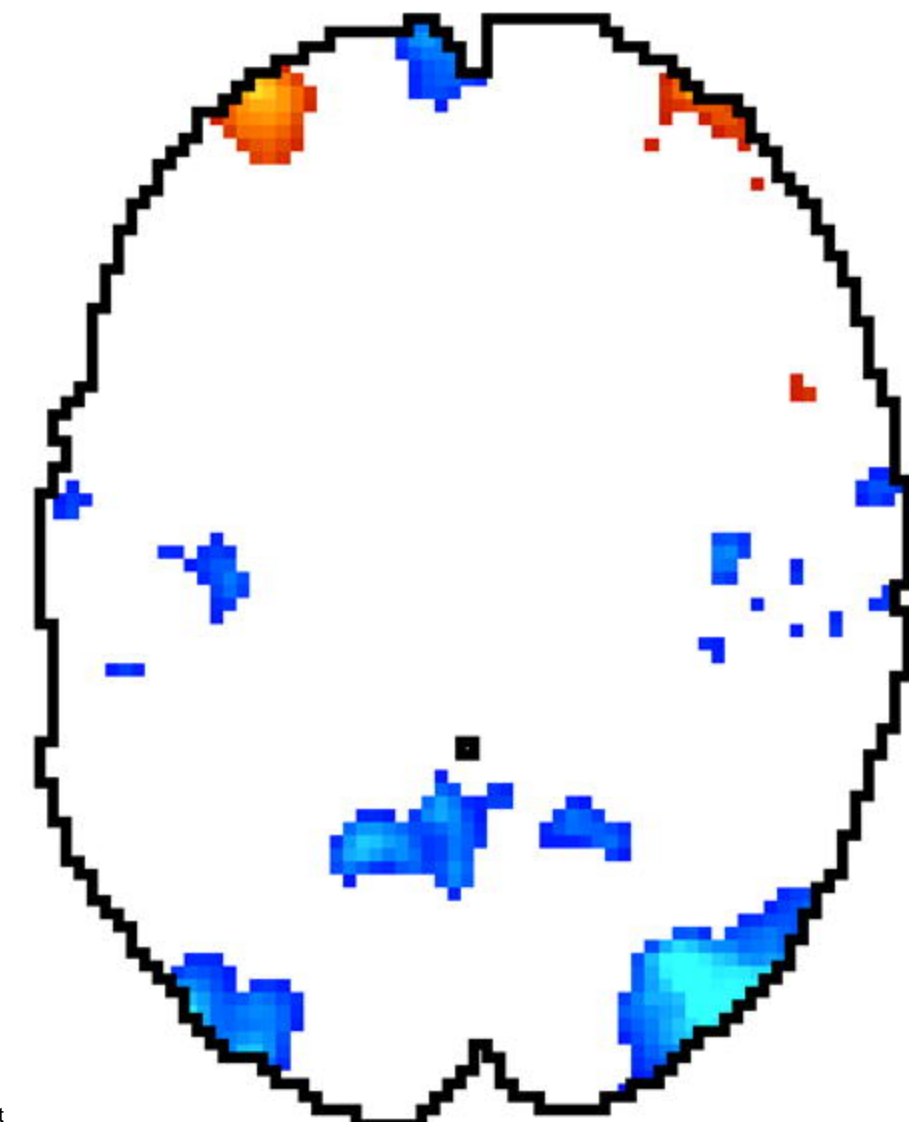
2-back > 0-back

Social > Random

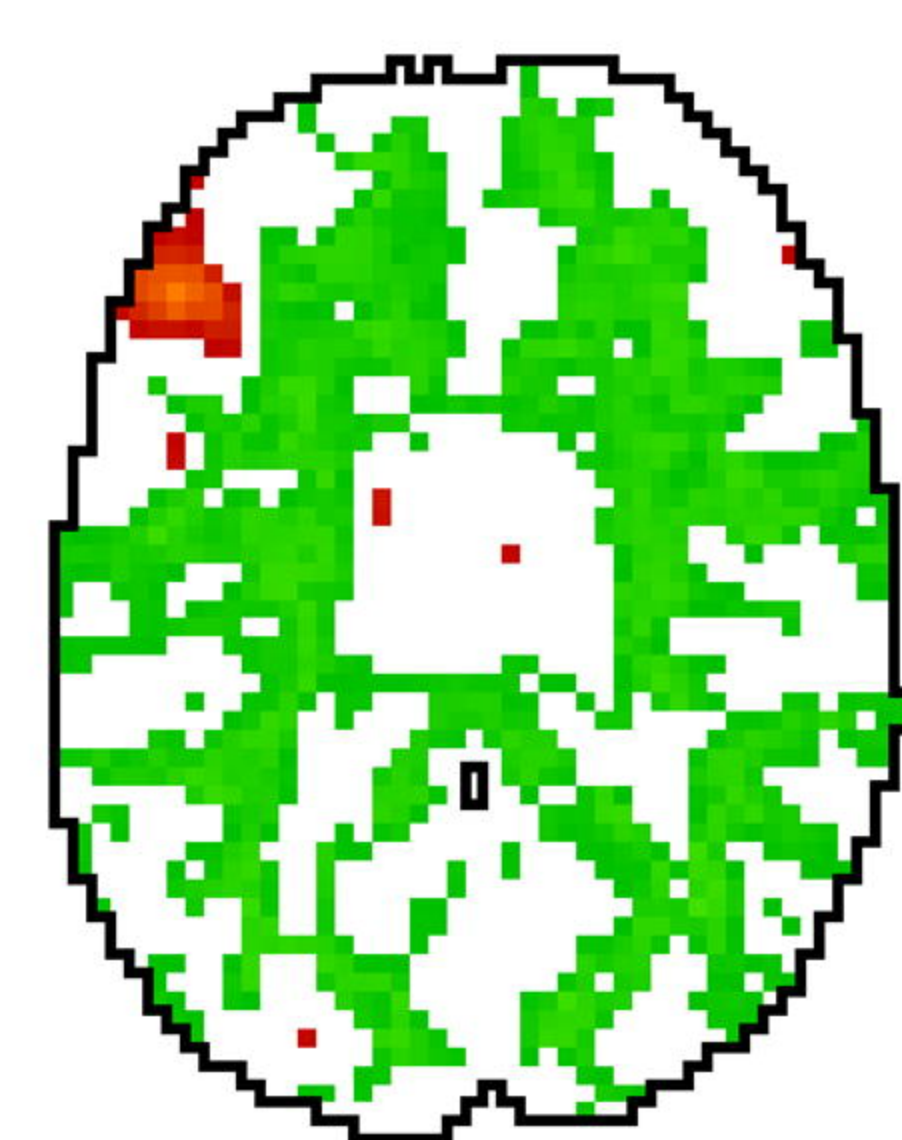
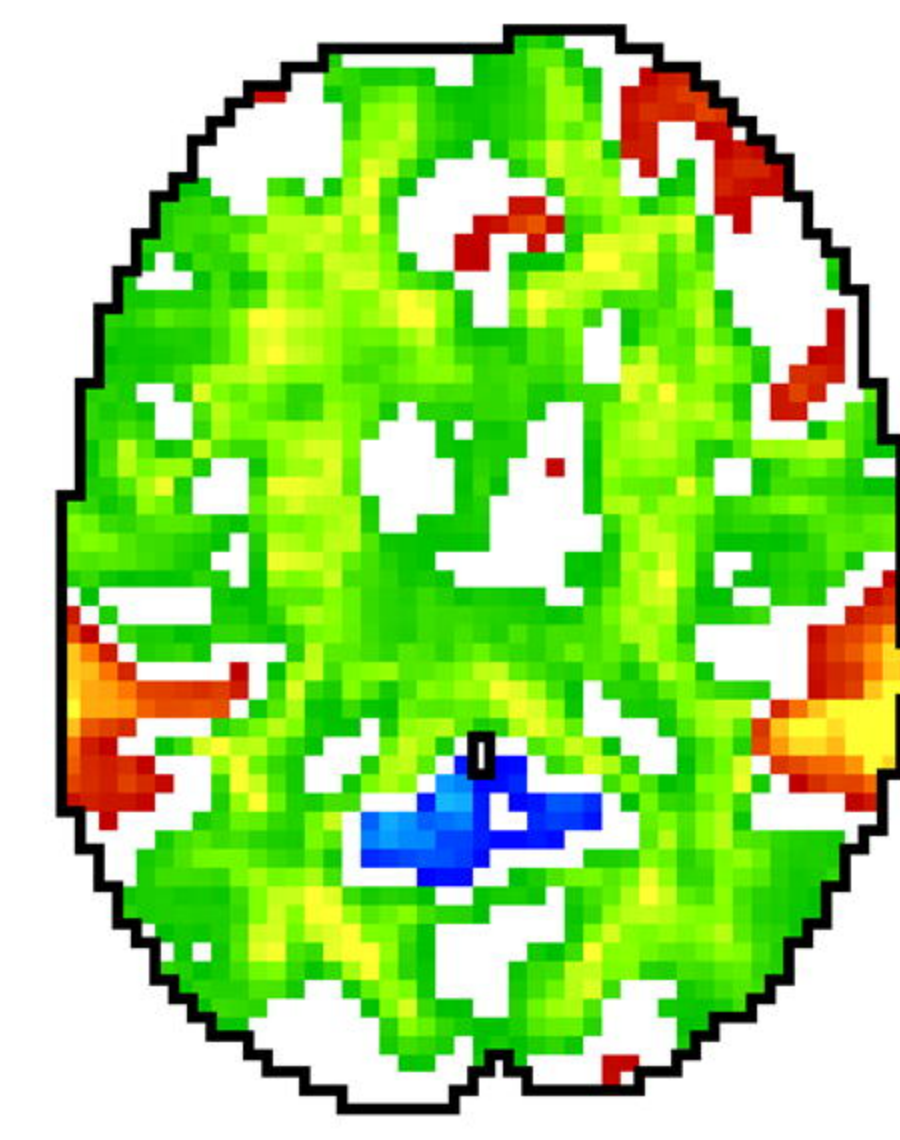
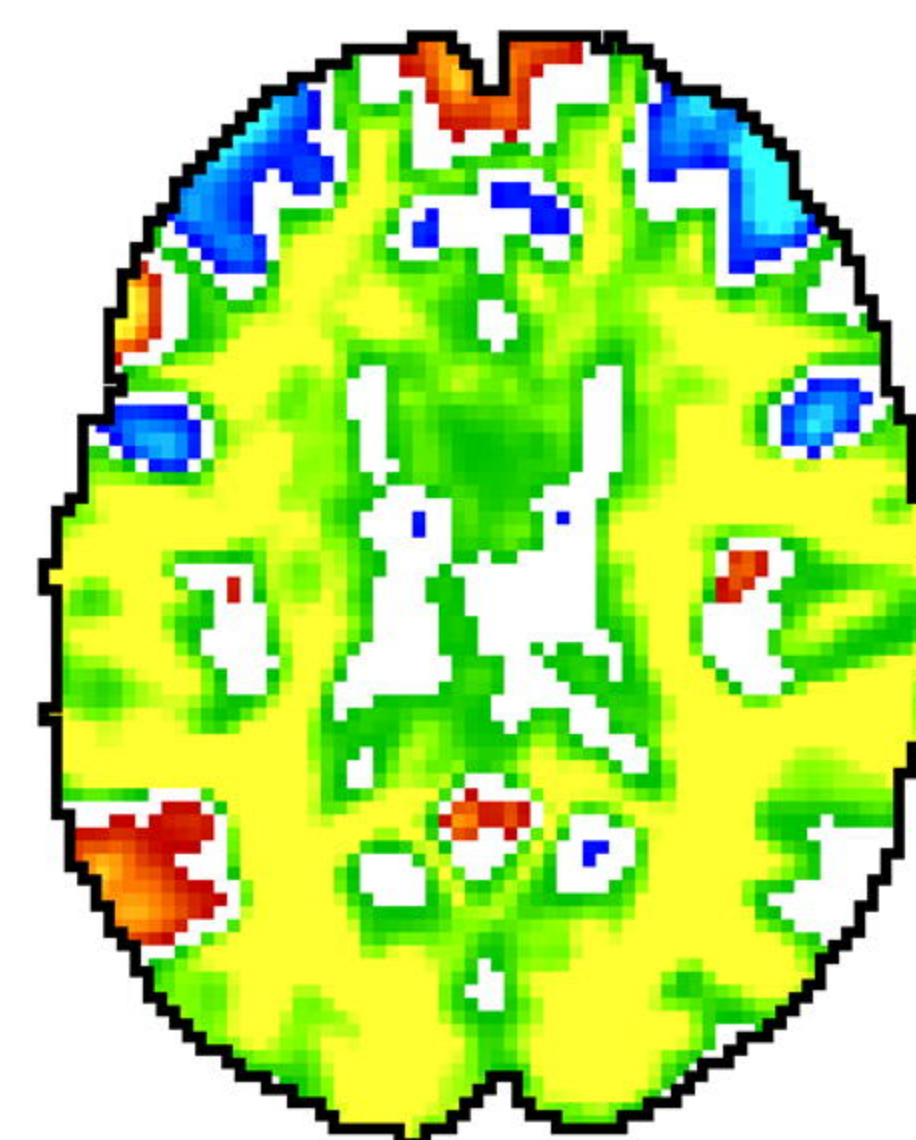
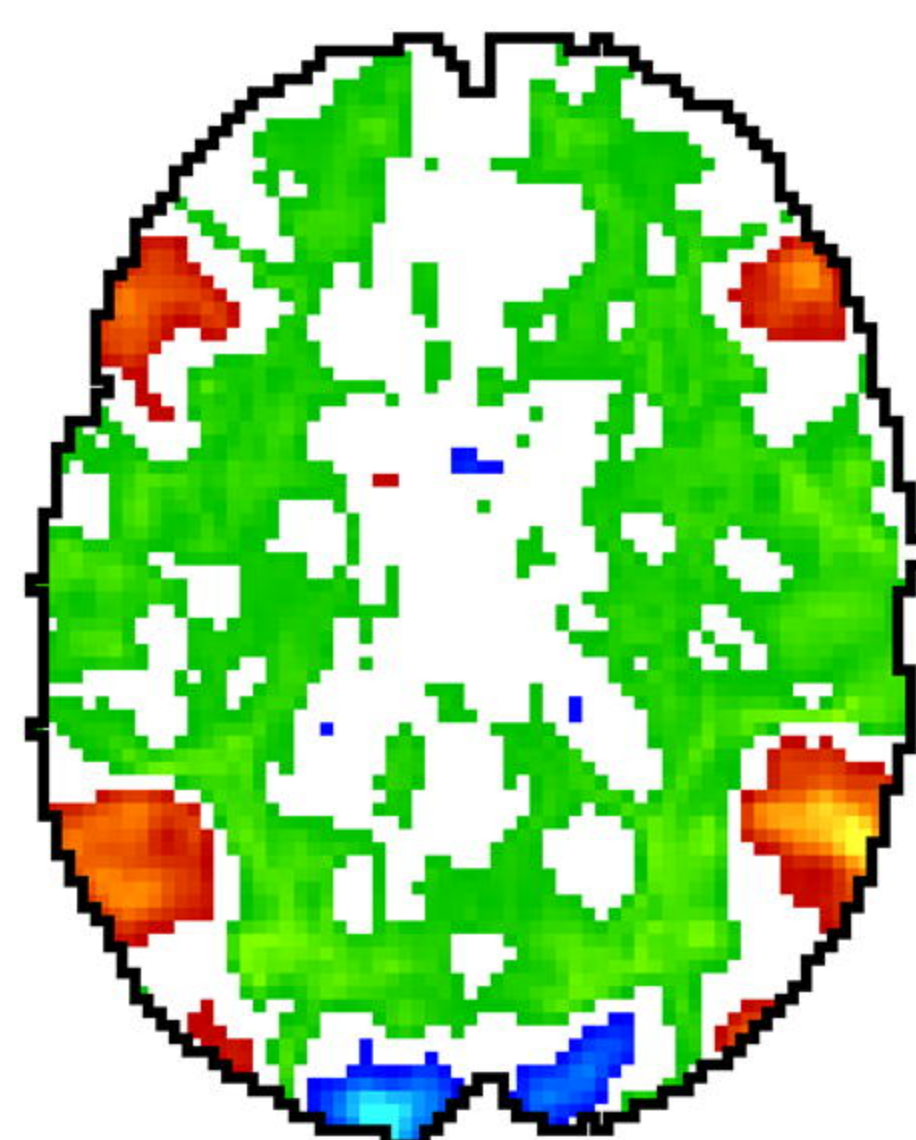
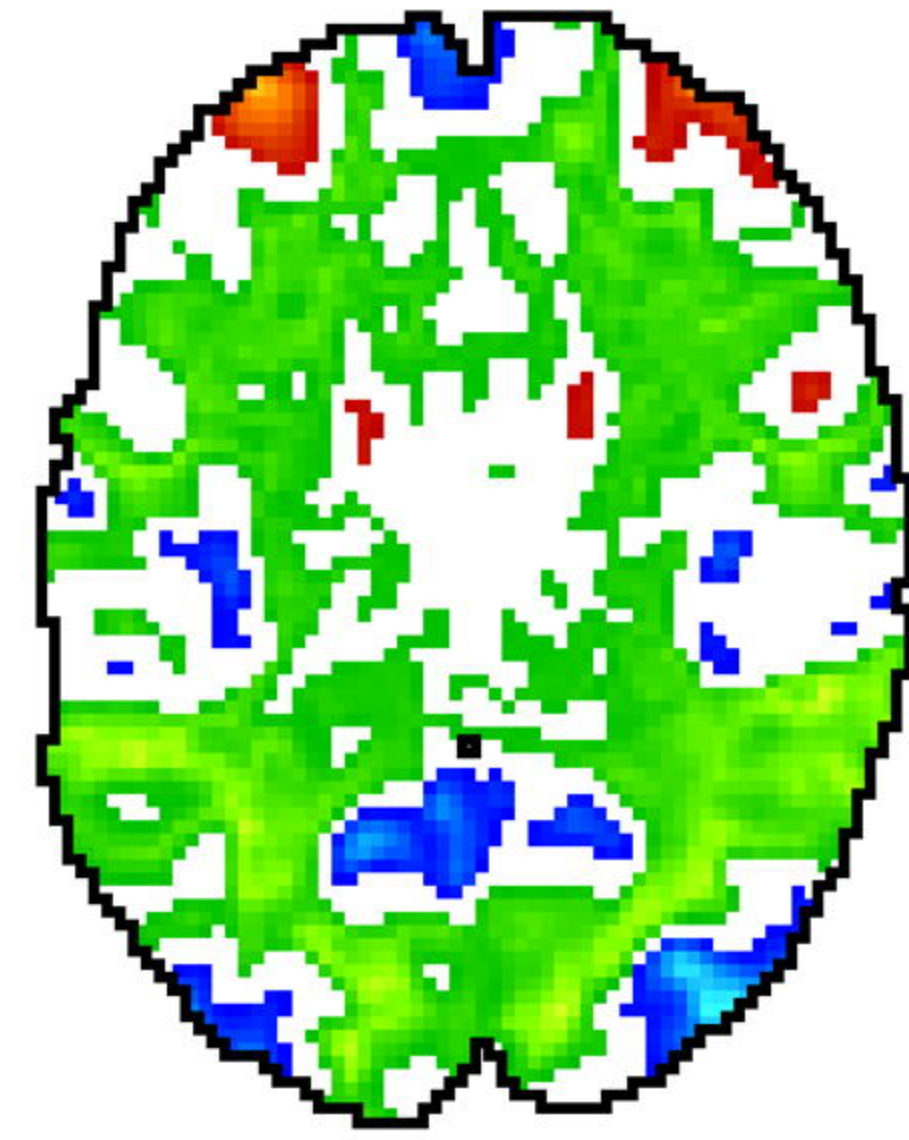
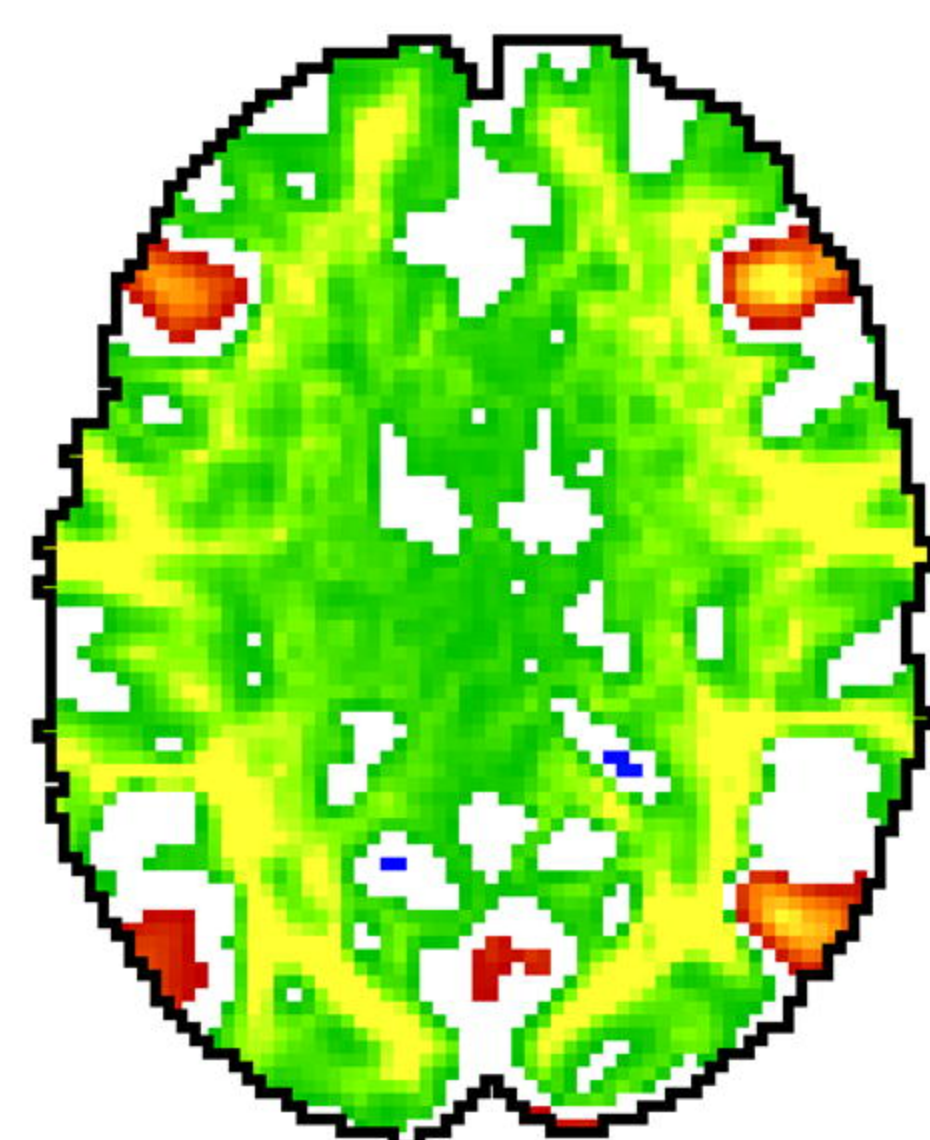
Story > Math

Stop > Go

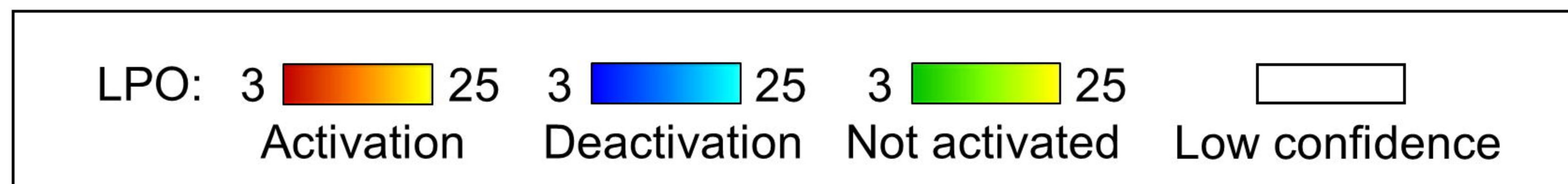
Switch > No Switch



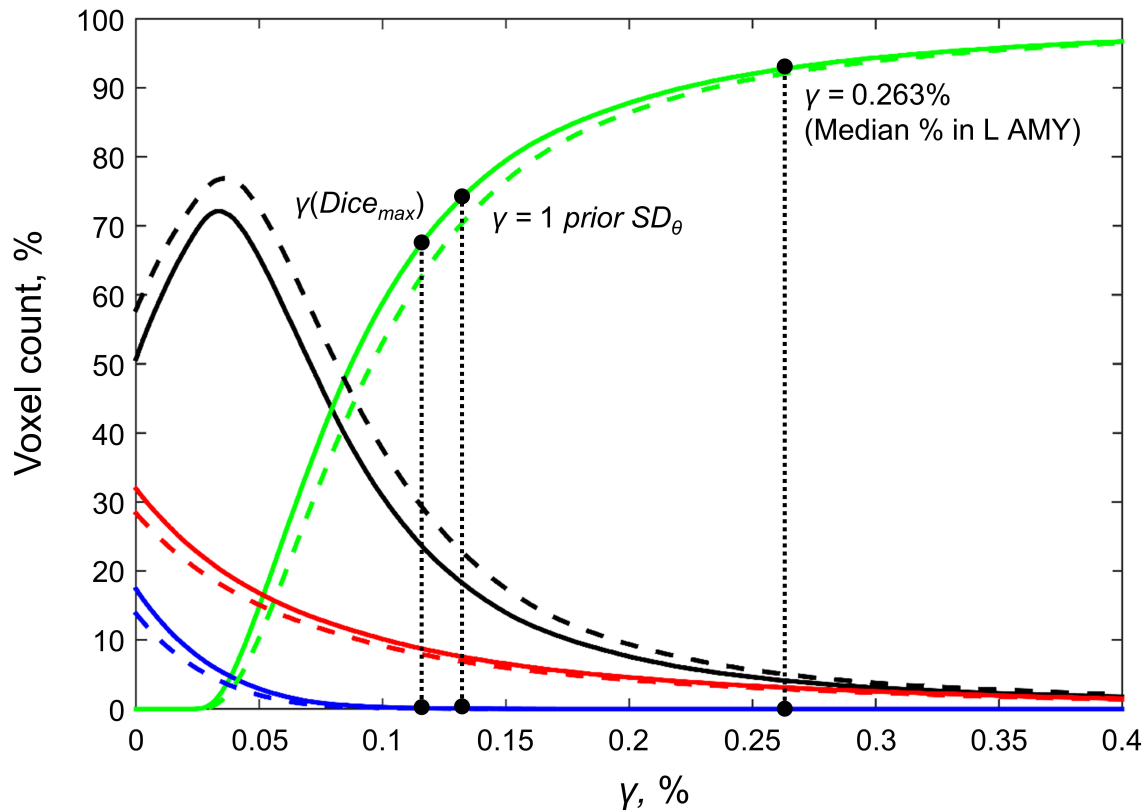
bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.13.442711>; this version posted June 2, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.



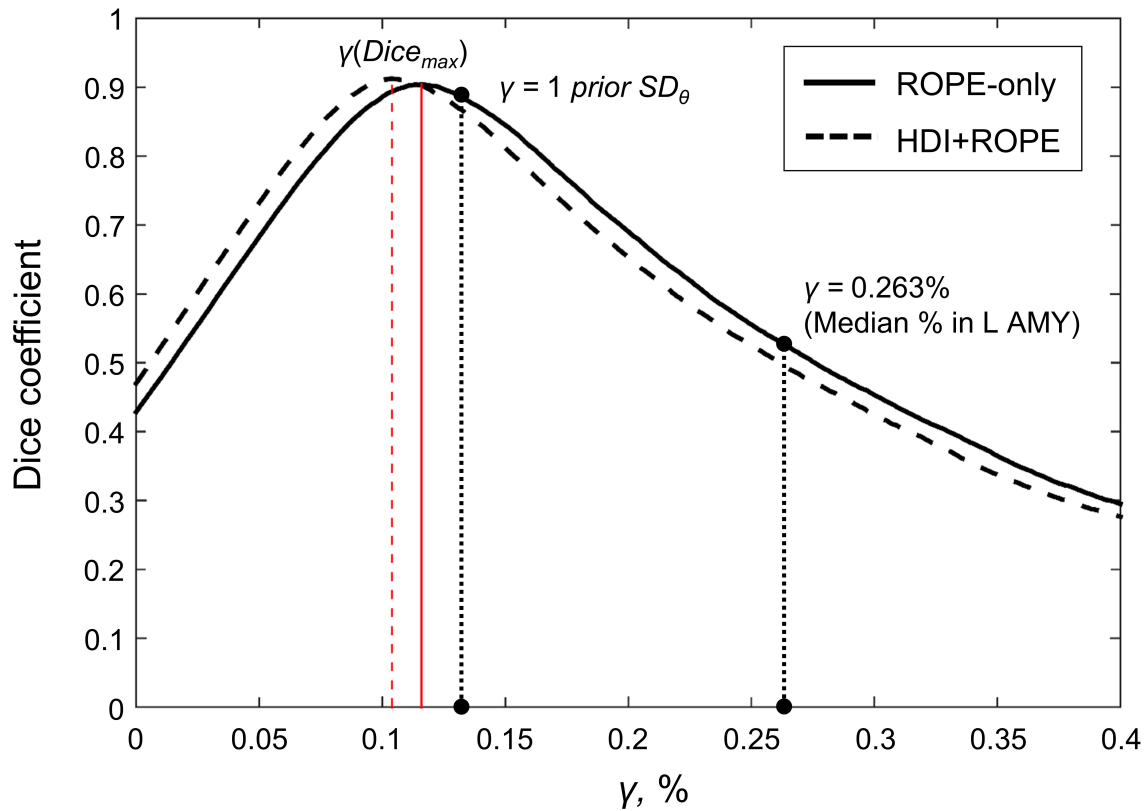
# Bayesian inference



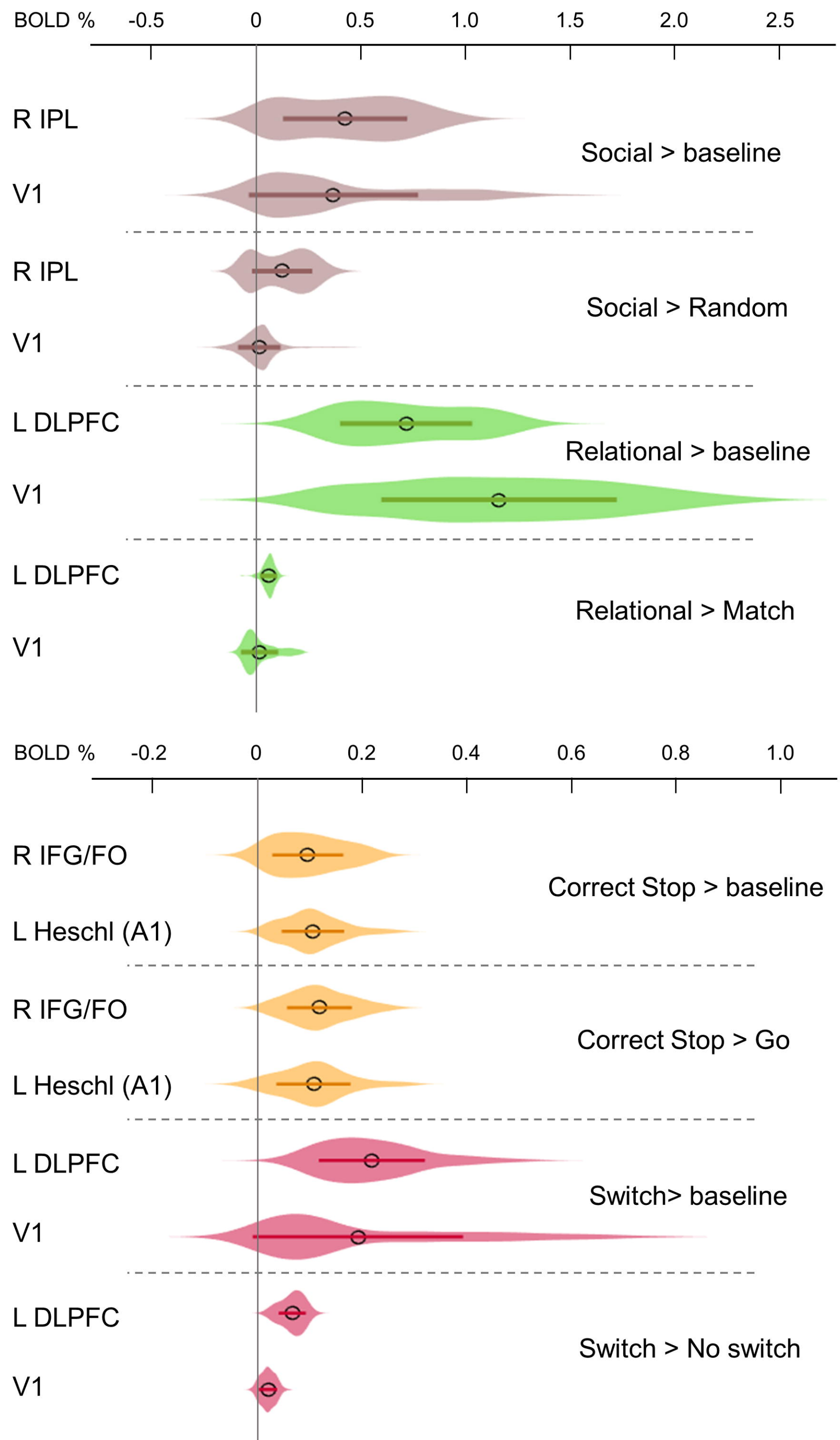
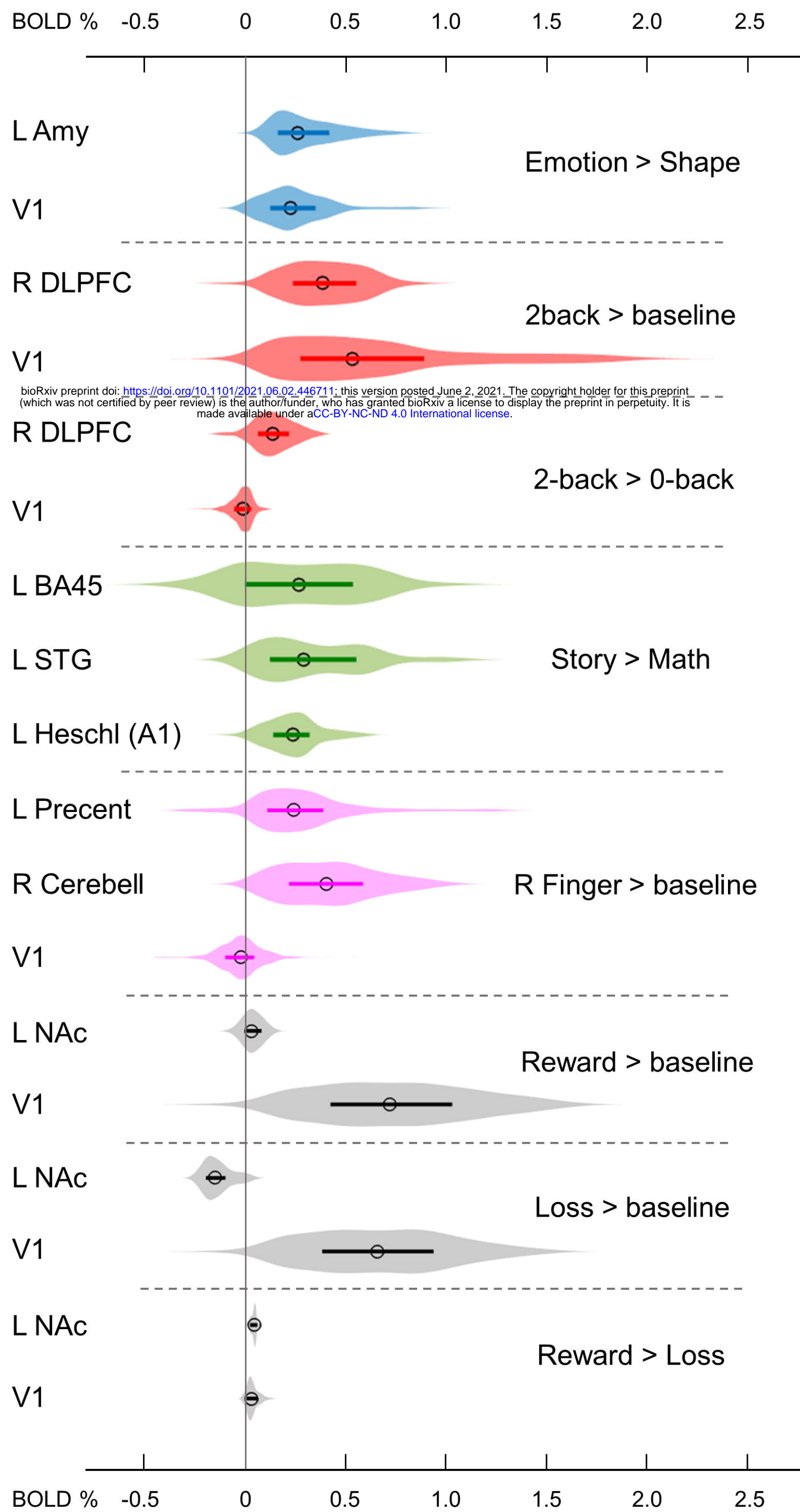




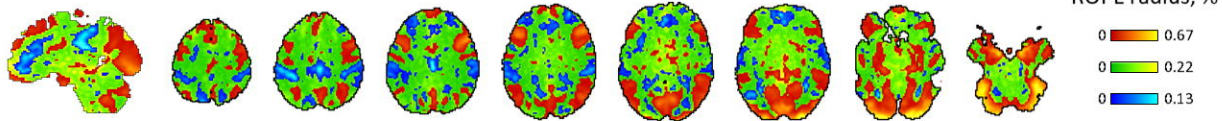
- Activated (ROPE-only)
- - - Activated (HDI+ROPE)
- Deactivated (ROPE-only)
- - - Deactivated (HDI+ROPE)
- Not activated (ROPE-only)
- - - Not activated (HDI+ROPE)
- Low confidence (ROPE-only)
- - - Low confidence (HDI+ROPE)



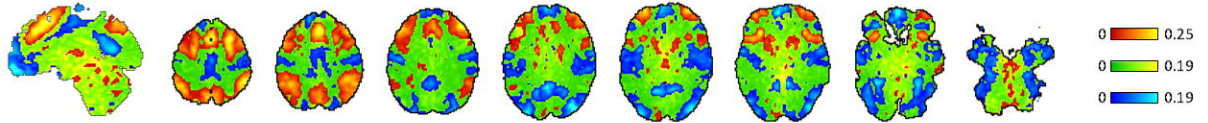




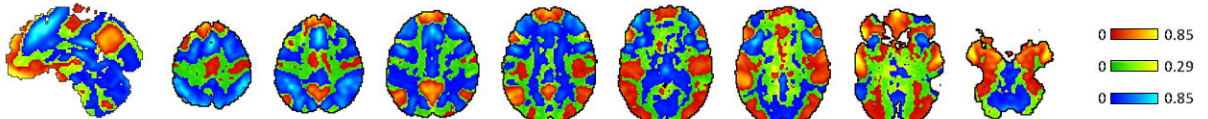
### Emotion > Shape



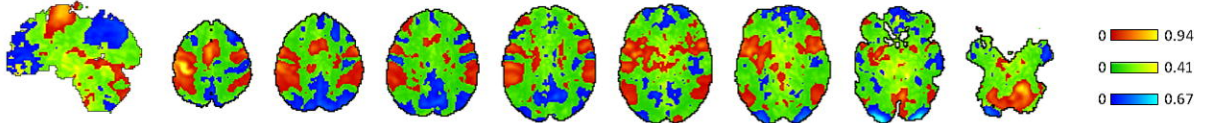
### 2-back > 0-back



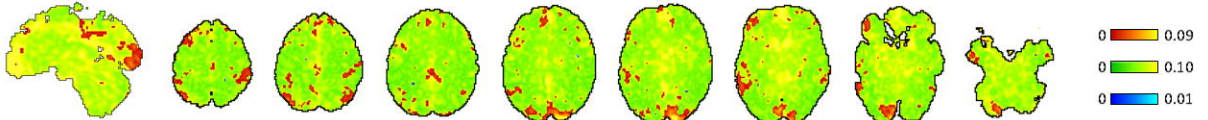
### Story > Math



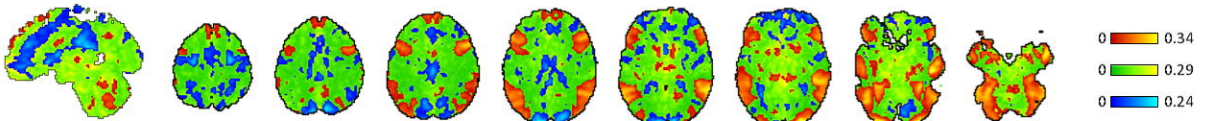
### Right finger > baseline



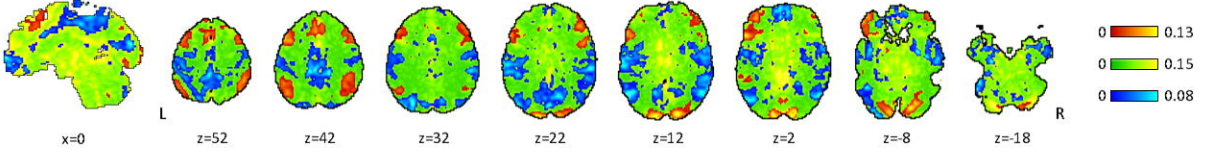
### Reward > Loss



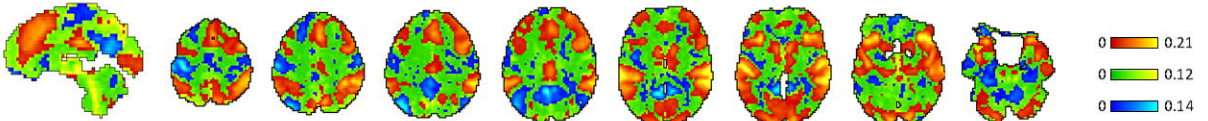
### Social > Random



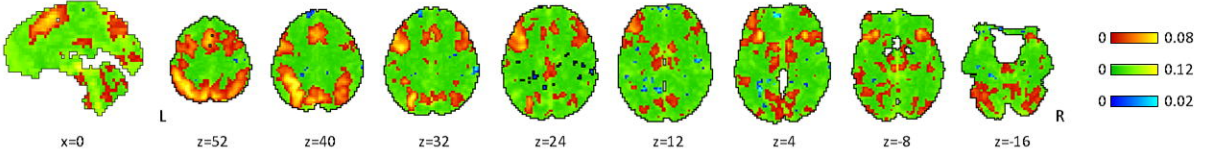
### Relational > Match



### Correct Stop > Go

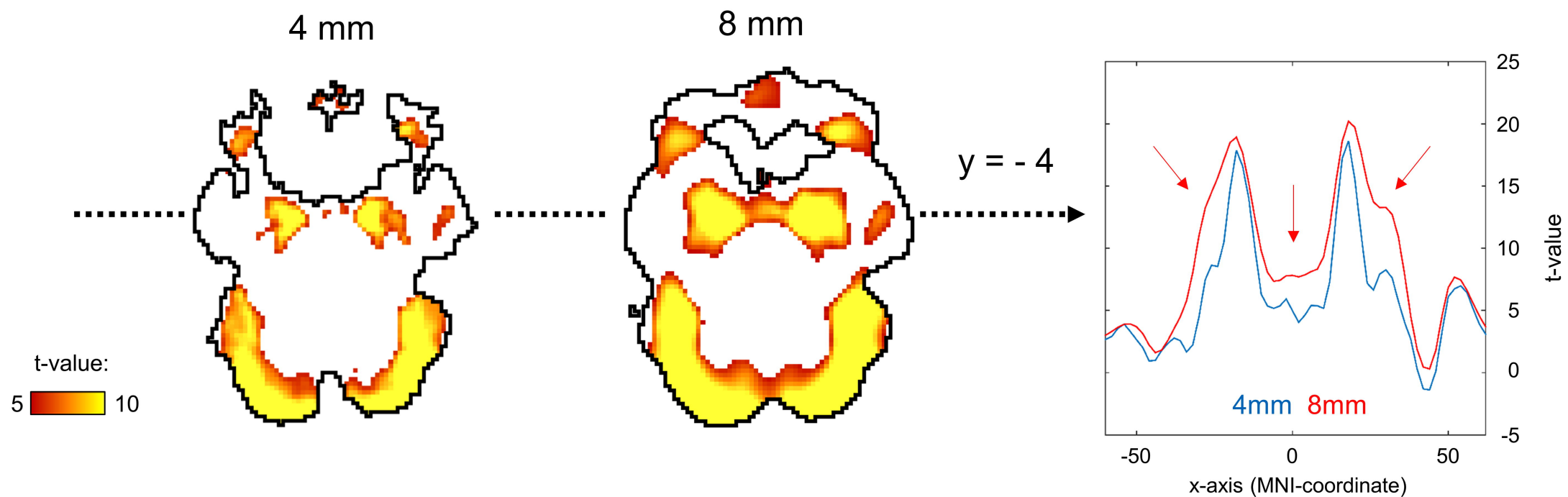


### Switch > No switch

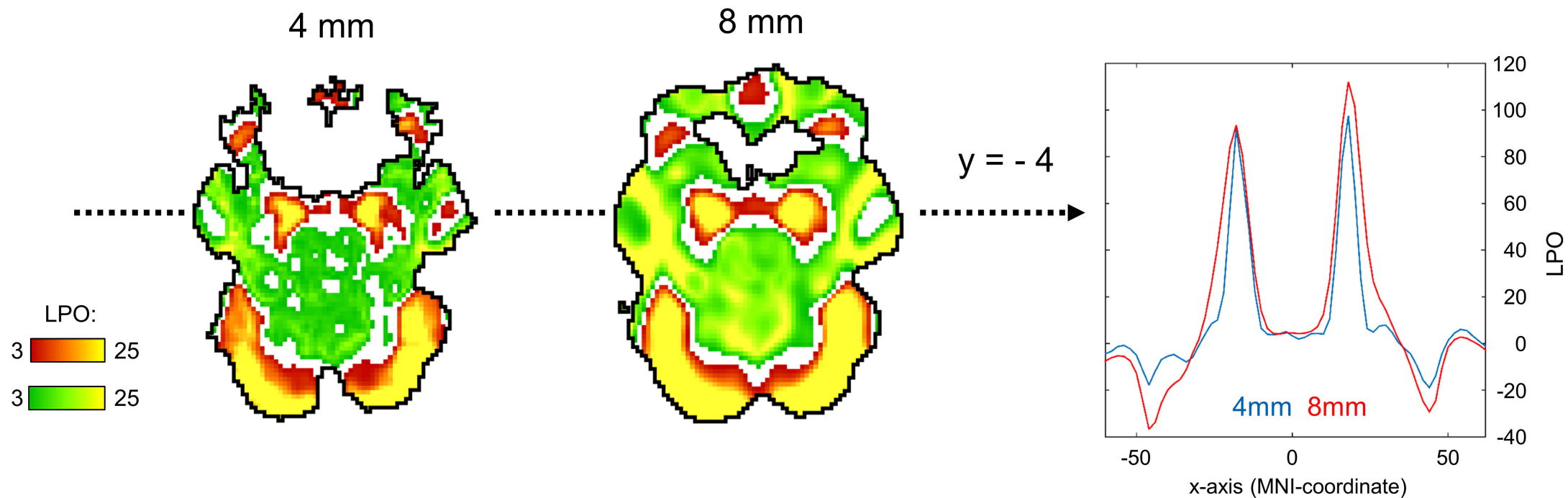




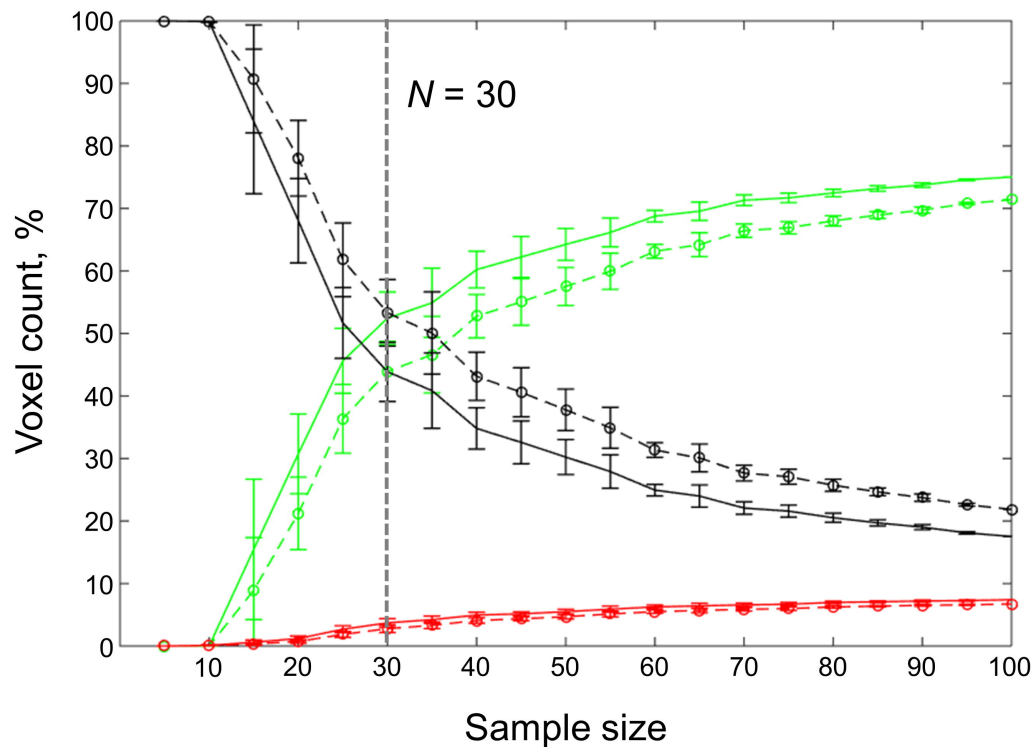
Classical NHST



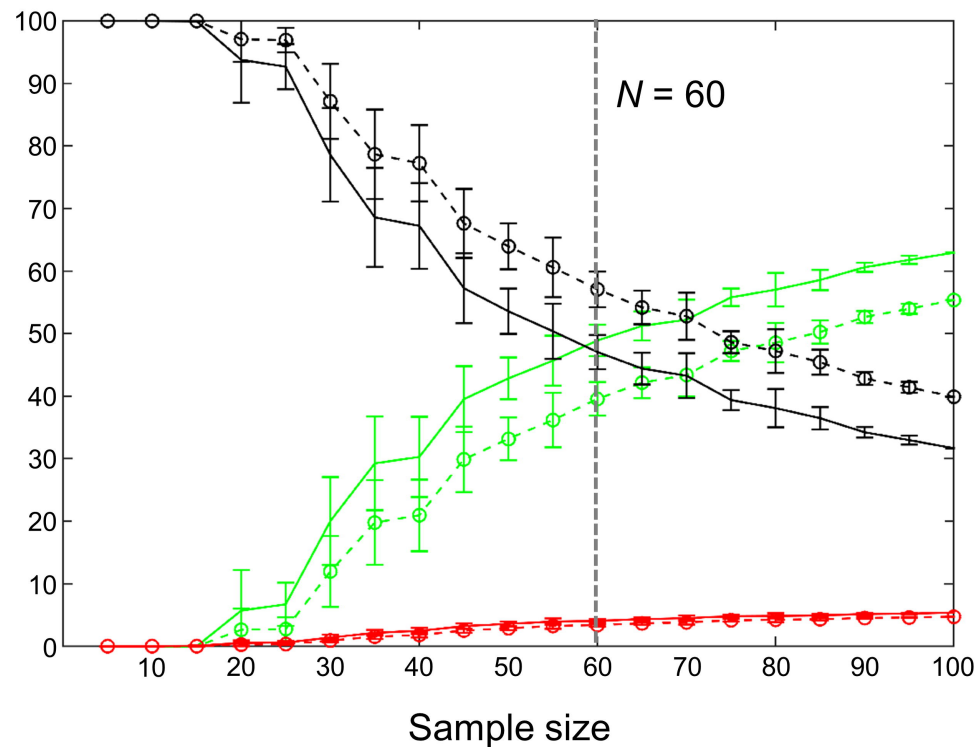
Bayesian inference



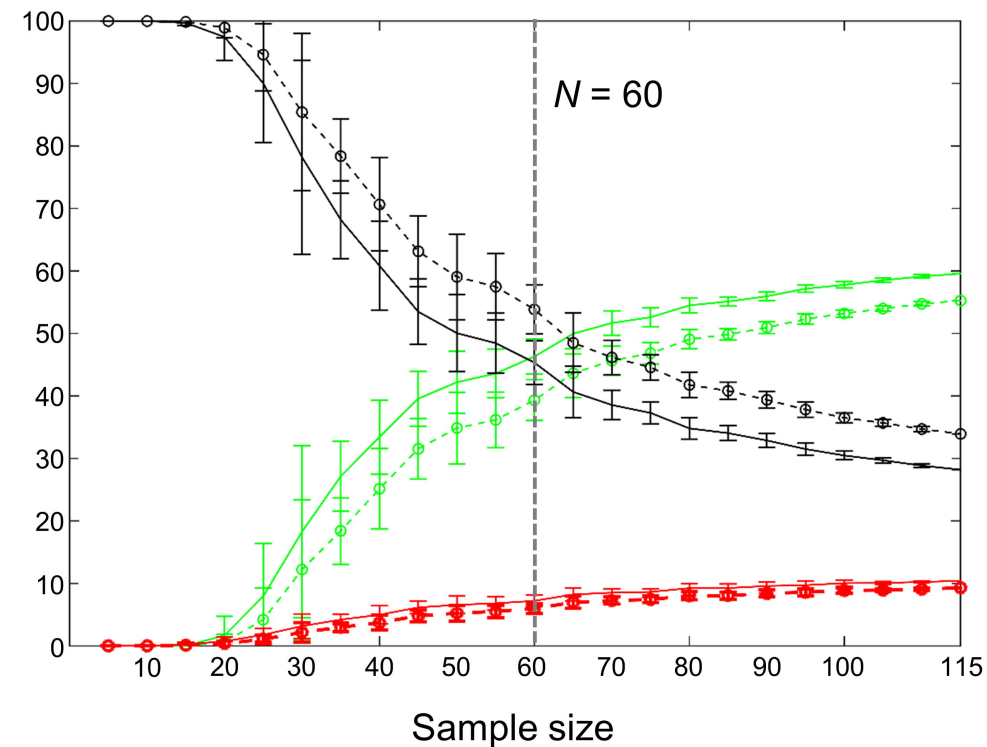
**A** Emotional processing task  
(block design, two sessions, 352 scans)



**B** Emotional processing task  
(block design, one session, 176 scans)



**C** Stop-signal task  
(event-related design, 184 scans)

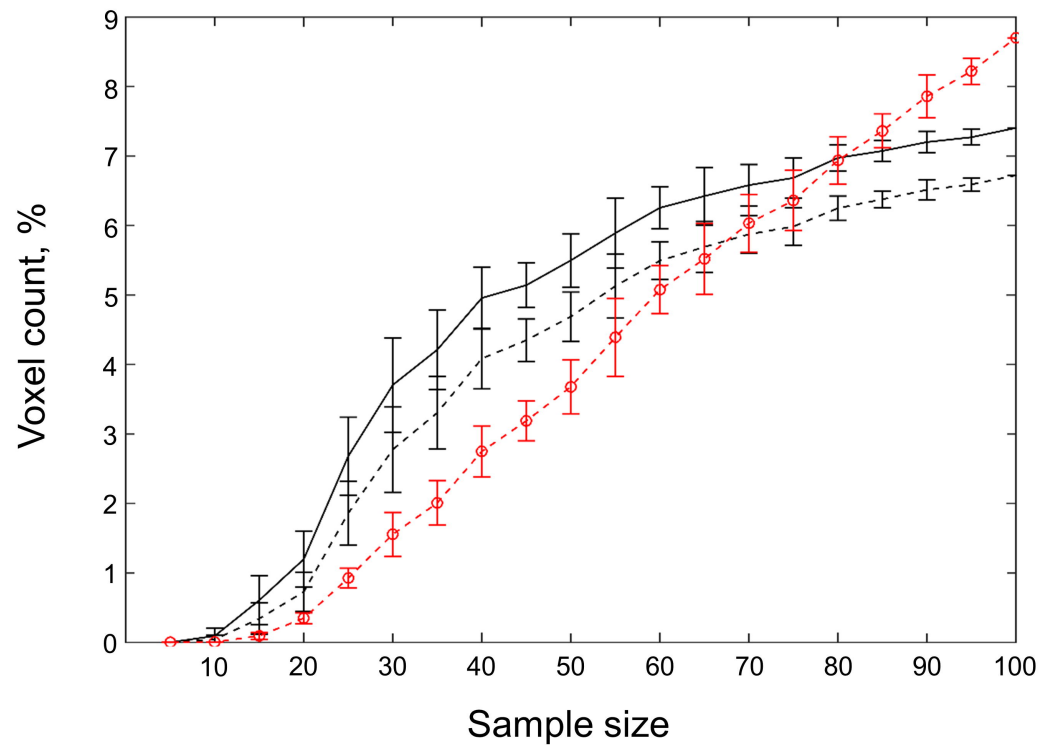


— Activated (ROPE-only)    — Not activated (ROPE-only)    — Low confidence (ROPE-only)  
- - Activated (HDI+ROPE)    - - Not activated (HDI+ROPE)    - - Low confidence (HDI+ROPE)

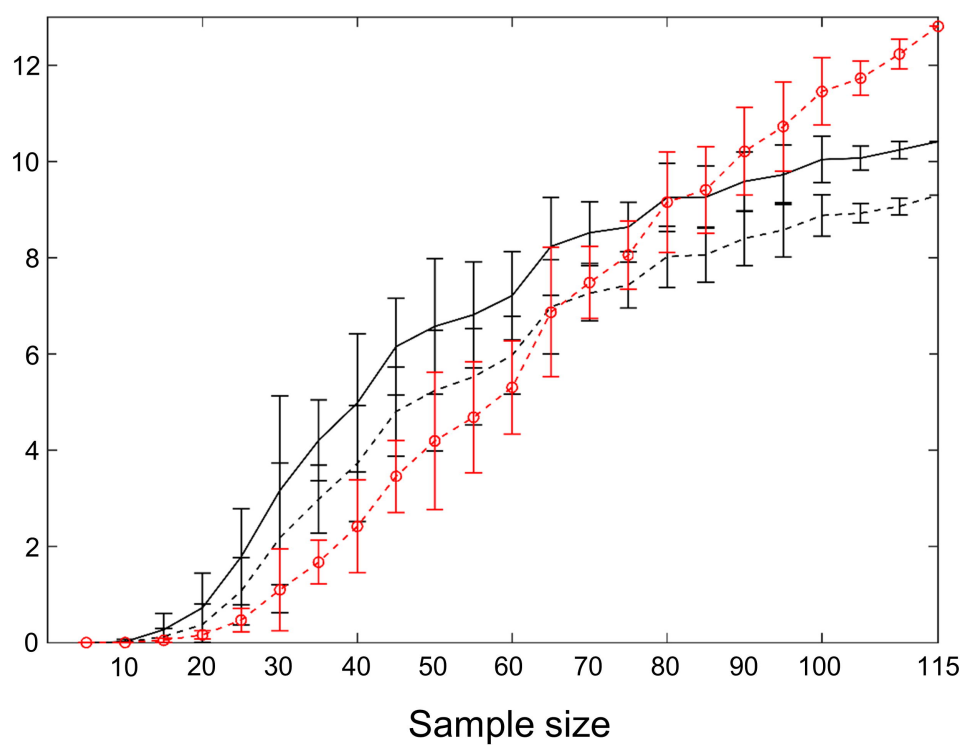


**A**

Emotional processing task

**B**

Stop-signal task



— ROPE-only

- - - HDI+ROPE

- - - Classic NHST

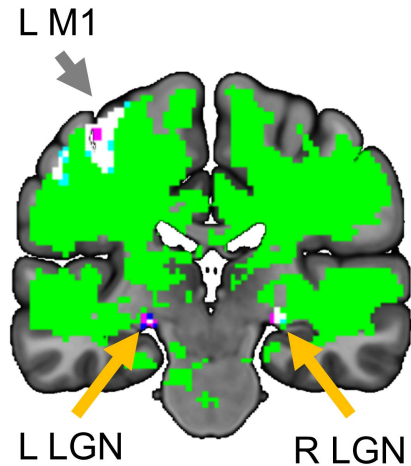
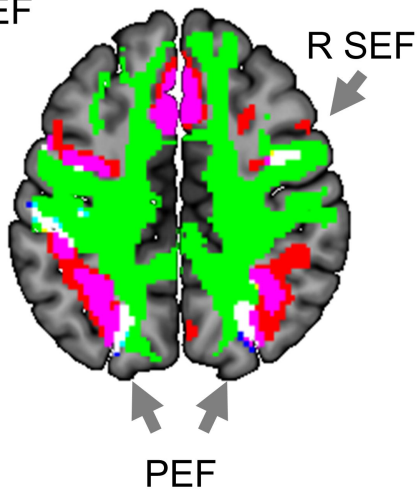
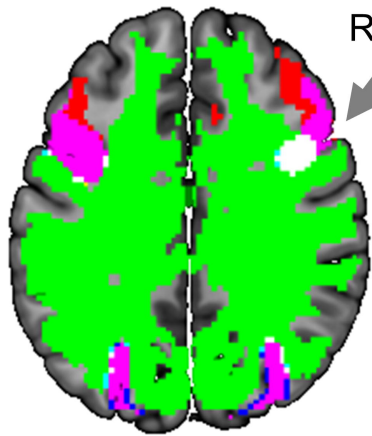
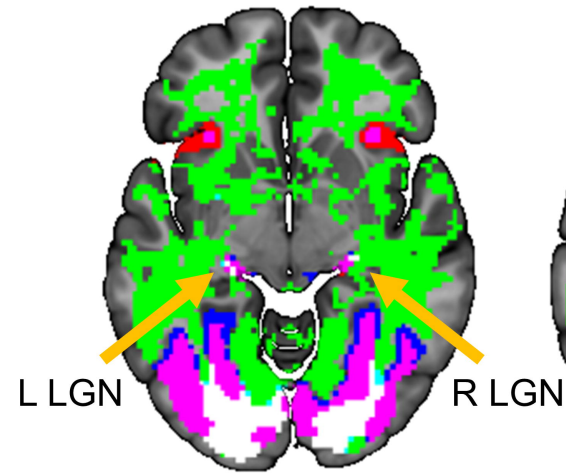
z = -6

z = 10

z = 30

z = 46

y = -24



2-back = 0-back

2-back > baseline



0-back > baseline