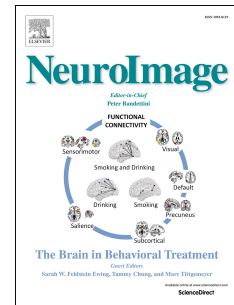


Accepted Manuscript

Probabilistic TFCE: A generalized combination of cluster size and voxel intensity to increase statistical power

Tamás Spisák, Zsófia Spisák, Matthias Zunhammer, Ulrike Bingel, Stephen Smith, Thomas Nichols, Tamás Kincses



PII: S1053-8119(18)31950-5

DOI: [10.1016/j.neuroimage.2018.09.078](https://doi.org/10.1016/j.neuroimage.2018.09.078)

Reference: YNIMG 15314

To appear in: *NeuroImage*

Received Date: 11 June 2018

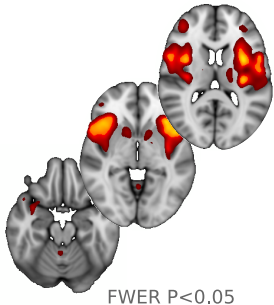
Revised Date: 30 August 2018

Accepted Date: 26 September 2018

Please cite this article as: Spisák, Tamás., Spisák, Zsó., Zunhammer, M., Bingel, U., Smith, S., Nichols, T., Kincses, Tamás., Probabilistic TFCE: A generalized combination of cluster size and voxel intensity to increase statistical power, *NeuroImage* (2018), doi: <https://doi.org/10.1016/j.neuroimage.2018.09.078>.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

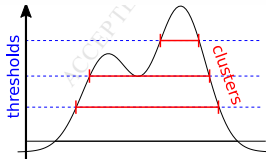
original image



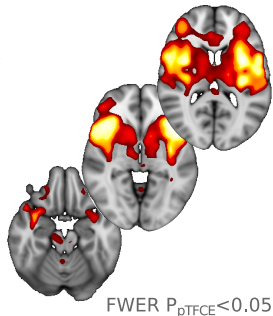
pTFCE



$$p(h|c) = \frac{p(c, |h)p(h)}{\int_{-\infty}^{\infty} p(c, |\xi)p(\xi) d\xi}$$



enhanced image



Probabilistic TFCE: a generalised combination of cluster size and voxel intensity to increase statistical power

Tamás Spisák^{1*}, Zsófia Spisák, Matthias Zunhammer¹, Ulrike Bingel¹, Stephen Smith², Thomas Nichols^{2,3}, Tamás Kincses⁴

1. Department of Neurology, University Hospital Essen, Essen, Germany
2. Wellcome Centre For Integrative Neuroimaging (FMRIB), University of Oxford, Oxford, United Kingdom
3. Department of Statistics, University of Warwick, Coventry, United Kingdom
4. Department of Neurology, University of Szeged, Szeged, Hungary

* Corresponding author: Tamás Spisák

e-mail: tamas.spisak@uk-essen.de

Address:

Department of Neurology
University Hospital Essen
University Duisburg-Essen
Hufelandstrasse 55
45147 Essen, Germany

Abstract

The threshold-free cluster enhancement (TFCE) approach integrates cluster information into voxel-wise statistical inference to enhance detectability of neuroimaging signal. Despite the significantly increased sensitivity, the application of TFCE is limited by several factors: (i) generalisation to data structures, like brain network connectivity data is not trivial, (ii) TFCE values are in an arbitrary unit, therefore, P-values can only be obtained by a computationally demanding permutation-test.

Here, we introduce a probabilistic approach for TFCE (pTFCE), that gives a simple general framework for topology-based belief boosting.

The core of pTFCE is a conditional probability, calculated based on Bayes' rule, from the probability of voxel intensity and the threshold-wise likelihood function of the measured cluster size. In this paper, we provide an estimation of these distributions based on Gaussian Random Field theory. The conditional probabilities are then aggregated across cluster-forming thresholds by a novel incremental aggregation method. pTFCE is validated on simulated and real fMRI data.

The results suggest that pTFCE is more robust to various ground truth shapes and provides a stricter control over cluster "leaking" than TFCE and, in many realistic cases, further improves its sensitivity.

Correction for multiple comparisons can be trivially performed on the enhanced P-values, without the need for permutation testing, thus pTFCE is well-suitable for the improvement of statistical inference in any neuroimaging workflow.

Implementation of pTFCE is available at <https://spisakt.github.io/pTFCE>.

Keywords

neuroimaging, statistics; inference; probabilistic; threshold free cluster enhancement;

Abbreviations

- TFCE threshold-free cluster enhancement
- FWER family-wise error rate
- PDF probability density function
- GRF Gaussian random field
- TPR, FPR true positive rate, false positive rate
- ROC receiver-operator characteristic
- AFROC alternative free-response ROC
- AUC area under the curve
- FWHM full width at half maximum

1 **Introduction**

2 Voxel-wise univariate statistical inference on neuroimaging signals is problematic due to the large
3 number of simultaneously performed statistical comparisons and the unknown, complex dependency
4 between tests. While correcting inferences for a brain-wide search is essential (Bennett et al., 2011;
5 Vul et al., 2009), the attempt to diminish Type I errors rendered most of the statistical thresholding
6 approaches overly conservative (Nichols and Hayasaka, 2003). This might result in increased Type II
7 errors (i.e. missing true effects)(Lohmann et al., 2017) and publication biases (Lieberman and
8 Cunningham, 2009) e.g. toward studying large rather than small effects.

9 Since the signal of interest is usually spatially more or less extended (that is, clustered), sensitivity
10 can be significantly boosted by relating the size of the activation cluster to the (empirical or
11 theoretical) distribution of random clusters in an image of given smoothness (Forman et al., 1995;
12 Friston et al., 1994b). This approach is called *cluster-wise inference*. It captures the spatial nature of
13 the signals, and thus suffers from less multiplicity than voxel-wise inference. However, it is not
14 always more sensitive, and its power depends on the spatial scale of the signal relative to the noise
15 smoothness. For instance focal, intense signals will be better detected by voxel-wise inference
16 (Nichols, 2012). Another serious pitfall when using cluster-level inference is to over-simplistically
17 interpret it at the voxel-wise level: spatial specificity in this case is low and the number of false
18 positive voxels is increased, especially when the significant clusters are large or the applied cluster-
19 forming thresholds are low (Woo et al., 2014). In this paper, we will refer to this phenomenon as
20 “cluster leaking” and discuss it in detail. Moreover, the dependence of results on the (initial cluster-
21 forming) hard-threshold is, on its own, problematic, since finding an optimal threshold is not trivial
22 and might even lead to another multiple comparison problem, if simultaneously testing many
23 different thresholds.

24 An improved way of making use of spatial neighbourhood information in order to boost belief in
25 extended areas of neuroimaging signals is the *threshold-free cluster enhancement* (TFCE) approach
26 (Smith and Nichols, 2009). In contrast to simple cluster-based inference, TFCE does not need a pre-
27 defined cluster-forming hard-threshold. Instead, it calculates the product of some powers of the
28 spatial extent of clusters and the corresponding image height threshold and aggregates this quantity
29 across multiple thresholds. The exponents of these powers are free parameters, but in practice they
30 are fixed to values justified by theory and empirical results.

31 The use of the TFCE method is limited by two factors. The first is that TFCE transforms statistical
32 images into an arbitrary value domain, which is then subject of permutation-testing to obtain P-
33 values. Therefore, it cannot be applied with parametric statistical approaches and a computationally
34 intensive statistical resampling step is always needed. Second, although fixing the free parameters of
35 TFCE provide robust results for three-dimensional images (Salimi-Khorshidi et al., 2011; Smith and
36 Nichols, 2009), generalization of the method for other data structures (e.g. brain connectivity
37 networks) is not trivial (Vinokur et al., 2015) and using suboptimal parameters might be even
38 “statistically dangerous”, that is, might result in an elevated number of false positives and false
39 negatives or non-nominal FWER-rates (Smith and Nichols, 2009).

40 Here, we introduce a method, which is similar to TFCE in its basic concept, but overcomes some of
41 these limitations by giving a *general, extendable probabilistic framework* for *integrating cluster- and*
42 *voxel-wise inference*. The introduced framework allows for converting P-values directly to enhanced
43 P-values and, in several cases, significantly improves the accuracy and the robustness of topology-
44 based belief boosting. The generalisability of the introduced method lies in the freedom of choice in
45 defining what a cluster is in various data structures. In the present study we apply the cluster
46 concept of Gaussian Random Filed theory which reduces the generalisation property of our core

47 formulation to three-dimensional images, but gives a fast implementation of the approach with clear
48 links to the current analysis practice in neuroimaging.

49 **Theory**

50 As many datasets with multiplicity, neuroimaging data (given its inherent spatial autocorrelation), is
51 massively multicollinear. This property of the data does not allow for fully exploiting the localising
52 performance of the typically-used mass-univariate statistical inference which handles the variables
53 largely as independent observations. In such datasets, making use of multicollinearity information in
54 order to boost belief in clusters of correlated signals during statistical inference can retrieve part of
55 the sensitivity that has been lost due to the massive multiplicity. In the case of neuroimages,
56 incorporating information about the spatial clustering properties of the images can be considered as
57 optimising the “localisation performance – sensitivity” trade-off by throwing out only that “part” of
58 localisation capacity which was “unutilised” due to image smoothness.

59 Such an approach, in practice, can be considered as a signal-enhancement, based on the integration
60 of the original voxel-level P-value (resulting from the mass-univariate voxel-wise analysis) with the
61 probability of the cluster it is part of. The resulting single “enhanced” P-value should exhibit the
62 following properties:

63 **Enhancement property: Enhancing sensitivity if data is spatially structured (clustered):** the original
64 P-value is enhanced so that it incorporates the information about the spatial topology of the
65 environment of the voxel. Practically speaking, the method enhances the P-value, if it is part of a
66 cluster-like structure (large enough to be unlikely to emerge when the null hypothesis is true).

67 **Control property: Controlling for false positives and multiplicity:** the enhancement of P-values does
68 not result in an undesirable accumulation of false positive voxels (e.g. due to “cluster leaking”), so
69 that the use of various statistical thresholding and *multiplicity correction* methods (like family-wise
70 error rate, FWER) remain approximately valid at the voxel-wise level.

71 In this section, we introduce a mathematical formulation of a novel candidate for such an
72 enhancement approach. Our basic concept is similar to that known as *threshold free cluster*
73 *enhancement* or, in short, *TFCE* (Smith and Nichols, 2009). Both methods are based on a threshold-
74 wise aggregation of a quantity, which implements a combination of spatial neighbourhood
75 information and intensity in the image at a given threshold. In the next section, these two methods,
76 together with the case of no enhancement, are evaluated and compared in terms of sensitivity on
77 simulated and real datasets (**Enhancement property**). On the same data, we demonstrate that the
78 methods provide an adequate control over false positive voxels (**Control property**, “cluster leaking”).
79 Moreover, our approach directly outputs P-values (without permutation test), and we show that it is
80 valid to correct these enhanced p-values for FWER based purely on the original unenhanced data
81 (**Control property**, multiplicity). This implies that, with our method, thresholds corrected for the
82 FWER in the original data (e.g. via parametric, GRF-based maximum height thresholding) remain
83 directly applicable on the cluster-enhanced data.

84 1.1 TFCE

85 The widely-used TFCE approach enhances areas of signal that exhibit some spatial contiguity without
 86 relying on hard-threshold-based clustering (Smith and Nichols, 2009). The image is thresholded at h ,
 87 and, for the voxel at position x , the single contiguous cluster containing x is used to define the score
 88 for that height h . The height threshold h is incrementally raised from a minimum value h_0 up to the
 89 height v_x (the signal intensity in voxel x , typically a Z-score), and each voxel's TFCE score is given by
 90 the sum of the scores of all "supporting sections" underneath it. Precisely, the TFCE output at voxel x
 91 is:

92

$$TFCE(x) = \int_{h=h_0}^{v_x} c_x(h)^E h^H dh$$

93

Eq. 1

94 , where h_0 will typically be zero (but see (Habib et al., 2017) for details), $c_x(h)$ is the size of the cluster
 95 containing x at threshold h and E and H are empirically set to 0.5 and 2, respectively. The cluster-
 96 enhanced output image can be turned into P-values (either uncorrected or fully corrected for
 97 multiple comparisons across space) via *permutation testing*. The values of parameters E and H were
 98 chosen so that the method gives good results over a wide range of signal and noise characteristics
 99 and, accordingly, can be pre-fixed in many cases. However, they have to be chosen differently for
 100 (largely two-dimensional) skeletonized data, like in the tract-based spatial statistics (Smith et al.,
 101 2006) approach ($E=1$, $H=2$, with 26-connectivity), and, interestingly, optimal parameter values
 102 *become strongly dependent on effect topology and effect magnitude* in the case of graphs (e.g.
 103 structural or functional brain connectivity data) (Vinokur et al., 2015).

104 1.2 pTFCE

105 Although TFCE is based on raw measures of image "height" and cluster extent, it is straightforward to
 106 hypothesise a close relationship to corresponding cluster occurrence probabilities (that is, the
 107 probability of the cluster extent, given the cluster forming threshold, as also used in cluster-level
 108 inference). In fact, in Appendix C of (Smith and Nichols, 2009) it is clarified that with a specific pair of
 109 exponent parameters ($H=2$ and $E=2/3$) $h^H c_x(h)^E$ is approximately proportional to the $-\log$ P-values of
 110 clusters found with different thresholds. This concept could directly link TFCE to cluster probability
 111 and allow for easy generalization, independent of data dimensionality and topology. However, as
 112 demonstrated by (Woo et al., 2014), the number of false positive voxels within an otherwise
 113 significant cluster is unknown and their proportion largely increases, if the cluster forming threshold
 114 is decreased. This leads to the phenomenon of "leaking" of positive observations into areas of
 115 background (if results are interpreted at the voxel-level). Therefore, one could expect that simply
 116 integrating the cluster occurrence probabilities and then using traditional voxel-wise thresholding
 117 approaches might comprise the effect of "cluster leaking" and, therefore, might easily lead to an
 118 accumulation of false positive voxels. In fact, this was confirmed by our preliminary analysis, where
 119 we computed the (negative logarithm of the) geometric mean of cluster occurrence probabilities
 120 across multiple thresholds: $TFCE_{Pclust}(x) = \int_{h=h_0}^{v_x} -\log P(C > c_x(h)) dh / N_h$, where N_h is the
 121 number of thresholds (See Supplementary Figure 1 for results).

122 To ensure a truly parameter-free, generalized cluster enhancement, and at the same time, reduce
 123 the issue of “cluster leaking”, here we introduce a method, called probabilistic TFCE (pTFCE). Instead
 124 of raw measures or probabilities of clusters, pTFCE aggregates the *conditional probability* of a voxel
 125 having an intensity greater or equal to the applied threshold, *given* the size of the corresponding
 126 cluster. By doing so, pTFCE projects the cluster-level neighbourhood information into the voxel-space
 127 and allows for voxel-level interpretation (which is not possible by conventional cluster-level
 128 inference). For an illustration of the proposed method and its relation to TFCE, see Figure 1.

129 **Core formula of pTFCE**

130 In the following, let us focus on a single voxel x and, if not relevant in the given context, neglect the x
 131 subscript in the notations. Accordingly, let V be the random variable modelling voxel intensity at an
 132 arbitrary voxel, and $p(v)$ denote the corresponding probability density function (PDF). For clarity, let
 133 us use $p(h)$ instead of $p(v)$, if we denote a suprathreshold voxel in a binary image thresholded at h . As
 134 conventionally, let $P(V \geq v)$ and $P(V \geq h)$ denote the probabilities for a given unthresholded or
 135 thresholded voxel value, respectively. Furthermore, let $p(c|h)$ denote the PDF of cluster size c , given
 136 that the image was thresholded at h . Next, let us, for a moment, consider that we were blinded on
 137 the actually applied threshold value h_i and are only informed on the measured cluster size $c_i = c(h_i)$. In
 138 that case, by using $p(c_i|h)$ as a *likelihood function*, the PDF for h , given c_i can be expressed using
 139 *Bayes' theorem*:

$$p(h|c_i) = \frac{p(c_i|h)p(h)}{\int_{-\infty}^{\infty} p(c_i|\xi)p(\xi) d\xi}$$

140

Eq. 2

141 The probability of $V \geq h_i$ at the applied threshold h_i , given the measured cluster size c_i is then simply:

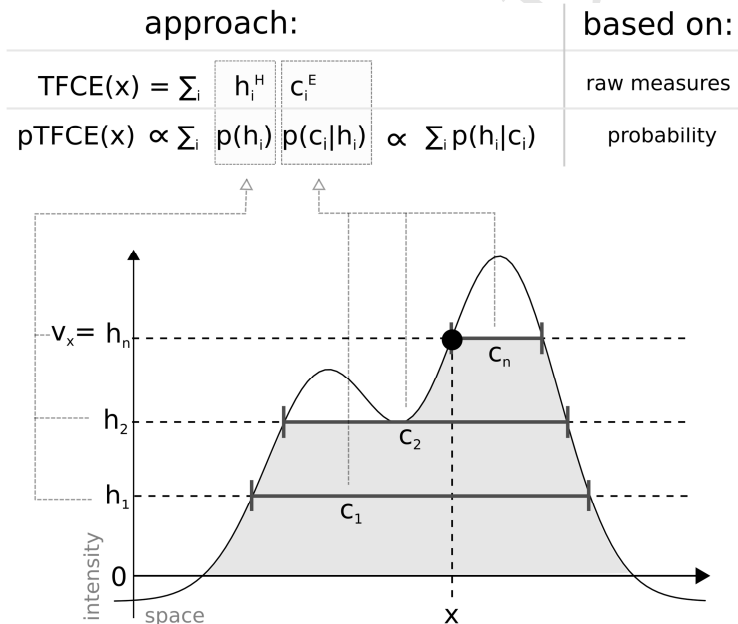


Figure 1. Illustration of the relation between the TFCE and the pTFCE approach. Both approaches are based on the integration of cluster-forming height threshold (h_1, h_2, \dots, h_n) and the supporting section or cluster size (c_1, c_2, \dots, c_n) at that given height. The difference is that, while TFCE combines raw measures of height and cluster size to an arbitrary unit, pTFCE realises the integration by constructing the conditional probability $p(h|c)$ based on Bayes' rule, thereby providing a natural adjustment for various signal topologies. Aggregating this probability across height thresholds provides enhanced P-values directly, without the need of permutation testing.

$$P(V \geq h_i | c_i) = \begin{cases} \int_{h=h_i}^{\infty} p(h|c_i) dh & , \text{if } v \geq h_i \\ 1 & , \text{otherwise} \end{cases}$$

142

Eq. 3

143 , where v is the voxel value (Z-score) in the voxel of interest and the $v < h_i$ branch simply covers the
144 case of *subthreshold* voxels.

145 Here we propose, that a proper aggregation of the $P(V \geq h_i | c_i)$ probabilities over a h_i series of
146 thresholds ($0 \leq h_i < v, i=1, \dots, n$), would satisfy our **Enhancement property and Control property** and
147 provide an appropriate way for integrating cluster information with voxel intensity. An illustration of
148 the “internal workings” of the proposed method can be seen in Figure 2 (parts A and B).

149 **Probability aggregation across thresholds**

150 The introduced conditional probability applies for the value of *the actual cluster-forming threshold*
151 instead of the value of *the actual voxel*. (As also suggested by the notation: $P(V \geq h_i | c_i)$ instead of
152 $P(V \geq v | c_i)$). That means that, when summarizing across thresholds, the aggregated probability should
153 be computed from an “incremental series” of probabilities (each of them corresponding a cluster-
154 forming threshold), rather than a pool of beliefs about the same event. Therefore, to aggregate the
155 $P(V \geq h_i | c_i)$ probability over all h_i thresholds, the common probability pooling methods (Genest and
156 Zidek, 1986; Stone, 1961), e.g., logarithmic pooling, are not suitable.

157 To overcome this, in the following, we give a solution for this scenario, hitherto referred to as
158 **equidistant incremental logarithmic probability aggregation**. Our aim can be formalized as finding
159 an aggregation function $Q(\cdot)$, which exhibits the following properties:

160 (i) $Q(\cdot)$ is interpreted on the *sum* of a series of negative log-probabilities P_i , which are
161 equidistantly distributed in the logarithmic domain, and returns the negative logarithm
162 of the aggregated probability \bar{P} (With the sum of logarithms, we implement a
163 multiplicative model, see Appendix B and C of (Smith and Nichols, 2009)).

$$Q: \sum_i -\log(P_i) \rightarrow -\log(\bar{P}),$$

$$i = 1, 2, \dots, n; \quad \forall i: \log(P_i) - \log(P_{i+1}) = \text{constant}$$

164

165

Eq. 4

166 (ii) for the sum of a series of *unenanced* $-\log$ P-values, $P_1=1$ and $P_n=P(V \geq v)$, $Q(\cdot)$ gives back
167 the negative logarithm of P_n , that is, the original voxel probability:

$$Q\left(\sum_i -\log(P_i)\right) = Q\left(\sum_i -\log(P(V \geq h_i))\right) = -\log(P_n) \cong -\log(P(V \geq v))$$

168

Eq. 5

170 (iii) $Q(\cdot)$ is monotonically increasing, in the sense that it ensures monotonicity about local
171 maxima in the image.

172 Let us note, that our assumption about the constant increment in log probabilities in property (i) is
173 not obligatory for an appropriate aggregation method, but, as detailed below, by allowing a
174 mathematical analogy to the equation of triangular numbers, it simplifies the construction of a
175 proper aggregation method. Furthermore, the uniform sampling of the $\log(P)$ space is a natural way
176 to address a greater accuracy at small p-values, which are typically of interest.

177 In the followings, we introduce a Q function, which fulfils properties i-iii. Moreover, below we
178 present how the value $Q(\cdot)$ returns with the series of *enhanced* $-\log P$ -values as input (instead of the
179 original unenhanced ones, like in property ii) can be considered as the negative logarithm of the
180 *pooled probability of interest*:

$$Q\left(\sum_i -\log(P(V \geq h_i | c_i)_x)\right) := pTFCE(x)$$

181 Eq. 6

182 As a starting point, let us consider the problem of finding the sum of the first n non-negative integer
183 numbers, which is $S_n = \sum_{k=0}^n k = \frac{n(n+1)}{2}$ (giving the so-called triangular numbers). It is easy to
184 generalize this, instead of non-negative integers, to equidistantly distributed non-negative real
185 numbers from 0 to $n\Delta_k$ and by an increment of Δ_k :

$$S_{\Delta_k, n\Delta_k} = \sum_{k=0}^n k\Delta_k = \frac{n\Delta_k(n+1)}{2} = \frac{n\Delta_k\left(\frac{n\Delta_k}{\Delta_k} + 1\right)}{2}$$

186 Eq. 7

187 If we use the notation $w=n\Delta_k$ and denote $S_{\Delta_k, n\Delta_k}$ as simply S , we can write Eq. 7 as:

$$0 = \frac{1}{2\Delta_k} w^2 + \frac{1}{2} w + S$$

188 Eq. 8

189

190 Solving Eq. 8 for w and taking the positive root gives:

$$w_+ = \frac{\sqrt{\Delta_k(8S + \Delta_k)} - \Delta_k}{2} := Q(S)$$

191 Eq. 9

192 In the context of our pTFCE approach, let $S(x)$ denote the sum of the $-\log P(V \geq h_i | c)$ enhanced log-
193 probabilities in voxel position x , so that the h_i thresholds change incrementally in the negative
194 logarithmic domain (to ensure a constant Δ_k):

$$pTFCE(x) = Q(S(x)) = \frac{\sqrt{\Delta_k(8S(x) + \Delta_k)} - \Delta_k}{2}$$

195 , where

$$S(x) = \sum_{i=0}^{N_h} -\log(P(V \geq h_i | c_i)) \quad \text{in voxel position } x$$

196

197 and

$$\forall i, 0 \leq i < N_h : \log(P(V \geq h_{i+1})) - \log(P(V \geq h_i)) = \Delta_k$$

198

Eq. 10

199 The proposed $Q(S(x))$ probability pooling for pTFCE clearly satisfies (i), (ii), (iii) of the above
 200 conditions. Let us note, that the formulation of the introduced probability aggregation method forces
 201 the starting $-\log P$ threshold to be zero ($P_1=1$) and, importantly, the *number of thresholds* and the
 202 *maximal threshold* does *only affect the accuracy* of the approximation as we normalise for Δ_k and
 203 the *enhanced* $-\log P$ -values are zero for subthreshold voxel.

204 The rationale of the proposed *equidistant incremental logarithmic probability aggregation* method
 205 can be understood as searching for the probability corresponding to a hypothetical v' voxel value for
 206 which the threshold-wise *unenanced* sum of negative log probabilities would be the same as the
 207 actual *enhanced* sum corresponding to the observed voxel value v .

208 Another, analogous and even more straightforward way to think about the method is that it links the
 209 series of enhanced probabilities produced by a “strongly clustered” voxel to another hypothetical
 210 series of probabilities produced by a higher intensity voxel, but with “average” clustering (relative to
 211 smoothness) where, accordingly, enhancement has no effect at all; and uses this higher intensity as
 212 the pooled (and enhanced) intensity value.

213 It can be also intuitive to consider a geometric meaning of the method as demonstrated on **Figure 3**.
 214 Computing the pooled log probability $pTFCE(x)$ of voxel x by the Q function is equal to finding the
 215 hypothetical voxel value, for which the sum of the *unenanced* (original) incremental log probability
 216 series ($S_{\Delta_k, v}$ according to Eq. 7, denoted as S'_{orig} and the blue area on the figure) is equal to the sum
 217 of the pTFCE *enhanced* log probability series belonging to the actual (observed) voxel intensity v
 218 (denoted as S'_{orig} and the red area on the figure).

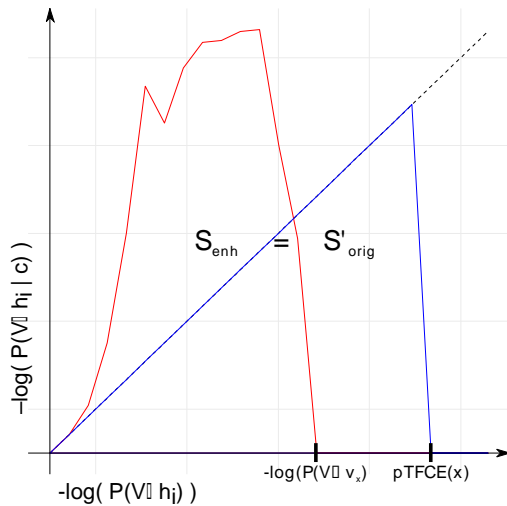


Figure 3. Geometric representation of the proposed equidistant incremental logarithmic probability pooling approach. The proposed *equidistant incremental logarithmic probability aggregation* method links the series of enhanced probabilities produced by a “strongly clustered” voxel (red area, S_{enh}) to another series of probabilities produced by a hypothetical, higher intensity voxel, but with “average” clustering (blue area, S'_{orig}), where enhancement has no effect at all; and uses this higher intensity as the pooled (and enhanced) intensity value.

219 Let us note that, up to this point, our formulation of the pTFCE approach does not contain *any detail*
 220 *about the data structure* and it can be easily *generalized* to various features of the data beside
 221 clustering. In the next section, we link the introduced formulation to the typical volumetric data
 222 structure in neuroimaging by estimating the appropriate PDFs based on Gaussian Random Field
 223 Theory.

224 ***Estimating the likelihood function***

225 We see at least two obvious ways for determining or approximating the probability density functions
 226 $p(h)$ and $p(c|h)$ in order to construct the $P(V \geq h)$, $P(C \geq c|h)$ and $P(V \geq h|c)$ probabilities: for n-
 227 dimensional images, the PDFs in question can be approximated based on *Gaussian Random Field*
 228 (GRF) theory, and in general, empirical estimation of the PDFs can be given by statistical resampling
 229 or permutation test. Since GRF theory (Bardeen et al., 1986; Nosko, 1969) is extensively used in
 230 statistical analysis of medical images (Friston et al., 1996; Nichols, 2012; Worsley et al., 2004), in the
 231 present work, we choose to investigate the use of GRF Theory to give an approximation of the PDFs
 232 and use simulations with a known ground truth to justify the validity of the approximation. Let us
 233 note however, that the use of GRF theory makes our approach specific to 3D images and implicitly
 234 introduces several assumptions not involved in the above discussed general formulation. While the
 235 GRF-based approach is well suited to establish the links of pTFCE to the existing practice in the field
 236 of neuroimaging, more general formulations (like the use of permutation-based empirical
 237 distributions) are subject of further investigation in order to fully exploit the generalisable
 238 formulation of the core pTFCE approach.

239 Determining $p(h)$ for a given threshold h (or $p(v)$ if v is the voxel value) for Gaussianised Z-score
 240 images is straightforward: it is simply the PDF of the normal distribution with zero mean and unit
 241 variance:

$$p(h) = \phi(h)$$

242 and

$$P(V \geq h) = 1 - \Phi(h)$$

243

Eq. 11

244 , where Φ (uppercase) is the cumulative distribution function and ϕ (lowercase) is the probability
 245 density function of the standard Gaussian. Note that, according to Eq. 11, we define the input image
 246 for pTFCE as an image of Z-scores, nevertheless, it is easy to convert any other statistical parametric
 247 images (like T- or P-values) to Z-scores.

248 For inference on clusters, we first utilize GRF theory to find the mean cluster size under the null at a
 249 given threshold, the distribution of cluster size about that mean and also, the P-value for a given
 250 cluster size.

251 Let N be the total supra-threshold volume (equivalently, the sum of all of the cluster sizes); let L be
 252 number of clusters observed. Assuming a large search region relative to the smoothness, and thus
 253 independence of number of clusters and cluster size, the mean cluster size under the null is

254 $E[C] = E[N]/E[L]$. Here, the numerator is easily obtained: $E[N] = V P(V \geq h)$, and specifically for
 255 a Gaussian image: $E[N] = V (1 - \Phi(h))$, where $\Phi(\cdot)$ is the CDF of a Gaussian and V is the total
 256 number of voxels (Friston et al., 1996).

257 The expected number of clusters is $E[L]$ approximated by the expected *Euler Characteristic* (Worsley
 258 et al., 1996). The Euler Characteristic (*EC*) of a D -dimensional random field thresholded at h is written
 259 χ_h , and is the number of clusters minus the number of holes ($D \geq 2$) plus the number of handles ($D \geq 3$).
 260 For sufficiently high h the probability of a hole or handle is small, and so the *EC* offers a good
 261 approximation of the number of clusters. Worsley's general results (Worsley et al., 1996) give a
 262 closed-form expression for $E[\chi_h]$ as a sum of D terms: $E[\chi_h] = \sum_{d=0}^D R_d \rho_d(h)$, where R_d is the *RESEL*
 263 count and ρ_d is the *EC* density. The *RESEL* count is a length, area or volume, depending on d , and it is
 264 the product of the spatial measure and a roughness measure $|\Lambda|^{1/2}$, where Λ is the $d \times d$ variance-
 265 covariance matrix of the partial derivatives of the data. Usually, only the $d = D$ term is appreciable, so
 266 for the $3D$ case we have $E[\chi_h] = R_3(h^2 - 1) e^{-h^2/2} (2\pi)^{-4/2}$. Thus, the expected cluster size for a
 267 $3D$ Gaussian (Z-score) image with cluster-forming threshold h is:

$$E[c_h] = \frac{V(1 - \Phi(h))}{V |\Lambda|^{1/2} (h^2 - 1) e^{-h^2/2} (2\pi)^{-2}} = \frac{V(1 - \Phi(h))}{R_D (h^2 - 1) e^{-h^2/2} (2\pi)^{-2}}$$

268

Eq. 12

269 Finally, using the result that cluster size to the $2/D$ power follows an exponential distribution (Nosko,
 270 1969), with

$$\lambda_h = \left(\frac{E[c]}{\Gamma\left(\frac{D}{2} + 1\right)} \right)^{-\frac{2}{D}}$$

271

Eq. 13

272 , the PDF of the cluster size at threshold h is:

$$p(c|h) = \frac{2\lambda_h e^{-\lambda_h c^{2/3}}}{3c^{1/3}}$$

273 , and

$$P(C > c|h) = c^{2/3} e^{-\lambda c^{2/3}}$$

274

Eq. 14

275 , where λ is the rate of the exponential distribution and $\Gamma(\cdot)$ is the gamma function.

276 There are several implicit assumptions in the Gaussian Random Field approach (Friston et al., 1996),
 277 and, importantly, GRF theory-based estimates become inaccurate or even undefined at low
 278 thresholds. Therefore, for the GRF-based implementation of pTFCE, here we propose that for any
 279 thresholds h smaller than a specific value h_{GRF} , the enhanced probability $P(H>h |c)$ is to be
 280 approximated simply by the unenhanced $P(H>h)$. Moreover, in Eq. 2, $p(c|h)$ should be truncated by
 281 setting it to 0 when $h < h_{GRF}$. Let us note, that this also truncates the resulting $p(h|c)$ distribution on
 282 the left side, thus slightly increases the enhanced probability $P(V \geq h |c)$, meaning that with this
 283 approximation, the cluster enhancement is *expected to be conservative*. That also means that, while
 284 this approximation might mean a loss in sensitivity, we can still expect our pTFCE approach to
 285 perform well in terms of our **Control property** (that is, it remains controllable e.g., for family wise
 286 error). Investigation of the effect of h_{GRF} for using GRF theory revealed that the pTFCE is robust to the
 287 choice of this value (see Supplementary Figure 5). Therefore, we fixed this value at $h_{GRF}=1.3$ (default
 288 parameter in the software implementations, as well), and tested the validity of this approximation in
 289 simulations with known ground truth and on real data.

290 **Evaluation Methods**

291 **1.3 Simulated data**

292 To assess the statistical validity of our methods, simulated data comprising seven 3D test image
 293 shapes (the same as in (Smith and Nichols, 2009)) were used to compare the pure voxel-level (a.k.a.
 294 unenhanced, hereinafter denoted simply as “VOXEL” method), the TFCE and the proposed pTFCE
 295 methods against each other, with ROC evaluations giving objective combined measures of specificity
 296 and sensitivity. Additionally, we tested whether the P-values output by the pTFCE method are valid
 297 for correction for multiple comparisons, by correcting based on the P-value distribution of the
 298 *unenhanced* voxel-based approach (See ROC methodology). This method we denote as pTFCE_{VOX}, to
 299 emphasize that it uses the thresholds computed for the VOXEL method and to distinguish it from the
 300 variant with the randomisation-based threshold (denoted simply as pTFCE).

301 **1.3.1 Test signal shapes**

302 In our simulations, we used the same seven 3D test signal shapes as in (Smith and Nichols, 2009).
 303 These are shown on **Error! Reference source not found.A**. These ground truth images cover a wide
 304 range of signal types, including small blobs, touching blobs and extended areas of activation. Each
 305 test signal has a background value of 0 and a peak height of 1. We then scaled the signal by a factor
 306 or 0.5, 1, 2 or 3 and added unsmoothed Gaussian white noise of standard deviation 1, to give a range
 307 of peak signal-to-noise (SNR) values: 0.5, 1, 2 and 3.

308 We evaluated also the effect of different Gaussian smoothing kernels, with full-width-at-half-
 309 maximum (FWHM) values of 1, 1.5, 2, 3 voxels applied. After smoothing, the data was scaled so as to
 310 keep the noise standard deviation equal to 1, so that the images were still analogous to T/Z images.
 311 For the TFCE method, we used the standard parameter values $E=1/2$, $H=2$.

312 **1.3.2 ROC methodology**

313 An ROC (receiver-operator characteristic) curve, given a signal+noise image and the known ground
 314 truth, plots true positive rate (TPR) against false positive rate (FPR), as one varies a threshold applied
 315 to binarise the image. An ideal algorithm gives perfect true positive rate at zero false positive rate,
 316 i.e., the perfect ROC curve jumps immediately up to TPR=1 (y-axis) for FPR=0 (x-axis) and stays at 1
 317 for all values of FPR. Hence a commonly-used single summary measure of the whole ROC curve is the
 318 AUC (area under curve); the higher the AUC, the better. To ease interpreting our results in relation to
 319 those in (Smith and Nichols, 2009), we use an analogous ROC methodology: we use AUC values for
 320 alternative free-response receiver-operator characteristic (AFROC) (Bunch et al., 1977; Chakraborty
 321 and Winter, 1990).

322 ***Alternative free-response ROC***

323 AFROC analysis plots the proportion of true positive tests (among all possible positive tests) on the y-
 324 axis and the probability of any false positive detections anywhere in the image on the x-axis (that is,
 325 the family-wise error rate, or FWER). Here we calculate FWER for the AFROC curves by counting the
 326 number of images with one or more false positive voxels among 1000 smoothed noise-only images.
 327 As neuroimaging analyses typically seek to control the FWER, we used this method to test the
 328 ***Enhancement property*** of different spatial enhancement/thresholding methods. For AFROC analysis,
 329 we define true positives based on the smoothed ground truth images.

330 Since ROC analysis is predominantly a binary concept, we threshold and binarise the smoothed
 331 ground truth images at $0.1/SNR$. This ensures, that voxels, in which a significant amount of signal was
 332 introduced by smoothing, count as true positive observation, if detected by any of the methods.

333 The above approach of calculating AFROC, by estimating FWER from processed pure-noise data,
 334 avoids the need to determine what is “real” background in the signal+noise data after passing
 335 through a given algorithm (Smith and Nichols, 2009). It is exactly what we want in the standard
 336 scenario of null-hypothesis testing which aims to explicitly control the FPR in the presence of no true
 337 signal; it tests sensitivity when the specificity is being controlled globally (that is among studies and
 338 not over voxels), in the way that we generally require in practice. This method of calculating FWER
 339 ignores the FP voxels in the signal+noise images that are spatially close to the true signal (as distinct
 340 from “real” FP voxels in the noise-only data), and in doing so does not weight, for example, against
 341 the smearing of estimated signal into neighbouring voxels due to smoothing.

342 ***“Negative” alternative free-response ROC***

343 Since FWER for AFROC is calculated from the *noise-only images*, aside from the effect of smoothing,
 344 AFROC is insensitive to “cluster leaking”, as well. This might be an undesirable property since
 345 “leaking” will possibly merge small blobs of activity with neighbouring random local maxima and
 346 present them as large clusters with lots of false positive voxels (Woo et al., 2014). Therefore, while
 347 AFROC is suitable for measuring signal detectability, to control for the voxel-level spatial specificity of
 348 the methods, here we introduce the “negative AFROC” method.

349 In the negative AFROC (short: nAFROC) analysis, we test our ***Control property*** for “cluster leaking”, by
 350 thresholding the smoothed ground truth images also at $0.001/SNR$ and then binarising and inverting
 351 it to define a region where the amount of signal can be neglected. Applying our AFROC method with
 352 this region as ground truth, we create ROC curves for the FPR, plotted against the FWER. An
 353 appropriate cluster enhancement algorithm should keep the area under the negative AFROC curve
 354 very close to zero, to minimise voxel-level Type I errors.

355 ***Correcting for multiple comparisons on enhanced probability values***

356 As the introduced pTFCE enhancement method *directly* outputs probability values, we also tested
 357 whether these enhanced probability values can indeed be interpreted in the range of the original,

358 unenhanced probability values (**Control property**: multiplicity) that is, thresholds corrected for
 359 multiple comparisons on the original, unenhanced data are directly applicable to the pTFCE
 360 enhanced data, as well. We did this by using the *unenhanced* noise-only images to calculate the
 361 FWER thresholds for the AFROC and negative AFROC curves. The rationale behind this analysis is that
 362 if thresholds computed from, and interpretable on *unenhanced* data are valid on pTFCE-enhanced
 363 probability values, as well (that is, they always guarantee a smaller or equal FWER), then a freedom
 364 in the choice of the correction method is granted, as long as it is valid on the original statistical
 365 image. In practice this would mean that no special correction methods for enhanced images (like
 366 permutation test-based maximum height correction) are needed and any threshold will give an equal
 367 (or close) FWER for the enhanced image than for the unenhanced.

368 We denote this analysis case, as $pTFCE_{vox}$, the subscript standing for **VOXEL** and referring to the way
 369 corrections for multiple comparisons were based on the P-values of the (unenhanced) voxel-level.

370 **1.3.3 Simulation details**

371 To summarize, our simulation and ROC analysis consisted of the following steps:

- 372 **1.** Generate the raw ground truth 3D image (“signal”).
- 373 **2.** Generate 1000 random Gaussian noise images (mean zero, unit variance) that, when added to the
 374 signal, give a specified SNR in the resulting 1000 signal+noise images. We applied exactly the same
 375 noise realizations as in (Smith and Nichols, 2009) and generated signal+noise images for SNRs 0.5, 1,
 376 2 and 3.
- 377 **3.** Pass all noise-only and signal+noise images through a smoothing stage (with FWHMs 1, 1.5, 2 and
 378 3) and afterwards, through the algorithm being tested: VOXEL (=no enhancement), TFCE, pTFCE,
 379 $pTFCE_{vox}$ (=pTFCE at this stage).
- 380 **4.** Pass also the ground-truth image through the smoothing stage, normalise its intensity to $\max=1$,
 381 and threshold it at $0.1/SNR$ and binarise, to define a ROI of true positive observations. The rationale
 382 behind smoothing the ground truth image is, that “de-smoothing” should not be the responsibility of
 383 any statistical inference method. Similarly, when smoothing is applied in a real data scenario, we can
 384 only expect to detect the smoothed version of the underlying activation pattern (which is also a
 385 strong argument against using excessive smoothing). Accordingly, areas of false positive observations
 386 are also defined based on the smoothed ground truth image, by thresholding it at $0.001/SNR$, and
 387 then binarising and inverting it.
- 388 **5.** Compute the traditional ROC curves for the processed signal+noise images and for the ground
 389 truth mask thresholded at $0.1/SNR$.
- 390 **6.** Compute AFROC and negative AFROC curves: threshold the appropriately processed noise-only
 391 and signal+noise images for the methods voxel, TFCE and pTFCE and $pTFCE_{vox}$ at the full range of
 392 possible threshold values. FWER, TPR and FPR at each threshold are computed as follows:
 - 393 • FWER: For each threshold level, count the number of processed noise-only images which contain
 394 any supra-threshold voxels. This count (divided by 1000) gives the family-wise FPR for this threshold
 395 level (i.e., achieves full correction for multiple comparisons across space). For $pTFCE_{vox}$, FWER is
 396 computed based on the *unenhanced* noise-only images (same as for the “voxel” method).
 - 397 • TPR: For each threshold level, use each of the 1000 processed signal+noise images, along with the
 398 ground truth mask thresholded at $0.1/SNR$, to obtain an estimate of the TPR. We use the raw voxel-
 399 wise TPR (fraction of non-background signal voxels correctly reported), averaged over the 1000

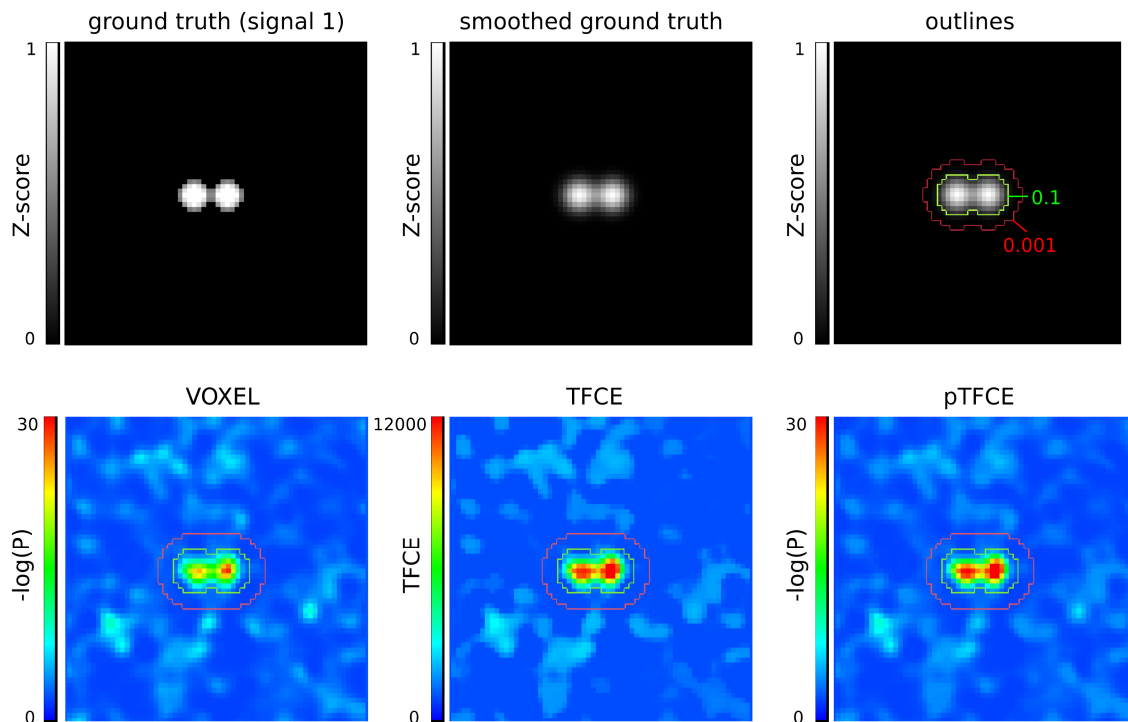


Figure 4. Representative images of the simulation and ROC ROIs for the investigated approaches. In the upper row, the unsmoothed and the smoothed (FWHM=1.5) versions of signal 2, and the outlines corresponding to the ROC (inside vs. outside the 0.1 contour), AFROC (inside the 0.1 contour) and negative AFROC (outside the 0.001 contour) analysis. In the lower row, results of the VOXEL, TFCE and pTFCE methods with the same noise realisation (noise 0001) are visualized with the contours as overlay.

400 signal+noise images (we also record the IQR of the TPR across the 1000 images, as a measure of the
401 stability of the various algorithms being tested).

402 • FPR: For each threshold level, use each of the 1000 processed signal+noise images, and use the
403 “negative” ground truth mask (thresholded at 0.01/SNR, binarised and inverted) to obtain an
404 estimate of the FPR. The raw voxel-wise FPR is then averaged across the 1000 images (again, we
405 record the IQR, as well).

406 6. Take the resulting ROC, AFROC and negative AFROC curves, and, using only the x-range of 0 to
407 0.05, calculate the AUC values (AUC ROC, AUC AFROC, AUC negative AFROC, respectively). Normalise
408 AUC by 0.05.

409 For an example of ground truth smoothing and the ROIs used for the ROC methodology along with
410 demonstrative results of the tested approaches, see Figure 4.

411 For estimating smoothness of the signal+noise images, for generating signal+noise images and for
412 AFROC analysis, we used FSL (Jenkinson et al., 2012; Smith et al., 2004). (We performed a
413 “traditional” ROC analysis, as well, see Supplementary Figure 2 and Supplementary Table 1) The
414 pTFCE algorithm was implemented in R, using the packages “mmad” (Clayden, 2014) (used for a fast
415 labelling of connected components in thresholded images) and “oro.nifti” (Whitcher et al., 2011) (to
416 manipulate nifty images). Calculation of pTFCE was performed with a fixed number of $\log(P)$ -
417 thresholds ($n=100$), ranging from 0 to the negative logarithm global image maximum $-\log(\max_x v_x)$,
418 and distributed equidistantly. Although this results in different deltas for various images, theory
419 suggests that this affects only the accuracy of the probability aggregation and our supplementary
420 analysis (Supplementary Figure 3) confirmed that the proposed *equidistant incremental logarithmic*
421 *probability aggregation* method is robust above a reasonable number ($n \geq 100$) of thresholds and
422 magnitude of delta values. Overflow problems were handled in most of the cases in the same way as

423 in the FSL source code (Jenkinson et al., 2012; Smith et al., 2004). Results were plotted with the
424 oro.nifti and the ggplot2 (Wickham, 2016) packages of R. We compare the AUC values at identical
425 parameter settings. Moreover, since it is a common practice to optimise neuroimaging pipelines in
426 terms of smoothing to achieve maximal sensitivity, and optimization often implicitly takes into
427 account the typical signal-to-noise level of the experimental design (through the inherent properties
428 of the data used to optimise the workflow), for each method, each test signal and each SNR, we
429 chose an optimal smoothing based on the best AUC values of the AFROC curves and compared
430 methods with their optimal settings.

431 **1.4 Testing on real data**

432 Since real neuroimaging data differs in many properties from the simplistic Gaussian model used in
433 the simulation, we use various fMRI datasets for the purposes of (1) evaluating the improvement in
434 sensitivity by investigating the dependence of results on sample size, (2) investigating whether pTFCE
435 maintains nominal family-wise error rates when corrected for multiple comparison, given that the
436 null hypothesis is true; and (3) illustrate the effect of pTFCE with enhancing the activation map
437 reflecting the pain matrix which is well known and complex enough to capture the advantages of our
438 method. We took advantage of the easy correction of pTFCE P-values for multiple comparisons and
439 did not apply the permutation-based FWER thresholds for pTFCE; only the pTFCE_{vox} method
440 (maximum-height thresholding based on GRF-theory) was tested in the real-data scenarios.

441 **1.4.1 Demonstration of the increased statistical power on real data**

442 For both the demonstration of increased statistical power and the evaluation of family-wise error
443 rates, we obtained data from the UCLA Consortium for Neuropsychiatric Phenomics LA5c Study (CNP)
444 (Gorgolewski et al., 2017; Poldrack et al., 2016) as shared via the OpenNeuro database (accession
445 number: ds000030, <https://openfmri.org/dataset/ds000030/>). Processed 1st-level activation maps
446 (contrast of parameter estimates, a.k.a “cope” images, as provided with the dataset) of N=119
447 healthy participants from the “switch-noswitch” contrast (cope39) of the task switching paradigm
448 were obtained and fed into FSL “randomise” (number of permutations: 5000) to create a group-
449 mean activation map. This activation map was then thresholded at a voxel-level (unenhanced) FWER-
450 corrected p<0.05 threshold and considered as “ground truth” for further analysis. As a next step, a
451 total of 900 activation maps were computed from random subgroups of the healthy population with
452 sample sizes N=5, 10, 20, 30, 40, 50, 60, 80 and 100 and with 100 random sampling per sample size.
453 Corrected (FWER, p<0.05) voxel-level and TFCE images and T-score maps were obtained. The latter
454 ones were converted to Z-score maps, fed into the pTFCE algorithm and thresholded based on GRF
455 theory, with a corrected threshold of p<0.05 (implementing the pTFCE_{vox} approach). True positive
456 rate was defined as the proportion of “ground truth” voxels found by the thresholded activation
457 maps of the subsamples. Mean, 0.025% and 0.975% percentiles of this true positive rate were
458 obtained for each of the investigated methods (VOXEL, TFCE and pTFCE_{vox}) and plotted against the
459 sample sizes.

460 **1.4.2 Evaluation of family-wise error rates on real data**

461 Evaluation of family-wise error rates in the case of a true null hypothesis was performed on the same
462 dataset as the demonstration of statistical power (CNP study, OpenNeuro database, ds000030). The
463 true null hypothesis (no group difference) was approximated by comparing 1st-level activation maps
464 (contrast of parameter estimates, a.k.a “cope” images, as provided with the dataset) of two random
465 samples (N=20 per group) from the group of healthy participants from the “switch-noswitch”
466 contrast (cope39) of the task switching paradigm. One thousand of these random control-to-control
467 comparisons was performed with FSL “randomise” (number of permutations: 5000). Corrected
468 (FWER, p<0.05) p-value images of the voxel-level and TFCE methods and simple voxel-wise T-score
469 maps were obtained. The latter ones were converted to Z-score maps, fed into the pTFCE algorithm

470 and thresholded based on GRF theory, with a corrected threshold of $p > 0.05$ (implementing the
471 pTFCE_{vox} approach). Family-wise error rate for each method was than estimated by the ratio of
472 thresholded images with any non-zero value across the 1000 random comparisons, for each method.

473

474 **1.4.3 Illustration of the cluster enhancement effect**

475 For the illustration of the cluster enhancement effect, we used data published as part of (Bingel et
476 al., 2011). This study investigated how divergent expectancies alter the analgesic efficacy of a potent
477 opioid in healthy volunteers by using fMRI. We aimed to evaluate enhancement approaches on the
478 published group-mean map of brain activation to painful stimulation, activating the pain matrix
479 (Figure 3. of the paper).

480 For TFCE, we fed the appropriate spatially standardized subject-level SPM contrast images into
481 randomise and estimated TFCE-enhanced, FWER-corrected p-values for the group mean activation
482 (number of permutations: $N=10000$). For pTFCE, we simply obtained the corresponding second-level
483 spmT image (see (Bingel et al., 2011) for details) and converted it to Z-score, based on the degrees of
484 freedom ($dof=63$). We estimated the smoothness of the map with the “smoothest” command-line
485 tool of FSL. The Z-score map and the smoothness estimate was then fed into the pTFCE algorithm
486 which computed the enhanced map (outputting negative log P-values). The working resolution of the
487 data in standard space was $2 \times 2 \times 2$ mm for both TFCE and pTFCE. For visualization, the original and
488 the enhanced negative log P-value maps were thresholded at 13.61, the $-\log(P)$ threshold
489 determined by the FWER correction ($P < 0.05$) in the original study. The FWER correction for TFCE was
490 based on the permutation test.

491 **Results**

492 **1.5 Implementation details**

493 The pTFCE algorithm is available as an R package called “pTFCE”, and also as an SPM (Friston et al.,
494 1994a) Toolbox with the same name. These packages, together with a fast C++ and FSL-based
495 implementation and nipy interfaces (Gorgolewski et al., 2011) are also available at
496 <https://spisakt.github.io/pTFCE>.

497 **1.6 Simulation results**

498 The full AFROC curves (See **Error! Reference source not found.C** for some representative curves)
499 reveal that our AUC results are not dependent on the applied x-axis threshold (FWER $P < 0.05$), since
500 all curves have a smooth rise in the 0-0.05 interval and differences between approaches are constant
501 at nearly all FWER thresholds. Therefore, AUC values provide valid description of these curves in the
502 following sections.

503 **1.6.1 Comparing enhancement methods with optimal parameter settings**

504 It is common practice to optimise neuroimaging pipelines in terms of smoothing to achieve maximal
505 sensitivity. Moreover, optimization often implicitly considers the typical signal-to-noise level of the
506 experimental design. Therefore, besides comparing all tested methods with identical parameter
507 settings (section 1.6.2), for each method, each test signal and each SNR, we chose an optimal
508 smoothing, based on the best AUC values of the AFROC curves. Results for SNR=1 and 2 are plotted
509 in Figure 6.

510 The mean(\pm sd) optimal smoothing FWHM across all test signal shapes and SNR values was
511 2.29(\pm 0.65), 1.86(\pm 0.84), 2.02(\pm 0.78) and 1.96(\pm 0.76) voxels for VOXEL, TFCE, pTFCE and pTFCE_{vox},
512 respectively. Although pooling across test signal shapes obviously does not necessary provide
513 summary statistics that are representative for real experimental settings, these results still strongly
514 suggest that the cluster enhancement methods generally require less smoothing. Not surprisingly, in
515 the case of spatially extended test signal shapes (signals 3, 6 and 7) a greater smoothing was
516 preferred by all methods. In the case of the other, spatially more restricted signal shapes (signals 1,
517 2, 4 and 5) an optimal smoothing of 1 or 1.5 was found. TFCE preferred in several instances an even
518 smaller smoothing than pTFCE and pTFCE_{vox}.

519 In general, when pooled over all ground truth images and all SNRs, pooled mean(\pm sd) for the mean
520 AUC of AFROC curves with optimal smoothing were 0.102(\pm 0.16), 0.141(\pm 0.2), 0.142(\pm 0.2) and
521 0.134(\pm 0.2) for VOXEL, TFCE, pTFCE and pTFCE_{vox}, respectively, clearly underpinning the
522 improvement in sensitivity in the case of cluster enhancement methods. Moreover, with optimal
523 smoothing, pTFCE and pTFCE_{vox} outperformed the *unenanced* inference for all SNR values and
524 ground truth shapes. TFCE did not manage to improve sensitivity in case of test signals 1, 2 and 5
525 with large SNR values. Summarising the inter-quartile ranges suggests that pTFCE and pTFCE_{vox} (mean
526 IQRs 0.018, 0.02) might be somewhat more robust than TFCE (0.022), but the unenhanced inference
527 displayed the strongest stability (0.013).

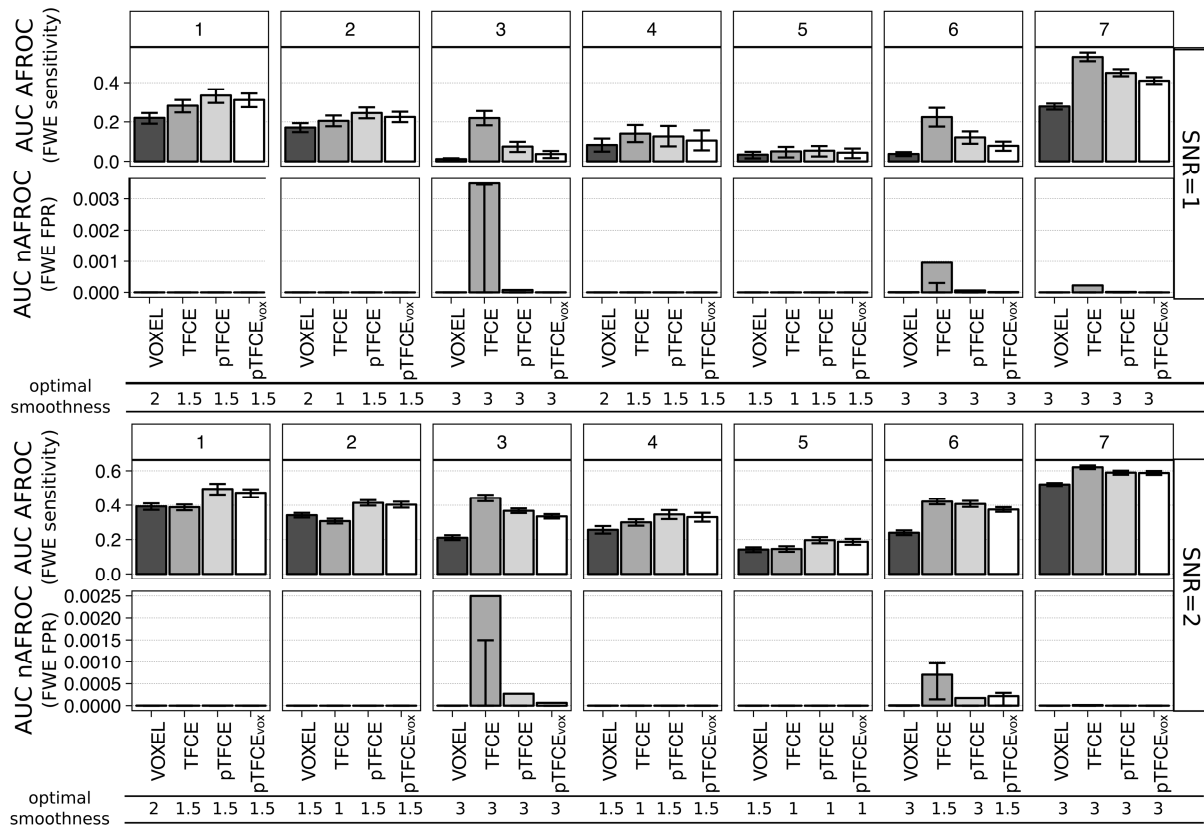


Figure 5. AUC values of the AFROC (representing FWER-level sensitivity) and negative AFROC (representing 1-FWER-level specificity) curves for all methods, ground truth shapes, SNRs, with the optimal smoothing. Barplots represent the mean across 1000 simulation runs, and whiskers represent the inter-quartile range.

528 When comparing the FWER-controlled sensitivity of cluster enhancement methods against each
 529 other, we found that pTFCE and pTFCE_{vox} *outperforms* TFCE in almost all cases with spatially more
 530 restricted test signals (signal 1, 2, 4 and 5). This difference is present already at SNR=1, but becomes
 531 more expressed with greater SNR. In contrast to pTFCE and pTFCE_{vox}, TFCE preferred ground truth
 532 images with spatially extended signal. With these ground truth images, it produced remarkably
 533 *greater* AUC values in AFROC analysis than pTFCE. However, we observed an *increased number* of
 534 probably “cluster leaking”-related *false positives* in these cases. Notably, a modest increase of false
 535 positives was experienced also in case of pTFCE and pTFCE_{vox}, but AUC values of nAFROC curves were
 536 significantly lower than for TFCE (second row in each panel of Figure 6). For the aforementioned
 537 spatially restricted test signals, all cluster enhancement methods give very strict control over cluster-
 538 leaking and corresponding false positives, as shown by the nAFROC curves and corresponding AUC
 539 values.

540 The AUC values of pTFCE_{vox} are systematically slightly below the mean AUC values of pTFCE, but,
 541 relative to the other methods, are not very much lower. This suggests that correcting GRF-based
 542 enhanced pTFCE P-values for multiple comparisons gives valid results even if it is based on the
 543 distribution of the unenhanced P-values (“VOXEL” noise images instead of “pTFCE noise images”).

544 **1.6.2 Comparing enhancement methods with identical parameter settings**

545 Another possible concept for contrasting the tested approaches is to investigate their performance
 546 with identical parameter settings, instead of using the optimized smoothing values for each. Separate
 547 comparison of the tested approaches with each parameter-setting showed, in general, very similar
 548 results to those revealed with optimized smoothing.

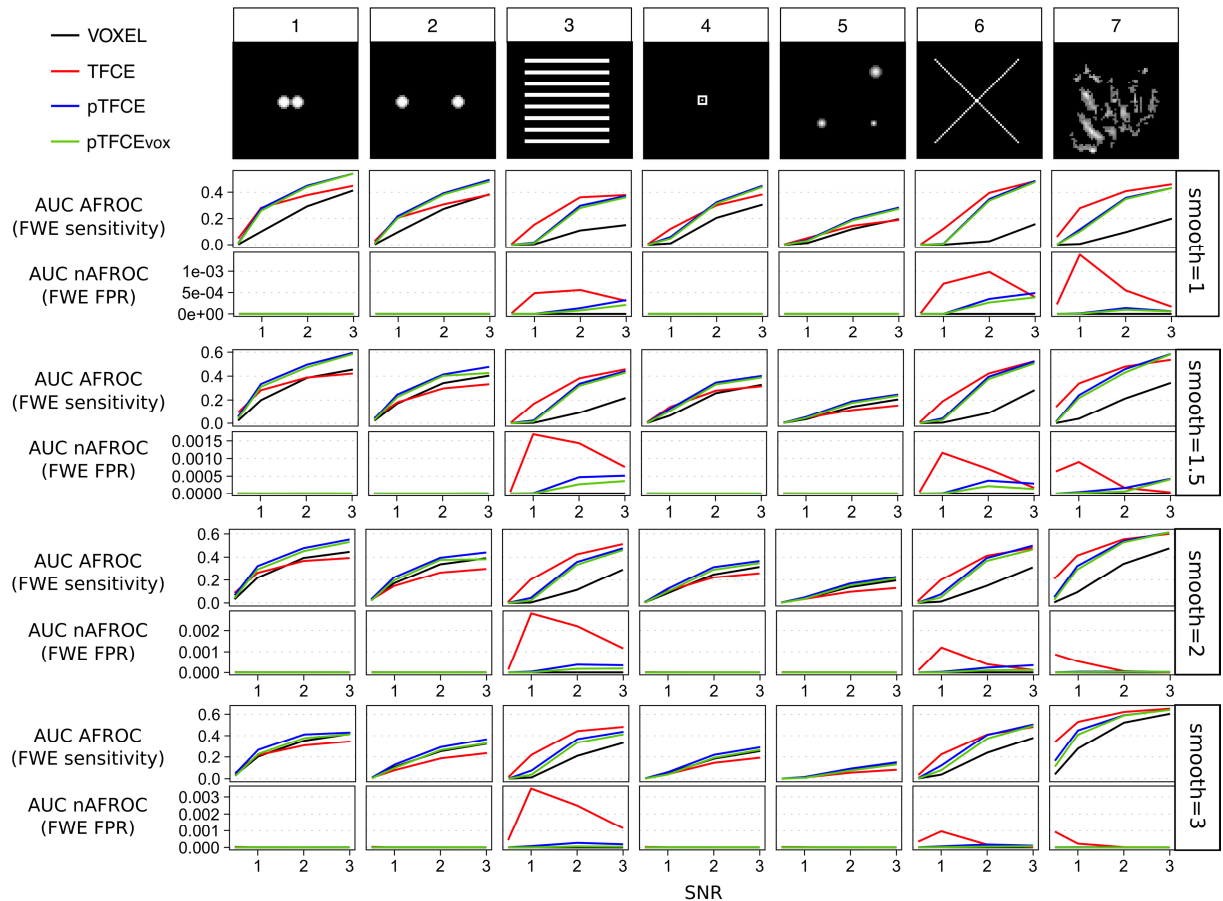


Figure 6. Mean AUC values of the AFROC (representing FWER-level sensitivity) and negative AFROC (representing 1-FWER-level specificity) curves for all methods, ground truth shapes, SNRs and smoothing kernels.

549 In these comparisons, as well, cluster enhancement methods tended to be superior to the
 550 unenhanced voxel-level inference (Figure 7, AUC AFROC: upper row for each smoothing value),
 551 although again, in a few cases (typically for high SNRs with test signals 1, 2 and 5) TFCE yielded lower
 552 mean AUC values than *unenhanced* voxel-level thresholding with the same simulation parameters.
 553 On the other hand, pTFCE *robustly outperformed the unenhanced thresholding* in *all* of the
 554 simulation parameter settings. The same is true for pTFCE_{vox} suggesting that thresholding the pTFCE
 555 enhanced image with thresholds originally interpreted on the unenhanced image also results in an
 556 enhanced sensitivity.

557 The mean area under the negative AFROC curve is, for most of the parameter settings, relatively
 558 close to zero, meaning an *appropriate control for false positives* and “cluster leaking”. However, in
 559 case of the cluster enhancement methods (TFCE, pTFCE and pTFCE_{vox}), false positives (defined by our
 560 nAFROC analysis, very conservatively, as regions where the noise is at least 1000 times greater than
 561 signal) slightly accumulate by ground truth images with extended area of signal (signals 3, 6 and 9,
 562 upper row of each smoothing parameter panel on Figure 7).

563 When comparing cluster enhancement methods against each other, we found a very similar pattern
 564 to that with optimized smoothing: pTFCE and pTFCE_{vox} tended to slightly outperform TFCE in terms of
 565 FWER-controlled sensitivity in almost all cases with spatially more restricted test signals (signals 1, 2,
 566 4 and 5). In these cases, all cluster enhancement methods give a very strict control over cluster-
 567 leaking and corresponding false positives. On the contrary, TFCE tends to outperform pTFCE and
 568 pTFCE_{vox} in simulations with spatially extended ground truth (signals 3, 6 and 7) and low SNRs;
 569 however, it does it at the price of decreased specificity, as suggested by the mean AUC of negative
 570 AFROC curves in these cases (second rows on each smoothing panel on Figure 7). All in all, in most of
 571

571 the cases pTFCE and pTFCE_{vox} seems to provide a similar or stricter control over cluster leaking than
 572 TFCE.

573 Importantly, low AUC values for the negative AFROC curves of the pTFCE_{vox} method imply that
 574 thresholding the pTFCE enhanced image with thresholds originally interpreted on the unenhanced
 575 image gives an appropriate control for FWER.

576 1.7 Real data evaluation

577

578

579

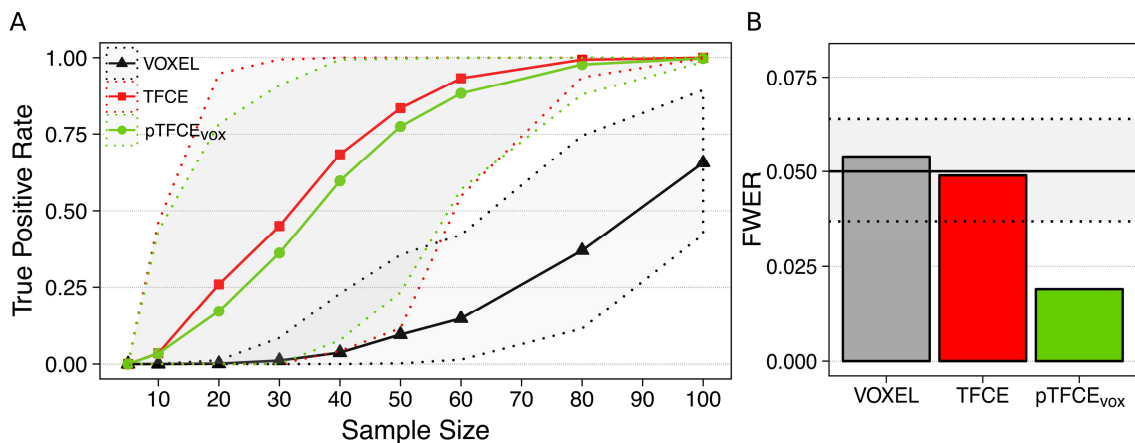
580

581

582

583

584



585 Figure 7. Performance of pTFCE on the “taskswitch” paradigm of the NCP dataset (Gorgolewski et al., 2017; Poldrack
 586 et al., 2016). (A) Counting the voxels of the full-sample (N=119) group-level FWER-corrected, unenhanced activation
 587 map also found by the investigated methods in various subsamples and plotting the mean (solid lines) and the 95%
 588 confidence intervals (dotted lines) as a function of sample sizes provides an estimate of the improvement in
 589 statistical power. (B) Family-wise error rates (FWER) of the investigated methods. Methods VOXEL and TFCE were
 590 thresholded using permutation test while pTFCE was thresholded at the same threshold as VOXEL, resulting in the
 591 pTFCE_{vox} approach. Dotted lines denote the 95% confidence interval for FWER for 1000 repetitions.

588 1.7.1 Demonstration of increased statistical power

589 Our real-data analysis demonstrates that both enhancement methods (TFCE and pTFCE_{vox}) result in a
 590 substantial increase of detected “true positive” voxels (see Methods for our assumption on true
 591 positives). For the investigated dataset (Figure 8A), the increase in statistical power introduced by
 592 the cluster enhancement methods allows for detecting about half of the “true activation” on average
 593 at a sample size of about N=35, as opposed to the sample size of about N=90 needed for the
 594 unenhanced statistical inference for the same average performance. Moreover, in 95% of the cases
 595 the true positive rate of 0.5 was already reached by the enhancement methods at a sample size 60.
 596 On the other hand, for the same sample size of N=20, the enhanced statistical inference might
 597 already detect 25% of the “true” activation while without enhancement no voxels at all were
 598 detected in most of the random samples. While the average true positive rate for pTFCE_{vox} is slightly
 599 lower than that of TFCE, across the random samples of sizes 30-60, pTFCE_{vox} yielded a narrower
 600 confidence interval and therefore a more pronounced separation (more robust improvement) from
 601 the unenhanced analysis.

602 1.7.2 Evaluation of family-wise error rate

603 Family-wise error rates in our real-data scenario with a true-null hypothesis were found to be in the
 604 nominal range for the VOXEL and TFCE methods (0.054 and 0.049, respectively) and, in accordance to
 605 the simulation results, below the nominal range (0.02) for pTFCE_{vox}.

606

607

608

609

610 ***1.7.3 Demonstration of enhancing an activation map with pTFCE***

611 Applying TFCE and pTFCE on real data resulted in an enhanced detection of pain-related activation in
612 both cases. The total number of FWER-corrected suprathreshold voxels approximately doubled for
613 both approaches (from 23665 voxels to 46722 for TFCE and to 50063 for pTFCE). The newly detected
614 areas, besides border regions of the unenhanced activation pattern, involve also completely new
615 activation maxima (See **Error! Reference source not found.** and Supplementary Table 2 and 3). New
616 local maxima emerge in the thalamus, brainstem, amygdala, caudate nuclei, pallidum, middle
617 cingulate cortex, middle frontal gyrus, postcentral gyrus, planum temporale and superior parietal
618 lobule.

619 Besides the striking similarity of the TFCE and the pTFCE map, several activations, among others, in
620 the cerebellum and the brainstem (and additionally, some border regions) were only detected by
621 pTFCE. In contrast, some blobs in white matter were detected only by TFCE. Notably, due to the
622 required permutation test, processing time of TFCE is longer than that of pTFCE by more than an
623 order magnitude.

624 ***Discussion***

625 In this paper, we have formulated and comprehensively evaluated a novel, generalized approach
626 which unifies the advantages of cluster- and voxel-wise statistical inference. Since the basic concepts
627 of our method are similar to those of the threshold-free cluster enhancement approach (TFCE)
628 (Smith and Nichols, 2009), we refer to our method as probabilistic TFCE or pTFCE.

629 In a pure theoretical sense, as we start to use spatial neighbourhood information to boost
630 neuroimaging signals we should inherently start to lose spatial localization accuracy. However, the
631 inherent smoothness of typical neuroimaging data (even without artificial smoothing) does not allow
632 for taking advantage of the high localising performance of the simple voxel-wise inference. On the
633 other hand, incorporating image smoothness (or clustering) information into the statistical inference
634 might also be considered as optimising the “localisation performance – sensitivity” trade-off by
635 throwing out only localisation capacity which is “unutilised” (due to smoothness). Both TFCE and
636 pTFCE take advantage of this property of neuroimages.

637 As opposed to the simple cluster-wise inference, both TFCE and pTFCE generate a voxel-wise output
638 image and maintain information about spatial detail within extended areas of signal. For example,
639 local maxima in the pTFCE output image can easily be identified, and separated from each other if a
640 “cluster” contains more than one maximum. These local maxima locations will be identical to those
641 in the original statistical image. This means that, similarly to TFCE, pTFCE provides rich and
642 interpretable output, retaining much more spatial information than the simple cluster-level
643 inference.

644 Our evaluation on simulated and real data underpins former results (Smith and Nichols, 2009) and
645 clearly justifies that incorporating spatial topological information when testing voxel-level differences
646 in neurological images provides a significant improvement in sensitivity over the simple, unenhanced
647 voxel-level inference (**Error! Reference source not found.**B, C). While the mathematical background
648 of pTFCE is significantly different from that of TFCE, importantly, pTFCE results are highly similar to

649 those of TFCE (**Error! Reference source not found.**B), suggesting that the spatial localisation
650 performances of the two methods are similar.

651 Our analysis with optimized values of smoothing reveals that pTFCE prefers smoothing extents
652 similar to TFCE, and typically smaller than the unenhanced inference. This can be considered as a
653 desirable property, because an extensive artificial smoothing changes the area of activation and the
654 positions of local maxima on the image, possibly leading to false conclusions.

655 The results also suggest that, in terms of family-wise sensitivity, pTFCE is *always superior* to the
656 uncorrected inference. Moreover, in doing so, pTFCE is *more stable* than TFCE, given the few cases of
657 the latter producing lower AUC values than the unenhanced inference (typically at high SNRs with
658 test signals 1 and 2; see Figure 7). The rationale behind these results might be that TFCE, with $E=0.5$,
659 $H=2$, is possibly better optimised for more extended signal shapes (where it performs very well),
660 while pTFCE does not have any free parameter to be optimised and, accordingly, it might implement
661 a more objective enhancement, unbiased regarding of the shape and extent of the true signal. TFCE
662 tends to produce higher AUC values than pTFCE for spatially extended test signals, mainly at low
663 signal-to-noise ratios. However, these cases were successively paired with an elevated number of
664 positive observations within background regions as revealed by our “negative AFROC” analysis.

665 Here we argue that this undesirable property is most probably the result of “cluster leaking”: when,
666 typically, low-threshold clusters containing many false positive voxels are integrated during
667 enhancement so that areas of no signal became enhanced enough to be detected even with FWER-
668 level correction. TFCE, with $E=0.5$, $H=2$, possibly implements a trade-off by allowing for some more
669 FPs but, at the same time, capturing the extended low-signal boundary regions, as well, thereby
670 being closer to what we could intuitively feel “true” for a big cluster. Nevertheless, in terms of this
671 issue, pTFCE seems to be more “pragmatic”, by maintaining an acceptable low-level of “leaked” false
672 positives even with morphologically complex and spatially extended ground truth signals. In other
673 words, the spatial localisation performance of pTFCE seems to be somewhat superior to that of TFCE
674 in case of large-extent signal shapes.

675 Combining the results of analysis with optimal and identical parameter-sets suggest that the above
676 discussed differences are not a consequence of suboptimal parameter settings for any of the
677 methods, but general differences in their overall performance. Therefore, we encourage the use of
678 pTFCE in any case, as an alternative to TFCE.

679 Importantly, our results also indicate that, when thresholding the pTFCE-enhanced image at FWER
680 corrected values completely based on the original unenhanced image, the thresholded image still
681 provides improved sensitivity, and gives a strict control for FWER, as well, comparable to that within
682 the original image. (see $pTFCE_{vox}$ on **Error! Reference source not found.**, Figure 6 and Figure 7).
683 Therefore, in contrast to TFCE, pTFCE does not require a permutation test, because the threshold
684 values obtained for the unenhanced image can be directly applied on the enhanced image. This
685 property renders pTFCE suitable for a wide range of studies, for instance with study designs where
686 permutation tests are not possible or not preferred.

687 We took advantage of this property when demonstrating pTFCE on real data. Our results
688 demonstrate that both TFCE and pTFCE result in a significant increase in statistical power (Figure 8A),
689 allowing for detecting the same activations at a lower sample size or more true positive voxels at the
690 same sample size as compared to the unenhanced statistical inference. While the tested $pTFCE_{vox}$
691 approach (that is using the GRF-theory based maximum-height threshold of the unenhanced values)
692 seems to perform slightly inferior to TFCE, it also seems to maintain narrower confidence intervals
693 across the random samples and therefore a more robust improvement over the unenhanced VOXEL
694 approach. While it was not investigated in this study on real data, analysis of simulated data suggests
695 that correcting pTFCE for multiple comparisons with a permutation-based approach (that is using

696 pTFCE instead of $pTFCE_{vox}$) would further improve sensitivity, and depending on the spatial topology
697 of the true signal, might even outperform TFCE.

698 Analysis of family-wise error rates (Figure 8B) also suggest, that $pTFCE_{vox}$ gives an overly strict control
699 of FWER. Based on our simulations, correcting pTFCE for multiple comparisons with a permutation-
700 based approach might yield a more liberal thresholding and result in nominal FWER values. This
701 aspect, together with a detailed evaluation of the effect of inhomogeneous spatial autocorrelation
702 on pTFCE is subject to further investigation.

703 Demonstrating the effect of enhancement approaches on an activation map (capturing response to
704 painful stimuli) reveals that both TFCE and pTFCE makes several extra activations detectable. Some
705 of these activation areas include completely new local maxima, meaning that the enhanced
706 activation map is not only a simply “boosted” version of the original image, with stronger activation
707 borders, but indeed, new activations are brought over the level of significance. The majority of the
708 newly detected activation maxima are considered to be part of the pain matrix, thus naturally fit into
709 the results of (Bingel et al., 2011). With the applied methods, pTFCE ($pTFCE_{vox}$) discovers slightly
710 more new voxels than TFCE and those seem to be in more relevant locations (cerebellum, brainstem)
711 than those specific for TFCE only (white matter). While analysis of the real-data example also
712 demonstrates that it is valid to threshold pTFCE images with values set up for the unenhanced image,
713 our simulation results still suggest, that, when it comes to correction for multiple comparisons,
714 permutation based maximum height thresholding using the empirical distribution of the enhanced
715 data gives a slightly better sensitivity. Nevertheless, the improvement over TFCE in the most realistic
716 simulation cases is present at the $pTFCE_{vox}$ approach (that is, without the need for permutation test),
717 as well. This improved sensitivity, together with the marginal processing time (compared to
718 permutation-test required for TFCE) renders pTFCE beneficial and easily applicable with any
719 neuroimaging workflow.

720

721 ***Limitations***

722 An obvious limitation of our study is that in the simulations we applied a stationary Gaussian noise
723 model, without any inhomogeneous spatial autocorrelation. Although, in that aspect, we did not
724 capture some relevant properties of real neuroimaging data, this is more than a reasonable
725 simplification: it is a standard first approximation of modelling neuroimaging data and as such, a
726 necessary first step in evaluating pTFCE. Therefore, the presented work should be considered as a
727 basis for further investigations aiming at the evaluation and adjustment of the pTFCE approach for
728 autocorrelation and nonstationarities in the data. Nevertheless, since pTFCE implements a concept
729 similar to TFCE, we could expect a similar robustness (Salimi-Khorshidi et al., 2011) to these
730 characteristics of real datasets. Indeed, our initial analysis (Figure 8B) revealed that, despite the
731 assumption of homogenous spatial autocorrelation in the GRF-based estimation of probabilities,
732 $pTFCE_{vox}$ still gives a very strict control of FWER on real data. Furthermore, the GRF theory based
733 implementation of TFCE can trivially be extended to consider local properties of smoothness,
734 analogous to (Salimi-Khorshidi et al., 2011), or alternatively, permutation-based nonparametric
735 estimation of the cluster-size distributions can be also used, yielding a “brute force” solution to the
736 issues of spatial inhomogeneity of image smoothness in real data.

737 ***Although, when introducing the theoretical background, we***
738 ***clarified that GRF theory is only one way to estimate the necessary***

739 **distributions for pTFCE, the main aim of this paper was to establish**
740 **and validate the links of pTFCE to existing, “de facto” standards in**
741 **neuroimaging analysis, and accordingly focused on the GRF-based**
742 **solution. Exploring the novel possibilities provided by non-GRF**
743 **based pTFCE solutions (e.g. cluster enhancement on graphs),**
744 **together with the comprehensive evaluation of the effect of spatial**
745 **inhomogeneities, will be the topic of upcoming research. Conclusion**

746 Here, we have proposed a novel approach for the integration of information about autocorrelation
747 into mass-univariate statistical analysis. Our solution, called pTFCE, can be considered as an
748 improvement and generalisation of the widely used threshold-free cluster enhancement (TFCE)
749 approach. While the theory behind pTFCE allows for generalising the approach for various data
750 structures, in this paper we have focused on a Gaussian Random Field-based implementation in
751 order to establish clear links to the standard volumetric analysis methodology in neuroimaging.

752 In our evaluation, we found that pTFCE is more robust to various ground truth shapes and provides a
753 stricter control over cluster “leaking” than TFCE and, in some realistic cases, further improves its
754 sensitivity. The fact that, as opposed to TFCE, pTFCE directly outputs (enhanced) P-values, makes it
755 well-suitable for the improvement of statistical inference in any neuroimaging workflow.

756 Importantly, the presented GRF-based likelihood function in the Bayesian formulation of pTFCE can
757 easily be exchanged, thus pTFCE is easy to adapt for data structures other than images (e.g.
758 skeletons, surfaces or graphs), and carries the potential to deploy the concept of topology-based
759 enhancement of statistical inference on a wider range of applications than ever before.

760 Various software implementations and documentation of the pTFCE approach are available at
761 <https://spisakt.github.io/pTFCE>.

762 **Acknowledgements**

763 We thank Prof. Irene Tracy (Nuffield Department of Clinical Neurosciences, University of Oxford, UK),
764 Krzysztof J. Gorgolewski and Russell A. Poldrack (Department of Psychology, Stanford University,
765 USA) for sharing the fMRI datasets used for the validation on real data. We are also thankful to Dr.
766 András Czürkó (Gedeon Richter Plc., Budapest, Hungary) for his support and valuable insights
767 regarding the possible applications of pTFCE.

768 **References**

- 769 Bardeen, J.M., Bond, J., Kaiser, N., Szalay, A., 1986. The statistics of peaks of Gaussian random fields.
770 The Astrophysical Journal 304, 15-61.
- 771 Bennett, C.M., Baird, A.A., Miller, M.B., Wolford, G.L., 2011. Neural correlates of interspecies
772 perspective taking in the post-mortem atlantic salmon: an argument for proper multiple comparisons
773 correction. Journal of Serendipitous and Unexpected Results 1, 1-5.
- 774 Bingel, U., Wanigasekera, V., Wiech, K., Ni Mhuircheartaigh, R., Lee, M.C., Ploner, M., Tracey, I., 2011.
775 The effect of treatment expectation on drug efficacy: imaging the analgesic benefit of the opioid
776 remifentanyl. Sci Transl Med 3, 70ra14.
- 777 Bunch, P.C., Hamilton, J.F., Sanderson, G.K., Simmons, A.H., 1977. A free response approach to the
778 measurement and characterization of radiographic observer performance. Application of Optical
779 Instrumentation in Medicine VI. International Society for Optics and Photonics, pp. 124-135.

- 780 Chakraborty, D.P., Winter, L., 1990. Free-response methodology: alternate analysis and a new
781 observer-performance experiment. *Radiology* 174, 873-881.
- 782 Clayden, J., 2014. *mmand: Mathematical Morphology in Any Number of Dimensions*. London, UK: R
783 package version 1.
- 784 Forman, S.D., Cohen, J.D., Fitzgerald, M., Eddy, W.F., Mintun, M.A., Noll, D.C., 1995. Improved
785 assessment of significant activation in functional magnetic resonance imaging (fMRI): use of a
786 cluster-size threshold. *Magnetic Resonance in medicine* 33, 636-647.
- 787 Friston, K.J., Holmes, A., Poline, J.B., Price, C.J., Frith, C.D., 1996. Detecting activations in PET and
788 fMRI: levels of inference and power. *Neuroimage* 4, 223-235.
- 789 Friston, K.J., Holmes, A.P., Worsley, K.J., Poline, J.P., Frith, C.D., Frackowiak, R.S., 1994a. Statistical
790 parametric maps in functional imaging: a general linear approach. *Human brain mapping* 2, 189-210.
- 791 Friston, K.J., Worsley, K.J., Frackowiak, R., Mazziotta, J.C., Evans, A.C., 1994b. Assessing the
792 significance of focal activations using their spatial extent. *Human brain mapping* 1, 210-220.
- 793 Genest, C., Zidek, J.V., 1986. Combining probability distributions: A critique and an annotated
794 bibliography. *Statistical Science*, 114-135.
- 795 Gorgolewski, K., Burns, C.D., Madison, C., Clark, D., Halchenko, Y.O., Waskom, M.L., Ghosh, S.S.,
796 2011. Nipype: a flexible, lightweight and extensible neuroimaging data processing framework in
797 python. *Frontiers in neuroinformatics* 5.
- 798 Gorgolewski, K.J., Durnez, J., Poldrack, R.A., 2017. Preprocessed Consortium for Neuropsychiatric
799 Phenomics dataset. *F1000Research* 6.
- 800 Habib, G., Steve, S., Thomas, N., 2017. Threshold-Free Cluster Enhancement revisited: Increasing
801 Power and Spatial specificity of TFCE.
- 802 Jenkinson, M., Beckmann, C.F., Behrens, T.E., Woolrich, M.W., Smith, S.M., 2012. *Fsl*. *Neuroimage* 62,
803 782-790.
- 804 Lieberman, M.D., Cunningham, W.A., 2009. Type I and Type II error concerns in fMRI research: re-
805 balancing the scale. *Social cognitive and affective neuroscience* 4, 423-428.
- 806 Lohmann, G., Stelzer, J., Mueller, K., Lacosse, E., Buschmann, T., Kumar, V.J., Grodd, W., Scheffler, K.,
807 2017. Inflated false negative rates undermine reproducibility in task-based fMRI. *bioRxiv*, 122788.
- 808 Nichols, T., Hayasaka, S., 2003. Controlling the familywise error rate in functional neuroimaging: a
809 comparative review. *Statistical methods in medical research* 12, 419-446.
- 810 Nichols, T.E., 2012. Multiple testing corrections, nonparametric methods, and random field theory.
811 *Neuroimage* 62, 811-815.
- 812 Nosko, V., 1969. LOCAL STRUCTURE OF GAUSSIAN RANDOM FIELDS IN VICINITY OF HIGH-LEVEL
813 SHINES. *DOKLADY AKADEMII NAUK SSSR* 189, 714-&.
- 814 Poldrack, R.A., Congdon, E., Triplett, W., Gorgolewski, K.J., Karlsgodt, K.H., Mumford, J.A., Sabb, F.W.,
815 Freimer, N.B., London, E.D., Cannon, T.D., Bilder, R.M., 2016. A phenome-wide examination of neural
816 and cognitive function. *Scientific data* 3, 160110.
- 817 Salimi-Khorshidi, G., Smith, S.M., Nichols, T.E., 2011. Adjusting the effect of nonstationarity in cluster-
818 based and TFCE inference. *Neuroimage* 54, 2006-2019.
- 819 Smith, S.M., Jenkinson, M., Johansen-Berg, H., Rueckert, D., Nichols, T.E., Mackay, C.E., Watkins, K.E.,
820 Ciccarelli, O., Cader, M.Z., Matthews, P.M., Behrens, T.E., 2006. Tract-based spatial statistics:
821 voxelwise analysis of multi-subject diffusion data. *Neuroimage* 31, 1487-1505.
- 822 Smith, S.M., Jenkinson, M., Woolrich, M.W., Beckmann, C.F., Behrens, T.E., Johansen-Berg, H.,
823 Bannister, P.R., De Luca, M., Drobnjak, I., Flitney, D.E., 2004. Advances in functional and structural
824 MR image analysis and implementation as FSL. *Neuroimage* 23, S208-S219.
- 825 Smith, S.M., Nichols, T.E., 2009. Threshold-free cluster enhancement: addressing problems of
826 smoothing, threshold dependence and localisation in cluster inference. *Neuroimage* 44, 83-98.
- 827 Stone, M., 1961. The opinion pool. *The Annals of Mathematical Statistics* 32, 1339-1342.
- 828 Vinokur, L., Zalesky, A., Raffelt, D., Smith, R., Connelly, A., 2015. A Novel Threshold-Free
829 Network-Based Statistics Method Demonstration using Simulated Pathology. *ISMRM*.
- 830 Vul, E., Harris, C., Winkielman, P., Pashler, H., 2009. Puzzlingly High Correlations in fMRI Studies of
831 Emotion, Personality, and Social Cognition. *Perspect Psychol Sci* 4, 274-290.

- 832 Whitcher, B., Schmid, V.J., Thornton, A., 2011. Working with the DICOM and NIfTI Data Standards in
833 R. *Journal of Statistical Software* 44, 1-28.
- 834 Wickham, H., 2016. *ggplot2: elegant graphics for data analysis*. Springer.
- 835 Woo, C.-W., Krishnan, A., Wager, T.D., 2014. Cluster-extent based thresholding in fMRI analyses:
836 pitfalls and recommendations. *Neuroimage* 91, 412-419.
- 837 Worsley, K.J., Marrett, S., Neelin, P., Vandal, A.C., Friston, K.J., Evans, A.C., 1996. A unified statistical
838 approach for determining significant signals in images of cerebral activation. *Hum Brain Mapp* 4, 58-
839 73.
- 840 Worsley, K.J., Taylor, J.E., Tomaiuolo, F., Lerch, J., 2004. Unified univariate and multivariate random
841 field theory. *Neuroimage* 23 Suppl 1, S189-195.

842