

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 GLC,
10 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 human
14 cognition. However, the orthodox method for analyzing such data – comparing means
15 between conditions – is known to conceal valuable information about the timeline of
16 psychological effects, such as their onset time and duration. The ability to reveal
17 finer-grained, “temporal states” of cognitive processes can have important consequences for
18 theory development by qualitatively changing the key inferences that are drawn from
19 psychological data. Luckily, well-established analytical approaches, such as event history
20 analysis (EHA), are able to evaluate the detailed shape of time-to-event distributions, and
21 thus characterize the time course of psychological states. One barrier to wider use of EHA,
22 however, is that the analytical workflow is typically more time-consuming and complex than
23 orthodox approaches. To help achieve broader uptake of EHA, in this paper we outline a set
24 of tutorials that detail one distributional method known as discrete-time EHA. We touch
25 upon several key aspects of the workflow, such as how to process raw data and specify
26 regression models, and we also consider the implications for experimental design, as well as
27 how to manage inter-individual differences. We finish the article by considering the benefits
28 of 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 means
30 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 distributional 36 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. However,
42 differences in mean RT conceal important pieces of information, such as when an
43 experimental effect starts, how it evolves with increasing waiting time, and whether its onset
44 is time-locked to other events (Panis, 2020; Panis, Moran, Wolkersdorfer, & Schmidt, 2020;
45 Panis & Schmidt, 2016, 2022; Panis, Torfs, Gillebert, Wagemans, & Humphreys, 2017; Panis
46 & Wagemans, 2009; Wolkersdorfer, Panis, & Schmidt, 2020). Such information is useful not
47 only for the interpretation of experimental effects under investigation, but also for cognitive
48 psychophysiology and computational model selection (Panis, Schmidt, Wolkersdorfer, &
49 Schmidt, 2020).

50

As a simple illustration, Figure 1 summarises simulated data that shows how
51 comparing means between two conditions can conceal the shapes of the underlying RT and
52 accuracy distributions. Indeed, compared to the aggregation of data across trials (Figure
53 1A), a distributional approach offers the possibility to reveal the time course of psychological
54 states (Figure 1B). For example, Figure 1B shows a first state (up to 400 ms after target
55 onset) for which the early upswing in hazard is equal for both conditions, and the emitted
56 responses are always correct in condition 1 and always incorrect in condition 2. In a second
57 state (400 to 500 ms), hazard is higher in condition 1, and conditional accuracies are close to
58 .5 in both conditions. In a third state (>500 ms), the effect disappears in hazard, and all
59 conditional accuracies are equal to 1.

60 Why does this matter for research in psychology? For many psychological questions,
61 the estimation of such “temporal states” information can be theoretically meaningful by
62 leading to more fine-grained understanding of psychological processes and by adding a
63 relatively under-used dimension – the passage of time – to the theory building toolkit. Thus,
64 a distributional approach permits different kinds of questions to be asked, different inferences
65 to be made, and it holds the potential to better discriminate between different theoretical
66 accounts of psychological and/or brain-based processes.

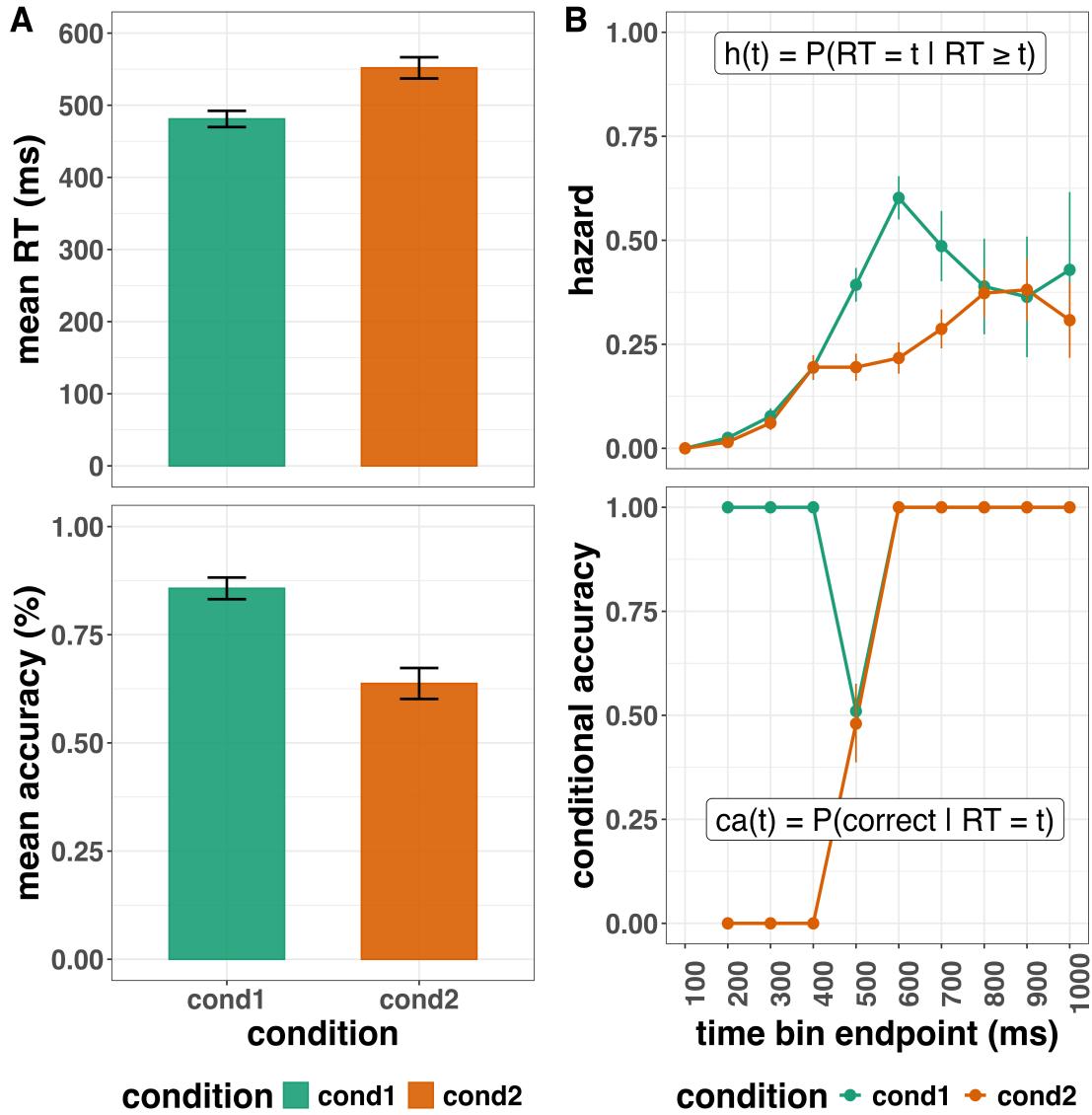


Figure 1. Simulated data showing mean performance versus distributional analyses. (A) The mean RT (top) and overall accuracy (bottom) for two conditions are plotted. (B) The discrete-time hazard functions (top) and conditional accuracy functions (bottom) are plotted for the same data. The first second 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).

67 1.2 Aims

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

73 Even though EHA is a widely used statistical tool and there already exist many
74 excellent reviews (Allison, 1982; Blossfeld & Rohwer, 2002; Box-Steffensmeier, 2004; Hosmer,
75 Lemeshow, & May, 2011; e.g., Singer & Willett, 2003; Teachman, 1983) and tutorials (e.g.,
76 Allison, 2010; Landes, Engelhardt, & Pelletier, 2020), we are not aware of any tutorials that
77 are aimed specifically at psychological RT (+ accuracy) data, and which provide worked
78 examples of the key data processing and Bayesian multilevel regression modelling steps. Set
79 within this context, our overall aim is to introduce a set of tutorials, which explain **how** to
80 do such analyses in the context of experimental psychology, rather than repeat in any detail
81 **why** you may do them. Therefore, we hope that our tutorials will provide a pathway for
82 research avenues in experimental psychology that have the potential to benefit from using
83 EHA in the future.

84 1.3 Structure

85 In what follows, the paper is organised in three main sections. In Section 2, we provide
86 a brief overview of EHA to orient the reader to the basic concepts that we will use
87 throughout the paper and why such an approach might be relevant for research in
88 experimental psychology. In Section 3, we outline a series of tutorials, which are written in
89 the R programming language and publicly available on our Github page
90 (https://github.com/sven-panis/Tutorial_Event_History_Analysis), along with all of the
91 other code and material associated with the project. The tutorials provide hands-on,
92 concrete examples of key parts of the analytical process, such as data wrangling, model

93 fitting and planning future studies, so that others can apply EHA to their own time-to-event
94 data measured in RT tasks. In Section 4, we discuss the strengths and weaknesses of the
95 approach for researchers in experimental psychology.

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

98 **2.1 A brief introduction to event history analysis**

99 EHA is a statistical approach to study the occurrence and timing of events, such as
100 disease onset, marriages, arrests, and job terminations (Allison, 2010). In this section, we
101 want to provide an intuition regarding how EHA works in general, as well as in the context
102 of experimental psychology. For those who want more detailed treatment of EHA and/or
103 regression equations, we refer the reader to several excellent textbooks on these topics
104 (Allison, 2010; Gelman, Hill, & Vehtari, 2020; Singer & Willett, 2003; Winter, 2019). We also
105 supply all relevant regression equations in section D of the Supplemental Material.

106 **2.1.1 Terminology and minimum requirements for EHA.** To avoid possible
107 confusion in terminology used, it is worth noting that EHA is frequently known by various
108 labels, such as survival analysis, hazard analysis, duration analysis, failure-time analysis, and
109 transition analysis (Singer & Willett, 2003). In this paper, we choose to use the term EHA
110 throughout.

111 In terms of minimum requirements to apply EHA, one must be able to:

- 112 1. define an event of interest that represents a qualitative change - a transition from one
113 discrete state to another - that can be situated in time (e.g., a button press, a saccade
114 onset, a fixation offset, etc.);
- 115 2. define time point zero (e.g., target stimulus onset, fixation onset, etc.);
- 116 3. measure the passage of time between time point zero and event occurrence in discrete

117 or continuous time units.

118 **2.1.2 Types of EHA.** There are different types of EHA. For example, the definition
119 of hazard and the type of models employed depend on whether one is using continuous or
120 discrete time units. As a lab, and mainly for practical reasons, we have much more
121 experience using discrete time EHA, and that is the approach that we describe and focus on
122 in this paper. This choice may seem counter-intuitive, given that RT is typically treated as a
123 continuous variable. However, continuous forms of EHA require much more data to estimate
124 the continuous-time hazard (rate) function well. Thus, by trading a bit of temporal
125 resolution for a lower number of trials, discrete-time methods seem ideal for dealing with
126 typical psychological time-to-event data sets for which there are less than ~200 trials per
127 condition per experiment. Moreover, as indicated by Allison (2010), learning discrete-time
128 EHA methods first will help in learning continuous-time methods, so it seems like a good
129 starting point.

130 To apply discrete time EHA, one divides time in discrete, contiguous time bins indexed
131 by t (e.g., $t = 1:10$ time bins; Figure 1B). Then let t be a discrete random variable denoting
132 the rank of the time bin in which a particular person's response occurs in a particular
133 experimental condition. For example, a response on a single trial might occur at 546 ms and
134 it would be in time bin 6 (any RTs from 501 ms to 600 ms). One then calculates the
135 discrete-time hazard function of event occurrence. The discrete-time hazard function gives
136 you, for each time bin, the probability that the event occurs (sometime) in bin t , given that
137 the event does not occur in previous bins. In other words, it reflects the instantaneous
138 likelihood that the event occurs in the current bin, given that it has not yet occurred in the
139 past, i.e., in one of the prior bins.

140 In the context of experimental psychology, it is often (but not always), the case that
141 responses can be classified as correct or incorrect. In those cases, one can also calculate the
142 conditional accuracy function (Figure B). The conditional accuracy function gives you for

143 each time bin the probability that a response is correct given that it is emitted in time bin t
144 (Allison, 2010; Kantowitz & Pachella, 2021; Wickelgren, 1977). The $ca(t)$ function is also
145 known as the micro-level speed-accuracy tradeoff (SAT) function. We refer to this extended
146 (hazard + conditional accuracy) analysis for choice RT data as EHA/SAT.

147 **2.2 Benefits of event history analysis for research in experimental psychology**

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

154 We highlight three benefits that we think are relevant to the domain of experimental
155 psychology. First, as illustrated in Figure 1, compared to averaging data across trials,
156 integrating results between hazard and conditional accuracy functions for choice RT data can
157 be informative for understanding psychological processes, in terms of inferences about
158 cognition and theoretical development. As such, the approach permits different kinds of
159 questions to be asked, different inferences to be made, and it holds the potential to
160 discriminate between theoretical accounts of psychological and/or brain-based processes. For
161 example, what kind of theory or theories could account for the effects reported in Figure 1B?
162 Are there new auxiliary assumptions that theories need to adopt? And are there new
163 experiments that need to be performed to test the novel predictions that follow from these
164 analyses?

165 Second, compared to more conventional analytical approaches, EHA uses more of the
166 data and deals with missing data differently. It is conventional with RT data to either (a)
167 use a response deadline and discard all trials without a response, or (b) wait in each trial

168 until a response occurs and then apply data trimming techniques, i.e., discarding too short
169 or too long RTs before calculating a mean RT (REF). Discarding data can introduce biases,
170 however. Rather than treat non-responses as missing data, EHA treats such trials as
171 *right-censored* observations on the variable RT, because all we know is that RT is greater
172 than some value. Right-censoring is a type of missing data problem and a nearly universal
173 feature of survival data including RT data. For example, if the censor time was 1 second,
174 then some trials result in observed event times (those with a RT below 1 second), while the
175 other trials result in response times that are right-censored at 1 second. The use of
176 right-censoring in EHA, therefore, presents a analytical strength of the approach compared
177 to many common approaches in experimental psychology.

178 Third, the approach is generalisable and applicable to many tasks that are commonly
179 used in experimental psychology, such as detection, discrimination and bistable perception
180 tasks (Supp Fig X). EHA can also be used with a range of common experimental
181 manipulations, such as stimulus-onset-asynchrony (Supp Fig X). The upshot is that one
182 general analytical approach, which holds several potential advantages, is widely applicable to
183 many substantive use-cases in the domain of experimental psychology.

184 **2.3 Implications for research design in experimental psychology**

185 Performing EHA in experimental psychology has implications for how experiments are
186 designed. More specifically, we consider two implications that researchers will need to
187 consider when using discrete time EHA.

188 First, since the number of trials per condition are spread across bins, it is important to
189 have a relatively large number of trial repetitions per participant and per condition.
190 Accordingly, experimental designs using this approach typically focus on factorial,
191 within-subject designs, in which a large number of observations are made on a relatively
192 small number of participants (so-called small-*N* designs). This approach emphasizes the

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

Second, the width of each time bin will need to be determined. For instance, in Figure 1B we chose 100ms in an arbitrary manner. In reality, however, bin width will need to be set by considering a number of factors simultaneously. For example, [[Sven - add a few sentences here]]

4. Tutorials

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

Tutorials 1a and 1b show how to calculate and plot the descriptive statistics of EHA/SAT when there are one or two independent variables, respectively. Tutorials 2a and 2b illustrate how to use Bayesian multilevel modeling to fit hazard and conditional accuracy models, respectively. Tutorials 3a and 3b show how to implement, respectively, multilevel models for hazard and conditional accuracy in the frequentist framework. Additionally, to further simplify the process for other users, the first two tutorials rely on a set of our own custom functions that make sub-processes easier to automate, such as data wrangling and plotting functions (see section B in the Supplemental Material for a list of the custom functions).

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

224 4.1 Tutorial 1a: Calculating descriptive statistics using a life table

225 **4.1.1 Data wrangling aims.** Our data wrangling procedures serve two related

226 purposes. First, we want to summarise and visualise descriptive statistics that relate to our
227 main research questions about the time course of psychological processes, using a life table.
228 A life table includes for each time bin, the risk set (i.e., the number of trials that are
229 event-free at the start of the bin), the number of observed events, and the estimates of $h(t)$,
230 $S(t)$, $P(t)$, possibly $ca(t)$, and their estimated standard errors (se).

231 Second, we want to produce two different data sets that can each be submitted to

232 different types of inferential modelling approaches. The two types of data structure we label
233 as ‘person-trial’ data and ‘person-trial-bin’ data. The ‘person-trial’ data (Table 1) will be
234 familiar to most researchers who record behavioural responses from participants, as it
235 represents the measured RT and accuracy per trial within an experiment. This data set is
236 used when fitting conditional accuracy models (Tutorials 2b and 3b).

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.

237 In contrast, the ‘person-trial-bin’ data (Table 2) has a different, more extended
 238 structure, which indicates in which bin a response occurred, if at all, in each trial. Therefore,
 239 the ‘person-trial-bin’ data generates a 0 in each bin until an event occurs and then it
 240 generates a 1 to signal an event has occurred in that bin. This data set is used when fitting
 241 hazard models (Tutorials 2a and 3a). It is worth pointing out that there is no requirement
 242 for an event to occur at all (in any bin), as maybe there was no response on that trial or the
 243 event occurred after the time window of interest. Likewise, when the event occurs in bin 1
 244 there would only be one row of data for that trial in the person-trial-bin data set.

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.

245 **4.1.2 A real data wrangling example.** To illustrate how to quickly set up life
 246 tables for calculating the descriptive statistics (functions of discrete time), we use a
 247 published data set on masked response priming from Panis and Schmidt (2016). In their first
 248 experiment, Panis and Schmidt (2016) presented a double arrow for 94 ms that pointed left
 249 or right as the target stimulus with an onset at time point zero in each trial. Participants
 250 had to indicate the direction in which the double arrow pointed using their corresponding
 251 index finger, within 800 ms after target onset. Response time and accuracy were recorded on
 252 each trial. Prime type (blank, congruent, incongruent) and mask type were manipulated.
 253 Here we focus on the subset of trials in which no mask was presented. The 13-ms prime
 254 stimulus was a double arrow presented 187 ms before target onset in the congruent (same

255 direction as target) and incongruent (opposite direction as target) prime conditions.

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

259 The required column names are as follows:

- 260 • “pid”, indicating unique participant IDs;
- 261 • “trial”, indicating each unique trial per participant;
- 262 • “condition”, a factor indicating the levels of the independent variable (1, 2, ...) and
 263 the corresponding labels;
- 264 • “rt”, indicating the response times in ms;
- 265 • “acc”, indicating the accuracies (1/0).

266 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")))
```

267 Next, we can set up the life tables and plots of the discrete-time functions $h(t)$, $S(t)$,
 268 $ca(t)$, and $P(t)$ – see section A of the Supplemental Material for their definitions. To do so
 269 using a functional programming approach, one has to nest the data within participants using
 270 the group_nest() function, and supply a user-defined censoring time and bin width to our
 271 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))

```

272 Note that the censoring time should be a multiple of the bin width (both in ms). The
 273 censoring time should be a time point after which no informative responses are expected
 274 anymore. In experiments that implement a response deadline in each trial the censoring time
 275 can equal that deadline time point. Trials with a RT larger than the censoring time, or trials
 276 in which no response is emitted during the data collection period, are treated as
 277 right-censored observations in EHA. In other words, these trials are not discarded, because
 278 they contain the information that the event did not occur before the censoring time.
 279 Removing such trials before calculating the mean event time will result in underestimation of
 280 the true mean.

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

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

286 the conditional accuracies using our custom function `calc_ca()`, add the $ca(t)$ estimates to

287 the life table using our custom function `join_lt_ca()`, and then plot the descriptive statistics

288 using our custom function `plot_eha()`. One can now inspect different aspects, including the

289 life table for a particular condition of a particular subject, and a plot of the different

290 functions for a particular participant. In general, it is important to visually inspect the

291 functions first for each participant, in order to identify individuals that may be guessing

292 (e.g., a flat conditional accuracy function at .5 indicates that someone is just guessing),

293 outlying individuals, and/or different groups with qualitatively different behavior.

294 Table 3 shows the life table for condition “blank” (no prime stimulus presented) for

295 participant 6.

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.

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

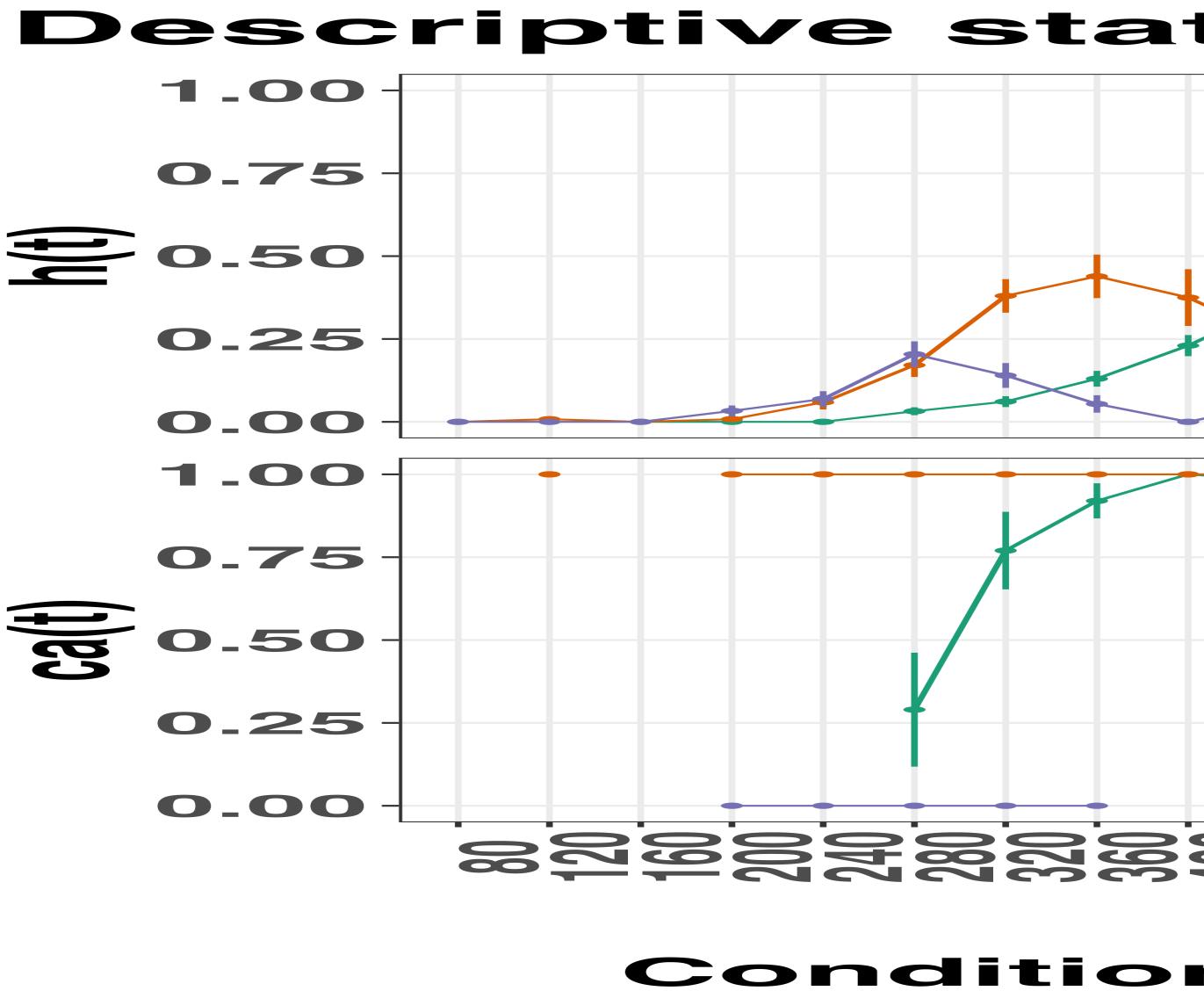


Figure 2. Estimated discrete-time hazard (h), survivor (S), conditional accuracy (ca) and probability mass (P) functions for participant 6. Vertical dotted lines indicate the estimated median RTs. Error bars represent +/- 1 standard error of the respective proportion.

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

302 0.17, and 0.20 for the blank, congruent, and incongruent prime conditions, respectively. In

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

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

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

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

307 0 for the blank, congruent, and incongruent prime conditions, respectively. In other words, if

308 a response is emitted in bin (240,280], then the probability that it is correct is estimated to

309 be 0.29, 1, and 0 for the blank, congruent, and incongruent prime conditions, respectively.

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

311 conditional probability of response occurrence in the next 40 ms is estimated to be 0.36, 0.25,

312 and 0.06 for the blank, congruent, and incongruent prime conditions, respectively. And when

313 a response does occur in bin (400,440], then the probability that it is correct is estimated to

314 be 0.98, 1, and 1 for the blank, congruent, and incongruent prime conditions, respectively.

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

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

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

318 responses are fully controlled by the target stimulus. Qualitatively similar results were

319 obtained for the other five participants. When participants show qualitatively similar

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

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

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

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

324 distributional data show that early responses are triggered exclusively by the prime stimulus,

325 while only later responses reflect primed responses to the target stimulus.

326 At this point, we have calculated, summarised and plotted descriptive statistics for the
327 key variables in EHA/SAT. As we will show in later Tutorials, statistical models for $h(t)$ and
328 $ca(t)$ can be implemented as generalized linear mixed regression models predicting event
329 occurrence (1/0) and conditional accuracy (1/0) in each bin of a selected time window for
330 analysis. But first we consider calculating the descriptive statistics for two independent
331 variables.

332 **4.2 Tutorial 1b: Generalising to a more complex design**

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

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

344 **4.3 Tutorial 2a: Fitting Bayesian hazard models to discrete time-to-event data**

345 In this third tutorial, we illustrate how to fit Bayesian multilevel regression models to
346 the RT data of the masked response priming data used in Tutorial 1a. Fitting (Bayesian or
347 non-Bayesian) regression models to time-to-event data is important when you want to study
348 how the shape of the hazard function depends on various predictors (Singer & Willett, 2003).

349 In general, when fitting regression models, our lab adopts an estimation approach to
350 multilevel regression (Kruschke & Liddell, 2018; Winter, 2019), which is heavily influenced

351 by the Bayesian framework as suggested by Richard McElreath (Kurz, 2023b; McElreath,
352 2020). We also use a “keep it maximal” approach to specifying varying (or random) effects
353 (Barr, Levy, Scheepers, & Tily, 2013). This means that wherever possible we include varying
354 intercepts and slopes per participant. To make inferences, we use two main approaches. We
355 compare models of different complexity, using information criteria (e.g., WAIC) and
356 cross-validation (e.g., LOO), to evaluate out-of-sample predictive accuracy (McElreath,
357 2020). We also take the most complex model and evaluate key parameters of interest using
358 point and interval estimates.

359 **4.3.1 Hazard model considerations.** There are several analytic decisions one has
360 to make when fitting a discrete-time hazard model. First, one has to select an analysis time
361 window, i.e., a contiguous set of bins for which there is enough data for each participant.
362 Second, given that the dependent variable (event occurrence) is binary, one has to select a
363 link function (see section C in the Supplemental Material). The cloglog link is preferred over
364 the logit link when events can occur in principle at any time point within a bin, which is the
365 case for RT data (Singer & Willett, 2003). Third, one has to choose whether to treat TIME
366 (i.e., the time bin index t) as a categorical or continuous predictor. And when you treat a
367 variable as a categorical predictor, you can choose between reference coding and index
368 coding. With reference coding, one defines the variable as a factor and selects one of the k
369 categories as the reference level. Brm() will then construct $k-1$ indicator variables (see model
370 M1d in Tutorial_2a.Rmd for an example). With index coding, one constructs an index
371 variable that contains integers that correspond to different categories (see models M0i and
372 M1i below). As explained by McElreath (2020), the advantage of index coding is that the
373 same prior can be assigned to each level of the index variable, so that each category has the
374 same prior uncertainty.

375 In the case of a large- N design without repeated measurements, the parameters of a
376 discrete-time hazard model can be estimated using standard logistic regression software after
377 expanding the typical person-trial data set into a person-trial-bin data set (Allison, 2010).

378 When there is clustering in the data, as in the case of a small- N design with repeated
 379 measurements, the parameters of a discrete-time hazard model can be estimated using
 380 population-averaged methods (e.g., Generalized Estimating Equations), and Bayesian or
 381 frequentist generalized linear mixed models (Allison, 2010).

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

387 In tutorial_2a.Rmd we fit several Bayesian multilevel models (i.e., generalized linear
 388 mixed models) that differ in complexity to the person-trial-bin oriented data set that we
 389 created in Tutorial 1a. We decided to select the analysis time window (200,600] and the
 390 cloglog link. Below, we shortly discuss two of these models. The person-trial-bin data set is
 391 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)))
```

392 **4.3.2 Prior distributions.** To get the posterior distribution of each model

393 parameter given the data, we need to specify prior distributions for the model parameters
394 which reflect our prior beliefs. In Tutorial_2a.Rmd we perform a few prior predictive checks
395 to make sure our selected prior distributions reflect our prior beliefs (Gelman, Vehtari, et al.,
396 2020).

397 The middle column of Supplementary Figure 2 (section E of the Supplemental

398 Material) shows six examples of prior distributions for an intercept on the logit and/or
399 cloglog scales. While a normal distribution with relatively large variance is often used as a
400 weakly informative prior for continuous dependent variables, rows A and B of Supplementary
401 Figure 2 show that specifying such distributions on the logit and cloglog scales actually leads
402 to rather informative distributions on the original probability scale, as most mass is pushed
403 to probabilities of 0 and 1.

404 **4.3.3 Model M0i: A null model with index coding.** When you do not want to

405 make assumptions about the shape of the hazard function, or its shape is not smooth but
406 irregular, then you can use a general specification of TIME, i.e., fit one grand intercept per
407 time bin. In this first model, we use a general specification of TIME using index coding, and
408 do not include experimental predictors. We call this model “M0i”.

409 Before we fit model M0i, we select the necessary columns from the data, and specify

410 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")
```

411 After selecting the bernoulli family and the cloglog link, the model formula is specified.

412 The specification “0 + …” removes the default intercept in brm(). The fixed effects include
 413 an intercept for each level of timebin. Each of these intercepts is allowed to vary across
 414 individuals (variable pid). We request 2000 samples from the posterior distribution for each
 415 of four chains. Estimating model M0i took about 30 minutes on a MacBook Pro (Sonoma
 416 14.6.1 OS, 18GB Memory, M3 Pro Chip).

417 **4.3.4 Model M1i: Adding the effects of prime-target congruency.** Previous
 418 research has shown that psychological effects typically change over time (Panis, 2020; Panis,
 419 Moran, et al., 2020; Panis & Schmidt, 2022; Panis et al., 2017; Panis & Wagemans, 2009). In
 420 the next model, therefore, we use index coding for both TIME (variable “timebin”) and the
 421 categorical predictor prime-target-congruency (variable “prime”), so that we get 30 grand
 422 intercepts, one for each combination of timebin level and prime level. Here is the model
 423 formula of this model that we call “M1i”.

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

424 Estimating model M1i took about 124 minutes.

425 **4.3.5 Compare the models.** We can compare the two models using the Widely

426 Applicable Information Criterion (WAIC) and Leave-One-Out (LOO) cross-validation, and
 427 look at model weights for both criteria (Kurz, 2023a; McElreath, 2020).

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

428 ## model_M0i model_M1i

```
429 ##          0          1  
  
model_weights(model_M0i, model_M1i, weights = "waic") %>% round(digits = 2)  
  
430 ## model_M0i model_M1i  
431 ##          0          1
```

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

434 **4.3.6 Evaluating parameter estimates in model M1i.** To make inferences from
435 the parameter estimates in model M1i, we first plot the densities of the draws from the
436 posterior distributions of its population-level parameters in Figure 5, together with point
437 (median) and interval estimates (80% and 95% credible intervals).

Posterior distributions for population-level effects in Model M1i

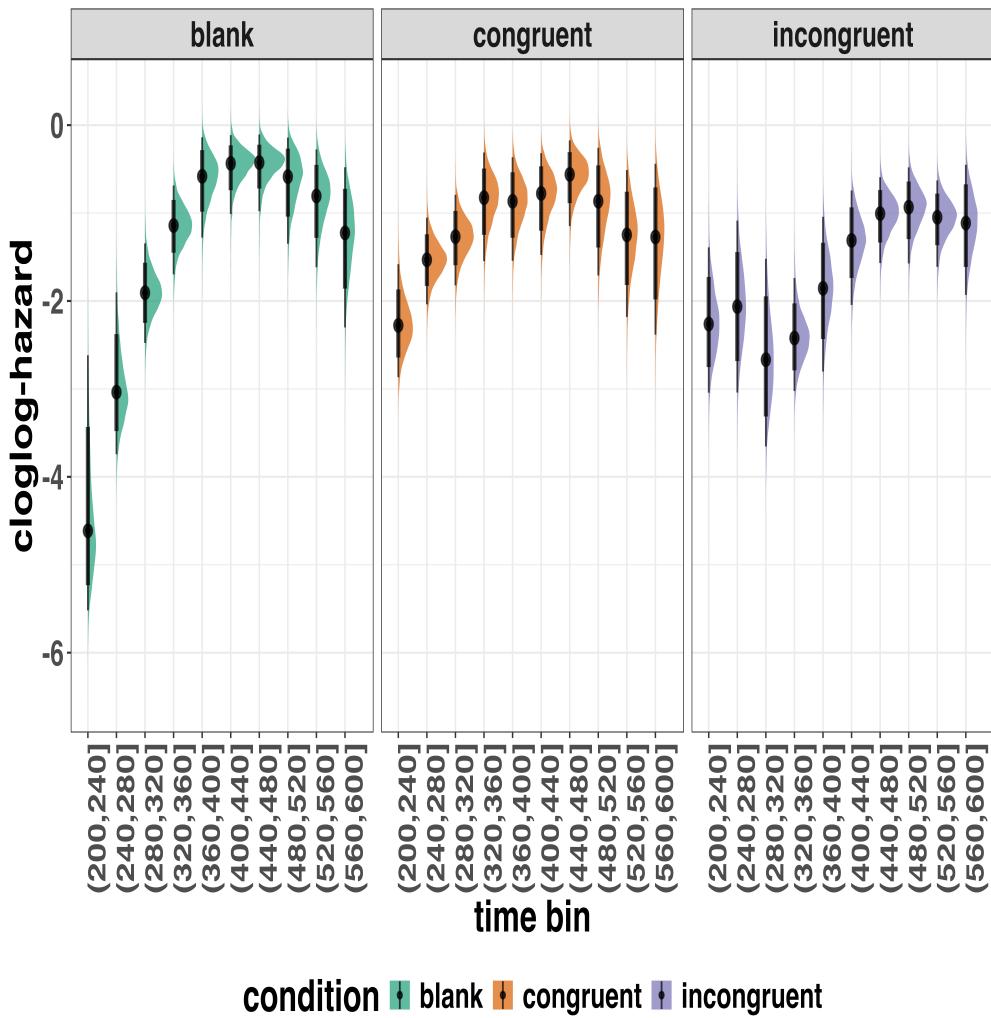


Figure 3. Medians and 80/95% credible intervals of the posterior distributions of the population-level parameters of model M1i.

Because the parameter estimates are on the cloglog-hazard scale, we can ease our

interpretation by plotting the expected value of the posterior predictive distribution – the predicted hazard values – at the population level (Figure 6A), and for each participant in the data set (Figure 6B).

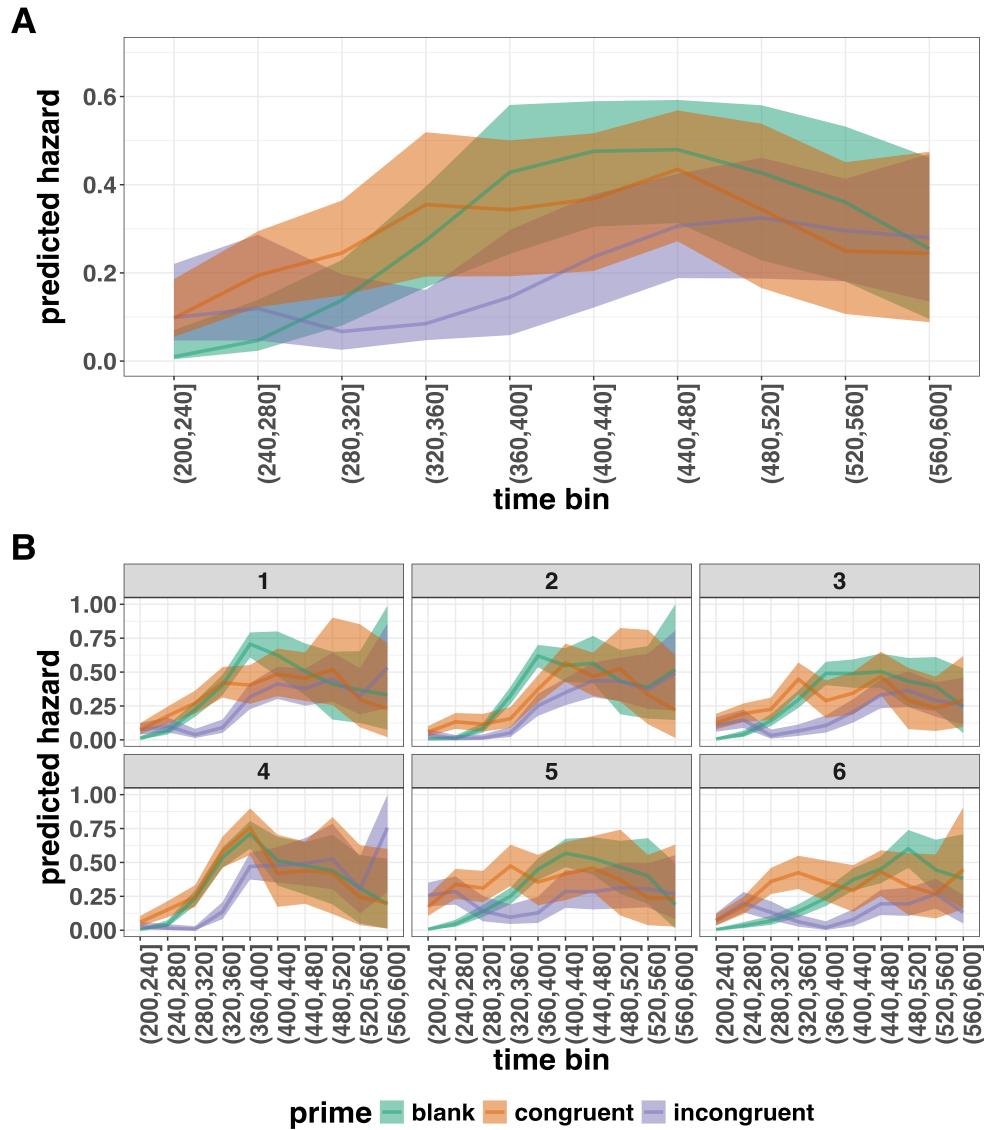


Figure 4. Point (median) and 80/95% credible interval summaries of the hazard estimates (expected values of the draws from the posterior predictive distributions) in each time bin at the population level (A), and for each participant (B).

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

⁴⁴⁶ (Figure 7B; subject-specific average marginal effect).

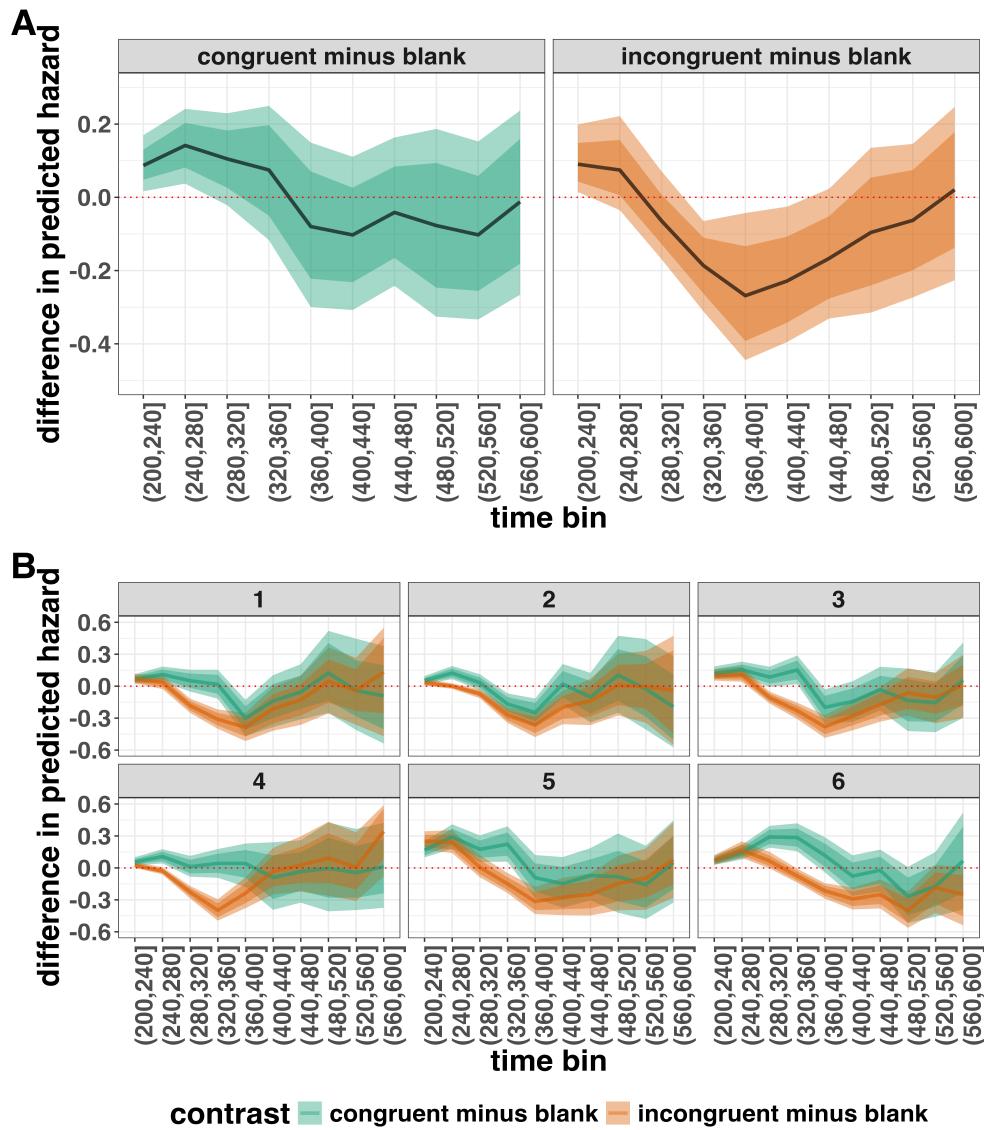


Figure 5. Point (mean) and 80/95% credible interval summaries of estimated differences in hazard in each time bin at the population level (A), and for each participant (B).

⁴⁴⁷ The point estimates and quantile intervals can be reported in a table (see
⁴⁴⁸ Tutorial_2a.Rmd for details).

⁴⁴⁹ **Example conclusions for M1i.** What can we conclude from model M1i about our
⁴⁵⁰ research question, i.e., the temporal dynamics of the effect of prime-target congruency on

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

454 If we want to estimate the population-level effect of prime type on hazard, we can base
455 our conclusion on Figure 7A. The contrast “congruent minus blank” was estimated to be
456 0.09 hazard units in bin (200,240] (95% CrI = [0.02, 0.17]), and 0.14 hazard units in bin
457 (240,280]) (95% CrI = [0.04, 0.25]). For the other bins, the 95% credible interval contained
458 zero. The contrast “incongruent minus blank” was estimated to be 0.09 hazard units in bin
459 (200,240] (95% CrI = [0.01, 0.21]), -0.19 hazard units in bin (320,360] (95% CrI = [-0.31,
460 -0.06]), -0.27 hazard units in bin (360,400] (95% CrI = [-0.45, -0.04]), and -0.23 hazard units
461 in bin (400,440] (95% CrI = [-0.40, -0.03]). For the other bins, the 95% credible interval
462 contained zero.

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

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

474 Future studies could (a) increase the number of participants to estimate the proportion
475 of “dippers” in the subject population, and/or (b) try to explain why this dip occurs. For

example, Panis and Schmidt (2016) concluded that active, top-down, task-guided response inhibition effects emerge around 360 ms after the onset of the stimulus following the prime (here: the target stimulus). Such a top-down inhibitory effect might exist in our priming data set, because after some time participants will learn that the first stimulus is not the one they have to respond to. To prevent a premature overt response to the prime they thus might gradually increase a global response threshold during the remainder of the experiment, which could result in a lower hazard in congruent trials compared to blank trials, for bins after \sim 360 ms, and towards the end of the experiment. This effect might be masked for incongruent primes by the response competition effect.

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

4.4 Tutorial 2b: Fitting Bayesian conditional accuracy models

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

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

501 intervals).

Posterior distributions for population-level effects in Model M1i_ca

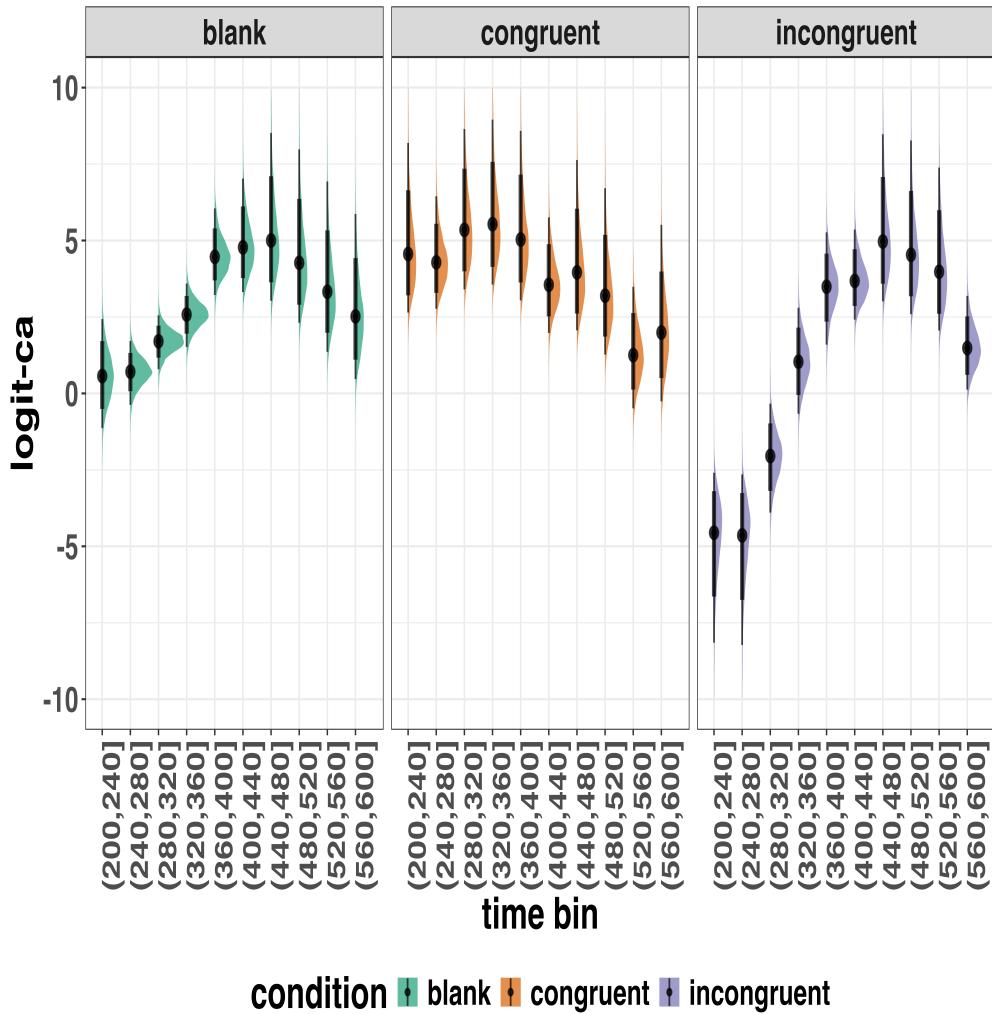


Figure 6. Medians and 80/95% credible intervals of the posterior distributions of the population-level parameters of model M1i_ca. ca = conditional accuracy.

502 Because the parameter estimates are on the logit-ca scale, we can ease our
503 interpretation by plotting the expected value of the posterior predictive distribution – the
504 predicted conditional accuracies – at the population level (Figure 9A), and for each
505 participant in the data set (Figure 9B).

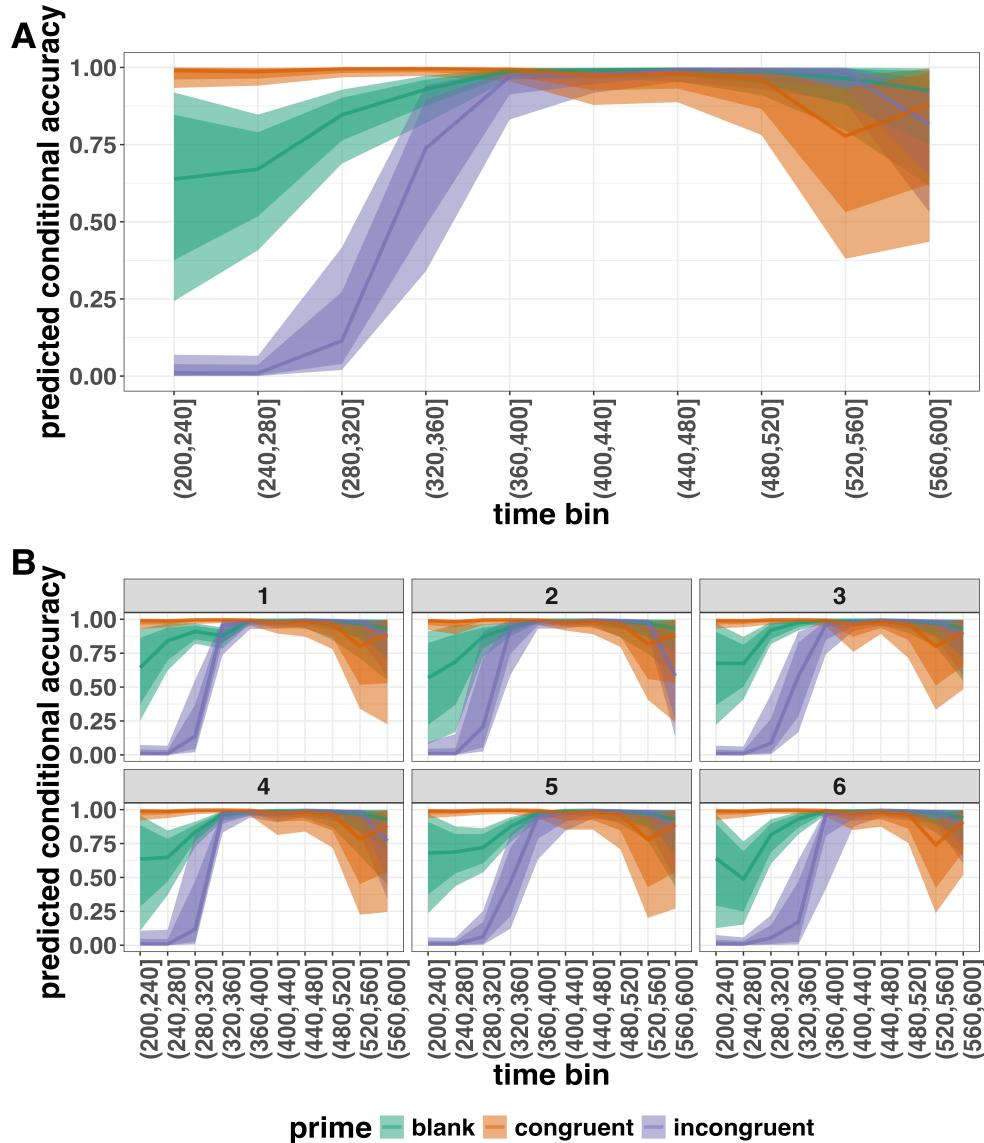


Figure 7. Point (median) and 80/95% credible interval summaries of the conditional accuracy estimates (expected values of the draws from the posterior predictive distributions) in each time bin at the population level (A), and for each participant (B).

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 at the

population level (Figure 10A) and for each participant (Figure 10B).

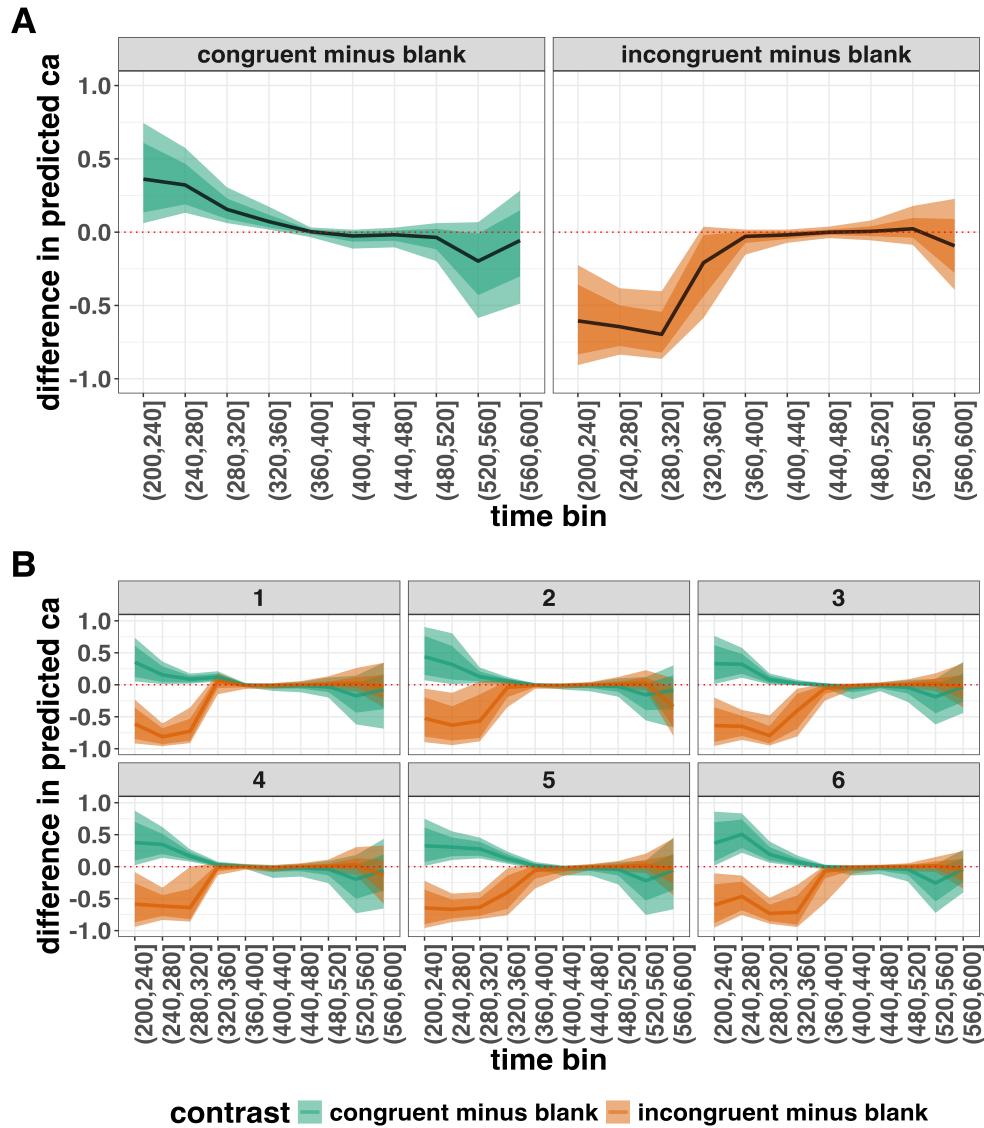


Figure 8. Point (mean) and 80/95% credible interval summaries of estimated differences in conditional accuracy in each time bin at the population level (A), and for each participant (B).

510 Based on Figure 10A we see that on the population level congruent primes have a
 511 positive effect on the conditional accuracy of emitted responses in time bins (200,240],
 512 (240,280], (280,320], and (320,360], relative to the estimates in the baseline condition (blank
 513 prime; red dashed lines in Figure 10A). Incongruent primes have a negative effect on the

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

516 **4.7 Tutorial 4: Planning**

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

518 **4.7.1 Background.** The general approach to planning that we adopt here involves

519 simulating reasonably structured data to help guide what you might be able to expect from
520 your data once you collect it (Gelman, Vehtari, et al., 2020). The basic structure and code
521 follows the examples outlined by Solomon Kurz in his ‘power’ blog posts
522 (<https://solomonkurz.netlify.app/blog/bayesian-power-analysis-part-i/>) and Lisa DeBruine’s
523 R package `faux{}` (<https://debruine.github.io/faux/>) as well as these related papers
524 (DeBruine & Barr, 2021; Pargent, Koch, Kleine, Lermer, & Gaube, 2024).

525 **4.7.2 Basic workflow.** The basic workflow is as follows:

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

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

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

539 **4.7.3 Fit a regression model and simulate one dataset.** We again use the data

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

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

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

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

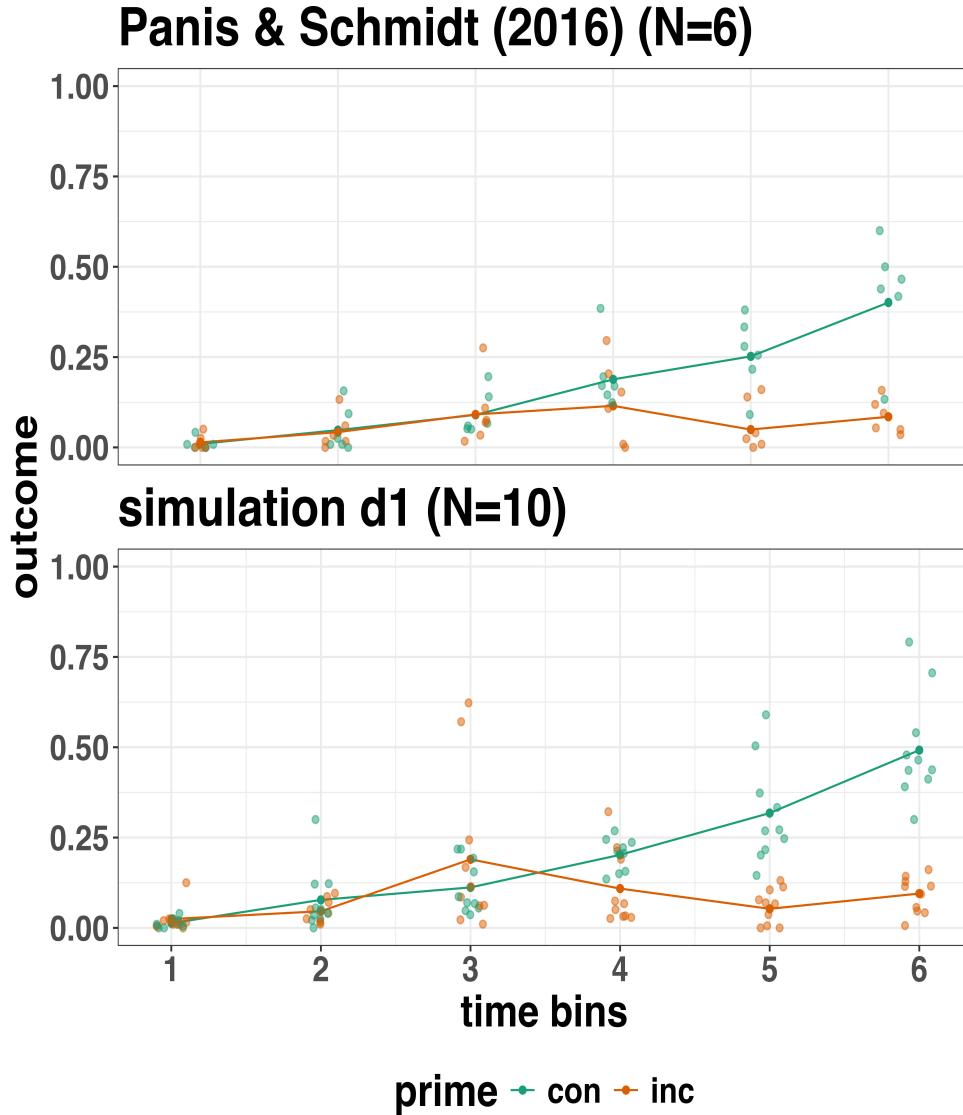


Figure 9. Raw data from Panis and Schmidt (2016) and simulated data from 10 participants.

549 4.7.4 Simulate and summarise data across a range of parameter values.

550 Here we use the same data simulation process as used above, but instead of simulating one
 551 dataset, we simulate 1000 datasets per variation in parameter values. Specifically, in
 552 Simulation 1, we vary the number of trials per condition (100, 200, and 400), as well as the
 553 effect size in bin 6. We focus on bin 6 only, in terms of varying the effect size, just to make
 554 things simpler and easier to understand. The effect size observed in bin 6 in this subsample
 555 of data was a 79% reduction in hazard value from the congruent prime (0.401 hazard value)

556 to the incongruent prime condition (0.085 hazard value). In other words, a hazard ratio of
557 0.21 (e.g., $0.085/0.401 = 0.21$). As a starting point, we chose three effect sizes, which covered
558 a fairly broad range of hazard ratios (0.25, 0.5, 0.75), which correspond to a 75%, 50% and
559 25% reduction in hazard value as a function of prime condition.

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

567 The results of Simulation 1 show that if we were targeting an effect size similar to the
568 one reported in the original study, then testing 10 participants and collecting 100 trials per
569 condition would be enough to provide over 95% power. However, we could not be as
570 confident about smaller effects, such as a hazard ratio of 50% or 25%. From this simulation,
571 we can see that somewhere between an effect size of a 50% and 75% reduction in hazard
572 value, power increases to a range that most researchers would consider acceptable (i.e.,
573 >95% power). To probe this space a little further, we decided to run a second simulation,
574 which varied different parameters.

575 In Simulation 2, we varied the effect size between a different range of values (0.5, 0.4,
576 0.3), which correspond to a 50%, 60% and 70% reduction in hazard value as a function of
577 prime condition. In addition, we varied the number of participants per experiment between
578 10, 15, and 20 participants. Given that trial count per condition made little difference to
579 power in Simulation 1, we fixed trial count at 200 trials per condition in Simulation 2.
580 Summary results from Simulation 2 are shown in Figure 13B. A summary of these power
581 calculations might be as follows (trial count = 200 per condition in all cases):

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

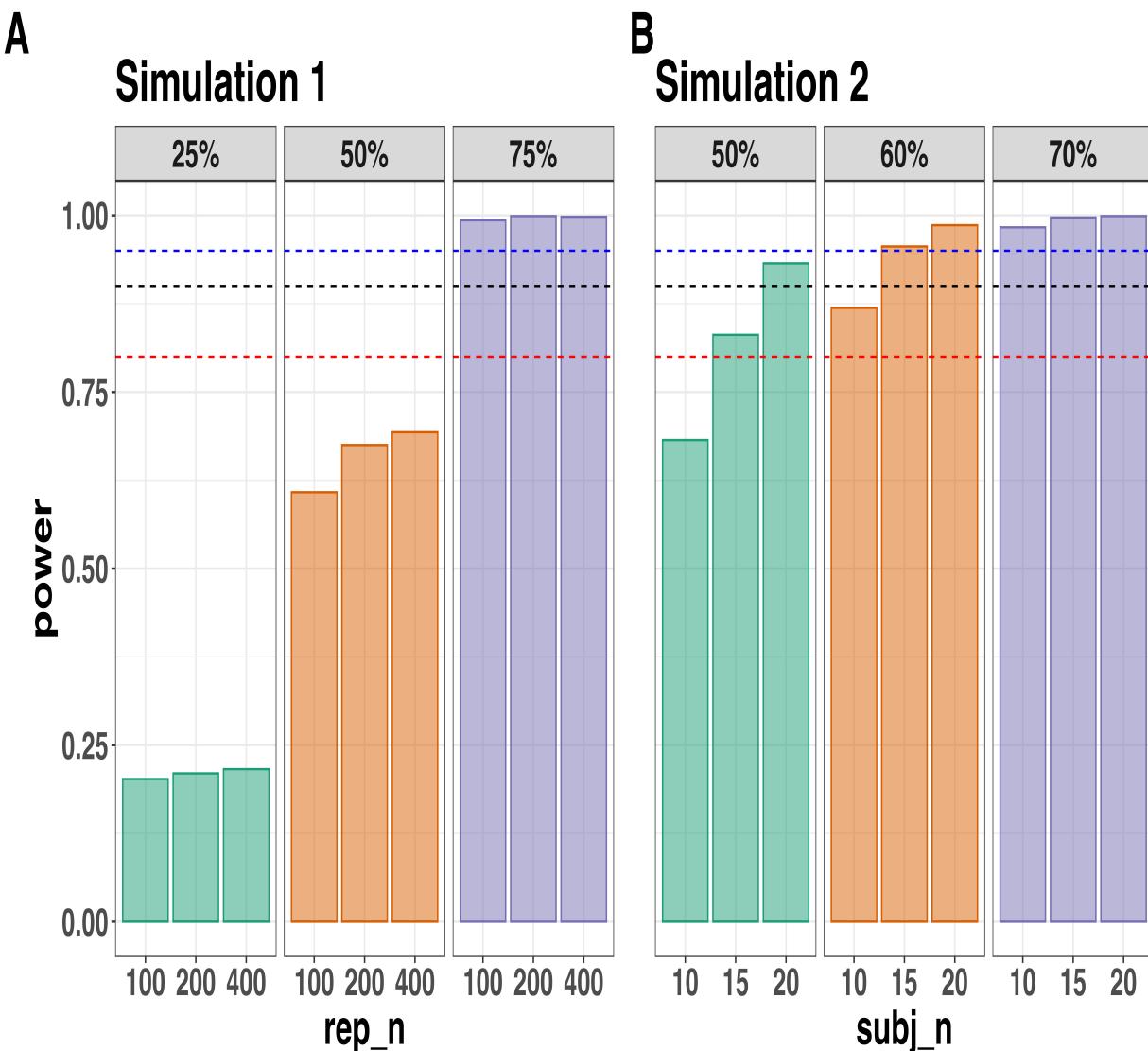


Figure 10. Statistical power across data Simulation 1 (A) and Simulation 2 (B). Power was calculated as the percentage of lower-bound 95% confidence intervals that exclude zero when the difference between prime condition is calculated (congruent - incongruent). In Simulation 1, the effect size was varied 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.

585 **4.7.5 Planning decisions.** Now that we have summarised our simulated data, what

586 planning decisions could we make about a future study? More concretely, how many trials
587 per condition should we collect and how many participants should we test? Like almost
588 always when planning future studies, the answer depends on your objectives, as well as the
589 available resources (Lakens, 2022). There is no straightforward and clear-cut answer. Some
590 considerations might be as follows:

- 591 • How much power or precision are you looking to obtain in this particular study?
- 592 • Are you running multiple studies that have some form of replication built in?
- 593 • What level of resources do you have at your disposal, such as time, money and
594 personnel?
- 595 • How easy or difficult is it to obtain the specific type of sample?

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

597 hazard ratio of 0.4 or 0.5 as a target effect size since this is much smaller than that observed
598 previously (Panis & Schmidt, 2016). Then we might pick the corresponding combination of
599 trial count per condition (e.g., 200) and participant sample size (e.g., N=10 or N=15) that
600 takes you over the 80% power mark. If we wanted to maximise power based on these
601 simulations, and we had the time and resources available, then we would test N=20
602 participants, which would provide >90% power for an effect size of 0.5.

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

604 what planning choices we made based on these data simulations, we would not solely rely on
605 data collected from one single study. Instead, we would run a follow-up experiment that
606 replicates and extends the initial result. By doing so, we would aim to avoid the Cult of the
607 Isolated Single Study (Nelder, 1999; Tong, 2019), and thus reduce the reliance on any one
608 type of planning tool, such as a power analysis. Then, we would look for common patterns
609 across two or more experiments, rather than trying to make the case that a single study on
610 its own has sufficient evidential value to hit some criterion mark.

611 4. Discussion

612 This main motivation for writing this paper is the observation that EHA and SAT
613 analysis remain under-used in psychological research. As a consequence, the field of
614 psychological research is not taking full advantage of the many benefits EHA/SAT provides
615 compared to more conventional analyses. By providing a freely available set of tutorials,
616 which provide step-by-step guidelines and ready-to-use R code, we hope that researchers will
617 feel more comfortable using EHA/SAT in the future. Indeed, we hope that our tutorials may
618 help to overcome a barrier to entry with EHA/SAT, which is that such approaches require
619 more analytical complexity compared to mean-average comparisons. While we have focused
620 here on within-subject, factorial, small- N designs, it is important to realize that EHA/SAT
621 can be applied to other designs as well (large- N designs with only one measurement per
622 subject, between-subject designs, etc.). As such, the general workflow and associated code
623 can be modified and applied more broadly to other contexts and research questions. In the
624 following, we discuss issues relating to model complexity and interpretability, individual
625 differences, as well as limitations of the approach and future extensions.

626 **627 4.1 What are the main use-cases of EHA for understanding cognition and brain
function?**

628 For those researchers, like ourselves, who are primarily interested in understanding
629 human cognitive and brain systems, we consider two broadly-defined, main use-cases of EHA.
630 First, as we hope to have made clear by this point, EHA is one way to investigating a
631 “temporal states” approach to cognitive processes. EHA provides one way to uncover when
632 cognitive states may start and stop, as well as what they may be tied to or interact with.
633 Therefore, if your research questions concern **when** and **for how long** psychological states
634 occur, our EHA tutorials could be useful tools for you to use.

635 Second, even if you are not primarily interested in studying the temporal states of
636 cognition, EHA could still be a useful tool to consider using, in order to qualify inferences

that are being made based on mean-average comparisons. Given that distinctly different inferences can be made from the same data based on whether one computes a mean-average across trials or a RT distribution of events (Figure 1), it may be important for researchers to supplement mean-average comparisons with EHA. One could envisage scenarios where the implicit assumption of an effect manifesting across all of the time bins measured would not be supported by EHA. Therefore, the conclusion of interest would not apply to all responses, but instead it would be restricted to certain aspects of time.

4.2 Model complexity versus interpretability

EHA can quickly become very complex when adding more than one time scale, due to the many possible higher-order interactions. For example, some of the models discussed in Tutorial 2a, which we did not focus on in the main text, contain two time scales as covariates: the passage of time on the within-trial time scale, and the passage of time on the across-trial (or within-experiment) time scale. However, when trials are presented in blocks, and blocks of trials within sessions, and when the experiment comprises three sessions, then four time scales can be defined (within-trial, within-block, within-session, and within-experiment). From a theoretical perspective, adding more than one time scale – and their interactions – can be important to capture plasticity and other learning effects that may play out on such longer time scales, and that are probably present in each experiment in general. From a practical perspective, therefore, some choices need to be made to balance the amount of data that is being collected per participant, condition and across the varying timescales. As one example, if there are several timescales of relevance, then it might be prudent for interpretational purposes to limit the number of experimental predictor variables (conditions). This is of course where planning and data simulation efforts would be important to provide a guide to experimental design choices (see Tutorial 4).

661 4.3 Limitations

662 Compared to the orthodox method – comparing mean-averages between conditions –
663 the most important limitation of multilevel hazard and conditional accuracy modeling is that
664 it might take a long time to estimate the parameters using Bayesian methods or the model
665 might have to be simplified significantly to use frequentist methods.

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

679 5. Conclusions

680 Estimating the temporal distributions of RT and accuracy provide a rich source of
681 information on the time course of cognitive processing, which have been largely undervalued
682 in the history of experimental psychology and cognitive neuroscience. We hope that by
683 providing a set of hands-on, step-by-step tutorials, which come with custom-built and freely
684 available code, researchers will feel more comfortable embracing EHA and investigating the
685 temporal profile of cognitive states. On a broader level, we think that wider adoption of such
686 approaches will have a meaningful impact on the inferences drawn from data, as well as the

687 development of theories regarding the structure of cognition.

688 Author contributions

689 Conceptualization: S. Panis and R. Ramsey; Software: S. Panis and R. Ramsey;

690 Writing - Original Draft Preparation: S. Panis; Writing - Review & Editing: S. Panis and R.

691 Ramsey; Supervision: R. Ramsey.

692 Conflicts of Interest

693 The author(s) declare that there were no conflicts of interest with respect to the

694 authorship or the publication of this article.

695 Prior versions

696 All of the submitted manuscript and Supplemental Material was previously posted to a

697 preprint archive: <https://doi.org/10.31234/osf.io/57bh6>

698 Supplemental Material**699 Disclosures****700 Data, materials, and online resources**

701 Link to public archive:

702 https://github.com/sven-panis/Tutorial_Event_History_Analysis

703 Supplemental Material: Panis_Ramsey_suppl_material.pdf

704 Ethical approval

705 Ethical approval was not required for this tutorial in which we reanalyze existing data

706 sets.

References

- 707
- 708 Allison, P. D. (1982). Discrete-Time Methods for the Analysis of Event Histories.
- 709 *Sociological Methodology*, 13, 61. <https://doi.org/10.2307/270718>
- 710 Allison, P. D. (2010). *Survival analysis using SAS: A practical guide* (2. ed.). Cary, NC: SAS
- 711 Press.
- 712 Baker, D. H., Vilidaite, G., Lygo, F. A., Smith, A. K., Flack, T. R., Gouws, A. D., &
- 713 Andrews, T. J. (2021). Power contours: Optimising sample size and precision in
- 714 experimental psychology and human neuroscience. *Psychological Methods*, 26(3),
- 715 295–314. <https://doi.org/10.1037/met0000337>
- 716 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for
- 717 confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*,
- 718 68(3), 10.1016/j.jml.2012.11.001. <https://doi.org/10.1016/j.jml.2012.11.001>
- 719 Blossfeld, H.-P., & Rohwer, G. (2002). *Techniques of event history modeling: New*
- 720 *approaches to causal analysis*, 2nd ed (pp. x, 310). Mahwah, NJ, US: Lawrence Erlbaum
- 721 Associates Publishers.
- 722 Box-Steffensmeier, J. M. (2004). Event history modeling: A guide for social scientists.
- 723 Cambridge: University Press.
- 724 DeBruine, L. M., & Barr, D. J. (2021). Understanding Mixed-Effects Models Through Data
- 725 Simulation. *Advances in Methods and Practices in Psychological Science*, 4(1),
- 726 2515245920965119. <https://doi.org/10.1177/2515245920965119>
- 727 Gelman, A., Hill, J., & Vehtari, A. (2020). Regression and Other Stories.
- 728 [https://www.cambridge.org/highereducation/books/regression-and-other-](https://www.cambridge.org/highereducation/books/regression-and-other-stories/DD20DD6C9057118581076E54E40C372C)
- 729 stories/DD20DD6C9057118581076E54E40C372C; Cambridge University Press.
- 730 <https://doi.org/10.1017/9781139161879>
- 731 Gelman, A., Vehtari, A., Simpson, D., Margossian, C. C., Carpenter, B., Yao, Y., ...
- 732 Modrák, M. (2020). *Bayesian Workflow*. arXiv.
- 733 <https://doi.org/10.48550/arXiv.2011.01808>

- 734 Heiss, A. (2021, November 10). A Guide to Correctly Calculating Posterior Predictions and
735 Average Marginal Effects with Multilevel Bayesian Models.
736 <https://doi.org/10.59350/wbn93-edb02>
- 737 Hosmer, D. W., Lemeshow, S., & May, S. (2011). *Applied Survival Analysis: Regression*
738 *Modeling of Time to Event Data* (2nd ed). Hoboken: John Wiley & Sons.
- 739 Kantowitz, B. H., & Pachella, R. G. (2021). The Interpretation of Reaction Time in
740 Information-Processing Research 1. *Human Information Processing*, 41–82.
741 <https://doi.org/10.4324/9781003176688-2>
- 742 Kruschke, J. K., & Liddell, T. M. (2018). The Bayesian New Statistics: Hypothesis testing,
743 estimation, meta-analysis, and power analysis from a Bayesian perspective. *Psychonomic*
744 *Bulletin & Review*, 25(1), 178–206. <https://doi.org/10.3758/s13423-016-1221-4>
- 745 Kurz, A. S. (2023a). *Applied longitudinal data analysis in brms and the tidyverse* (version
746 0.0.3). Retrieved from <https://bookdown.org/content/4253/>
- 747 Kurz, A. S. (2023b). *Statistical rethinking with brms, ggplot2, and the tidyverse: Second*
748 *edition* (version 0.4.0). Retrieved from <https://bookdown.org/content/4857/>
- 749 Lakens, D. (2022). Sample Size Justification. *Collabra: Psychology*, 8(1), 33267.
750 <https://doi.org/10.1525/collabra.33267>
- 751 Landes, J., Engelhardt, S. C., & Pelletier, F. (2020). An introduction to event history
752 analyses for ecologists. *Ecosphere*, 11(10), e03238. <https://doi.org/10.1002/ecs2.3238>
- 753 McElreath, R. (2020). *Statistical Rethinking: A Bayesian Course with Examples in R and*
754 *STAN* (2nd ed.). New York: Chapman and Hall/CRC.
755 <https://doi.org/10.1201/9780429029608>
- 756 Nelder, J. A. (1999). From Statistics to Statistical Science. *Journal of the Royal Statistical*
757 *Society. Series D (The Statistician)*, 48(2), 257–269. Retrieved from
758 <https://www.jstor.org/stable/2681191>
- 759 Panis, S. (2020). How can we learn what attention is? Response gating via multiple direct
760 routes kept in check by inhibitory control processes. *Open Psychology*, 2(1), 238–279.

- 761 https://doi.org/10.1515/psych-2020-0107
- 762 Panis, S., Moran, R., Wolkersdorfer, M. P., & Schmidt, T. (2020). Studying the dynamics of
763 visual search behavior using RT hazard and micro-level speed–accuracy tradeoff
764 functions: A role for recurrent object recognition and cognitive control processes.
765 *Attention, Perception, & Psychophysics*, 82(2), 689–714.
- 766 https://doi.org/10.3758/s13414-019-01897-z
- 767 Panis, S., Schmidt, F., Wolkersdorfer, M. P., & Schmidt, T. (2020). Analyzing Response
768 Times and Other Types of Time-to-Event Data Using Event History Analysis: A Tool for
769 Mental Chronometry and Cognitive Psychophysiology. *I-Perception*, 11(6),
770 2041669520978673. https://doi.org/10.1177/2041669520978673
- 771 Panis, S., & Schmidt, T. (2016). What Is Shaping RT and Accuracy Distributions? Active
772 and Selective Response Inhibition Causes the Negative Compatibility Effect. *Journal of*
773 *Cognitive Neuroscience*, 28(11), 1651–1671. https://doi.org/10.1162/jocn_a_00998
- 774 Panis, S., & Schmidt, T. (2022). When does “inhibition of return” occur in spatial cueing
775 tasks? Temporally disentangling multiple cue-triggered effects using response history and
776 conditional accuracy analyses. *Open Psychology*, 4(1), 84–114.
- 777 https://doi.org/10.1515/psych-2022-0005
- 778 Panis, S., Torfs, K., Gillebert, C. R., Wagemans, J., & Humphreys, G. W. (2017).
779 Neuropsychological evidence for the temporal dynamics of category-specific naming.
780 *Visual Cognition*, 25(1-3), 79–99. https://doi.org/10.1080/13506285.2017.1330790
- 781 Panis, S., & Wagemans, J. (2009). Time-course contingencies in perceptual organization and
782 identification of fragmented object outlines. *Journal of Experimental Psychology: Human*
783 *Perception and Performance*, 35(3), 661–687. https://doi.org/10.1037/a0013547
- 784 Pargent, F., Koch, T. K., Kleine, A.-K., Lermer, E., & Gaube, S. (2024). A Tutorial on
785 Tailored Simulation-Based Sample-Size Planning for Experimental Designs With
786 Generalized Linear Mixed Models. *Advances in Methods and Practices in Psychological*
787 *Science*, 7(4), 25152459241287132. https://doi.org/10.1177/25152459241287132

- 788 Singer, J. D., & Willett, J. B. (2003). *Applied Longitudinal Data Analysis: Modeling Change
789 and Event Occurrence*. Oxford, New York: Oxford University Press.
- 790 Smith, P. L., & Little, D. R. (2018). Small is beautiful: In defense of the small-N design.
791 *Psychonomic Bulletin & Review*, 25(6), 2083–2101.
792 <https://doi.org/10.3758/s13423-018-1451-8>
- 793 Teachman, J. D. (1983). Analyzing social processes: Life tables and proportional hazards
794 models. *Social Science Research*, 12(3), 263–301.
795 [https://doi.org/10.1016/0049-089X\(83\)90015-7](https://doi.org/10.1016/0049-089X(83)90015-7)
- 796 Tong, C. (2019). Statistical Inference Enables Bad Science; Statistical Thinking Enables
797 Good Science. *The American Statistician*, 73(sup1), 246–261.
798 <https://doi.org/10.1080/00031305.2018.1518264>
- 799 Wickelgren, W. A. (1977). Speed-accuracy tradeoff and information processing dynamics.
800 *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)
- 801 Winter, B. (2019). *Statistics for Linguists: An Introduction Using R*. New York: Routledge.
802 <https://doi.org/10.4324/9781315165547>
- 803 Wolkersdorfer, M. P., Panis, S., & Schmidt, T. (2020). Temporal dynamics of sequential
804 motor activation in a dual-prime paradigm: Insights from conditional accuracy and
805 hazard functions. *Attention, Perception, & Psychophysics*, 82(5), 2581–2602.
806 <https://doi.org/10.3758/s13414-020-02010-5>