

1 Precise temporal localisation of M/EEG effects with
2 Bayesian generalised additive multilevel models

3 Ladislas Nalborczyk¹ and Paul Bürkner²

4 ¹Aix Marseille Univ, CNRS, LPL

5 ²TU Dortmund University, Department of Statistics

6 Abstract

7 Time-resolved electrophysiological measurements such as those obtained through magneto- or electro-encephalography (M/EEG) offer a unique window into the neural activity underlying cognitive processes. Researchers are often interested in determining whether and when these signals differ across experimental conditions or participant groups. The conventional approach involves mass-univariate statistical testing across time and space, followed by corrections for multiple comparisons such as cluster-based inference. While effective for controlling error rates at the cluster-level, cluster-based inference comes with a significant limitation: by shifting the focus of inference from individual time points to clusters, it makes difficult to draw precise conclusions about the onset or offset of observed effects. Here, we introduce a *model-based* approach for analysing M/EEG timeseries such as event-related potentials (ERPs) or decoding performance over time. Our approach leverages Bayesian generalised additive multilevel models, providing posterior probabilities that an effect is above zero (or above chance) at each time point, while naturally accounting for temporal dependencies and between-subject variability. Using both simulated and actual M/EEG datasets, we demonstrate that this approach substantially outperforms conventional methods in estimating the onset and offset of neural effects, yielding more precise and reliable results. We provide an R package implementing the method and describe how it can be integrated into M/EEG analysis pipelines using MNE-Python.

Keywords: EEG, MEG, cluster-based inference, simulation, multiple comparisons, generalised additive models, mixed-effects models, multilevel models, Bayesian statistics, brms

8 **Table of contents**

9 Introduction	3
10 Introduction	3
11 Problem statement	3

12	Statistical errors and cluster-based inference	4
13	Previous work on modelling M/EEG data	5
14	Generalised additive models	6
15	Bayesian generalised additive multilevel models	7
16	Objectives	8
17	Methods	8
18	M/EEG data simulation	8
19	Model description and model fitting	8
20	Error properties of the proposed approach	12
21	Comparing the onsets/offsets estimates from other approaches	13
22	Simulation study	14
23	Application to actual MEG data	14
24	Results	16
25	Simulation study (bias and variance)	16
26	Application to actual MEG data (reliability)	17
27	Discussion	20
28	Summary of the proposed approach	20
29	Limitations and future directions	20
30	Data and code availability	22
31	Packages	22
32	Acknowledgements	22
33	References	23
34	Application to 2D time-resolved decoding results (cross-temporal generalisation)	29
35	Alternative to GAMs: Approximate Gaussian Process regression	32
36	How to choose the GAM basis dimension?	34
37	R package and integration with MNE-Python	35

Ladislas Nalborczyk  <https://orcid.org/0000-0002-7419-9855>

Paul Bürkner  <https://orcid.org/0000-0001-5765-8995>

The authors have no conflicts of interest to disclose.

Correspondence concerning this article should be addressed to Ladislas Nalborczyk, Aix Marseille Univ, CNRS, LPL, 5 avenue Pasteur, 13100 Aix-en-Provence, France, email: ladislas.nalborczyk@cnrs.fr

38 **Precise temporal localisation of M/EEG effects with Bayesian generalised additive
39 multilevel models**

1 **Introduction**

2 **Problem statement**

3 Understanding the temporal dynamics of cognitive processes requires methods that can
4 capture fast-changing neural activity with high temporal resolution. Magnetoencephalography
5 and electroencephalography (M/EEG) are two such methods, widely used in cognitive neuro-
6 science for their ability to track brain activity at the millisecond scale. These techniques provide
7 rich time series data that reflect how neural responses unfold in response to stimuli or tasks. A
8 central goal in many M/EEG studies is to determine whether, when, and where neural responses
9 differ across experimental conditions or participant groups.

10 The conventional approach involves mass-univariate statistical testing through time
11 and/or space followed by some form or correction for multiple comparisons with the goal of
12 maintaining the familywise error rate (FWER) or the false discovery rate (FDR) at the nomi-
13 nal level (e.g., 5%). Cluster-based inference is the most common way of achieving this sort of
14 error control in the M/EEG literature, being the recommended approach in several software
15 programs (e.g., [EEGLab](#), [Delorme & Makeig, 2004](#); [MNE-Python](#), [Gramfort, 2013](#)). While ef-
16 fective for controlling error rates, cluster-based inference comes with a significant limitation:
17 by shifting the focus of inference from individual datapoints (e.g., timesteps, sensors, voxels) to
18 clusters, it prevents the ability to draw precise conclusions about the spatiotemporal localisation
19 of such effects ([Maris & Oostenveld, 2007](#); [Sassenhagen & Draschkow, 2019](#)). As pointed out
20 by Maris & Oostenveld (2007): “there is a conflict between this interest in localized effects and
21 our choice for a global null hypothesis: by controlling the FA [false alarm] rate under this global
22 null hypothesis, one cannot quantify the uncertainty in the spatiotemporal localization of the
23 effect”. Even worse, Rosenblatt et al. (2018) note that cluster-based inference suffers from low
24 spatial resolution: “Since discovering a cluster means that ‘there exists at least one voxel with
25 an evoked response in the cluster’, and not that ‘all the voxels in the cluster have an evoked
26 response’, it follows that the larger the detected cluster, the less information we have on the
27 location of the activation.” As a consequence, cluster-based inference is expected to perform
28 poorly for localising the onset of M/EEG effects; a property that was later demonstrated in
29 simulations studies (e.g., [Rousselet, 2025](#); [Sassenhagen & Draschkow, 2019](#)).

30 To overcome the limitations of cluster-based inference, we introduce a novel *model-based*
31 approach for precisely localising M/EEG effects in time, space, and other dimensions. The pro-
32 posed approach, based on Bayesian generalised additive multilevel models, allows quantifying
33 the posterior probability of effects being above chance at the level of timesteps, sensors, voxels,
34 etc, while naturally taking into account spatiotemporal dependencies present in M/EEG data.
35 We compare the performance of the proposed approach to well-established alternative meth-
36 ods using both simulated and actual M/EEG data and show that it significantly outperforms
37 alternative methods in estimating the onset and offset of M/EEG effects.

38 **Statistical errors and cluster-based inference**

39 The issues with multiple comparisons represent a common and well-recognised danger in
 40 neuroimaging and M/EEG research where the collected data allows for a multitude of potential
 41 hypothesis tests and is characterised by complex structures of spatiotemporal dependencies. The
 42 probability of obtaining at least one false positive in an ensemble (family) of m tests (i.e., the
 43 FWER) is computed as $1 - (1 - \alpha)^m$ (for $m = 10$ tests and $\alpha = 0.05$, it is approximately equal
 44 to 0.4). Different methods exist to control the FWER (i.e., bring it back to α). Most methods
 45 apply a simple correction to series of “raw” p -values issued from univariate statistical tests (e.g.,
 46 t-tests). For instance, the Bonferroni correction (Dunn, 1961) consists in setting the significance
 47 threshold to α/m , or equivalently, multiplying the “raw” p -values by m and using the standard
 48 α significance threshold. This method is generally overconservative (i.e., under-powered) as
 49 it assumes statistical independence of the tests, an assumption that is clearly violated in the
 50 context of M/EEG timeseries characterised by massive spatiotemporal dependencies. Some
 51 alternative methods aims at controlling the FDR, defined as the proportion of false positive
 52 *among positive tests* (e.g., Benjamini & Hochberg, 1995; Benjamini & Yekutieli, 2001). However,
 53 a major limitation of both types of corrections is that they do not take into account the spatial
 54 and temporal information contained in M/EEG data.

55 A popular technique to account for spatiotemporal dependencies while controlling the
 56 FWER is cluster-based inference (Bullmore et al., 1999; Maris & Oostenveld, 2007). A typ-
 57 ical cluster-based inference consists of two successive steps (for more details on cluster-based
 58 inference, see for instance Frossard & Renaud, 2022; Maris, 2011; Maris & Oostenveld, 2007;
 59 Sassenhagen & Draschkow, 2019). First, clusters are defined as sets of contiguous timesteps,
 60 sensors, voxels, etc, whose activity, summarised by some test statistic (e.g., a t -value), exceeds
 61 a predefined threshold (e.g., the 95th percentile of the parametric null distribution). Clusters
 62 are then characterised by their height (i.e., maximal value), extent (number of constituent ele-
 63 ments), or some combination of both, for instance by summing the statistics within a cluster, an
 64 approach referred to as “cluster mass” (Maris & Oostenveld, 2007; Pernet et al., 2015). Then,
 65 the null hypothesis is tested by computing a p -value for each identified cluster by comparing
 66 its mass with the null distribution of cluster masses (obtained via permutation). As alluded
 67 previously, a significant cluster is a cluster which contains *at least one* significant time-point. As
 68 such, it would be incorrect to conclude, for instance, that the timestep of a significant cluster
 69 is the first moment at which some conditions differ (Frossard & Renaud, 2022; Sassenhagen
 70 & Draschkow, 2019). In other words, because the inference is performed at the second step
 71 (i.e., once clusters have been formed), it prevents any conclusion to be made about individual
 72 datapoints (e.g., timesteps, sensors, etc).

73 As different cluster-forming thresholds lead to clusters with different spatial or temporal
 74 extent, this initial threshold modulates the sensitivity of the subsequent permutation test. The
 75 threshold-free cluster enhancement (TFCE) method was introduced by S. Smith & Nichols
 76 (2009) to overcome this choice of an arbitrary threshold. In brief, the TFCE method works
 77 as follows. Instead of picking an arbitrary cluster-forming threshold (e.g., $t = 2$), the methods
 78 consist in trying all (or many) possible thresholds in a given range and checking whether a given
 79 datapoint (e.g., timestep, sensor, voxel) belongs to a significant cluster under any of the set of

80 thresholds. Then, instead of using cluster mass, one uses a weighted average between the cluster
 81 extend (e , how broad is the cluster, that is, how many connected samples it contains) and the
 82 cluster height (h , how high is the cluster, that is, how large is the test statistic). The TFCE
 83 score at each timestep t is given by:

$$\text{TFCE}(t) = \int_{h=h_0}^{h=h_t} e(h)^E h^H dh$$

84 where h_0 is typically 0 and parameters E and H are set a priori (typically to 0.5 and
 85 2, respectively) and control the influence of the extend and height on the TFCE. Note that in
 86 practise, this integral is approximated by a sum over small h increments. Then, a p-value for
 87 each timestep t is computed by comparing its TFCE with the null distribution of TFCE values
 88 (obtained via permutation). For each permuted signal, we keep the maximal value over the whole
 89 signal for the null distribution of the TFCE. The TFCE combined with permutation (assuming
 90 a large enough number of permutations) has been shown to provide accurate FWER control
 91 (e.g., Pernet et al., 2015). However, further simulation work showed that cluster-based methods
 92 (including TFCE) perform poorly in localising the onset of M/EEG effects (e.g., Rousselet,
 93 2025; Sassenhagen & Draschkow, 2019).

94 To sum up, the main limitation of cluster-based inference is that it allows for inference
 95 at the cluster level only, not allowing inference at the level of timesteps, sensors, etc. As a
 96 consequence, it does not allow inferring the precise spatial and temporal localisation of effects.
 97 In the following, we briefly review previous M/EEG modelling work. Then, we provide a
 98 short introduction to generalised additive models (GAMs) and Bayesian generalised additive
 99 multilevel models (BGAMMs) to illustrate how these models can be used to precisely localise
 100 the onset and offset of M/EEG effects.

101 Previous work on modelling M/EEG data

102 Recent example of GLM for EEG (Fischer & Ullsperger, 2013; Wüllhorst et al., 2025)...
 103 See also (Hauk et al., 2006; Rousselet et al., 2008)... Example of two-stage regression analysis
 104 (i.e., individual-level then group-level, Dunagan et al., 2024)...

105 As put by Rousselet (2025), concluding on the onset of effect based on a series of univari-
 106 ate tests... commits to three fallacies... here, we want to avoid these by introducing a model-based
 107 approach, which naturally take into account the temporal dependencies in the data to output
 108 a series of posterior probabilities...

109 See also the rERP framework (N. J. Smith & Kutas, 2014a, 2014b) and Tremblay &
 110 Newman (2014)...

111 From Dimigen & Ehinger (2021): Recently, spline regression has been applied to ERPs
 112 (e.g., Hendrix et al., 2017; Kryuchkova et al., 2012)... GAMMs for EEG data (Abugaber et
 113 al., 2023; Meulman et al., 2015, 2023)... See also the UNFOLD toolbox (Ehinger & Dimigen,
 114 2019)...

115 Disentangling overlapping processes (Skukies et al., 2024; Skukies & Ehinger, 2021)...
 116 Weighting single trials (Pernet, 2022)... The LIMO toolbox (Pernet et al., 2011) using linear
 117 and 2-stage regression (not a proper multilevel model)...

Recently, Teichmann (2022) provided a detailed tutorial on using Bayes factors (BFs) to analyse the 1D or 2D output from MVPA, that is, for testing, at every timestep, whether decoding performance is above chance level. However, this approach provides timeseries of BFs that ignores temporal dependencies...

Generalised additive models

See for instance these tutorials (Sóskuthy, 2017; Winter & Wieling, 2016) or application to phonetic data (Sóskuthy, 2021; Wieling, 2018) or this introduction (Baayen & Linke, 2020) or these reference books (Hastie & Tibshirani, 2017; Wood, 2017a)... application to pupillometry (Rij et al., 2019)... GAMLSS for neuroimaging data (Dinga et al., 2021)... Modelling auto-correlation in GAMMs + EEG example (Baayen et al., 2018)...

In generalised additive models (GAMs), the functional relationship between predictors and response variable is decomposed into a sum of low-dimensional non-parametric functions. A typical GAM has the following form:

$$y_i \sim \text{EF}(\mu_i, \phi)$$

$$g(\mu_i) = \underbrace{\mathbf{A}_i \gamma}_{\text{parametric part}} + \underbrace{\sum_{j=1}^J f_j(x_{ij})}_{\text{non-parametric part}}$$

where $y_i \sim \text{EF}(\mu_i, \phi)$ denotes that the observations y_i are distributed as some member of the exponential family of distributions (e.g., Gaussian, Gamma, Beta, Poisson) with mean μ_i and scale parameter ϕ ; $g(\cdot)$ is the link function, \mathbf{A}_i is the i th row of a known parametric model matrix, γ is a vector of parameters for the parametric terms (to be estimated), f_j is a smooth function of covariate x_j (to be estimated as well). The smooth functions f_j are represented in the model via penalised splines basis expansions of the covariates, that are a weighted sum of K simpler, basis functions:

$$f_j(x_{ij}) = \sum_{k=1}^K \beta_{jk} b_{jk}(x_{ij})$$

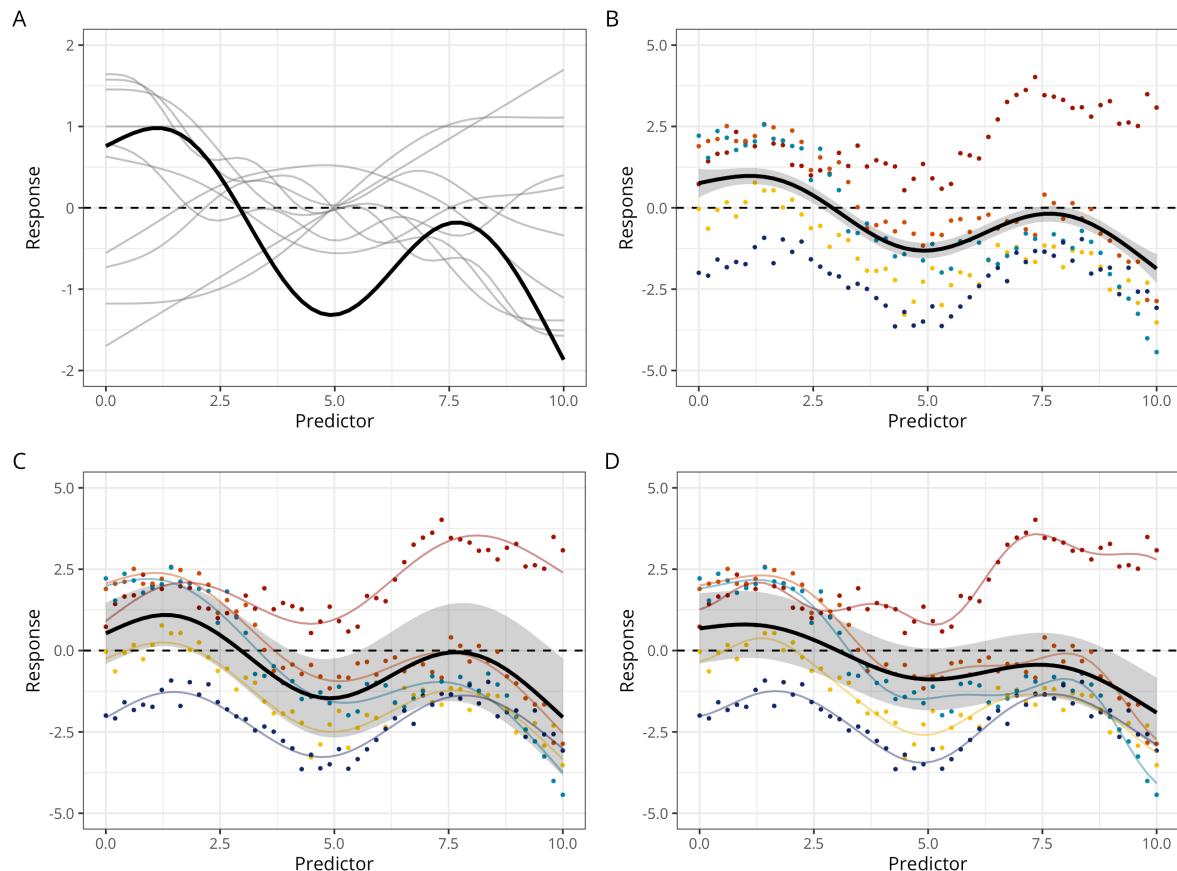
where β_{jk} is the weight (coefficient) associated with the k th basis function $b_{jk}()$ evaluated at the covariate value x_{ij} for the j th smooth function f_j . To clarify the terminology at this stage: *splines* are functions composed of simpler functions. These simpler functions are basis functions (e.g., cubic polynomial, thin-plate) and the set of basis functions is a *basis*. Each basis function gets its coefficient and the resultant spline is the sum of these weighted basis functions (Figure 1). Splines coefficients are penalised (usually through the squared of the smooth functions' second derivative) in a way that can be interpreted, in Bayesian terms, as a prior on the “wigginess” of the function. In other words, more complex (wiggly) basis functions are penalised.

¹⁴⁷ **Bayesian generalised additive multilevel models**

¹⁴⁸ The Bayesian approach to statistical modelling is characterised by its reliance on probability theory to (Gelman et al., 2020)... In this framework, all unknown entities are assigned probability distributions reflecting the uncertainty... These probability distributions are commonly referred to as “priors” and represent some state of knowledge about unknown quantities before seeing any data. There are debates among Bayesian practitioners as to whether prior distributions should encode subjective (personal) beliefs or... but these debates are outside the scope of the present paper and we therefore leave the interested reader to dedicated work (e.g., XX; YY)... In practice, weakly informative priors are often used as default priors in situations in which subjective priors are difficult to define/elicit... Bayesian models are then fitted on empirical (actual or simulated) data to update prior states of knowledge to posterior states of knowledge using Bayes theorem, or in practise, sampling-based approximations of the posterior distribution...

Figure 1

Different types of GAM(M)s. **A:** GAMs predictions are computed as the weighted sum (in black) of basis functions (here thin-plate basis functions, in grey). **B:** Constant-effect GAM, with 5 participants in colours and the group-level prediction in black. **C:** Varying-intercept + varying-slope GAMM (with common smoother). **D:** Varying-intercept + varying-slope + varying-smoother GAMM. In this model, each participant gets its own intercept, slope, and degree of ‘wiggleness’ (smoother).



¹⁶⁰ See Figure 1... Introduction to multilevel GAMs (E. J. Pedersen et al., 2019)... Now

161 describe the Bayesian GAMM (Miller, 2025)... Proper inclusion of varying/random effects in the
 162 model specification protects against overly wiggly curves (Baayen & Linke, 2020)... Generalising
 163 to scale and shape or “distributional GAMs” (Rigby & Stasinopoulos, 2005; Umlauf et al., 2018)
 164 and applied to neuroimaging data (Dinga et al., 2021)...

165 Instead of averaging, obtain the smooth ERP signal from multilevel GAM... less suscep-
 166 tible to outliers (Meulman et al., 2023)...

167 Objectives

168 Given the previously reported limitations of conventional methods to precisely identify
 169 the onset and offset of M/EEG effects (e.g., ERPs, decoding performance), we developed a
 170 model-based approach for estimating the onset and offset of such effects. To achieve this, we
 171 leveraged Bayesian generalised additive multilevel models (BGAMMs) fitted in R via the **brms**
 172 package and compared the performance of this approach to conventional methods on both
 173 simulated and actual M/EEG data.

174 Methods

175 M/EEG data simulation

176 Following the approach of Sassenhagen & Draschkow (2019) and Rousselet (2025), we
 177 simulated EEG data stemming from two conditions, one with noise only, and the other with
 178 noise + signal. As in previous studies, the noise was generated by superimposing 50 sinusoids
 179 at different frequencies, following an EEG-like spectrum (see code in the online supplementary
 180 materials and details in Yeung et al., 2004). As in Rousselet (2025), the signal was generated
 181 from a truncated Gaussian distribution with an objective onset at 160 ms, a peak at 250 ms,
 182 and an offset at 342 ms. We simulated this signal for 250 timesteps between 0 and 0.5s, akin to
 183 a 500 Hz sampling rate. We simulated data for a group of 20 participants (with variable true
 184 onset) with 50 trials per participant and condition (Figure 2).

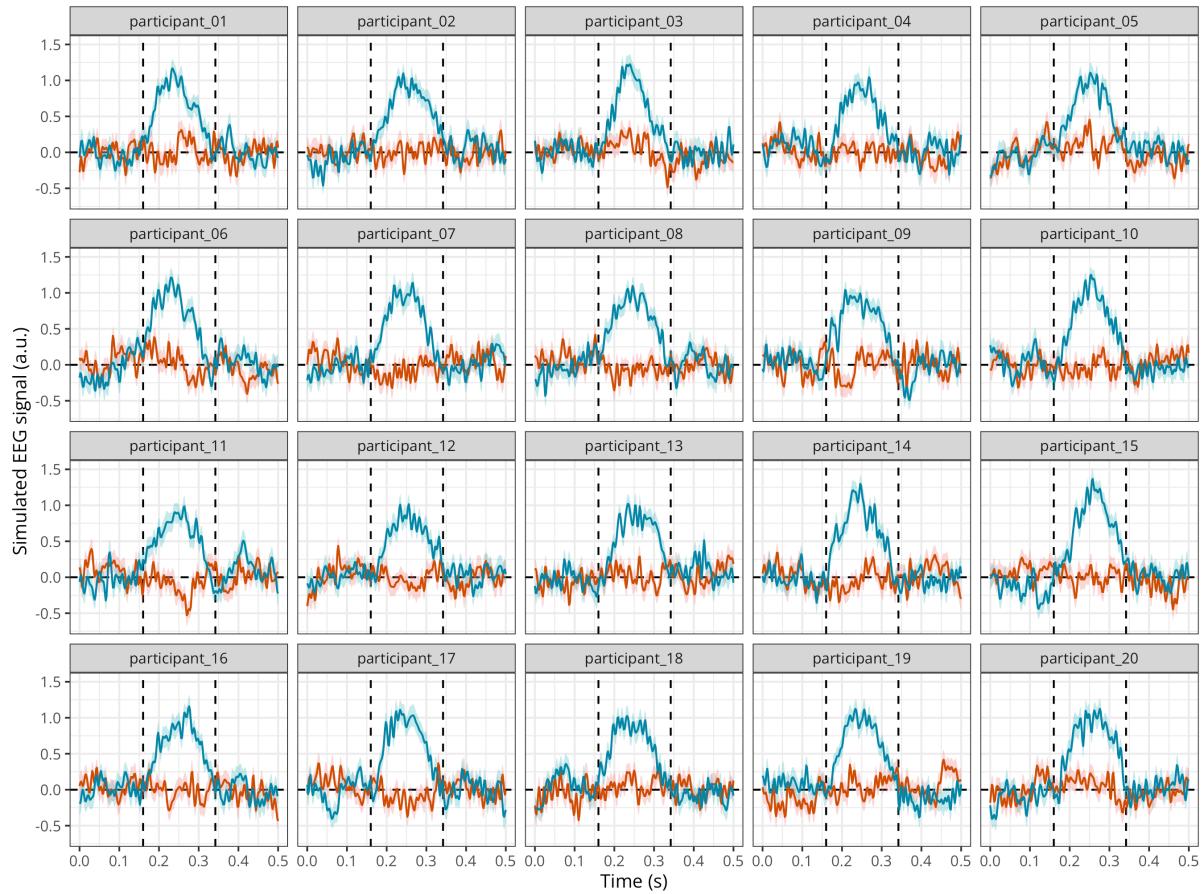
185 Model description and model fitting

186 We then fitted a Bayesian GAM (BGAM) using the **brms** package (Bürkner, 2017,
 187 2018; Nalborczyk et al., 2019). We used the default priors in **brms** (i.e., weakly informative
 188 priors). We ran eight Markov Chain Monte-Carlo (MCMC) to approximate the posterior dis-
 189 tribution, including each 5000 iterations and a warmup of 2000 iterations, yielding a total of
 190 $8 \times (5000 - 2000) = 24000$ posterior samples to use for inference. Posterior convergence was
 191 assessed examining trace plots as well as the Gelman–Rubin statistic \hat{R} (Gabry et al., 2019;
 192 Gelman et al., 2020). The **brms** package uses the same syntax as the R package **mgcv** v 1.9-
 193 1 (Wood, 2017b) for specifying smooth effects. Figure 3 shows the predictions of this model
 194 together with the raw data.

```
# averaging across participants
ppt_df <- raw_df %>%
  group_by(participant, condition, time) %>%
  summarise(eeg = mean(eeg) ) %>%
```

Figure 2

Mean simulated EEG activity in two conditions with 50 trials each, for a group of 20 participants. The error band represents the mean +/- 1 standard error of the mean.



```
ungroup()
```

```
# defining a contrast for condition
contrasts(ppt_df$condition) <- c(-0.5, 0.5)

# fitting the BGAM
gam <- brm(
  # thin-plate regression splines with k-1 basis functions
  eeg ~ condition + s(time, bs = "tp", k = 20, by = condition),
  data = ppt_df,
  family = gaussian(),
  warmup = 2000,
  iter = 5000,
  chains = 8,
  cores = 8,
  file = "models/gam.rds"
)
```

195 However, the previous model only included constant (fixed) effects, thus not prop-
 196 erly accounting for between-participant variability. We next fit a multilevel version of the
 197 BGAM (BGAMM, for an introduction to Bayesian multilevel models in `brms`, see [Nalborczyk](#)
 198 [et al., 2019](#)). Although it is possible to fit a BGAMM at the single-trial level, we present a
 199 computationally-lighter version of the model that is fitted directly on by-participant summary
 200 statistics (mean and SD), similar to what is done in meta-analysis.

```
# averaging across participants
summary_df <- raw_df %>%
  summarise(
    eeg_mean = mean(eeg),
    eeg_sd = sd(eeg),
    .by = c(participant, condition, time)
  )

# defining a contrast for condition
contrasts(summary_df$condition) <- c(-0.5, 0.5)

# fitting the BGAMM
meta_gam <- brm(
  # using by-participant SD of ERPs across trials
  eeg_mean | se(eeg_sd) ~
    condition + s(time, bs = "cr", k = 20, by = condition) +
    (1 | participant),
  data = summary_df,
  family = gaussian(),
  warmup = 2000,
  iter = 5000,
  chains = 8,
  cores = 8,
  file = "models/meta_gam.rds"
)

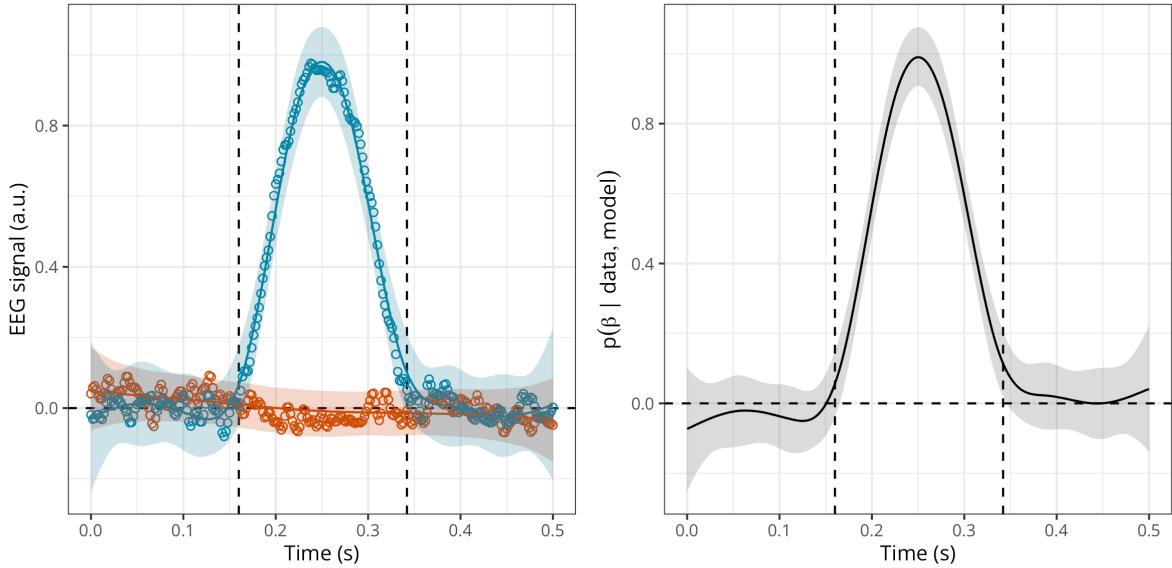
# fitting the BGAMM
# meta_gamm <- brm(
#   # using by-participant SD of ERPs across trials
#   eeg_mean | se(eeg_sd) ~ condition +
#   s(participant, bs = "re") +
#   s(time, participant, bs = "fs", by = condition, m = 1),
#   data = summary_df,
#   family = gaussian(),
#   warmup = 2000,
#   iter = 5000,
```

```
#     chains = 8,
#     cores = 8,
#     control = list(adapt_delta = 0.95),
#     file = "models/meta_gamm.rds"
#   )
```

201 We depict the posterior predictions together with the posterior estimate of the slope
 202 for **condition** at each timestep (Figure 3). This figure suggests that the BGAMM provides an
 203 adequate description of the simulated data (see further posterior predictive checks in Section C).

Figure 3

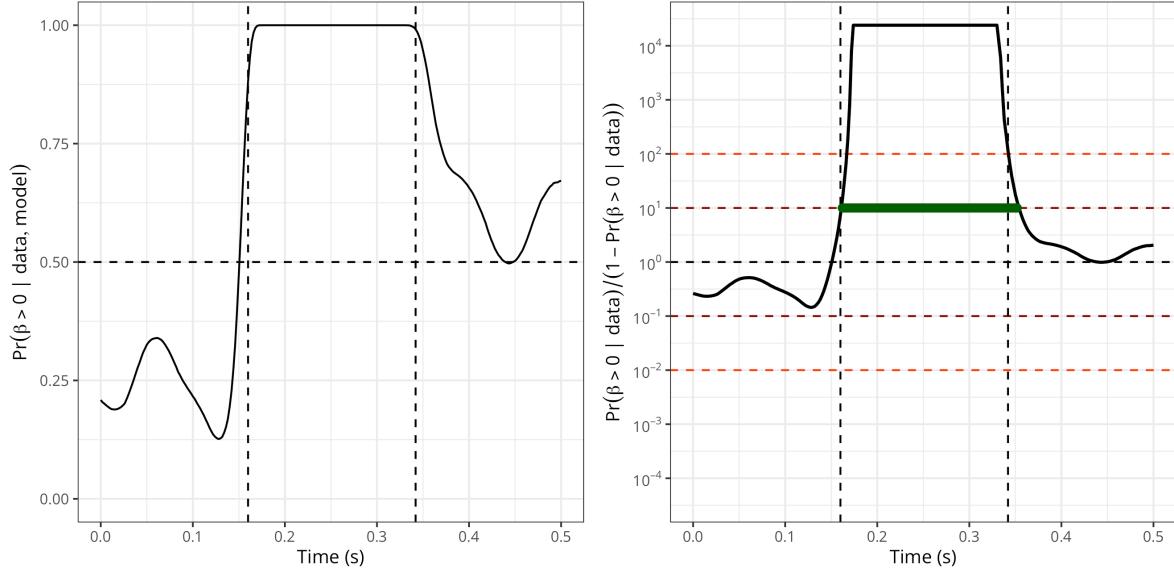
Posterior estimate of the EEG activity in each condition (left) and posterior estimate of the difference in EEG activity (right) according to the BGAMM.



204 We then compute the posterior probability of the slope for **condition** being above
 205 0 (Figure 4, left). From this quantity, we then compute the ratio of posterior probabilities
 206 (i.e., $p/(1-p)$) and visualise the timecourse of this ratio superimposed with the conventional
 207 thresholds on evidence ratios (Figure 4, right). Note that a ratio of 10 means that the probability
 208 of the difference being above 0 is 10 times higher than the probability of the difference not being
 209 above 0, given the data, the priors, and other model's assumptions. Thresholding the posterior
 210 probability ratio thus provides a model-based approach for estimating the onset and offset of
 211 M/EEG effects. An important advantage is that the proposed approach can be extended to
 212 virtually any model structure. Moreover, the approach provides both group-level and individual-
 213 level onset and offset estimates of M/EEG effects.

Figure 4

Left: Posterior probability of the EEG difference (slope) being above 0 according to the BGAMM. Right: Ratio of posterior probability according to the BGAMM (on a log10 scale). Timesteps above threshold (10) are highlighted in green. NB: the minimum and maximum possible ratio values are determined (bounded) by the number of available posterior samples in the model.

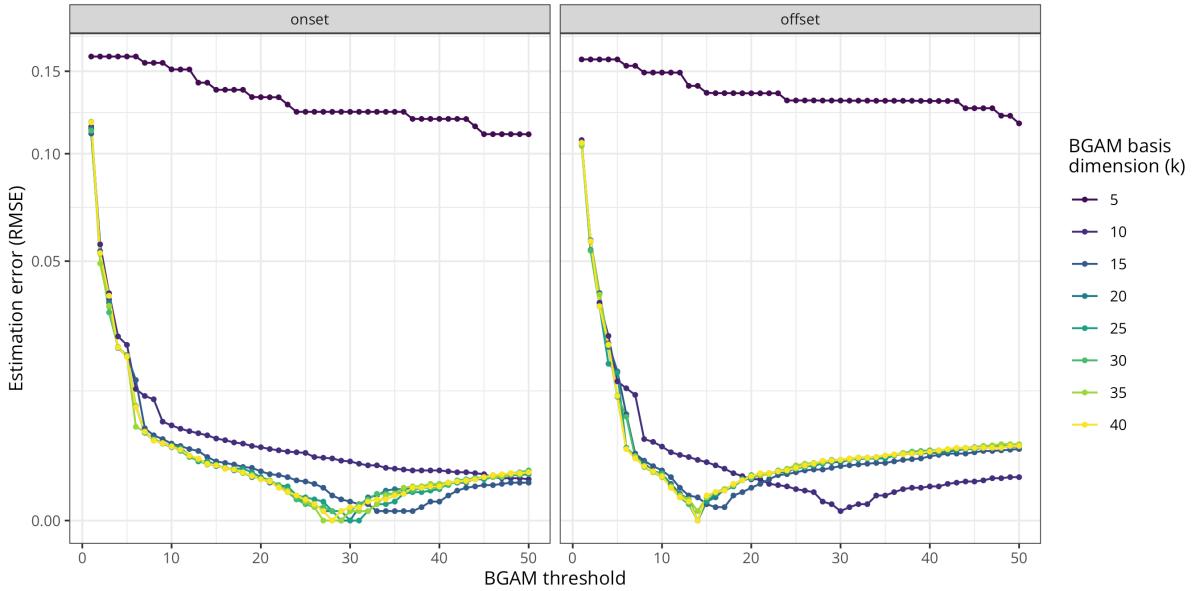


214 Error properties of the proposed approach

215 We then assess the performance of the proposed approach by computing the differ-
 216 ence between the true and estimated onset/offset of the EEG difference according to various
 217 **k** (BGAM basis dimension) and **threshold** values. Remember that the EEG signal was gen-
 218 erated from a truncated Gaussian with an objective onset at 160 ms, a maximum at 250 ms,
 219 and an offset at 342 ms. Figure 5 shows that the multilevel GAM can almost exactly recover
 220 the true onset and offset values, given some reasonable choice of **k** and **threshold** values. We
 221 provide more detailed recommendations on how to set **k** in Section C. This figure further reveals
 222 that the optimal **k** and **threshold** values may differ for the onset and offset values, and there
 223 there seems to exist a trade-off between these two parameters: lower **k** values lead to poorer
 224 estimations, but these poor estimations can be compensated (only to some extent) by higher
 225 **threshold** values (and reciprocally).

Figure 5

Average estimation error (RMSE) for the onset (left) and offset (right) according to various basis dimension and threshold values for the BGAM (computed from 100 simulated datasets).

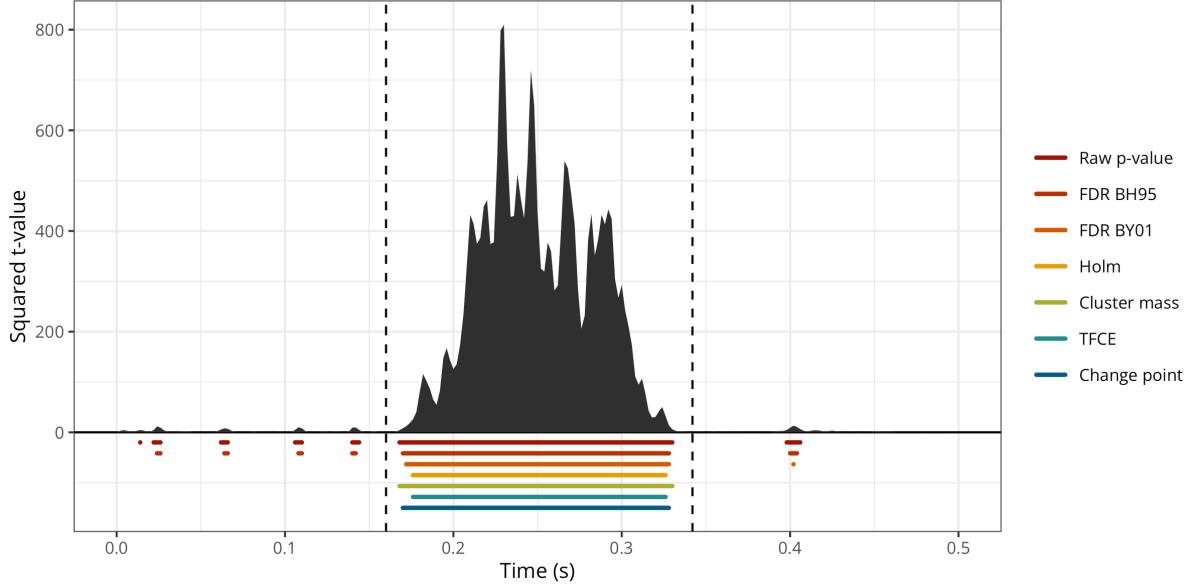


226 Comparing the onsets/offsets estimates from other approaches

227 We then compared the ability of the BGAMM to accurately estimate the onset and offset
 228 of the ERP difference to other widely-used methods. First, we conducted mass-univariate t-
 229 tests (thus treating each timestep independently) and identified the onset and offset of the ERP
 230 difference as the first and last values crossing an arbitrary significance threshold ($\alpha = 0.05$). We
 231 then followed the same approach but after applying different forms of multiplicity correction
 232 to the p -values. We compared two methods that control the FDR (i.e., BH95, Benjamini &
 233 Hochberg, 1995; and BY01, Benjamini & Yekutieli, 2001), one method that controls the FWER
 234 (i.e., Holm–Bonferroni method, Holm, 1979), and two cluster-based permutation methods (per-
 235 mutation with a single cluster-forming threshold and threshold-free cluster enhancement, TFCE,
 236 S. Smith & Nichols, 2009). The BH95, BY01, and Holm corrections were applied to the p-
 237 values using the `p.adjust()` function in R. The cluster-based inference was implemented using
 238 a cluster-sum statistic of squared t -values, as implemented in MNE-Python (Gramfort, 2013),
 239 called via the R package `reticulate` v 1.42.0 (Ushey et al., 2024). We also compared these
 240 estimates to the onset and offset as estimated using the binary segmentation algorithm, as im-
 241 plemented in the R package `changepoint` v 2.3 (Killick et al., 2022), and applied directly to
 242 the squared t -values (as in Rousselet, 2025). Figure 6 illustrates the onsets and offsets esti-
 243 mated by each method on a single simulated dataset and shows that all methods systematically
 244 overestimate the true onset and underestimate the true offset.

Figure 6

Exemplary timecourse of squared t-values with true onset and offset (vertical black dashed lines) and onsets/offsets identified using the raw p-values, the corrected p-values (BH95, BY01, Holm), the cluster-based methods (Cluster mass, TFCE), or using the binary segmentation method (Change point).



245 Simulation study

246 To assess the accuracy of group-level onset estimation, the various methods were com-
 247 pared using the bias (i.e., median(estimated-true)), median absolute error (MAE), root mean
 248 square error (RMSE), variance, and median absolute deviation (MAD) of onset/offset estimates
 249 computed on 1,000 simulated datasets. As in Rousselet (2025), each participant was assigned a
 250 random onset between 150 and 170ms. Whereas the present article focuses on one-dimensional
 251 signals (e.g., one M/EEG channel), we provide a 2D application in Section A.

252 Application to actual MEG data

253 To complement the simulation study, we evaluated the performance of the various meth-
 254 ods on actual MEG data (decoding results from Nalborczyk et al., in preparation). In this study,
 255 we conducted time-resolved multivariate pattern analysis (MVPA, also known as decoding) of
 256 MEG data during reading tasks. As a result, we obtain a timecourse of decoding performance
 257 (ROC AUC), bounded between 0 and 1, for each participant ($N = 32$). Next, we wanted to test
 258 whether the group-level average decoding accuracy is above chance (i.e., 0.5) at each timestep
 259 (Figure 7). To achieve this, we fitted a BGAM as introduced previously, but we replaced the
 260 Normal likelihood function by a Beta one to account for the bounded nature of AUC values
 261 (between 0 and 1) (for a tutorial on Beta regression, see Coretta & Bürkner, 2025).

262 Note that although we chose a basis dimension of $k = 50$, which seems appropriate for
 263 the present data, this choice should be adapted according to the properties of the modelled
 264 data (e.g., signal-to-noise ratio, prior low-pass filtering, sampling rate) and should be assessed
 265 by the usual model checking tools (e.g., posterior predictive checks, see also Section C). To

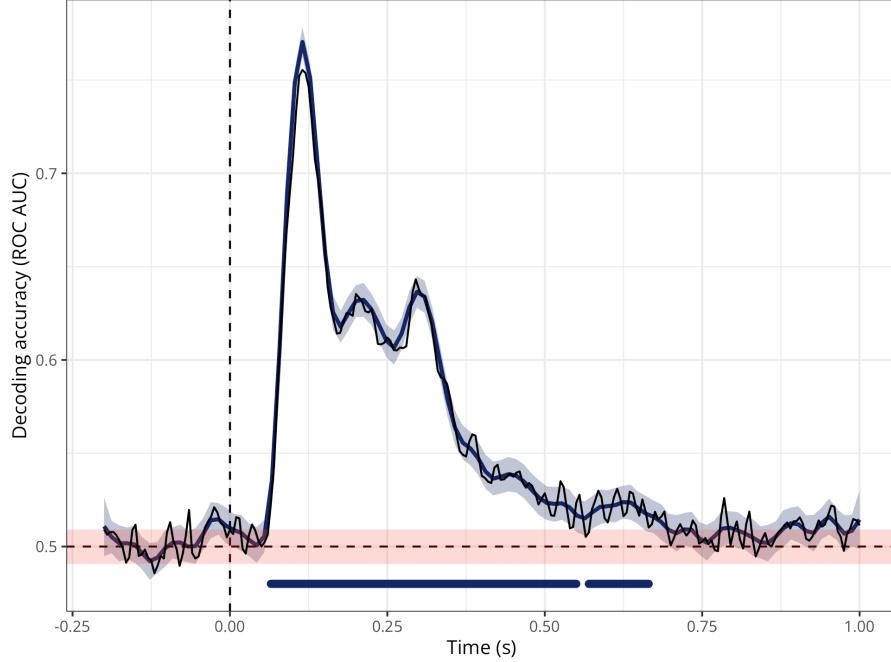
266 better distinguish signal from noise, we also defined a region of practical equivalence (ROPE,
 267 [Kruschke & Liddell, 2017](#)), defined as the chance level plus the standard deviation of the (group-
 268 level average) decoding performance during the baseline period.

```
# fitting the Beta GAM
meg_decoding_gam <- brm(
  auc ~ s(time, bs = "cr", k = 50),
  data = decoding_df,
  family = Beta(),
  warmup = 2000,
  iter = 5000,
  chains = 4,
  cores = 4
)
```

269 We assessed the reliability of the proposed approach using a form of permutation-based
 270 split-half reliability (as for instance in [Rosenblatt et al., 2018](#)), which consisted of the following
 271 steps. First, we created 1,000 split halves of the data (i.e., with half the participants in the
 272 original data, that is, 16 participants). For each split, we estimated the onset/offset using all
 273 methods described previously. Third, we summarised the distribution of onset/offset estimates
 274 using the median “error” (i.e., difference between the split estimate and the estimate obtained
 275 using the full dataset) and the variance across splits. This approach allows assessing how similar
 276 the estimate of each half split is to the full dataset (thus acting as a proxy for the population)
 277 and how variable the estimates are across split halves.

Figure 7

Group-level average decoding performance ($N=32$) superimposed with the GAM predictions (in blue) and the region of practical equivalence (ROPE, in orange) computed from the baseline period (data from Nalborczyk et al., in preparation). The blue horizontal markers indicate the timesteps at which the posterior probability ratio exceeds 20.



278

Results

279 This section is divided in two parts. First, we present the results from the simulation
 280 study, assessing the bias and variance of each method when applied to simulated data in which
 281 the ground truth is known. Second, we present the results obtained when applying the different
 282 methods to actual MEG data (decoding performance through time), assessing the reliability of
 283 the estimates provided by each method.

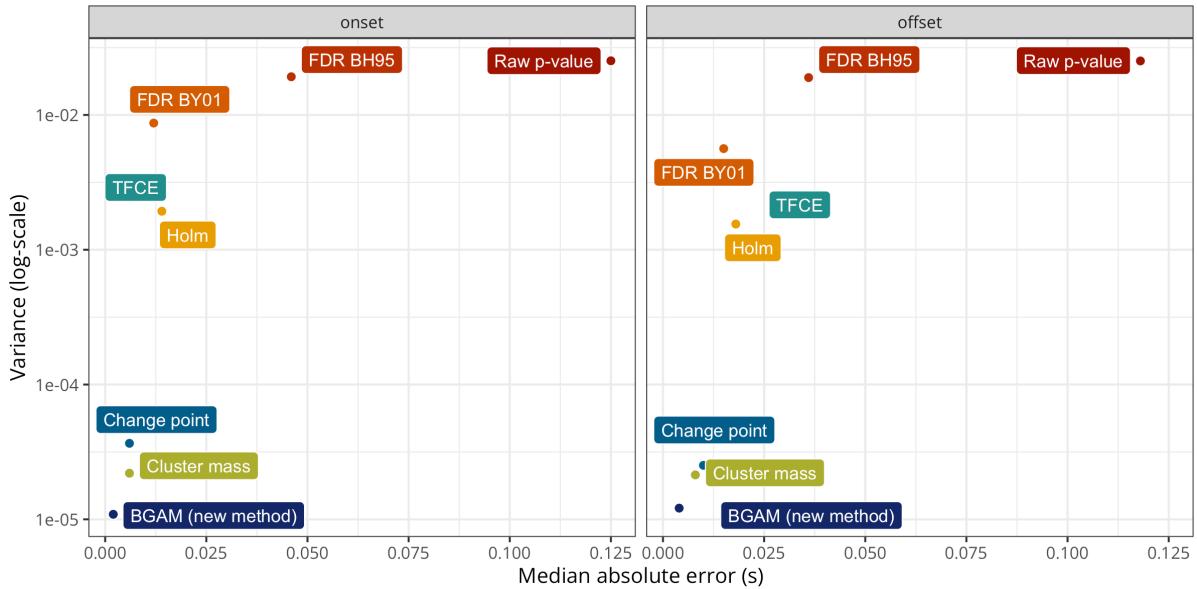
284 **Simulation study (bias and variance)**

285 Figure 8 shows a summary of the simulation results, revealing that the proposed ap-
 286 proach (BGAM) has the lowest median absolute error (MAE) and variance for both the onset
 287 and offset estimates. The Cluster mass and Change point also have good performance, but
 288 surprisingly, the TFCE method has relatively bad performance for estimating the effect offset
 289 (similar performance to the Holm and FDR BY01 methods). Unsurprisingly, the Raw p-value
 290 and FDR BH95 methods show the worst performance.

291 Results are further summarised in Table 1, which shows that the BGAM is perfectly
 292 unbiased (i.e., it has a bias of 0s) for the onset and almost exactly unbiased for the onset (with
 293 a bias of approximately -2ms). The Bias column shows that all other methods tend to estimate
 294 the onset later than the true onset and to estimate the offset earlier than the true offset. As can
 295 be seen from this table, the BGAM has the best performance on all included metrics.

Figure 8

Median absolute error and variance of onset and offset estimates for each method. Variance is plotted on a log10 scale for visual purposes.



296 Application to actual MEG data (reliability)

297 Figure 9 shows the group-level average decoding performance through time with onset
 298 and offset estimates from each method. Overall, this figure shows that both the Raw p-value
 299 and FDR BH95 methods are extremely lenient, considering that the decoding performance is
 300 above chance before the onset of the stimulus (false positive) and until the end of the trial. The
 301 Change point and Cluster mass methods seem the most conservative methods, identifying
 302 a time window from approximately +60ms to +500ms. The Holm, TFCE, and BGAM methods
 303 produce similar estimates of onset and offset, ranging from approximately +60ms to +650ms,
 304 although the BGAM method seems to result in fewer clusters.¹

305 Figure 10 shows the median difference between the onset and offset estimates from
 306 each data split and the onset and offset estimates from the full dataset (x-axis) along with
 307 the variance of its onset and offset estimates across data splits (error bar). This figure reveals
 308 that the BGAM onset and offset estimates on each split are the closest to the estimates from the
 309 full dataset on average (0ms difference for the onset estimate and 5ms difference for the offset
 310 estimate). The Raw p-value method has similar performance, but given the aberrant estimates
 311 it produces (cf. Figure 9), the fact that it is consistent between data splits and the full dataset
 312 is not convincing on its own. The Change point method also has a very good performance
 313 (i.e., very low difference between split estimates and full estimates), but produces too short
 314 cluster of significant decoding performance (cf. Figure 9).² Overall, the figure reveals that for
 315 all other methods, split datasets produce later onset estimates and earlier offset estimates (as
 316 compared to the estimates from the model fitted on the full dataset). These results highlight

¹It should be noted that although each method can produce several “clusters” of timesteps, we only considered the first (onset) and last (offset) timesteps identified by each method to compute the estimation error.

²As in Rousselet (2025), we fixed the number of expected change points to two in the binary segmentation algorithm, thus producing always one cluster.

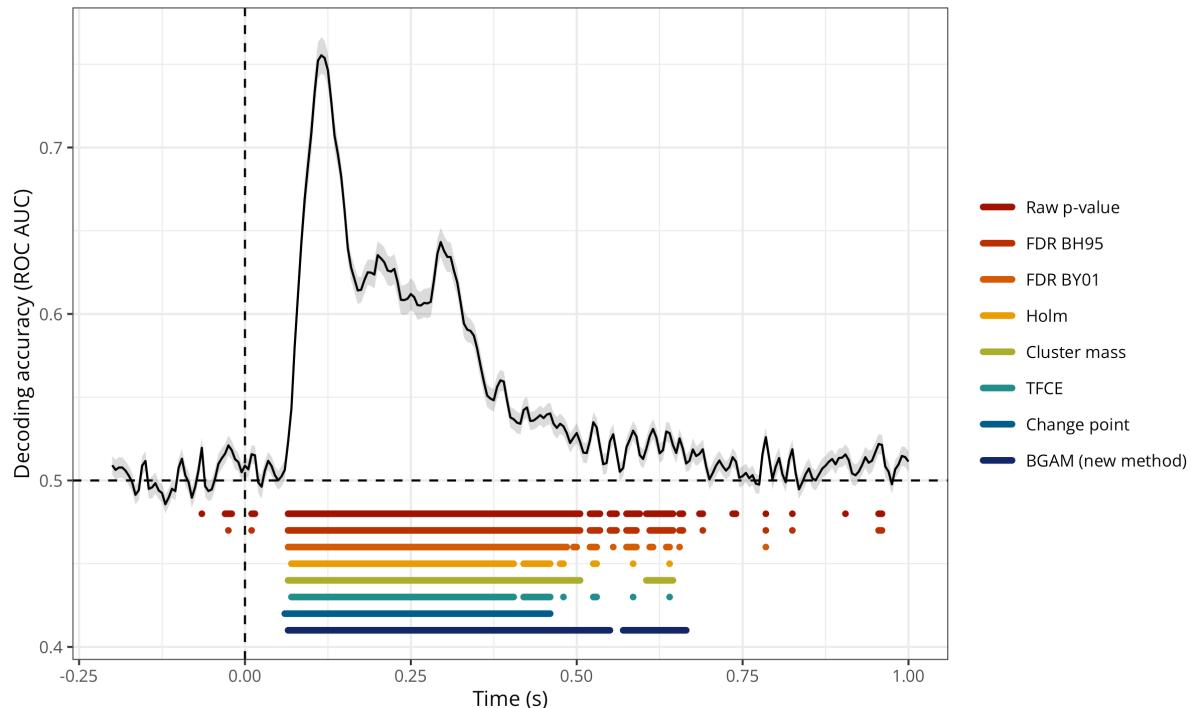
Table 1

Summary statistics of the onset and offset estimates for each method (ordered by the absolute value of the bias).

	Bias	MAE	RMSE	Variance	MAD
onset					
BGAM (new method)	0.0000	0.0020	0.0001	0.0000	0.0030
Cluster mass	0.0060	0.0060	0.0057	0.0000	0.0044
Change point	0.0060	0.0060	0.0051	0.0000	0.0059
Raw p-value	0.0060	0.1250	0.0711	0.0252	0.1942
FDR BH95	0.0080	0.0460	0.0599	0.0192	0.0801
FDR BY01	0.0120	0.0120	0.0443	0.0087	0.0059
TFCE	0.0140	0.0140	0.0270	0.0019	0.0059
Holm	0.0140	0.0140	0.0270	0.0019	0.0059
offset					
BGAM (new method)	-0.0020	0.0040	0.0021	0.0000	0.0030
Cluster mass	-0.0080	0.0080	0.0084	0.0000	0.0059
Raw p-value	-0.0080	0.1180	0.0725	0.0252	0.1868
Change point	-0.0100	0.0100	0.0096	0.0000	0.0059
FDR BH95	-0.0100	0.0360	0.0578	0.0189	0.0682
FDR BY01	-0.0140	0.0150	0.0204	0.0056	0.0059
TFCE	-0.0180	0.0180	0.0261	0.0016	0.0030
Holm	-0.0180	0.0180	0.0262	0.0016	0.0030

Figure 9

Group-level average decoding performance through time with onset and offset estimates for each method (data from Nalborczyk et al., in preparation).



317 some desirable properties for a method aiming to precisely and reliably estimate the onset and
 318 offset of M/EEG effects, namely, it should i) have good asymptotic properties on simulated data,
 319 ii) provide sensible identified clusters in actual data, and iii) provide reliable/stable estimates

320 on actual data.

Figure 10

Median error and median absolute deviation of the error for the onset (left) and offset (right) estimates according to each method. Methods are ordered from lowest (top) to highest (bottom) median absolute error (separately for the onset and offset estimates).

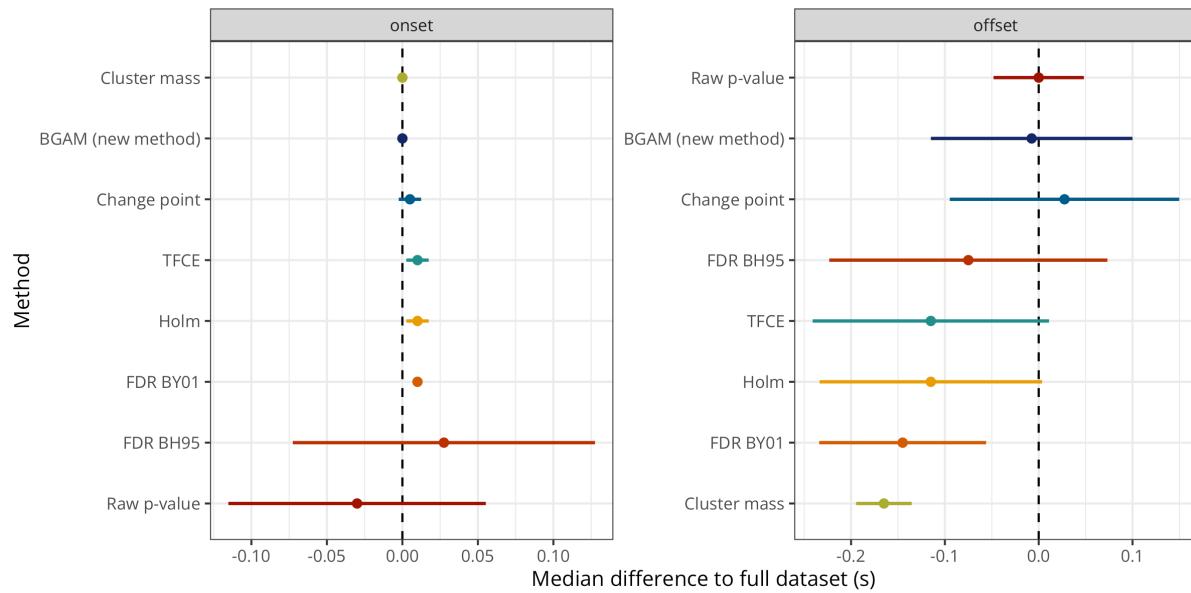


Table 2

Summary statistics of the onset and offset estimates for each method (ordered by the absolute value of the bias).

	Bias	MAE	RMSE	Variance	MAD
onset					
BGAM (new method)	0.0000	0.0000	0.0106	0.0008	0.0000
Cluster mass	0.0000	0.0000	0.0019	0.0000	0.0000
Change point	0.0050	0.0050	0.0011	0.0001	0.0074
TFCE	0.0100	0.0100	0.0120	0.0001	0.0074
FDR BY01	0.0100	0.0100	0.0065	0.0023	0.0000
Holm	0.0100	0.0100	0.0118	0.0001	0.0074
FDR BH95	0.0250	0.0900	0.0101	0.0076	0.0964
Raw p-value	-0.0300	0.0650	0.0201	0.0062	0.0815
offset					
Raw p-value	0.0000	0.0300	0.0256	0.0047	0.0445
BGAM (new method)	-0.0050	0.0750	0.0716	0.0274	0.1038
Change point	0.0250	0.0650	0.0566	0.0090	0.1186
FDR BH95	-0.0750	0.0750	0.1053	0.0153	0.1483
TFCE	-0.1150	0.1450	0.1094	0.0146	0.1260
Holm	-0.1150	0.1450	0.1091	0.0133	0.1186
FDR BY01	-0.1450	0.1450	0.1197	0.0154	0.0890
Cluster mass	-0.1650	0.1650	0.1448	0.0051	0.0297

321

Discussion

322 Summary of the proposed approach

323 Overall, before concluding on the onset/offset of effect based on the model, we need to
 324 ensure that the model provides a faithful description of the data-generating process (e.g., via
 325 posterior predictive checks etc)...

326 The TFCE performs worse than the cluster-sum approach, which was anticipated by
 327 Rousselet (2025) based on the initial results of S. Smith & Nichols (2009)... we did not include
 328 the cluster-depth algorithm (Frossard & Renaud, 2022), as Rousselet (2025) already showed its
 329 performance were worse than the cluster mass algorithm...

330 Limitations and future directions

331 As in previous simulation work (e.g., Rousselet et al., 2008; Sassenhagen & Draschkow,
 332 2019), the present simulation results depend on various choices such as the specific cluster-
 333 forming algorithm and threshold, signal-to-noise ratio, negative impact of preprocessing steps
 334 (e.g., low-pass filter) on temporal resolution... note however, that the same caveats apply to all
 335 methods...

336 Discussion about the nature of the model: we modelled the surface M/EEG signals,
 337 however, the true interests (probably) lie in the brain, that is, in the source space... we could
 338 build a “full” Bayesian model of the generated EEG signal (i.e., including hypotheses about the
 339 source and a forward model), but this model would become computationally heavier...

340 The error properties depend on the threshold parameter, a value of 10 or 20 seems to
 341 be a reasonable default, but the optimal threshold parameter can be adjusted using split-half
 342 reliability assessment... also depends on k...

343 Can be applied to any 1D timeseries (e.g., pupillometry, electromyography)... Extending

344 the approach to spatiotemporal data (i.e., time + sensors) or spatiotemporal time-frequency 4D

345 data...

346 We kept the exemplary models simple, but can be extended by adding varying/random

347 effects (intercept and slope) for item (e.g., word)... but also continuous predictors at the trial

348 level?

349

Data and code availability

350

The simulation results as well as the R code to reproduce the simulations are available on GitHub: https://github.com/lNALBORCZYK/brms_meeg. The `neurogam` R package is available at <https://github.com/lNALBORCZYK/neurogam>.

353

Packages

354

We used R version 4.4.3 ([R Core Team, 2025](#)) and the following R packages: assertthat v. 0.2.1 ([Wickham, 2019](#)), brms v. 2.22.0 ([Bürkner, 2017, 2018, 2021](#)), changepoint v. 2.3 ([Killick et al., 2024; Killick & Eckley, 2014](#)), doParallel v. 1.0.17 ([Corporation & Weston, 2022](#)), easystats v. 0.7.4 ([Lüdecke et al., 2022](#)), foreach v. 1.5.2 ([Microsoft & Weston, 2022](#)), furrr v. 0.3.1 ([Vaughan & Dancho, 2022](#)), future v. 1.34.0 ([Bengtsson, 2021](#)), ggrepel v. 0.9.6 ([Slowikowski, 2024](#)), glue v. 1.8.0 ([Hester & Bryan, 2024](#)), grateful v. 0.2.11 ([Rodriguez-Sanchez & Jackson, 2024](#)), gt v. 1.0.0 ([Iannone et al., 2025](#)), knitr v. 1.50 ([Xie, 2014, 2015, 2025](#)), MetBrewer v. 0.2.0 ([Mills, 2022](#)), neurogam v. 0.0.1 ([Nalborczyk, 2025](#)), pakret v. 0.2.2 ([Gallou, 2024](#)), patchwork v. 1.3.0 ([T. L. Pedersen, 2024](#)), rmarkdown v. 2.29 ([Allaire et al., 2024; Xie et al., 2018, 2020](#)), scales v. 1.3.0 ([Wickham et al., 2023](#)), scico v. 1.5.0 ([T. L. Pedersen & Cramer, 2023](#)), tictoc v. 1.2.1 ([Izrailev, 2024](#)), tidybayes v. 3.0.7 ([Kay, 2024](#)), tidytext v. 0.4.2 ([Silge & Robinson, 2016](#)), tidyverse v. 2.0.0 ([Wickham et al., 2019](#)).

366

Acknowledgements

367

Centre de Calcul Intensif d’Aix-Marseille is acknowledged for granting access to its high performance computing resources.

369

References

- 370 Abugaber, D., Finestrat, I., Luque, A., & Morgan-Short, K. (2023). Generalized additive mixed
 371 modeling of EEG supports dual-route accounts of morphosyntax in suggesting no word
 372 frequency effects on processing of regular grammatical forms. *Journal of Neurolinguistics*,
 373 67, 101137. <https://doi.org/10.1016/j.jneuroling.2023.101137>
- 374 Allaire, J., Xie, Y., Dervieux, C., McPherson, J., Luraschi, J., Ushey, K., Atkins, A., Wickham,
 375 H., Cheng, J., Chang, W., & Iannone, R. (2024). *rmarkdown: Dynamic documents for r*.
 376 <https://github.com/rstudio/rmarkdown>
- 377 Baayen, R. H., & Linke, M. (2020). *Generalized Additive Mixed Models* (pp. 563–591). Springer
 378 International Publishing. https://doi.org/10.1007/978-3-030-46216-1_23
- 379 Baayen, R. H., Rij, J. van, Cat, C. de, & Wood, S. (2018). *Autocorrelated errors in ex-
 380 perimental data in the language sciences: Some solutions offered by generalized additive
 381 mixed models* (pp. 49–69). Springer International Publishing. [https://doi.org/10.1007/978-3-319-69830-4_4](https://doi.org/10.1007/

 382 978-3-319-69830-4_4)
- 383 Bengtsson, H. (2021). A unifying framework for parallel and distributed processing in r using
 384 futures. *The R Journal*, 13(2), 208–227. <https://doi.org/10.32614/RJ-2021-048>
- 385 Benjamini, Y., & Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical and
 386 Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society Series B:
 387 Statistical Methodology*, 57(1), 289–300. <https://doi.org/10.1111/j.2517-6161.1995.tb02031.x>
- 388 x
- 389 Benjamini, Y., & Yekutieli, D. (2001). The control of the false discovery rate in multiple
 390 testing under dependency. *The Annals of Statistics*, 29(4). <https://doi.org/10.1214/aos/1013699998>
- 392 Bullmore, E. T., Suckling, J., Overmeyer, S., Rabe-Hesketh, S., Taylor, E., & Brammer, M. J.
 393 (1999). Global, voxel, and cluster tests, by theory and permutation, for a difference between
 394 two groups of structural MR images of the brain. *IEEE Transactions on Medical Imaging*,
 395 18(1), 32–42. <https://doi.org/10.1109/42.750253>
- 396 Bürkner, P.-C. (2017). brms: An R package for Bayesian multilevel models using Stan. *Journal
 397 of Statistical Software*, 80(1), 1–28. <https://doi.org/10.18637/jss.v080.i01>
- 398 Bürkner, P.-C. (2018). Advanced Bayesian multilevel modeling with the R package brms. *The
 399 R Journal*, 10(1), 395–411. <https://doi.org/10.32614/RJ-2018-017>
- 400 Bürkner, P.-C. (2021). Bayesian item response modeling in R with brms and Stan. *Journal of
 401 Statistical Software*, 100(5), 1–54. <https://doi.org/10.18637/jss.v100.i05>
- 402 Coretta, S., & Bürkner, P.-. C. (2025). *Bayesian beta regressions with brms in r: A tutorial for
 403 phoneticians*. http://dx.doi.org/10.31219/osf.io/f9rqp_v1
- 404 Corporation, M., & Weston, S. (2022). doParallel: Foreach parallel adaptor for the “parallel”
 405 package. <https://CRAN.R-project.org/package=doParallel>
- 406 Delorme, A., & Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-
 407 trial EEG dynamics including independent component analysis. *Journal of Neuroscience
 408 Methods*, 134(1), 9–21. <https://doi.org/10.1016/j.jneumeth.2003.10.009>
- 409 Dimigen, O., & Ehinger, B. V. (2021). Regression-based analysis of combined EEG and eye-
 410 tracking data: Theory and applications. *Journal of Vision*, 21(1), 3. <https://doi.org/10.1111/jovi.13333>

- 411 1167/jov.21.1.3
- 412 Dinga, R., Fraza, C. J., Bayer, J. M. M., Kia, S. M., Beckmann, C. F., & Marquand, A.
413 F. (2021). *Normative modeling of neuroimaging data using generalized additive models of
414 location scale and shape*. <http://dx.doi.org/10.1101/2021.06.14.448106>
- 415 Dunagan, D., Jordan, T., Hale, J. T., Pylkkänen, L., & Chacón, D. A. (2024). *Evaluating
416 the timecourses of morpho-orthographic, lexical, and grammatical processing following rapid
417 parallel visual presentation: An EEG investigation in english*. <http://dx.doi.org/10.1101/2024.04.10.588861>
- 418 Dunn, O. J. (1961). Multiple Comparisons among Means. *Journal of the American Statistical
419 Association*, 56(293), 52–64. <https://doi.org/10.1080/01621459.1961.10482090>
- 420 Ehinger, B. V., & Dimigen, O. (2019). Unfold: An integrated toolbox for overlap correction,
421 non-linear modeling, and regression-based EEG analysis. *PeerJ*, 7, e7838. <https://doi.org/10.7717/peerj.7838>
- 422 Fischer, Adrian G., & Ullsperger, M. (2013). Real and Fictive Outcomes Are Processed Dif-
423 ferently but Converge on a Common Adaptive Mechanism. *Neuron*, 79(6), 1243–1255.
424 <https://doi.org/10.1016/j.neuron.2013.07.006>
- 425 Frossard, J., & Renaud, O. (2022). The cluster depth tests: Toward point-wise strong con-
426 trol of the family-wise error rate in massively univariate tests with application to M/EEG.
427 *NeuroImage*, 247, 118824. <https://doi.org/10.1016/j.neuroimage.2021.118824>
- 428 Gabry, J., Simpson, D., Vehtari, A., Betancourt, M., & Gelman, A. (2019). Visualization in
429 Bayesian workflow. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*,
430 182(2), 389–402. <https://doi.org/10.1111/rssa.12378>
- 431 Gallou, A. (2024). *pakret: Cite “R” packages on the fly in “R Markdown” and “Quarto”*. <https://CRAN.R-project.org/package=pakret>
- 432 Gelman, A., Vehtari, A., Simpson, D., Margossian, C. C., Carpenter, B., Yao, Y., Kennedy,
433 L., Gabry, J., Bürkner, P.-C., & Modrák, M. (2020). Bayesian workflow. *arXiv:2011.01808
434 /Stat*. <http://arxiv.org/abs/2011.01808>
- 435 Gramfort, A. (2013). MEG and EEG data analysis with MNE-python. *Frontiers in Neuro-
436 science*, 7. <https://doi.org/10.3389/fnins.2013.00267>
- 437 Hastie, T. J., & Tibshirani, R. J. (2017). *Generalized Additive Models*. Routledge. <https://doi.org/10.1201/9780203753781>
- 438 Hauk, O., Davis, M. H., Ford, M., Pulvermüller, F., & Marslen-Wilson, W. D. (2006). The
439 time course of visual word recognition as revealed by linear regression analysis of ERP data.
440 *NeuroImage*, 30(4), 1383–1400. <https://doi.org/10.1016/j.neuroimage.2005.11.048>
- 441 Hendrix, P., Bolger, P., & Baayen, H. (2017). Distinct ERP signatures of word frequency, phrase
442 frequency, and prototypicality in speech production. *Journal of Experimental Psychology:
443 Learning, Memory, and Cognition*, 43(1), 128–149. <https://doi.org/10.1037/a0040332>
- 444 Hester, J., & Bryan, J. (2024). *glue: Interpreted string literals*. <https://CRAN.R-project.org/package=glue>
- 445 Holm, S. (1979). A simple sequentially rejective multiple test procedure. *Scandinavian Journal
446 of Statistics*, 6(2), 65–70. <http://www.jstor.org/stable/4615733>
- 447 Iannone, R., Cheng, J., Schloerke, B., Hughes, E., Lauer, A., Seo, J., Brevoort, K., & Roy, O.

- 453 (2025). *gt*: Easily create presentation-ready display tables. [https://CRAN.R-project.org/
454 package=gt](https://CRAN.R-project.org/package=gt)
- 455 Izrailev, S. (2024). *tictoc*: Functions for timing r scripts, as well as implementations of “Stack”
456 and “StackList” structures. <https://CRAN.R-project.org/package=tictoc>
- 457 Kay, M. (2024). *tidybayes*: Tidy data and geoms for Bayesian models. <https://doi.org/10.5281/zenodo.1308151>
- 458 Killick, R., & Eckley, I. A. (2014). changepoint: An R package for changepoint analysis. *Journal
459 of Statistical Software*, 58(3), 1–19. <https://www.jstatsoft.org/article/view/v058i03>
- 460 Killick, R., Haynes, K., & Eckley, I. A. (2022). *changepoint*: An R package for changepoint
461 analysis. <https://CRAN.R-project.org/package=changepoint>
- 462 Killick, R., Haynes, K., & Eckley, I. A. (2024). *changepoint*: An R package for changepoint
463 analysis. <https://CRAN.R-project.org/package=changepoint>
- 464 King, J.-R., & Dehaene, S. (2014). Characterizing the dynamics of mental representations:
465 the temporal generalization method. *Trends in Cognitive Sciences*, 18(4), 203–210. <https://doi.org/10.1016/j.tics.2014.01.002>
- 466 Kruschke, J. K., & Liddell, T. M. (2017). The Bayesian New Statistics: Hypothesis testing,
467 estimation, meta-analysis, and power analysis from a Bayesian perspective. *Psychonomic
468 Bulletin & Review*, 25(1), 178–206. <https://doi.org/10.3758/s13423-016-1221-4>
- 469 Kryuchkova, T., Tucker, B. V., Wurm, L. H., & Baayen, R. H. (2012). Danger and usefulness
470 are detected early in auditory lexical processing: Evidence from electroencephalography.
471 *Brain and Language*, 122(2), 81–91. <https://doi.org/10.1016/j.bandl.2012.05.005>
- 472 Lüdecke, D., Ben-Shachar, M. S., Patil, I., Wiernik, B. M., Bacher, E., Thériault, R., &
473 Makowski, D. (2022). easystats: Framework for easy statistical modeling, visualization,
474 and reporting. CRAN. <https://doi.org/10.32614/CRAN.package.easystats>
- 475 Maris, E. (2011). Statistical testing in electrophysiological studies. *Psychophysiology*, 49(4),
476 549–565. <https://doi.org/10.1111/j.1469-8986.2011.01320.x>
- 477 Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and MEG-data.
478 *Journal of Neuroscience Methods*, 164(1), 177–190. <https://doi.org/10.1016/j.jneumeth.2007.03.024>
- 479 Meulman, N., Sprenger, S. A., Schmid, M. S., & Wieling, M. (2023). GAM-based individual
480 difference measures for L2 ERP studies. *Research Methods in Applied Linguistics*, 2(3),
481 100079. <https://doi.org/10.1016/j.rmal.2023.100079>
- 482 Meulman, N., Wieling, M., Sprenger, S. A., Stowe, L. A., & Schmid, M. S. (2015). Age Effects
483 in L2 Grammar Processing as Revealed by ERPs and How (Not) to Study Them. *PLOS
484 ONE*, 10(12), e0143328. <https://doi.org/10.1371/journal.pone.0143328>
- 485 Microsoft, & Weston, S. (2022). *foreach*: Provides foreach looping construct. <https://CRAN.R-project.org/package=foreach>
- 486 Miller, D. L. (2025). Bayesian views of generalized additive modelling. *Methods in Ecology and
487 Evolution*. <https://doi.org/10.1111/2041-210x.14498>
- 488 Mills, B. R. (2022). *MetBrewer*: Color palettes inspired by works at the metropolitan museum
489 of art. <https://CRAN.R-project.org/package=MetBrewer>
- 490 Nalborczyk, L. (2025). *neurogam*: Precise temporal localisation of m/EEG effects with bayesian
491

- 495 generalised additive multilevel models. <https://github.com/lmalborczyk/neurogam>
- 496 Nalborczyk, L., Batailler, C., Loevenbruck, H., Vilain, A., & Bürkner, P.-C. (2019). An In-
497 troduction to Bayesian Multilevel Models Using brms: A Case Study of Gender Effects
498 on Vowel Variability in Standard Indonesian. *Journal of Speech, Language, and Hearing
499 Research*, 62(5), 1225–1242. https://doi.org/10.1044/2018_jslhr-s-18-0006
- 500 Nalborczyk, L., Hauw, F., Torcy, H. de, Dehaene, S., & Cohen, L. (in preparation). *Neural and
501 representational dynamics of tickertape synesthesia*.
- 502 Pedersen, E. J., Miller, D. L., Simpson, G. L., & Ross, N. (2019). Hierarchical generalized
503 additive models in ecology: An introduction with mgcv. *PeerJ*, 7, e6876. <https://doi.org/10.7717/peerj.6876>
- 504 Pedersen, T. L. (2024). *patchwork: The composer of plots*. <https://CRAN.R-project.org/package=patchwork>
- 505 Pedersen, T. L., & Crameri, F. (2023). *scico: Colour palettes based on the scientific colour-maps*.
506 <https://CRAN.R-project.org/package=scico>
- 507 Pernet, C. R. (2022). Electroencephalography robust statistical linear modelling using a single
508 weight per trial. *Aperture Neuro*, 2, 1–22. <https://doi.org/10.52294/apertureneuro.2022.2. seoo9435>
- 509 Pernet, C. R., Chauveau, N., Gaspar, C., & Rousselet, G. A. (2011). LIMO EEG: A Tool-
510 box for Hierarchical LInear MOdeling of ElectroEncephaloGraphic Data. *Computational
511 Intelligence and Neuroscience*, 2011, 1–11. <https://doi.org/10.1155/2011/831409>
- 512 Pernet, C. R., Latinus, M., Nichols, T. E., & Rousselet, G. A. (2015). Cluster-based com-
513 putational methods for mass univariate analyses of event-related brain potentials/fields: A
514 simulation study. *Journal of Neuroscience Methods*, 250, 85–93. <https://doi.org/10.1016/j.jneumeth.2014.08.003>
- 515 R Core Team. (2025). *R: A language and environment for statistical computing*. R Foundation
516 for Statistical Computing. <https://www.R-project.org/>
- 517 Rasmussen, C. E., & Williams, C. K. I. (2005). *Gaussian Processes for Machine Learning*.
518 <https://doi.org/10.7551/mitpress/3206.001.0001>
- 519 Rigby, R. A., & Stasinopoulos, D. M. (2005). Generalized Additive Models for Location, Scale
520 and Shape. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 54(3),
521 507–554. <https://doi.org/10.1111/j.1467-9876.2005.00510.x>
- 522 Rij, J. van, Hendriks, P., Rijn, H. van, Baayen, R. H., & Wood, S. N. (2019). Analyzing
523 the Time Course of Pupillometric Data. *Trends in Hearing*, 23. <https://doi.org/10.1177/2331216519832483>
- 524 Riutort-Mayol, G., Bürkner, P.-C., Andersen, M. R., Solin, A., & Vehtari, A. (2023). Practi-
525 cal Hilbert space approximate Bayesian Gaussian processes for probabilistic programming.
526 *Statistics and Computing*, 33(1), 17. <https://doi.org/10.1007/s11222-022-10167-2>
- 527 Rodriguez-Sanchez, F., & Jackson, C. P. (2024). *grateful: Facilitate citation of R packages*.
528 <https://pakillo.github.io/grateful/>
- 529 Rosenblatt, J. D., Finos, L., Weeda, W. D., Solari, A., & Goeman, J. J. (2018). All-
530 Resolutions Inference for brain imaging. *NeuroImage*, 181, 786–796. <https://doi.org/10.1016/j.neuroimage.2018.07.060>

- 537 Rousselet, G. A. (2025). Using cluster-based permutation tests to estimate MEG/EEG onsets:
538 How bad is it? *European Journal of Neuroscience*, 61(1), e16618. <https://doi.org/10.1111/ejn.16618>
- 540 Rousselet, G. A., Pernet, C. R., Bennett, P. J., & Sekuler, A. B. (2008). Parametric study
541 of EEG sensitivity to phase noise during face processing. *BMC Neuroscience*, 9(1). <https://doi.org/10.1186/1471-2202-9-98>
- 543 Sassenhagen, J., & Draschkow, D. (2019). Cluster-based permutation tests of MEG/EEG data
544 do not establish significance of effect latency or location. *Psychophysiology*, 56(6). <https://doi.org/10.1111/psyp.13335>
- 546 Silge, J., & Robinson, D. (2016). tidytext: Text mining and analysis using tidy data principles
547 in r. *JOSS*, 1(3). <https://doi.org/10.21105/joss.00037>
- 548 Skukies, R., & Ehinger, B. (2021). Modelling event duration and overlap during EEG analysis.
549 *Journal of Vision*, 21(9), 2037. <https://doi.org/10.1167/jov.21.9.2037>
- 550 Skukies, R., Schepers, J., & Ehinger, B. (2024, December 9). *Brain responses vary in duration*
551 - modeling strategies and challenges. <https://doi.org/10.1101/2024.12.05.626938>
- 552 Slowikowski, K. (2024). *ggrepel: Automatically position non-overlapping text labels with “gg-*
553 *plot2”*. <https://CRAN.R-project.org/package=ggrepel>
- 554 Smith, N. J., & Kutas, M. (2014a). Regression-based estimation of ERP waveforms: I. The
555 rERP framework. *Psychophysiology*, 52(2), 157–168. <https://doi.org/10.1111/psyp.12317>
- 556 Smith, N. J., & Kutas, M. (2014b). Regression-based estimation of ERP waveforms: II. Non-
557 linear effects, overlap correction, and practical considerations. *Psychophysiology*, 52(2),
558 169–181. <https://doi.org/10.1111/psyp.12320>
- 559 Smith, S., & Nichols, T. (2009). Threshold-free cluster enhancement: Addressing problems of
560 smoothing, threshold dependence and localisation in cluster inference. *NeuroImage*, 44(1),
561 83–98. <https://doi.org/10.1016/j.neuroimage.2008.03.061>
- 562 Sóskuthy, M. (2017). *Generalised additive mixed models for dynamic analysis in linguistics: A*
563 *practical introduction*. <https://doi.org/10.48550/ARXIV.1703.05339>
- 564 Sóskuthy, M. (2021). Evaluating generalised additive mixed modelling strategies for dynamic
565 speech analysis. *Journal of Phonetics*, 84, 101017. <https://doi.org/10.1016/j.wocn.2020.101017>
- 566 Teichmann, L. (2022). An empirically driven guide on using bayes factors for m/EEG decoding.
567 *Aperture Neuro*, 2, 1–10. <https://doi.org/10.52294/apertureneuro.2022.2.maoc6465>
- 569 Tremblay, A., & Newman, A. J. (2014). Modeling nonlinear relationships in ERP data using
570 mixed-effects regression with R examples. *Psychophysiology*, 52(1), 124–139. <https://doi.org/10.1111/psyp.12299>
- 572 Umlauf, N., Klein, N., & Zeileis, A. (2018). BAMLSS: Bayesian Additive Models for Location,
573 Scale, and Shape (and Beyond). *Journal of Computational and Graphical Statistics*, 27(3),
574 612–627. <https://doi.org/10.1080/10618600.2017.1407325>
- 575 Ushey, K., Allaire, J., & Tang, Y. (2024). *Reticulate: Interface to ‘python’*. <https://CRAN.R-project.org/package=reticulate>
- 577 Vaughan, D., & Dancho, M. (2022). *furrr: Apply mapping functions in parallel using futures*.
578 <https://CRAN.R-project.org/package=furrr>

- 579 Wickham, H. (2019). *assertthat: Easy pre and post assertions.* [https://CRAN.R-project.org/
package=assertthat](https://CRAN.R-project.org/package=assertthat)
- 580
- 581 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Grolemund,
582 G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L., Miller, E., Bache, S. M.,
583 Müller, K., Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., ... Yutani, H. (2019). Welcome
584 to the tidyverse. *Journal of Open Source Software*, 4(43), 1686. [https://doi.org/10.21105/
joss.01686](https://doi.org/10.21105/joss.01686)
- 585
- 586 Wickham, H., Pedersen, T. L., & Seidel, D. (2023). *scales: Scale functions for visualization.*
587 <https://CRAN.R-project.org/package=scales>
- 588 Wieling, M. (2018). Analyzing dynamic phonetic data using generalized additive mixed model-
589 ing: A tutorial focusing on articulatory differences between L1 and L2 speakers of English.
590 *Journal of Phonetics*, 70, 86–116. <https://doi.org/10.1016/j.wocn.2018.03.002>
- 591 Winter, B., & Wieling, M. (2016). How to analyze linguistic change using mixed models, Growth
592 Curve Analysis and Generalized Additive Modeling. *Journal of Language Evolution*, 1(1),
593 7–18. <https://doi.org/10.1093/jole/lzv003>
- 594 Wood, S. N. (2003). Thin Plate Regression Splines. *Journal of the Royal Statistical Society
Series B: Statistical Methodology*, 65(1), 95–114. <https://doi.org/10.1111/1467-9868.00374>
- 595
- 596 Wood, S. N. (2017a). *Generalized Additive Models*. Chapman; Hall/CRC. <https://doi.org/10.1201/9781315370279>
- 597
- 598 Wood, S. N. (2017b). *Generalized additive models: An introduction with r* (2nd ed.). Chapman;
599 Hall/CRC.
- 600 Wüllhorst, V., Wüllhorst, R., Overmeyer, R., & Endrass, T. (2025). Comprehensive Analysis
601 of Event-Related Potentials of Response Inhibition: The Role of Negative Urgency and
602 Compulsivity. *Psychophysiology*, 62(2). <https://doi.org/10.1111/psyp.70000>
- 603 Xie, Y. (2014). knitr: A comprehensive tool for reproducible research in R. In V. Stodden, F.
604 Leisch, & R. D. Peng (Eds.), *Implementing reproducible computational research*. Chapman;
605 Hall/CRC.
- 606 Xie, Y. (2015). *Dynamic documents with R and knitr* (2nd ed.). Chapman; Hall/CRC. <https://yihui.org/knitr/>
- 607
- 608 Xie, Y. (2025). knitr: A general-purpose package for dynamic report generation in R. <https://yihui.org/knitr/>
- 609
- 610 Xie, Y., Allaire, J. J., & Grolemund, G. (2018). *R markdown: The definitive guide*. Chapman;
611 Hall/CRC. <https://bookdown.org/yihui/rmarkdown>
- 612 Xie, Y., Dervieux, C., & Riederer, E. (2020). *R markdown cookbook*. Chapman; Hall/CRC.
613 <https://bookdown.org/yihui/rmarkdown-cookbook>
- 614 Yeung, N., Bogacz, R., Holroyd, C. B., & Cohen, J. D. (2004). Detection of synchronized os-
615 cillations in the electroencephalogram: An evaluation of methods. *Psychophysiology*, 41(6),
616 822–832. <https://doi.org/10.1111/j.1469-8986.2004.00239.x>

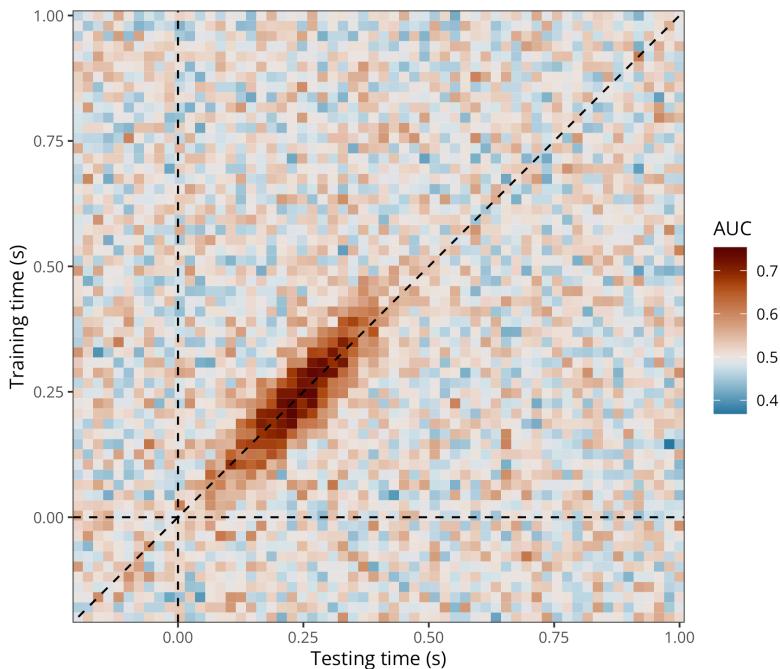
Appendix A

Application to 2D time-resolved decoding results (cross-temporal generalisation)

617 Assume we have M/EEG data and we have conducted cross-temporal generalisation analyses
 618 ([King & Dehaene, 2014](#)). As a result, we have a 2D matrix where each element contains the
 619 decoding accuracy (e.g., ROC AUC) of a classifier trained at timestep training_i and tested at
 620 timestep testing_j (Figure A1).

Figure A1

Exemplary (simulated) group-level average cross-temporal generalisation matrix of decoding performance (ROC AUC).



621 To model cross-temporal generalisation matrices of decoding performance (ROC AUC),
 622 we extended the initial (decoding) GAM to take into account the bivariate temporal distribution
 623 of AUC values, thus producing naturally smoothed estimates (timecourses) of AUC values and
 624 posterior probabilities. This model can be written as follows:

$$\begin{aligned} \text{AUC}_i &\sim \text{Beta}(\mu_i, \phi) \\ g(\mu_i) &= f(\text{train}_i, \text{test}_i) \end{aligned}$$

625 where we assume that AUC values come from a Beta distribution with two parameters
 626 μ and ϕ . We can think of $f(\text{train}_i, \text{test}_i)$ as a surface (a smooth function of two variables) that
 627 we can model using a 2-dimensional splines. Let $\mathbf{s}_i = (\text{train}_i, \text{test}_i)$ be some pair of training and
 628 testing samples, and let $\mathbf{k}_m = (\text{train}_m, \text{test}_m)$ denote the m^{th} knot in the domain of train_i and
 629 test_i . We can then express the smooth function as:

$$f(\text{train}_i, \text{test}_i) = \alpha + \sum_{m=1}^M \beta_m b_m(\tilde{s}_i, \tilde{k}_m)$$

630 Note that $b_m(,)$ is a basis function that maps $R \times R \rightarrow R$. A popular bivariate basis

631 function uses *thin-plate splines* (Wood, 2003), which extend to $\mathbf{s}_i \in \mathbb{R}^d$ and ∂l_g penalties. These
 632 splines are designed to interpolate and approximate smooth surfaces over two dimensions (hence
 633 the “bivariate” term). For $d = 2$ dimensions and $l = 2$ (smoothness penalty involving second
 634 order derivative):

$$f(\tilde{s}_i) = \alpha + \beta_1 x_i + \beta_2 z_i + \sum_{m=1}^M \beta_{2+m} b_m(\tilde{s}_i, \tilde{k}_m)$$

635 using the the radial basis function given by:

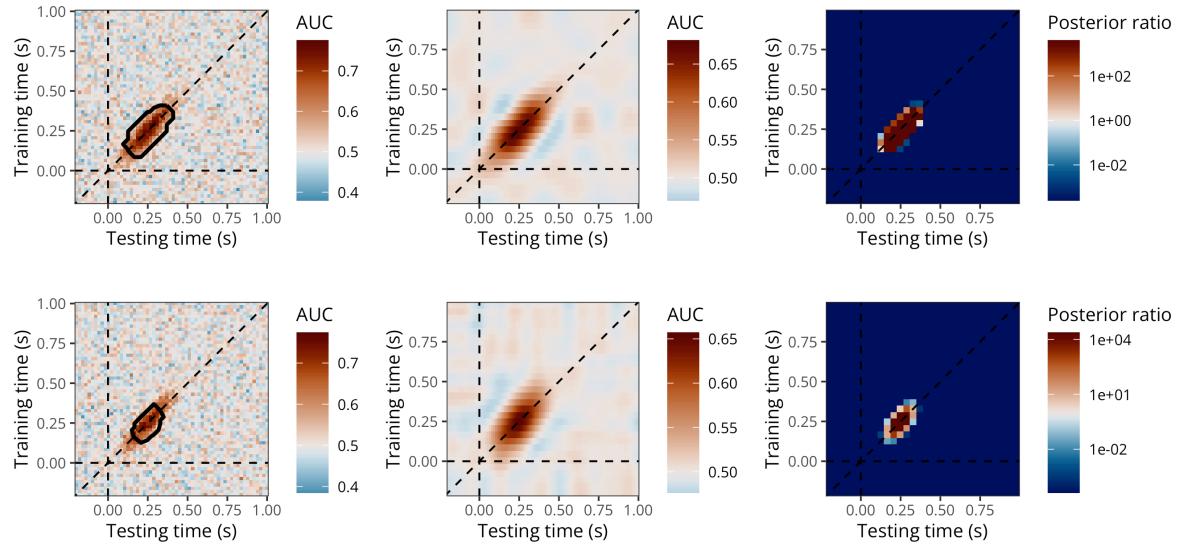
$$b_m(\tilde{s}_i, \tilde{k}_m) = \|\tilde{s}_i - \tilde{k}_m\|^2 \log \|\tilde{s}_i - \tilde{k}_m\|$$

636 where $\|\mathbf{s}_i - \mathbf{k}_m\|$ is the Euclidean distance between the covariate \mathbf{s}_i and the knot location
 637 \mathbf{k}_m . We fitted this model using `brms...`

```
# fitting a GAM with two temporal dimensions
timegen_gam <- brm(
  # 2D thin-plate spline (tp)
  # auc ~ t2(train_time, test_time, bs = "tp", k = 10),
  auc ~ t2(train_time, test_time, bs = "tp", k = 20),
  data = timegen_data,
  family = Beta(),
  warmup = 2000,
  iter = 5000,
  chains = 8,
  cores = 8,
  # file = "models/timegen_gam_t2.rds" # k = 10
  file = "models/timegen_gam_t2_k20.rds" # k = 20
)
```

Figure A2

Simulated data with thresholds derived from the BGAM (left), predicted AUC values (middle), and ratios of posterior probabilities of decoding accuracy being above chance level (right) according to the bivariate BGAM (for $k=10$, top row, and $k=20$, bottom row).



638 Could be extended to spatial and temporal dimensions with formulas such as `te(x, y,`
 639 `Time, d = c(2, 1) ...`

Appendix B

Alternative to GAMs: Approximate Gaussian Process regression

640 A Gaussian process (GP) is a stochastic process that defines the distribution over a collection
 641 of random variables indexed by a continuous variable, that is $\{f(t) : t \in \mathcal{T}\}$ for some index
 642 set \mathcal{T} (Rasmussen & Williams, 2005; Riutort-Mayol et al., 2023). Whereas Bayesian linear
 643 regression outputs a distribution over the parameters of some predefined parametric model, the
 644 GP approach, in contrast, is a non-parametric approach, in that it finds a distribution over the
 645 possible functions that are consistent with the observed data. However, note that nonparametric
 646 does not mean there aren't parameters, it means that there are infinitely many parameters.

647 From brms documentation: A GP is a stochastic process, which describes the relation
 648 between one or more predictors $x = (x_1, \dots, x_d)$ and a response $f(x)$, where d is the number
 649 of predictors. A GP is the generalization of the multivariate normal distribution to an infinite
 650 number of dimensions. Thus, it can be interpreted as a prior over functions. The values of $f()$
 651 at any finite set of locations are jointly multivariate normal, with a covariance matrix defined
 652 by the covariance kernel $k_p(x_i, x_j)$, where p is the vector of parameters of the GP:

$$(f(x_1), \dots, f(x_n)) \sim \text{MVN}\left(0, (k_p(x_i, x_j))_{i,j=1}^n\right)$$

653 The smoothness and general behaviour of the function f depends only on the choice of
 654 covariance kernel, which ensures that values that are close together in the input space will be
 655 mapped to similar output values...

656 From this perspective, f is a realisation of an infinite dimensional normal distribution:

$$f \sim \text{Normal}(0, C(\lambda))$$

657 where C is a covariance kernel with hyperparameters λ that defines the covariance
 658 between two function values $f(t_1)$ and $f(t_2)$ for two time points t_1 and t_2 (Rasmussen &
 659 Williams, 2005). Similar to the different choices of the basis function for splines, different
 660 choices of the covariance kernel lead to different GPs. In this article, we consider the squared-
 661 exponential (a.k.a. radial basis function) kernel, which computes the squared distance between
 662 points and converts it into a measure of similarity. It is defined as:

$$C(\lambda) := C(t_1, t_2, \sigma, \gamma) := \sigma^2 \exp\left(-\frac{\|t_1 - t_2\|^2}{2\gamma^2}\right)$$

663 with hyperparameters $\lambda = (\sigma, \gamma)$, expressing the overall scale of GP and the length-
 664 scale, respectively (Rasmussen & Williams, 2005). The advantages of this kernel are that it is
 665 computationally efficient and (infinitely) smooth making it a reasonable choice for the purposes
 666 of the present article. Here again, λ hyperparameters are estimated from the data, along with
 667 all other model parameters.

668 Taken from <https://michael-franke.github.io/Bayesian-Regression/practice-sheets/>
 669 10c-Gaussian-processes.html: For a given vector \mathbf{x} , we can use the kernel to construct finite
 670 multi-variate normal distribution associated with it like so:

$$\mathbf{x} \mapsto_{GP} \text{MVNormal}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}))$$

671 where m is a function that specifies the mean for the distribution associated with \mathbf{x} . This
672 mapping is essentially the Gaussian process: a systematic association of vectors of arbitrary
673 length with a suitable multi-variate normal distribution.

674 Low-rank approximate Gaussian processes are of main interest in machine learning and
675 statistics due to the high computational demands of exact Gaussian process models ([Riutort-](#)
676 [Mayol et al., 2023](#))...

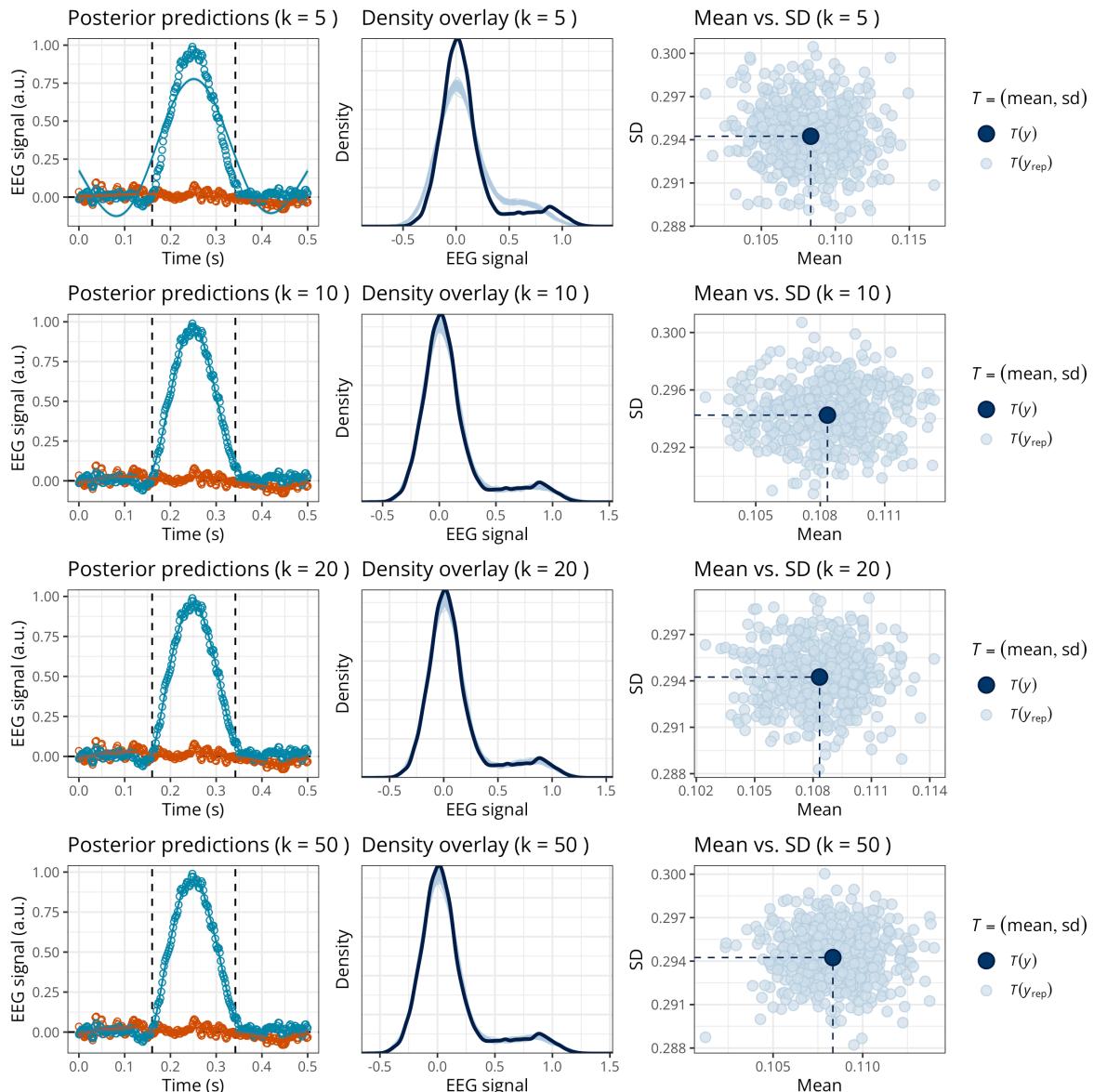
Appendix C

How to choose the GAM basis dimension?

677 Here provide recommendation about how to define k . An option is to vary k and examine the
 678 predictions and posterior predictive checks (PPCs) of each model... In this example (Figure C1)...
 679 However, it is not possible to provide general recommendations, as the optimal k depends on
 680 the sampling rate, the preprocessing steps (e.g., signal-to-noise ratio, low-pass filtering, etc),
 681 and the neural dynamics of the phenomenon under study.

Figure C1

Posterior predictions and posterior predictive checks for the GAM with varying k (in rows).



Appendix D

R package and integration with MNE-Python

682 For users who are already familiar with `brms`, the recommended pipeline is to import ERPs
 683 or decoding results in R and analyse these data using the code provided in the main paper.
 684 However, it is also possible to call functions from the `neurogam` R package (available at <https://github.com/lNALBORCZYK/neurogam>), which come with sensible defaults.
 685

```
# installing (if needed) and loading the neurogam R package
# remotes::install_github("https://github.com/lNALBORCZYK/neurogam")
library(neurogam)

# using the testing_through_time() function from the neurogam package
# this may take a few minutes (or hours depending the machine's
# performance and data size)...
gam_onset_offset <- testing_through_time(
  # dataframe with M/EEG data in long format
  data = raw_df,
  # threshold for defining clusters (20 by default)
  threshold = 20,
  # the *_id arguments are used to specify the relevant columns in data
  participant_id = "participant", meeg_id = "eeg",
  time_id = "time", predictor_id = "condition",
  # number of warmup MCMC iterations
  warmup = 1000,
  # total number of MCMC iterations
  iter = 5000,
  # number of MCMCs
  chains = 4,
  # number of parallel cores to use for running the MCMCs
  cores = 4
)

# displaying the results
gam_onset_offset$clusters
```

686 The `neurogam` package can also be called from Python using the `rpy2` module, and can
 687 easily be integrated into MNE-Python pipelines. For example, we use it below to estimate the
 688 onset and offset of effects for one EEG channel from a MNE evoked object. Note that the code
 689 used to reshape the `sample` MNE dataset is available in the online supplementary materials, and
 690 we refer to the [MNE documentation](#) about converting MNE epochs to Pandas dataframes in long
 691 format (i.e., with one observation per row).

```
# loading the Python modules
import rpy2.robj as robj
from rpy2.robj.packages import importr
from rpy2.robj import pandas2ri
from rpy2.robj.conversion import localconverter

# importing the "neurogam" R package
neurogam = importr("neurogam")

# activating automatic pandas-R conversion
pandas2ri.activate()

# assuming reshaped_df is some M/EEG data reshaped in long format
with localconverter(robj.default_converter + pandas2ri.converter):

    reshaped_df_r = robj.conversion.py2rpy(reshaped_df)

# using the testing_through_time() function from the neurogam R package
gam_onset_offset = neurogam.testing_through_time(
    data=reshaped_df_r,
    threshold=10,
    multilevel=False
)

# displaying the results
print(list(gam_onset_offset) )
```