

No need to be indirect: On the role of data in the validation of theory

Niels Vanhasbroeck¹, Kenny Yu², and Sigert Ariens²

¹Psychological Methods, University of Amsterdam, Amsterdam, the Netherlands

²Quantitative Psychology and Individual Differences, KU Leuven, Leuven, Belgium

Author Note

All authors contributed equally. Authors can be contacted through their respective email addresses: niels.vanhasbroeck@gmail.com, kenny.yu@kuleuven.be, sigert.ariens@kuleuven.be

Abstract

In recent years, the case has been made to move from narrative accounts of psychological theories to rigorous mathematical descriptions thereof. Yet, a persistent tension remains in how the models that result from the formalization of psychological theories are to be evaluated – whether by their theoretical elegance, their ability to reproduce established phenomena, or their quantitative fit to empirical data. This paper challenges an exclusive use of validation paradigms that are based on phenomena or data alone as evaluative unit by reexamining the relationships between theory, phenomena, and data. We argue that testing the tenability of a theory or model while placing too much emphasis on only phenomena or data is inherently fallible. Instead, we argue that both phenomena and data should occupy a central role in theory testing. This perspective demands that models make substantive and explanatory claims while engaging with the full structure of empirical observations. This perspective not only clarifies the dynamic interdependence among theory, phenomena, and data, but also establishes a principled basis for cumulative, iterative scientific progress in psychological science.

Keywords: computational modeling; theory evaluation; phenomena; simulation; data

No need to be indirect: On the role of data in the validation of theory

Introduction

1

2 Theories are the engines of scientific progress; They offer an explanation for how
3 a system behaves and help predict as of yet unknown behavior that may occur in the
4 future (Popper, 1959). While psychological researchers have a history of using natural
5 language to communicate their theories, they are increasingly turning towards the
6 language of mathematics for this purpose (Borsboom & Haslbeck, 2024; Farrell &
7 Lewandowsky, 2010; Smaldino, 2017). While not uncontroversial (Eronen & Romeijn,
8 2020; Oude Maatman, preprint), we align with scholars who view this evolution as
9 advantageous in several respects. For example, the use of mathematics allows for a
10 greater precision in articulating explanatory mechanisms and in generating testable
11 predictions. This increased precision, together with the tools of mathematics, provides
12 researchers with rigorous methods to evaluate the adequacy of their theoretical
13 frameworks and derive novel predictions that might otherwise have remained obscure in
14 purely verbal accounts.

15 As efforts to formalize theories continue, so too must our efforts for the
16 development of methods for evaluating them. Specifically, we should ask what
17 constitutes as evidence for or against a theory, and how we can decide whether a model
18 genuinely captures the psychological processes it purports to explain. Fortunately, a
19 vast literature on model evaluation procedures exists, providing researchers with various
20 tools to evaluate a model's explanatory and predictive power (see, for example Navarro
21 et al., 2004; D. R. Roberts et al., 2017; Schwarz, 1978; Tashman, 2000; Wagenmakers
22 et al., 2004; Wehrens et al., 2000).

23 However, researchers increasingly express skepticism towards the practice of
24 evaluating a model's fit to empirical data, with some scholars arguing that such an
25 analysis, when conducted in isolation, yields little meaningful insight (S. Roberts &
26 Pashler, 2000). Instead of investigating fit, we should investigate a theory's ability to
27 make "risky predictions", that is predictions that would be unlikely to hold if the theory
28 were false, and use this as the primary basis of evaluating a theory's validity (Lakatos,

29 1978; Navarro, 2019; S. Roberts & Pashler, 2000). This perspective has been voiced
30 persuasively by Borsboom and colleagues, who, across a series of recent publications,
31 have argued that psychological phenomena provide a more robust foundation for the
32 development and evaluation of theories (Borsboom & Haslbeck, 2024; Borsboom et al.,
33 2021; Haslbeck et al., 2022; van Dongen et al., 2025). According to this framework,
34 (formal) theories should address the mechanisms underlying one or more psychological
35 phenomena, and the testing of such theories should consequently focus on the
36 predictions of phenomena under a pre-specified set of conditions. Importantly, they
37 contrast this phenomenon-based approach to the dominant data-driven approach in
38 which a model’s fit to data is taken as evidence for or against the theory that has been
39 formalized to the model. The validity of this approach is called into question on
40 grounds that the data that serve as a basis for testing theory are only unreliably so
41 (Borsboom & Haslbeck, 2024; van Dongen et al., 2025).

42 We wholeheartedly agree that a crucial aspect of building and reasoning about
43 the validity of a theory concerns its ability to explain phenomena. However, when
44 focusing on phenomena, one may start to underappreciate the role empirical data play
45 in *validating* theory, and in our opinion unrightfully so. In this article, we therefore aim
46 to clarify the role data can and should play when validating a formal theory.

47 In what follows, we will first discuss the data-driven and phenomenon-based
48 approaches to theory validation and show where a pure application of either of these
49 approaches may go wrong. We then place data alongside phenomena in an integrated
50 approach to theory testing, highlighting the complimentary role both approaches play
51 in evaluating theory. We end this paper with a short discussion on the integrated
52 approach an a historical example.

53 In this article, we focus primarily on theories that are expressed as mathematical
54 equations, guided by our shared expertise in statistics and mathematical psychology.
55 We furthermore want to emphasize that we fundamentally agree with many aspects of
56 the phenomenon-based approach outlined in previous literature. We endorse the
57 construction of theory through phenomena and the testing of theory through “risky”

58 predictions, positions advocated by numerous philosophers and psychologists (Haslbeck
59 et al., 2022; Lakatos, 1978; Navarro, 2019). Our contribution therefore does not aim to
60 discredit the consideration of phenomena, but rather aims to provide a crucial missing
61 element to the ongoing discussion on the evaluation of theory.

62 Definitions

63 Before we present our argument, it is useful to define the primary concepts of
64 interest, namely those of *data*, *phenomena*, *theory*, *model*, and *theory validation*.

65 Data

66 In this paper, we define data as a collection of observations that are
67 systematically structured and recorded (a definition aligned with van Dongen et al.,
68 2025). These observations are typically encoded as variables that serve as the basis for
69 subsequent analysis. Examples include self-reported ratings (Cloos et al., 2023; Krosnick
70 & Fabrigar, 1997; Russel et al., 1989), spatial position recordings (Bastida-Castillo
71 et al., 2019; Gamble et al., 2023; Hedrick, 2008), and psychophysiological measurements
72 (Dawson et al., 2000; Mauss & Robinson, 2009; Read, 2017).

73 A fundamental characteristic of data is its particularity – each data point
74 represents an observation obtained from a specific individual, within a specific context,
75 at a specific moment in time. In their raw form, data lack inherent structure or
76 generalizability: The process of identifying regularities within data constitutes the
77 discovery of *phenomena*. Given that theories aim to explain regularities rather than
78 idiosyncrasies, this particularity presents a challenge for using raw data directly to
79 validate formal theories. The gap between particular observations and generalizable
80 theoretical principles requires bridging mechanisms that can connect these distinct
81 epistemological domains.

82 Phenomena

83 Given the inherent limitations of data in their particularity, researchers have
84 advocated for phenomena as more appropriate targets for validating theory (Bogen &
85 Woodward, 1988; Haig, 2005; Woodward, 1989). Phenomena are conceptualized as
86 stable, recurring patterns that manifest consistently across diverse datasets and

87 contexts. Examples of phenomena are the speed-accuracy tradeoff (Brown & Heathcote,
88 2008; Dutilh et al., 2011; Ratcliff & Rouder, 1998), lowered mood at the end of
89 multi-trial experiments (Jangraw et al., 2023), and the power law for learning (Logan,
90 1992; Newell and Rosenbloom, 1980, but see Heathcote et al., 2000). Critically,
91 phenomena exist at a higher level of abstraction than individual data points or specific
92 datasets – they represent generalizable regularities rather than context-dependent
93 particulars. This generalizability makes phenomena especially suitable targets for
94 theoretical explanation, as theories themselves aim to capture general principles rather
95 than idiosyncratic observations (van Dongen et al., 2025).

96 **Theory**

97 In this paper, we define theory as a systematic and explanatory account of how
98 aspects of the world work (following van Dongen et al., 2025). As convincingly argued
99 by several authors (Bogen & Woodward, 1988; Borsboom et al., 2021; Guest & Martin,
100 2021), theories should be constructed to explain phenomena, either by describing how a
101 collection of phenomena interrelate or, preferably, by targeting the mechanisms that
102 underlie them. These mechanisms often focus on causal relationships between variables
103 which, together, account for the manifestation of the phenomena. Examples of such
104 theories are the kinetic theory of gases (Pathria & Beale, 2022; Schroeder, 2021), the
105 theory of evolution (Darwin, 1859), and cognitive dissonance theory (Festinger, 1957).

106 An important characteristic of theories is that they can be used to predict under
107 which conditions a given phenomenon should occur, thereby providing an opportunity
108 to test a theory’s adequacy to explain the phenomenon. However, deriving precise
109 predictions from theories often presents significant challenges, particularly when
110 theories are formulated with insufficient specificity (Oude Maatman, preprint). This
111 limitation has led researchers to advocate for the development of *formal theories* or
112 *models* (Borsboom et al., 2021; Guest & Martin, 2021; van Rooij & Baggio, 2021),
113 which we define in the next section.

114 **Model**

115 For the purposes of this article, we distinguish between a theory and its
116 formalization as a mathematical or computational *model*. While a theory requires the
117 specification of all theoretically relevant mechanisms that may give rise to a particular
118 set of phenomena, a formal model attempts to convey these mechanisms by translating
119 them, as accurately as possible, into mathematical equations. Formalizing a theory into
120 a model requires both making implicit theoretical assumptions explicit as well as the
121 specification of auxiliary assumptions that may not be relevant for the theory itself
122 (Smaldino, 2020; Suppes, 1962; van Dongen et al., 2025). Examples of such auxiliary
123 assumptions concern assumptions of statistical independence or normality, which are
124 typically imposed to allow for a model’s estimation on data. Throughout this article, we
125 use “model” to refer to what we previously called a “formal theory”, therefore
126 considering models to represent a theory that has been formalized into mathematical
127 language.

128 Interestingly, the formalization of theory has been argued to greatly increase the
129 precision with which the mechanisms that generate a particular set of phenomena are
130 described (Borsboom & Haslbeck, 2024; van Dongen et al., 2025; van Rooij & Baggio,
131 2021). Importantly, this statement is not uncontroversial. For example, some authors
132 have argued that psychology is not ready for formalization yet, first requiring us to
133 precisely define our constructs before we can meaningfully formulate a theory, let alone
134 formalize one (Eronen & Bringmann, 2021; Kellen et al., 2021). Additionally, the
135 translation of a theory into a model does not imply an increased precision in its
136 specification (Oude Maatman, preprint): If a theory lacks specificity, then a model
137 based on this theory will likely suffer the same fate. We agree with both of these
138 arguments, but at the same time believe that the act of formalization forces researchers
139 to think their theory through. The fact that, in formalization, researchers have to
140 engage with all aspects of a theory, greatly constrains researcher degrees of freedom and
141 will, ultimately, lead to more rigorous and transparent scientific discourse.

142 **Theory validation**

143 Theory validation concerns the adequacy with which a theory captures the
144 phenomena it aims to explain. We distinguish between two related aspects of theory
145 validation that will clarify the argument presented in this article.

146 The first aspect of theory validation concerns *theory building*. During this phase,
147 the researcher should identify the phenomena to-be-explained and construct a theory
148 that can account for these phenomena. When formalizing a theory into a model, we can
149 validate the model through examining whether it can adequately produce the
150 phenomena that it has to explain (Borsboom & Haslbeck, 2024; Borsboom et al., 2021).
151 Taking the recently proposed model of panic disorder as an example (Robinaugh et al.,
152 2024), the validity of this model depends, in part, on its success in producing
153 phenomena such as the avoidance of panic attacks through engaging in escape behavior
154 and the efficacy of cognitive behavioral therapy in treating a patient with panic
155 disorder. Ideally, a valid theory should explain all relevant phenomena for the construct
156 of interest.

157 The second aspect of theory validation concerns *theory testing*. During this
158 phase, the researcher is concerned with whether the theory's predictions are in
159 accordance with actual observation. For this, one typically sets up a study to test
160 whether a theory's predictions hold in empirical practice. Again taking the model of
161 Robinaugh et al. (2024) as an example, one can validate this model by examining this
162 model's predictions with regard to the success rate of different types of therapy and
163 comparing these predictions to empirical observations. In these types of studies, the
164 strength of evidence is often taken to depend on the "riskiness" of the prediction, that is
165 on the a priori likelihood of the prediction and the prediction's uniqueness among
166 competing theories (Vanpaemel, 2020).

167 **Phenomenon-based Approach**

168 As its name implies, the phenomenon-based approach places phenomena at the
169 center of theory validation (Bogen & Woodward, 1988; Borsboom et al., 2021; Guest &
170 Martin, 2021). On the one hand, this implies that observed phenomena inform theory,

171 so that a theory should be constructed in such a way as to explain phenomena that are
172 already known. For example, a new theory of gravity needs to explain all currently
173 observed phenomena that are relevant to it, going from a simple apple falling down the
174 tree to the bending of light around a large object. On the other hand, theory can also
175 predict new, currently unobserved phenomena that, if discovered in the way the theory
176 predicted, grants credence to the theory. Staying within the gravitation example, the
177 discovery of Neptune due to irregularities in the orbit of Uranus is often cited as
178 showing this principle in action (Gershman, 2019).

179 In summary, the phenomenon-based approach holds that a theory constrains
180 phenomena and that these phenomena should therefore serve as primary targets for
181 theory validation. For theory building, this implies that one should assess the theory's
182 success in reproducing a set of phenomena that are deemed interesting to the
183 researcher. When using a model, the primary analysis then consists of repeatedly
184 simulating the model's behavior under different types of conditions and different values
185 for the parameters, providing one with a comprehensive overview of what phenomena
186 the model is able to produce (Pitt et al., 2006). Here, we agree with proponents of the
187 phenomenon-based approach in saying that this step is invaluable for those who aim to
188 construct (and formalize) theory.

189 For theory testing, the phenomenon-based approach implies that one should test
190 a theory's validity through an assessment of its predicted behavior. Again bringing this
191 to the model level, one should derive a model's predicted behavior through either
192 analytic means or through simulation studies and assess whether these predictions hold
193 in empirical studies (e.g., Ariens et al., 2023; Navarro, 2019) Proponents of the
194 phenomenon-based approach argue that phenomena provide sufficient, albeit indirect,
195 evidence for theory testing (Borsboom & Haslbeck, 2024; Haslbeck et al., 2022).
196 However, we do not agree with this statement and believe that phenomena are not a
197 sufficient basis to achieve this type of inference. In what follows, we will formulate three
198 challenges to the phenomenon-based approach that show the insufficiency of phenomena
199 for the testing of theory.

200 **Challenges to the Phenomenon-based Approach**

201 We now discuss three challenges when using phenomena to test a model. We will
202 illustrate each of these challenges with an example that was chosen to be as simple as
203 possible yet remaining as relevant as we can to real situations researchers may
204 encounter in their investigations.

205 *Phenomena Selection*

206 This first issue concerns the problem of selecting phenomena for the validation of
207 a particular model. It can be very difficult to find a set of *sufficient phenomena* – that
208 is, a set of phenomena which can validate a theory in its entirety. To see this, we draw
209 an instructive parallel to the concept of sufficient statistics in statistical theory. In
210 statistics, a sufficient statistic contains all the information in a sample that is relevant
211 for estimating the parameters of a model. For instance, when sampling from a normal
212 distribution, the sample mean and variance together form the sufficient statistics for
213 estimating the population parameters. However, the mean alone is insufficient: Multiple
214 normal distributions can produce the same mean through variations in their variance.
215 Similarly, in theory testing, a single phenomenon is often insufficient to discriminate
216 between competing theories, as multiple theoretical mechanisms could generate the
217 same phenomenon. Practically, this means researchers should be able to identify those
218 phenomena that are sufficient in order to validate that theory. While we believe that
219 there are sufficient phenomena to be found, identifying them becomes troublesome in
220 more complex psychological theories where numerous underlying mechanisms could
221 generate similar surface-level phenomena (e.g., Yu et al., 2023).

222 **Example 1:** To make this concrete, consider an example where researchers are
223 testing competing theories of decision making. Through manipulation of an
224 experimental variable (e.g., difficulty of an item), they want to investigate how different
225 models compare with regard to the capturing mean differences in response times
226 between the conditions. Because such differences are consistently found over data
227 particulars, the researchers consider these mean differences to be the crucial
228 phenomenon-of-interest that any theory must be able to explain.

229 Now imagine four competing models, each specifying different expected response
230 time distributions based on distinct underlying cognitive processes. Model A posits that
231 response times follow a normal distribution, with differences between experimental
232 conditions arising from a shift in the mean of this distribution—representing a simple
233 additive effect. Model B posits a bimodal distribution, where the observed mean
234 difference results from a strategic mixture between two response modes (e.g., fast
235 guessing versus deliberate responding). Model C posits a skewed distribution, reflecting
236 a threshold-based evidence accumulation process (such as in drift-diffusion models),
237 where condition differences emerge from changes in the decision threshold or drift rate.
238 Model D posits a heavy-tailed distribution, where rare but extreme response times play
239 a substantial role, and mean differences are driven by changes in the tail behavior (e.g.,
240 increased likelihood of lapses or extreme delays). Through applying the
241 phenomenon-based approach in isolation, we found that despite the fundamentally
242 different assumptions and predictions of the four theories, we are not able to distinguish
243 between them based on only this single phenomenon (see Figure 1). This shows how a
244 selective focus on a single phenomenon can lead researchers to incorrectly confirm their
245 preferred theory.

246 It is useful to repeat that we do think that there are sufficient phenomena to be
247 found: In this example, we may have looked at other distributional parameters such as
248 the skewness and kurtosis to distinguish between the four competing models (see Panel
249 (c) in Figure 1). However, identifying such sufficient phenomena becomes increasingly
250 difficult with more complex models that share prediction of many different phenomena.

251 *Model Complexity*

252 A single phenomenon can be implied by many formal theories, but the converse
253 is also true, so that one model can imply a wide range of different phenomena. Indeed,
254 while it is typically argued that verbal theories lack specificity, we should appreciate
255 that mathematical models can suffer from the same limitation. In extreme cases,
256 models can predict virtually any phenomenon if their parameters or boundary
257 conditions are “tweaked” in a certain way. More generally, one can always come up with

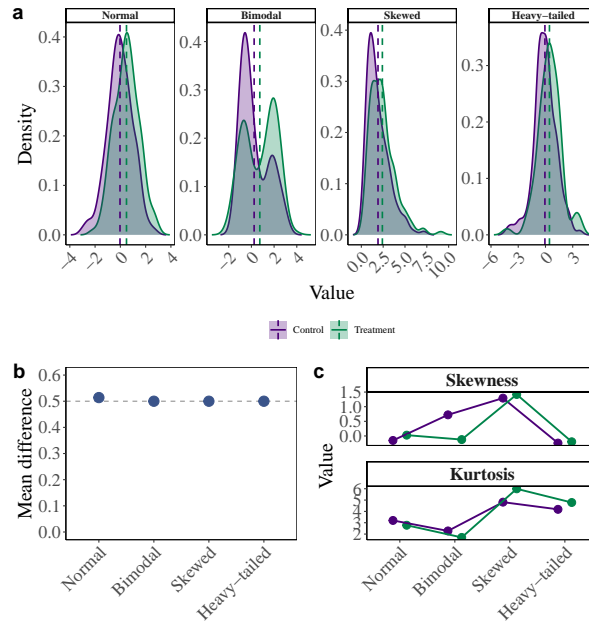


Figure 1

Illustration of the selective focus issue. Using the models described in the section, we simulated 150 data points for each model and each condition separately. Panel (a) shows the four distributions of the model, each having a distinctly different shapes representing different psychological processes: Normal (simple additive effects), bimodal (strategic choice between options), skewed (threshold-based accumulation), and heavy-tailed (processes with occasional extreme responses). Panel (b) shows that all four models produce identical mean differences of 0.5 between control and treatment conditions, which serves as the phenomenon-of-interest. Panel (c) visualizes the different skewness and kurtosis values of the models, providing clear evidence that they represent different underlying mechanisms.

258 a more complex model that generates several phenomena equally well as long as it is
 259 used in a particular way, as we illustrate in the following example.

260 **Example 2:** In the field of affect dynamics, an important phenomenon is
 261 “emotional inertia” or the extent to which affective states linger on over time. Imagine
 262 that there are two models that attempt to formalize the functional form of this inertia,
 263 so that Model A assumes that affect decays exponentially over time (Ariens et al., 2020;
 264 Rutledge et al., 2014) while Model B states that affect decays in a quasi-hyperbolic

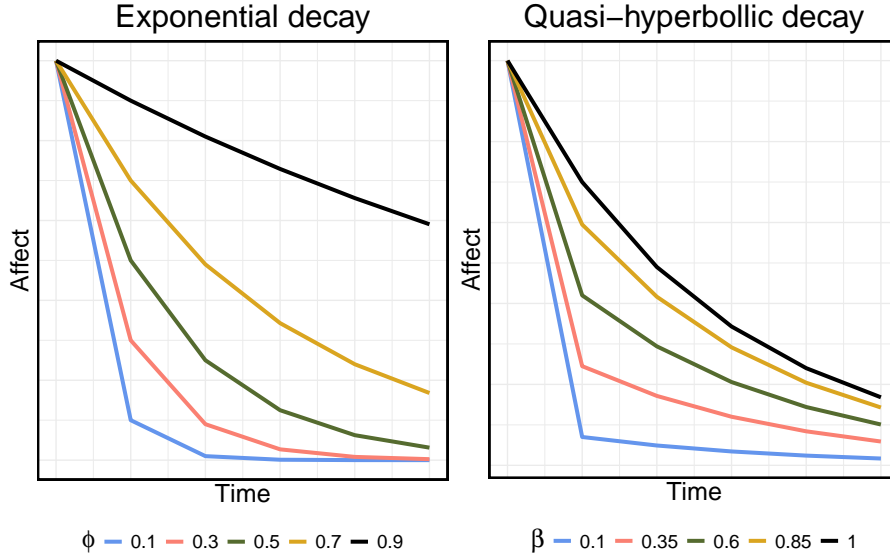


Figure 2

Illustration of the functional form of emotional inertia according to Model A (exponential, left) and Model B (quasi-hyperbolic, right). For the exponential model, we vary the defining parameter $\phi \in \{0.10, 0.30, 0.50, 0.70, 0.90\}$ to show the strength of the regulation that is expected with each of these parameter settings. For the quasi-hyperbolic model, we fixed the parameter $\delta = 0.70$ and vary the value of the parameter $\beta \in \{0.10, 0.35, 0.60, 0.85, 1.00\}$. Observe that the quasi-hyperbolic expectation for the value of parameter $\beta = 1$ corresponds to the expectation of the exponential model for the value of $\phi = 0.70$.

265 fashion (Laibson, 1997; Vanhasbroeck et al., 2024, see Figure 2)). In symbols, these
 266 models can be written down as:

$$\text{A: } y_t = \phi^j y_{t-j} \quad \phi \in [0, 1)$$

$$\text{B: } y_t = \beta \delta^j y_{t-j} \quad \beta \in [0, 1], \delta \in [0, 1)$$

267 where $j \in \{1, \dots, T\}$ represents the lag in the variable y .

268 Following the phenomenon-based approach, we have to find out which of the two
 269 models best reproduces the observed decay rates. However, if we imagine that in reality
 270 affect decays exponentially, then we will not be able to distinguish between these two
 271 models, as Model B is able to produce both exponential and quasi-hyperbolic decay
 272 through variation in its parameter β . In other words, reproducing the
 273 phenomenon-of-interest is not enough to differentiate between these two competitor

274 models.

275 To get around this issue, one could consider the parsimony of a model to
276 distinguish between alternatives (Borsboom et al., 2021). However, to our knowledge,
277 the phenomenon-based approach provides no principled way of doing so, meaning there
278 is no inherent limit to the complexity of the models we consider. Exacerbating this
279 issue is that parsimony is difficult to define when based on phenomena alone, as even
280 relatively simple models can produce complicated ranges of phenomena. One example is
281 the linear differential equation which, under the right specification of the eigenvalues of
282 its drift matrix, allows not only for the typical exponential decay function, but also for
283 more complex oscillatory behavior (Strogatz, 2018). Should this be considered a
284 “simple” model due to its limited number of parameters, or as a “complex” model due
285 to the wide range of behavior that the model predicts? At this moment, the
286 phenomenon-based approach does not answer this question.

287 The relevance of the complexity debate for modeling is further illustrated in an
288 anecdote involving physicist Enrico Fermi, as retold by Freeman Dyson (Dyson, 2004).
289 At the time, Dyson was investigating the strong nuclear force using mathematical
290 equations that he had previously used to understand the weak nuclear force and which
291 had their roots in the theory of quantum electrodynamics, which states that electrons
292 and protons interact through weak electromagnetic forces. After finishing his
293 calculations on the strong nuclear force, Dyson asked Fermi for his opinion on his work.
294 Fermi responded by asking “How many arbitrary parameters did you use for your
295 calculations?”, afterwards concluding by quoting an adage of his friend John von
296 Neumann: “With four parameters I can fit an elephant, and with five I can make him
297 wiggle his trunk” (p. 297, Dyson, 2004). Initially disheartened, Dyson abandoned the
298 project and only later appreciated the depth of Fermi’s critique with the discovery of
299 the quarks, which enhanced the understanding of the strong nuclear force more so than
300 his complex model ever could have.

301 *Circularity*

302 The final limitation of the phenomenon-based approach is the fact that it takes
 303 phenomena as a given. If a phenomenon is indeed present in nature, a valid model
 304 should necessarily recover that phenomenon. However, this recovery becomes
 305 problematic if a phenomenon turns out to be non-existent (see also van Dongen et al.,
 306 2025) The strength of inference within the phenomenon-based approach therefore
 307 depends heavily on the phenomenon's existence.

308 **Example 3:** To make this issue concrete, imagine a research process where
 309 researchers attempt to uncover the true relationship between two quantities by building
 310 formal models and testing them using the phenomenon-based approach. Imagine that
 311 the true relationship between two variables X and Y is given by the following
 312 third-degree polynomial:

$$Y = \alpha + \beta X + \gamma X^2 + \delta X^3.$$

313 where $\alpha, \beta, \gamma \in \mathbb{R}$ and $\delta > 0$. Now imagine that researchers routinely fit a linear
 314 regression model to their data in order to estimate the relationship between X and Y :

$$Y = \alpha + \beta X + v_t$$

$$v_t \stackrel{iid}{\sim} N(0, \sigma^2).$$

315 Across many studies, researchers would discover a seemingly linear relationship
 316 between X and Y , in the sense that the slope $\hat{\beta}$ would consistently be estimated to be
 317 positive. Because of the robustness of this phenomenon, theory builders conclude that
 318 their formal theories should accommodate the expected positive linear relationship
 319 between X and Y .

320 When it comes to theory testing, however, the ability of the model to reproduce
 321 a curve of the form $Y = \alpha + \beta X$ is clearly irrelevant, if not counterproductive: The
 322 more accurately the formal theory reproduces the phenomenon, the less likely it will be
 323 to conform to the true relationship $Y = \alpha + \beta X + \gamma X^2 + \delta X^3$. Reproducing the
 324 phenomena does not bring the literature closer to discovering, understanding, or

325 explaining the true relationship between the variables Y and X . If the practice of
326 validating a theory using the phenomena it was designed to explain becomes
327 mainstream, it may be to the detriment of new ideas or findings which challenge the
328 current way of thinking. More strongly, a pure application of the phenomenon-based
329 approach creates an echo chamber where theory prescribes which phenomena to model
330 without questioning whether these phenomena were correctly identified or without
331 allowing for data to disagree.

332 **Data-driven Approach**

333 As an alternative to the phenomenon-based approach, researchers can instead
334 turn to the dominant data-driven approach to theory validation. This approach holds
335 that data serve as the primary targets for theory validation and that theory testing can
336 be achieved through the use of model estimation and model selection techniques. The
337 fit of the model to the data then serves as a tool to test the model and the theory it
338 formalizes. Beyond simple fitting, one can also validate the model through prediction of
339 unseen data. The cross-validation method levies on these types of inferences, for
340 example, where a model is estimated on a particular part of the data called the
341 “training set” and then used to predict the unseen data contained in the “test set”
342 (Bates et al., 2023; Hastie et al., 2011; Stone, 1974) Importantly, both methods directly
343 use empirical data to test the theories that are formalized with the models.

344 Note that these data-driven methods are not very useful when applied to only a
345 singular model (S. Roberts & Pashler, 2000). Ideally, researchers compare the
346 performance of multiple models that encompass different theories on the same dataset.
347 Such a comparison is thought to inform one on the tenability of the theoretical
348 assumptions that differ between models, therefore serving as a comparative form of
349 theory testing.

350 Given our critique of the phenomenon-based approach, readers might expect us
351 to champion the application of these pure data-driven methods. This would be
352 mistaken. While the data-driven approach provides us with some formal tools that lack
353 in the phenomenon-based approach, it faces its own critical limitations, which we turn

354 to now.

355 **Challenges to the Data-driven Approach**

356 *Auxiliary Assumptions*

357 As we previously mentioned, researchers are required to make some auxiliary
358 assumptions when formalizing a theory into a model (Suppes, 1962). These assumptions
359 concern many different aspects, including the properties of the measurements (Krosnick
360 & Presser, 2010; Stevens, 1946), structures that relate theoretical constructs to
361 observables (Kellen et al., 2021; Regenwetter & Robinson, 2017; Robinson et al., 2025),
362 or the structure of the error (Westfall & Yarkoni, 2016). Critically, when we use models
363 to test theory, we never test theories in isolation. Rather, we test a theory plus
364 auxiliary assumptions. This creates a fundamental attribution problem. When a model
365 fits data poorly, this failure could come from multiple sources: The core theory could be
366 wrong, but so could any of the auxiliary assumptions of the model. Unfortunately, data
367 alone cannot indicate the source of the misfit, meaning that current judgment of the
368 performance of the model is largely determined by its theory-irrelevant auxiliary
369 assumptions.

370 *Estimation*

371 A more practical issue concerns that of estimation. To be able to leverage on the
372 data-driven approach, one requires their model to be estimable. This can be difficult to
373 achieve for more complex models, especially given that classical estimation methods
374 require the specification of a likelihood function (Spanos, 2024; Viscardi et al., 2025).
375 Recent developments have tried to address this problem, using for example
376 simulation-based methods (Mestdagh et al., 2019) or neural networks (Radev et al.,
377 2020, 2023) to allow for model estimation without the explicit derivation of its
378 likelihood function. Such developments are very useful, and allow researchers to focus
379 on the structure of the models rather than how one should go about estimating and
380 performing inference on the parameters. However, even with these developments, it still
381 holds that more complex models will be more difficult to estimate, hindering the
382 application of these models in research.

412 complimentary. To test a theory, we need to consider both phenomena and data as they
413 provide important counterweights to each approach's limitations. The way forward is
414 thus to combine both approaches in an integrative approach of theory validation.

415 In the integrative approach, we distinguish between a top-down
416 (phenomenon-based) and a bottom-up (data-driven) stream of analysis. The top-down
417 stream is concerned with making and validating predictions from theory to phenomena.
418 It starts at the theoretical framework and makes predictions about the patterns one
419 should observe in data as well as about the conditions under which a given set of
420 phenomena should occur. Key questions are therefore “whether a model can reproduce
421 theoretically interesting patterns in the data” and “whether a model is able to make
422 useful predictions”?

423 The bottom-up stream is concerned with evaluating how well a model fits the
424 full structure of observed data. In other words, it starts at the data and examines how
425 well a model can capture the statistical properties of the data. This allows not only for
426 the comparison of different models with regard to their relative fit to a given dataset,
427 but also for the discovery of new phenomena through the assessment of absolute fit to
428 the data and model assumptions. Key questions are therefore “how closely a model fits
429 the data relative to another model” and “whether the model captures all systematic
430 patterns in the data, or whether there are signs of systematic misfit”.

431 Importantly, both streams of analysis are necessary for theory testing. The
432 top-down stream ensures that our tests remain theoretically grounded and produce
433 theoretically interesting predictions. Meanwhile, the bottom-up approach ensures that
434 we do not get stuck in an echo chamber, allowing us to more easily see what the model
435 can and cannot explain yet.

436 **Counterweight to the Challenges**

437 While the integrative approach cannot fully resolve the fundamental limitations
438 of the two approaches it is based on, it does provide an important counterweight to
439 these challenges. We will illustrate this by discussing how the integrative approach deals
440 with the previously identified challenges to the phenomenon-based approach.

Table 1

Model comparison results using information criteria to solve the problem in Example 1.

Best-fitting models for each dataset are indicated in bold.

Model	Information Criteria (AIC/BIC)			
	Normal	Bimodal	Skewed	Heavy-tailed
Normal	850 / 865	1026 / 1040	1040 / 1054	1006 / 1021
Bimodal	856 / 884	909 / 938	969 / 997	1006 / 1035
Skewed	956 / 978	1009 / 1030	938 / 960	1182 / 1204
Heavy-tailed	854 / 876	1030 / 1051	1019 / 1040	993 / 1014

441 *Phenomena selection*

442 In *Example 1*, we showed that four different formal theories produced identical
 443 mean differences of 0.5 between conditions in a decision-making task, making them
 444 indistinguishable based on this phenomenon alone. As we have explained in the section
 445 around this issue, one could try to alleviate this issue by trying to find another set of
 446 phenomena that is sufficient to distinguish between different theories. However, one
 447 could also address the issue more directly, namely through assessing the fit of the formal
 448 theories model to the data. Indeed, by fitting the formal theory, we are imposing all of
 449 its assumptions on the data and we need not select individual phenomena.

450 In Table 1, we show the relative fit as assessed through the information criteria
 451 AIC and BIC. The results paint a clear picture: When a model was fit to data that it
 452 itself generated, it outperformed the alternative models. By considering data, we are
 453 now able to distinguish between the different models because all implications of the
 454 models, including their differences with respect to phenomena other than the mean,
 455 have been imposed on the data.

456 Note that model fit alone would not be sufficient to assess the validity of the
 457 theory. Instead, it is the combination of the data-driven model fit together with the
 458 phenomenon-based predictions that lead us to believe one model is a better
 459 representation of reality than the other.

460 *Model complexity*

461 One of the primary advantages of the integrative approach is that it allows for a
462 principled way to handle the complexity of a model. In *Example 2*, we described a
463 situation in which the reproduction of a phenomenon alone did not allow us to find out
464 whether affect decays in an exponential or quasi-hyperbolic fashion. However, by
465 assessing fit more directly and penalizing the models according to the number of
466 parameters they have, we are able to get around this issue.

467 To illustrate this, we report the results of a simulation-based model comparison
468 of the situation described in *Example 2*. We simulated 20 time-series of 5 datapoints
469 with the exponential model, setting its parameter $\phi = 0.25$ and its error variance $\sigma^2 = 1$
470 (see Figure 2). When we then assess the AIC and BIC of the fit of the exponential and
471 quasi-hyperbolic models, we find that the former ($AIC = -15.44$, $BIC = -9.63$)
472 outperforms the latter ($AIC = -14.40$, $BIC = -5.98$). This again illustrates the value
473 of using the integrative approach rather than an approach based on only phenomena.

474 Like in the previous example, we do not wish to suggest that the data-driven fit
475 of the model is sufficient to assess the models' validity. If either Model A or Model B
476 would have been unable to reproduce the phenomenon-of-interest (in this case,
477 exponential decay), then we follow the phenomenon-based approach in saying that this
478 would count as evidence against that particular model. Instead, our argument is that
479 only looking at this reproduction is not enough to validate a model, and that one
480 should additionally use data-driven approaches to have a more comprehensive test of
481 the model and the theory that it formalizes.

482 *Circularity*

483 In *Example 3*, we showed that the the ability of a model to reproduce a
484 phenomenon can be counterproductive if the phenomenon itself is incorrectly specified.
485 This is the case due to the absence of a principled correcting mechanism. Within the
486 integrative approach, one is able to avoid this issue through a systematic study of
487 misfit. Specifically, when one fits a model to data, one can examine the residuals of this
488 model to identify missing or incorrectly specified phenomena (see Figure 3). Conversely,

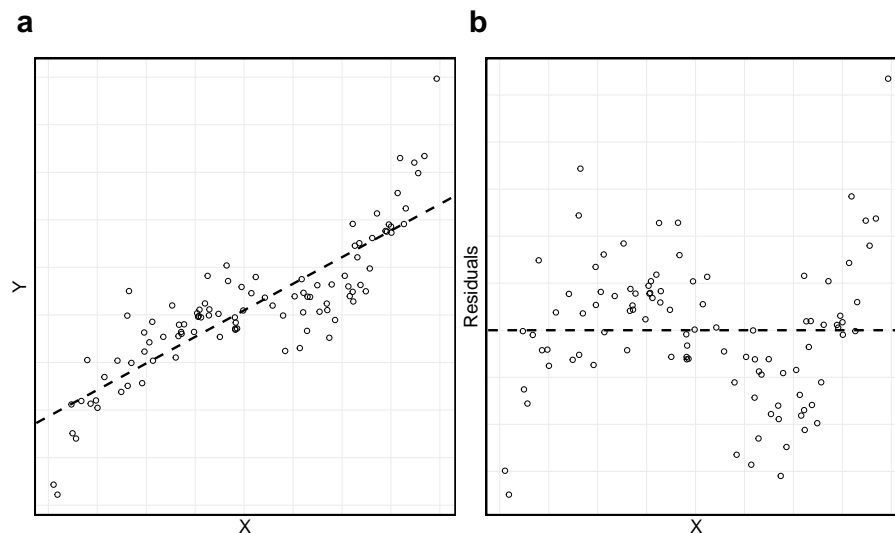


Figure 3

Inspection of the residuals of the linear regression model when fitted to $N = 100$ data points generated from a third-degree polynomial with $\alpha = 0, \beta = \gamma = \delta = 1$; The situation we described in Example 3. Panel (a) shows the fit of the linear regression to the raw data. Panel (b) zooms in on the residuals of this fit, showing that residuals are not independently and identically distributed, which signals model misfit.

489 one may also spot overfit of a model to the data through the observation of extremely
 490 small residual variances, therefore also accommodating the issue of model complexity.
 491 Finally, one can check the values of the estimated parameters of the formal theory,
 492 which may lie outside of the boundaries that are imposed by the theory.

493 Note that the use of data for the validation of a model could also be argued to
 494 be circular, as the models one fits to the data in principle need to be correctly specified
 495 to retrieve unbiased estimates of the model parameters (Ariens et al., 2023). However, a
 496 key advantage of the integrative approach is that it allows for the data to speak against
 497 the merits of the model we fit to it. These assumption violations are therefore important
 498 sources of information rather than nuisances and allow for researchers to detect that
 499 something has gone wrong. Unfortunately, this kind of correction mechanism is often
 500 neglected (for exceptions, see Revol et al., in press; Vanhasbroeck et al., 2022).

501 We must admit, though, that if a researcher encounters an assumption violation,
 502 this might indicate a failure of an auxiliary assumption rather than a theoretical one. In

503 this case, the task of the model builder will be to reconsider either the method with
504 which data has been collected, the auxiliary assumptions of the model, or the theory
505 itself in light of this mismatch.

506 **Strengths and Limitations**

507 We end this section with a discussion of several strengths and limitations of the
508 integrative approach. Here, we focus on those features that have not been mentioned
509 earlier in this article but which show how the phenomenon-based and data-driven
510 approaches enrich each other.

511 *Strengths*

512 **Data-inspired Parameter Values.** A first strength of the integrated approach
513 is the direct testability of a model to data and, critically, the ability to estimate a
514 realistic set of parameters from those data. The importance of a realistic set of
515 parameters for evaluating a model can not be understated: Parameters govern the
516 behavior of the model. The ability to estimate the parameters of a model is therefore
517 useful, as it provides us with a realistic parameter set that can be used to predict
518 phenomena that can be realistically expected to occur – the hallmark of strong scientific
519 theories. Additionally, it allows for an empirical check of whether the parameter values
520 satisfy theoretical constraints, rather than merely assuming that the parameters do.
521 This type of analysis corresponds to what Pitt et al. (2006) referred to as local analysis,
522 which focuses on evaluating a model’s behavior at specific parameter values. In
523 contrast, global analysis explores the full range of qualitative behaviors a model can
524 produce across its entire parameter space. Critically, Pitt et al. (2006) argue that while
525 global analysis offers valuable theoretical insights into the range of behaviors a model
526 can produce, it should not stand alone. Instead, it must be complemented by local
527 analysis based on parameter estimates derived from actual data. We fully agree with
528 this position.

529 Parameter recoverability and identifiability can furthermore serve as critical
530 diagnostics for whether theoretical constructs in a model can be empirically grounded.
531 When different combinations of parameter values generate the same phenomena, or

532 when parameter estimates fluctuate erratically across similar datasets, the
533 correspondence between parameters and psychological constructs becomes vague,
534 despite mathematical formalization. Pushing this argument a bit further, well-designed
535 models may use their parameters as indirect measurements of psychological constructs,
536 extending the reach of psychological science beyond directly observable behavior (e.g.,
537 Brown & Heathcote, 2008; Heathcote et al., in preparation; Kahneman & Tversky,
538 1979; Yu et al., 2023, 2025). When a model includes parameters which represent
539 mechanisms like learning rates, generalization gradients, or attentional focus, estimating
540 these parameters allows researchers to investigate how specific psychological processes
541 vary across individuals, developmental stages, or experimental conditions. These
542 parameter-based insights provide a deeper understanding of psychological functioning
543 than would be possible through analysis of behavioral measures alone.

544 **Formal Theory as Data Model.** The integrative approach questions the
545 utility of the dichotomy between theoretical and *data models*. Based on Suppes’
546 proposition that one takes many auxiliary assumptions when moving from theory to
547 data (Suppes, 1962), data models have been defined as general representations of data,
548 going from a simple mean to statistical models such as linear regressions and
549 autoregressive models (Haslbeck et al., 2022). However, we believe that the distinction
550 between theoretical and data model is, to a degree, arbitrary.¹ and data models, that is
551 distinguishing between the two types of models on the basis of their purpose (i.e.,
552 communicate a theory vs fitting data). However, this distinction has nothing to do with
553 the mathematical structure of the models and as such can not be used to classify them.
554 There does not seem to exist a set of necessary nor sufficient conditions that allows for
555 the distinction between one or the other. For example, one may argue that, according
556 to the current definition, data models are those we can apply or fit to data directly.
557 However, if we apply this criterion, each theoretical model that can be fit to data should
558 be classified as a data model. Contrary to this conclusion, many theoretical models in
559 the cognitive and affective sciences can be fit to data (e.g., Brown & Heathcote, 2008;

¹ Note that there one may maintain a teleological distinction between theoretical

560 Dalege et al., 2016; Laibson, 1997; Loossens et al., 2020; Verdonck & Tuerlinckx, 2014;
561 Yu et al., 2023) with only few exceptions to this rule (e.g., Robinaugh et al., 2024).
562 Would this mean that only the latter qualify as theoretical models? We don't think so.

563 *Limitations*

564 **Combined Inference.** While we have focused on how the data-driven and
565 phenomenon-based approaches compliment each other, we have not discussed what
566 happens when the two methods fail to fully solve the limitations of each approach
567 separately. For example, model fit may prefer models that do not necessarily reproduce
568 all phenomena the researcher is interested in, leaving it to the researcher to determine
569 which of the two pieces of information is more crucial. Additionally, there is some
570 subjectivity in the metric that is chosen to assess model fit, a type of subjectivity that
571 is also present in operationalizing the phenomena to be reproduced according to the
572 phenomenon-based approach. Unfortunately, this subjectivity is not resolved in the
573 integrative approach. In other words, while the integrative approach provides the
574 researcher with a broader set of methods, it is not a catch-all remedy against problems
575 in inferential reasoning.

576 Instead, one should view the use of the integrative approach as providing more
577 direct feedback towards the modeler through both the data and the phenomena. To
578 leverage on the strength of both approaches, one therefore has to consider this valuable
579 feedback and integrate it into future iterations of the model.

580 **Estimability.** Like the data-driven approach, the integrative approach only
581 applies to the testing of models that are estimable. Building a model that is both
582 theoretically valid and parsimonious enough to be fit to data is no easy task, but we
583 believe a necessary one. In cases where this is infeasible, however, one can use the
584 methods of the phenomenon-based approach to perform theory building. However, we
585 caution readers to closely consider the limitations of this approach with regard to
586 theory testing, asking them to remain cautious when claiming validity of a model.

Discussion

587

588 In this paper, we have provided a critical evaluation of the strengths and
589 limitations of the phenomenon-based approach to theory testing. Additionally, we have
590 argued that the inclusion of data in validating a model can mitigate the issues
591 associated to the use of the phenomenon-based approach alone. Instead, we have argued
592 for combining the phenomenon-based and data-driven approaches to theory validation
593 into an integrative approach, which we believe provides the most promising avenue for
594 future research.

595 It is important to reiterate that we do not disregard the importance of
596 phenomena when building and testing a theory. The phenomenon-based approach
597 represents a fresh breath of air in an otherwise data-driven field, allowing for the initial
598 validation of theories that do not necessarily allow for estimation. However, we believe
599 that there is a crucial piece missing from the methodology proposed by Borsboom et al.
600 (2021) with regard to testing theory, and that data represents this missing piece.
601 Similarly, we believe that purely data-driven methods have their own place in the
602 literature, but that a focus on only data may equally harm theory validation attempts.
603 Given the complimentary nature of the strengths and weaknesses of both approaches
604 when used in isolation, a combination of these approaches may yield the best results.

605 We close by returning to the previously mentioned discovery of Neptune, as it
606 provides a beautiful historical example of the topics we have discussed in this article.
607 As mentioned before, the discovery of Neptune is sometimes cited as an example of the
608 phenomenon-based approach in action (Haslbeck et al., 2022). However, it is easy to
609 overlook the decades of observation, prediction, and theory testing that led up to this
610 discovery and which, when taken together, are indicative of the integrative approach to
611 theory validation.

612 Uranus displayed “residual perturbations”, that is deviations from the elliptic
613 orbit predicted by the law of gravitation, which became apparent due to a systematic
614 comparison of decades of carefully gathered observations (data) with the orbits
615 (phenomena) implied by the theory of gravity (theory). To explain why these deviations

616 were observed, astronomers put forward a few theories (Sheehan et al., 2021). First,
617 some astronomers questioned the quality of the data, stating that old observations
618 might be unreliable, while others argued that the theory of gravity was incorrect and
619 put forward alternative theories that accommodated the phenomenon of residual
620 perturbations (e.g., through selective attraction, see Sheehan et al., 2021). Yet other
621 astronomers were less quick to seek to modify the theory of gravitation, writing that
622 “the law of gravitation was too firmly established to be doubted till every other
623 hypothesis had failed” (quote of John Adams taken from Bamford, 1996, p. 216).

624 Accepting both the validity of the model and the data, astronomers put forward
625 the possibility that an undiscovered planet caused the deviations in the orbit of Uranus,
626 indicating a problem with our understanding of the solar system at the time.
627 Importantly, this hypothesis was supported through *fitting various orbits to the*
628 *observed data*. For example, Urbain Leverrier attempted to explain the irregularities in
629 the orbit of Uranus by correcting the relevant observations for the gravitational pull of
630 other, known planets, concluding that “no ellipse would satisfy the range of
631 observations, ancient and modern, even on the most favorable distribution of errors in
632 them” (Bamford, 1996, p. 215). Following this result, Leverrier then approximated the
633 location of Neptune by assuming the existence of another planet so that, when
634 correcting the observations in Uranus for the pull of such a planet, the errors in the
635 orbit of Uranus would no longer be systematic (Grant, 1852).

636 It is clear that the discovery of Neptune arose due to both a top-down stream,
637 where theoretical predictions are derived from a theory mathematically, and a
638 bottom-up stream, where data is used to judge the validity of these predictions. It is
639 difficult to see how any stream in isolation would have led to the discovery of Neptune.
640 Indeed citing Sheehan et al. (2021): “The discovery of Neptune was a story of two
641 parts. The first one was theoretical, belonging to what was then called ‘mathematical
642 astronomy’ (...) The other part was empirical, and consisted in preparing the ground,
643 such as charting the skies to facilitate location of transiting bodies, and the actual
644 observations necessary for any genuine discovery of astronomical objects.” (p. 189).

References

645

646 Ariens, S., Ceulemans, E., & Adolf, J. K. (2020). Time series analysis of intensive
647 longitudinal data in psychosomatic research: A methodological overview. *Journal*
648 *of Psychosomatic Research*, *137*, 110191.

649 <https://doi.org/10.1016/j.jpsychores.2020.110191>

650 Ariens, S., Adolf, J. K., & Ceulemans, E. (2023). One does not simply correct for serial
651 dependence. *Structural Equation Modeling: A Multidisciplinary Journal*, *30*(5),
652 807–821.

653 Bamford, G. (1996). Popper and his commentators on the discovery of Neptune: A close
654 shave for the law of gravitation? *Studies in History and Philosophy of Science*
655 *Part A*, *27*(2), 207–232.

656 Bastida-Castillo, A., Gómez-Carmona, C. D., De La Cruz Sánchez, E., &

657 Pino-Ortega, J. (2019). Comparing accuracy between global positioning systems
658 and ultra-wideband-based position tracking systems used for tactical analyses in
659 soccer. *European Journal of Sport Science*, *19*, 1157–1165.

660 <https://doi.org/10.1080/17461391.2019.1584248>

661 Bates, S., Hastie, T., & Tibshirani, R. (2023). Cross-validation: What does it estimate
662 and how well does it do it? *Journal of the American Statistical Association*, *119*,
663 1434–1445. <https://doi.org/10.1080/01621459.2023.2197686>

664 Bogen, J., & Woodward, J. (1988). Saving the phenomena. *Philosophical Review*, *97*(3),
665 303–352. <https://doi.org/10.2307/2185445>

666 Borsboom, D., & Haslbeck, J. M. B. (2024). Integrating intra- and interindividual
667 phenomena in psychological theories. *Multivariate Behavioral Research*, *59*,
668 1290–1309. <https://doi.org/10.1080/00273171.2024.2336178>

669 Borsboom, D., van der Maas, H. L. J., Dalege, J., Kievit, R. A., & Haig, B. D. (2021).

670 Theory construction methodology: A practical framework for building theories in
671 psychology. *Perspectives on Psychological Science*, *16*, 756–766.

672 <https://doi.org/10.1177/1745691620969647>

- 673 Bramson, A., Grim, P., Singer, D. J., Berger, W. J., Sack, G., Fisher, S., Flocken, C., &
674 Holman, B. (2017). Understanding polarization: Meanings, measures, and model
675 evaluation. *Philosophy of Science*, *84*, 115–159. <https://doi.org/10.1086/688938>
- 676 Brown, S. D., & Heathcote, A. (2008). The simplest complete model of choice response
677 time: Linear ballistic accumulation. *Cognitive Psychology*, *57*, 153–178.
678 <https://doi.org/10.1016/j.cogpsych.2007.12.002>
- 679 Campanella, M., Hoogendoorn, S., & Daamen, W. (2014). The Nomad model: Theory,
680 developments and applications. *Transportation Research Procedia*, *2*, 462–467.
681 <https://doi.org/10.1016/j.trpro.2014.09.061>
- 682 Cloos, L., Ceulemans, E., & Kuppens, P. (2023). Development, validation, and
683 comparison of self-report measures for positive and negative affect in intensive
684 longitudinal research. *Psychological Assessment*, *35*, 189–204.
685 <https://doi.org/10.1037/pas0001200>
- 686 Dalege, J., Borsboom, D., van Harreveld, F., van den Berg, H., M., C., &
687 van der Maas, H. L. J. (2016). Toward a formalized account of attitudes: The
688 Causal Attitude Network (CAN) model. *Psychological Review*, *123*, 2–22.
689 <https://doi.org/10.1037/a0039802>
- 690 Darwin, C. (1859). *On the origin of species by means of natural selection, or, the*
691 *preservation of favoured races in the struggle for life*. John Murray.
- 692 Dawson, M. E., Schell, A. M., & Filion, D. L. (2000). The electrodermal system. In
693 J. T. Cacioppo, L. G. Tassinary, & G. G. Berntson (Eds.), *Handbook of*
694 *psychophysiology* (2nd ed., pp. 200–223). Cambridge University Press.
- 695 Dutilh, G., Wagenmakers, E.-J., Visser, I., & van der Maas, H. L. J. (2011). A phase
696 transition model for the speed-accuracy trade-off in response time. *Cognitive*
697 *Science*, *35*, 211–250. <https://doi.org/10.1111/j.1551-6709.2010.01147.x>
- 698 Dyson, F. (2004). A meeting with Enrico Fermi. *Nature*, *427*, 297.
699 <https://doi.org/10.1038/427297a>

- 700 Eronen, M. I., & Bringmann, L. F. (2021). The theory crisis in psychology: How to
701 move forward. *Perspectives on Psychological Science, 16*, 779–788.
702 <https://doi.org/10.1177/1745691620970586>
- 703 Eronen, M. I., & Romeijn, J.-W. (2020). Philosophy of science and the formalization of
704 psychological theory. *Theory & Psychology, 30*, 786–799.
705 <https://doi.org/10.1177/0959354320969876>
- 706 Farrell, S., & Lewandowsky, S. (2010). Computational models as aids to better
707 reasoning in psychology. *Current Directions in Psychological Science, 19*,
708 329–335. <https://doi.org/10.1177/0963721410386677>
- 709 Festinger, L. (1957). *A theory of cognitive dissonance*. Stanford University Press.
- 710 Gamble, A. S. D., Bigg, J. L., Pignanelli, C., Nyman, D. L. E., Burr, J. F., &
711 Spriet, L. L. (2023). Reliability and validity of an indoor local positioning system
712 for measuring external load in ice hockey players. *European Journal of Sport
713 Science, 23*, 311–318. <https://doi.org/10.1080/17461391.2022.2032371>
- 714 Gelman, A., Hwang, J., & Vehtari, A. (2014). Understanding predictive information
715 criteria for Bayesian models. *Statistics and Computing, 24*(6), 997–1016.
716 <https://doi.org/10.1007/s11222-013-9416-2>
- 717 Gershman, S. J. (2019). How to never be wrong. *Psychonomic Bulletin and Review, 26*,
718 13–28. <https://doi.org/10.3758/s13423-018-1488-8>
- 719 Grant, R. (1852). *History of physical astronomy from the earliest ages to the middle of
720 the nineteenth century*. Bohn.
- 721 Guest, O., & Martin, A. E. (2021). How computational modeling can force theory
722 building in psychological science. *Perspectives on Psychological Science, 16*,
723 789–802. <https://doi.org/10.1177/1745691620970585>
- 724 Haig, B. D. (2005). An abductive theory of scientific method. *Psychological Methods,*
725 *10*, 371–388. <https://doi.org/10.1037/1082-989X.10.4.371>
- 726 Haslbeck, J. M. B., Ryan, O., Robinaugh, D. J., Waldorp, L. J., & Borsboom, D.
727 (2022). Modeling psychopathology: From data models to formal theories.
728 *Psychological Methods, 27*, 930–957. <https://doi.org/10.1037/met0000303>

- 729 Hastie, T., Tibshirani, R., & Friedman, J. (2011). *The elements of statistical learning:*
730 *Data mining, inference, and prediction* (2nd ed.). Springer.
- 731 Heathcote, A., Brown, S., & Mewhort, D. J. K. (2000). The power law repealed: The
732 case for an exponential law of practice. *Psychonomic Bulletin & Review*, *7*,
733 185–207. <https://doi.org/10.3758/BF03212979>
- 734 Heathcote, A., Vanhasbroeck, N., Anderson, A., Blanken, T., Borsboom, D., &
735 Matzke, D. (in preparation). A psychological modeling framework for individual
736 pedestrian decisions in complex environments.
- 737 Hedrick, T. L. (2008). Software techniques for two- and three-dimensional kinematic
738 measurements of biological and biomimetic systems. *Bioinspiration &*
739 *Biomimetics*, *3*, 034001. <https://doi.org/10.1088/1748-3182/3/3/034001>
- 740 Jangraw, D. C., Keren, H., Sun, H., Bedder, R. L., Rutledge, R. B., Pereira, F.,
741 Thomas, A. G., Pine, D. S., Zheng, C., Nielson, D. M., & Stringaris, A. (2023).
742 A highly replicable decline in mood during rest and simple tasks. *Nature Human*
743 *Behaviour*, *7*, 596–610. <https://doi.org/10.1038/s41562-023-01519-7>
- 744 Kahneman, D., & Tversky, A. (1979). Prospect theory: An analysis of decision under
745 risk. *Econometrica*, *47*, 263–292. <https://doi.org/10.2307/1914185>
- 746 Kellen, D., Davis-Stober, C. P., Dunn, J. C., & Kalish, M. L. (2021). The problem of
747 coordination and the pursuit of structural constraints in psychology. *Perspectives*
748 *on Psychological Science*, *16*(4), 767–778.
749 <https://doi.org/10.1177/1745691620974771>
- 750 Krosnick, J. A., & Fabrigar, L. R. (1997). Designing rating scales for effective
751 measurement in surveys. In L. Lyberg, P. Biemer, M. Collins, E. De Leeuw,
752 C. Dippo, N. Schwarz, & D. Trewin (Eds.), *Survey measurement and process*
753 *quality*. John Wiley & Sons, Inc.
- 754 Krosnick, J. A., & Presser, A. (2010). Question and questionnaire design. In
755 P. V. Marsden & J. D. Wright (Eds.), *Handbook of survey research* (2nd ed.).
756 Emerald.

- 757 Laibson, D. (1997). Golden eggs and hyperbolic discounting. *The Quarterly Journal of*
758 *Economics*, *2*, 443–477.
- 759 Lakatos, I. (1978). *The methodology of scientific research programmes*. Cambridge
760 University Press.
- 761 Logan, G. D. (1992). Shapes of reaction-time distributions and shapes of learning
762 curves: A test of the instance theory of automaticity. *Journal of Experimental*
763 *Psychology: Learning, Memory, and Cognition*, *18*, 883–914.
764 <https://doi.org/10.1037//0278-7393.18.5.883>
- 765 Loossens, T., Mestdagh, M., Dejonckheere, E., Kuppens, P., Tuerlinckx, F., &
766 Verdonck, S. (2020). The Affective Ising Model: A computational account of
767 human affect dynamics. *PLoS Computational Biology*, *16*, e1007860.
768 <https://doi.org/10.1371/journal.pcbi.1007860>
- 769 Mauss, I. B., & Robinson, M. D. (2009). Measures of emotion: A review. *Cognition and*
770 *Emotion*, *23*, 209–237. <https://doi.org/10.1080/02699930802204677>
- 771 Mestdagh, M., Verdonck, S., Meers, K., Loossens, T., & Tuerlinckx, F. (2019). Prepaid
772 parameter estimation without likelihoods. *PLoS Computational Biology*, *15*,
773 e1007181. <https://doi.org/10.1371/journal.pcbi.1007181>
- 774 Navarro, D. J. (2019). Between the devil and the deep blue sea: Tensions between
775 scientific judgement and statistical model selection. *Computational Brain &*
776 *Behavior*, *2*, 28–34. <https://doi.org/10.1007/s42113-018-0019-z>
- 777 Navarro, D. J., Pitt, M. A., & Myung, I. J. (2004). Assessing the distinguishability of
778 models and the informativeness of data. *Cognitive Psychology*, *49*, 47–84.
779 <https://doi.org/10.1016/j.cogpsych.2003.11.001>
- 780 Newell, A., & Rosenbloom, P. S. (1980). Mechanisms of skill acquisition and the law of
781 practice [Carnegie Mellon University].
- 782 Oude Maatman, F. J. W. (preprint). *Psychology's theory crisis, and why formal*
783 *modelling cannot solve it* [Retrieved from
784 https://osf.io/preprints/psyarxiv/puqvs_v1].
- 785 Pathria, R. K., & Beale, P. D. (2022). *Statistical mechanics* (4th ed.). Academic Press.

- 786 Piironen, J., & Vehtari, A. (2017). Comparison of Bayesian predictive methods for
787 model selection. *Statistics and Computing*, *27*(3), 711–735.
788 <https://doi.org/10.1007/s11222-016-9649-y>
- 789 Pitt, M. A., Kim, W., Navarro, D. J., & Myung, J. I. (2006). Global model analysis by
790 parameter space partitioning. *Psychological Review*, *113*, 57–83.
791 <https://doi.org/10.1037/0033-295X.113.1.57>
- 792 Popper, K. (1959). *The logic of scientific discovery*. Routledge.
- 793 Radev, S. T., Mertens, U. K., Voss, A., Ardizzone, L., & Köthe, U. (2020). BayesFlow:
794 Learning complex stochastic models with invertible neural networks. *IEEE*
795 *transactions on neural networks and learning systems*, *33*(4), 1452–1466.
- 796 Radev, S. T., Schmitt, M., Schumacher, L., Else Müller, L., Pratz, V., Schälte, Y.,
797 Köthe, U., & Bürkner, P.-C. (2023). BayesFlow: Amortized Bayesian workflows
798 with neural networks. *Journal of Open Source Software*, *8*(89), 5702.
- 799 Ratcliff, R., & Rouder, J. (1998). Modeling response times for two-choice decisions.
800 *Psychological Science*, *9*, 347–356. <https://doi.org/10.1111/1467-9280.00067>
- 801 Read, G. L. (2017). Facial electromyography (EMG). In J. Matthes, C. S. Davis, &
802 R. F. Potter (Eds.), *The international encyclopedia of communication research*
803 *methods* (pp. 1–10). John Wiley & Sons, Ltd.
804 <https://doi.org/10.1002/9781118901731.iecrm0100>
- 805 Regenwetter, M., & Robinson, M. M. (2017). The construct–behavior gap in behavioral
806 decision research: A challenge beyond replicability. *Psychological Review*, *124*(5),
807 533–550. <https://doi.org/10.1037/rev0000067>
- 808 Revol, J., Ariens, S., Lafit, G., Adolf, J. K., & Ceulemans, E. (in press).
809 Episode-contingent experience-sampling designs for accurate estimates of
810 autoregressive dynamics [Accepted for publication in *Psychological Methods*].
- 811 Reynolds, C. W. (1987). Flocks, herds, and schools: A distributed behavioral model.
812 *Computer Graphics*, *21*, 25–34.
- 813 Roberts, D. R., Bahn, V., Ciuti, S., Boyce, M. S., Elith, J., Guillerá-Arroita, G.,
814 Hauenstein, S., LahozMonfort, J. J., Schröder, B., Thuiller, W., Warton, D. I.,

- 815 Wintle, B. A., Hartig, F., & Dormann, C. F. (2017). Cross-validation strategies
816 for data with temporal, spatial, hierarchical, or phylogenetic structure.
817 *Ecography*, *40*, 913–929. <https://doi.org/10.1111/ecog.02881>
- 818 Roberts, S., & Pashler, H. (2000). How persuasive is a good fit? A comment on theory
819 testing. *Psychological Review*, *107*, 358–367.
820 <https://doi.org/10.1037/0033-295X.107.2.358>
- 821 Robinaugh, D. J., Haslbeck, J. M. B., Waldorp, L. J., Kossakowski, J. J., Fried, E. I.,
822 Millner, A. J., McNally, R. J., R., O., de Ron, J., van der Maas, H. L. J., van
823 Nes, E. H., Scheffer, M., Kendler, K. S., & Borsboom, D. (2024). Advancing the
824 network theory of mental disorders: A computational model of panic disorder.
825 *Psychological Review*, *131*, 1482–1508. <https://doi.org/10.1037/rev0000515>
- 826 Robinson, M. M., Williams, J. R., Wixted, J. T., & Brady, T. F. (2025). Zooming in on
827 what counts as core and auxiliary: A case study on recognition models of visual
828 working memory. *Psychonomic Bulletin & Review*, *32*(2), 547–569.
829 <https://doi.org/10.3758/s13423-024-02562-9>
- 830 Russel, J. A., Weiss, A., & Mendelsohn, G. A. (1989). Affect grid: A single-item scale of
831 pleasure and arousal. *Journal of Personality and Social Psychology*, *57*, 493–502.
832 <https://doi.org/10.1037/0022-3514.57.3.493>
- 833 Rutledge, R. B., Skandali, N., Dayan, P., & Dolan, R. J. (2014). A computational and
834 neural model of momentary subjective well-being. *Proceedings of the National
835 Academy of Sciences*, *111*, 12252–12257.
836 <https://doi.org/10.1073/pnas.1407535111>
- 837 Schroeder, D. V. (2021). *An introduction to thermal physics*. Oxford University Press.
- 838 Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, *6*,
839 461–464.
- 840 Sheehan, W., Bell, T. E., Kennett, C., & Smith, R. W. (2021). Neptune: From grand
841 discovery to a world revealed. <https://doi.org/10.1007/978-3-030-54218-4>

- 842 Smaldino, P. E. (2017). Models are stupid, and we need more of them. In
843 R. R. Vallacher, S. J. Read, & A. Nowak (Eds.), *Computational social psychology*
844 (pp. 311–331). Routledge.
- 845 Smaldino, P. E. (2020). How to translate a verbal theory into a formal model. *Social*
846 *Psychology*, *51*, 207–218. <https://doi.org/10.1027/1864-9335/a000425>
- 847 Spanos, A. (2024). How the post-data severity converts testing results into evidence for
848 or against pertinent inferential claims. *Entropy*, *26*(1), 95.
849 <https://doi.org/10.3390/e26010095>
- 850 Stevens, S. S. (1946). On the theory of scales of measurement. *Science*, *103*, 677–680.
851 <https://doi.org/10.1126/science.103.2684.677>
- 852 Stone, M. (1974). Cross-validatory choice and assessment of statistical predictions.
853 *Journal of the Royal Statistical Society, Series B (Methodological)*, *38*, 111–147.
- 854 Strogatz, S. H. (2018). *Nonlinear dynamics and chaos: With applications to physics,*
855 *biology, chemistry, and engineering* (2nd ed.). CRC Press.
- 856 Suppes, P. (1962). Models of data. In *Logic, methodology and the philosophy of science:*
857 *Proceedings of the 1960 international congress* (pp. 252–261, Vol. 44). Stanford
858 University Press.
- 859 Susko, E., & Roger, A. J. (2020). On the use of information criteria for model selection
860 in phylogenetics (N. Saitou, Ed.). *Molecular Biology and Evolution*, *37*(2),
861 549–562. <https://doi.org/10.1093/molbev/msz228>
- 862 Tashman, L. J. (2000). Out-of-sample tests of forecasting accuracy: An analysis and
863 review. *International Journal of Forecasting*, *16*, 437–450.
- 864 van Dongen, N., van Bork, R., Finnemann, A., Haslbeck, J. M. B.,
865 van der Maas, H. L. J., Robinaugh, D. J., de Ron, J., Sprenger, J., &
866 Borsboom, D. (2025). Productive explanation: A framework for evaluating
867 explanations in psychological science. *Psychological Review*, *132*, 311–329.
868 <https://doi.org/10.1037/rev0000479>

- 869 van Rooij, I., & Baggio, G. (2021). Theory before the test: How to build
870 high-verisimilitude explanatory theories in psychological science. *Perspectives on*
871 *Psychological Science*, *16*, 682–697. <https://doi.org/10.1177/1745691620970604>
- 872 Vanhasbroeck, N., Loossens, T., Anarat, N., Ariens, S., Vanpaemel, W., Moors, A., &
873 Tuerlinckx, F. (2022). Stimulus-driven affective change: Evaluating
874 computational models of affect dynamics in conjunction with input. *Affective*
875 *Science*, *3*, 559–576. <https://doi.org/10.1007/s42761-022-00118-5>
- 876 Vanhasbroeck, N., Loossens, T., & Tuerlinckx, F. (2024). Two peas in a pod:
877 Discounting models as a special case of the VARMAX. *Journal of Mathematical*
878 *Psychology*, *120-121*, 102856. <https://doi.org/10.1016/j.jmp.2024.102856>
- 879 Vanpaemel, W. (2020). Strong theory testing using the prior predictive and the data
880 prior. *Psychological Review*, *127*, 136–145. <https://doi.org/10.1037/rev0000167>
- 881 Verdonck, S., & Tuerlinckx, F. (2014). The Ising Decision Maker: A binary stochastic
882 network for choice response time. *Psychological Review*, *121*, 422–462.
883 <https://doi.org/10.1037/a0037012>
- 884 Viscardi, C., Lachi, A., & Baccini, M. (2025). Discrete-time compartmental models with
885 partially observed data: A comparison among frequentist and Bayesian
886 approaches for addressing likelihood intractability. *Epidemiologic Methods*,
887 *14*(1), 20240032. <https://doi.org/10.1515/em-2024-0032>
- 888 Wagenmakers, E.-J., Ratcliff, R., Gomez, P., & Iverson, G. J. (2004). Assessing model
889 mimicry using the parametric bootstrap. *Journal of Mathematical Psychology*,
890 *48*, 28–50. <https://doi.org/10.1016/j.jmp.2003.11.004>
- 891 Wehrens, R., Putter, H., & Buydens, L. M. C. (2000). The bootstrap: A tutorial.
892 *Chemometrics and Intelligent Laboratory Systems*, *54*, 35–52.
- 893 Westfall, J., & Yarkoni, T. (2016). Statistically controlling for confounding constructs is
894 harder than you think (U. S. Tran, Ed.). *PLOS ONE*, *11*(3), e0152719.
895 <https://doi.org/10.1371/journal.pone.0152719>
- 896 Woodward, J. (1989). Data and phenomena. *Synthese*, *79*, 393–472.

- 897 Yarkoni, T., & Westfall, J. (2017). Choosing prediction over explanation in psychology:
898 Lessons from machine learning. *Perspectives on Psychological Science*, *12*(6),
899 1100–1122. <https://doi.org/10.1177/1745691617693393>
- 900 Yu, K., Tuerlinckx, F., Vanpaemel, W., & Zaman, J. (2023). Humans display
901 interindividual differences in the latent mechanisms underlying fear
902 generalization behaviour. *Communications Psychology*, *1*(1), 5.
903 <https://doi.org/10.1038/s44271-023-00005-0>
- 904 Yu, K., Vanpaemel, W., Tuerlinckx, F., & Zaman, J. (2025). The probabilistic and
905 dynamic nature of perception in human generalization behavior. *iScience*, *28*(4),
906 112228. <https://doi.org/10.1016/j.isci.2025.112228>

907 **Acknowledgements:** The authors would like to thank Prof. Denny Borsboom
908 and Kyra Evers for critical and constructive discussions on the contents of this paper.
909 We furthermore thank the reviewers for their valuable comments.

910 **Conflict of interest:** The authors report no conflict of interest.

911 **Contribution:** All authors contributed equally.

912 **Funding:** N.V. is supported by funding from the European Research Council
913 (ERC) under the European Union's Horizon 2020 Excellent Science program (Grant
914 agreement No. 101053880). K.Y. is supported by a grant from the Fund for Scientific
915 Research - Flanders (FWO; G079520N) and in part by the Research Fund of KU
916 Leuven (C14/23/062). S.A. is supported by a grant from FWO (1278525N).