

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 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
23 than orthodox approaches. To help achieve broader uptake of EHA, in this paper we
24 outline a set of tutorials that detail one distributional method known as discrete-time
25 EHA. We touch upon several key aspects of the workflow, such as how to process raw data
26 and specify regression models, and we also consider the implications for experimental
27 design, as well as how to manage inter-individual differences. We finish the article by
28 considering the benefits of the approach for understanding psychological states, as well as
29 the limitations and future directions of this work. Finally, the project is written in R and
30 freely available, which 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 mean-average comparisons 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 long it lasts, how it evolves with increasing waiting time,
44 and whether its onset is time-locked to other events (Panis, 2020; Panis, Moran,
45 Wolkersdorfer, & Schmidt, 2020; Panis & Schmidt, 2016, 2022; Panis, Torfs, Gillebert,
46 Wagemans, & Humphreys, 2017; Panis & Wagemans, 2009; Wolkersdorfer, Panis, &
47 Schmidt, 2020). Such information is useful not only for the interpretation of experimental
48 effects under investigation, but also for cognitive psychophysiology and computational
49 model selection (Panis, Schmidt, Wolkersdorfer, & Schmidt, 2020).

50 As a simple illustration, Figure 1 shows the results of several simulated RT data sets,
51 which show how mean-average comparisons between two conditions can conceal the shape
52 of the underlying RT distributions. For instance, in examples 1-3, mean RT is always
53 comparable between two conditions, while the distributions differ (Figure 1, left). In
54 contrast, in examples 4-6, mean RT is lower in condition 2 compared to condition 1, but
55 the RT distributions differ in each case (Figure 1, right). Therefore, a comparison of means
56 would lead to a similar conclusion in examples 1-3, as well as examples 4-6, whereas a
57 comparison of the distributions would lead to a different conclusion in every case.

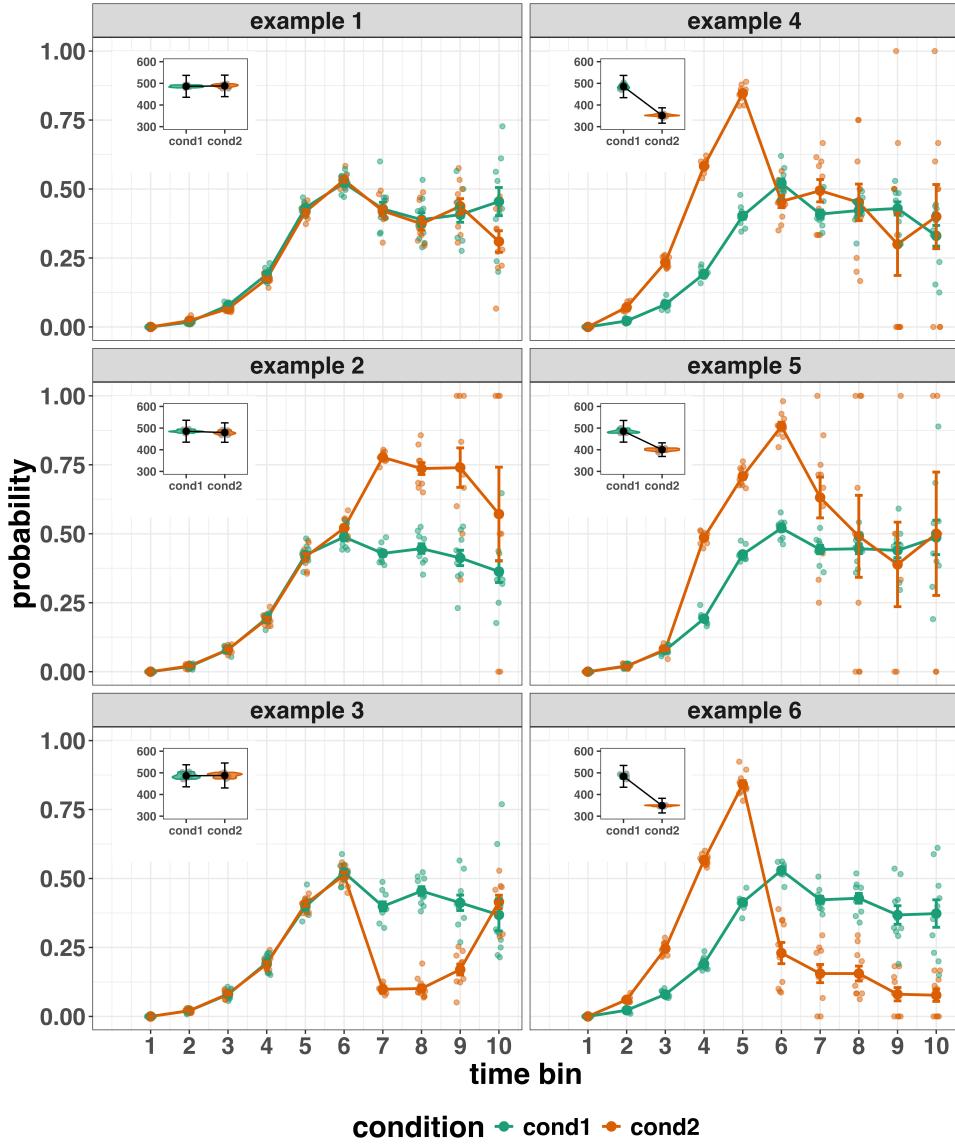


Figure 1. Means versus distributional shapes for six different simulated data set examples.

The first second after stimulus onset is divided in ten bins of 100 ms. Time bin indicates the bin rank. The first bin is (0,100], the last bin is (900,1000]. For our purposes here, it is enough to know that the distributions plotted represent the probability of an event occurring in that time bin, given that it has not yet occurred. Insets show mean response time per condition.

59 data across trials, a distributional approach offers the possibility to reveal the time course
60 of psychological states. As such, the approach permits different kinds of questions to be
61 asked, different inferences to be made, and it holds the potential to discriminate between
62 different theoretical accounts of psychological and/or brain-based processes. For example,
63 the distributions in example 4 show that the effect starts between 100 and 200 ms (in time
64 bin 2) and is gone when the waiting time reaches 500 ms or more. In contrast, in example
65 5, the effect starts around 300 ms and is gone by 700 ms. And in example 6, the effect
66 reverses between 500 and 600 ms. What kind of theory or theories could account for such
67 effects? Are there new auxiliary assumptions that theories need to adopt? And are there
68 new experiments that need to be performed to test the novel predictions that follow from
69 these analyses? As we show later using published examples, for many psychological
70 questions, such “temporal states” information can be theoretically meaningful by leading to
71 more fine-grained understanding of psychological processes, as well as adding a relatively
72 under-used dimension – the passage of time – to the theory building toolkit.

73 From a historical perspective, it is worth noting that the development of analytical
74 tools that can estimate or predict whether and when events will occur is not a new
75 endeavour. Indeed, hundreds of years ago, analytical methods were developed to predict
76 the duration of time until people died (e.g., Halley, 1693; Makeham, 1860). The same logic
77 has been applied to psychological time-to-event data, as previously demonstrated (Panis,
78 Schmidt, et al., 2020).

79 1.2 Aims and structure of the paper

80 In this paper, we focus on a distributional method for time-to-event data known as
81 discrete-time Event History Analysis (EHA), a.k.a. survival analysis, hazard analysis,
82 duration analysis, failure-time analysis, and transition analysis (Singer & Willett, 2003).
83 We hope to show the added value of EHA for knowledge and theory building in cognitive
84 psychology and related areas of research, such as cognitive neuroscience. Most importantly,

85 we provide tutorials that provide step-by-step code and instructions in the hope that we
86 can enable others to use EHA in a more routine, efficient and effective manner.

87 We first provide a brief overview of EHA to orient the reader to the basic concepts
88 that we will use throughout the paper. However, this will remain relatively short, as this
89 has been covered in detail before (Allison, 1982, 2010; Singer & Willett, 2003). Indeed, our
90 primary aim here is to introduce the set of tutorials, which explain **how** to do such
91 analyses, rather than repeat in any detail **why** you may do them.

92 We provide seven different tutorials, which are written in the R programming
93 language and publicly available on our Github page ([https://github.com/sven-panis/
94 Tutorial_Event_History_Analysis](https://github.com/sven-panis/Tutorial_Event_History_Analysis)), along with all of the other code and material
95 associated with the project. The tutorials provide hands-on, concrete examples of key parts
96 of the analytical process, so that others can apply EHA to their own time-to-event data.
97 Each tutorial is provided as an RMarkdown file, so that others can download and adapt
98 the code to fit their own purposes. Additionally, each tutorial is made available as a .html
99 file, so that it can be viewed by any web browser, and thus available to those that do not
100 use R. Finally, the manuscript itself is written in R using the papaja package (Aust &
101 Barth, 2024a), which makes it computationally reproducible, in terms of the underlying
102 data and figures.

103 In Tutorial 1a, we illustrate how to process or “wrangle” a previously published RT +
104 accuracy data set to calculate descriptive statistics when there is one independent variable.
105 The descriptive statistics are plotted, and we comment on their interpretation. In Tutorial
106 1b we provide a generalisation of this approach to illustrate how one can calculate the
107 descriptive statistics when using a more complex design, such as when there are two
108 independent variables.

109 In Tutorial 2a, we illustrate how one can fit Bayesian multilevel regression models to
110 RT data using the R package brms. We perform prior predictive checks, compare models,

and interpret the plots of the predicted hazard functions for the selected model, and the posterior distributions of our contrasts of interest. In Tutorial 2b we fit Bayesian multilevel regression models to *timed* accuracy data to perform a micro-level speed-accuracy tradeoff (SAT) analysis, which complements the EHA of RT data for choice RT data.

In Tutorial 3a, we shortly illustrate how to fit similar multilevel regression models for RT data in a frequentist framework using the R package lme4. We then briefly compare and contrast these inferential frameworks when applied to EHA. In Tutorial 3b, we illustrate how to perform the SAT analysis in a frequentist framework.

In tutorial 4, we illustrate one approach to planning how much data to collect in an experiment using EHA. We use data simulation techniques to vary sample size and trial count per condition until a certain degree of statistical power or precision is reached.

In summary, even though EHA is a widely used statistical tool and there already exist many excellent reviews (e.g., Blossfeld & Rohwer, 2002; Box-Steffensmeier, 2004; Hosmer, Lemeshow, & May, 2011; Teachman, 1983) and tutorials (e.g., Allison, 2010; Landes, Engelhardt, & Pelletier, 2020) on its general use-cases, we are not aware of any tutorials that are aimed at psychological time-to-event data, and which provide worked examples of the key data processing and multilevel regression modelling steps. Therefore, our ultimate goal is twofold: first, we want to convince readers of the many benefits of using EHA when dealing with time-to-event data with a focus on psychological time-to-event data, and second, we want to provide a set of practical tutorials, which provide step-by-step instructions on how you actually perform a discrete-time EHA on time-to-event data such as RT data, as well as a complementary discrete-time SAT analysis on timed accuracy data.

2. A brief introduction to event history analysis

We recommend several excellent textbooks for a comprehensive background context to EHA (Allison, 2010; Singer & Willett, 2003) and for a more general introduction to

136 understanding regression equations (Gelman, Hill, & Vehtari, 2020; Winter, 2019). Our
137 focus here is not on providing a detailed account of the underlying regression equations,
138 since this topic has been comprehensively covered many times before. Instead, we want to
139 provide an intuition regarding how EHA works in general, as well as in the context of
140 experimental psychology. As such, we only supply regression equations in section D of the
141 Supplemental Material.

142 **2.1 Basic features of event history analysis**

143 To apply EHA, one must be able to:

- 144 1. define an event of interest that represents a qualitative change that can be situated in
145 time (e.g., a button press, a saccade onset, a fixation offset, etc.);
- 146 2. define time point zero (e.g., target stimulus onset, fixation onset, etc.);
- 147 3. measure the passage of time between time point zero and event occurrence in discrete
148 or continuous time units.

149 In EHA, the definition of hazard and the type of models employed depend on
150 whether one is using continuous or discrete time units. Since our focus here is on hazard
151 models that use discrete time units, we describe that approach. After dividing time in
152 discrete, contiguous time bins indexed by t (e.g., $t = 1:10$ time bins), let RT be a discrete
153 random variable denoting the rank of the time bin in which a particular person's response
154 occurs in a particular experimental condition. For example, the first response might occur
155 at 546 ms and it would be in time bin 6 (any RTs from 501 ms to 600 ms).

156 Discrete-time EHA focuses on the discrete-time hazard function of event occurrence
157 and the discrete-time survivor function (Figure 2). The equations that define both of these
158 functions are reported in section A of the Supplemental Material. The discrete-time hazard

159 function gives you, for each time bin, the probability that the event occurs (sometime) in
160 bin t, given that the event does not occur in previous bins. In other words, it reflects the
161 instantaneous likelihood that the event occurs in the current bin, given that it has not yet
162 occurred in the past, i.e., in one of the prior bins. In contrast, the discrete-time survivor
163 function cumulates the bin-by-bin risks of event *nonoccurrence* to obtain the survival
164 probability, the probability that the event occurs after bin t. In other words, the survivor
165 function gives you for each time bin the likelihood that the event occurs in the future, i.e.,
166 in one of the subsequent time bins.

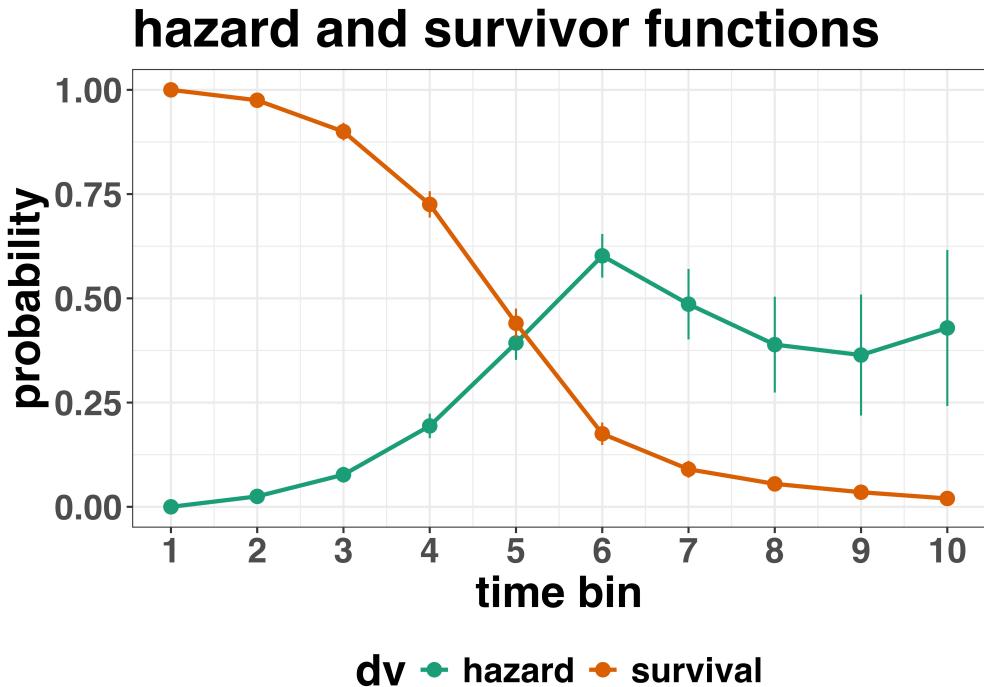


Figure 2. Discrete-time hazard and survivor functions. Discrete time-to-event data were simulated for 200 trials of 1 experimental condition. Error bars represent ± 1 standard error of the respective proportion. While the hazard function is the vehicle for inferring the time course of cognitive processes, the survival probability $S(t-1)$ can help to qualify or provide context to the interpretation of the hazard probability $h(t)$. For example, the high hazard of $.60 = h(t=6)$ is experienced only by 44 percent of the trials, as $S(t=5) = .44$. Because the survivor function is a decreasing function of time, the error bars in later parts of the hazard function will always be wider and less precise compared to earlier parts.

¹⁶⁷ 2.2 Benefits of event history analysis

¹⁶⁸ Statisticians and mathematical psychologists recommend focusing on the hazard
¹⁶⁹ function when analyzing time-to-event data for various reasons. We do not cover these
¹⁷⁰ benefits in detail here, as these are more general topics that have been covered elsewhere in
¹⁷¹ textbooks. Instead, we briefly list the benefits below, and refer the reader to section F of
¹⁷² the Supplemental Material for more detailed coverage of the benefits. The benefits include:

- 173 1. Hazard functions are more diagnostic than density functions when one is interested in
174 studying the detailed shape of a RT distribution (Holden et al., 2009).
- 175 2. RT distributions may differ from each other in multiple ways, and hazard functions
176 allow one to capture these differences which mean-average comparisons may conceal
177 (Townsend, 1990).
- 178 3. EHA takes account of more of the data collected in a typical speeded response
179 experiment, by virtue of not discarding right-censored observations. Trials with very
180 long RTs are not discarded, but instead contribute to the risk set in each time bin
181 (see section 4.1.2 below).
- 182 4. Hazard modeling allows one to incorporate time-varying explanatory covariates, such
183 as heart rate, electroencephalogram (EEG) signal amplitude, gaze location, etc.
184 (Allison, 2010). This is useful for linking physiological effects to behavioral effects
185 when performing cognitive psychophysiology (Meyer, Osman, Irwin, & Yantis, 1988).
- 186 5. EHA can help to solve the problem of model mimicry, i.e., the fact that different
187 computational models can often predict the same mean RTs as observed in the
188 empirical data, but not necessarily the detailed shapes of the empirical RT hazard
189 distributions. As such, EHA can be a tool to help distinguish between competing
190 theories of cognition and brain function.

191 **2.3 Event history analysis in the context of experimental psychology**

192 To make EHA more relevant to researchers studying cognitive psychology and

193 cognitive neuroscience, in this section we provide a relevant worked example and consider
194 implications that are relevant to that domain of research.

195 **2.3.1 A worked example.** In the context of experimental psychology, it is

196 common for participants to be presented with either a 1-button detection task or a

197 discrimination task. For example, a task may involve choosing between two response
198 options with only one of them being correct. For such two-choice RT data, the
199 discrete-time EHA of the RT data (hazard and survivor functions) can be extended with a
200 discrete-time SAT analysis of the timed accuracy data. Specifically, the hazard function of
201 event occurrence can be extended with the discrete-time conditional accuracy function,
202 which gives you the probability that a response is correct given that it is emitted in time
203 bin t (Allison, 2010; Kantowitz & Pachella, 2021; Wickelgren, 1977). We refer to this
204 extended (hazard + conditional accuracy) analysis for choice RT data as EHA/SAT.

205 Integrating results between hazard and conditional accuracy functions for choice RT
206 data can be informative for understanding psychological processes. To illustrate, we
207 consider a hypothetical choice RT example that is inspired by real data (Panis & Schmidt,
208 2016), but simplified to make the main point clearer (Figure 3). In a standard priming
209 paradigm, there is a prime stimulus (e.g., an arrow pointing left or right) followed by a
210 target stimulus (another arrow pointing left or right). The prime can then be congruent or
211 incongruent with the target.

212 Figure 3 shows that the early upswing in hazard is equal for both priming conditions
213 (Figure 3, upper panel), and that early emitted responses are always correct in the
214 congruent condition and always incorrect in the incongruent condition (Figure 3, lower
215 panel). These results show that for short waiting times (< bin 6), responses always follow
216 the prime (and not the target, as instructed). During time bin 6 the target-triggered
217 response channel is activated and causes response competition – $ca(6) = .5$ – and a lower
218 hazard probability in the incongruent condition. For waiting times of 600 ms or more, the
219 hazard of response occurrence is lower in incongruent compared to congruent trials, and all
220 responses emitted in these late bins are correct.

221 This joint pattern of results is interesting because it can provide meaningfully
222 different conclusions about psychological processes compared to conventional analyses, such

as computing mean-average RT and accuracy across trials. Mean-average RT would only represent the overall ability of cognition to overcome interference, on average, across trials. For instance, if mean-average RT was higher in incongruent than congruent trials, one may conclude that cognitive mechanisms that support interference control are working as expected across trials, and are indexed by each recorded response. But such a conclusion is not supported when the effects are explored over a timeline. Instead, the psychological conclusion is much more nuanced and suggests that multiple states start, stop and possibly interact over a particular temporal window.

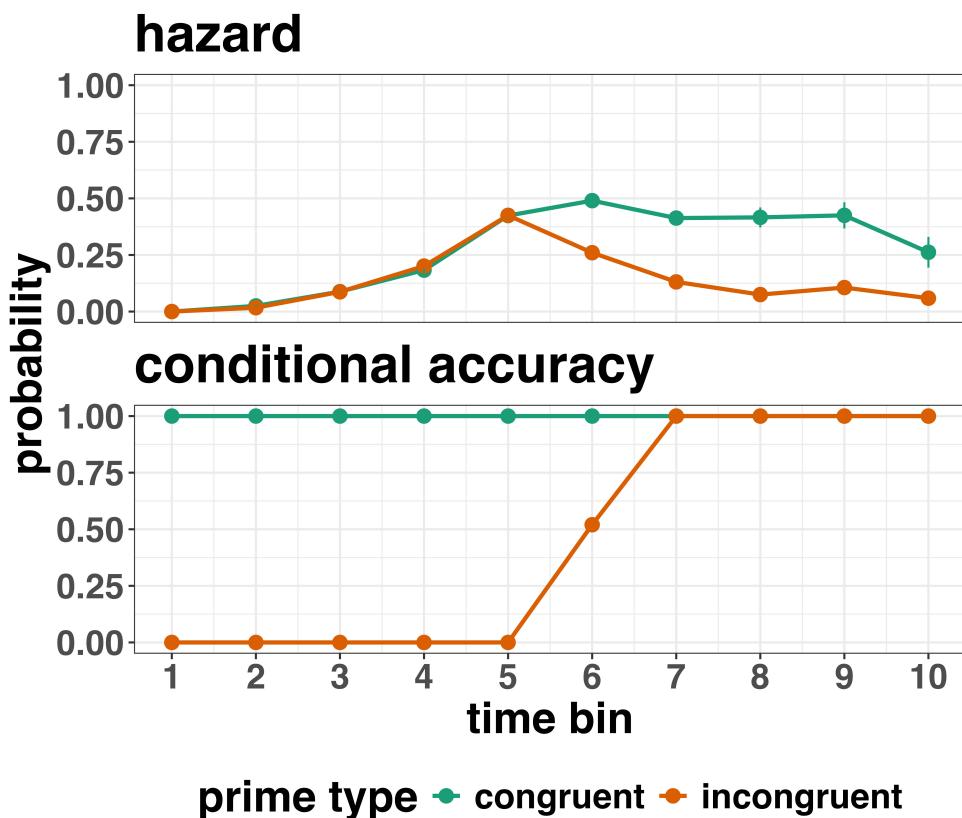


Figure 3. Discrete-time hazard and conditional accuracy functions. Discrete time-to-event and conditional accuracy data were simulated for 1000 trials for each of two priming conditions (congruent and incongruent prime stimuli). Error bars represent ± 1 standard error of the respective proportion. Bin width equals 100 ms.

Unlocking the temporal states of cognitive processes can be revealing for theory development and the understanding of basic psychological processes. Possibly more importantly, however, is that it simultaneously opens the door to address many new and previously unanswered questions. Do all participants show similar temporal states or are there individual differences? Do such individual differences extend to those individuals that have been diagnosed with some form of psychopathology? How do temporal states relate to brain-based mechanisms that might be studied using other methods from cognitive neuroscience? And how much of theory in cognitive psychology would be in need of revision if mean-average comparisons were supplemented with a temporal states approach?

2.3.2 Implications for designing experiments. Performing EHA in experimental psychology has implications for how experiments are designed. Indeed, if trials are categorised as a function of when responses occur, then each time bin will only include a subset of the total number of trials. For example, let's consider an experiment where each participant performs 2 conditions and there are 100 trial repetitions per condition. Those 100 trials must be distributed in some manner across the chosen number of bins.

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

In contrast to the large- N design that typically average across many participants without being able to scrutinize individual data patterns, small- N designs retain crucial information about the data patterns of individual observers. This can be advantageous

258 whenever participants differ systematically in their strategies or in the time courses of their
259 effects, so that averaging them would lead to misleading data patterns. Note that because
260 statistical power derives both from the number of participants and from the number of
261 repeated measures per participant and condition, small- N designs can still achieve what
262 are generally considered acceptable levels of statistical power, if they have a sufficient
263 amount of data overall (Baker et al., 2021; Smith & Little, 2018).

264 **3. An overview of the general analytical workflow**

265 Although the focus is on EHA/SAT, we also want to briefly comment on broader
266 aspects of our general analytical workflow, which relate more to data science and data
267 analysis workflows.

268 **3.1 Data science workflow and descriptive statistics**

269 We perform data wrangling following tidyverse principles and a functional
270 programming approach (Wickham, Çetinkaya-Rundel, & Grolemund, 2023). In short,
271 functional programming means that you avoid writing your own loops and instead use
272 functions that have been built and tested by others. In addition, we also supply a set of
273 custom-built functions, which make the process of data wrangling in the context of data
274 preparation and descriptive statistics a lot quicker and more efficient.

275 **3.2 Inferential statistical approach**

276 Our lab adopts an estimation approach to multilevel regression (Kruschke & Liddell,
277 2018; Winter, 2019), which is heavily influenced by the Bayesian framework as suggested
278 by Richard McElreath (Kurz, 2023b; McElreath, 2020). We also use a “keep it maximal”
279 approach to specifying varying (or random) effects (Barr, Levy, Scheepers, & Tily, 2013).
280 This means that wherever possible we include varying intercepts and slopes per participant.

281 To make inferences, we use two main approaches. We compare models of different
 282 complexity, using information criteria (e.g., WAIC) and cross-validation (e.g., LOO), to
 283 evaluate out-of-sample predictive accuracy (McElreath, 2020). We also take the most
 284 complex model and evaluate key parameters of interest using point and interval estimates.

285 **3.3 Implementation**

286 We used R (Version 4.4.0; R Core Team, 2024)¹ for all reported analyses. The
 287 content of the tutorials, in terms of EHA and multilevel regression modelling, is mainly
 288 based on Allison (2010), Singer and Willett (2003), McElreath (2020), Heiss (2021), Kurz
 289 (2023a), and Kurz (2023b).

290 **4. Tutorials**

291 Tutorials 1a and 1b show how to calculate and plot the descriptive statistics of
 292 EHA/SAT when there are one or two independent variables, respectively. Tutorials 2a and
 293 2b illustrate how to use Bayesian multilevel modeling to fit hazard and conditional

¹ We, furthermore, used the R-packages *bayesplot* (Version 1.11.1; Gabry, Simpson, Vehtari, Betancourt, & Gelman, 2019), *brms* (Version 2.22.0; Bürkner, 2017, 2018, 2021), *citr* (Version 0.3.2; Aust, 2019), *cmdstanr* (Version 0.8.1.9000; Gabry, Češnovar, Johnson, & Broder, 2024), *dplyr* (Version 1.1.4; Wickham, François, Henry, Müller, & Vaughan, 2023), *forcats* (Version 1.0.0; Wickham, 2023a), *ggplot2* (Version 3.5.1; Wickham, 2016), *lme4* (Version 1.1.35.5; Bates, Mächler, Bolker, & Walker, 2015), *lubridate* (Version 1.9.3; Grolemund & Wickham, 2011), *Matrix* (Version 1.7.1; Bates, Maechler, & Jagan, 2024), *nlme* (Version 3.1.166; Pinheiro & Bates, 2000), *papaja* (Version 0.1.3; Aust & Barth, 2024b), *patchwork* (Version 1.3.0; Pedersen, 2024), *purrr* (Version 1.0.2; Wickham & Henry, 2023), *RColorBrewer* (Version 1.1.3; Neuwirth, 2022), *Rcpp* (Eddelbuettel & Balamuta, 2018; Version 1.0.13.1; Eddelbuettel & François, 2011), *readr* (Version 2.1.5; Wickham, Hester, & Bryan, 2024), *RJ-2021-048* (Bengtsson, 2021), *rstan* (Version 2.32.6; Stan Development Team, 2024), *standist* (Version 0.0.0.9000; Girard, 2024), *StanHeaders* (Version 2.32.10; Stan Development Team, 2020), *stringr* (Version 1.5.1; Wickham, 2023b), *tibble* (Version 3.2.1; Müller & Wickham, 2023), *tidybayes* (Version 3.0.7; Kay, 2024), *tidyverse* (Version 1.3.1; Wickham, Vaughan, & Girlich, 2024), *tidyverse* (Version 2.0.0; Wickham et al., 2019) and *tinylabels* (Version 0.2.4; Barth, 2023).

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

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

4.1.1 Data wrangling aims.

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

Second, we want to produce two different data sets that can each be submitted to different types of inferential modelling approaches. The two types of data structure we label as ‘person-trial’ data and ‘person-trial-bin’ data. The ‘person-trial’ data (Table 1) will be familiar to most researchers who record behavioural responses from participants, as

³¹⁹ it represents the measured RT and accuracy per trial within an experiment. This data set
³²⁰ is 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.

³²¹ In contrast, the ‘person-trial-bin’ data (Table 2) has a different, more extended
³²² structure, which indicates in which bin a response occurred, if at all, in each trial.
³²³ Therefore, the ‘person-trial-bin’ data generates a 0 in each bin until an event occurs and
³²⁴ then it generates a 1 to signal an event has occurred in that bin. This data set is used
³²⁵ when fitting hazard models (Tutorials 2a and 3a). It is worth pointing out that there is no
³²⁶ requirement for an event to occur at all (in any bin), as maybe there was no response on
³²⁷ that trial or the event occurred after the time window of interest. Likewise, when the event

328 occurs in bin 1 there would only be one row of data for that trial in the person-trial-bin
329 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.

330 **4.1.2 A real data wrangling example.** To illustrate how to quickly set up life
331 tables for calculating the descriptive statistics (functions of discrete time), we use a
332 published data set on masked response priming from Panis and Schmidt (2016). In their
333 first experiment, Panis and Schmidt (2016) presented a double arrow for 94 ms that
334 pointed left or right as the target stimulus with an onset at time point zero in each trial.
335 Participants had to indicate the direction in which the double arrow pointed using their
336 corresponding index finger, within 800 ms after target onset. Response time and accuracy

337 were recorded on each trial. Prime type (blank, congruent, incongruent) and mask type
 338 were manipulated. Here we focus on the subset of trials in which no mask was presented.
 339 The 13-ms prime stimulus was a double arrow presented 187 ms before target onset in the
 340 congruent (same direction as target) and incongruent (opposite direction as target) prime
 341 conditions.

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

345 The required column names are as follows:

- 346 • “pid”, indicating unique participant IDs;
- 347 • “trial”, indicating each unique trial per participant;
- 348 • “condition”, a factor indicating the levels of the independent variable (1, 2, ...) and
 349 the corresponding labels;
- 350 • “rt”, indicating the response times in ms;
- 351 • “acc”, indicating the accuracies (1/0).

352 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")))
```

353 Next, we can set up the life tables and plots of the discrete-time functions $h(t)$, $S(t)$,
 354 $ca(t)$, and $P(t)$ – see section A of the Supplemental Material for their definitions. To do so

355 using a functional programming approach, one has to nest the data within participants
 356 using the group_nest() function, and supply a user-defined censoring time and bin width
 357 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))
```

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

367 The person-trial-bin oriented data set is created by our custom function ptb(), and it
 368 has one row for each time bin (of each trial) that is at risk for event occurrence. The

369 variable “event” in the person-trial-bin oriented data set indicates whether a response
370 occurs (1) or not (0) for each bin.

371 The next step is to set up the life table using our custom function `setup_lt()`,
372 calculate the conditional accuracies using our custom function `calc_ca()`, add the `ca(t)`
373 estimates to the life table using our custom function `join_lt_ca()`, and then plot the
374 descriptive statistics using our custom function `plot_eha()`. When creating the plots, some
375 warning messages will likely be generated, like these:

- 376 • Removed 2 rows containing missing values or values outside the scale range
377 (`geom_line()`).
378 • Removed 2 rows containing missing values or values outside the scale range
379 (`geom_point()`).
380 • Removed 2 rows containing missing values or values outside the scale range
381 (`geom_segment()`).

382 The warning messages are generated because some bins have no hazard and `ca(t)`
383 estimates, and no error bars. They can thus safely be ignored. One can now inspect
384 different aspects, including the life table for a particular condition of a particular subject,
385 and a plot of the different functions for a particular participant. In general, it is important
386 to visually inspect the functions first for each participant, in order to identify individuals
387 that may be guessing (e.g., a flat conditional accuracy function at .5 indicates that
388 someone is just guessing), outlying individuals, and/or different groups with qualitatively
389 different behavior.

390 Table 3 shows the life table for condition “blank” (no prime stimulus presented) for
391 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.

Figure 4 displays the discrete-time hazard, survivor, conditional accuracy, and

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

Descriptive stats for subject 6

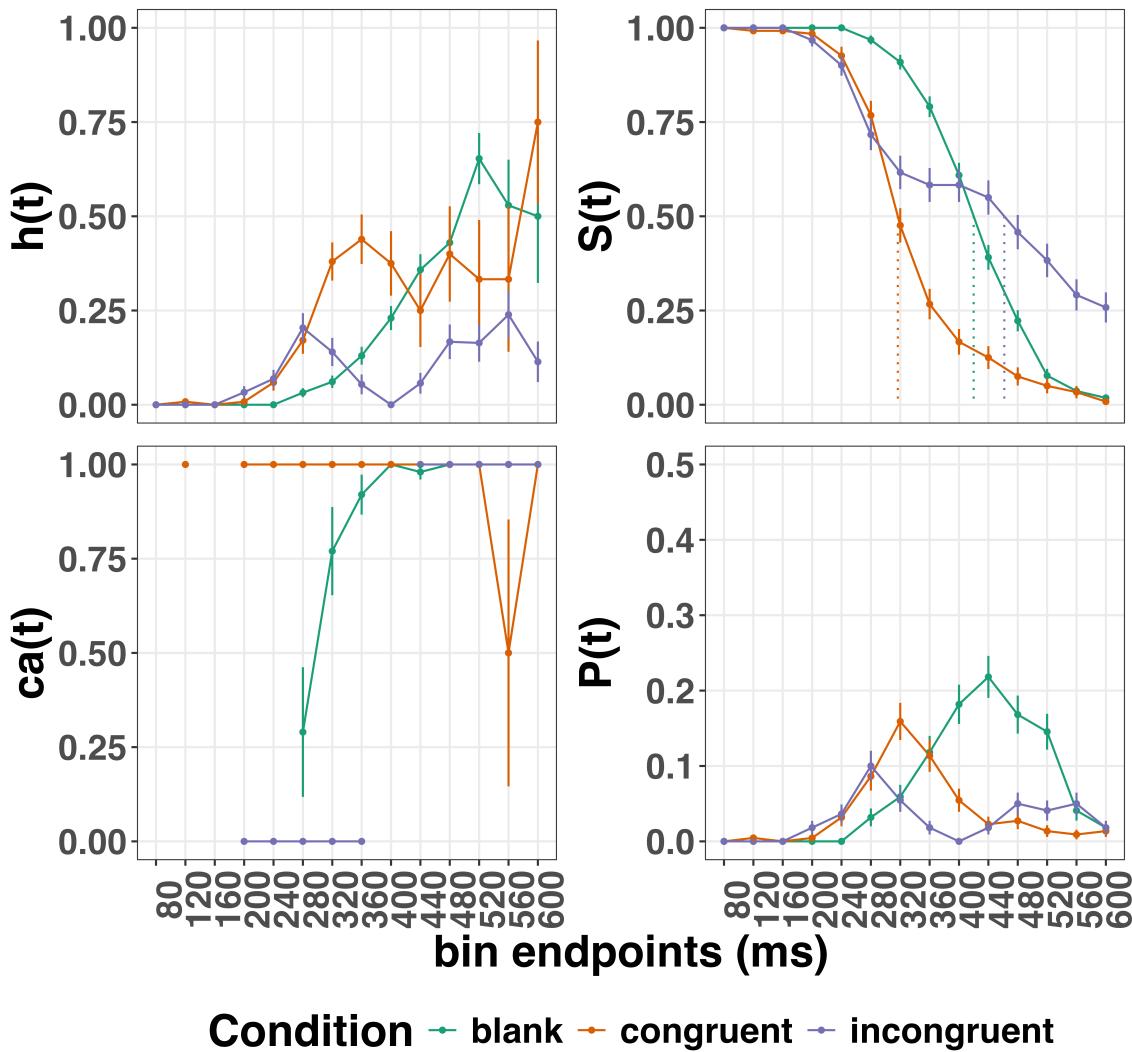


Figure 4. Estimated discrete-time hazard (h), survivor (S), conditional accuracy (ca) and probability mass (P) functions for participant 6. Vertical dotted lines indicate the estimated median RTs. Error bars represent ± 1 standard error of the respective proportion.

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

399 0.17, and 0.20 for the blank, congruent, and incongruent prime conditions, respectively. In

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

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

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

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

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

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

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

407 respectively.

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

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

410 0.36, 0.25, and 0.06 for the blank, congruent, and incongruent prime conditions,

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

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

413 conditions, respectively.

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

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

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

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

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

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

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

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

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

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

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

425 At this point, we have calculated, summarised and plotted descriptive statistics for
426 the key variables in EHA/SAT. As we will show in later Tutorials, statistical models for
427 $h(t)$ and $ca(t)$ can be implemented as generalized linear mixed regression models predicting
428 event occurrence (1/0) and conditional accuracy (1/0) in each bin of a selected time
429 window for analysis. But first we consider calculating the descriptive statistics for two
430 independent variables.

431 **4.2 Tutorial 1b: Generalising to a more complex design**

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

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

443 **4.3 Tutorial 2a: Fitting Bayesian hazard models to discrete time-to-event data**

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

448 Willett, 2003).

449 **4.3.1 Hazard model considerations.** There are several analytic decisions one
450 has to make when fitting a discrete-time hazard model. First, one has to select an analysis
451 time window, i.e., a contiguous set of bins for which there is enough data for each
452 participant. Second, given that the dependent variable (event occurrence) is binary, one
453 has to select a link function (see section C in the Supplemental Material). The cloglog link
454 is preferred over the logit link when events can occur in principle at any time point within
455 a bin, which is the case for RT data (Singer & Willett, 2003). Third, one has to choose
456 whether to treat TIME (i.e., the time bin index t) as a categorical or continuous predictor.
457 And when you treat a variable as a categorical predictor, you can choose between reference
458 coding and index coding. With reference coding, one defines the variable as a factor and
459 selects one of the k categories as the reference level. `Brm()` will then construct $k-1$
460 indicator variables (see model M1d in Tutorial_2a.Rmd for an example). With index
461 coding, one constructs an index variable that contains integers that correspond to different
462 categories (see models M0i and M1i below). As explained by McElreath (2020), the
463 advantage of index coding is that the same prior can be assigned to each level of the index
464 variable, so that each category has the same prior uncertainty.

465 In the case of a large- N design without repeated measurements, the parameters of a
466 discrete-time hazard model can be estimated using standard logistic regression software
467 after expanding the typical person-trial data set into a person-trial-bin data set (Allison,
468 2010). When there is clustering in the data, as in the case of a small- N design with
469 repeated measurements, the parameters of a discrete-time hazard model can be estimated
470 using population-averaged methods (e.g., Generalized Estimating Equations), and Bayesian
471 or frequentist generalized linear mixed models (Allison, 2010).

472 In general, there are three assumptions one can make or relax when adding
473 experimental predictor variables and other covariates: The linearity assumption for
474 continuous predictors (the effect of a 1 unit change is the same anywhere on the scale), the

475 additivity assumption (predictors do not interact), and the proportionality assumption
 476 (predictors do not interact with TIME).

477 In tutorial_2a.Rmd we fit several Bayesian multilevel models (i.e., generalized linear
 478 mixed models) that differ in complexity to the person-trial-bin oriented data set that we
 479 created in Tutorial 1a. We decided to select the analysis time window (200,600] and the
 480 cloglog link. Below, we shortly discuss two of these models. The person-trial-bin data set is
 481 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)))
```

482 **4.3.2 Prior distributions.** To get the posterior distribution of each model
 483 parameter given the data, we need to specify prior distributions for the model parameters
 484 which reflect our prior beliefs. In Tutorial_2a.Rmd we perform a few prior predictive
 485 checks to make sure our selected prior distributions reflect our prior beliefs (Gelman,
 486 Vehtari, et al., 2020).

487 The middle column of Supplementary Figure 2 (section E of the Supplemental
 488 Material) shows six examples of prior distributions for an intercept on the logit and/or
 489 cloglog scales. While a normal distribution with relatively large variance is often used as a

490 weakly informative prior for continuous dependent variables, rows A and B of
 491 Supplementary Figure 2 show that specifying such distributions on the logit and cloglog
 492 scales actually leads to rather informative distributions on the original probability scale, as
 493 most mass is pushed to probabilities of 0 and 1.

494 **4.3.3 Model M0i: A null model with index coding.** When you do not want to
 495 make assumptions about the shape of the hazard function, or its shape is not smooth but
 496 irregular, then you can use a general specification of TIME, i.e., fit one grand intercept per
 497 time bin. In this first model, we use a general specification of TIME using index coding,
 498 and do not include experimental predictors. We call this model “M0i”.

499 Before we fit model M0i, we select the necessary columns from the data, and specify
 500 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")
```

501 After selecting the bernoulli family and the cloglog link, the model formula is
 502 specified. The specification “ $0 + \dots$ ” removes the default intercept in brm(). The fixed
 503 effects include an intercept for each level of timebin. Each of these intercepts is allowed to

504 vary across individuals (variable pid). We request 2000 samples from the posterior
 505 distribution for each of four chains. Estimating model M0i took about 30 minutes on a
 506 MacBook Pro (Sonoma 14.6.1 OS, 18GB Memory, M3 Pro Chip).

507 **4.3.4 Model M1i: Adding the effects of prime-target congruency.** Previous
 508 research has shown that psychological effects typically change over time (Panis, 2020;
 509 Panis, Moran, et al., 2020; Panis & Schmidt, 2022; Panis et al., 2017; Panis & Wagemans,
 510 2009). In the next model, therefore, we use index coding for both TIME (variable
 511 “timebin”) and the categorical predictor prime-target-congruency (variable “prime”), so
 512 that we get 30 grand intercepts, one for each combination of timebin level and prime level.
 513 Here is the model formula of this model that we call “M1i”.

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

514 Estimating model M1i took about 124 minutes.

515 **4.3.5 Compare the models.** We can compare the two models using the Widely
 516 Applicable Information Criterion (WAIC) and Leave-One-Out (LOO) cross-validation, and
 517 look at model weights for both criteria (Kurz, 2023a; McElreath, 2020).

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

518 ## model_M0i model_M1i
 519 ## 0 1

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

520 ## model_M0i model_M1i
 521 ## 0 1

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

524 **4.3.6 Evaluating parameter estimates in model M1i.** To make inferences

525 from the parameter estimates in model M1i, we first plot the densities of the draws from
 526 the posterior distributions of its population-level parameters in Figure 5, together with
 527 point (median) and interval estimates (80% and 95% credible intervals).

Posterior distributions for population-level effects in Model M1i

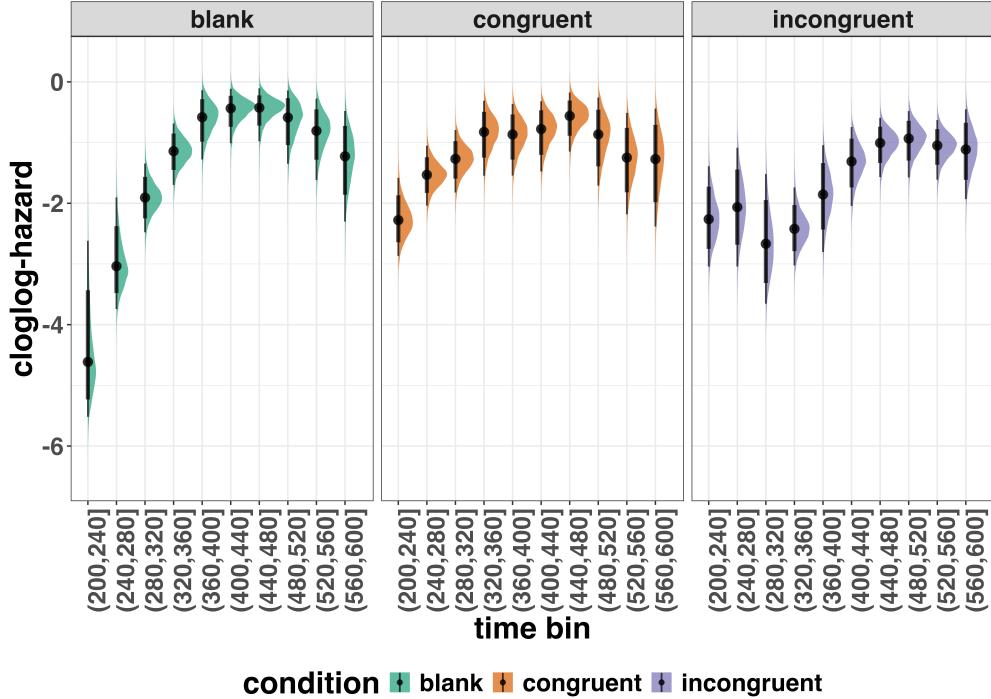


Figure 5. Medians and 80/95% credible intervals of the posterior distributions of the population-level parameters of model M1i.

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

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

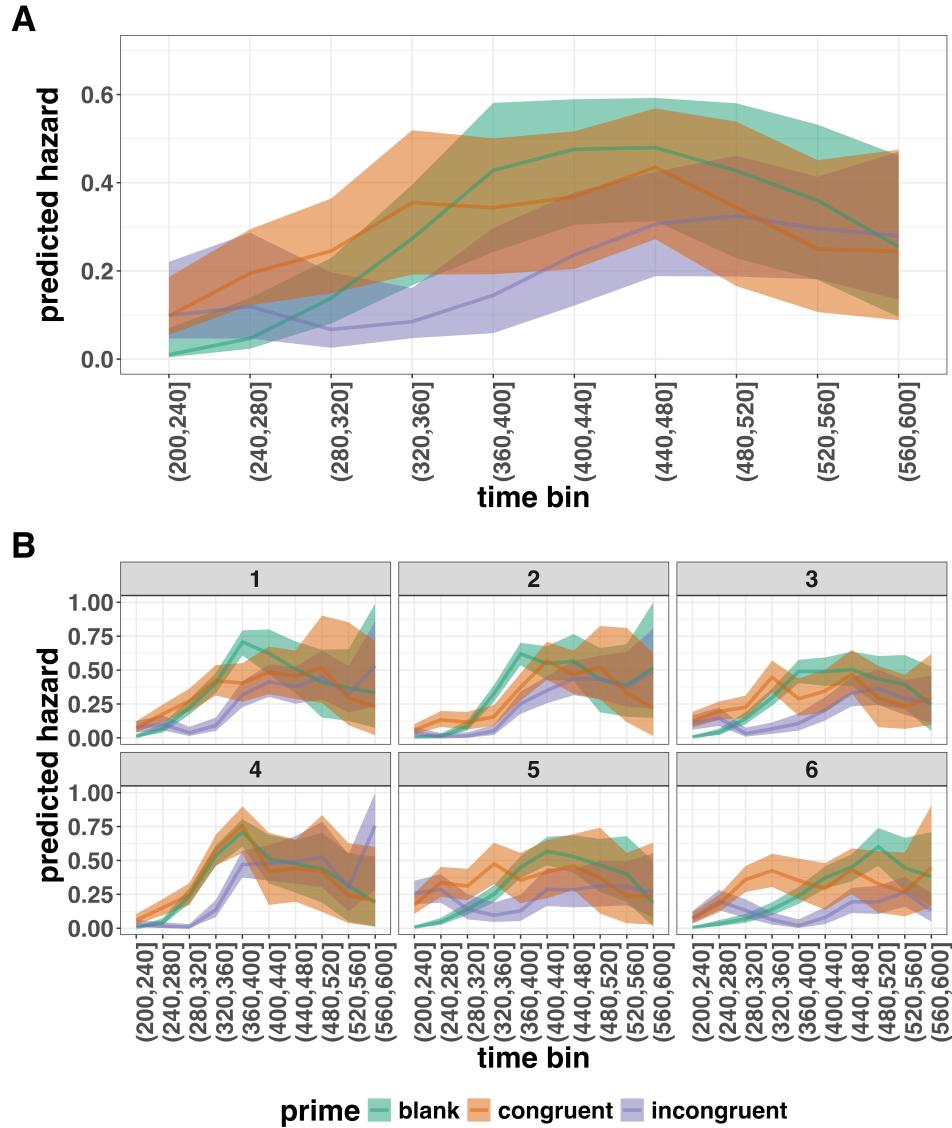


Figure 6. Point (median) and 80/95% credible interval summaries of the hazard estimates (expected values of the draws from the posterior predictive distributions) in each time bin at the population level (A), and for each participant (B).

532 As we are actually interested in the effects of congruent and incongruent primes,
 533 relative to the blank prime condition, we can construct two contrasts (congruent-blank,
 534 incongruent-blank), and plot the posterior distributions of these contrast effects, both at
 535 the population level (Figure 7A; grand average marginal effect) and at the participant level

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

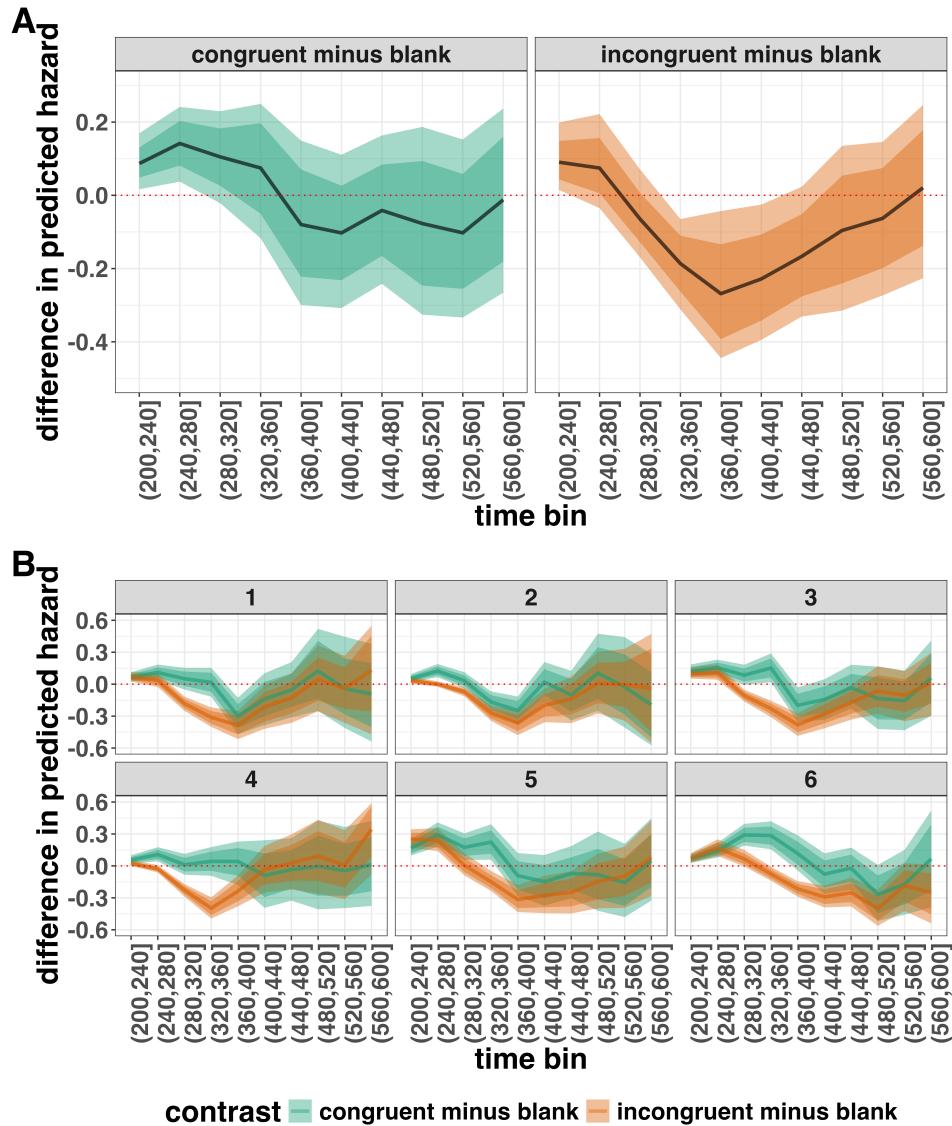


Figure 7. Point (mean) and 80/95% credible interval summaries of estimated differences in hazard in each time bin at the population level (A), and for each participant (B).

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

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

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

544 If we want to estimate the population-level effect of prime type on hazard, we can
545 base our conclusion on Figure 7A. The contrast “congruent minus blank” was estimated to
546 be 0.09 hazard units in bin (200,240] (95% CrI = [0.02, 0.17]), and 0.14 hazard units in bin
547 (240,280]) (95% CrI = [0.04, 0.25]). For the other bins, the 95% credible interval contained
548 zero. The contrast “incongruent minus blank” was estimated to be 0.09 hazard units in bin
549 (200,240] (95% CrI = [0.01, 0.21]), -0.19 hazard units in bin (320,360] (95% CrI = [-0.31,
550 -0.06]), -0.27 hazard units in bin (360,400] (95% CrI = [-0.45, -0.04]), and -0.23 hazard
551 units in bin (400,440] (95% CrI = [-0.40, -0.03]). For the other bins, the 95% credible
552 interval contained zero.

553 There are thus two phases of performance for the average person between 200 and
554 600 ms after target onset. In the first phase, the addition of a congruent or incongruent
555 prime stimulus increases the hazard of response occurrence compared to blank prime trials
556 in the time period (200, 240]. In the second phase, only the incongruent prime decreases
557 the hazard of response occurrence compared to blank primes, in the time period (320,440].
558 The sign of the effect of incongruent primes on the hazard of response occurrence thus
559 depends on how much waiting time has passed since target onset.

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

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

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

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

581 4.4 Tutorial 2b: Fitting Bayesian conditional accuracy models

582 In this fourth tutorial, we illustrate how to fit a Bayesian multilevel regression model
583 to the timed accuracy data from the masked response priming data used in Tutorial 1a.
584 The general process is similar to Tutorial 2a, except that (a) we use the person-trial data,
585 (b) we use the logit link function, and (c) we change the priors. To keep the tutorial short,
586 we only fit one conditional accuracy model, which was based on model M1i from Tutorial
587 2a and labelled M1i_ca.

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

Posterior distributions for population-level effects in Model M1i_ca

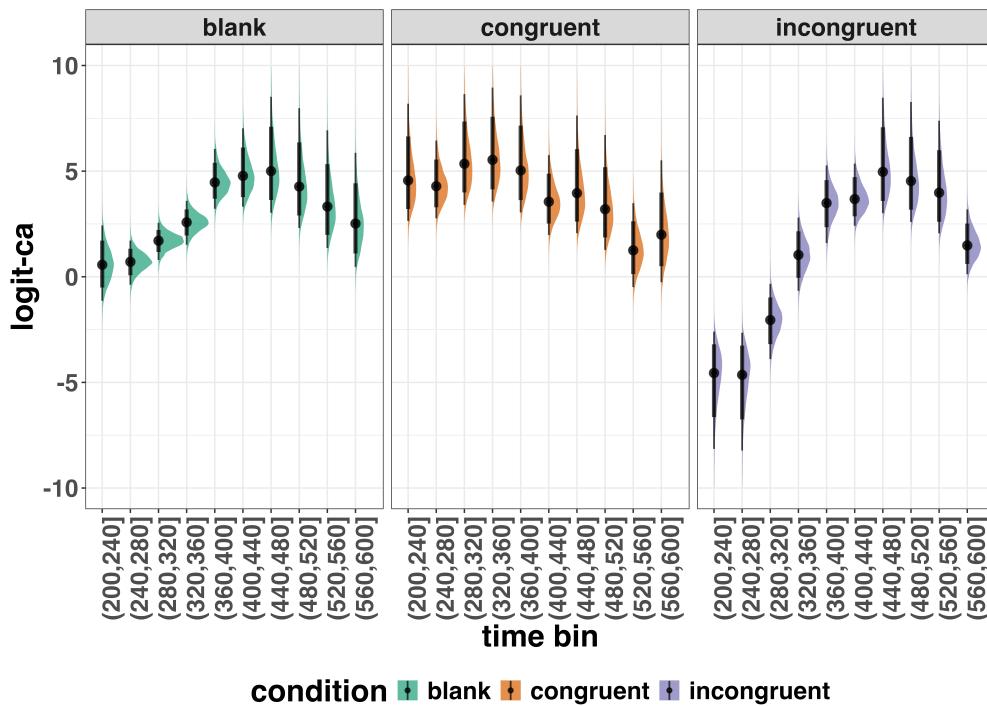


Figure 8. Medians and 80/95% credible intervals of the posterior distributions of the population-level parameters of model M1i_ca. ca = conditional accuracy.

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).

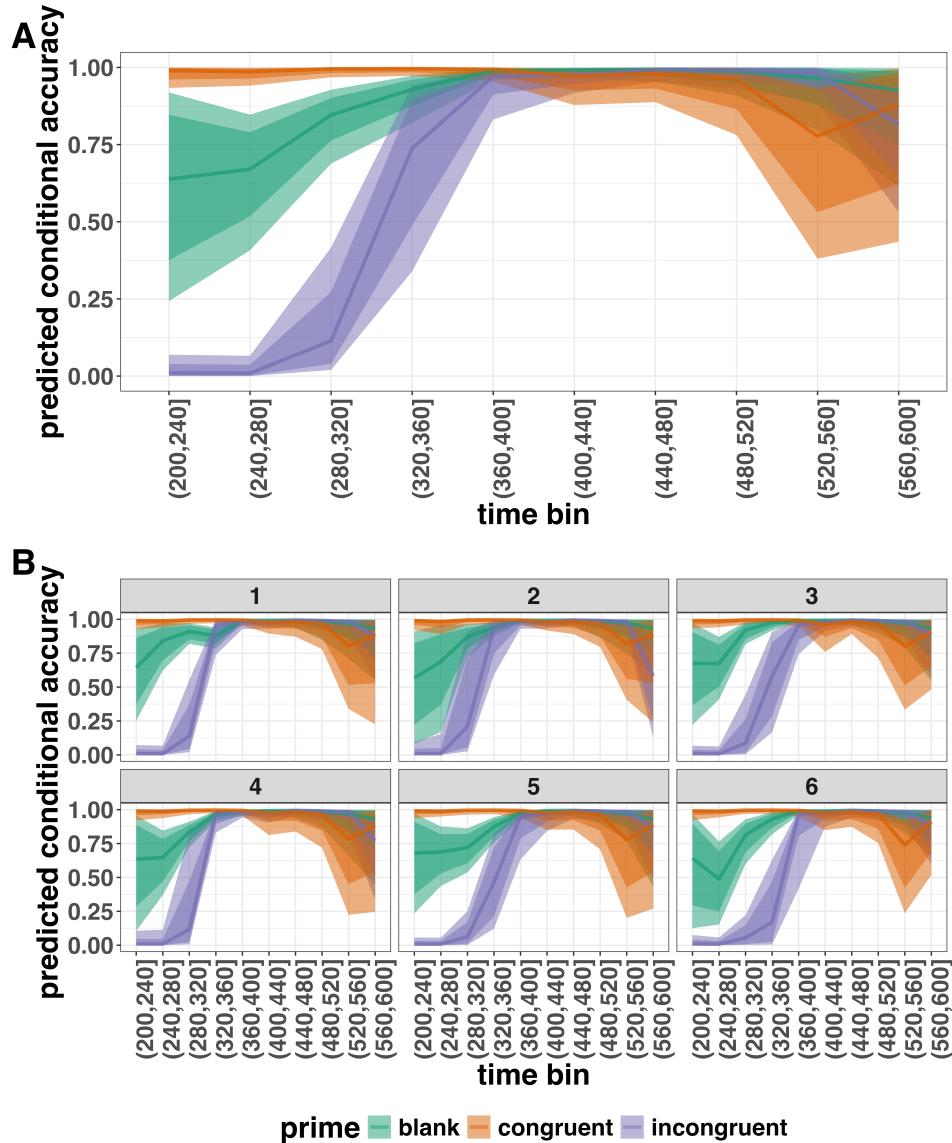


Figure 9. Point (median) and 80/95% credible interval summaries of the conditional accuracy estimates (expected values of the draws from the posterior predictive distributions) in each time bin at the population level (A), and for each participant (B).

596 As we are actually interested in the effects of congruent and incongruent primes,
 597 relative to the blank prime condition, we can construct two contrasts (congruent-blank,
 598 incongruent-blank), and plot the posterior distributions of these contrast effects at the
 599 population level (Figure 10A) and for each participant (Figure 10B).

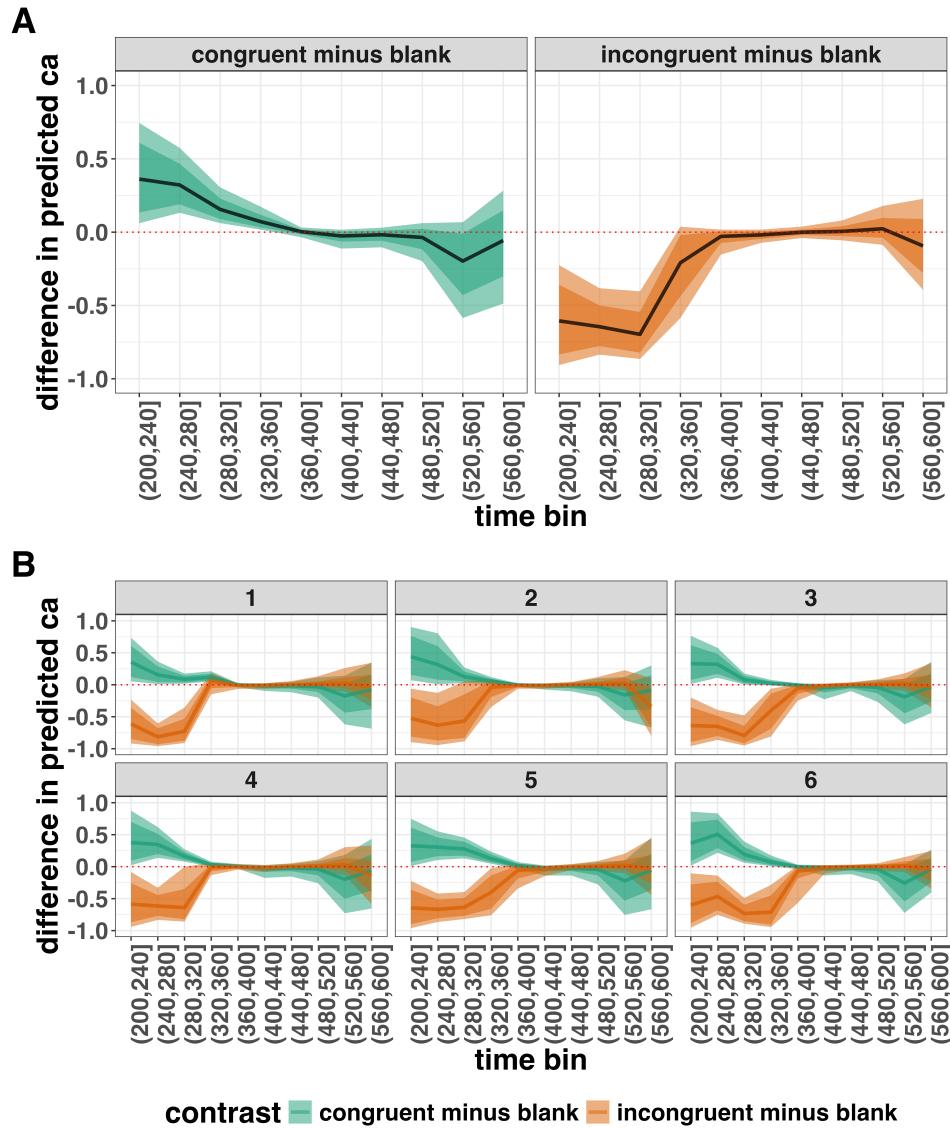


Figure 10. Point (mean) and 80/95% credible interval summaries of estimated differences in conditional accuracy in each time bin at the population level (A), and for each participant (B).

Based on Figure 10A we see that on the population level congruent primes have a

positive effect on the conditional accuracy of emitted responses in time bins (200,240],

(240,280], (280,320], and (320,360], relative to the estimates in the baseline condition

(blank prime; red dashed lines in Figure 10A). Incongruent primes have a negative effect on

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

606 **4.5 Tutorial 3a: Fitting Frequentist hazard models**

607 In this fifth tutorial we illustrate how to fit a multilevel regression model to RT data
 608 in the frequentist framework, for the data used in Tutorial 1a. The general process is
 609 similar to that in Tutorial 2a, except that there are no priors to set.

610 Again, to keep the tutorial concise, we only fit model M1i (see Tutorial 2a) using the
 611 function `glmer()` from the R package `lme4`. Alternatively, one could also use the function
 612 `glmmPQL()` from the R package `MASS` (Ripley et al., 2024). The resulting hazard model
 613 is called `M1i_f` with the appended “`_f`” denoting a frequentist model.

614 In Figure 11 we compare the parameter estimates from the Bayesian regression model
 615 `M1i` with those from the frequentist model `M1i_f`.

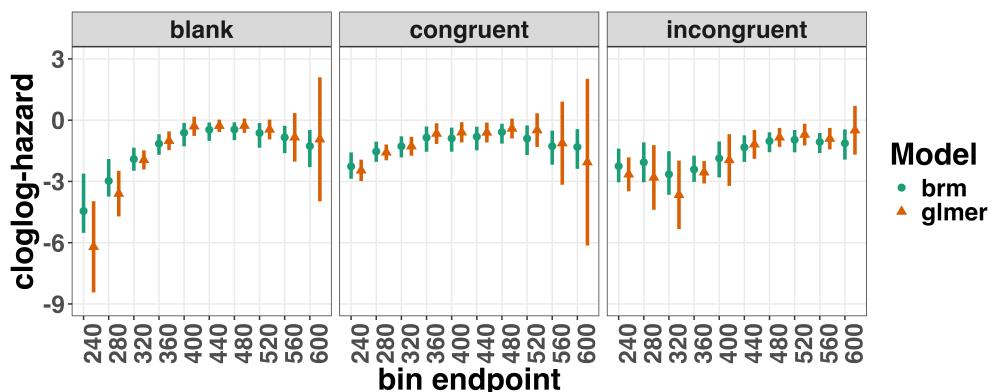


Figure 11. Parameter estimates for model M1i from `brm()` – means and 95% credible intervals – and model M1i_f from `glmer()` – maximum likelihood estimates and 95% confidence intervals.

616 Figure 11 confirms that the parameter estimates from both Bayesian and frequentist
 617 models are pretty similar, which makes sense given the close similarity in model structure.

618 However, model M1i_f did not converge and resulted in a singular fit. This is of course one
619 of the reasons why Bayesian modeling has become so popular in recent years. But the price
620 you pay for being able to fit models with more complex varying effects structures via a
621 Bayesian framework is increased computation time. In other words, as we have noted
622 throughout, some of the Bayesian models in Tutorials 2a took several hours to build.

623 **4.6 Tutorial 3b: Fitting Frequentist conditional accuracy models**

624 In this sixth tutorial we illustrate how to fit a multilevel regression model to the
625 timed accuracy data in the frequentist framework, for the data used in Tutorial 1a. To be
626 concise, we only fit effects from model M1i_ca (see Tutorial 2b) using the function glmer()
627 from the R package lme4. Alternatively, one could also use the function glmmPQL() from
628 the R package MASS (Ripley et al., 2024). The resulting conditional accuracy model,
629 which we labelled M1i_ca_f, did not converge and resulted in a singular fit. Again, this
630 just highlights some of the difficulties in fitting reasonably complex varying/random effects
631 structures in frequentist workflows.

632 **4.7 Tutorial 4: Planning**

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

634 **4.7.1 Background.** The general approach to planning that we adopt here involves
635 simulating reasonably structured data to help guide what you might be able to expect from
636 your data once you collect it (Gelman, Vehtari, et al., 2020). The basic structure and code
637 follows the examples outlined by Solomon Kurz in his ‘power’ blog posts
638 (<https://solomonkurz.netlify.app/blog/bayesian-power-analysis-part-i/>) and Lisa
639 DeBruine’s R package faux{} (<https://debruine.github.io/faux/>) as well as these related
640 papers (DeBruine & Barr, 2021; Pargent, Koch, Kleine, Lermer, & Gaube, 2024).

641 **4.7.2 Basic workflow.** The basic workflow is as follows:

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

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

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

655 **4.7.3 Fit a regression model and simulate one dataset.** We again use the
 656 data from Panis and Schmidt (2016) to provide a worked example. We fit an index coding
 657 model on a subset of time bins (six time bins in total) and for two prime conditions
 658 (congruent and incongruent). We chose to focus on a subsample of the data to ease the
 659 computational burden. We also used a full varying effects structure, with the model
 660 formula as follows:

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

661 We then took parameters from this model and used them to create a single dataset
 662 with 200 trials per condition for 10 individual participants. The raw data and the
 663 simulated data are plotted in Figure 12 and show quite close correspondence, which is
 664 re-assuring. But, this is only one dataset. What we really want to do is simulate many
 665 datasets and vary parameters of interest, which is what we turn to in the next section.

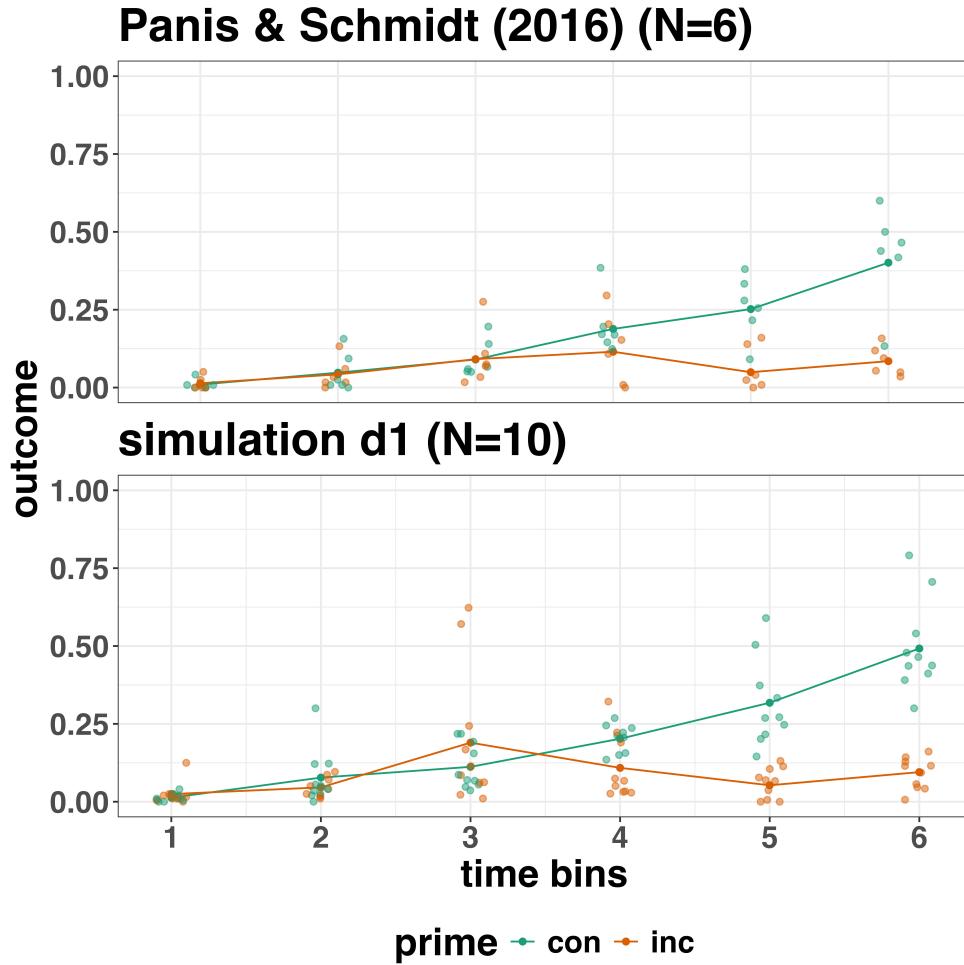


Figure 12. Raw data from Panis and Schmidt (2016) and simulated data from 10 participants.

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

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

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

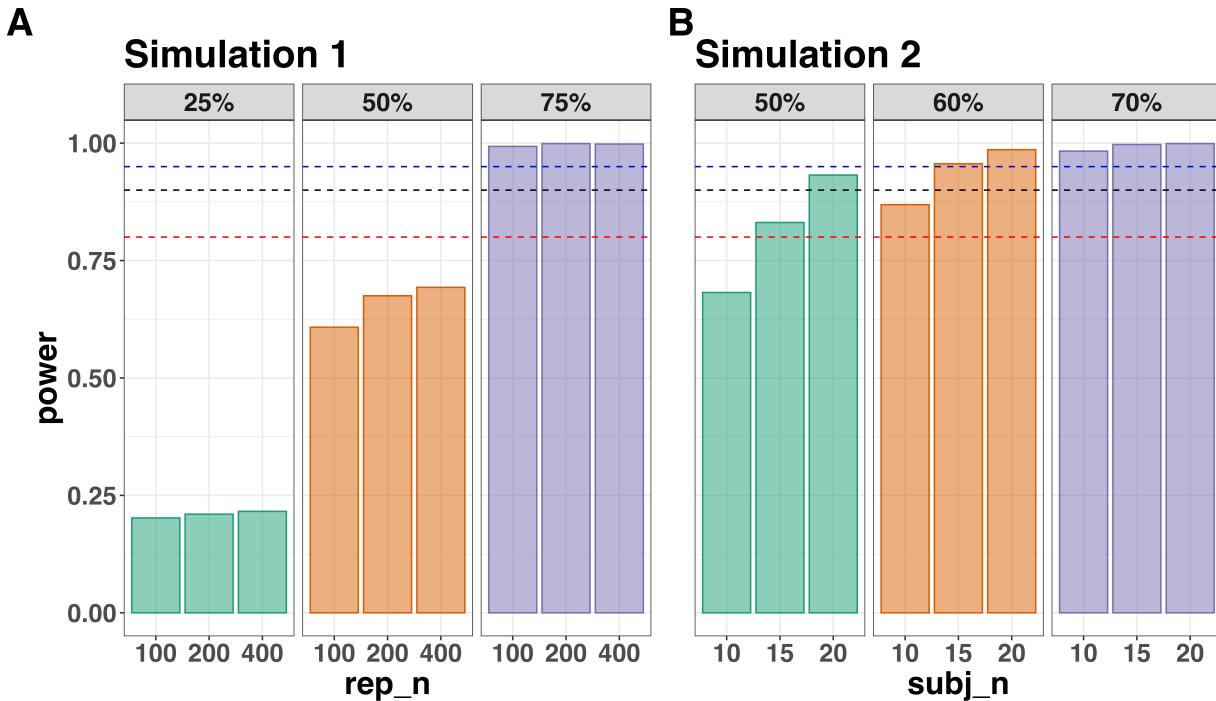


Figure 13. Statistical power across data Simulation 1 (A) and Simulation 2 (B). Power was calculated as the percentage of lower-bound 95% confidence intervals that exclude zero when the difference between prime condition is calculated (congruent - incongruent). In Simulation 1, the effect size was varied 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.

702 **4.7.5 Planning decisions.** Now that we have summarised our simulated data,

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

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

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

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

707 Some considerations might be as follows:

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

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

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

711 personnel?

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

728

5. Discussion

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

744 **5.1 What are the main use-cases of EHA for understanding cognition and brain
745 function?**

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

753 Second, even if you are not primarily interested in studying the temporal states of
754 cognition, EHA could still be a useful tool to consider using, in order to qualify inferences
755 that are being made based on mean-average comparisons. Given that distinctly different
756 inferences can be made from the same data based on whether one computes a
757 mean-average across trials or a RT distribution of events (Figure 1), it may be important
758 for researchers to supplement mean-average comparisons with EHA. One could envisage
759 scenarios where the implicit assumption of an effect manifesting across all of the time bins
760 measured would not be supported by EHA. Therefore, the conclusion of interest would not
761 apply to all responses, but instead it would be restricted to certain aspects of time.

762 5.2 Model complexity versus interpretability

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

779 4).

780 **5.3 Individual differences**

781 One important issue is that of possible individual differences in the overall location of
782 the distribution, and the time course of psychological effects. For example, when you wait
783 for a response of the participant on each trial, you allow the participant to have control
784 over the trial duration, and some participants might respond only when they are confident
785 that their emitted response will be correct. These issues can be avoided by introducing a
786 (relatively short) response deadline in each trial, e.g., 500 ms for simple detection tasks,
787 800 ms for more difficult discrimination tasks, or 2 s for tasks requiring extended high-level
788 processing. Because EHA can deal in a straightforward fashion with right-censored
789 observations (i.e., trials without an observed response in the analysis time window),
790 introducing a response deadline is recommended when designing RT experiments.
791 Furthermore, introducing a response deadline and asking participants to respond before the
792 deadline as much as possible, will also lead to individual distributions that overlap in time,
793 which is important when selecting a common analysis time window when fitting hazard
794 and conditional accuracy models.

795 But even when using a response deadline, participants can differ qualitatively in the
796 effects they display (see Panis, 2020). One way to deal with this is to describe and
797 interpret the different patterns. Another way is to run a clustering algorithm on the
798 individual hazard estimates across all bins and conditions. The obtained dendrogram can
799 then be used to identify a (hopefully big) cluster of participants that behave similarly, and
800 to identify a (hopefully small) cluster of participants with different behavioral patterns.
801 One might then exclude the smaller sub-group of participants before fitting a hazard model
802 or consider the possibility that different cognitive processes may be at play during task
803 performance across the different sub-groups.

804 Another approach to deal with individual differences is Bayesian prevalence (Ince,

805 Paton, Kay, & Schyns, 2021), which is a form of small- N approach (Smith & Little, 2018).

806 This method looks at effects within each individual in the study and asks how likely it

807 would be to see the same result if the experiment was repeated with a new person chosen

808 from the wider population at random. This approach allows one to quantify how typical or

809 uncommon an observed effect is in the population, and the uncertainty around this

810 estimate.

811 5.4 Limitations

812 Compared to the orthodox method – comparing mean-averages between conditions –

813 the most important limitation of multilevel hazard and conditional accuracy modeling is

814 that it might take a long time to estimate the parameters using Bayesian methods or the

815 model might have to be simplified significantly to use frequentist methods.

816 Another issue is that you need a relatively large number of trials per condition to

817 estimate the hazard function with high temporal resolution, which is required when testing

818 predictions of process models of cognition. Indeed, in general, there is a trade-off between

819 the number of trials per condition and the temporal resolution (i.e., bin width) of the

820 hazard function. Therefore, we recommend researchers to collect as many trials as possible

821 per experimental condition, given the available resources and considering the participant

822 experience (e.g., fatigue and boredom). For instance, if the maximum session length

823 deemed reasonable is between 1 and 2 hours, what is the maximum number of trials per

824 condition that you could reasonably collect? After consideration, it might be worth

825 conducting multiple testing sessions per participant and/or reducing the number of

826 experimental conditions. Finally, there is a user-friendly online tool for calculating

827 statistical power as a function of the number of trials as well as the number of participants,

828 and this might be worth consulting to guide the research design process (Baker et al., 2021).

We did not discuss continuous-time EHA, nor continuous-time SAT analysis. As indicated by Allison (2010), learning discrete-time EHA methods first will help in learning continuous-time methods. Given that RT is typically treated as a continuous variable, it is possible that continuous-time methods will ultimately prevail. However, they require much more data to estimate the continuous-time hazard (rate) function well. Thus, by trading a bit of temporal resolution for a lower number of trials, discrete-time methods seem ideal for dealing with typical psychological time-to-event data sets for which there are less than ~200 trials per condition per experiment.

5.5 Extensions

The hazard models in this tutorial assume that there is one event of interest. For RT data, this button-press event constitutes a single transition between an “idle” state and a “responded” state. However, in certain situations, more than one event of interest might exist. For example, in a medical or health-related context, an individual might transition back and forth between a “healthy” state and a “depressed” state, before being absorbed into a final “death” state. When you have data on the timing of these transitions, one can apply multi-state hazard models, which generalize EHA to transitions between three or more states (Steele, Goldstein, & Browne, 2004). Also, the predictor variables in this tutorial are time-invariant, i.e., their value did not change over the course of a trial. Thus, another extension is to include time-varying predictors, i.e., predictors whose value can change across the time bins within a trial (Allison, 2010). For example, when gaze position is tracked during a visual search trial, the gaze-target distance will vary during a trial when the eyes move around before a manual response is given; shorter gaze-target distances should be associated with a higher hazard of response occurrence. Note that the effect of a time-varying predictor (e.g., an occipital EEG signal) can itself vary over time.

853

6. Conclusions

854 Estimating the temporal distributions of RT and accuracy provide a rich source of
855 information on the time course of cognitive processing, which have been largely
856 undervalued in the history of experimental psychology and cognitive neuroscience. We hope
857 that by providing a set of hands-on, step-by-step tutorials, which come with custom-built
858 and freely available code, researchers will feel more comfortable embracing EHA and
859 investigating the temporal profile of cognitive states. On a broader level, we think that
860 wider adoption of such approaches will have a meaningful impact on the inferences drawn
861 from data, as well as the development of theories regarding the structure of cognition.

862

Author contributions

863 Conceptualization: S. Panis and R. Ramsey; Software: S. Panis and R. Ramsey;
864 Writing - Original Draft Preparation: S. Panis; Writing - Review & Editing: S. Panis and
865 R. Ramsey; Supervision: R. Ramsey.

866

Conflicts of Interest

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

869

Prior versions

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

872

Supplemental Material

873

Disclosures**874 Data, materials, and online resources**

875 Link to public archive:
876 https://github.com/sven-panis/Tutorial_Event_History_Analysis
877 Supplemental Material: Panis_Ramsey_suppl_material.pdf

878 Ethical approval

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

881

References

- 882 Allison, P. D. (1982). Discrete-Time Methods for the Analysis of Event Histories.
883 *Sociological Methodology*, 13, 61. <https://doi.org/10.2307/270718>
- 884 Allison, P. D. (2010). *Survival analysis using SAS: A practical guide* (2. ed). Cary, NC:
885 SAS Press.
- 886 Aust, F. (2019). *Citr: 'RStudio' add-in to insert markdown citations*. Retrieved from
887 <https://github.com/crsh/citr>
- 888 Aust, F., & Barth, M. (2024a). *papaja: Prepare reproducible APA journal articles with R*
889 *Markdown*. <https://doi.org/10.32614/CRAN.package.papaja>
- 890 Aust, F., & Barth, M. (2024b). *papaja: Prepare reproducible APA journal articles with R*
891 *Markdown*. <https://doi.org/10.32614/CRAN.package.papaja>
- 892 Baker, D. H., Vilidaite, G., Lygo, F. A., Smith, A. K., Flack, T. R., Gouws, A. D., &
893 Andrews, T. J. (2021). Power contours: Optimising sample size and precision in
894 experimental psychology and human neuroscience. *Psychological Methods*, 26(3),
895 295–314. <https://doi.org/10.1037/met0000337>
- 896 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for
897 confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*,
898 68(3), 10.1016/j.jml.2012.11.001. <https://doi.org/10.1016/j.jml.2012.11.001>
- 899 Barth, M. (2023). *tinylabes: Lightweight variable labels*. Retrieved from
900 <https://cran.r-project.org/package=tinylabes>
- 901 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects
902 models using lme4. *Journal of Statistical Software*, 67(1), 1–48.
903 <https://doi.org/10.18637/jss.v067.i01>
- 904 Bates, D., Maechler, M., & Jagan, M. (2024). *Matrix: Sparse and dense matrix classes and*
905 *methods*. Retrieved from <https://Matrix.R-forge.R-project.org>
- 906 Bengtsson, H. (2021). A unifying framework for parallel and distributed processing in r
907 using futures. *The R Journal*, 13(2), 208–227. <https://doi.org/10.32614/RJ-2021-048>

- 908 Blossfeld, H.-P., & Rohwer, G. (2002). *Techniques of event history modeling: New
909 approaches to causal analysis, 2nd ed* (pp. x, 310). Mahwah, NJ, US: Lawrence
910 Erlbaum Associates Publishers.
- 911 Box-Steffensmeier, J. M. (2004). Event history modeling: A guide for social scientists.
912 Cambridge: University Press.
- 913 Bürkner, P.-C. (2017). brms: An R package for Bayesian multilevel models using Stan.
914 *Journal of Statistical Software*, 80(1), 1–28. <https://doi.org/10.18637/jss.v080.i01>
- 915 Bürkner, P.-C. (2018). Advanced Bayesian multilevel modeling with the R package brms.
916 *The R Journal*, 10(1), 395–411. <https://doi.org/10.32614/RJ-2018-017>
- 917 Bürkner, P.-C. (2021). Bayesian item response modeling in R with brms and Stan. *Journal
918 of Statistical Software*, 100(5), 1–54. <https://doi.org/10.18637/jss.v100.i05>
- 919 DeBruine, L. M., & Barr, D. J. (2021). Understanding Mixed-Effects Models Through
920 Data Simulation. *Advances in Methods and Practices in Psychological Science*, 4(1),
921 2515245920965119. <https://doi.org/10.1177/2515245920965119>
- 922 Eddelbuettel, D., & Balamuta, J. J. (2018). Extending R with C++: A Brief Introduction
923 to Rcpp. *The American Statistician*, 72(1), 28–36.
924 <https://doi.org/10.1080/00031305.2017.1375990>
- 925 Eddelbuettel, D., & François, R. (2011). Rcpp: Seamless R and C++ integration. *Journal
926 of Statistical Software*, 40(8), 1–18. <https://doi.org/10.18637/jss.v040.i08>
- 927 Gabry, J., Češnovar, R., Johnson, A., & Broder, S. (2024). *Cmdstanr: R interface to
928 'CmdStan'*. Retrieved from <https://github.com/stan-dev/cmdstanr>
- 929 Gabry, J., Simpson, D., Vehtari, A., Betancourt, M., & Gelman, A. (2019). Visualization
930 in bayesian workflow. *J. R. Stat. Soc. A*, 182, 389–402.
931 <https://doi.org/10.1111/rssa.12378>
- 932 Gelman, A., Hill, J., & Vehtari, A. (2020). Regression and Other Stories.
933 <https://www.cambridge.org/highereducation/books/regression-and-other-stories/DD20DD6C9057118581076E54E40C372C>; Cambridge University Press.

- 935 https://doi.org/10.1017/9781139161879
- 936 Gelman, A., Vehtari, A., Simpson, D., Margossian, C. C., Carpenter, B., Yao, Y., ...
- 937 Modrák, M. (2020). *Bayesian Workflow*. arXiv.
- 938 https://doi.org/10.48550/arXiv.2011.01808
- 939 Girard, J. (2024). *Standist: What the package does (one line, title case)*. Retrieved from
- 940 https://github.com/jmgirard/standist
- 941 Grolemund, G., & Wickham, H. (2011). Dates and times made easy with lubridate.
- 942 *Journal of Statistical Software*, 40(3), 1–25. Retrieved from
- 943 https://www.jstatsoft.org/v40/i03/
- 944 Halley, E. (1693). VI. An estimate of the degrees of the mortality of mankind; drawn from
- 945 curious tables of the births and funerals at the city of breslaw; with an attempt to
- 946 ascertain the price of annuities upon lives. *Philosophical Transactions of the Royal*
- 947 *Society of London*, 17(196), 596–610. https://doi.org/10.1098/rstl.1693.0007
- 948 Heiss, A. (2021, November 10). A Guide to Correctly Calculating Posterior Predictions
- 949 and Average Marginal Effects with Multilevel Bayesian Models.
- 950 https://doi.org/10.59350/wbn93-edb02
- 951 Hosmer, D. W., Lemeshow, S., & May, S. (2011). *Applied Survival Analysis: Regression*
- 952 *Modeling of Time to Event Data* (2nd ed). Hoboken: John Wiley & Sons.
- 953 Ince, R. A., Paton, A. T., Kay, J. W., & Schyns, P. G. (2021). Bayesian inference of
- 954 population prevalence. *eLife*, 10, e62461. https://doi.org/10.7554/eLife.62461
- 955 Kantowitz, B. H., & Pachella, R. G. (2021). The Interpretation of Reaction Time in
- 956 Information-Processing Research 1. *Human Information Processing*, 41–82.
- 957 https://doi.org/10.4324/9781003176688-2
- 958 Kay, M. (2024). *tidybayes: Tidy data and geoms for Bayesian models*.
- 959 https://doi.org/10.5281/zenodo.1308151
- 960 Kruschke, J. K., & Liddell, T. M. (2018). The Bayesian New Statistics: Hypothesis testing,
- 961 estimation, meta-analysis, and power analysis from a Bayesian perspective.

- 962 *Psychonomic Bulletin & Review*, 25(1), 178–206.
- 963 <https://doi.org/10.3758/s13423-016-1221-4>
- 964 Kurz, A. S. (2023a). *Applied longitudinal data analysis in brms and the tidyverse* (version 0.0.3). Retrieved from <https://bookdown.org/content/4253/>
- 965 Kurz, A. S. (2023b). *Statistical rethinking with brms, ggplot2, and the tidyverse: Second edition* (version 0.4.0). Retrieved from <https://bookdown.org/content/4857/>
- 966 Lakens, D. (2022). Sample Size Justification. *Collabra: Psychology*, 8(1), 33267.
- 967 <https://doi.org/10.1525/collabra.33267>
- 968 Landes, J., Engelhardt, S. C., & Pelletier, F. (2020). An introduction to event history analyses for ecologists. *Ecosphere*, 11(10), e03238. <https://doi.org/10.1002/ecs2.3238>
- 969 Makeham, W. M. (1860). *On the Law of Mortality and the Construction of Annuity Tables*. The Assurance Magazine, and Journal of the Institute of Actuaries.
- 970 McElreath, R. (2020). *Statistical Rethinking: A Bayesian Course with Examples in R and STAN* (2nd ed.). New York: Chapman and Hall/CRC.
- 971 <https://doi.org/10.1201/9780429029608>
- 972 Meyer, D. E., Osman, A. M., Irwin, D. E., & Yantis, S. (1988). Modern mental chronometry. *Biological Psychology*, 26(1-3), 3–67.
- 973 [https://doi.org/10.1016/0301-0511\(88\)90013-0](https://doi.org/10.1016/0301-0511(88)90013-0)
- 974 Müller, K., & Wickham, H. (2023). *Tibble: Simple data frames*. Retrieved from <https://CRAN.R-project.org/package=tibble>
- 975 Nelder, J. A. (1999). From Statistics to Statistical Science. *Journal of the Royal Statistical Society. Series D (The Statistician)*, 48(2), 257–269. Retrieved from <https://www.jstor.org/stable/2681191>
- 976 Neuwirth, E. (2022). *RColorBrewer: ColorBrewer palettes*. Retrieved from <https://CRAN.R-project.org/package=RColorBrewer>
- 977 Panis, S. (2020). How can we learn what attention is? Response gating via multiple direct routes kept in check by inhibitory control processes. *Open Psychology*, 2(1), 238–279.
- 978
- 979
- 980
- 981
- 982
- 983
- 984
- 985
- 986
- 987
- 988

- 989 https://doi.org/10.1515/psych-2020-0107
- 990 Panis, S., Moran, R., Wolkersdorfer, M. P., & Schmidt, T. (2020). Studying the dynamics
991 of visual search behavior using RT hazard and micro-level speed–accuracy tradeoff
992 functions: A role for recurrent object recognition and cognitive control processes.
993 *Attention, Perception, & Psychophysics*, 82(2), 689–714.
- 994 https://doi.org/10.3758/s13414-019-01897-z
- 995 Panis, S., Schmidt, F., Wolkersdorfer, M. P., & Schmidt, T. (2020). Analyzing Response
996 Times and Other Types of Time-to-Event Data Using Event History Analysis: A Tool
997 for Mental Chronometry and Cognitive Psychophysiology. *I-Perception*, 11(6),
998 2041669520978673. https://doi.org/10.1177/2041669520978673
- 999 Panis, S., & Schmidt, T. (2016). What Is Shaping RT and Accuracy Distributions? Active
1000 and Selective Response Inhibition Causes the Negative Compatibility Effect. *Journal of*
1001 *Cognitive Neuroscience*, 28(11), 1651–1671. https://doi.org/10.1162/jocn_a_00998
- 1002 Panis, S., & Schmidt, T. (2022). When does “inhibition of return” occur in spatial cueing
1003 tasks? Temporally disentangling multiple cue-triggered effects using response history
1004 and conditional accuracy analyses. *Open Psychology*, 4(1), 84–114.
- 1005 https://doi.org/10.1515/psych-2022-0005
- 1006 Panis, S., Torfs, K., Gillebert, C. R., Wagemans, J., & Humphreys, G. W. (2017).
1007 Neuropsychological evidence for the temporal dynamics of category-specific naming.
1008 *Visual Cognition*, 25(1-3), 79–99. https://doi.org/10.1080/13506285.2017.1330790
- 1009 Panis, S., & Wagemans, J. (2009). Time-course contingencies in perceptual organization
1010 and identification of fragmented object outlines. *Journal of Experimental Psychology:*
1011 *Human Perception and Performance*, 35(3), 661–687.
- 1012 https://doi.org/10.1037/a0013547
- 1013 Pargent, F., Koch, T. K., Kleine, A.-K., Lermer, E., & Gaube, S. (2024). A Tutorial on
1014 Tailored Simulation-Based Sample-Size Planning for Experimental Designs With
1015 Generalized Linear Mixed Models. *Advances in Methods and Practices in Psychological*

- 1016 *Science*, 7(4), 25152459241287132. <https://doi.org/10.1177/25152459241287132>
- 1017 Pedersen, T. L. (2024). *Patchwork: The composer of plots*. Retrieved from
1018 <https://patchwork.data-imaginist.com>
- 1019 Pinheiro, J. C., & Bates, D. M. (2000). *Mixed-effects models in s and s-PLUS*. New York:
1020 Springer. <https://doi.org/10.1007/b98882>
- 1021 R Core Team. (2024). *R: A language and environment for statistical computing*. Vienna,
1022 Austria: R Foundation for Statistical Computing. Retrieved from
1023 <https://www.R-project.org/>
- 1024 Ripley, B., Venables, B., Bates, D. M., ca 1998), K. H. (partial. port, ca 1998), A. G.
1025 (partial. port, & polr), D. F. (support. functions for. (2024). *MASS: Support Functions
1026 and Datasets for Venables and Ripley's MASS*.
- 1027 Singer, J. D., & Willett, J. B. (2003). *Applied Longitudinal Data Analysis: Modeling
1028 Change and Event Occurrence*. Oxford, New York: Oxford University Press.
- 1029 Smith, P. L., & Little, D. R. (2018). Small is beautiful: In defense of the small-N design.
1030 *Psychonomic Bulletin & Review*, 25(6), 2083–2101.
1031 <https://doi.org/10.3758/s13423-018-1451-8>
- 1032 Stan Development Team. (2020). *StanHeaders: Headers for the R interface to Stan*.
1033 Retrieved from <https://mc-stan.org/>
- 1034 Stan Development Team. (2024). *RStan: The R interface to Stan*. Retrieved from
1035 <https://mc-stan.org/>
- 1036 Steele, F., Goldstein, H., & Browne, W. (2004). A general multilevel multistate competing
1037 risks model for event history data, with an application to a study of contraceptive use
1038 dynamics. *Statistical Modelling*, 4(2), 145–159.
1039 <https://doi.org/10.1191/1471082X04st069oa>
- 1040 Teachman, J. D. (1983). Analyzing social processes: Life tables and proportional hazards
1041 models. *Social Science Research*, 12(3), 263–301.
1042 [https://doi.org/10.1016/0049-089X\(83\)90015-7](https://doi.org/10.1016/0049-089X(83)90015-7)

- 1043 Tong, C. (2019). Statistical Inference Enables Bad Science; Statistical Thinking Enables
1044 Good Science. *The American Statistician*, 73(sup1), 246–261.
1045 <https://doi.org/10.1080/00031305.2018.1518264>
- 1046 Wickelgren, W. A. (1977). Speed-accuracy tradeoff and information processing dynamics.
1047 *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)
- 1048 Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New
1049 York. Retrieved from <https://ggplot2.tidyverse.org>
- 1050 Wickham, H. (2023a). *Forcats: Tools for working with categorical variables (factors)*.
1051 Retrieved from <https://forcats.tidyverse.org/>
- 1052 Wickham, H. (2023b). *Stringr: Simple, consistent wrappers for common string operations*.
1053 Retrieved from <https://stringr.tidyverse.org>
- 1054 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., ...
1055 Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43),
1056 1686. <https://doi.org/10.21105/joss.01686>
- 1057 Wickham, H., Çetinkaya-Rundel, M., & Grolemund, G. (2023). *R for data science: Import,
1058 tidy, transform, visualize, and model data* (2nd edition). Beijing Boston Farnham
1059 Sebastopol Tokyo: O'Reilly.
- 1060 Wickham, H., François, R., Henry, L., Müller, K., & Vaughan, D. (2023). *Dplyr: A
1061 grammar of data manipulation*. Retrieved from <https://dplyr.tidyverse.org>
- 1062 Wickham, H., & Henry, L. (2023). *Purrr: Functional programming tools*. Retrieved from
1063 [https://purrr.tidyverse.org/](https://purrr.tidyverse.org)
- 1064 Wickham, H., Hester, J., & Bryan, J. (2024). *Readr: Read rectangular text data*. Retrieved
1065 from <https://readr.tidyverse.org>
- 1066 Wickham, H., Vaughan, D., & Girlich, M. (2024). *Tidyr: Tidy messy data*. Retrieved from
1067 <https://tidyr.tidyverse.org>
- 1068 Winter, B. (2019). *Statistics for Linguists: An Introduction Using R*. New York:
1069 Routledge. <https://doi.org/10.4324/9781315165547>

- 1070 Wolkersdorfer, M. P., Panis, S., & Schmidt, T. (2020). Temporal dynamics of sequential
1071 motor activation in a dual-prime paradigm: Insights from conditional accuracy and
1072 hazard functions. *Attention, Perception, & Psychophysics*, 82(5), 2581–2602.
1073 <https://doi.org/10.3758/s13414-020-02010-5>