

1 **Why does preregistration increase the persuasiveness of evidence? A Bayesian**
2 **rationalization**

3 Aaron Peikert^{1,2,3}, Maximilian S. Ernst^{1, 4}, and & Andreas M. Brandmaier^{1, 3, 5}

4 ¹ Center for Lifespan Psychology

5 Max Planck Institute for Human Development

6 Berlin

7 Germany

8 ² Department of Imaging Neuroscience

9 University College London

10 London

11 UK

12 ³ Max Planck UCL Centre for Computational Psychiatry and Ageing Research

13 Berlin

14 Germany

15 ⁴ Max Planck School of Cognition

16 Leipzig

17 Germany

18 ⁵ Department of Psychology

19 MSB Medical School Berlin

20 Berlin

21 Germany

22 The materials for this article are available on [GitHub](#) (Peikert & Brandmaier, 2023a). This
23 version was created from git commit `9f4ccd7`. The manuscript is available as [preprint](#)
24 (Peikert & Brandmaier, 2023b).

25 Submitted to *Meta-Psychology*. Participate in open peer review by sending an email to
26 open.peer.reviewer@gmail.com. The full editorial process of all articles under review at
27 Meta-Psychology can be found following this link:

28 <https://tinyurl.com/mp-submissions>

29 You will find this preprint by searching for the first author's name.

30 Author Note

31
32 The authors made the following contributions. Aaron Peikert: Conceptualization,
33 Writing—Original Draft Preparation, Writing—Review & Editing, Methodology, Formal
34 analysis, Software, Visualization, Project administration; Maximilian S. Ernst:
35 Writing—Review & Editing, Formal analysis, Validation; Andreas M. Brandmaier:
36 Writing—Review & Editing, Supervisions.

37 Correspondence concerning this article should be addressed to Aaron Peikert,
38 Center for Lifespan Psychology, Max Planck Institute for Human Development, Lentzeallee
39 94, 14195 Berlin, Germany. E-mail: peikert@mpib-berlin.mpg.de

Abstract

40

41 The replication crisis has led many researchers to preregister their hypotheses and data
42 analysis plans before collecting data. A widely held view is that preregistration is supposed
43 to limit the extent to which data may influence the hypotheses to be tested. Only if data
44 have no influence an analysis is considered confirmatory. Consequently, many researchers
45 believe that preregistration is only applicable in confirmatory paradigms. In practice,
46 researchers may struggle to preregister their hypotheses because of vague theories that
47 necessitate data-dependent decisions (aka exploration). We argue that preregistration
48 benefits any study on the continuum between confirmatory and exploratory research. To
49 that end, we formalize a general objective of preregistration and demonstrate that
50 exploratory studies also benefit from preregistration. Drawing on Bayesian philosophy of
51 science, we argue that preregistration should primarily aim to reduce uncertainty about the
52 inferential procedure used to derive results. This approach provides a principled
53 justification of preregistration, separating the procedure from the goal of ensuring strictly
54 confirmatory research. We acknowledge that knowing the extent to which a study is
55 exploratory is central, but certainty about the inferential procedure is a prerequisite for
56 persuasive evidence. Finally, we discuss the implications of these insights for the practice of
57 preregistration.

58 *Keywords:* preregistration; confirmation; exploration; hypothesis testing; Bayesian;

59 Open Science

60 Word count: 8390

61 **Why does preregistration increase the persuasiveness of evidence? A Bayesian**
62 **rationalization**

63 The scientific community has long pondered the vital distinction between
64 exploration and confirmation, discovery and justification, hypothesis generation and
65 hypothesis testing, or prediction and postdiction (Hoyningen-Huene, 2006; Nosek et al.,
66 2018; Shmueli, 2010; Tukey, 1980). Despite the different names, it is fundamentally the
67 same dichotomy that is at stake here. There is a broad consensus that both approaches are
68 necessary for science to progress; exploration, to make new discoveries and confirmation, to
69 expose these discoveries to potential falsification, and assess empirical support for the
70 theory. However, mistaking exploratory findings for empirically confirmed results is
71 dangerous. It inflates the likelihood of believing that there is evidence supporting a given
72 hypothesis, even if it is false. A variety of problems, such as researchers' degrees of freedom
73 together with researchers' hindsight bias or naive p-hacking have led to such mistakes
74 becoming commonplace yet unnoticed for a long time. Recognizing them has led to a crisis
75 of confidence in the empirical sciences (Ioannidis, 2005), and psychology in particular
76 (Open Science Collaboration, 2015). As a response to the crisis, evermore researchers
77 preregister their hypotheses and their data collection and analysis plans in advance of their
78 studies (Nosek et al., 2018). They do so to stress the predictive nature of their registered
79 statistical analyses, often with the hopes of obtaining a label that marks the study as
80 "confirmatory". Indeed, rigorous application of preregistration prevents researchers from
81 reporting a set of results produced by an arduous process of trial and error as a simple
82 confirmatory story (Wagenmakers et al., 2012) while keeping low false-positive rates. This
83 promise of a clear distinction between confirmation and exploration has obvious appeal to
84 many who have already accepted the practice. Still, the majority of empirical researchers
85 do not routinely preregister their studies. One reason may be that some do not find that
86 the theoretical advantages outweigh the practical hurdles, such as specifying every aspect of
87 a theory and the corresponding analysis in advance. We believe that we can reach a greater

88 acceptance of preregistration by explicating a more general objective of preregistration that
89 benefits all kinds of studies, even those that allow data-dependent decisions.

90 One goal of preregistration that has received widespread attention is to clearly
91 distinguish confirmatory from exploratory research (Bakker et al., 2020; Mellor & Nosek,
92 2018; Nosek et al., 2018; Simmons et al., 2021; Wagenmakers et al., 2012). In such a
93 narrative, preregistration is justified by a confirmatory research agenda. However, two
94 problems become apparent under closer inspection. First, many researchers do not
95 subscribe to a purely confirmatory research agenda (Baumeister, 2016; Brandmaier et al.,
96 2013; Finkel et al., 2017; Tukey, 1972). Second, there is no strict mapping of the categories
97 preregistered vs. non-preregistered onto the categories confirmatory vs. exploratory
98 research.

99 Obviously, researchers can conduct confirmatory research without preregistration —
100 though it might be difficult to convince other researchers of the confirmatory nature of
101 their research, that is, that they were free of cognitive biases, made no data-dependent
102 decisions, and so forth. The opposite, that is, preregistered but not strictly confirmatory
103 studies, are also becoming more commonplace (Chan et al., 2004; Dwan et al., 2008; Silagy
104 et al., 2002).

105 This is the result of researchers applying one of two strategies to evade the
106 self-imposed restrictions of preregistrations: writing a loose preregistration to begin with
107 (Stefan & Schönbrodt, 2023) or deviating from the preregistration afterward (Lakens, 2024).
108 The latter is a frequent occurrence and, perhaps more worryingly, often remains
109 undisclosed (Akker et al., 2023; Claesen et al., 2021). Both strategies may be used for
110 sensible scientific reasons or with the self-serving intent of generating desirable results.
111 Thus, insisting on equating preregistration and confirmation has led to the criticism that,
112 all things considered, preregistration is actually harmful and neither sufficient nor
113 necessary for doing good science (Pham & Oh, 2021; Szollosi et al., 2020).

114 We argue that such criticism is not directed against preregistration itself but against
115 a justification through a confirmatory research agenda (Wagenmakers et al., 2012). When
116 researchers criticize preregistration as being too inflexible to fit their research question,
117 they often simply acknowledge that their research goals are not strictly confirmatory.
118 Forcing researchers into adopting a strictly confirmatory research agenda does not only
119 imply changing *how* they investigate a phenomenon but also *what* research questions they
120 pose. However reasonable such a move is, changing the core beliefs of a large community is
121 much harder than convincing them that a method is well justified. We, therefore, attempt
122 to disentangle the *methodological* goals of preregistration from the *ideological* goals of
123 confirmatory science. It might well be the case that psychology needs more confirmatory
124 studies to progress as a science. However, independently of such a goal, preregistration can
125 be useful for any kind of study on the continuum between strictly confirmatory and fully
126 exploratory.

127 To form such an objective for preregistration, we first introduce some tools of
128 Bayesian philosophy of science and map the exploration/confirmation distinction onto a
129 dimensional quantity we call “theoretical risk” (a term borrowed from Meehl, 1978, but
130 formalized as the probability of proving a hypothesis wrong if it does not hold).

131 We are interested in why preregistrations should change researchers’ evaluation of
132 evidence. Applying a Bayesian framework allows us to investigate our research question
133 most straightforwardly because it directly deals with what we ought to believe, given the
134 evidence presented. Specifically, it allows us to model changes in subjective degrees of
135 belief due to preregistration or, more simply, “persuasion”. Please note that our decision to
136 adopt a Bayesian philosophy of science does not make assumptions about the statistical
137 methods researchers use. In fact, this conceptualization is intentionally as minimal as
138 possible to be compatible with a wide range of philosophies of science and statistical
139 methods researchers might subscribe to. One feature of the Bayesian framework, is the

140 strong emphasis on subjective yet rational judgement. Therefore, we assume that
141 researchers will differ significantly in how they value evidence but that by making
142 assumptions about the general process, we can make general statements that apply to all
143 these subjective evaluations. However, we should note that Popperians would be appalled
144 that we are content with positive inductive inferences (but we regard “failing to disprove”
145 as too limited), and Neopopperians would flinch that we assign probabilities to beliefs (we
146 are fond of calculating things). While the latter move is not strictly necessary it allows us
147 to connect the more abstract considerations more closely with what researchers believe.

148 Now, we outline two possible perspectives on the utility of preregistration. The first
149 one corresponds to the traditional application of preregistration to research paradigms that
150 focus on confirmation by maximizing the theoretical risk or, equivalently, by limiting type-I
151 error (when dichotomous decisions about theories are an inferential goal). We argue that
152 this view on the utility of preregistration can be interpreted as maximizing theoretical risk,
153 which otherwise may be reduced by researchers’ degrees of freedom, p-hacking, and suchlike.
154 The second interpretation is our main contribution: We argue that contrary to the classic
155 view, the objective of preregistration is *not* the maximization of theoretical risk but rather
156 the minimization of uncertainty about the theoretical risk. This interpretation leads to a
157 broad applicability of preregistration to both exploratory and confirmatory studies.

158 To arrive at this interpretation, we rely on three arguments. The first is that
159 theoretical risk is vital for judging evidential support for theories. The second argument is
160 that the theoretical risk for a given study is generally uncertain. The third and last
161 argument is that this uncertainty is reduced by applying preregistration. We conclude that
162 because preregistration decreases uncertainty about the theoretical risk, which in turn
163 increases the amount of knowledge we gain from a particular study, preregistration is
164 potentially useful for any kind of study, no matter where it falls on the
165 exploratory-confirmatory continuum.

Persuasion and the Bayesian rationale

166

167 If researchers plan to conduct a study, they usually hope that it will change their
168 assessment of some theory's verisimilitude (Niiniluoto, 1998). Moreover, they hope to
169 convince other researchers can be persuaded to change their believe in a theory as well.
170 Beforehand, researchers cannot know what evidence a study will provide but still must form
171 an expectation in order to decide about the specifics of a planned study, including if they
172 should preregister it. If they can expect that preregistration helps them to persuade other
173 researchers to change their believe, it is only rational to employ preregistration. To make
174 our three arguments, we must assume three things about what an ideal estimation process
175 entails and how it relates to what studies (preregistered vs not preregistered) to conduct.

176

1. Researchers judge the evidence for or against a hypothesis rationally.

177

2. They expect other researchers to apply a similar rational process.

178

3. Researchers try to maximize the expected persuasiveness for *other* researchers.

179

180 The assumption of rationality can be connected to Bayesian reasoning and leads to
181 our adoption of the framework. Our rationale is as follows. Researchers who decide to
182 conduct a certain study are actually choosing a study to bet on. They have to “place the
183 bet” by conducting the study by investing resources and stand to gain evidence for or
184 against a theory with some probability. This conceptualization of choosing a study as a
185 betting problem allows us to apply a “Dutch book” argument (Christensen, 1991). This
186 argument states that any better must follow the axioms of probability to avoid being
187 “irrational,” i.e., accepting bets that lead to sure losses. Fully developing a Dutch book
188 argument for this problem requires careful consideration of what kind of studies to include
189 as possible bets, defining a conversion rate from the stakes to the reward, and modeling
190 what liberties researchers have in what studies to conduct. Without deliberating these
191 concepts further, we find it reasonable that researchers should not violate the axioms of
probability if they have some expectation about what they stand to gain with some

192 likelihood from conducting a study. The axioms of probability are sufficient to derive the
 193 Bayes formula, on which we will heavily rely for our further arguments. The argument is
 194 not sufficient, however, to warrant conceptualizing persuasiveness in terms of posterior
 195 probability; that remains a leap of faith. In fact, persuasiveness depends on how other
 196 researchers weigh evidence which differs between individuals.

197 However, the argument applies to any reward function that satisfies the “statistical
 198 relevancy condition” (Fetzer, 1974; Salmon, 1970), that is, evidence only increases believe
 199 for a theory if the evidence is more likely to be observed under the theory than under the
 200 alternative. In particular, “diagnosticity” (Fiedler, 2017; Oberauer & Lewandowsky, 2019),
 201 a concept highlighted in recent psychological literature, seems to adhere to the statistical
 202 relevancy condition.

203 **Theoretical risk**

204 Our first argument is that theoretical risk is crucial for judging the persuasiveness of
 205 evidence. Put simply, risky predictions create persuasive evidence if they turn out to be
 206 correct. This point is crucial because we attribute much of the appeal of a confirmatory
 207 research agenda to this notion.

208 Let us make some simplifying assumptions and define our notation. To keep the
 209 notation simple, we restrict ourselves to evidence of a binary nature (either it was observed
 210 or not). We denote the probability of a hypothesis before observing evidence as $P(H)$ and
 211 its complement as $P(\neg H) = 1 - P(H)$. The probability of observing evidence under some
 212 hypothesis is $P(E|H)$. We can calculate the probability of the hypothesis after observing
 213 the evidence with the help of the Bayes formula:

$$P(H|E) = \frac{P(H)P(E|H)}{P(E)} \quad (1)$$

214 The posterior probability $P(H|E)$ is of great relevance since it is often used directly

215 or indirectly as a measure of confirmation of a hypothesis. In the tradition of Carnap, in its
216 direct use, it is called *confirmation as firmness*; in its relation to the a priori probability
217 $P(H)$, it is called *increase in firmness* (Carnap, 1950, preface to the 1962 edition). We
218 concentrate on the posterior probability because of its simplicity but take it only as one
219 example of a possible measure. In reality, researchers surely differ in what function they
220 apply to judge evidence and it is often most fruitful to compare more than two competing
221 hypotheses. The goal is therefore to reason about the space of possible measures
222 researchers might apply. However, since any measure fulfilling the statistical relevancy
223 condition increases monotonically with an increase in posterior probability $P(H|E)$, we
224 might well take it to illustrate our reasoning.

225 In short, we want to increase posterior probability $P(H|E)$. Increases in posterior
226 probability $P(H|E)$ are associated with increases in persuasiveness, of which we want to
227 maximize the expectation. So how can we increase posterior probability? The Bayes
228 formula yields three components that influence confirmation, namely $P(H)$, $P(E|H)$ and
229 $P(E)$. The first option leads us to the unsurprising conclusion that higher a priori
230 probability $P(H)$ leads to higher posterior probability $P(H|E)$. If a hypothesis is more
231 probable to begin with, observing evidence in its favor will result in a hypothesis that is
232 more strongly confirmed, all else being equal. However, the prior probability of a
233 hypothesis is nothing our study design can change. The second option is equally
234 reasonable; that is, an increase in $P(E|H)$ leads to a higher posterior probability $P(H|E)$.
235 $P(E|H)$ is the probability of obtaining evidence for a hypothesis when it holds. We call
236 this probability of detecting evidence, given that the hypothesis holds “detectability.”
237 Consequently, researchers should ensure that their study design allows them to find
238 evidence for their hypothesis, in case it is true. When applied strictly within the bounds of
239 null hypothesis testing, detectability is equivalent to power (or the complement of type-II
240 error rate). However, while detectability is of great importance for study design, it is not
241 directly relevant to what a preregistration is communicating to other researchers. We later

242 discuss how issues of detectability must be considered in a preregistration. Thus, $P(E)$
 243 remains to be considered. Since $P(E)$ is the denominator, decreasing it can increase the
 244 posterior probability. In other words, high risk, high reward.

245 If we equate riskiness with a low probability of obtaining evidence (when the
 246 hypothesis is false), the Bayesian rationale perfectly aligns with the observation that risky
 247 predictions lead to persuasive evidence. This tension between high risk leading to high gain
 248 is central to our consideration of preregistration. A high-risk, high-gain strategy is bound
 249 to result in many losses that are eventually absorbed by the high gains. Sustaining many
 250 “failed” studies is not exactly aligned with the incentive structure under which many, if not
 251 most, researchers operate. Consequently, researchers are incentivized to appear to take
 252 more risks than they actually do, which misleads their readers to give their claims more
 253 credence than they deserve. It is at this juncture that the practice and mispractice of
 254 preregistration comes into play. We argue that the main function of preregistration is to
 255 enable proper judgment of the riskiness of a study.

256 To better understand how preregistrations can achieve that, let us take a closer look
 257 at the factors contributing to $P(E)$. Using the law of total probability, we can split $P(E)$
 258 into two terms:

$$P(E) = P(H)P(E|H) + P(\neg H)P(E|\neg H) \quad (2)$$

259 We have already noted that there is not much to be done about prior probability
 260 ($P(H)$, and hence its counter probability $P(\neg H)$), and that it is common sense to increase
 261 detectability $P(E|H)$. The real lever to pull is therefore $P(E|\neg H)$. This probability tells
 262 us how likely it is that we find evidence in favor of the theory when in fact, the theory is
 263 not true. Its counter probability $P(\neg E|\neg H) = 1 - P(E|\neg H)$ is what we call “theoretical
 264 risk”, because it is the risk a theory takes on in predicting the occurrence of particular

265 evidence in its favor. We borrow the term from Meehl (1978), though he has not assigned
266 it to the probability $P(\neg E|\neg H)$. Kukla (1990) argued that the core arguments in Meehl
267 (1990) can be reconstructed in a purely Bayesian framework. However, while he did not
268 mention $P(\neg E|\neg H)$ he suggested that Meehl (1978) used the term “very strange
269 coincidence” for a small $P(E|\neg H)$ which would imply, that $P(\neg E|\neg H)$ can be related to or
270 even equated to theoretical risk.

271 Let us note some interesting properties of theoretical risk $P(\neg E|\neg H)$. First,
272 increasing theoretical risk leads to higher posterior probability $P(H|E)$, our objective.
273 Second, if the theoretical risk is smaller than detectability $P(E|H)$ it follows that the
274 posterior probability must decrease when observing the evidence. If detectability exceeds
275 theoretical risk, the evidence is less likely under the theory than it is when the theory does
276 not hold (the inverse of statistical relevancy). Third, if the theoretical risk equals zero, then
277 posterior probability is at best equal to prior probability but only if detectability is perfect
278 ($P(H|E) = 1$). In other words, observing a sure fact does not lend credence to a hypothesis.

279 The last statement sounds like a truism but is directly related to Popper’s seminal
280 criterion of demarcation. He stated that if it is impossible to prove that a hypothesis is false
281 ($P(\neg E|\neg H) = 0$, theoretical risk is zero), it cannot be considered a scientific hypothesis
282 (Popper, 2002, p. 18). We note these relations to underline that the Bayesian rationale we
283 apply here is able to reconstruct many commonly held views on how “risky” predictions are
284 valued (but we of course differ from Popper on the central role of induction in science).

285 Both theoretical risk $P(\neg E|\neg H)$ and detectability $P(E|H)$ aggregate countless
286 influences; otherwise, they could not model the process of evidential support for theories.
287 To illustrate the concepts we have introduced here, consider the following example of a
288 single theory and three experiments that may test it. The experiments were created to
289 illustrate how they may differ in their theoretical risk and detectability. Suppose the
290 primary theory is about the cognitive phenomenon of “insight.” For the purpose of

291 illustration, we define it, with quite some hand-waving, as a cognitive abstraction that
292 allows agents to consistently solve a well-defined class of problems. We present the
293 hypothesis that the following problem belongs to such a class of insight problems:

294 Use five matches (IIIII) to form the number eight.

295 We propose three experiments that differ in theoretical risk and detectability. All
296 experiments take a sample of ten psychology students. We present the students with the
297 problem for a brief span of time. After that, the three experiments differ as follows:

- 298 1. The experimenter gives a hint that the problem is easy to solve when using Roman
299 numerals; if all students come up with the solution, she records it as evidence for the
300 hypothesis.
- 301 2. The experimenter shows the solution “VIII” and explains it; if all students come up
302 with the solution, she records it as evidence for the hypothesis.
- 303 3. The experimenter does nothing; if all students come up with the solution, she records
304 it as evidence for the hypothesis.

305 We argue that experiment 1 has high theoretical risk $P(\neg E_1|\neg H)$ and high
306 detectability $P(E_1|H)$. If “insight” has nothing to do with solving the problem ($\neg H$), then
307 presenting the insight that Roman numerals can be used should not lead to all students
308 solving the problem ($\neg E_1$); the experiment, therefore, has high theoretical risk
309 $P(\neg E_1|\neg H)$. Conversely, if insight is required to solve the problem (H), then it is likely to
310 help all students to solve the problem (E_1), the experiment, therefore, has high
311 detectability $P(E_1|H)$. The second experiment, on the other hand, has low theoretical risk
312 $P(\neg E_2|\neg H)$. Even if “insight” has nothing to do with solving the problem ($\neg H$), there are
313 other plausible reasons for observing the evidence (E_2), because the students could simply
314 copy the solution without having any insight. With regard to detectability, experiments 1
315 and 2 differ in no obvious way. Experiment 3, however, also has low detectability. It is

316 unlikely that all students will come up with the correct solution in a short time (E_3), even
317 if insight is required (H); experiment 3 therefore has low detectability $P(E_3|H)$. The
318 theoretical risk, however, is also low in absolute terms, but high compared to the
319 detectability (statistical relevancy condition is satisfied). In the unlikely event that all 10
320 students place their matches to form the Roman numeral VIII (E_3), it is probably due to
321 insight (H) and not by chance $P(\neg E_3|\neg H)$. Of course, in practice, we would allow the
322 evidence to be probabilistic, e.g., relax the requirement of “all students” to nine out of ten
323 students, more than eight, and so forth.

324 As mentioned earlier, we restrict ourselves to binary evidence, to keep the
325 mathematical notation as simple as possible. We discuss the relation between statistical
326 methods and theoretical risk in the [Statistical Methods](#) section.

327 Preregistration as a means to increase theoretical risk?

328 Having discussed that increasing the theoretical risk will increase the persuasiveness,
329 it is intuitive to task preregistration with maximizing theoretical risk, i.e., a confirmatory
330 research agenda. Indeed, limiting the type-I error rate is commonly stated as *the* central
331 goal of preregistration (Nosek et al., 2018; Oberauer, 2019; Rubin, 2020). We argue that
332 while such a conclusion is plausible, we must first consider at least two constraints that
333 place an upper bound on the theoretical risk.

334 First, the theory itself limits theoretical risk: Some theories simply do not make
335 risky predictions, and preregistration will not change that. Consider the case of a
336 researcher contemplating the relation between two sets of variables. Suppose each set is
337 separately well studied, and strong theories tell the researcher how the variables within the
338 set relate. However, our imaginary researcher now considers the relation between these two
339 sets. For lack of a better theory, they assume that some relation between any variables of
340 the two sets exists. This is not a risky prediction to make in psychology (Orben & Lakens,
341 2020). However, we would consider it a success if the researcher would use the evidence

368 risk into our framework.

369 **Statistical methods**

370 One widely known factor is the contribution of statistical methods to theoretical
371 risk. Theoretical risk $P(\neg E|\neg H)$ is deeply connected with statistical methods, because it is
372 related to the type-I error rate in statistical hypothesis testing $P(E|\neg H)$ by
373 $P(\neg E|\neg H) = 1 - P(E|\neg H)$, if you consider the overly simplistic case where the research
374 hypothesis is equal to the statistical alternative-hypothesis because then the null-hypothesis
375 is $\neg H$. Because many researchers are familiar with the type-I error rate, it can be helpful
376 to remember this connection to theoretical risk. Researchers who choose a smaller type-I
377 error rate can be more sure of their results, if significant, because the theoretical risk is
378 higher. However, this connection should not be overinterpreted for two reasons. First,
379 according to most interpretations of null hypothesis testing, the absence of a significant
380 result should not generally be interpreted as evidence against the hypothesis (Mayo, 2018,
381 p. 5.3). Second, the research hypothesis rarely equals the statistical alternative hypothesis
382 (most research hypothesis are more specific than “any value except zero”). In fact, it is
383 entirely possible to assume the null hypothesis as a research hypothesis, as is commonly
384 done in e.g., structural equation modelling, where the roles of detectability, theoretical risk
385 and type-I/II error rate switch. We argue that theoretical risk (and hence its complement,
386 $P(E|\neg H)$) also encompasses factors outside the statistical realm, most notably the study
387 design and broader analytical strategies. Type-I error rate is the property of a statistical
388 test under some assumptions, whereas theoretical risk is a researchers’ belief. One may
389 take such theoretical properties as a first starting point to form a substantive belief but
390 surely researchers ought to take other factors into consideration. For example, if a
391 researcher believes that there might be confounding variables at play for the relation
392 between two variables, this should decrease theoretical risk; after all they might find an
393 association purely on account of the confounders (Fiedler, 2017).

394 Statistical methods stand out among these factors because we have a large and
395 well-understood toolbox for assessing and controlling their contribution to theoretical risk.
396 Examples of our ability to exert this control are the choice of type-I error rate, adjustments
397 for multiple testing, the use of corrected fit measures (i.e., adjusted R^2), information
398 criteria, or cross-validation in machine learning. These tools help us account for biases in
399 statistical methods that influence theoretical risk (and hence, $P(E|\neg H)$).

400 The point is that the contribution of statistical methods to theoretical risk can be
401 formally assessed. For many statistical models it can be analytically computed under some
402 assumptions. For those models or assumptions where this is impossible, one can employ
403 Monte Carlo simulation to estimate the contribution to theoretical risk. The precision with
404 which statisticians can discuss contributions to theoretical risk has lured the community
405 concerned with research methods into ignoring other factors that are much more uncertain.
406 We cannot hope to resolve this uncertainty; but we have to be aware of its implications.
407 These are presented in the following.

408 **Sources of uncertainty**

409 As we have noted, it is possible to quantify how statistical models affect the
410 theoretical risk based on mathematical considerations and simulation. However, other
411 factors in the broader context of a study are much harder to quantify. If one chooses to
412 focus only on the contribution of statistical methods to theoretical risk, one is bound to
413 overestimate it. Take, for example, a t-test of mean differences in two samples. Under ideal
414 circumstances (assumption of independence, normality of residuals, equal variance), it
415 stays true to its type-I error rate. However, researchers may do many very reasonable
416 things in the broader context of the study that affect theoretical risk: They might exclude
417 outliers, choose to drop an item before computing a sum score, broaden their definition of
418 the population to be sampled, translate their questionnaires into a different language,
419 impute missing values, switch between different estimators of the pooled variance, or any

420 number of other things. All of these decisions carry a small risk that they will increase the
 421 likelihood of obtaining evidence despite the underlying research hypothesis being false.
 422 Even if the t-test itself perfectly maintains its type I error rate, these factors influence
 423 $P(E|\neg H)$. While, in theory, these factors may leave $P(E|\neg H)$ unaffected or even decrease
 424 it, we argue that this is not the case in practice. Whether researchers want to or not, they
 425 continuously process information about how the study is going, except under strict
 426 blinding. While one can hope that processing this information does not affect their
 427 decision-making either way, this cannot be ascertained. Therefore, we conclude that
 428 statistical properties only guarantee a lower bound for theoretical risk. The only thing we
 429 can conclude with some certainty is that theoretical risk is not higher than what the
 430 statistical model guarantees without knowledge about the other factors at play.

431 **The effects of uncertainty**

432 Before we ask how preregistration influences this uncertainty, we must consider the
 433 implications of being uncertain about the theoretical risk. Within the Bayesian framework,
 434 this is both straightforward and insightful. Let us assume a researcher is reading a study
 435 from another lab and tries to decide whether and how much the presented results confirm
 436 the hypothesis. As the researcher did not conduct the study (and the study is not
 437 preregistered), they can not be certain about the various factors influencing theoretical risk
 438 (researcher degrees of freedom). We therefore express this uncertainty about the theoretical
 439 risk as a probability distribution Q of $P(E|\neg H)$ (remember that $P(E|\neg H)$ is related to
 440 theoretical risk by $P(E|\neg H) = 1 - P(\neg E|\neg H)$, so it does not matter whether we consider
 441 the distribution of theoretical risk or $P(E|\neg H)$). To get the expected value of $P(H|E)$
 442 that follows from the researchers' uncertainty about the theoretical risk, we can compute
 443 the expectation using Bayes theorem:

$$\mathbb{E}_Q[P(H|E)] = \mathbb{E}_Q \left[\frac{P(H)P(E|H)}{P(H)P(E|H) + P(\neg H)P(E|\neg H)} \right] \quad (3)$$

444 Of course, the assigned probabilities and the distribution Q vary from study to
 445 study and researcher to researcher (and even the measure of confirmation), but we can
 446 illustrate the effect of uncertainty with an example. Assuming $P(E|H) = 0.8$ (relative of
 447 the typically strived for power of 80%). Let us further assume that the tested hypothesis is
 448 considered unlikely to be true by the research community before the study is conducted
 449 ($P(H) = 0.1$) and assign a uniform distribution for $P(E|\neg H) \sim U([1 - \tau, 1])$ where τ is set
 450 to $1 - \alpha$, reflecting our assumption that this term gives an upper bound for theoretical risk
 451 $P(\neg E|\neg H)$. We chose this uniform distribution as it is the maximum entropy distribution
 452 with support $[1 - \tau, 1]$ and hence conforms to our Bayesian framework (Giffin & Caticha,
 453 2007).

With this, we derive the expected value of $P(H|E)$ as

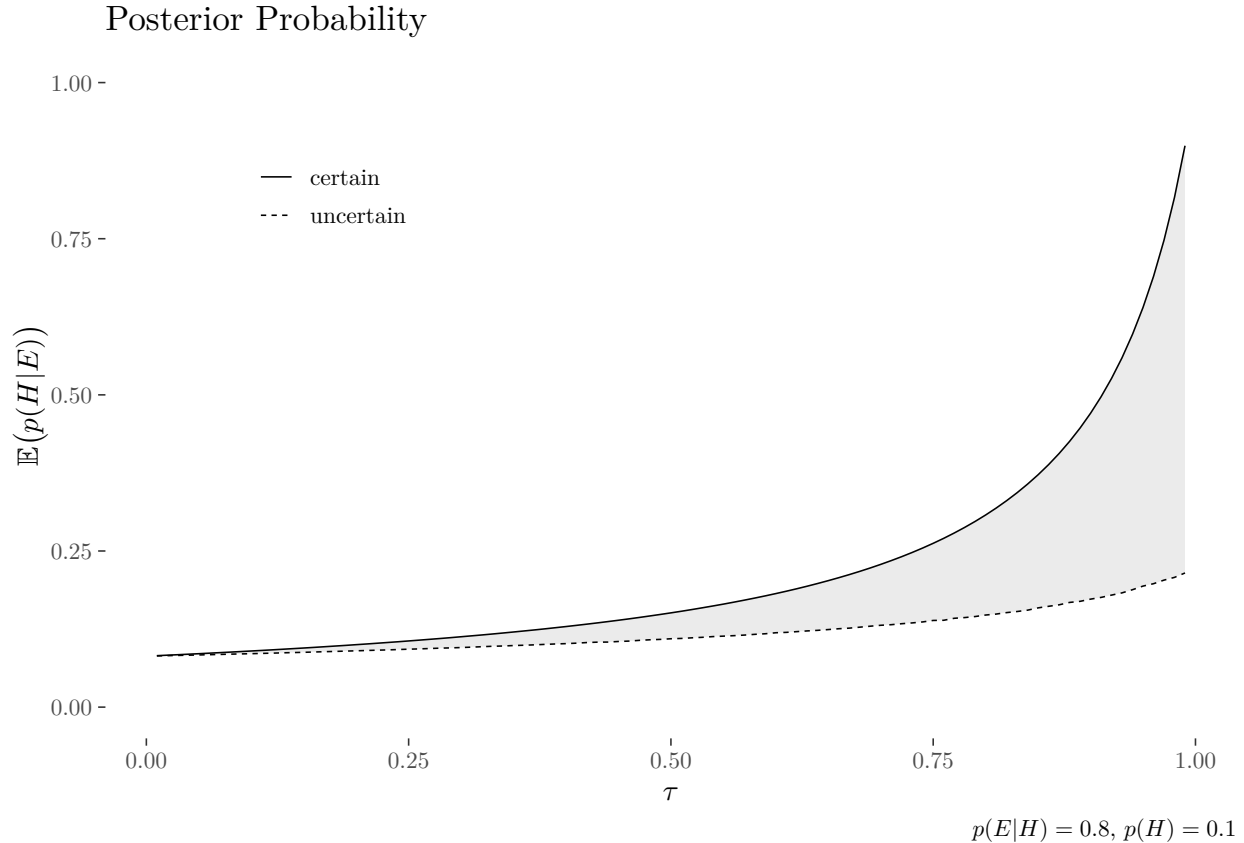
$$\mathbb{E}_Q[P(H|E)] = \mathbb{E}_Q \left[\frac{P(H)P(E|H)}{P(H)P(E|H) + P(\neg H)P(E|\neg H)} \right] \quad (4)$$

$$= \int_{[1-\tau, 1]} \tau^{-1} \frac{P(H)P(E|H)}{P(H)P(E|H) + P(\neg H)P(E|\neg H)} dP(E|\neg H) \quad (5)$$

$$= \frac{P(H)P(E|H)}{P(\neg H)\tau} \ln \left(\frac{P(H)P(E|H) + P(\neg H)}{P(H)P(E|H) + P(\neg H)(1 - \tau)} \right) \quad (6)$$

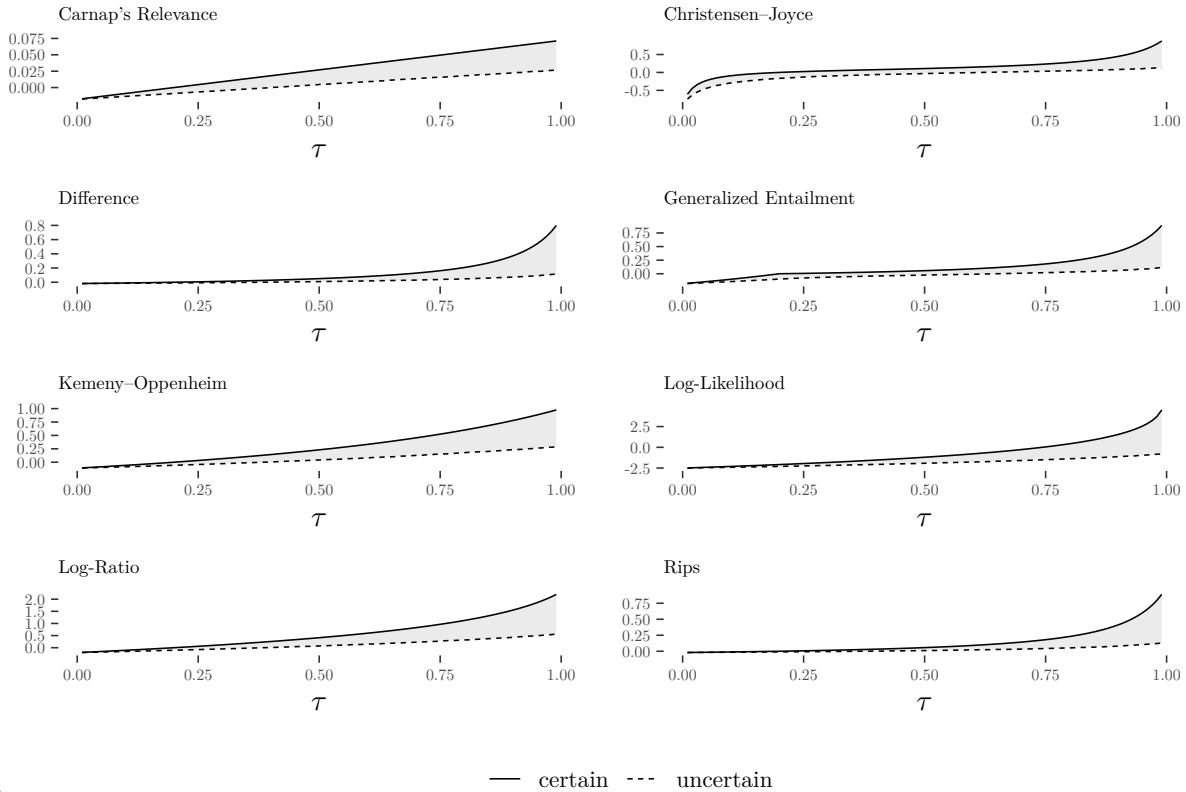
454 Figure 1 shows exemplary the effect of theoretical risk (x-axis) on the posterior
 455 probability (y-axis) being certain (solid line) or uncertain (dashed line) about the
 456 theoretical risk of a study. Our expectation of the persuasiveness varies considerably
 457 depending on how uncertain we are about the theoretical risk a study took on.
 458 Mathematically, uncertainty about theoretical risk is expressed through the variance (or
 459 rather entropy) of the distribution. The increase in uncertainty (expressed as more entropic
 460 distributions) leads to a decreased expected persuasiveness.

461 The argument for a confirmatory research agenda is that by increasing theoretical
 462 risk we increase expected persuasiveness, i.e., moving to the right on the x-axis in Figure 1

**Figure 1**

Posterior probability (confirmation as firmness) as a function of theoretical risk τ , where τ is either certain (solid line) or maximally uncertain (dotted line).

463 increases posterior probability (on the y-axis). However, if a hypothesis in a certain study
 464 has low theoretical risk, there is not much researchers can do about it. However, studies do
 465 not only differ by how high the theoretical risk is but also by how certain the recipient is
 466 about the theoretical risk. A study that has a very high theoretical risk (e.g., 1.00% chance
 467 that if the hypothesis is wrong, evidence in its favor will be observed,) but has also
 468 maximum uncertainty will result in a posterior probability of 21%, while the same study
 469 with maximum certainty will result in 90% posterior probability. The other factors
 470 (detectability, prior beliefs, measure of confirmation) and, therefore, the extent of the
 471 benefit varies, of course, with the specifics of the study. Crucially, even studies with some
 472 exploratory aspects benefit from preregistration, e.g., in this scenario with a $\tau = 0.80$ (false

**Figure 2**

Several measures for confirmation as an increase in firmness as a function of τ , where τ is either certain (solid line) or maximally uncertain (dotted line). Measures taken from Sprenger and Hartmann (2019), Table 1.3, p. 51.

473 positive rate of 0.20) moving from uncertain to certain increases the posterior from 0.15 to
 474 0.31. We find it helpful to calculate an example because of the nonlinear nature of the
 475 evidence functions.

476 Preregistration as a means to decrease uncertainty about the theoretical risk

477 We hope to have persuaded the reader to accept two arguments: First, the
 478 theoretical risk is important for judging evidential support for theories. Second, the
 479 theoretical risk is inherently uncertain, and the degree of uncertainty diminishes the
 480 persuasiveness of the gathered evidence. The third and last argument is that
 481 preregistrations reduce this uncertainty. Following the last argument, a preregistered study
 482 is represented by the solid line (certainty about theoretical risk), and a study that was not

483 preregistered is more similar to the dashed line (maximally uncertain about theoretical
484 risk) in Figure 1 and Figure 2.

485 Let us recall our three assumptions:

- 486 1. Researchers judge the evidence for or against a hypothesis rationally.
- 487 2. They expect other researchers to apply a similar rational process.
- 488 3. Researchers try to maximize the expected persuasiveness for other researchers.

489 The point we make with these assumptions is that researchers aim to persuade
490 other researchers, for example, the readers of their articles. Not only the original authors
491 are concerned with the process of weighing evidence for or against a theory but really the
492 whole scientific community the study authors hope to persuade. Unfortunately, readers of a
493 scientific article (or, more generally, any consumer of a research product) will likely lack
494 insight into the various factors that influence theoretical risk. While the authors
495 themselves may have a clear picture of what they did and how it might have influenced the
496 theoretical risk they took, their readers have much greater uncertainty about these factors.
497 In particular, they never know which relevant factors the authors of a given article failed to
498 disclose, be it intentionally or not. From the perspective of the ultimate skeptic, they may
499 claim maximum uncertainty.

500 Communicating clearly how authors of a scientific report collected their data and
501 consequently analyzed it to arrive at the evidence they present is crucial for judging the
502 theoretical risk they took. Preregistrations are ideal for communicating just that because
503 any description after the fact is prone to be incomplete. For instance, the authors could
504 have opted for selective reporting, that is, they decided to exclude a number of analytic
505 strategies they tried out. That is not to say that every study that was not-preregistered
506 was subjected to practices of questionable research practices. The point is that we cannot
507 exclude it with certainty. This uncertainty is drastically reduced if the researchers have

508 described what they intended to do beforehand and then report that they did exactly that.
509 In that case, readers can be certain they received a complete account of the situation.
510 They still might be uncertain about the actual theoretical risk the authors took, but to a
511 much smaller extent than if the study would not have been preregistered.

512 The remaining sources of uncertainty might be unfamiliarity with statistical
513 methods or experimental paradigms used, the probability of an implementation error in the
514 statistical analyses, a bug in the software used for analyses, etc. To further reduce the
515 uncertainty about theoretical risk, researchers must therefore publish code and ideally data.
516 After all, computational reproducibility is only possible if the data analytic procedure was
517 communicated clearly enough to allow others to retrace the computational steps (Peikert &
518 Brandmaier, 2021).

519 In any case, a well-written preregistration should aim to reduce the uncertainty
520 about the theoretical risk and hence increase the persuasiveness of evidence. Therefore, a
521 study that perfectly adhered to its preregistration will resemble the solid line in Figure 1/2.
522 Crucially, perfect means here that the theoretical risk can be judged with low uncertainty,
523 not that the theoretical risk is necessarily high.

524 **Hacking, harking, and other harms**

525 The importance of distinguishing between low and highly uncertain theoretical risk
526 becomes perhaps clearer if we consider a few hypothetical cases for illustration.

- 527 1. We know with absolute certainty that researchers will revert to p-hacking to create
528 evidence that is favorable for the theory.
- 529 2. A hypothesis was picked to explain reported results after the fact (HARKing, Kerr,
530 1998).
- 531 3. We cannot exclude the possibility of p-hacking having led to the reported results.
- 532 4. Reported results were obtained by planned exploration.

533 5. Reported results were obtained by unplanned exploration.

534 In case 1, there is no theoretical risk ($P(\neg E|\neg H) = 0$). If we know that the results
535 will be engineered to support the hypothesis no matter what, there is no reason to collect
536 data. A prime example of this case is the p_{ointless} metric (Hussey, 2021). Case 2 has a
537 similar problem. After all, the hypothesis that it had to happen the way it did happen is
538 irrefutable. In fact, both cases should be problematic to anyone who subscribes to the
539 statistical relevancy condition because if we choose the hypothesis in accordance with the
540 data or vice versa, without restrictions, they are not related anymore (i.e., observing the
541 data does not tell us anything about the hypothesis and the other way around). Case 3 is
542 different since here the theoretical risk is not necessarily low but simply uncertain (and
543 perhaps best represented by the dotted line in Figure 1/2). In case 4, the theoretical risk is
544 neither zero (unless the researcher plans to do run variations of analyses until a favourable
545 outcome is obtained, then we have a particular instance case of 1) nor high (as this is the
546 nature of exploratory approaches). However, we can take advantage of computational
547 reproducibility, use statistical properties, simulation or resampling methods, together with
548 scientific reasoning, to get a reasonably certain evaluation of the theoretical risk. Low
549 uncertainty about high theoretical risk is a somewhat favourable position (i.e., close to the
550 solid line in Figure 1/2). This favorable position leads us to recommend preregistration of
551 exploratory studies. Case 5 shares the neither zero nor high theoretical risk of case 4 but
552 has additional uncertainty about how much exploration was going on (how hard exactly
553 did the researchers try to come up with favourable results). Its low *and uncertain*
554 theoretical risk make it difficult to produce compelling evidence.

555

Discussion

556 To summarize, we showed that both higher theoretical risk and lower uncertainty
557 about theoretical risk lead to higher persuasiveness across a variety of measures. The
558 former result that increasing theoretical risk leads to higher expected persuasiveness

559 reconstructs the appeal and central goal of preregistration of confirmatory research
560 agendas. However, theoretical risk is something researchers have only limited control over.
561 For example, theories are often vague and ill-defined, resources are limited, and increasing
562 theoretical risk usually decreases detectability of a hypothesized effect (a special instance of
563 this trade-off is the well-known tension between type-I error and statistical power). While
564 we believe that preregistration is always beneficial, it might be counterproductive to pursue
565 high theoretical risk if the research context is inappropriate for strictly confirmatory
566 research. Specifically, appropriateness here entails the development of precise theories and
567 the availability of necessary resources (often, large enough sample size, but also see
568 Brandmaier et al. (2015)) to adequately balance detectability against theoretical risk.

569 In terms of preparing the conditions for confirmatory research, preregistration may
570 at most help to invest some time into developing more specific, hence riskier, implications
571 of a theory. But for a confirmatory science, it will not be enough to preregister all studies.
572 This undertaking requires action from the whole research community (Lishner, 2015).
573 Incentive structures must be created to evaluate not the outcomes of a study but the rigor
574 with which it was conducted (Cagan, 2013; Schönbrodt et al., 2022). Journal editors could
575 encourage theoretical developments that allow for precise predictions that will be tested by
576 other researchers and be willing to accept registered reports (Fried, 2020a, 2020b; van
577 Rooij & Baggio, 2021, 2020). Funding agencies should demand an explicit statement about
578 theoretical risk in relation to detectability and must be willing to provide the necessary
579 resources to reach adequate levels of both (Koole & Lakens, 2012).

580 Theoretical risk may conceptually be related to the framework of “severity” (Mayo,
581 2018; Mayo & Spanos, 2011). Severity, is a Neopopperian view which asserts that there is
582 evidence for a hypothesis just to the extent that it survives stringent scrutiny. However,
583 there are crucial differences between the two. First, our perspective on theoretical risk is
584 not primarily concerned with avoiding inductive reasoning but with subjective changes of

585 belief. This is important because, while severity is calculable, it remains unclear how
586 severity should be valued, e.g. if an increase in severity from .80 to .81 should be as
587 impressive as from .99 to .999. Second, severity considerations are mainly after the fact.
588 Severity, a measure with which we can rule out alternative explanations, can only be
589 calculated after evidence was observed. This makes it difficult to guide a priori decisions in
590 planning a study, after all severity disregards power, if we observe evidence, and disregards
591 Type I error rate when we do not. This implies that for a priori balancing Type I and Type
592 II error rate, a researcher must assign a priori probabilities to, for example, the size of an
593 effect. Since such a move is not in line with frequentist rationale there is no guidelines
594 available on how to do this. Third, we would argue that severity considerations assume full
595 information about how the evidence came about and hence imply axiomatically the need
596 for perfect preregistration. This comes down to frequentist understanding of probability as
597 the outcome of a well defined random experiments. When judging a particular study, a
598 frequentist, and hence a severe tester, may not assign probability to the event that the
599 researchers did, for example, p-hack. The lack of knowledge on the readers side does not
600 turn the p hacking into a random event of which we can calculate the long run frequency
601 aka frequentist probability. A severe test, hence, must assume that they know the Type I
602 and Type II error rate precisely. Full transparency, is hence assumed, and we can not
603 imagine many ways except preregistration that get close to this ideal. This assumptions
604 also makes it difficult to deal with less than perfect preregistrations and post-hoc changes
605 without appealing to principles outside the core philosophy of severity. One such approach
606 is Lakens (2024)' introduction of validity as an additional consideration to severity when
607 evaluating deviations from preregistrations. Interestingly, in this work he unconventionally
608 defines high severity as high $P(E|H)$ and high $P(\neg E|\neg H)$, which is closer to definitions of
609 "diagnosticity" (Fiedler, 2017; Oberauer & Lewandowsky, 2019) and falls under the broad
610 class of measures for evaluating evidence we consider here. Notably, the original definition
611 of severity does not satisfy the statistical relevancy condition and is not such a measure;

612 Mayo (2018), p. 14:

613 Severity Principle (strong): We have evidence for a claim C just to the extent it
614 survives a stringent scrutiny. If C passes a test that was highly capable of
615 findings flaws or discrepancies from C, and yet none or few are found, the
616 passing result, x, is evidence for C.

617 However, there also are communalities between our approach and severity, like the
618 strong emphasis on counterfactual consideration (imagining the hypothesis was false), and
619 there are even proposals to reconcile Bayesian and severity considerations (van Dongen et
620 al., 2023).

621 Our latter result, on the importance of preregistration for minimizing uncertainty,
622 has two important implications. The first is, that even if all imaginable actions regarding
623 promoting higher theoretical risk are taken, confirmatory research should be preregistered.
624 Otherwise, the uncertainty about the theoretical risk will diminish the advantage of
625 confirmatory research. Second, even under less-than-ideal circumstances for confirmatory
626 research, preregistration is beneficial. Preregistering exploratory studies increases the
627 expected persuasiveness by virtue of reducing uncertainty about theoretical risk.
628 Nevertheless, exploratory studies will have a lower expected persuasiveness than a more
629 confirmatory study if both are preregistered and have equal detectability.

630 Focusing on uncertainty reduction also explains two common practices of
631 preregistration that do not align with a confirmatory research agenda. First, researchers
632 seldomly predict precise numerical outcomes, instead they use preregistrations to describe
633 the process that generates the results. Precise predictions would have very high theoretical
634 risk (they are likely incorrect if the theory is wrong). A statistical procedure may have high
635 or low theoretical risk depending on the specifics of the model used. Specifying the process,
636 therefore, is in line with the rationale we propose here, but is less reasonable when the goal

637 of preregistration is supposed to be a strictly confirmatory research agenda.

638 Second, researchers often have to deviate from the preregistration and make
639 data-dependent decisions after the preregistration. If the only goal of preregistration is to
640 ensure confirmatory research, such changes are not justifiable. However, under our rational,
641 some changes may be justified. Any change increases the uncertainty about the theoretical
642 risk and may even decrease the theoretical risk. The changes still may be worthwhile if the
643 negative outcomes may be offset by an increase in detectability due to the change.
644 Consider a preregistration that failed to specify how to handle missing values, and
645 researchers subsequently encountering missing values. In such case, detectability becomes
646 zero because the data cannot be analyzed without a post-hoc decision about how to handle
647 the missing data. Any such decision would constitute a deviation from the preregistration,
648 which is possible under our proposed objective. Note that a reader cannot rule out that the
649 researchers leveraged the decision to decrease theoretical risk, i.e., picking among all
650 options the one that delivers the most beneficial results for the theory (in the previous
651 example, choosing between various options of handling missing values). Whatever decision
652 they make, increased uncertainty about the theoretical risk is inevitable and the expected
653 persuasiveness is decreased compared to a world where they anticipated the need to deal
654 with missing data. However, it is still justified to deviate. After all they have not
655 anticipated the case and are left with a detectability of zero. Any decision will increase
656 detectability to a non-zero value offsetting the increase in uncertainty. The researchers also
657 may do their best to argue that the deviation was not motivated by increasing theoretical
658 risk, thereby, decreasing the uncertainty. Ideally, there is a default decision that fits well
659 with the theory or with the study design. Or, if there is no obvious candidate, the
660 researchers could conduct a multiverse analysis of the available options to deal with
661 missings to show the influence of the decision (Steen et al., 2016). In any case, deviations
662 must be transparently reported and we applaud recent developments to standardize and
663 normalize this process (Willroth & Atherton, 2023).

664 As explained above, reduction in uncertainty as the objective for preregistration
665 does not only explain some existing practice, that does not align with confirmation as a
666 goal, it also allows to form recommendations to improve the practice of preregistration.
667 Importantly, we now have a theoretical measure to gauge the functionality of
668 preregistrations, which can only help increase its utility. In particular, a preregistration
669 should be specific about the procedure that is intended to generate evidence for a theory.
670 Such a procedure may accommodate a wide range of possible data, i.e., it may be
671 exploratory. The theoretical risk, however low, must be communicated clearly. Parts of the
672 process left unspecified imply uncertainty, which preregistration should reduce. However,
673 specifying procedures that can be expected to fail will lead to deviation and, subsequently,
674 to larger uncertainty.

675 Our emphasis on transparency aligns with other justifications of preregistration,
676 especially those put forth by Lakens (2019)'s, although based on quite different
677 philosophical foundations. Our goal is to contribute a rationale that more comprehensively
678 captures the spectrum of exploration and confirmation in relation to preregistrations,
679 post-hoc changes of preregistrations, and subjective evaluations of evidence. We find it
680 difficult to content ourselves with vague terms like “control” or “transparency” if they
681 ultimately remain unconnected to how much researchers believe in a theory. Within our
682 framework, researchers have the ability to input their assumptions regarding the
683 perspectives of other researchers and calculate the potential impact of their actions on their
684 readership, whether these actions relate to study design, to the preregistration itself, or
685 subsequent deviations from it. We put subjective evaluations at the center of our
686 considerations; we deal explicitly with researchers who are proponents of some theory (they
687 have higher priors for the theory being true), researchers who suspect confounding variables
688 (they assume lower theoretical risk), or those who remain doubtful if everything relevant
689 was reported (they have higher uncertainty about theoretical risk) or even those who place
690 greater value on incongruent evidence than others (they differ in their confirmation

691 function). We, therefore, hope to not only provide a rationale for preregistration for those
692 who subscribe to a Bayesian philosophy of science but also a framework to navigate the
693 complicated questions that arise in the practice of preregistration.

694 At the same time, approaching the evaluation of evidence using a Bayesian
695 formalism is far from novel (Fiedler, 2017; e.g., Kukla, 1990; Sprenger & Hartmann, 2019).
696 To our knowledge, it was not yet applied to the problem of preregistration. However,
697 Oberauer and Lewandowsky (2019) made use of the formalism to model the relation
698 between theory, hypothesis, and evidence. In the context of this conceptualization, they
699 discussed the usefulness of preregistration, though without applying the formalism there.
700 Most importantly, they are rather critical of the idea that preregistration has tangible
701 benefits. Instead, they prefer multiverse analyses but contend that those could be
702 preregistered if one fancies it. Their reasoning is based on two intuitions about what
703 should *not* influence the evaluation of evidence: temporal order and the mental state of the
704 originator. In our opinion, they disregard the temporal order a bit too hastily, as it is a
705 long-standing issue in Bayesian philosophy of science known as the “problem of old
706 evidence” (Chihara, 1987). However, we agree that not the temporal order is decisive but if
707 the researchers incorporated the information into the hypothesis the evidence is supposed
708 to confirm. For the other, we argue that the mental state of the originator does matter.
709 Suppose there are $k = 1, 2, \dots, K$ ways to analyze data, where each k has a $P(E_k|\neg H) > 0$.
710 If they intend to try each way after another but happen to be “lucky” on the first try and
711 stop, should we then apply $P(E|\neg H) = P(E_1|\neg H)$ or $P(E|\neg H) = P(E_1 \vee \dots \vee E_k|\neg H)$?
712 We think the latter. However, this “Defeatist” intuition is not universally warranted and
713 depends on what we take H to mean specifically (Kotzen, 2013). Addressing, this problem
714 might benefit from combining Oberauer and Lewandowsky (2019)’s idea of updating on
715 two nested levels (theory-hypothesis layered on top of hypothesis-evidence) with our
716 approach to modelling uncertainty.

717 Whatever the difference in evaluating preregistration as a tool, maybe conceptually
718 more profound is that Oberauer and Lewandowsky (2019) conceptualizes
719 “discovery-oriented research” differently than we do “exploratory”. They assume the same
720 theoretical risk ($P(\neg E|\neg H) = .05$) and detectability ($P(E|H) = .8$) in their calculation
721 example as we do but assign different prior probabilities, namely .06 for discovery versus .6
722 for theory testing. Then, they conclude that discovery-oriented researcher requires a much
723 lower type-I error rate to control false positive in light of the low prior probability. This
724 runs counter to our definition of exploratory research having low theoretical risk. Of course,
725 we agree that low priors require more persuasive evidence; our disagreement, therefore, lies
726 mainly in terminology. They imagine discovery-oriented researchers to conduct
727 experiments where they have low expectations that they obtain positive evidence
728 ($.06 \cdot .8 + .94 \cdot .05 = 0.095$), but if they do, it raises the posterior significantly (from .06 to
729 .51) In our view, researchers who set out to explore a data set often find “something” (due
730 to low $P(\neg E|\neg H)$); therefore, it should only slightly raise your posterior if they do. On a
731 substantive matter, we believe both kinds of research are common in psychology. It is,
732 therefore, mostly a disagreement on terminology. This disagreement only highlights why
733 using a mathematical framework to investigate such things is so useful and ultimately
734 indispensable because we can clearly see where and how we differ in our reasoning.

735 We believe that our reasoning is quite similar to Höfler et al. (2022), who call for
736 transparent exploration using preregistration. We could be more sure of our agreement, if
737 they had formulated their arguments within a mathematical framework, which would also
738 have helped to dissolve an apparent conflict in their definitions of confirmation, exploration,
739 and transparency. On the one hand, they define “The principle difference between
740 confirmation and exploration is that confirmation adheres to an evidential norm for the
741 test of a hypothesis to pass.”, but then suggest that transparent exploration can be
742 conducted using inferences tests as a filtering mechanism. Their distinction between
743 confirmation, intransparent and transparent exploration are otherwise just as well placed

744 along the dimensions, theoretical risk and uncertainty about theoretical risk.

745 With the goal to facilitate rigorous exploration, we have proposed a workflow for
746 preregistration called *preregistration as code* (PAC) elsewhere (Peikert et al., 2021). In a
747 PAC, researchers use computer code for the planned analysis as well as a verbal description
748 of theory and methods for the preregistration. This combination is facilitated by dynamic
749 document generation, where the results of the code, such as numbers, figures, and tables,
750 are inserted automatically into the document. The idea is that the preregistration already
751 contains “mock results” based on simulated or pilot data, which are replaced after the
752 actual study data becomes available. Such an approach dissolves the distinction between
753 the preregistration document and the final scientific report. Instead of separate documents,
754 preregistration, and final report are different versions of the same underlying dynamic
755 document. Deviations from the preregistration can therefore be clearly (and if necessary,
756 automatically) isolated, highlighted, and inspected using version control. Crucially, because
757 the preregistration contains code, it may accommodate many different data patterns, i.e., it
758 may be exploratory. However, while a PAC does not limit the extent of exploration, it is
759 very specific about the probability to generate evidence even when the theory does not
760 hold (theoretical risk). Please note that while PAC is ideally suited to reduce uncertainty
761 about theoretical risk, other more traditional forms of preregistration are also able to
762 advance this goal.

763 Contrary to what is widely assumed about preregistration, a preregistration is not
764 necessarily a seal of confirmatory research. Confirmatory research would almost always be
765 less persuasive without preregistration, but in our view, preregistration primarily
766 communicates the extent of confirmation, i.e., theoretical risk, of a study. Clearly
767 communicating theoretical risk is important because it reduces the uncertainty and hence
768 increases expected persuasiveness.

769

Acknowledgement

770

771

772

773

We thank Leo Richter, Caspar van Lissa, Felix Schönbrodt, the discussants at the DGPS2022 conference and Open Science Center Munich, and many more for the insightful discussions about disentangling preregistration and confirmation. We are grateful to Julia Delius for her helpful assistance in language and style editing.

774

Declarations

775

776

777

778

All code and materials required to reproduce this article are available under <https://github.com/aaronpeikert/bayes-prereg> (Peikert & Brandmaier, 2023a). The authors have no competing interests to declare that are relevant to the content of this article.

References

- 779 Akker, O. van den, Bakker, M., Assen, M. A. L. M. van, Pennington, C. R., Verweij, L.,
780 Elsherif, M., Claesen, A., Gaillard, S. D. M., Yeung, S. K., Frankenberger, J.-L.,
781 Krautter, K., Cockcroft, J. P., Kreuer, K. S., Evans, T. R., Heppel, F., Schoch, S. F.,
782 Korbmacher, M., Yamada, Y., Albayrak-Aydemir, N., ... Wicherts, J. (2023, May 10).
783 *The effectiveness of preregistration in psychology: Assessing preregistration strictness*
784 *and preregistration-study consistency*. <https://doi.org/10.31222/osf.io/h8xjw>
785
- 786 Bakker, M., Veldkamp, C. L. S., Assen, M. A. L. M. van, Cromptoets, E. A. V., Ong, H.
787 H., Nosek, B. A., Soderberg, C. K., Mellor, D., & Wicherts, J. M. (2020). Ensuring the
788 quality and specificity of preregistrations. *PLOS Biology*, *18*(12), e3000937.
789 <https://doi.org/10.1371/journal.pbio.3000937>
- 790 Baumeister, R. F. (2016). Charting the future of social psychology on stormy seas:
791 Winners, losers, and recommendations. *Journal of Experimental Social Psychology*, *66*,
792 153–158. <https://doi.org/10.1016/j.jesp.2016.02.003>
- 793 Brandmaier, A. M., Oertzen, T. von, Ghisletta, P., Hertzog, C., & Lindenberger, U. (2015).
794 LIFESPAN: A tool for the computer-aided design of longitudinal studies. *Frontiers in*
795 *Psychology*, *6*, 272.
- 796 Brandmaier, A. M., von Oertzen, T., McArdle, J. J., & Lindenberger, U. (2013). Structural
797 equation model trees. *Psychological Methods*, *18*(1), 71–86.
798 <https://doi.org/10.1037/a0030001>
- 799 Cagan, R. (2013). San Francisco Declaration on Research Assessment. *Disease Models &*
800 *Mechanisms*, dmm.012955. <https://doi.org/10.1242/dmm.012955>
- 801 Carnap, R. (1950). *Logical Foundations of Probability*. Chicago, IL, USA: Chicago
802 University of Chicago Press.
- 803 Chan, A.-W., Hróbjartsson, A., Haahr, M. T., Gøtzsche, P. C., & Altman, D. G. (2004).
804 Empirical Evidence for Selective Reporting of Outcomes in Randomized
805 TrialsComparison of Protocols to Published Articles. *JAMA*, *291*(20), 2457–2465.

806 <https://doi.org/10.1001/jama.291.20.2457>

807 Chihara, C. S. (1987). Some Problems for Bayesian Confirmation Theory. *The British*
808 *Journal for the Philosophy of Science*, 38(4), 551–560.

809 <https://doi.org/10.1093/bjps/38.4.551>

810 Christensen, D. (1991). Clever Bookies and Coherent Beliefs. *The Philosophical Review*,
811 100(2), 229–247. <https://doi.org/10.2307/2185301>

812 Claesen, A., Gomes, S., Tuerlinckx, F., & Vanpaemel, W. (2021). Comparing dream to
813 reality: An assessment of adherence of the first generation of preregistered studies.

814 *Royal Society Open Science*, 8(10), 211037. <https://doi.org/10.1098/rsos.211037>

815 Dwan, K., Altman, D. G., Arnaiz, J. A., Bloom, J., Chan, A.-W., Cronin, E., Decullier, E.,
816 Easterbrook, P. J., Elm, E. V., Gamble, C., Ghersi, D., Ioannidis, J. P. A., Simes, J., &
817 Williamson, P. R. (2008). Systematic Review of the Empirical Evidence of Study
818 Publication Bias and Outcome Reporting Bias. *PLOS ONE*, 3(8), e3081.

819 <https://doi.org/10.1371/journal.pone.0003081>

820 Fetzer, J. H. (1974). Statistical Explanations. In K. F. Schaffner & R. S. Cohen (Eds.),
821 *PSA 1972: Proceedings of the 1972 Biennial Meeting of the Philosophy of Science*
822 *Association* (pp. 337–347). Springer Netherlands.

823 https://doi.org/10.1007/978-94-010-2140-1_23

824 Fiedler, K. (2017). What Constitutes Strong Psychological Science? The (Neglected) Role
825 of Diagnosticity and A Priori Theorizing. *Perspectives on Psychological Science*, 12(1),
826 46–61. <https://doi.org/10.1177/1745691616654458>

827 Finkel, E. J., Eastwick, P. W., & Reis, H. T. (2017). Replicability and other features of a
828 high-quality science: Toward a balanced and empirical approach. *Journal of Personality*
829 *and Social Psychology*, 113(2), 244–253. <https://doi.org/10.1037/pspi0000075>

830 Fried, E. I. (2020a). Lack of Theory Building and Testing Impedes Progress in The Factor
831 and Network Literature. *Psychological Inquiry*, 31(4), 271–288.

832 <https://doi.org/10.1080/1047840X.2020.1853461>

- 833 Fried, E. I. (2020b). Theories and Models: What They Are, What They Are for, and What
834 They Are About. *Psychological Inquiry*, 31(4), 336–344.
835 <https://doi.org/10.1080/1047840X.2020.1854011>
- 836 Giffin, A., & Caticha, A. (2007). Updating Probabilities with Data and Moments. *AIP*
837 *Conference Proceedings*, 954, 74–84. <https://doi.org/10.1063/1.2821302>
- 838 Höfler, M., Scherbaum, S., Kanske, P., McDonald, B., & Miller, R. (2022). Means to
839 valuable exploration: I. The blending of confirmation and exploration and how to
840 resolve it. *Meta-Psychology*, 6. <https://doi.org/10.15626/MP.2021.2837>
- 841 Hoyningen-Huene, P. (2006). Context of Discovery Versus Context of Justification and
842 Thomas Kuhn. In J. Schickore & F. Steinle (Eds.), *Revisiting Discovery and*
843 *Justification: Historical and philosophical perspectives on the context distinction* (pp.
844 119–131). Springer Netherlands. https://doi.org/10.1007/1-4020-4251-5_8
- 845 Hussey, I. (2021). A method to streamline p-hacking. *Meta-Psychology*, 5.
846 <https://doi.org/10.15626/MP.2020.2529>
- 847 Ioannidis, J. P. A. (2005). Why Most Published Research Findings Are False. *PLOS*
848 *Medicine*, 2(8), e124. <https://doi.org/10.1371/journal.pmed.0020124>
- 849 Kerr, N. L. (1998). HARKing: Hypothesizing After the Results are Known. *Personality*
850 *and Social Psychology Review*, 2(3), 196–217.
851 https://doi.org/10.1207/s15327957pspr0203_4
- 852 Koole, S. L., & Lakens, D. (2012). Rewarding Replications: A Sure and Simple Way to
853 Improve Psychological Science. *Perspectives on Psychological Science*, 7(6), 608–614.
854 <https://doi.org/10.1177/1745691612462586>
- 855 Kotzen, M. (2013). Multiple Studies and Evidential Defeat. *Noûs*, 47(1), 154–180.
856 <http://www.jstor.org/stable/43828821>
- 857 Kukla, A. (1990). Clinical Versus Statistical Theory Appraisal. *Psychological Inquiry*, 1(2),
858 160–161. https://doi.org/10.1207/s15327965pli0102_9
- 859 Lakens, D. (2024). When and How to Deviate From a Preregistration. *Collabra*:

- 860 *Psychology*, 10(1), 117094. <https://doi.org/10.1525/collabra.117094>
- 861 Lakens, D. (2019). The value of preregistration for psychological science: A conceptual
862 analysis. *Psychological Science*, 62(3), 221–230. https://doi.org/10.24602/sjpr.62.3_221
- 863 Lishner, D. A. (2015). A Concise Set of Core Recommendations to Improve the
864 Dependability of Psychological Research. *Review of General Psychology*, 19(1), 52–68.
865 <https://doi.org/10.1037/gpr0000028>
- 866 Mayo, D. G. (2018). *Statistical Inference as Severe Testing: How to Get Beyond the*
867 *Statistics Wars* (First). Cambridge University Press.
868 <https://doi.org/10.1017/9781107286184>
- 869 Mayo, D. G., & Spanos, A. (2011). Error Statistics. In *Philosophy of Statistics* (pp.
870 153–198). Elsevier. <https://doi.org/10.1016/B978-0-444-51862-0.50005-8>
- 871 Meehl, P. E. (1990). Appraising and Amending Theories: The Strategy of Lakatosian
872 Defense and Two Principles that Warrant It. *Psychological Inquiry*, 1(2), 108–141.
873 https://doi.org/10.1207/s15327965pli0102_1
- 874 Meehl, P. E. (1978). Theoretical risks and tabular asterisks: Sir Karl, Sir Ronald, and the
875 slow progress of soft psychology. *Journal of Consulting and Clinical Psychology*, 46(4),
876 806–834. <https://doi.org/10.1037/0022-006X.46.4.806>
- 877 Mellor, D. T., & Nosek, B. A. (2018). Easy preregistration will benefit any research.
878 *Nature Human Behaviour*, 2(2), 98–98. <https://doi.org/10.1038/s41562-018-0294-7>
- 879 Niiniluoto, I. (1998). Verisimilitude: The Third Period. *The British Journal for the*
880 *Philosophy of Science*, 49(1), 1–29. <https://doi.org/10.1093/bjps/49.1.1>
- 881 Nosek, B. A., Ebersole, C. R., DeHaven, A. C., & Mellor, D. T. (2018). The preregistration
882 revolution. *Proceedings of the National Academy of Sciences*, 115(11), 2600–2606.
883 <https://doi.org/10.1073/pnas.1708274114>
- 884 Oberauer, K. (2019). Preregistration of a forking path – What does it add to the garden of
885 evidence? In *Psychonomic Society Featured Content*.
- 886 Oberauer, K., & Lewandowsky, S. (2019). Addressing the theory crisis in psychology.

- 887 *Psychonomic Bulletin & Review*, 26(5), 1596–1618.
888 <https://doi.org/10.3758/s13423-019-01645-2>
- 889 Open Science Collaboration. (2015). Estimating the reproducibility of psychological
890 science. *Science*, 349(6251), aac4716. <https://doi.org/10.1126/science.aac4716>
- 891 Orben, A., & Lakens, D. (2020). Crud (Re)Defined. *Advances in Methods and Practices in*
892 *Psychological Science*, 3(2), 238–247. <https://doi.org/10.1177/2515245920917961>
- 893 Peikert, A., & Brandmaier, A. M. (2023a). *Supplemental materials for preprint: Why does*
894 *preregistration increase the persuasiveness of evidence? A Bayesian rationalization.*
895 Zenodo. <https://doi.org/10.5281/zenodo.7648471>
- 896 Peikert, A., & Brandmaier, A. M. (2023b). *Why does preregistration increase the*
897 *persuasiveness of evidence? A Bayesian rationalization.* PsyArXiv; PsyArXiv.
898 <https://doi.org/10.31234/osf.io/cs8wb>
- 899 Peikert, A., & Brandmaier, A. M. (2021). A Reproducible Data Analysis Workflow With R
900 Markdown, Git, Make, and Docker. *Quantitative and Computational Methods in*
901 *Behavioral Sciences*, 1–27. <https://doi.org/10.5964/qcmb.3763>
- 902 Peikert, A., van Lissa, C. J., & Brandmaier, A. M. (2021). Reproducible Research in R: A
903 Tutorial on How to Do the Same Thing More Than Once. *Psych*, 3(4), 836–867.
904 <https://doi.org/10.3390/psych3040053>
- 905 Pham, M. T., & Oh, T. T. (2021). Preregistration Is Neither Sufficient nor Necessary for
906 Good Science. *Journal of Consumer Psychology*, 31(1), 163–176.
907 <https://doi.org/10.1002/jcpy.1209>
- 908 Popper, K. R. (2002). *The logic of scientific discovery.* Routledge.
- 909 Rubin, M. (2020). Does preregistration improve the credibility of research findings? *The*
910 *Quantitative Methods for Psychology*, 16(4), 376–390.
911 <https://doi.org/10.20982/tqmp.16.4.p376>
- 912 Salmon, W. C. (1970). Statistical Explanation. In *The Nature & function of scientific*
913 *theories: Essays in contemporary science and philosophy* (pp. 173–232). University of

- 914 Pittsburgh Press.
- 915 Schönbrodt, F., Gärtner, A., Frank, M., Gollwitzer, M., Ihle, M., Mischkowski, D., Phan, L.
916 V., Schmitt, M., Scheel, A. M., Schubert, A.-L., Steinberg, U., & Leising, D. (2022).
917 *Responsible Research Assessment I: Implementing DORA for hiring and promotion in*
918 *psychology*. PsyArXiv. <https://doi.org/10.31234/osf.io/rgh5b>
- 919 Shmueli, G. (2010). To Explain or to Predict? *Statistical Science*, 25(3), 289–310.
920 <https://doi.org/10.1214/10-STS330>
- 921 Silagy, C. A., Middleton, P., & Hopewell, S. (2002). Publishing Protocols of Systematic
922 Reviews Comparing What Was Done to What Was Planned. *JAMA*, 287(21),
923 2831–2834. <https://doi.org/10.1001/jama.287.21.2831>
- 924 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2021). Pre-registration: Why and How.
925 *Journal of Consumer Psychology*, 31(1), 151–162. <https://doi.org/10.1002/jcpy.1208>
- 926 Sprenger, J., & Hartmann, S. (2019). *Bayesian Philosophy of Science*. Oxford University
927 Press. <https://doi.org/10.1093/oso/9780199672110.001.0001>
- 928 Steegen, S., Tuerlinckx, F., Gelman, A., & Vanpaemel, W. (2016). Increasing Transparency
929 Through a Multiverse Analysis. *Perspectives on Psychological Science*, 11(5), 702–712.
930 <https://doi.org/10.1177/1745691616658637>
- 931 Stefan, A. M., & Schönbrodt, F. D. (2023). Big little lies: A compendium and simulation
932 of p-hacking strategies. *Royal Society Open Science*, 10(2).
933 <https://doi.org/10.1098/rsos.220346>
- 934 Szollosi, A., Kellen, D., Navarro, D. J., Shiffrin, R., Rooij, I. van, Zandt, T. V., & Donkin,
935 C. (2020). Is Preregistration Worthwhile? *Trends in Cognitive Sciences*, 24(2), 94–95.
936 <https://doi.org/10.1016/j.tics.2019.11.009>
- 937 Tukey, J. W. (1972). Exploratory data analysis: As part of a larger whole. *Proceedings of*
938 *the 18th Conference on Design of Experiments in Army Research and Development and*
939 *Training*, 1–18. <https://apps.dtic.mil/sti/tr/pdf/AD0776910.pdf>
- 940 Tukey, J. W. (1980). We Need Both Exploratory and Confirmatory. *The American*

- 941 *Statistician*, 34(1), 23–25. <https://doi.org/10.2307/2682991>
- 942 van Dongen, N., Sprenger, J., & Wagenmakers, E.-J. (2023). A Bayesian perspective on
943 severity: Risky predictions and specific hypotheses. *Psychonomic Bulletin & Review*,
944 30(2), 516–533. <https://doi.org/10.3758/s13423-022-02069-1>
- 945 van Rooij, I., & Baggio, G. (2021). Theory Before the Test: How to Build
946 High-Verisimilitude Explanatory Theories in Psychological Science. *Perspectives on*
947 *Psychological Science*, 16(4), 682–697. <https://doi.org/10.1177/1745691620970604>
- 948 van Rooij, I., & Baggio, G. (2020). Theory Development Requires an Epistemological Sea
949 Change. *Psychological Inquiry*, 31(4), 321–325.
950 <https://doi.org/10.1080/1047840X.2020.1853477>
- 951 Wagenmakers, E.-J., Wetzels, R., Borsboom, D., van der Maas, H. L. J., & Kievit, R. A.
952 (2012). An Agenda for Purely Confirmatory Research. *Perspectives on Psychological*
953 *Science*, 7(6), 632–638. <https://doi.org/10.1177/1745691612463078>
- 954 Willroth, E. C., & Atherton, O. E. (2023). *Best Laid Plans: A Guide to Reporting*
955 *Preregistration Deviations* [Preprint]. PsyArXiv. <https://doi.org/10.31234/osf.io/dwx69>