

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

³Ladislas Nalborczyk¹ and Paul Bürkner²

⁴¹Aix Marseille Univ, CNRS, LPL

⁵²TU Dortmund University, Department of Statistics

⁶Abstract

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 method leverages Bayesian generalised additive multilevel models (GAMMs), 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 show how it can be seamlessly 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

⁸Table of contents

⁹ Introduction	³
¹⁰ Introduction	³

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

11	Problem statement	3
12	Cluster-based inference	3
13	Previous work on modelling M/EEG data	4
14	Generalised additive models	5
15	Bayesian generalised additive multilevel models	5
16	Objectives	6
17	Methods	7
18	M/EEG data simulation	7
19	Model description and model fitting	7
20	Error properties of the proposed approach	10
21	Comparing the identified onsets/offsets to other approaches	11
22	Simulation study	12
23	Application to actual MEG data	12
24	Results	14
25	Simulation study (bias and variance)	14
26	Application to actual MEG data (reliability)	15
27	Discussion	16
28	Summary of the proposed approach	16
29	Limitations and future directions	17
30	Data and code availability	18
31	Packages	18
32	Acknowledgements	18
33	References	19
34	Application to 2D time-resolved decoding results (cross-temporal generalisation)	24
35	Alternative to GAMs: Approximate Gaussian Process regression	27
36	How to choose the GAM basis dimension?	28
37	R package and integration with MNE-Python	29

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 with a unique
11 threshold, followed by corrections for multiple comparisons with the goal of maintaining the
12 familywise error rate (FWER) or false discovery rate (FDR) at the nominal level. Cluster-based
13 inference is the most common way of achieving this sort of error control in the M/EEG literature,
14 being the default option in several software programs. While effective for controlling error rates,
15 cluster-based inference comes with a significant limitation: they shift the focus of inference from
16 individual time points to clusters, making it difficult to draw precise conclusions about the onset
17 or offset of observed effects (Maris & Oostenveld, 2007; Sassenhagen & Draschkow, 2019). As
18 pointed by Maris & Oostenveld (2007); “there is a conflict between this interest in localized effects
19 and our choice for a global null hypothesis: by controlling the FA [false alarm] rate under this
20 global null hypothesis, one cannot quantify the uncertainty in the spatiotemporal localization
21 of the effect”. Even worse, as noted by Rosenblatt et al. (2018), cluster-based inference suffers
22 from low spatial resolution: “Since discovering a cluster means that ‘there exists at least one
23 voxel with an evoked response in the cluster’, and not that ‘all the voxels in the cluster have an
24 evoked response’, it follows that the larger the detected cluster, the less information we have on
25 the location of the activation.” As a consequence, cluster-based inferenc has been shown...

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

34 **Cluster-based inference**

35 Different methods exist to control the family-wiser error rate (FWER), defined as the
36 type-1 error rate (false positive) over an ensemble (family) of tests... for instance the Bonferroni
37 correction (Dunn, 1961), however this method is generally overconservative as it assumes sta-
38 tistical independence of tests, an assumptions that is clearly violated in the context of M/EEG
39 timeseries with spatiotemporal dependencies... or the false discovery rate (FDR), defined as the
40 proportion of false positive among positive tests (e.g., Benjamini & Hochberg, 1995; Benjamini
41 & Yekutieli, 2001)...

42 On popular technique to account for spatiotemporal dependencies while controlling the
43 FWER is cluster-based inference. Description of cluster-based approaches (see Sassenhagen
44 & Draschkow, 2019)... A typical cluster-based inference consists of two successive steps. First,
45 clusters are defined as sets of contiguous voxels, channels, and/or timesteps, whose intensity/ac-
46 tivity exceeds some predefined threshold. Clusters are then characterised by their height (i.e.,

maximal value), extent (number of constituent elements), or some combination of both, such at its “mass”, for instance by summing the statistics within a cluster, an approach refereed to as “cluster mass” (Maris & Oostenveld, 2007; Pernet et al., 2015). Then, the null hypothesis is tested by assessing the probability.. As different cluster-forming thresholds lead to clusters with different spatial or temporal extent, this initial threshold modulates the sensitivity of the subsequent permutation test. The threshold-free cluster enhancement method (TFCE) was introduced by S. Smith & Nichols (2009) to overcome this choice of an arbitrary threshold.

In brief, the TFCE method works as follows. Instead of picking an arbitrary cluster-forming threshold (e.g., $t = 2$), we try all (or many) possible thresholds in a given range and check whether a given timestep/voxel belongs to a significant cluster under any of the set of thresholds. Then, instead of using cluster mass (e.g., the sum of squared t-values within the cluster), we use a weighted average between the cluster extend (e , how broad is the cluster, that is, how many connected samples it contains) and the cluster height (h , how high is the cluster, that is, how large is the test statistic). The TFCE score at each timestep/voxel t is given by:

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

where h_0 is typically 0 and parameters E and H are set a priori (typically to 0.5 and 2, respectively) and control the influence of the extend and height on the TFCE. Then, p-value for timestep/voxel t is computed by comparing it TFCE with the null distribution of TFCE values. For each permuted signal, we keep the maximal value over the whole signal for the null distribution of the TFCE. The TFCE combined with permutation (assuming a large enough number of permutations) has been shown to provide accurate type 1 FWER (e.g., Pernet et al., 2015). However, previous simulation work showed that cluster-based methods have indeed poor performance in localising the onset of M/EEG effects (e.g., Rousselet, 2025; Sassenhagen & Draschkow, 2019).

To sum up this section, cluster-based inference main limitations is that it provides inference at the cluster level only, not allowing inference at the level of timesteps, sensors, etc. Thus, it does not allow inferring the localisation of effects (Maris, 2011; for more details on cluster-based inference, see for instance Maris & Oostenveld, 2007; Sassenhagen & Draschkow, 2019). In the following, we briefly review previous modelling work. Then, we provide a brief introduction to Bayesian statistical modelling, before moving to generalised additive models (GAMs) and generalised additive multilevel models (GAMMs) to illustrate how they can be used to precisely localise the onset and offset of M/EEG effects.

Previous work on modelling M/EEG data

Recent example of GLM for EEG (Fischer & Ullsperger, 2013; Wüllhorst et al., 2025)...

See also (Hauk et al., 2006; Rousselet et al., 2008)... Example of two-stage regression analysis (i.e., individual-level then group-level, Dunagan et al., 2024)...

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

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

From Dimigen & Ehinger (2021): Recently, spline regression has been applied to ERPs (e.g., Hendrix et al., 2017; Kryuchkova et al., 2012)... GAMMs for EEG data (Abugaber et al., 2023; Meulman et al., 2015, 2023)...

Disentangling overlapping processes (Skukies et al., 2024; Skukies & Ehinger, 2021)...

Weighting single trials (Pernet, 2022)... The LIMO toolbox (Pernet et al., 2011)...

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 “wiggliness” of the function. In other words, more complex (wiggly) basis function are penalised.

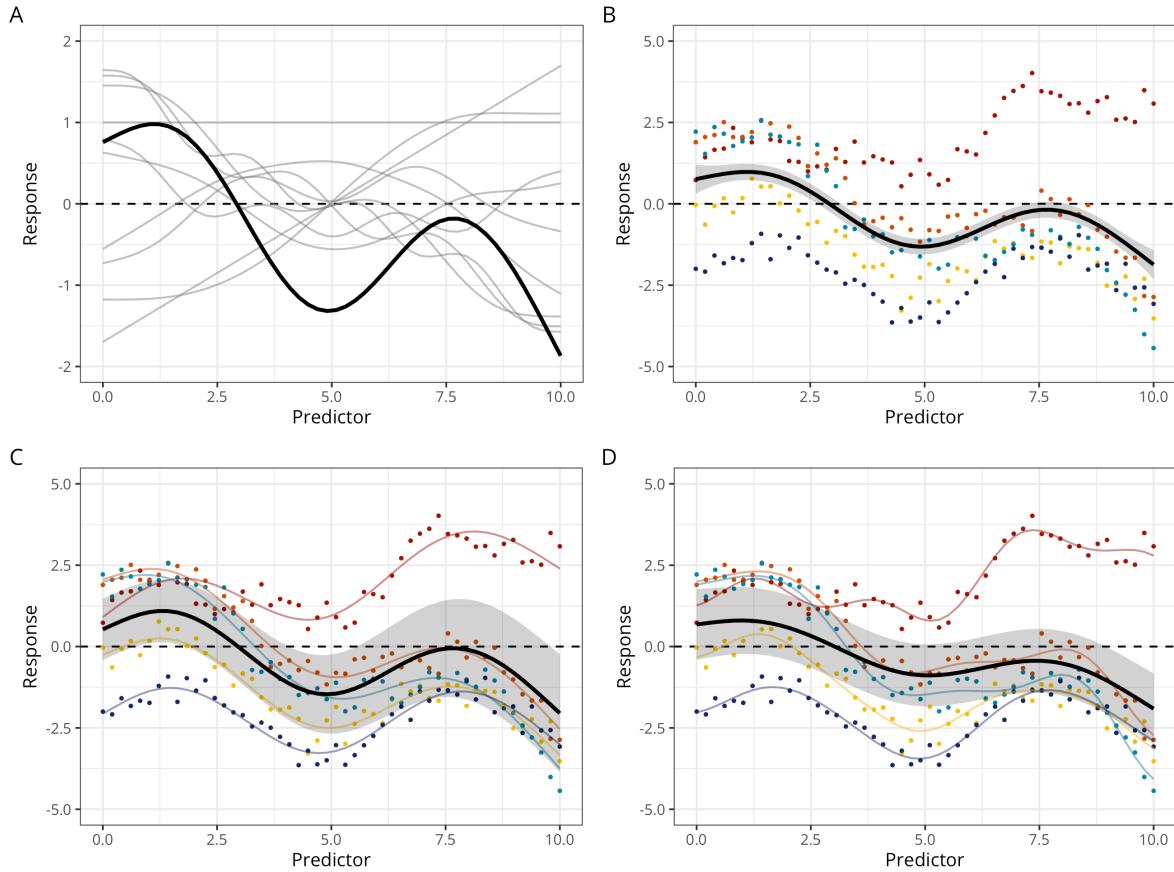
Bayesian generalised additive multilevel models

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 enclose subjective (personal) beliefs or... but these debates are outside the scope of the present paper and we therefore 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 ‘wiggliness’ (smoother).



See Figure 1... Introduction to multilevel GAMs ([E. J. Pedersen et al., 2019](#))... Now describe the Bayesian GAMM ([Miller, 2025](#))... Proper inclusion of varying/random effects in the model specification protects against overly wiggly curves ([Baayen & Linke, 2020](#))... Generalising to scale and shape or “distributional GAMs” ([Rigby & Stasinopoulos, 2005](#); [Umlauf et al., 2018](#)) and applied to neuroimaging data ([Dinga et al., 2021](#))...

Instead of averaging, obtain the smooth ERP signal from multilevel GAM... less susceptible to outliers ([Meulman et al., 2023](#))...

142 Objectives

Given the previously reported limitations of conventional methods to precisely identify the onset and offset of M/EEG effects (e.g., ERPs, decoding performance), we developed a model-based approach for estimating the onset and offset of such effects. To achieve this, we

¹⁴⁶ leveraged Bayesian generalised additive multilevel models (BGAMMs) fitted in R via the `brms`
¹⁴⁷ package and compared the performance of this approach to conventional methods on both
¹⁴⁸ simulated and actual M/EEG data.

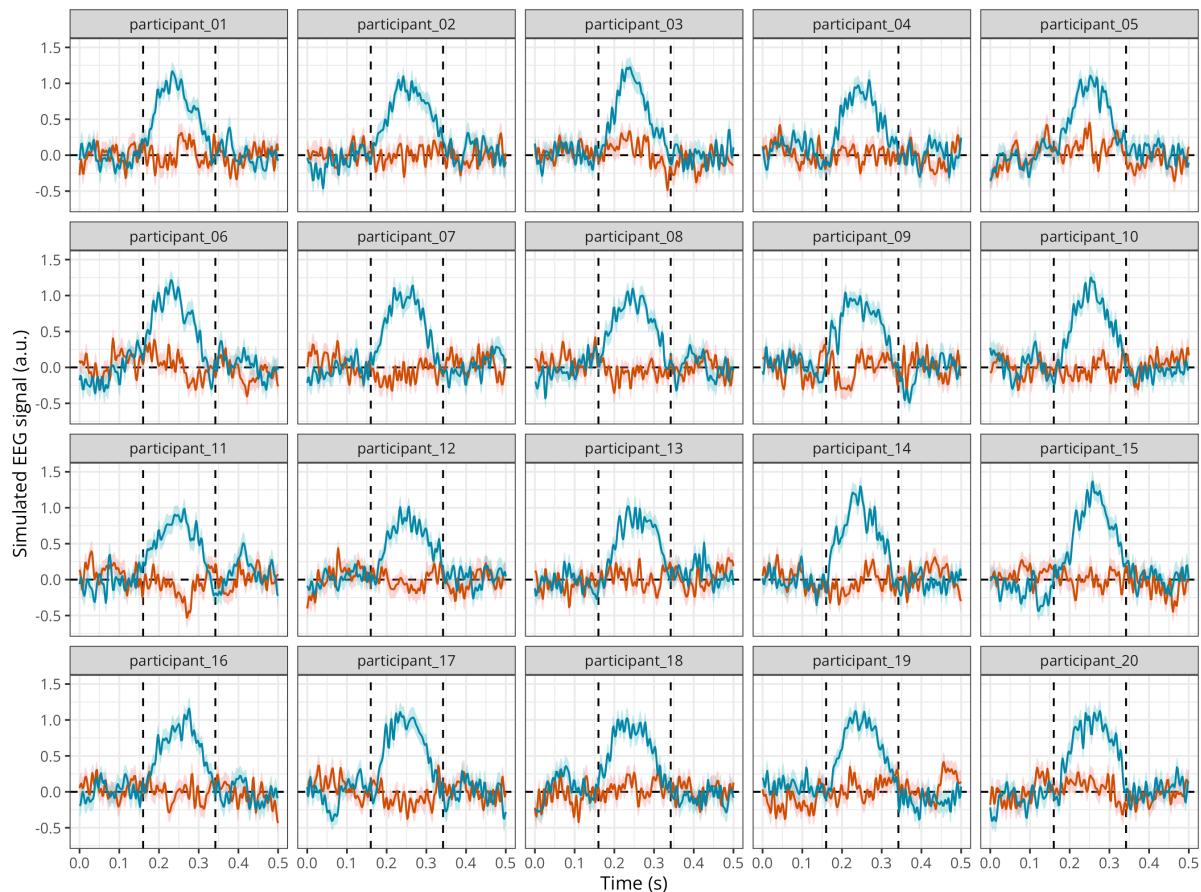
¹⁴⁹ **Methods**

¹⁵⁰ **M/EEG data simulation**

¹⁵¹ Following the approach of Sassenhagen & Draschkow (2019) and Rousselet (2025), we
¹⁵² simulated EEG data stemming from two conditions, one with noise only, and the other with
¹⁵³ noise + signal. As in previous studies, the noise was generated by superimposing 50 sinusoids
¹⁵⁴ at different frequencies, following an EEG-like spectrum (see code in the online supplementary
¹⁵⁵ materials and details in Yeung et al., 2004). As in Rousselet (2025), the signal was generated
¹⁵⁶ from a truncated Gaussian distribution with an objective onset at 160 ms, a peak at 250 ms,
¹⁵⁷ and an offset at 342 ms. We simulated this signal for 250 timesteps between 0 and 0.5s, akin to
¹⁵⁸ a 500 Hz sampling rate. We simulated data for a group of 20 participants (with variable true
¹⁵⁹ onset) with 50 trials per participant and condition (Figure 2).

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.*



¹⁶⁰ **Model description and model fitting**

¹⁶¹ We then fitted a Bayesian GAM (BGAM) using the `brms` package (Bürkner, 2017,
¹⁶² 2018; Nalborczyk et al., 2019). We used the default priors in `brms` (i.e., weakly informative

¹⁶³ priors). We ran eight Markov Chain Monte-Carlo (MCMC) to approximate the posterior distribution, including each 5000 iterations and a warmup of 2000 iterations, yielding a total of $8 \times (5000 - 2000) = 24000$ posterior samples to use for inference. Posterior convergence was ¹⁶⁴ assessed examining trace plots as well as the Gelman–Rubin statistic \hat{R} (Gabry et al., 2019; ¹⁶⁵ Gelman et al., 2020). The `brms` package uses the same syntax as the R package `mgcv` v 1.9-¹⁶⁶ 1 (Wood, 2017b) for specifying smooth effects. Figure 3 shows the predictions of this model ¹⁶⁷ together with the raw data.

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

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

# fitting the BGAM
gam <- brm(
  # cubic regression splines with k-1 basis functions
  # eeg ~ condition + s(time, bs = "cr", k = 20, by = condition),
  # 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"
)
```

¹⁷⁰ However, the previous model only included constant (fixed) effects, thus not properly ¹⁷¹ accounting for between-participant variability. We next fit a multilevel version of the BGAM ¹⁷² (BGAMM). Although it is possible to fit a BGAMM at the single-trial level, we present a ¹⁷³ resource-efficient version of the model that is fitted directly on the by-participant summary ¹⁷⁴ 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
```

```

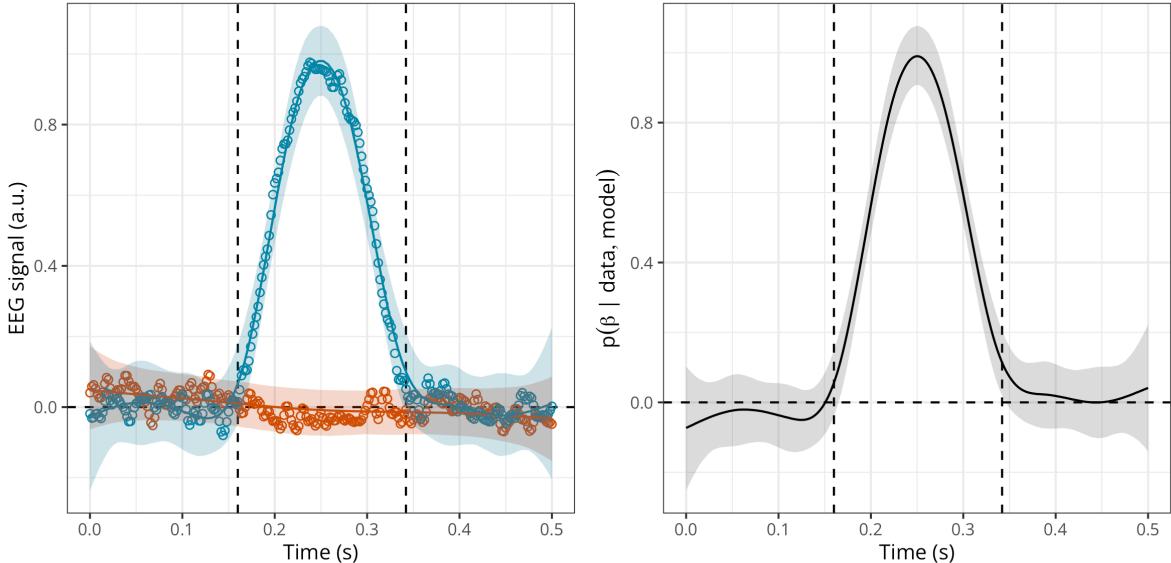
eeg_mean | se(eeg_sd) ~
  condition + s(time, bs = "cr", k = 20, by = condition) +
  (1 | participant),
data = summary_df,
family = gaussian(),
warmup = 2000,
iter = 5000,
chains = 8,
cores = 8,
file = "models/meta_gam.rds"
)

```

175 We depict the posterior predictions together with the posterior estimate of the slope for
 176 `condition` at each timestep (Figure 3). This figure suggests that the BGAMM provides an
 177 adequate description of the simulated data.

Figure 3

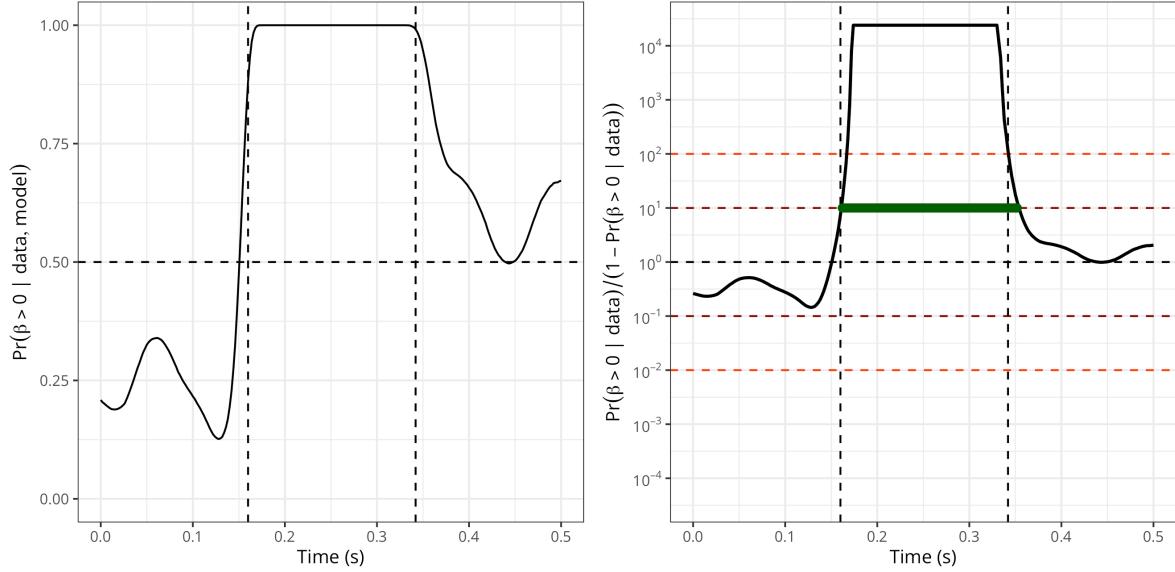
Posterior estimate of the EEG activity in each condition (left) and posterior estimate of the difference in EEG activity (right) according to the BGAMM.



178 We then compute the posterior probability of the slope for `condition` being above 0
 179 (Figure 4, left). We can also express this as the ratio of posterior probabilities (i.e., $p/(1 - p)$)
 180 and visualise the timecourse of this ratio superimposed with the conventional thresholds on
 181 evidence ratios (Figure 4, right). Note that a ratio of 10 means that the probability of the
 182 difference being above 0 is 10 times higher than the probability of the difference not being
 183 above 0, given the data, the priors, and other model's assumptions. Thresholding the posterior
 184 probability ratio thus provides a model-based approach for estimating the onset and offset of
 185 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.

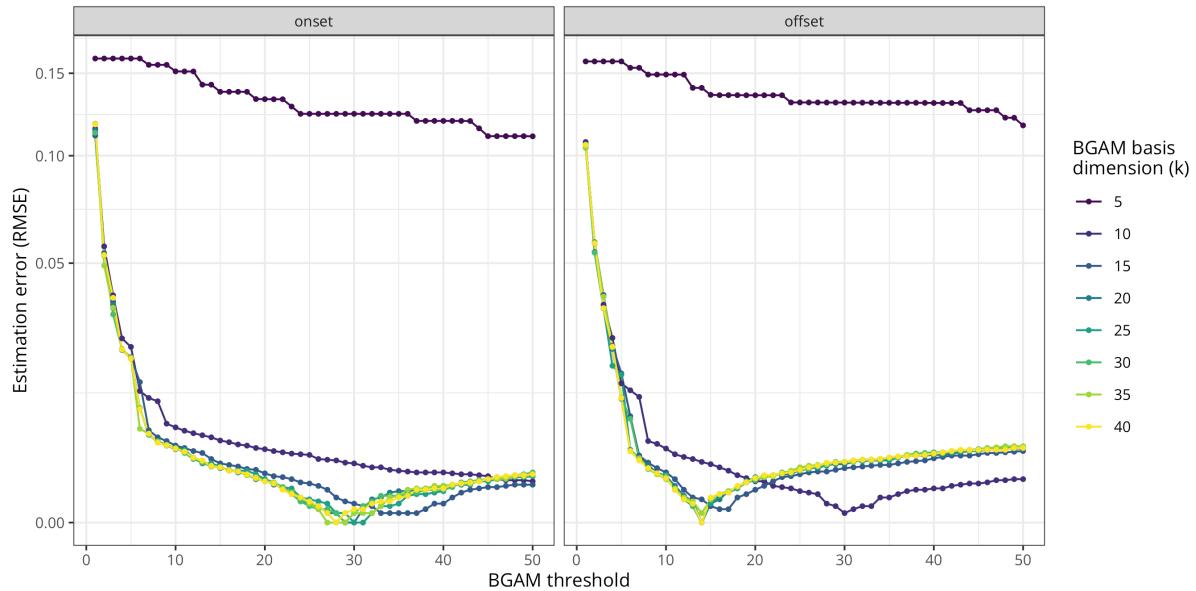


186 Error properties of the proposed approach

187 We then assess the performance of the proposed approach by computing the differ-
 188 ence between the true and estimated onset/offset of the EEG difference according to various
 189 k (BGAM basis dimension) and **threshold** values. Remember that the EEG signal was gen-
 190 erated from a truncated Gaussian with an objective onset at 160 ms, a maximum at 250 ms,
 191 and an offset at 342 ms. Figure 5 shows that the multilevel GAM can almost exactly recover
 192 the true onset and offset values, given some reasonable choice of k and **threshold** values. We
 193 provide more detailed recommendations on how to set k in Section C. This figure further reveals
 194 that the optimal k and **threshold** values may differ for the onset and offset values, and there
 195 there seems to exist a trade-off between these two parameters: lower k values lead to poorer
 196 estimations, but these poor estimations can be compensated (only to some extent) by higher
 197 **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).

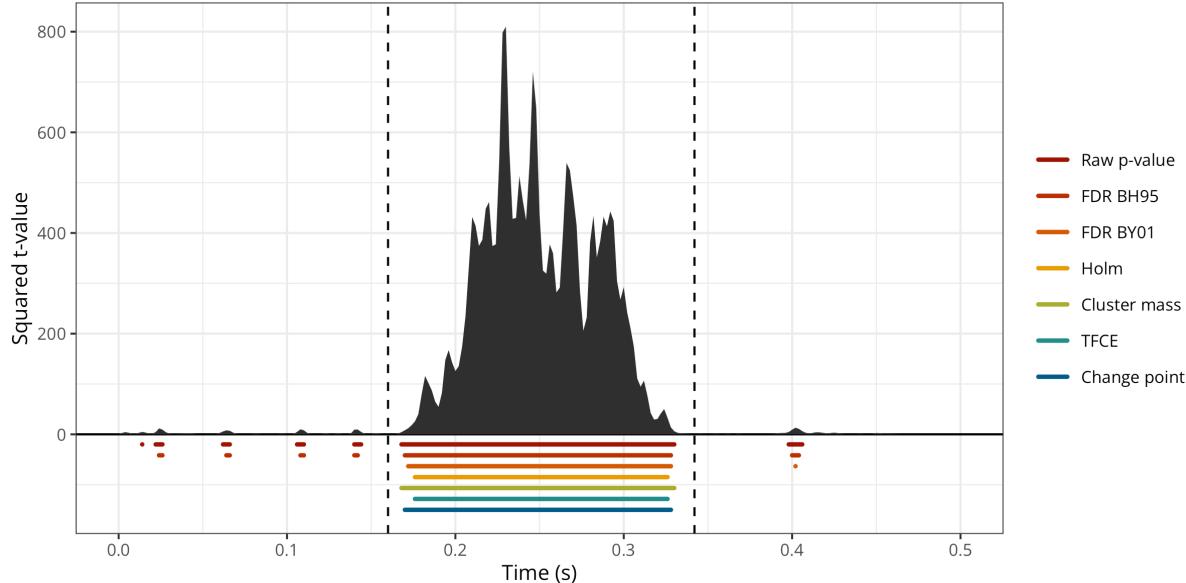


198 Comparing the identified onsets/offsets to other approaches

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



218 Simulation study

219 To assess the accuracy of group-level onset estimation, the various methods were com-
 220 pared using the bias (i.e., median(estimated-true)), median absolute error (MAE), root mean
 221 square error (RMSE), variance, and median absolute deviation (MAD) of onset/offset estimates
 222 computed on 1000 simulated datasets. As in Rousselet (2025), each participant was assigned a
 223 random onset between 150 and 170ms.

224 Application to actual MEG data

225 To complement the simulation study, we evaluated the performance of the various meth-
 226ods on actual MEG data (decoding results from Nalborczyk et al., in preparation). In this study,
 227 we conducted time-resolved multivariate pattern analysis (MVPA, also known as decoding) of
 228 MEG data during reading tasks. As a result, we obtain a timecourse of decoding performance
 229 (ROC AUC), bounded between 0 and 1, for each participant (for a total of 32 participants).
 230 Next, we wanted to test whether the group-level average decoding accuracy is above chance (i.e.,
 231 0.5) at each timestep (Figure 7). To achieve this, we fitted a BGAM as introduced previously,
 232 but we replaced the Normal likelihood function by a Beta one to account for the bounded nature
 233 of AUC values (between 0 and 1) (for a tutorial on Beta regression, see Coretta & Bürkner,
 234 2025).

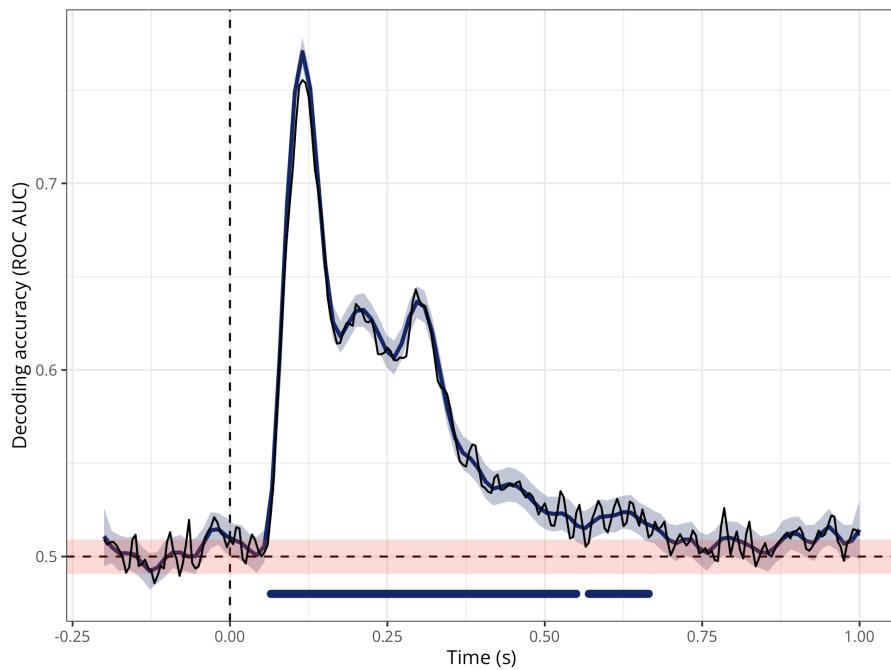
235 Note that although we chose a basis dimension of $k = 50$, which seems appropriate for
 236 the present data, this choice should be adapted according to the properties of the modelled data
 237 (e.g., signal-to-noise ratio, prior low-pass filtering, sampling rate, etc) and should be assessed
 238 by the usual model checking tools (e.g., posterior predictive checks, see also Section C). To
 239 better distinguish signal from noise, we also defined a region of practical equivalence (ROPE,
 240 Kruschke & Liddell, 2017), defined as the chance level plus the standard deviation of the (group-
 241 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
)
```

242 We assessed the reliability of the proposed approach using a form of permutation-based
 243 split-half reliability (as for instance in [Rosenblatt et al., 2018](#)), which consisted of the following
 244 steps. First, we created 1000 split halves of the data (i.e., with half the participants in the
 245 original data, that is, 16 participants). For each split, we estimated the onset/offset using all
 246 methods described previously. Third, we summarised the distribution of onset/offset estimates
 247 using the median “error” (i.e., difference between the split estimate and the estimate obtained
 248 using the full dataset) and the variance across splits. This approach allows assessing how similar
 249 the estimate of each half split is to the full dataset (thus acting as a proxy for the population)
 250 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.



251

Results

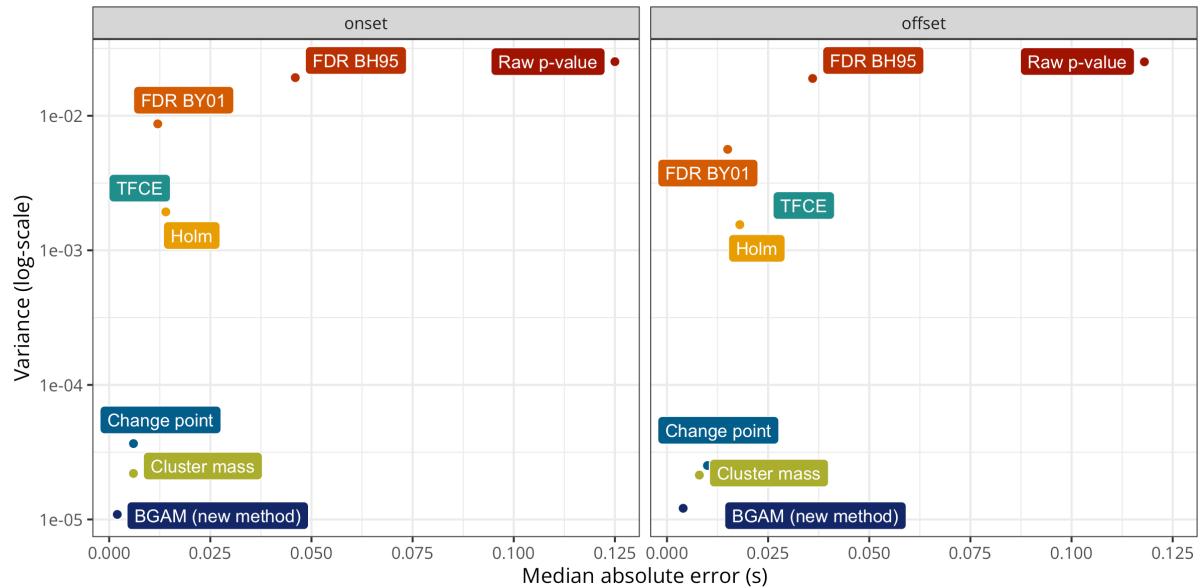
252 This section is divided in two parts. First, we present the results from the simulation
 253 study, assessing the bias and variance of each method when applied to simulated data in which
 254 the ground truth is known. Second, we present the results obtained when applying the different
 255 methods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

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

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

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.



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

Table 1

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

	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
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
FDR BH95	0.0080	0.0460	0.0599	0.0192	0.0801
Raw p-value	0.0060	0.1250	0.0711	0.0252	0.1942
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
Change point	-0.0100	0.0100	0.0096	0.0000	0.0059
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
FDR BH95	-0.0100	0.0360	0.0578	0.0189	0.0682
Raw p-value	-0.0080	0.1180	0.0725	0.0252	0.1868

269 Application to actual MEG data (reliability)

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

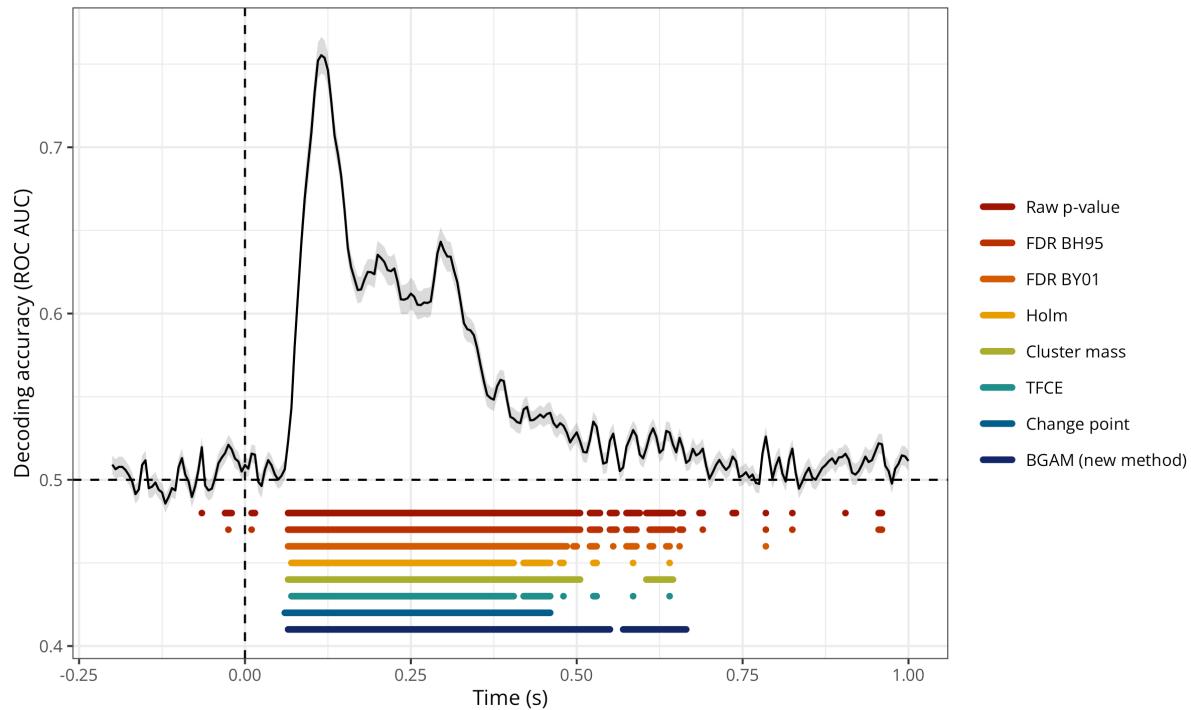
278 Figure 10 shows the median difference between the onset and offset estimates from
 279 each data split and the onset and offset estimates from the full dataset (x-axis) along with
 280 the variance of its onset and offset estimates across data splits (error bar). This figure reveals
 281 that the BGAM onset and offset estimates on each split are the closest to the estimates from the
 282 full dataset on average (0ms difference for the onset estimate and 5ms difference for the offset
 283 estimate). The Raw p-value method has similar performance, but given the aberrant estimates
 284 it produces (cf. Figure 9), the fact that it is consistent between data splits and the full dataset
 285 is not convincing on its own. The Change point method also has a very good performance (i.e.,
 286 very low difference between split estimates and full estimates), but produces too short cluster
 287 of significant decoding performance (cf. Figure 9).² Overall, the figure reveals that for all other
 288 methods, split datasets produce later onset estimates and earlier offset estimates (as compared
 289 to the estimates from the model fitted on the full dataset). These results highlight that the
 290 results of any method should be assessed using both i) good asymptotic properties on simulated
 291 data, ii) sensible identified clusters in actual data, and iii) reliable/stable estimates on actual
 292 data.

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.

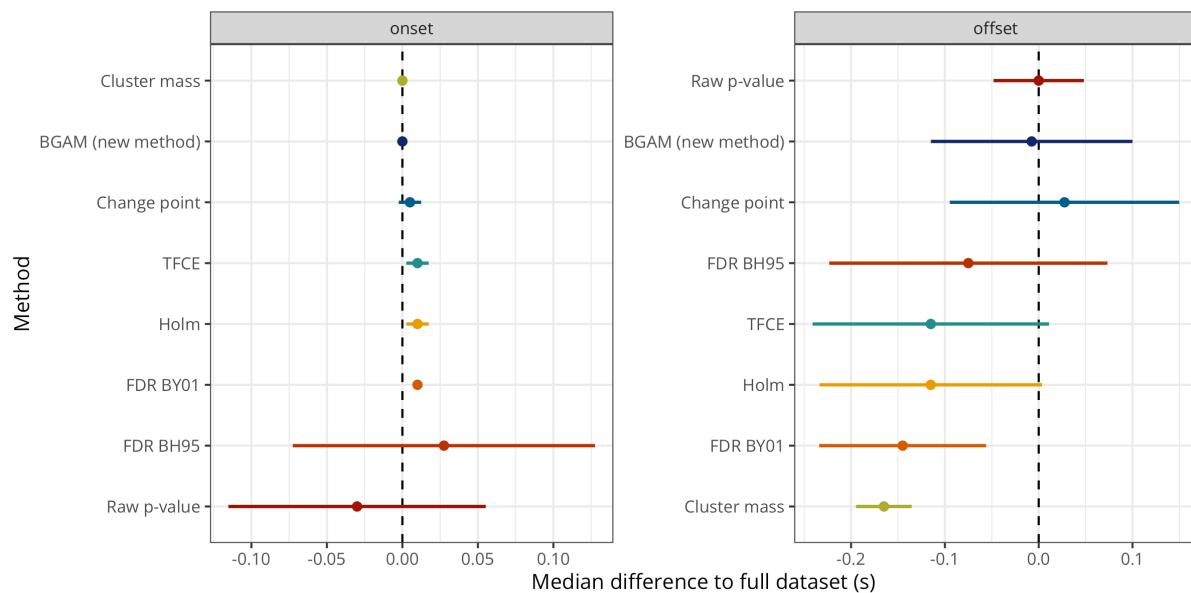
²As in Rousselet (2025), we fixed the number of expected change points to two in the binary segmentation algorithm, thus producing always one cluster.

Figure 9

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

**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).



293

Discussion

294 **Summary of the proposed approach**

295 Overall, before concluding on the onset/offset of effect based on the model, we need to
 296 ensure that the model provides a faithful description of the data-generating process (e.g., via

297 posterior predictive checks etc)...

298 The TFCE performs worse than the cluster-sum approach, which was anticipated by
299 Rousselet (2025) based on the initial results of S. Smith & Nichols (2009)...

300 **Limitations and future directions**

301 As in previous simulation work (e.g., Rousselet et al., 2008; Sassenhagen & Draschkow,
302 2019), the present simulation results depend on various choices such as the specific cluster-
303 forming algorithm and threshold, signal-to-noise ratio, negative impact of preprocessing steps
304 (e.g., low-pass filter) on temporal resolution... note however, that the same caveats apply to all
305 methods...

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

309 Can be applied to any 1D timeseries (e.g., pupillometry, electromyography)... Extending
310 the approach to spatiotemporal data (i.e., time + sensors) or spatiotemporal time-frequency 4D
311 data...

312 We kept the exemplary models simple, but can be extended by adding varying/random
313 effects (intercept and slope) for item (e.g., word)... but also continuous predictors at the trial
314 level?

315

Data and code availability

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

319

Packages

320 We used R version 4.4.3 (R Core Team, 2025) and the following R packages: assertthat
321 v. 0.2.1 (Wickham, 2019), brms v. 2.22.0 (Bürkner, 2017, 2018, 2021), changepoint v. 2.3
322 (Killick et al., 2024; Killick & Eckley, 2014), doParallel v. 1.0.17 (Corporation & Weston,
323 2022), easystats v. 0.7.4 (Lüdecke et al., 2022), foreach v. 1.5.2 (Microsoft & Weston, 2022),
324 furrr v. 0.3.1 (Vaughan & Dancho, 2022), future v. 1.34.0 (Bengtsson, 2021), ggrepel v. 0.9.6
325 (Slowikowski, 2024), glue v. 1.8.0 (Hester & Bryan, 2024), grateful v. 0.2.11 (Rodriguez-
326 Sanchez & Jackson, 2024), gt v. 1.0.0 (Iannone et al., 2025), knitr v. 1.50 (Xie, 2014, 2015,
327 2025), MetBrewer v. 0.2.0 (Mills, 2022), neurogam v. 0.0.1 (Nalborczyk, 2025), pakret v. 0.2.2
328 (Gallou, 2024), patchwork v. 1.3.0 (T. L. Pedersen, 2024), rmarkdown v. 2.29 (Allaire et al.,
329 2024; Xie et al., 2018, 2020), scales v. 1.3.0 (Wickham et al., 2023), scico v. 1.5.0 (T. L.
330 Pedersen & Cramer, 2023), tictoc v. 1.2.1 (Izrailev, 2024), tidybayes v. 3.0.7 (Kay, 2024),
331 tidytext v. 0.4.2 (Silge & Robinson, 2016), tidyverse v. 2.0.0 (Wickham et al., 2019).

332

Acknolwedgements

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

335

References

- 336 Abugaber, D., Finestrat, I., Luque, A., & Morgan-Short, K. (2023). Generalized additive mixed
 337 modeling of EEG supports dual-route accounts of morphosyntax in suggesting no word
 338 frequency effects on processing of regular grammatical forms. *Journal of Neurolinguistics*,
 339 67, 101137. <https://doi.org/10.1016/j.jneuroling.2023.101137>
- 340 Allaire, J., Xie, Y., Dervieux, C., McPherson, J., Luraschi, J., Ushey, K., Atkins, A., Wickham,
 341 H., Cheng, J., Chang, W., & Iannone, R. (2024). *rmarkdown: Dynamic documents for r*.
 342 <https://github.com/rstudio/rmarkdown>
- 343 Baayen, R. H., & Linke, M. (2020). *Generalized Additive Mixed Models* (pp. 563–591). Springer
 344 International Publishing. https://doi.org/10.1007/978-3-030-46216-1_23
- 345 Baayen, R. H., Rij, J. van, Cat, C. de, & Wood, S. (2018). *Autocorrelated errors in ex-
 346 perimental data in the language sciences: Some solutions offered by generalized additive
 347 mixed models* (pp. 49–69). Springer International Publishing. [https://doi.org/10.1007/978-3-319-69830-4_4](https://doi.org/10.1007/

 348 978-3-319-69830-4_4)
- 349 Bengtsson, H. (2021). A unifying framework for parallel and distributed processing in r using
 350 futures. *The R Journal*, 13(2), 208–227. <https://doi.org/10.32614/RJ-2021-048>
- 351 Benjamini, Y., & Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical and
 352 Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society Series B:
 353 Statistical Methodology*, 57(1), 289–300. <https://doi.org/10.1111/j.2517-6161.1995.tb02031.x>
- 355 Benjamini, Y., & Yekutieli, D. (2001). The control of the false discovery rate in multiple
 356 testing under dependency. *The Annals of Statistics*, 29(4). <https://doi.org/10.1214/aos/1013699998>
- 358 Bürkner, P.-C. (2017). brms: An R package for Bayesian multilevel models using Stan. *Journal
 359 of Statistical Software*, 80(1), 1–28. <https://doi.org/10.18637/jss.v080.i01>
- 360 Bürkner, P.-C. (2018). Advanced Bayesian multilevel modeling with the R package brms. *The
 361 R Journal*, 10(1), 395–411. <https://doi.org/10.32614/RJ-2018-017>
- 362 Bürkner, P.-C. (2021). Bayesian item response modeling in R with brms and Stan. *Journal of
 363 Statistical Software*, 100(5), 1–54. <https://doi.org/10.18637/jss.v100.i05>
- 364 Coretta, S., & Bürkner, P.-. C. (2025). *Bayesian beta regressions with brms in r: A tutorial for
 365 phoneticians*. http://dx.doi.org/10.31219/osf.io/f9rqp_v1
- 366 Corporation, M., & Weston, S. (2022). doParallel: Foreach parallel adaptor for the “parallel”
 367 package. <https://CRAN.R-project.org/package=doParallel>
- 368 Dimigen, O., & Ehinger, B. V. (2021). Regression-based analysis of combined EEG and eye-
 369 tracking data: Theory and applications. *Journal of Vision*, 21(1), 3. <https://doi.org/10.1167/jov.21.1.3>
- 371 Dinga, R., Fraza, C. J., Bayer, J. M. M., Kia, S. M., Beckmann, C. F., & Marquand, A.
 372 F. (2021). *Normative modeling of neuroimaging data using generalized additive models of
 373 location scale and shape*. <http://dx.doi.org/10.1101/2021.06.14.448106>
- 374 Dunagan, D., Jordan, T., Hale, J. T., Pylkkänen, L., & Chacón, D. A. (2024). *Evaluating
 375 the timecourses of morpho-orthographic, lexical, and grammatical processing following rapid
 376 parallel visual presentation: An EEG investigation in english*. <http://dx.doi.org/10.1101/2024.04.10.588861>
- 378 Dunn, O. J. (1961). Multiple Comparisons among Means. *Journal of the American Statistical
 379 Association*, 56(293), 52–64. <https://doi.org/10.1080/01621459.1961.10482090>
- 380 Fischer, Adrian G., & Ullsperger, M. (2013). Real and Fictive Outcomes Are Processed Dif-
 381 ferently but Converge on a Common Adaptive Mechanism. *Neuron*, 79(6), 1243–1255.
 382 <https://doi.org/10.1016/j.neuron.2013.07.006>
- 383 Gabry, J., Simpson, D., Vehtari, A., Betancourt, M., & Gelman, A. (2019). Visualization in
 384 Bayesian workflow. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*,
 385 182(2), 389–402. <https://doi.org/10.1111/rssa.12378>

- 386 Gallou, A. (2024). *pakret: Cite “R” packages on the fly in “R Markdown” and “Quarto”*. <https://CRAN.R-project.org/package=pakret>
- 387
- 388 Gelman, A., Vehtari, A., Simpson, D., Margossian, C. C., Carpenter, B., Yao, Y., Kennedy, L., Gabry, J., Bürkner, P.-C., & Modrák, M. (2020). Bayesian workflow. *arXiv:2011.01808 [Stat]*. <http://arxiv.org/abs/2011.01808>
- 389
- 390
- 391 Gramfort, A. (2013). MEG and EEG data analysis with MNE-python. *Frontiers in Neuroscience*, 7. <https://doi.org/10.3389/fnins.2013.00267>
- 392
- 393 Hastie, T. J., & Tibshirani, R. J. (2017). *Generalized Additive Models*. Routledge. <https://doi.org/10.1201/9780203753781>
- 394
- 395 Hauk, O., Davis, M. H., Ford, M., Pulvermüller, F., & Marslen-Wilson, W. D. (2006). The time course of visual word recognition as revealed by linear regression analysis of ERP data. *NeuroImage*, 30(4), 1383–1400. <https://doi.org/10.1016/j.neuroimage.2005.11.048>
- 396
- 397 Hendrix, P., Bolger, P., & Baayen, H. (2017). Distinct ERP signatures of word frequency, phrase frequency, and prototypicality in speech production. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 43(1), 128–149. <https://doi.org/10.1037/a0040332>
- 398
- 399
- 400
- 401 Hester, J., & Bryan, J. (2024). *glue: Interpreted string literals*. <https://CRAN.R-project.org/package=glue>
- 402
- 403 Holm, S. (1979). A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics*, 6(2), 65–70. <http://www.jstor.org/stable/4615733>
- 404
- 405 Iamone, R., Cheng, J., Schloerke, B., Hughes, E., Lauer, A., Seo, J., Brevoort, K., & Roy, O. (2025). *gt: Easily create presentation-ready display tables*. <https://CRAN.R-project.org/package=gt>
- 406
- 407
- 408 Izrailev, S. (2024). *tictoc: Functions for timing r scripts, as well as implementations of “Stack” and “StackList” structures*. <https://CRAN.R-project.org/package=tictoc>
- 409
- 410 Kay, M. (2024). *tidybayes: Tidy data and geoms for Bayesian models*. <https://doi.org/10.5281/zenodo.1308151>
- 411
- 412 Killick, R., & Eckley, I. A. (2014). changepoint: An R package for changepoint analysis. *Journal of Statistical Software*, 58(3), 1–19. <https://www.jstatsoft.org/article/view/v058i03>
- 413
- 414 Killick, R., Haynes, K., & Eckley, I. A. (2022). *changepoint: An R package for changepoint analysis*. <https://CRAN.R-project.org/package=changepoint>
- 415
- 416 Killick, R., Haynes, K., & Eckley, I. A. (2024). *changepoint: An R package for changepoint analysis*. <https://CRAN.R-project.org/package=changepoint>
- 417
- 418 King, J.-R., & Dehaene, S. (2014). Characterizing the dynamics of mental representations: the temporal generalization method. *Trends in Cognitive Sciences*, 18(4), 203–210. <https://doi.org/10.1016/j.tics.2014.01.002>
- 419
- 420
- 421 Kruschke, J. K., & Liddell, T. M. (2017). The Bayesian New Statistics: Hypothesis testing, estimation, meta-analysis, and power analysis from a Bayesian perspective. *Psychonomic Bulletin & Review*, 25(1), 178–206. <https://doi.org/10.3758/s13423-016-1221-4>
- 422
- 423
- 424 Kryuchkova, T., Tucker, B. V., Wurm, L. H., & Baayen, R. H. (2012). Danger and usefulness are detected early in auditory lexical processing: Evidence from electroencephalography. *Brain and Language*, 122(2), 81–91. <https://doi.org/10.1016/j.bandl.2012.05.005>
- 425
- 426
- 427 Lüdecke, D., Ben-Shachar, M. S., Patil, I., Wiernik, B. M., Bacher, E., Thériault, R., & Makowski, D. (2022). easystats: Framework for easy statistical modeling, visualization, and reporting. CRAN. <https://doi.org/10.32614/CRAN.package.easystats>
- 428
- 429
- 430 Maris, E. (2011). Statistical testing in electrophysiological studies. *Psychophysiology*, 49(4), 549–565. <https://doi.org/10.1111/j.1469-8986.2011.01320.x>
- 431
- 432 Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and MEG-data. *Journal of Neuroscience Methods*, 164(1), 177–190. <https://doi.org/10.1016/j.jneumeth.2007.03.024>
- 433
- 434
- 435 Meulman, N., Sprenger, S. A., Schmid, M. S., & Wieling, M. (2023). GAM-based individual difference measures for L2 ERP studies. *Research Methods in Applied Linguistics*, 2(3),
- 436

- 437 100079. <https://doi.org/10.1016/j.rmal.2023.100079>
- 438 Meulman, N., Wieling, M., Sprenger, S. A., Stowe, L. A., & Schmid, M. S. (2015). Age Effects
439 in L2 Grammar Processing as Revealed by ERPs and How (Not) to Study Them. *PLOS
440 ONE*, 10(12), e0143328. <https://doi.org/10.1371/journal.pone.0143328>
- 441 Microsoft, & Weston, S. (2022). *foreach*: Provides foreach looping construct. <https://CRAN.R-project.org/package=foreach>
- 442 Miller, D. L. (2025). Bayesian views of generalized additive modelling. *Methods in Ecology and
443 Evolution*. <https://doi.org/10.1111/2041-210x.14498>
- 444 Mills, B. R. (2022). *MetBrewer*: Color palettes inspired by works at the metropolitan museum
445 of art. <https://CRAN.R-project.org/package=MetBrewer>
- 446 Nalborczyk, L. (2025). *neurogam*: Precise temporal localisation of m/EEG effects with bayesian
447 generalised additive multilevel models. <https://github.com/lhalborczyk/neurogam>
- 448 Nalborczyk, L., Batailler, C., Lœvenbruck, H., Vilain, A., & Bürkner, P.-C. (2019). An In-
449 troduction to Bayesian Multilevel Models Using brms: A Case Study of Gender Effects
450 on Vowel Variability in Standard Indonesian. *Journal of Speech, Language, and Hearing
451 Research*, 62(5), 1225–1242. https://doi.org/10.1044/2018_jslhr-s-18-0006
- 452 Nalborczyk, L., Hauw, F., Torcy, H. de, Dehaene, S., & Cohen, L. (in preparation). *Neural and
453 representational dynamics of tickertape synesthesia*.
- 454 Pedersen, E. J., Miller, D. L., Simpson, G. L., & Ross, N. (2019). Hierarchical generalized
455 additive models in ecology: An introduction with mgcv. *PeerJ*, 7, e6876. <https://doi.org/10.7717/peerj.6876>
- 456 Pedersen, T. L. (2024). *patchwork*: The composer of plots. <https://CRAN.R-project.org/package=patchwork>
- 457 Pedersen, T. L., & Crameri, F. (2023). *scico*: Colour palettes based on the scientific colour-maps.
458 <https://CRAN.R-project.org/package=scico>
- 459 Pernet, C. R. (2022). Electroencephalography robust statistical linear modelling using a single
460 weight per trial. *Aperture Neuro*, 2, 1–22. <https://doi.org/10.52294/apertureneuro.2022.2.seeo9435>
- 461 Pernet, C. R., Chauveau, N., Gaspar, C., & Rousselet, G. A. (2011). LIMO EEG: A Tool-
462 box for Hierarchical LInear MOdeling of ElectroEncephaloGraphic Data. *Computational
463 Intelligence and Neuroscience*, 2011, 1–11. <https://doi.org/10.1155/2011/831409>
- 464 Pernet, C. R., Latinus, M., Nichols, T. E., & Rousselet, G. A. (2015). Cluster-based com-
465 putational methods for mass univariate analyses of event-related brain potentials/fields: A
466 simulation study. *Journal of Neuroscience Methods*, 250, 85–93. <https://doi.org/10.1016/j.jneumeth.2014.08.003>
- 467 R Core Team. (2025). *R: A language and environment for statistical computing*. R Foundation
468 for Statistical Computing. <https://www.R-project.org/>
- 469 Rasmussen, C. E., & Williams, C. K. I. (2005). *Gaussian Processes for Machine Learning*.
470 <https://doi.org/10.7551/mitpress/3206.001.0001>
- 471 Rigby, R. A., & Stasinopoulos, D. M. (2005). Generalized Additive Models for Location, Scale
472 and Shape. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 54(3),
473 507–554. <https://doi.org/10.1111/j.1467-9876.2005.00510.x>
- 474 Rij, J. van, Hendriks, P., Rijn, H. van, Baayen, R. H., & Wood, S. N. (2019). Analyzing
475 the Time Course of Pupillometric Data. *Trends in Hearing*, 23. <https://doi.org/10.1177/2331216519832483>
- 476 Riutort-Mayol, G., Bürkner, P.-C., Andersen, M. R., Solin, A., & Vehtari, A. (2023). Practical
477 Hilbert space approximate Bayesian Gaussian processes for probabilistic programming.
478 *Statistics and Computing*, 33(1), 17. <https://doi.org/10.1007/s11222-022-10167-2>
- 479 Rodriguez-Sanchez, F., & Jackson, C. P. (2024). *grateful*: Facilitate citation of R packages.
480 <https://pakillo.github.io/grateful/>
- 481 Rosenblatt, J. D., Finos, L., Weeda, W. D., Solari, A., & Goeman, J. J. (2018). All-

- 488 Resolutions Inference for brain imaging. *NeuroImage*, 181, 786–796. <https://doi.org/10.1016/j.neuroimage.2018.07.060>
- 489 Rousselet, G. A. (2025). Using cluster-based permutation tests to estimate MEG/EEG onsets:
490 How bad is it? *European Journal of Neuroscience*, 61(1), e16618. <https://doi.org/10.1111/ejn.16618>
- 491 Rousselet, G. A., Pernet, C. R., Bennett, P. J., & Sekuler, A. B. (2008). Parametric study
492 of EEG sensitivity to phase noise during face processing. *BMC Neuroscience*, 9(1). <https://doi.org/10.1186/1471-2202-9-98>
- 493 Sassenhagen, J., & Draschkow, D. (2019). Cluster-based permutation tests of MEG/EEG data
494 do not establish significance of effect latency or location. *Psychophysiology*, 56(6). <https://doi.org/10.1111/psyp.13335>
- 495 Silge, J., & Robinson, D. (2016). tidytext: Text mining and analysis using tidy data principles
496 in r. *JOSS*, 1(3). <https://doi.org/10.21105/joss.00037>
- 497 Skukies, R., & Ehinger, B. (2021). Modelling event duration and overlap during EEG analysis.
498 *Journal of Vision*, 21(9), 2037. <https://doi.org/10.1167/jov.21.9.2037>
- 499 Skukies, R., Schepers, J., & Ehinger, B. (2024, December 9). *Brain responses vary in duration*
500 - modeling strategies and challenges. <https://doi.org/10.1101/2024.12.05.626938>
- 501 Slowikowski, K. (2024). *ggrepel: Automatically position non-overlapping text labels with “gg-*
502 *plot2”*. <https://CRAN.R-project.org/package=ggrepel>
- 503 Smith, N. J., & Kutas, M. (2014a). Regression-based estimation of ERP waveforms: I. The
504 rERP framework. *Psychophysiology*, 52(2), 157–168. <https://doi.org/10.1111/psyp.12317>
- 505 Smith, N. J., & Kutas, M. (2014b). Regression-based estimation of ERP waveforms: II. Non-
506 linear effects, overlap correction, and practical considerations. *Psychophysiology*, 52(2),
507 169–181. <https://doi.org/10.1111/psyp.12320>
- 508 Smith, S., & Nichols, T. (2009). Threshold-free cluster enhancement: Addressing problems of
509 smoothing, threshold dependence and localisation in cluster inference. *NeuroImage*, 44(1),
510 83–98. <https://doi.org/10.1016/j.neuroimage.2008.03.061>
- 511 Sóskuthy, M. (2017). *Generalised additive mixed models for dynamic analysis in linguistics: A*
512 *practical introduction*. <https://doi.org/10.48550/ARXIV.1703.05339>
- 513 Sóskuthy, M. (2021). Evaluating generalised additive mixed modelling strategies for dynamic
514 speech analysis. *Journal of Phonetics*, 84, 101017. <https://doi.org/10.1016/j.wocn.2020.101017>
- 515 Teichmann, L. (2022). An empirically driven guide on using bayes factors for m/EEG decoding.
516 *Aperture Neuro*, 2, 1–10. <https://doi.org/10.52294/apertureneuro.2022.2.maoc6465>
- 517 Tremblay, A., & Newman, A. J. (2014). Modeling nonlinear relationships in ERP data using
518 mixed-effects regression with R examples. *Psychophysiology*, 52(1), 124–139. <https://doi.org/10.1111/psyp.12299>
- 519 Umlauf, N., Klein, N., & Zeileis, A. (2018). BAMLSS: Bayesian Additive Models for Location,
520 Scale, and Shape (and Beyond). *Journal of Computational and Graphical Statistics*, 27(3),
521 612–627. <https://doi.org/10.1080/10618600.2017.1407325>
- 522 Ushey, K., Allaire, J., & Tang, Y. (2024). *Reticulate: Interface to ‘python’*. <https://CRAN.R-project.org/package=reticulate>
- 523 Vaughan, D., & Dancho, M. (2022). *furrr: Apply mapping functions in parallel using futures*.
524 <https://CRAN.R-project.org/package=furrr>
- 525 Wickham, H. (2019). *assertthat: Easy pre and post assertions*. <https://CRAN.R-project.org/package=assertthat>
- 526 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Grolemund,
527 G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L., Miller, E., Bache, S. M.,
528 Müller, K., Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., ... Yutani, H. (2019). Welcome
529 to the tidyverse. *Journal of Open Source Software*, 4(43), 1686. <https://doi.org/10.21105/joss.01686>

- 539 Wickham, H., Pedersen, T. L., & Seidel, D. (2023). *scales: Scale functions for visualization.*
540 <https://CRAN.R-project.org/package=scales>
- 541 Wieling, M. (2018). Analyzing dynamic phonetic data using generalized additive mixed model-
542 ing: A tutorial focusing on articulatory differences between L1 and L2 speakers of English.
543 *Journal of Phonetics*, 70, 86–116. <https://doi.org/10.1016/j.wocn.2018.03.002>
- 544 Winter, B., & Wieling, M. (2016). How to analyze linguistic change using mixed models, Growth
545 Curve Analysis and Generalized Additive Modeling. *Journal of Language Evolution*, 1(1),
546 7–18. <https://doi.org/10.1093/jole/lzv003>
- 547 Wood, S. N. (2003). Thin Plate Regression Splines. *Journal of the Royal Statistical Society
548 Series B: Statistical Methodology*, 65(1), 95–114. <https://doi.org/10.1111/1467-9868.00374>
- 549 Wood, S. N. (2017a). *Generalized Additive Models*. Chapman; Hall/CRC. <https://doi.org/10.1201/9781315370279>
- 550 Wood, S. N. (2017b). *Generalized additive models: An introduction with r* (2nd ed.). Chapman;
551 Hall/CRC.
- 552 Wüllhorst, V., Wüllhorst, R., Overmeyer, R., & Endrass, T. (2025). Comprehensive Analysis
553 of Event-Related Potentials of Response Inhibition: The Role of Negative Urgency and
554 Compulsivity. *Psychophysiology*, 62(2). <https://doi.org/10.1111/psyp.70000>
- 555 Xie, Y. (2014). knitr: A comprehensive tool for reproducible research in R. In V. Stodden, F.
556 Leisch, & R. D. Peng (Eds.), *Implementing reproducible computational research*. Chapman;
557 Hall/CRC.
- 558 Xie, Y. (2015). *Dynamic documents with R and knitr* (2nd ed.). Chapman; Hall/CRC. <https://yihui.org/knitr/>
- 559 Xie, Y. (2025). knitr: A general-purpose package for dynamic report generation in R. <https://yihui.org/knitr/>
- 560 Xie, Y., Allaire, J. J., & Grolemund, G. (2018). *R markdown: The definitive guide*. Chapman;
561 Hall/CRC. <https://bookdown.org/yihui/rmarkdown>
- 562 Xie, Y., Dervieux, C., & Riederer, E. (2020). *R markdown cookbook*. Chapman; Hall/CRC.
563 <https://bookdown.org/yihui/rmarkdown-cookbook>
- 564 Yeung, N., Bogacz, R., Holroyd, C. B., & Cohen, J. D. (2004). Detection of synchronized os-
565 cillations in the electroencephalogram: An evaluation of methods. *Psychophysiology*, 41(6),
566 822–832. <https://doi.org/10.1111/j.1469-8986.2004.00239.x>
- 567
- 568
- 569

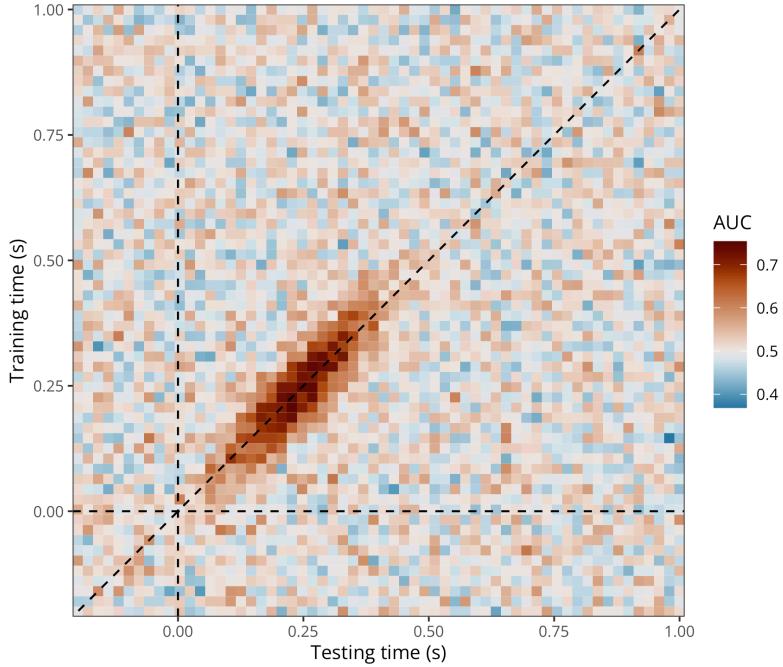
Appendix A

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

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

Figure A1

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



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

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

where we assume that AUC values come from a Beta distribution with two parameters μ and ϕ . We can think of $f(\text{train}_i, \text{test}_i)$ as a surface (a smooth function of two variables) that we can model using a 2-dimensional splines. Let $\mathbf{s}_i = (\text{train}_i, \text{test}_i)$ be some pair of training and 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 test_i . We can then express the smooth function as:

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

Note that $b_m(\cdot)$ is a basis function that maps $R \times R \rightarrow R$. A popular bivariate basis function uses *thin-plate splines* (Wood, 2003), which extend to $\mathbf{s}_i \in \mathbb{R}^d$ and ∂l_g penalties. These splines are designed to interpolate and approximate smooth surfaces over two dimensions (hence the “bivariate” term). For $d = 2$ dimensions and $l = 2$ (smoothness penalty involving second order derivative):

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

588 using the the radial basis function given by:

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

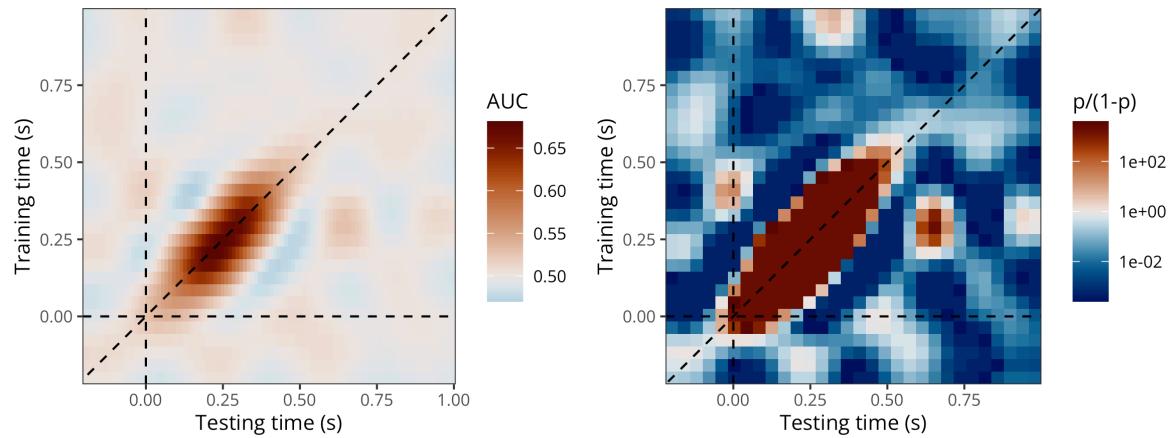
589 where $\|\mathbf{s}_i - \mathbf{k}_m\|$ is the Euclidean distance between the covariate \mathbf{s}_i and the knot location
590 \mathbf{k}_m . We fitted this model using `brms`...

```
# fitting a GAM with two temporal dimensions
timegen_gam <- brm(
  # 2D thin-plate spline (tp)
  auc ~ t2(train_time, test_time, bs = "tp", k = 20),
  # auc ~ t2(train_time, test_time, bs = "tp", k = 30),
  data = timegen_data,
  family = Beta(),
  warmup = 1000,
  iter = 2000,
  chains = 8,
  cores = 8,
  file = "models/timegen_gam_t2.rds"
  # file = "models/timegen_gam_t2_k30.rds"
)

# fitting a GP with two temporal dimensions
# timegen_gp <- brm(
#   auc ~ gp(train_time, test_time, k = 20),
#   data = timegen_data,
#   family = Beta(),
#   control = list(adapt_delta = 0.95),
#   iter = 2000,
#   chains = 4,
#   cores = 4,
#   file = "models/timegen_gp.rds"
# )
```

Figure A2

Predicted AUC values (left) and posterior probabilities of decoding accuracy being above chance level (right) according to the bivariate GAM.



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

Appendix B

Alternative to GAMs: Approximate Gaussian Process regression

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

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

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

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

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

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

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

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

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

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

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

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

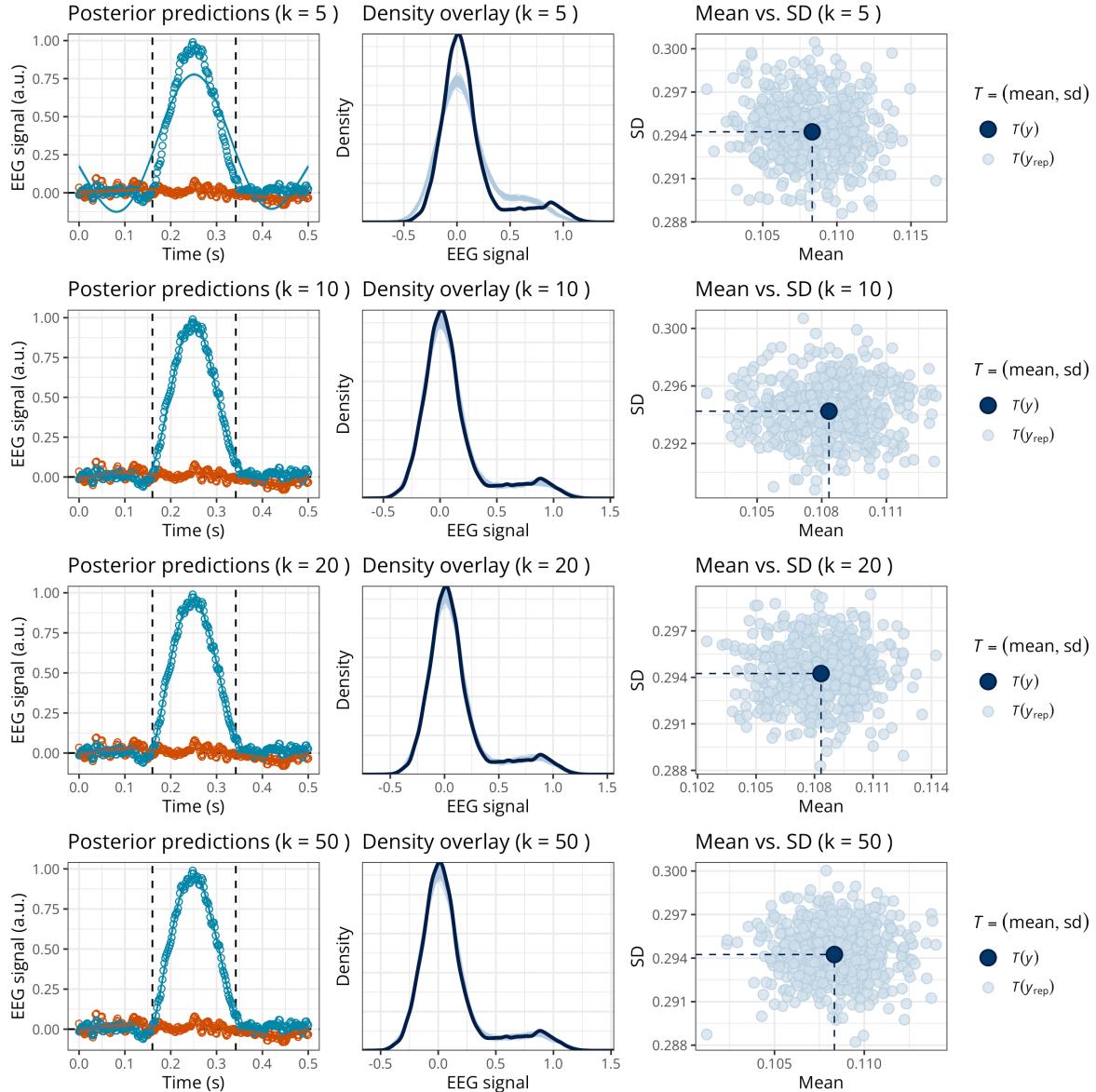
627 Low-rank approximate Gaussian processes are of main interest in machine learning and
 628 statistics due to the high computational demands of exact Gaussian process models (Riutort-
 629 Mayol et al., 2023)...

Appendix C How to choose the GAM basis dimension?

630 Here provide recommendation about how to define k . An option is to vary k and examine the
 631 predictions and posterior predictive checks (PPCs) of each model... In this example (Figure C1)...
 632 However, it is not possible to provide general recommendations, as the optimal k depends on
 633 the sampling rate, the preprocessing steps (e.g., signal-to-noise ratio, low-pass filtering, etc),
 634 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

635 For users who are already familiar with `brms`, the recommended pipeline is to import ERPs
 636 or decoding results in R and analyse these data using the code provided in the main paper.
 637 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.
 638

```
# 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 use 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
```

639 The `neurogam` package can also be called from Python using the `rpy2` module, and
 640 can easily be integrated into MNE-Python pipelines. For example, we use it below to estimate
 641 the onset and offset of effects for one EEG channel from a MNE evoked object. Note that
 642 the code used to reshape the `sample` MNE dataset is provided later, and we refer to the [MNE](#)
 643 [documentation](#) about converting MNE epochs to Pandas dataframes in long format.

```
# loading the Python modules
import rpy2.robj as robjects
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(robjects.default_converter + pandas2ri.converter):

    reshaped_df_r = robjects.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) )
```

```
# loading the Python modules
import mne
from mne.datasets import sample
import numpy as np
import pandas as pd

# defining the path to the "sample" dataset
path = sample.data_path()

# loading the evoked data
evokeds = mne.read_evokeds(path / "MEG" / "sample" / "sample_audvis-ave.fif")

# defining a function to reshape the data
def reshape_eeg_channels_as_participants(evokeds, condition_labels):

    all_dfs = []

    for evoked, condition in zip(evokeds, condition_labels):

        # selecting only EEG channels
        picks = mne.pick_types(evoked.info, meg=False, eeg=True)
        data = evoked.data[picks, :]
        channel_names = np.array(evoked.ch_names)[picks]

        n_channels, n_times = data.shape
        # repeating time for each channel
        times = np.tile(evoked.times, n_channels)
        meeg_values = data.flatten()
        pseudo_participant_ids = np.repeat(channel_names, n_times)

        # converting to dataframe
        df = pd.DataFrame({
            "participant": pseudo_participant_ids,
            "time": times,
```

```
        "eeg": meeg_values,
        "condition": condition
    })

    all_dfs.append(df)

    return pd.concat(all_dfs, ignore_index=True)

# picking two evoked conditions
faces = evokeds[0]
scrambled = evokeds[1]

# reshaping data by pretending each EEG channel is a participant
reshaped_df = reshape_eeg_channels_as_participants(
    evokeds=[faces, scrambled],
    condition_labels=["faces", "scrambled"]
)
```