

1 Event History Analysis for psychological time-to-event data: A tutorial in R with examples
2 in Bayesian and frequentist workflows

3 Sven Panis¹ & Richard Ramsey¹

4 ¹ ETH Zürich

5 Author Note

6 Neural Control of Movement lab, Department of Health Sciences and Technology
7 (D-HEST). Social Brain Sciences lab, Department of Humanities, Social and Political
8 Sciences (D-GESS).

9 Correspondence concerning this article should be addressed to Sven Panis, ETH
10 GLC, room G16.2, Gloriastrasse 37/39, 8006 Zürich. E-mail: sven.panis@hest.ethz.ch

11

Abstract

12 Time-to-event data such as response times and saccade latencies form a cornerstone of
13 experimental psychology, and have had a widespread impact on our understanding of
14 human cognition. However, the orthodox method for analyzing such data – comparing
15 means between conditions – is known to conceal valuable information about the timeline of
16 psychological effects, such as their onset time and how they evolve with increasing waiting
17 time. The ability to reveal finer-grained, “temporal states” of cognitive processes can have
18 important consequences for theory development by qualitatively changing the key
19 inferences that are drawn from psychological data. Luckily, well-established analytical
20 approaches, such as event history analysis (EHA), are able to evaluate the detailed shape
21 of time-to-event distributions, and thus characterize the time course of psychological states.
22 One barrier to wider use of EHA, however, is that the analytical workflow is typically more
23 time-consuming and complex than orthodox approaches. To help achieve broader uptake of
24 EHA, in this paper we outline a set of tutorials that detail one distributional method
25 known as discrete-time EHA. We touch upon several key aspects of the workflow, such as
26 how to process raw data and specify regression models, and we also consider the
27 implications for experimental design. We finish the article by considering the benefits of
28 the approach for understanding psychological states, as well as the limitations and future
29 directions of this work. Finally, the project is written in R and freely available, which
30 means the approach can easily be adapted to other data sets.

31 *Keywords:* response times, event history analysis, Bayesian multilevel regression
32 models, experimental psychology, cognitive psychology

33 Word count: 11664 (body) + 1593 (references) + 2394 (supplemental material)

34

1. Introduction

35 1.1 Motivation and background context: Comparing means versus 36 distributional shapes

37 In experimental psychology, it is standard practice to analyse response times (RTs),
38 saccade latencies, and fixation durations by calculating average performance across a series
39 of trials. Such comparisons between means have been the workhorse of experimental
40 psychology over the last century, and have had a substantial impact on theory development
41 as well as our understanding of the structure of cognition and brain function. Indeed, the
42 view that mean values are truth and variations around the mean are error is deeply
43 ingrained in experimental psychology (Bolger et al., 2019). However, differences in mean
44 RT conceal important pieces of information, such as when an experimental effect starts,
45 how it evolves with increasing waiting time, and whether its onset is time-locked to other
46 events (Panis, 2020; Panis, Moran, Wolkersdorfer, & Schmidt, 2020; Panis & Schmidt,
47 2016, 2022; Panis, Torfs, Gillebert, Wagemans, & Humphreys, 2017; Panis & Wagemans,
48 2009; Wolkersdorfer, Panis, & Schmidt, 2020). Such information is useful not only for the
49 interpretation of experimental effects under investigation, but also for cognitive
50 psychophysiology and computational model selection (Panis, Schmidt, Wolkersdorfer, &
51 Schmidt, 2020).

52 As a simple illustration, Figure 1 summarises simulated single-subject data (200 trials
53 per condition) that shows how comparing means between two conditions can conceal the
54 shapes of the underlying RT and accuracy distributions. Indeed, compared to the
55 aggregation of data across trials (Figure 1A), a distributional approach offers the
56 possibility to reveal the time course of psychological states (Figure 1B).

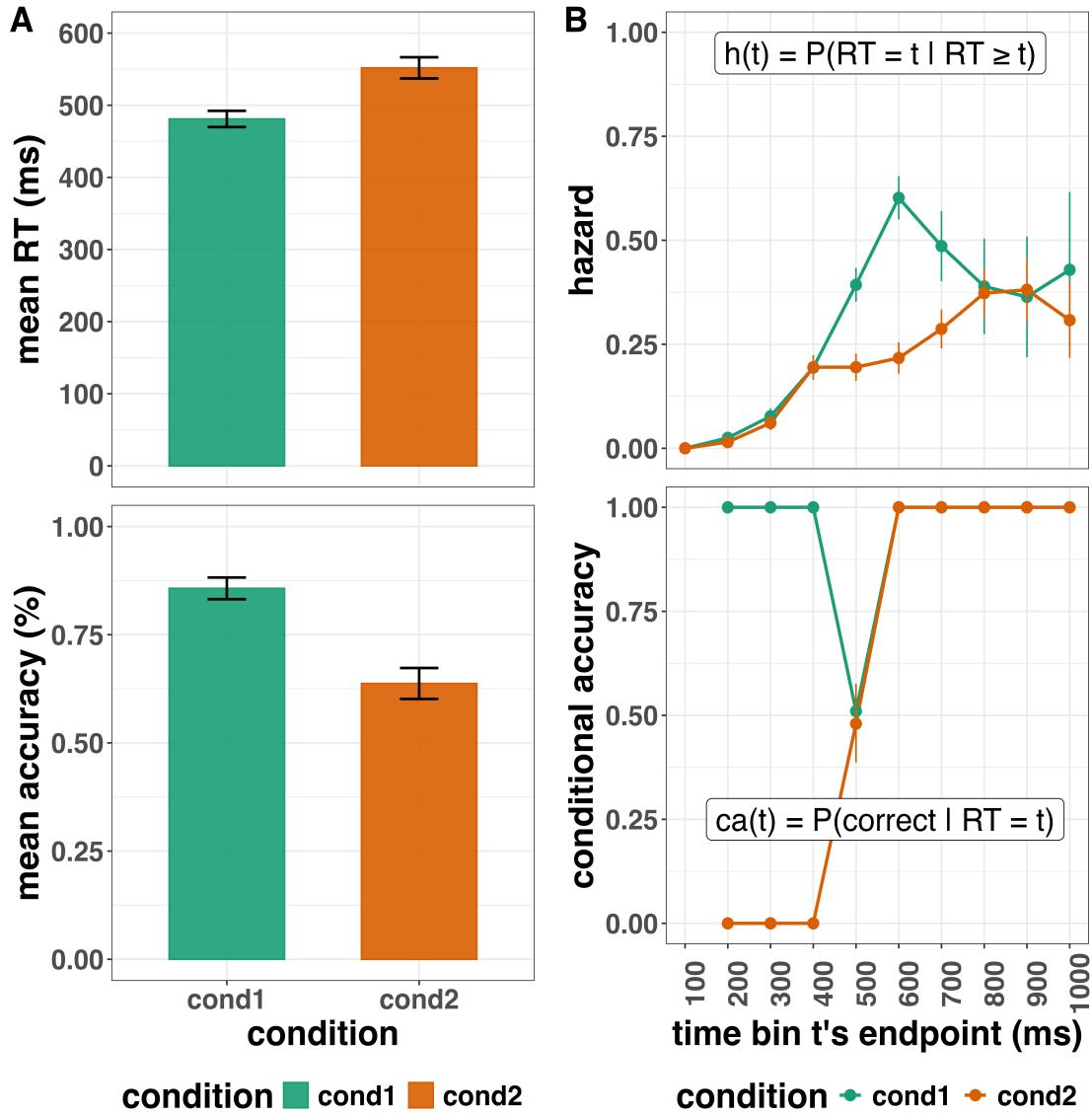


Figure 1. Simulated data showing mean performance versus distributional (EHA/SAT) analyses. (A) The mean RT (top) and overall accuracy (bottom) for two conditions are plotted. (B) The discrete-time hazard functions (top) and conditional accuracy functions (bottom) are plotted for the same data. The first second after target stimulus onset (time zero) is divided in ten bins of 100 ms. The first bin is (0,100], the last bin is (900,1000]. Note that the hazard and conditional accuracy estimates are plotted at the endpoint of each time bin. The definitions of discrete-time hazard and conditional accuracy are further explained in section 2. Error bars represent ± 1 standard error of the mean (A) or proportion (B).

57 For example, Figure 1B shows a first state (up to 400 ms after target onset) for which
58 the early upswing in hazard is equal for both conditions, and the emitted responses are
59 always correct in condition 1 and always incorrect in condition 2. In a second state (400 to
60 500 ms), hazard is higher in condition 1, and conditional accuracies are close to .5 in both
61 conditions. In a third state (>500 ms), the effect disappears in hazard, and all conditional
62 accuracies are equal to 1 (see also Panis & Schmidt, 2016).

63 Why does this matter for research in psychology? For many psychological questions,
64 the estimation of such “temporal states” information can be theoretically meaningful by
65 leading to more fine-grained understanding of psychological processes. Because EHA adds
66 a relatively under-used but ever-present dimension – the passage of time – to the theory
67 building toolkit, it provides one possible answer to the recent call for better temporal
68 methods (ref1, ref2).

69 1.2 Aims

70 Our ultimate aim in this paper is twofold: first, we want to convince readers of the
71 many benefits of using EHA when dealing with psychological RT data, and second, we
72 want to provide a set of practical tutorials, which provide step-by-step instructions on how
73 you actually perform a discrete-time EHA on RT data, as well as a complementary
74 discrete-time speed-accuracy tradeoff (SAT) analysis on timed accuracy data in case of
75 choice RT data.

76 Even though EHA is a widely used statistical tool and there already exist many
77 excellent reviews (Allison, 1982; Blossfeld & Rohwer, 2002; Box-Steffensmeier, 2004;
78 Hosmer, Lemeshow, & May, 2011; e.g., Singer & Willett, 2003; Teachman, 1983) and
79 tutorials (e.g., Allison, 2010; Landes, Engelhardt, & Pelletier, 2020), we are not aware of
80 any tutorials that are aimed specifically at psychological RT (+ accuracy) data, and which
81 provide worked examples of the key data processing and Bayesian multilevel regression

82 modelling steps. Set within this context, our overall aim is to introduce a set of tutorials,
83 which explain **how** to do such analyses in the context of experimental psychology, rather
84 than repeat in any detail **why** you may do them. Therefore, we hope that our tutorials will
85 provide a pathway for research avenues in experimental psychology that have the potential
86 to benefit from using EHAf in the future.

87 **1.3 Structure**

88 In what follows, the paper is organised in three main sections. In Section 2, we
89 provide a brief overview of EHA to orient the reader to the basic concepts that we will use
90 throughout the paper and why such an approach might be relevant for research in
91 experimental psychology. In Section 3, we outline a series of tutorials, which are written in
92 the R programming language and publicly available on our Github page
93 (https://github.com/sven-panis/Tutorial_Event_History_Analysis), along with all of the
94 other code and material associated with the project. The tutorials provide hands-on,
95 concrete examples of key parts of the analytical process, such as data wrangling, model
96 fitting and planning future studies, so that others can apply EHA to their own
97 time-to-event data measured in RT tasks. In Section 4, we discuss the strengths and
98 weaknesses of the approach for researchers in experimental psychology.

99 **2. What is event history analysis and why is it relevant to research in**
100 **experimental psychology?**

101 **2.1 A brief introduction to event history analysis**

102 EHA is a class of statistical approaches to study the occurrence and timing of events,
103 such as disease onset, marriages, arrests, and job terminations (Allison, 2010). In this
104 section, we want to provide an intuition regarding how EHA works in general, as well as in
105 the context of experimental psychology. For those who want more detailed treatment of

106 EHA and/or regression equations, we refer the reader to several excellent textbooks on
107 these topics (Allison, 2010; Gelman, Hill, & Vehtari, 2020; Singer & Willett, 2003; Winter,
108 2019). We also visualize and discuss the types of time-to-event data that are obtained in
109 typical RT tasks in section A of the Supplemental Material, and supply relevant regression
110 equations in section E of the Supplemental Material.

111 **2.1.1 Terminology and minimum requirements for EHA.** To avoid possible
112 confusion in terminology used, it is worth noting that EHA is known by various labels,
113 such as survival analysis, hazard analysis, duration analysis, failure-time analysis, and
114 transition analysis (Singer & Willett, 2003). In this paper, we choose to use the term EHA
115 throughout.

116 In terms of minimum requirements to apply a single-event EHA, one must be able to:

- 117 1. define an event of interest that represents a qualitative change - a transition from one
118 discrete state to another - that can be situated in time (e.g., a button press, a
119 saccade onset, a fixation offset, etc.);
- 120 2. define time point zero in each trial (e.g., target stimulus onset, fixation onset, etc.);
- 121 3. measure the passage of time between time point zero and event occurrence in discrete
122 or continuous time units in each trial.

123 **2.1.2 Types of EHA.** There are different types of EHA. For example, the
124 definition of hazard and the type of models employed depend on whether one is using
125 continuous or discrete time units. As a lab, and mainly for practical reasons, we have much
126 more experience using discrete time EHA, and that is the approach that we describe and
127 focus on in this paper. This choice may seem counter-intuitive, given that RT is typically
128 treated as a continuous variable. However, continuous forms of EHA require much more
129 data to estimate the continuous-time hazard (rate) function well (REFS). Thus, by trading
130 a bit of temporal resolution for a lower number of trials, discrete-time methods seem ideal

131 for dealing with typical psychological RT data sets for which there are less than ~200 trials
132 per condition per experiment (REF PANIS). Moreover, as indicated by Allison (2010),
133 learning discrete-time EHA methods first will help in learning continuous-time methods, so
134 it seems like a good starting point.

135 To apply discrete-time EHA, one divides the within-trial time in discrete, contiguous
136 time bins indexed by t (e.g., $t = 1:10$ time bins; Figure 1B). Then let RT be a discrete
137 random variable denoting the rank of the time bin in which a particular person's response
138 occurs in a particular trial (i.e., repeated measure). For example, a response in one trial
139 might occur at 546 ms and it would be in time bin 6 (any RTs from 501 ms to 600 ms).
140 One then calculates the sample-based estimate of the discrete-time hazard function of
141 event occurrence for each experimental condition (Figure 1B top). The discrete-time
142 hazard function gives you, for each time bin, the conditional probability that the event
143 occurs (sometime) in bin t , given that the event does not occur in previous bins. In other
144 words, it reflects the instantaneous risk that the event occurs in the current bin t , given
145 that it has not yet occurred in the past, i.e., in one of the prior bins ($t-1, t-2, \dots, 1$).

146 In the context of experimental psychology, it is often (but not always), the case that
147 responses can be classified as correct or incorrect. In those cases, one can also calculate the
148 conditional accuracy function (Figure 1B bottom). The conditional accuracy function gives
149 you for each time bin the conditional probability that a response is correct given that it is
150 emitted in time bin t (Allison, 2010; Kantowitz & Pachella, 2021; Wickelgren, 1977). The
151 $ca(t)$ function is also known as the micro-level speed-accuracy tradeoff (SAT) function. We
152 refer to this extended (hazard + conditional accuracy) analysis for choice RT data as
153 EHA/SAT.

154 The definitions of these and other discrete-time functions is given in section B of the
155 Supplemental Material.

156 2.2 Benefits of event history analysis for research in experimental psychology

157 Statisticians and mathematical psychologists recommend focusing on the hazard
158 function when analyzing time-to-event data for various reasons (REF?? - Sven please cite a
159 relevant REF here). We do not cover these benefits in detail here, as these are more
160 general topics that have been covered elsewhere in textbooks (see also section F of the
161 Supplemental Material). Instead, here we focus on the benefits as we see them for common
162 research programmes in experimental psychology.

163 We highlight three benefits that we think are relevant to the domain of experimental
164 psychology. First, as illustrated in Figure 1, compared to averaging data across trials,
165 integrating results between hazard functions and their associated conditional accuracy
166 functions for choice RT data can be informative for understanding psychological processes,
167 in terms of inferences about the microgenesis and temporal organization of cognition and
168 theoretical development. As such, the approach permits different kinds of questions to be
169 asked, different inferences to be made, and it holds the potential to discriminate between
170 theoretical accounts of psychological and/or brain-based processes. For example, what kind
171 of theory or set of mechanisms could account for the shape of the functions and the
172 temporally localized effects reported in Figure 1B? Are there new auxiliary assumptions
173 that computational models need to adopt (ref search)? Will the temporal effect patterns
174 align nicely with EEG findings (ref IOR)? And are there new experiments that need to be
175 performed to test the novel predictions that follow from these analyses?

176 Second, compared to more conventional analytical approaches, EHA uses more of the
177 data because it deals with missing data differently. It is conventional with RT data to
178 either (a) use a response deadline and discard all trials without a response, or (b) wait in
179 each trial until a response occurs and then apply data trimming techniques, i.e., discarding
180 too short or too long RTs (and perhaps also erroneous responses) before calculating a mean
181 RT (REF). Discarding data can introduce biases, however. Rather than treat

182 non-responses as missing data, EHA treats such trials as *right-censored* observations on the
183 variable RT, because all we know is that RT is greater than some value. Right-censoring is
184 a type of missing data problem and a nearly universal feature of survival data including RT
185 data. For example, if the censoring time was 1 second, then some trials result in observed
186 event times (those with a RT below 1 second), while the other trials result in response
187 times that are right-censored at 1 second. The fact that EHA can deal with
188 right-censoring, therefore, presents a analytical strength of the approach compared to many
189 common approaches in experimental psychology (ANOVA, linear regression, delta plots).

190 Third, the approach is generalisable and applicable to many tasks that are commonly
191 used in experimental psychology, such as detection, discrimination and bistable perception
192 tasks, and to a range of common experimental manipulations, such as
193 stimulus-onset-asynchrony (see section A of the Supplemental Material). The upshot is
194 that one general analytical approach, which holds several potential advantages, is widely
195 applicable to many substantive use-cases in the domain of experimental psychology,
196 irrespective of the analyst's current view on the nature of cognition (REFS).

197 **2.3 Implications for research design in experimental psychology**

198 Performing EHA in experimental psychology has implications for how experiments
199 are designed. More specifically, we consider three implications that researchers will need to
200 consider when using discrete-time EHA.

201 First, one can use a response deadline in each trial because EHA deals with
202 right-censored observations.

203 Second, since the number of trials per condition are spread across bins, it is
204 important to have a relatively large number of trial repetitions per participant and per
205 condition. Accordingly, experimental designs using this approach typically focus on
206 factorial, within-subject designs, in which a large number of observations are made on a

207 relatively small number of participants (so-called small- N designs). This approach
208 emphasizes the precision and reproducibility of data patterns at the individual participant
209 level to increase the inferential validity of the design (Baker et al., 2021; Smith & Little,
210 2018). Note that because statistical power derives both from the number of participants
211 and from the number of repeated measures per participant and condition, small- N designs
212 can still achieve what are generally considered acceptable levels of statistical power, if they
213 have a sufficient amount of data overall (Baker et al., 2021; Smith & Little, 2018).

214 Third, the width of each time bin will need to be determined. For instance, in Figure
215 1B we chose 100ms in an arbitrary manner. In reality, however, bin width will need to be
216 set by considering a number of factors simultaneously. The optimal bin width will depend
217 on (a) the length of the observation period in each trial, (b) the rarity of event occurrence,
218 (c) the number of repeated measures (or trials) per condition per participant, and (d) the
219 shape of the hazard function. Finding an appropriate bin width in a given user case before
220 fitting models will require testing a number of options, when calculating and plotting the
221 descriptive statistics (see section 3.1). The goal is to find the smallest bin width that is
222 supported by the amount of data available. Based on our experience, a bin width of 50 ms
223 is a good starting value when the number of repeated measures is 100 or less. Too small
224 bin widths will result in erratic hazard functions as many bins will have no events, and
225 thus hazard estimates of zero. Interestingly, the time bins do not need to have the same
226 width. For example, Panis (2020) used larger bins towards the end of the observation
227 period, as fewer events occurred there.

228 3. Tutorials

229 We used r my_r_citation\$r for all reported analyses. The content of the tutorials, in
230 terms of EHA and multilevel regression modelling, is mainly based on Allison (2010), Singer
231 and Willett (2003), McElreath (2020), Heiss (2021), Kurz (2023a), and Kurz (2023b).

232 Tutorials 1a and 1b show how to calculate and plot the descriptive statistics of

233 EHA/SAT when there are one or two independent variables, respectively. Tutorials 2a and

234 2b illustrate how to use Bayesian multilevel modeling to fit hazard and conditional

235 accuracy models, respectively. Tutorials 3a and 3b show how to implement, respectively,

236 multilevel models for hazard and conditional accuracy in the frequentist framework.

237 Additionally, to further simplify the process for other users, the first two tutorials rely on a

238 set of our own custom functions that make sub-processes easier to automate, such as data

239 wrangling and plotting functions (see section C in the Supplemental Material for a list of

240 the custom functions).

241 Our list of tutorials is as follows:

- 1a. Wrangle raw data and calculate descriptive stats for one independent variable

- 1b. Wrangle raw data and calculate descriptive stats for two independent variables

- 2a. Bayesian multilevel modeling for $h(t)$

- 2b. Bayesian multilevel modeling for $ca(t)$

- 3a. Frequentist multilevel modeling for $h(t)$

- 3b. Frequentist multilevel modeling for $ca(t)$

- 4. Simulation and power analysis for planning experiments

249 **3.1 Tutorial 1a: Calculating descriptive statistics using a life table**

250 **3.1.1 Data wrangling aims.** Our data wrangling procedures serve two related

251 purposes. First, we want to summarise and visualise descriptive statistics using a life table.

252 A life table includes for each time bin, the risk set (i.e., the number of trials that are

253 event-free at the start of the bin), the number of observed events, and the estimates of the

254 discrete-time hazard function $h(t)$, survivor function $S(t)$, probability mass function $P(t)$,

255 possibly the conditional accuracy function $ca(t)$, and their estimated standard errors (se).

256 The definitions of these functions are provided in section A of the Supplemental Material.

257 Second, we want to produce two different data sets that can each be submitted to
 258 different types of inferential modelling approaches. The two types of data structure we
 259 label as ‘person-trial’ data and ‘person-trial-bin’ data. The ‘person-trial’ data (Table 1)
 260 will be familiar to most researchers who record behavioural responses from participants, as
 261 it represents the measured RT and accuracy per trial within an experiment. This data set
 262 is used when fitting conditional accuracy models (Tutorials 2b and 3b).

```
263 ## Warning in attr(x, "align"): 'xfun::attr()' is deprecated.  

264 ## Use 'xfun::attr2()' instead.  

265 ## See help("Deprecated")
```

Table 1

Data structure for ‘person-trial’ data

pid	trial	condition	rt	accuracy
1	1	congruent	373.49	1
1	2	incongruent	431.31	1
1	3	congruent	455.43	0
1	4	incongruent	622.41	1
1	5	incongruent	535.98	1
1	6	incongruent	540.08	1
1	7	congruent	511.07	1
1	8	incongruent	444.42	1
1	9	congruent	678.69	1
1	10	congruent	549.79	1

Note. The first 10 trials for participant 1 are shown. These data are simulated and for illustrative purposes only.

266 In contrast, the ‘person-trial-bin’ data (Table 2) has a different, more extended

267 structure, which indicates in which bin a response occurred, if at all, in each trial.

268 Therefore, the ‘person-trial-bin’ data generates a 0 in each bin until an event occurs and

269 then it generates a 1 to signal an event has occurred in that bin. This data set is used

270 when fitting hazard models (Tutorials 2a and 3a). It is worth pointing out that there is no

271 requirement for an event to occur at all (in any bin), as maybe there was no response on

272 that trial or the event occurred after the time window of interest. Likewise, when the event

273 occurs in bin 1 there would only be one row of data for that trial in the person-trial-bin

274 data set.

```
275 ## Warning in attr(x, "align"): 'xfun::attr()' is deprecated.
```

```
276 ## Use 'xfun::attr2()' instead.
```

```
277 ## See help("Deprecated")
```

Table 2
Data structure for ‘person-trial-bin’ data

pid	trial	condition	timebin	event
1	1	congruent	1	0
1	1	congruent	2	0
1	1	congruent	3	0
1	1	congruent	4	1
1	2	incongruent	1	0
1	2	incongruent	2	0
1	2	incongruent	3	0
1	2	incongruent	4	0
1	2	incongruent	5	1

Note. The first 2 trials for participant 1 from Table 1 are shown. The width of the time bins is 100 ms. These data are simulated and for illustrative purposes only.

278 **3.1.2 A real data wrangling example.** To illustrate how to quickly set up life
 279 tables for calculating the descriptive statistics (functions of discrete time), we use a
 280 published data set on masked response priming from Panis and Schmidt (2016). In their
 281 first experiment, Panis and Schmidt (2016) presented a double arrow for 94 ms that
 282 pointed left or right as the target stimulus with an onset at time point zero in each trial.
 283 Participants had to indicate the direction in which the double arrow pointed using their
 284 corresponding index finger, within 800 ms after target onset. Response time and accuracy
 285 were recorded on each trial. Prime type (blank, congruent, incongruent) and mask type
 286 were manipulated. Here we focus on the subset of trials in which no mask was presented.

287 The 13-ms prime stimulus was a double arrow presented 187 ms before target onset in the
 288 congruent (same direction as target) and incongruent (opposite direction as target) prime
 289 conditions.

290 There are several data wrangling steps to be taken. First, we need to load the data
 291 before we (a) supply required column names, and (b) specify the factor condition with the
 292 correct levels and labels.

293 The required column names are as follows:

- 294 • “pid”, indicating unique participant IDs;
- 295 • “trial”, indicating each unique trial per participant;
- 296 • “condition”, a factor indicating the levels of the independent variable (1, 2, ...) and
 the corresponding labels;
- 298 • “rt”, indicating the response times in ms;
- 299 • “acc”, indicating the accuracies (1/0).

300 In the code of Tutorial 1a, this is accomplished as follows.

```
data_wr<-read_csv("../Tutorial_1_descriptive_stats/data/DataExp1_6subjects_wrangled.csv")
data_wr <- data_wr %>%
  rename(pid = vp, condition = prime_type, acc = respac, trial = TrialNr) %>%
  mutate(condition = condition + 1, # original levels were 0, 1, 2,
         condition = factor(condition,
                               levels=c(1,2,3),
                               labels=c("blank","congruent","incongruent")))
```

301 Next, we can set up the life tables and plots of the discrete-time functions $h(t)$, $S(t)$,
 302 $ca(t)$, and $P(t)$. To do so using a functional programming approach, one has to nest the
 303 data within participants using the `group_nest()` function, and supply a user-defined
 304 censoring time and bin width to our custom function “`censor()`”, as follows.

```

data_nested <- data_wr %>% group_nest(pid)

data_final <- data_nested %>%
  # ! user input: censoring time, and bin width
  mutate(censored = map(data, censor, 600, 40)) %>%
  # create person-trial-bin data set
  mutate(ptb_data = map(censored, ptb)) %>%
  # create life tables without ca(t)
  mutate(lifetable = map(ptb_data, setup_lt)) %>%
  # calculate ca(t)
  mutate(condacc = map(censored, calc_ca)) %>%
  # create life tables with ca(t)
  mutate(lifetable_ca = map2(lifetable, condacc, join_lt_ca)) %>%
  # create plots
  mutate(plot = map2(.x = lifetable_ca, .y = pid, plot_eha,1))

```

305 Note that the censoring time (here: 600 ms) should be a multiple of the bin width
 306 (here: 40 ms). The censoring time should be a time point after which no informative
 307 responses are expected anymore. In experiments that implement a response deadline in
 308 each trial the censoring time can equal that deadline time point. Trials with a RT larger
 309 than the censoring time, or trials in which no response is emitted during the data collection
 310 period, are treated as right-censored observations in EHA. In other words, these trials are
 311 not discarded, because they contain the information that the event did not occur before the
 312 censoring time. Removing such trials before calculating the mean event time will result in
 313 underestimation of the true mean.

314 The person-trial-bin oriented data set is created by our custom function ptb(), and it
 315 has one row for each time bin (of each trial) that is at risk for event occurrence. The
 316 variable “event” in the person-trial-bin oriented data set indicates whether a response
 317 occurs (1) or not (0) for each bin.

318 The next step is to set up the life table using our custom function setup_lt(),

319 calculate the conditional accuracies using our custom function calc_ca(), add the ca(t)
320 estimates to the life table using our custom function join_lt_ca(), and then plot the
321 descriptive statistics using our custom function plot_eha(). One can now inspect different
322 aspects, including the life table for a particular condition of a particular subject, and a plot
323 of the different functions for a particular participant.

324 In general, it is important to visually inspect the functions first for each participant,
325 in order to identify individuals that may be guessing (e.g., a flat conditional accuracy
326 function at .5 indicates that someone is just guessing), outlying individuals, and/or
327 different groups with qualitatively different behavior. Also, to select a bin width for fitting
328 models, one should test and compare various bin widths in the censor function, and select
329 the smallest one that is supported by the data. Too small bin widths will result in erratic
330 hazard functions because many bins will have estimates equal to zero.

331 Table 3 shows the life table for condition “blank” (no prime stimulus presented) for
332 participant 6.

```
333 ## Warning in attr(x, "align"): 'xfun::attr()' is deprecated.  
334 ## Use 'xfun::attr2()' instead.  
335 ## See help("Deprecated")
```

Table 3

The life table for the blank prime condition of participant 6.

bin	risk_set	events	hazard	se_haz	survival	se_surv	ca	se_ca
0	220	NA	NA	NA	1.00	0.00	NA	NA
40	220	0	0.00	0.00	1.00	0.00	NA	NA
80	220	0	0.00	0.00	1.00	0.00	NA	NA
120	220	0	0.00	0.00	1.00	0.00	NA	NA
160	220	0	0.00	0.00	1.00	0.00	NA	NA
200	220	0	0.00	0.00	1.00	0.00	NA	NA
240	220	0	0.00	0.00	1.00	0.00	NA	NA
280	220	7	0.03	0.01	0.97	0.01	0.29	0.17
320	213	13	0.06	0.02	0.91	0.02	0.77	0.12
360	200	26	0.13	0.02	0.79	0.03	0.92	0.05
400	174	40	0.23	0.03	0.61	0.03	1.00	0.00
440	134	48	0.36	0.04	0.39	0.03	0.98	0.02
480	86	37	0.43	0.05	0.22	0.03	1.00	0.00
520	49	32	0.65	0.07	0.08	0.02	1.00	0.00
560	17	9	0.53	0.12	0.04	0.01	1.00	0.00
600	8	4	0.50	0.18	0.02	0.01	1.00	0.00

Note. The column named “bin” indicates the endpoint of each time bin (in ms), and includes time point zero. For example the first bin is (0,40] with the starting point excluded and the endpoint included. At time point zero, no events can occur and therefore $h(t=0)$ and $ca(t=0)$ are undefined. $se =$ standard error. $ca =$ conditional accuracy. $NA =$ undefined.

337 probability mass functions for each prime condition for participant 6. By using
 338 discrete-time hazard functions of event occurrence – in combination with conditional
 339 accuracy functions for two-choice tasks – one can provide an unbiased, time-varying, and
 340 probabilistic description of the latency and accuracy of responses based on all trials of any
 341 data set.

Descriptive stats for subject 6

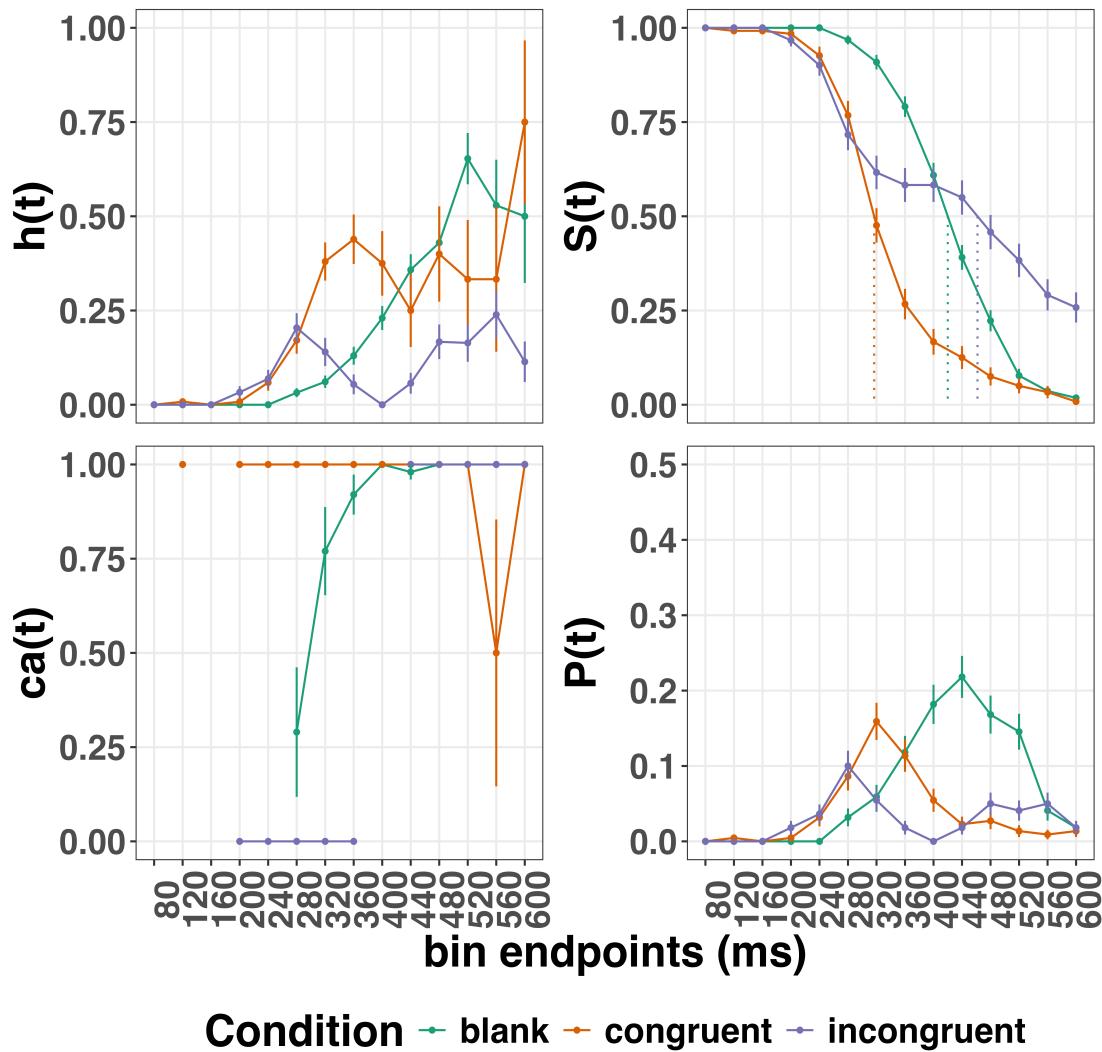


Figure 2. Estimated discrete-time hazard (h), survivor (S), conditional accuracy (ca) and probability mass (P) functions for participant 6. Vertical dotted lines indicate the estimated median RTs. Error bars represent ± 1 standard error of the respective proportion.

342 For example, for participant 6, the estimated hazard values in bin (240,280] are 0.03,

343 0.17, and 0.20 for the blank, congruent, and incongruent prime conditions, respectively. In

344 other words, when the waiting time has increased until *240 ms* after target onset, then the

345 conditional probability of response occurrence in the next 40 ms is more than five times

346 larger for both prime-present conditions, compared to the blank prime condition.

347 Furthermore, the estimated conditional accuracy values in bin (240,280] are 0.29, 1,

348 and 0 for the blank, congruent, and incongruent prime conditions, respectively. In other

349 words, if a response is emitted in bin (240,280], then the probability that it is correct is

350 estimated to be 0.29, 1, and 0 for the blank, congruent, and incongruent prime conditions,

351 respectively.

352 However, when the waiting time has increased until *400 ms* after target onset, then

353 the conditional probability of response occurrence in the next 40 ms is estimated to be

354 0.36, 0.25, and 0.06 for the blank, congruent, and incongruent prime conditions,

355 respectively. And when a response does occur in bin (400,440], then the probability that it

356 is correct is estimated to be 0.98, 1, and 1 for the blank, congruent, and incongruent prime

357 conditions, respectively.

358 These distributional results suggest that participant 6 is initially responding to the

359 prime even though (s)he was instructed to only respond to the target, that response

360 competition emerges in the incongruent prime condition around 300 ms, and that only

361 slower responses are fully controlled by the target stimulus. Qualitatively similar results

362 were obtained for the other five participants. When participants show qualitatively similar

363 distributional patterns, one might consider aggregating their data and plotting the

364 group-average distribution per condition (see Tutorial_1a.Rmd).

365 In general, these results go against the (often implicit) assumption in research on

366 priming that all observed responses are primed responses to the target stimulus. Instead,

367 the distributional data show that early responses are triggered exclusively by the prime

368 stimulus, while only later responses reflect primed responses to the target stimulus.

369 At this point, we have calculated, summarised and plotted descriptive statistics for
370 the key variables in EHA/SAT. As we will show in later Tutorials, statistical models for
371 $h(t)$ and $ca(t)$ can be implemented as generalized linear mixed regression models predicting
372 event occurrence (1/0) and conditional accuracy (1/0) in each bin of a selected time
373 window for analysis. But first we consider calculating the descriptive statistics for two
374 independent variables.

375 **3.2 Tutorial 1b: Generalising to a more complex design**

376 So far in this paper, we have used a simple experimental design, which involved one
377 condition with three levels. But psychological experiments are often more complex, with
378 crossed factorial designs and/or conditions with more than three levels. The purpose of
379 Tutorial 1b, therefore, is to provide a generalisation of the basic approach, which extends
380 to a more complicated design. We felt that this might be useful for researchers in
381 experimental psychology that typically use crossed factorial designs.

382 To this end, Tutorial 1b illustrates how to calculate and plot the descriptive statistics
383 for the full data set of Experiment 1 of Panis and Schmidt (2016), which includes two
384 independent variables: mask type and prime type. As we use the same functional
385 programming approach as in Tutorial 1a, we simply present the sample-based functions for
386 each participant as part of Tutorial_1b.Rmd for those that are interested.

387 **3.3 Tutorial 2a: Fitting Bayesian hazard models to discrete time-to-event data**

388 In this third tutorial, we illustrate how to fit Bayesian multilevel regression models to
389 the RT data of the masked response priming data used in Tutorial 1a. Fitting (Bayesian or
390 non-Bayesian) regression models to time-to-event data is important when you want to
391 study how the shape of the hazard function depends on various predictors (Singer &

392 Willett, 2003).

393 In general, when fitting regression models, our lab adopts an estimation approach to
394 multilevel regression (Kruschke & Liddell, 2018; Winter, 2019), which is heavily influenced
395 by the Bayesian framework as suggested by Richard McElreath (Kurz, 2023b; McElreath,
396 2020). We also use a “keep it maximal” approach to specifying varying (or random) effects
397 (Barr, Levy, Scheepers, & Tily, 2013). This means that wherever possible we include
398 varying intercepts and slopes per participant. To make inferences, we use two main
399 approaches. We compare models of different complexity, using information criteria (e.g.,
400 WAIC) and cross-validation (e.g., LOO), to evaluate out-of-sample predictive accuracy
401 (McElreath, 2020). We also take the most complex model and evaluate key parameters of
402 interest using point and interval estimates.

403 **3.3.1 Hazard model considerations.** There are several analytic decisions one
404 has to make when fitting a discrete-time hazard model. First, one has to select an analysis
405 time window, i.e., a contiguous set of bins for which there is enough data for each
406 participant. Second, given that the dependent variable (event occurrence) is binary, one
407 has to select a link function (see section C in the Supplemental Material). The cloglog link
408 is preferred over the logit link when events can occur in principle at any time point within
409 a bin, which is the case for RT data (Singer & Willett, 2003). Third, one has to choose
410 whether to treat TIME (i.e., the time bin index t) as a categorical or continuous predictor.
411 And when you treat a variable as a categorical predictor, you can choose between reference
412 coding and index coding. With reference coding, one defines the variable as a factor and
413 selects one of the k categories as the reference level. Brm() will then construct $k-1$
414 indicator variables (see model M1d in Tutorial_2a.Rmd for an example). With index
415 coding, one constructs an index variable that contains integers that correspond to different
416 categories (see models M0i and M1i below). As explained by McElreath (2020), the
417 advantage of index coding is that the same prior can be assigned to each level of the index
418 variable, so that each category has the same prior uncertainty.

In the case of a large- N design without repeated measurements, the parameters of a discrete-time hazard model can be estimated using standard logistic regression software after expanding the typical person-trial data set into a person-trial-bin data set (Allison, 2010). When there is clustering in the data, as in the case of a small- N design with repeated measurements, the parameters of a discrete-time hazard model can be estimated using population-averaged methods (e.g., Generalized Estimating Equations), and Bayesian or frequentist generalized linear mixed models (Allison, 2010).

In general, there are three assumptions one can make or relax when adding experimental predictor variables and other covariates: The linearity assumption for continuous predictors (the effect of a 1 unit change is the same anywhere on the scale), the additivity assumption (predictors do not interact), and the proportionality assumption (predictors do not interact with TIME).

In tutorial_2a.Rmd we fit several Bayesian multilevel models (i.e., generalized linear mixed models) that differ in complexity to the person-trial-bin oriented data set that we created in Tutorial 1a. We decided to select the analysis time window (200,600] and the cloglog link. Below, we shortly discuss two of these models. The person-trial-bin data set is prepared as follows.

```
# read in the file we saved in tutorial 1a
ptb_data <- read_csv("Tutorial_1_descriptive_stats/data/inputfile_hazard_modeling.csv")

ptb_data <- ptb_data %>%
  # select analysis time range: (200,600] with 10 bins (time bin ranks 6 to 15)
  filter(period > 5) %>%
  # define categorical predictor TIME as index variable named timebin
  mutate(timebin = factor(period, levels = c(6:15)),
  # factor "condition" using reference coding, with "blank" as the reference level
  condition = factor(condition, labels = c("blank", "congruent", "incongruent")),
  # categorical predictor "prime" with index coding
```

```
prime = ifelse(condition=="blank", 1, ifelse(condition=="congruent", 2, 3)),
prime = factor(prime, levels = c(1,2,3)))
```

436 3.3.2 Prior distributions. To get the posterior distribution of each model

437 parameter given the data, we need to specify prior distributions for the model parameters
438 which reflect our prior beliefs. In Tutorial_2a.Rmd we perform a few prior predictive
439 checks to make sure our selected prior distributions reflect our prior beliefs (Gelman,
440 Vehtari, et al., 2020).

441 The middle column of Supplementary Figure 2 (section E of the Supplemental
442 Material) shows six examples of prior distributions for an intercept on the logit and/or
443 cloglog scales. While a normal distribution with relatively large variance is often used as a
444 weakly informative prior for continuous dependent variables, rows A and B of
445 Supplementary Figure 2 show that specifying such distributions on the logit and cloglog
446 scales actually leads to rather informative distributions on the original probability scale, as
447 most mass is pushed to probabilities of 0 and 1.

448 3.3.3 Model M0i: A null model with index coding. When you do not want to
449 make assumptions about the shape of the hazard function, or its shape is not smooth but
450 irregular, then you can use a general specification of TIME, i.e., fit one grand intercept per
451 time bin. In this first model, we use a general specification of TIME using index coding,
452 and do not include experimental predictors. We call this model “M0i”.

453 Before we fit model M0i, we select the necessary columns from the data, and specify
454 our priors. In the code of Tutorial 2a, model M0i is specified as follows.

```
model_M0i <-
  brm(data = data_M0i,
       family = bernoulli(link="cloglog"),
       formula = event ~ 0 + timebin + (0 + timebin | pid),
```

```

prior = priors_M0i,
chains = 4, cores = 4,
iter = 3000, warmup = 1000,
control = list(adapt_delta = 0.999,
                step_size = 0.04,
                max_treedepth = 12),
seed = 12, init = "0",
file = "Tutorial_2_Bayesian/models/model_M0i")

```

455 After selecting the bernoulli family and the cloglog link, the model formula is
 456 specified. The specification “0 + …” removes the default intercept in brm(). The fixed
 457 effects include an intercept for each level of timebin. Each of these intercepts is allowed to
 458 vary across individuals (variable pid). We request 2000 samples from the posterior
 459 distribution for each of four chains. Estimating model M0i took about 30 minutes on a
 460 MacBook Pro (Sonoma 14.6.1 OS, 18GB Memory, M3 Pro Chip).

461 **3.3.4 Model M1i: Adding the effects of prime-target congruency.** Previous
 462 research has shown that psychological effects typically change over time (Panis, 2020;
 463 Panis, Moran, et al., 2020; Panis & Schmidt, 2022; Panis et al., 2017; Panis & Wagemans,
 464 2009). In the next model, therefore, we use index coding for both TIME (variable
 465 “timebin”) and the categorical predictor prime-target-congruency (variable “prime”), so
 466 that we get 30 grand intercepts, one for each combination of timebin level and prime level.
 467 Here is the model formula of this model that we call “M1i”.

```
event ~ 0 + timebin:prime + (0 + timebin:prime | pid)
```

468 Estimating model M1i took about 124 minutes.

469 **3.3.5 Compare the models.** We can compare the two models using the Widely
 470 Applicable Information Criterion (WAIC) and Leave-One-Out (LOO) cross-validation, and

471 look at model weights for both criteria (Kurz, 2023a; McElreath, 2020).

```
model_weights(model_M0i, model_M1i, weights = "loo") %>% round(digits = 2)
```

472 ## model_M0i model_M1i

473 ## 0 1

```
model_weights(model_M0i, model_M1i, weights = "waic") %>% round(digits = 2)
```

474 ## model_M0i model_M1i

475 ## 0 1

476 Clearly, both the loo and waic weighting schemes assign a weight of 1 to model M1i,
477 and a weight of 0 to the other simpler model.

478 **3.3.6 Evaluating parameter estimates in model M1i.** To make inferences
479 from the parameter estimates in model M1i, we first plot the densities of the draws from
480 the posterior distributions of its population-level parameters in Figure 5, together with
481 point (median) and interval estimates (80% and 95% credible intervals).

Posterior distributions for population-level effects in Model M1i

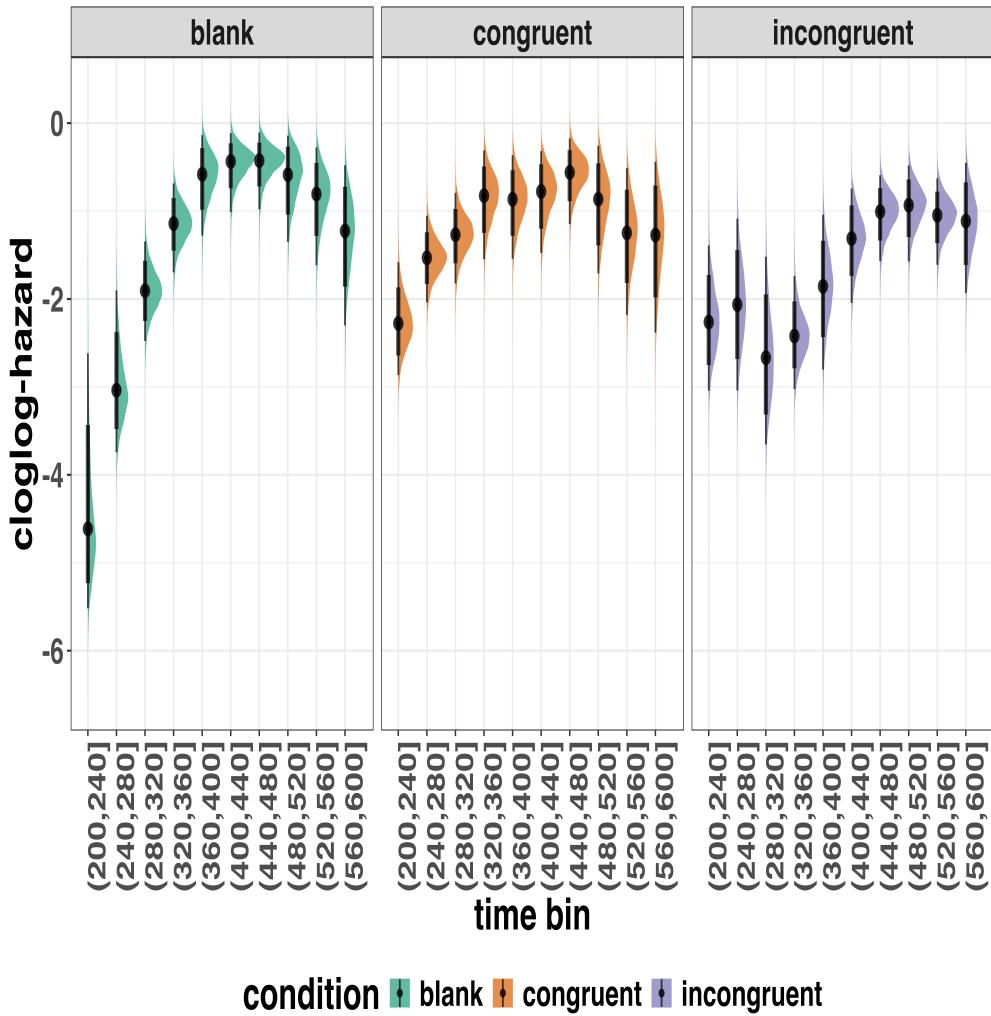


Figure 3. Medians and 80/95% credible intervals of the posterior distributions of the population-level parameters of model M1i.

482 Because the parameter estimates are on the cloglog-hazard scale, we can ease our
 483 interpretation by plotting the expected value of the posterior predictive distribution – the
 484 predicted hazard values – at the population level (Figure 6A), and for each participant in
 485 the data set (Figure 6B).

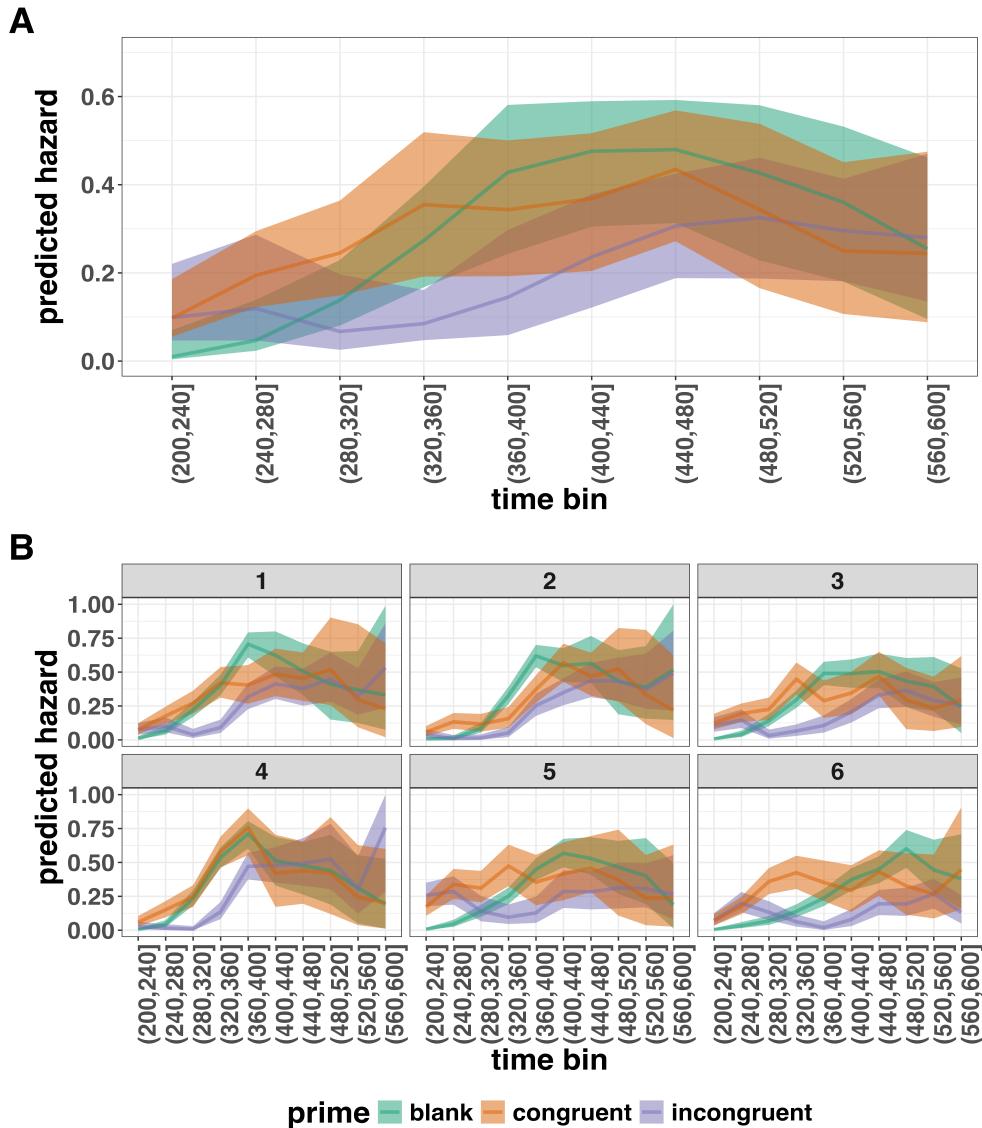


Figure 4. Point (median) and 80/95% credible interval summaries of the hazard estimates (expected values of the draws from the posterior predictive distributions) in each time bin at the population level (A), and for each participant (B).

As we are actually interested in the effects of congruent and incongruent primes,

relative to the blank prime condition, we can construct two contrasts (congruent-blank, incongruent-blank), and plot the posterior distributions of these contrast effects, both at the population level (Figure 7A; grand average marginal effect) and at the participant level

490 (Figure 7B; subject-specific average marginal effect).

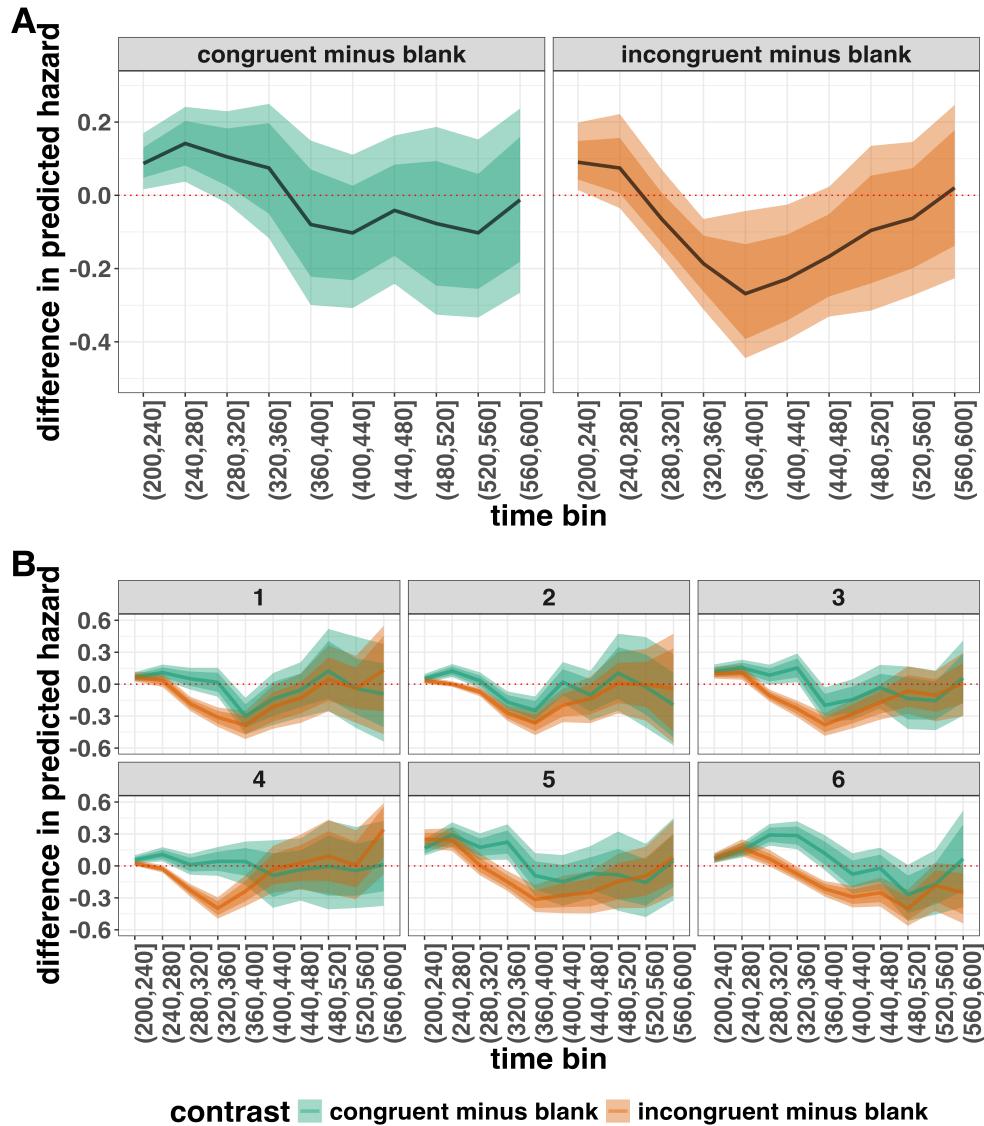


Figure 5. Point (mean) and 80/95% credible interval summaries of estimated differences in hazard in each time bin at the population level (A), and for each participant (B).

491 The point estimates and quantile intervals can be reported in a table (see
 492 Tutorial_2a.Rmd for details).

493 **Example conclusions for M1i.** What can we conclude from model M1i about
 494 our research question, i.e., the temporal dynamics of the effect of prime-target congruency

495 on RT? In other words, in which of the 40-ms time bins between 200 and 600 ms after
496 target onset does changing the prime from blank to congruent or incongruent affect the
497 hazard of response occurrence (for a prime-target stimulus-onset-asynchrony of 187 ms)?

498 If we want to estimate the population-level effect of prime type on hazard, we can
499 base our conclusion on Figure 7A. The contrast “congruent minus blank” was estimated to
500 be 0.09 hazard units in bin (200,240] (95% CrI = [0.02, 0.17]), and 0.14 hazard units in bin
501 (240,280]) (95% CrI = [0.04, 0.25]). For the other bins, the 95% credible interval contained
502 zero. The contrast “incongruent minus blank” was estimated to be 0.09 hazard units in bin
503 (200,240] (95% CrI = [0.01, 0.21]), -0.19 hazard units in bin (320,360] (95% CrI = [-0.31,
504 -0.06]), -0.27 hazard units in bin (360,400] (95% CrI = [-0.45, -0.04]), and -0.23 hazard
505 units in bin (400,440] (95% CrI = [-0.40, -0.03]). For the other bins, the 95% credible
506 interval contained zero.

507 There are thus two phases of performance for the average person between 200 and
508 600 ms after target onset. In the first phase, the addition of a congruent or incongruent
509 prime stimulus increases the hazard of response occurrence compared to blank prime trials
510 in the time period (200, 240]. In the second phase, only the incongruent prime decreases
511 the hazard of response occurrence compared to blank primes, in the time period (320,440].
512 The sign of the effect of incongruent primes on the hazard of response occurrence thus
513 depends on how much waiting time has passed since target onset.

514 If we want to focus more on inter-individual differences, we can study the
515 subject-specific hazard functions in Figure 7B. Note that three participants (1, 2, and 3)
516 show a negative difference for the contrast “congruent minus incongruent” in bin (360,400]
517 – subject 2 also in bin (320,360].

518 Future studies could (a) increase the number of participants to estimate the
519 proportion of “dippers” in the subject population, and/or (b) try to explain why this dip
520 occurs. For example, Panis and Schmidt (2016) concluded that active, top-down,

521 task-guided response inhibition effects emerge around 360 ms after the onset of the stimulus
522 following the prime (here: the target stimulus). Such a top-down inhibitory effect might
523 exist in our priming data set, because after some time participants will learn that the first
524 stimulus is not the one they have to respond to. To prevent a premature overt response to
525 the prime they thus might gradually increase a global response threshold during the
526 remainder of the experiment, which could result in a lower hazard in congruent trials
527 compared to blank trials, for bins after ~360 ms, and towards the end of the experiment.
528 This effect might be masked for incongruent primes by the response competition effect.

529 Interestingly, all subjects show a tendency in their mean difference (congruent minus
530 blank) to “dip” around that time (Figure 7B). Therefore, future modeling efforts could
531 incorporate the trial number into the model formula, in order to also study how the effects
532 of prime type on hazard change on the long experiment-wide time scale, next to the short
533 trial-wide time scale. In Tutorial_2a.Rmd we provide a number of model formulae that
534 should get you going.

535 3.4 Tutorial 2b: Fitting Bayesian conditional accuracy models

536 In this fourth tutorial, we illustrate how to fit a Bayesian multilevel regression model
537 to the timed accuracy data from the masked response priming data used in Tutorial 1a.
538 The general process is similar to Tutorial 2a, except that (a) we use the person-trial data,
539 (b) we use the logit link function, and (c) we change the priors. To keep the tutorial short,
540 we only fit one conditional accuracy model, which was based on model M1i from Tutorial
541 2a and labelled M1i_ca.

542 To make inferences from the parameter estimates in model M1i_ca, we first plot the
543 densities of the draws from the posterior distributions of its population-level parameters in
544 Figure 8, together with point (median) and interval estimates (80% and 95% credible
545 intervals).

Posterior distributions for population-level effects in Model M1i_ca

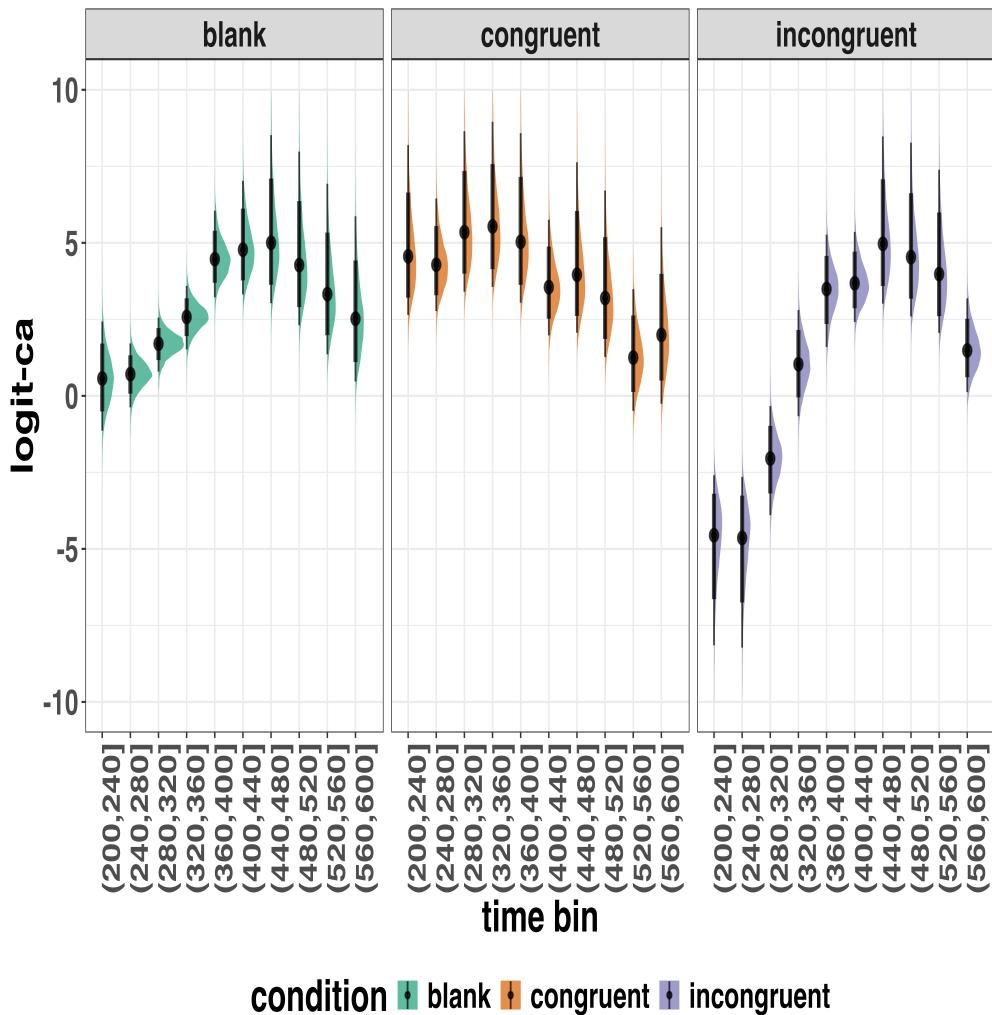


Figure 6. Medians and 80/95% credible intervals of the posterior distributions of the population-level parameters of model M1i_ca. ca = conditional accuracy.

Because the parameter estimates are on the logit-ca scale, we can ease our

interpretation by plotting the expected value of the posterior predictive distribution – the predicted conditional accuracies – at the population level (Figure 9A), and for each participant in the data set (Figure 9B).

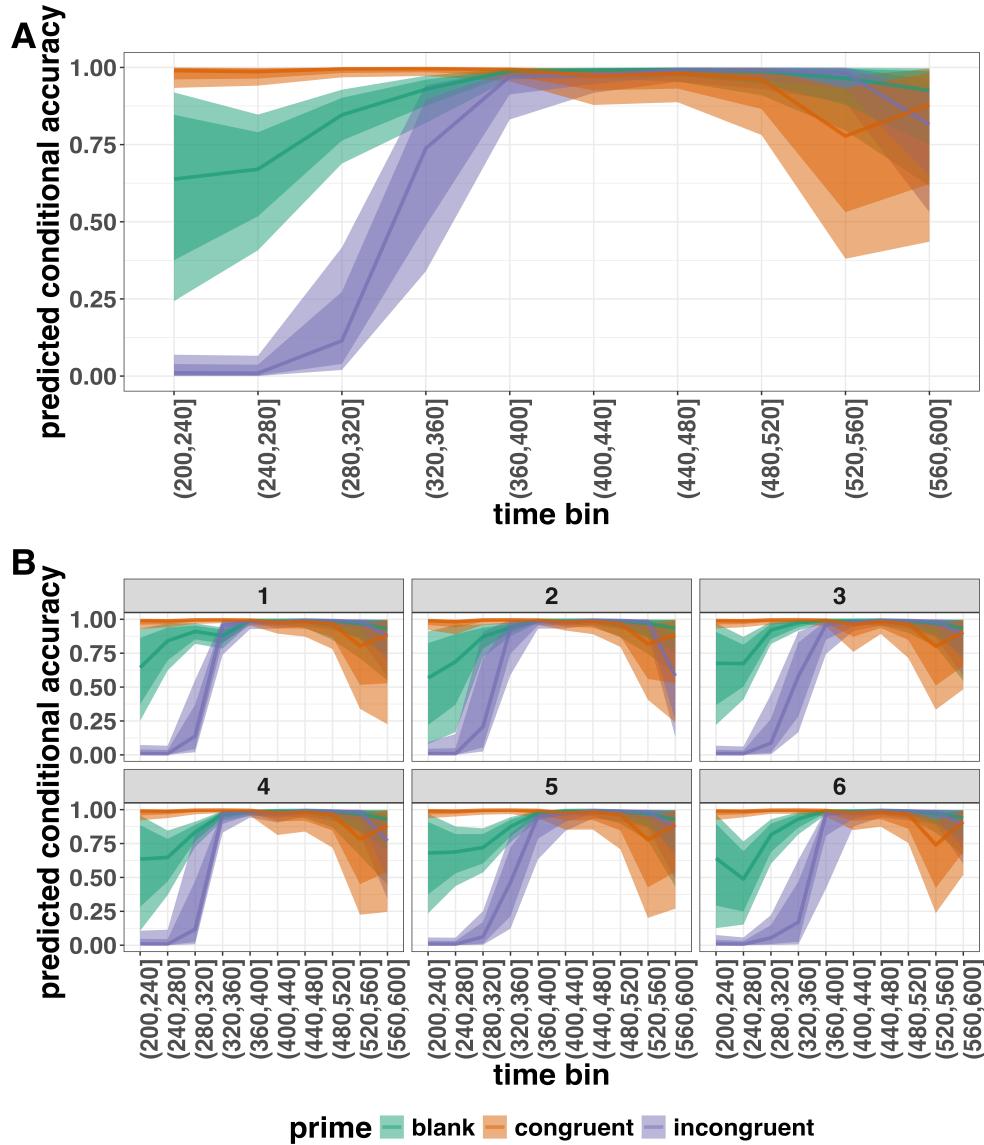


Figure 7. Point (median) and 80/95% credible interval summaries of the conditional accuracy estimates (expected values of the draws from the posterior predictive distributions) in each time bin at the population level (A), and for each participant (B).

550 As we are actually interested in the effects of congruent and incongruent primes,

551 relative to the blank prime condition, we can construct two contrasts (congruent-blank,
 552 incongruent-blank), and plot the posterior distributions of these contrast effects at the
 553 population level (Figure 10A) and for each participant (Figure 10B).

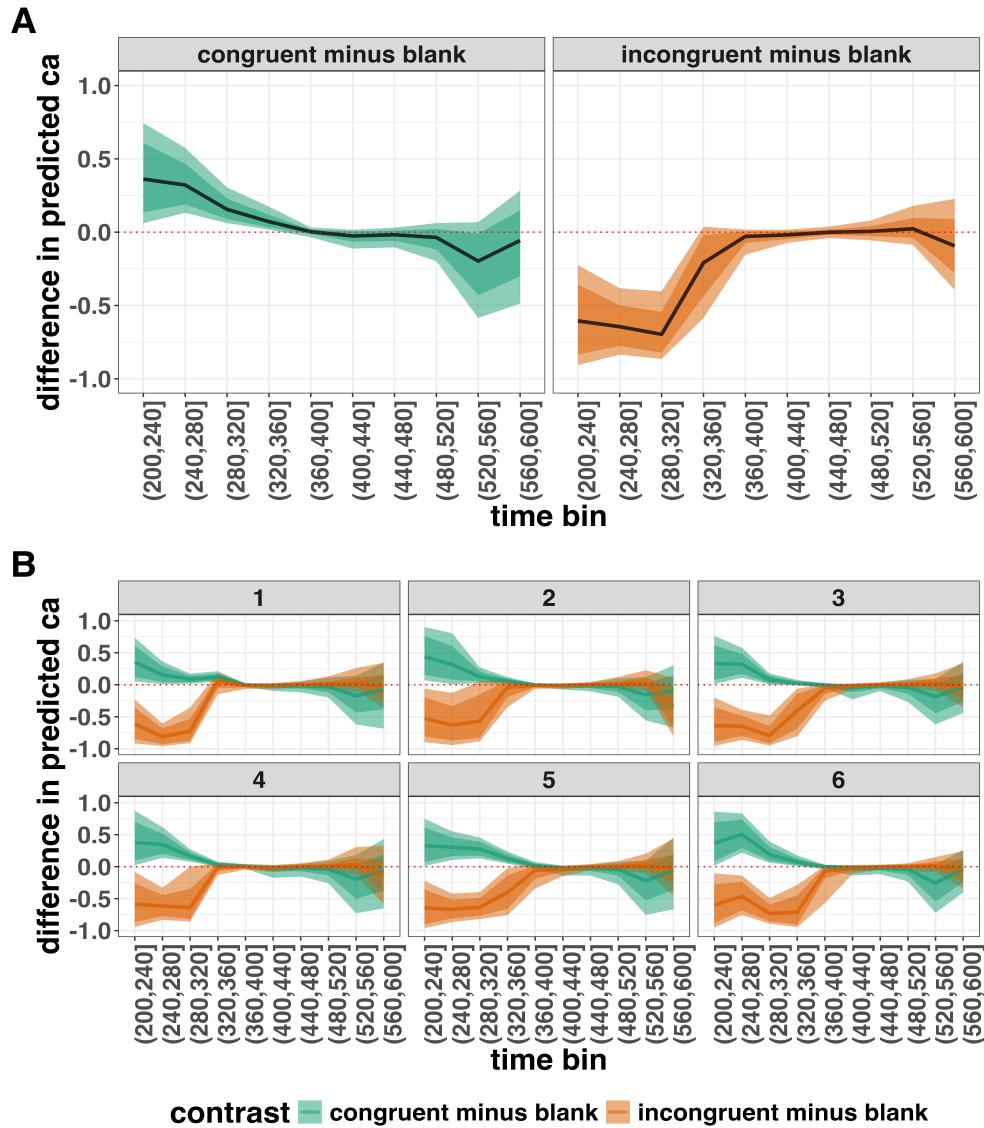


Figure 8. Point (mean) and 80/95% credible interval summaries of estimated differences in conditional accuracy in each time bin at the population level (A), and for each participant (B).

554 Based on Figure 10A we see that on the population level congruent primes have a
 555 positive effect on the conditional accuracy of emitted responses in time bins (200,240],
 556 (240,280], (280,320], and (320,360], relative to the estimates in the baseline condition
 557 (blank prime; red dashed lines in Figure 10A). Incongruent primes have a negative effect on

558 the conditional accuracy of emitted responses in the first time bins, relative to the
559 estimates in the baseline condition.

560 **3.5 Tutorial 4: Planning**

561 In the final tutorial, we look at planning a future experiment, which uses EHA.

562 **3.5.1 Background.** The general approach to planning that we adopt here involves
563 simulating reasonably structured data to help guide what you might be able to expect from
564 your data once you collect it (Gelman, Vehtari, et al., 2020). The basic structure and code
565 follows the examples outlined by Solomon Kurz in his ‘power’ blog posts
566 (<https://solomonkurz.netlify.app/blog/bayesian-power-analysis-part-i/>) and Lisa
567 DeBruine’s R package `faux{}` (<https://debruine.github.io/faux/>) as well as these related
568 papers (DeBruine & Barr, 2021; Pargent, Koch, Kleine, Lermer, & Gaube, 2024).

569 **3.5.2 Basic workflow.** The basic workflow is as follows:

- 570 1. Fit a regression model to existing data.
- 571 2. Use the regression model parameters to simulate new data.
- 572 3. Write a function to create 1000s of datasets and vary parameters of interest (e.g.,
573 sample size, trial count, effect size).
- 574 4. Summarise the simulated data to estimate likely power or precision of the research
575 design options.

576 Ideally, in the above workflow, we would also fit a model to each dataset and
577 summarise the model output, rather than the raw data. However, when each model takes
578 several hours to build, and we may want to simulate many 1000s of datasets, it can be
579 computationally demanding for desktop machines. So, for ease, here we just use the raw
580 simulated datasets to guide future expectations.

581 In the below, we only provide a high-level summary of the process and let readers
582 dive into the details within the tutorial should they feel so inclined.

3.5.3 Fit a regression model and simulate one dataset.

We again use the data from Panis and Schmidt (2016) to provide a worked example. We fit an index coding model on a subset of time bins (six time bins in total) and for two prime conditions (congruent and incongruent). We chose to focus on a subsample of the data to ease the computational burden. We also used a full varying effects structure, with the model formula as follows:

```
event ~ 0 + timebin:prime + (0 + timebin:prime | pid)
```

We then took parameters from this model and used them to create a single dataset

with 200 trials per condition for 10 individual participants. The raw data and the simulated data are plotted in Figure 12 and show quite close correspondence, which is re-assuring. But, this is only one dataset. What we really want to do is simulate many datasets and vary parameters of interest, which is what we turn to in the next section.

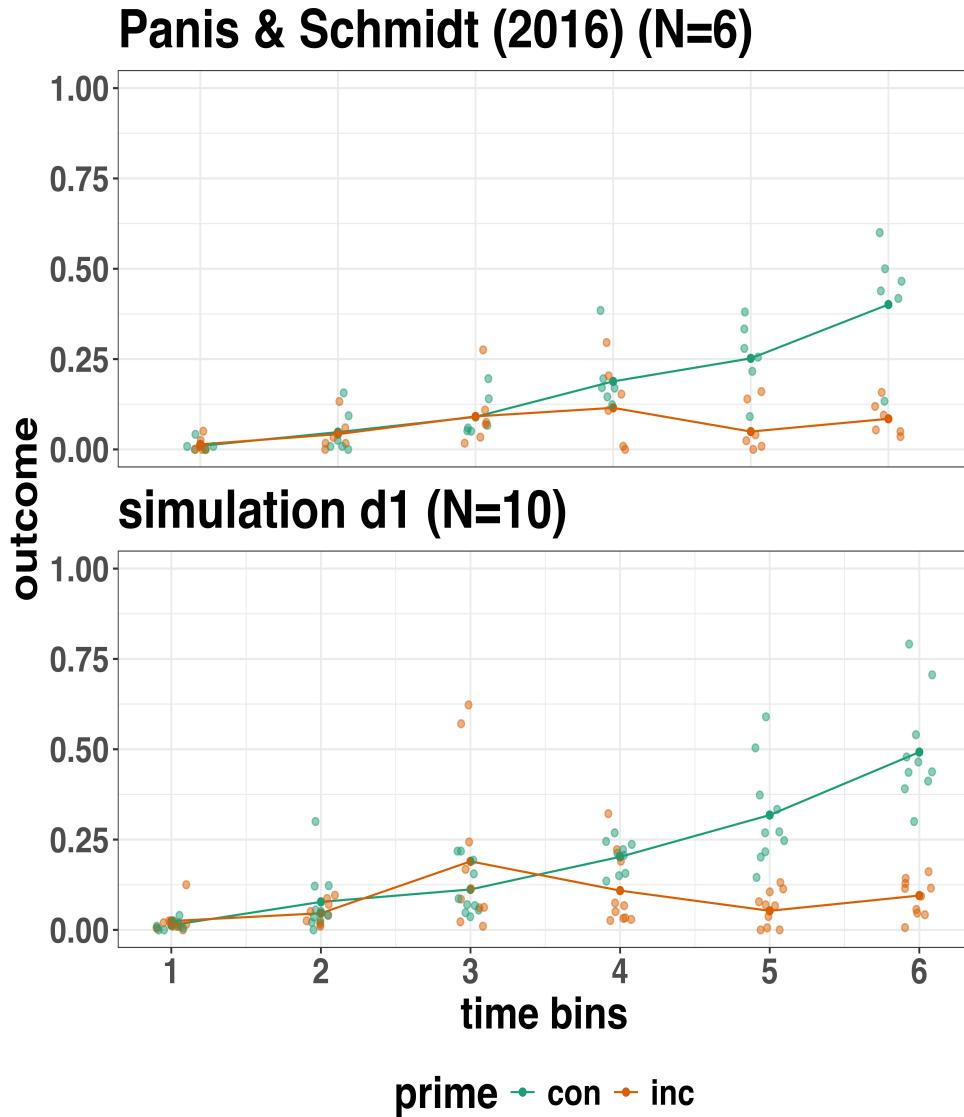


Figure 9. Raw data from Panis and Schmidt (2016) and simulated data from 10 participants.

594 3.5.4 Simulate and summarise data across a range of parameter values.

595 Here we use the same data simulation process as used above, but instead of simulating one
 596 dataset, we simulate 1000 datasets per variation in parameter values. Specifically, in
 597 Simulation 1, we vary the number of trials per condition (100, 200, and 400), as well as the
 598 effect size in bin 6. We focus on bin 6 only, in terms of varying the effect size, just to make
 599 things simpler and easier to understand. The effect size observed in bin 6 in this subsample

of data was a 79% reduction in hazard value from the congruent prime (0.401 hazard value) to the incongruent prime condition (0.085 hazard value). In other words, a hazard ratio of 0.21 (e.g., $0.085/0.401 = 0.21$). As a starting point, we chose three effect sizes, which covered a fairly broad range of hazard ratios (0.25, 0.5, 0.75), which correspond to a 75%, 50% and 25% reduction in hazard value as a function of prime condition.

Summary results from Simulation 1 are shown in Figure 13A. Figure 13A depicts statistical “power” as calculated by the percentage of lower-bound 95% confidence intervals that exclude zero when the difference between prime condition is calculated (congruent - incongruent). In other words, what fraction of the simulated datasets generated an effect of prime that excludes the criterion mark of zero. We are aware that “power” is not part of a Bayesian analytical workflow, but we choose to include it here, as it is familiar to most researchers in experimental psychology.

The results of Simulation 1 show that if we were targeting an effect size similar to the one reported in the original study, then testing 10 participants and collecting 100 trials per condition would be enough to provide over 95% power. However, we could not be as confident about smaller effects, such as a hazard ratio of 50% or 25%. From this simulation, we can see that somewhere between an effect size of a 50% and 75% reduction in hazard value, power increases to a range that most researchers would consider acceptable (i.e., >95% power). To probe this space a little further, we decided to run a second simulation, which varied different parameters.

In Simulation 2, we varied the effect size between a different range of values (0.5, 0.4, 0.3), which correspond to a 50%, 60% and 70% reduction in hazard value as a function of prime condition. In addition, we varied the number of participants per experiment between 10, 15, and 20 participants. Given that trial count per condition made little difference to power in Simulation 1, we fixed trial count at 200 trials per condition in Simulation 2. Summary results from Simulation 2 are shown in Figure 13B. A summary of these power

626 calculations might be as follows (trial count = 200 per condition in all cases):

- 627 • For a 70% reduction (0.3 hazard ratio), N=10 would give nearly 100% power.
- 628 • For a 60% reduction (0.4 hazard ratio), N=10 would give nearly 90% power.
- 629 • For a 50% reduction (0.5 hazard ratio), N=15 would give over 80% power.

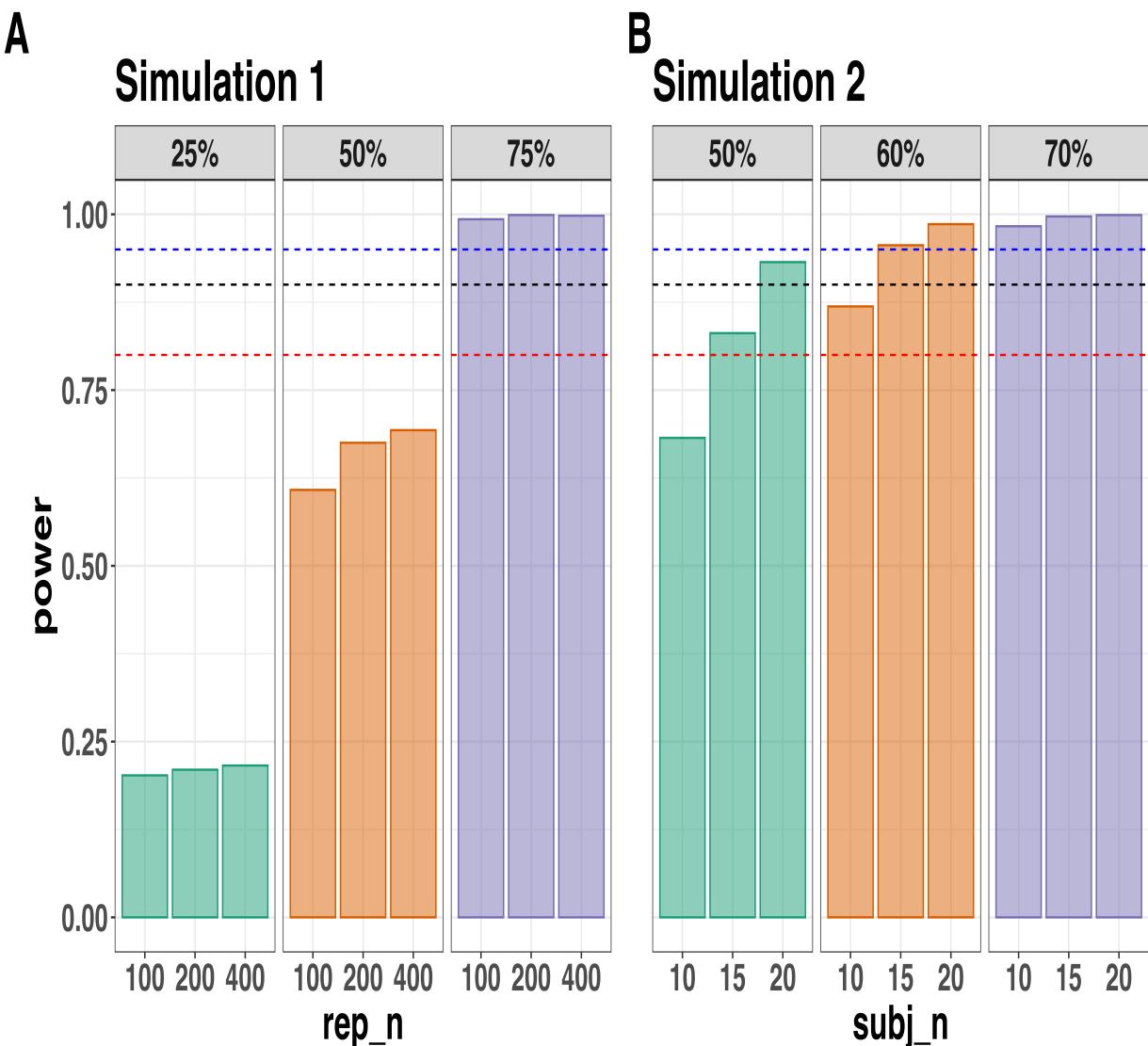


Figure 10. Statistical power across data Simulation 1 (A) and Simulation 2 (B). Power was calculated as the percentage of lower-bound 95% confidence intervals that exclude zero when the difference between prime condition is calculated (congruent - incongruent). In Simulation 1, the effect size was varied between a 25%, 50% and 75% reduction in hazard value, whereas the trial count was varied between 100, 200 and 400 trials per condition (the number of participants was fixed at N=10). In Simulation 2, the effect size was varied between a 50%, 60% and 70% reduction in hazard value, whereas the number of participants was varied between N=10, 15 and 20 (the number of trials per condition was fixed at 200). The dashed lines represent 80% (red), 90% (black) and 95% (blue) power. Abbreviations: rep_n = the number of trials per experimental condition; subj_n = the number of participants per simulated experiment.

630 **3.5.5 Planning decisions.** Now that we have summarised our simulated data,

631 what planning decisions could we make about a future study? More concretely, how many

632 trials per condition should we collect and how many participants should we test? Like

633 almost always when planning future studies, the answer depends on your objectives, as well

634 as the available resources (Lakens, 2022). There is no straightforward and clear-cut answer.

635 Some considerations might be as follows:

- 636 • How much power or precision are you looking to obtain in this particular study?

- 637 • Are you running multiple studies that have some form of replication built in?

- 638 • What level of resources do you have at your disposal, such as time, money and

639 personnel?

- 640 • How easy or difficult is it to obtain the specific type of sample?

641 If we were running this kind of study in our lab, what would we do? We might pick a

642 hazard ratio of 0.4 or 0.5 as a target effect size since this is much smaller than that

643 observed previously (Panis & Schmidt, 2016). Then we might pick the corresponding

644 combination of trial count per condition (e.g., 200) and participant sample size (e.g., N=10

645 or N=15) that takes you over the 80% power mark. If we wanted to maximise power based

646 on these simulations, and we had the time and resources available, then we would test

647 N=20 participants, which would provide >90% power for an effect size of 0.5.

648 **But**, and this is an important “but”, unless there are unavoidable reasons, no matter

649 what planning choices we made based on these data simulations, we would not solely rely

650 on data collected from one single study. Instead, we would run a follow-up experiment that

651 replicates and extends the initial result. By doing so, we would aim to avoid the Cult of

652 the Isolated Single Study (Nelder, 1999; Tong, 2019), and thus reduce the reliance on any

653 one type of planning tool, such as a power analysis. Then, we would look for common

654 patterns across two or more experiments, rather than trying to make the case that a single

655 study on its own has sufficient evidential value to hit some criterion mark.

656

4. Discussion

657 This main motivation for writing this paper is the observation that EHA and SAT
658 analysis remain under-used in psychological research. As a consequence, the field of
659 psychological research is not taking full advantage of the many benefits EHA/SAT provides
660 compared to more conventional analyses. By providing a freely available set of tutorials,
661 which provide step-by-step guidelines and ready-to-use R code, we hope that researchers
662 will feel more comfortable using EHA/SAT in the future. Indeed, we hope that our
663 tutorials may help to overcome a barrier to entry with EHA/SAT, which is that such
664 approaches require more analytical complexity compared to mean-average comparisons.
665 While we have focused here on within-subject, factorial, small- N designs, it is important to
666 realize that EHA/SAT can be applied to other designs as well (large- N designs with only
667 one measurement per subject, between-subject designs, etc.). As such, the general workflow
668 and associated code can be modified and applied more broadly to other contexts and
669 research questions. In the following, we discuss the main use-cases, issues relating to model
670 complexity and interpretability, as well as limitations of the approach and future
671 extensions.

672 **4.1 What are the main use-cases of EHA for understanding cognition and brain**
673 **function?**

674 For those researchers, like ourselves, who are primarily interested in understanding
675 human cognitive and brain systems, we consider two broadly-defined, main use-cases of
676 EHA. First, as we hope to have made clear by this point, EHA is one way to investigating
677 a “temporal states” approach to cognitive processes. EHA provides one way to uncover the
678 microgenesis of cognitive effects, by revealing when cognitive states may start and stop,
679 how states are replaced with others, as well as what they may be tied to or interact with.
680 Therefore, if your research questions concern **when psychological states occur, and**

681 **how they are temporally organized**, our EHA tutorials could be useful tools for you to
682 use.

683 Second, even if you are not primarily interested in studying the temporal
684 organization of cognitive states, EHA could still be a useful tool to consider using, in order
685 to qualify inferences that are being made based on comparisons between means. Given that
686 distinctly different inferences can be made from the same data based on whether one
687 computes a mean across trials or a RT distribution of events (Figure 1), it may be
688 important for researchers to supplement comparisons between means with EHA. One could
689 envisage scenarios where the implicit assumption of an effect in mean RTs manifesting
690 across all of the time bins measured would not be supported by EHA. Therefore, the
691 conclusion of interest would not apply to all responses, but instead it would be restricted to
692 certain periods of within-trial time.

693 **4.2 Model complexity versus interpretability**

694 EHA can quickly become very complex when adding more than one time scale, due to
695 the many possible higher-order interactions. For example, some of the models discussed in
696 Tutorial 2a, which we did not focus on in the main text, contain two time scales as
697 covariates: the passage of time on the within-trial time scale, and the passage of time on
698 the across-trial (or within-experiment) time scale. However, when trials are presented in
699 blocks, and blocks of trials within sessions, and when the experiment comprises three
700 sessions, then four time scales can be defined (within-trial, within-block, within-session,
701 and within-experiment). From a theoretical perspective, adding more than one time scale –
702 and their interactions – can be important to capture plasticity and other learning effects
703 that may play out on such longer time scales, and that are probably present in each
704 experiment in general (REF DFT). From a practical perspective, therefore, some choices
705 need to be made to balance the amount of data that is being collected per participant,
706 condition and across the varying timescales. As one example, if there are several timescales

707 of relevance, then it might be prudent for interpretational purposes to limit the number of
708 experimental predictor variables (conditions). This is of course where planning and data
709 simulation efforts would be important to provide a guide to experimental design choices
710 (see Tutorial 4).

711 **4.3 Limitations**

712 Compared to the orthodox method – comparing means between conditions – the
713 most important limitation of multilevel hazard and conditional accuracy modeling is that it
714 might take a long time to estimate the parameters using Bayesian methods or the model
715 might have to be simplified significantly to use frequentist methods.

716 Another issue is that you need a relatively large number of trials per condition to
717 estimate the hazard function with high temporal resolution, which is required when testing
718 predictions of process models of cognition. Indeed, in general, there is a trade-off between
719 the number of trials per condition and the temporal resolution (i.e., bin width) of the
720 discrete-time hazard function. Therefore, we recommend researchers to collect as many
721 trials as possible per experimental condition, given the available resources and considering
722 the participant experience (e.g., fatigue and boredom). For instance, if the maximum
723 session length deemed reasonable is between 1 and 2 hours, what is the maximum number
724 of trials per condition that you could reasonably collect? After consideration, it might be
725 worth conducting multiple testing sessions per participant and/or reducing the number of
726 experimental conditions. Finally, there is a user-friendly online tool for calculating
727 statistical power as a function of the number of trials as well as the number of participants,
728 and this might be worth consulting to guide the research design process (Baker et al., 2021).

729

5. Conclusions

730 Estimating the temporal distributions of RT and accuracy provide a rich source of
731 information on the time course of cognitive processing, which have been largely
732 undervalued in the history of experimental psychology and cognitive neuroscience. We
733 hope that by providing a set of hands-on, step-by-step tutorials, which come with
734 custom-built and freely available code, researchers will feel more comfortable embracing
735 EHA and investigating the shape of empirical hazard functions and the temporal profile of
736 cognitive states. On a broader level, we think that wider adoption of such approaches will
737 have a meaningful impact on the inferences drawn from data, as well as the development of
738 theories regarding the structure of cognition.

739

Author contributions

740 Conceptualization: S. Panis and R. Ramsey; Software: S. Panis and R. Ramsey;
741 Writing - Original Draft Preparation: S. Panis; Writing - Review & Editing: S. Panis and
742 R. Ramsey; Supervision: R. Ramsey.

743

Conflicts of Interest

744 The author(s) declare that there were no conflicts of interest with respect to the
745 authorship or the publication of this article.

746

Prior versions

747 All of the submitted manuscript and Supplemental Material was previously posted to
748 a preprint archive: <https://doi.org/10.31234/osf.io/57bh6>

749

Supplemental Material

750

Disclosures**751 Data, materials, and online resources**

752 Link to public archive:
753 https://github.com/sven-panis/Tutorial_Event_History_Analysis
754 Supplemental Material: Panis_Ramsey_suppl_material.pdf

755 Ethical approval

756 Ethical approval was not required for this tutorial in which we reanalyze existing
757 data sets.

758

References

- 759 Allison, P. D. (1982). Discrete-Time Methods for the Analysis of Event Histories.
760 *Sociological Methodology*, 13, 61. <https://doi.org/10.2307/270718>
- 761 Allison, P. D. (2010). *Survival analysis using SAS: A practical guide* (2. ed). Cary, NC:
762 SAS Press.
- 763 Baker, D. H., Vilidaite, G., Lygo, F. A., Smith, A. K., Flack, T. R., Gouws, A. D., &
764 Andrews, T. J. (2021). Power contours: Optimising sample size and precision in
765 experimental psychology and human neuroscience. *Psychological Methods*, 26(3),
766 295–314. <https://doi.org/10.1037/met0000337>
- 767 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for
768 confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*,
769 68(3), 10.1016/j.jml.2012.11.001. <https://doi.org/10.1016/j.jml.2012.11.001>
- 770 Blossfeld, H.-P., & Rohwer, G. (2002). *Techniques of event history modeling: New*
771 *approaches to causal analysis*, 2nd ed (pp. x, 310). Mahwah, NJ, US: Lawrence
772 Erlbaum Associates Publishers.
- 773 Box-Steffensmeier, J. M. (2004). Event history modeling: A guide for social scientists.
774 Cambridge: University Press.
- 775 DeBruine, L. M., & Barr, D. J. (2021). Understanding Mixed-Effects Models Through
776 Data Simulation. *Advances in Methods and Practices in Psychological Science*, 4(1),
777 2515245920965119. <https://doi.org/10.1177/2515245920965119>
- 778 Gelman, A., Hill, J., & Vehtari, A. (2020). Regression and Other Stories.
779 <https://www.cambridge.org/highereducation/books/regression-and-other-stories/DD20DD6C9057118581076E54E40C372C>; Cambridge University Press.
780 <https://doi.org/10.1017/9781139161879>
- 782 Gelman, A., Vehtari, A., Simpson, D., Margossian, C. C., Carpenter, B., Yao, Y., ...
783 Modrák, M. (2020). *Bayesian Workflow*. arXiv.
784 <https://doi.org/10.48550/arXiv.2011.01808>

- 785 Heiss, A. (2021, November 10). A Guide to Correctly Calculating Posterior Predictions
786 and Average Marginal Effects with Multilevel Bayesian Models.
787 <https://doi.org/10.59350/wbn93-edb02>
- 788 Hosmer, D. W., Lemeshow, S., & May, S. (2011). *Applied Survival Analysis: Regression*
789 *Modeling of Time to Event Data* (2nd ed). Hoboken: John Wiley & Sons.
- 790 Kantowitz, B. H., & Pachella, R. G. (2021). The Interpretation of Reaction Time in
791 Information-Processing Research 1. *Human Information Processing*, 41–82.
792 <https://doi.org/10.4324/9781003176688-2>
- 793 Kruschke, J. K., & Liddell, T. M. (2018). The Bayesian New Statistics: Hypothesis testing,
794 estimation, meta-analysis, and power analysis from a Bayesian perspective.
795 *Psychonomic Bulletin & Review*, 25(1), 178–206.
796 <https://doi.org/10.3758/s13423-016-1221-4>
- 797 Kurz, A. S. (2023a). *Applied longitudinal data analysis in brms and the tidyverse* (version
798 0.0.3). Retrieved from <https://bookdown.org/content/4253/>
- 799 Kurz, A. S. (2023b). *Statistical rethinking with brms, ggplot2, and the tidyverse: Second*
800 *edition* (version 0.4.0). Retrieved from <https://bookdown.org/content/4857/>
- 801 Lakens, D. (2022). Sample Size Justification. *Collabra: Psychology*, 8(1), 33267.
802 <https://doi.org/10.1525/collabra.33267>
- 803 Landes, J., Engelhardt, S. C., & Pelletier, F. (2020). An introduction to event history
804 analyses for ecologists. *Ecosphere*, 11(10), e03238. <https://doi.org/10.1002/ecs2.3238>
- 805 McElreath, R. (2020). *Statistical Rethinking: A Bayesian Course with Examples in R and*
806 *STAN* (2nd ed.). New York: Chapman and Hall/CRC.
807 <https://doi.org/10.1201/9780429029608>
- 808 Nelder, J. A. (1999). From Statistics to Statistical Science. *Journal of the Royal Statistical*
809 *Society. Series D (The Statistician)*, 48(2), 257–269. Retrieved from
810 <https://www.jstor.org/stable/2681191>
- 811 Panis, S. (2020). How can we learn what attention is? Response gating via multiple direct

- 812 routes kept in check by inhibitory control processes. *Open Psychology*, 2(1), 238–279.
- 813 <https://doi.org/10.1515/psych-2020-0107>
- 814 Panis, S., Moran, R., Wolkersdorfer, M. P., & Schmidt, T. (2020). Studying the dynamics
815 of visual search behavior using RT hazard and micro-level speed–accuracy tradeoff
816 functions: A role for recurrent object recognition and cognitive control processes.
817 *Attention, Perception, & Psychophysics*, 82(2), 689–714.
818 <https://doi.org/10.3758/s13414-019-01897-z>
- 819 Panis, S., Schmidt, F., Wolkersdorfer, M. P., & Schmidt, T. (2020). Analyzing Response
820 Times and Other Types of Time-to-Event Data Using Event History Analysis: A Tool
821 for Mental Chronometry and Cognitive Psychophysiology. *I-Perception*, 11(6),
822 2041669520978673. <https://doi.org/10.1177/2041669520978673>
- 823 Panis, S., & Schmidt, T. (2016). What Is Shaping RT and Accuracy Distributions? Active
824 and Selective Response Inhibition Causes the Negative Compatibility Effect. *Journal of*
825 *Cognitive Neuroscience*, 28(11), 1651–1671. https://doi.org/10.1162/jocn_a_00998
- 826 Panis, S., & Schmidt, T. (2022). When does “inhibition of return” occur in spatial cueing
827 tasks? Temporally disentangling multiple cue-triggered effects using response history
828 and conditional accuracy analyses. *Open Psychology*, 4(1), 84–114.
829 <https://doi.org/10.1515/psych-2022-0005>
- 830 Panis, S., Torfs, K., Gillebert, C. R., Wagemans, J., & Humphreys, G. W. (2017).
831 Neuropsychological evidence for the temporal dynamics of category-specific naming.
832 *Visual Cognition*, 25(1-3), 79–99. <https://doi.org/10.1080/13506285.2017.1330790>
- 833 Panis, S., & Wagemans, J. (2009). Time-course contingencies in perceptual organization
834 and identification of fragmented object outlines. *Journal of Experimental Psychology:*
835 *Human Perception and Performance*, 35(3), 661–687.
836 <https://doi.org/10.1037/a0013547>
- 837 Pargent, F., Koch, T. K., Kleine, A.-K., Lermer, E., & Gaube, S. (2024). A Tutorial on
838 Tailored Simulation-Based Sample-Size Planning for Experimental Designs With

- 839 Generalized Linear Mixed Models. *Advances in Methods and Practices in Psychological*
840 *Science*, 7(4), 25152459241287132. <https://doi.org/10.1177/25152459241287132>
- 841 Singer, J. D., & Willett, J. B. (2003). *Applied Longitudinal Data Analysis: Modeling*
842 *Change and Event Occurrence*. Oxford, New York: Oxford University Press.
- 843 Smith, P. L., & Little, D. R. (2018). Small is beautiful: In defense of the small-N design.
844 *Psychonomic Bulletin & Review*, 25(6), 2083–2101.
845 <https://doi.org/10.3758/s13423-018-1451-8>
- 846 Teachman, J. D. (1983). Analyzing social processes: Life tables and proportional hazards
847 models. *Social Science Research*, 12(3), 263–301.
848 [https://doi.org/10.1016/0049-089X\(83\)90015-7](https://doi.org/10.1016/0049-089X(83)90015-7)
- 849 Tong, C. (2019). Statistical Inference Enables Bad Science; Statistical Thinking Enables
850 Good Science. *The American Statistician*, 73(sup1), 246–261.
851 <https://doi.org/10.1080/00031305.2018.1518264>
- 852 Wickelgren, W. A. (1977). Speed-accuracy tradeoff and information processing dynamics.
853 *Acta Psychologica*, 41(1), 67–85. [https://doi.org/10.1016/0001-6918\(77\)90012-9](https://doi.org/10.1016/0001-6918(77)90012-9)
- 854 Winter, B. (2019). *Statistics for Linguists: An Introduction Using R*. New York:
855 Routledge. <https://doi.org/10.4324/9781315165547>
- 856 Wolkersdorfer, M. P., Panis, S., & Schmidt, T. (2020). Temporal dynamics of sequential
857 motor activation in a dual-prime paradigm: Insights from conditional accuracy and
858 hazard functions. *Attention, Perception, & Psychophysics*, 82(5), 2581–2602.
859 <https://doi.org/10.3758/s13414-020-02010-5>