



# Bayesian statistical approaches to evaluate cognitive models

Jeffrey Annis\* and Thomas J. Palmeri

Cognitive models aim to explain complex human behavior in terms of hypothesized mechanisms of the mind. These mechanisms can be formalized in terms of mathematical structures containing parameters that are theoretically meaningful. For example, in the case of perceptual decision making, model parameters might correspond to theoretical constructs like response bias, evidence quality, response caution, and the like. Formal cognitive models go beyond verbal models in that cognitive mechanisms are instantiated in terms of mathematics and they go beyond statistical models in that cognitive model parameters are psychologically interpretable. We explore three key elements used to formally evaluate cognitive models: parameter estimation, model prediction, and model selection. We compare and contrast traditional approaches with Bayesian statistical approaches to performing each of these three elements. Traditional approaches rely on an array of seemingly *ad hoc* techniques, whereas Bayesian statistical approaches rely on a single, principled, internally consistent system. We illustrate the Bayesian statistical approach to evaluating cognitive models using a running example of the Linear Ballistic Accumulator model of decision making (Brown and Heathcote 2008). © 2017 Wiley Periodicals, Inc.

How to cite this article:

*WIREs Cogn Sci* 2017, e1458. doi: 10.1002/wcs.1458

## INTRODUCTION

Cognitive models aim to explain behavior by positing mechanisms that underlie perception, memory, decision making, and other fundamental aspects of cognition. Formal cognitive models instantiate these hypothesized mechanisms in terms of mathematics, computations, and simulations, and models are fitted, evaluated, and compared based on tools and techniques from statistics. Formal cognitive models go beyond verbal theories in that they are precisely defined and make explicit predictions. They go beyond statistical models that describe patterns of behavior in that they attempt to explain patterns of behavior in terms of hypothesized mechanisms of the mind. While statistical models make parametric assumptions about observed data, such as linearity

or a particular distributional form, cognitive models make further assumptions about the underlying cognitive processes hypothesized to cause observed behavior, allowing differences across stimuli, conditions, groups, or individuals to be quantitatively characterized (e.g., Ref 1).

Such differences are often reflected in the values of free parameters in cognitive models. In the case of decision making, model parameters might describe the quality of evidence driving the decision process, how decisions are made more or less cautiously, and whether decisions are biased or not (e.g., Ref 2). In the case of object categorization, parameters might describe attention weights to relevant or irrelevant features, the relative strength of stored category representations, and biases to chose particular categorization responses over others (e.g., Refs 3,4). And in the case of memory recall, parameters might reflect how well items are stored in memory, how strongly stored items are integrated with current context, and when recall will terminate.<sup>5</sup>

Given a particular cognitive model and a set of observed data, a first step is often *parameter*

\*Correspondence to: jeff.annis@vanderbilt.edu

Department of Psychology, Vanderbilt University, Nashville, TN, USA

Conflict of interest: The authors have declared no conflicts of interest for this article.

1 *estimation*, fitting a model to observed data, estimat- 57  
2 ing model parameters that maximize the correspon- 58  
3 dence between model predictions and observed data. 59  
4 Once parameters are estimated, they might be related 60  
5 to brain measures,<sup>6</sup> used to explain individual 61  
6 differences,<sup>7,8</sup> or characterize psychiatric disorders.<sup>9</sup> 62  
7 Best-fitting model parameters can be used to generate 63  
8 *model predictions*. Does a model actually provide an 64  
9 adequate account of observed data and does it gener- 65  
10 ate novel predictions about yet-unseen data? 66

11 Often there are multiple competing models of 67  
12 any particular aspect of cognition, and researchers 68  
13 might do *model selection* based on how well or how 69  
14 poorly different models account for observed data. 70  
15 Parameter estimation and model prediction are often 71  
16 repeated for competing models to determine which 72  
17 model best explains the data. ‘Best’ usually means a 73  
18 model that provides the best fit to data with the low- 74  
19 est model complexity, a formal application of 75  
20 Occam’s razor.<sup>10</sup> This model comparison, which 76  
21 leads to model selection, might pit a more general 77  
22 model against a simpler special case of that general 78  
23 model, or pit one model against a different model 79  
24 making different mechanistic assumptions, or pit a 80  
25 mechanistic model against a ‘saturated’ statistical 81  
26 model to ask whether the mechanistic model leaves 82  
27 any observed variance unexplained. 83

28 This article focuses on these three key elements 84  
29 of formally testing a cognitive model: parameter esti- 85  
30 mation, model prediction, and model selection. To 86  
31 maintain a tight focus, we assume throughout that a 87  
32 researcher has a formal cognitive model already in 88  
33 hand, and is ready to fit that model to observed data 89  
34 and evaluate its quality. There are many interesting 90  
35 and important aspect of cognitive modeling that we 91  
36 will not address, such as a discussion of why formal 92  
37 cognitive models are developed and how they are 93  
38 used to advance theory (e.g., Refs 11–13), how cog- 94  
39 nitive models are initially developed or extended 95  
40 (e.g., Ref 14), a survey different kinds of cognitive 96  
41 models (e.g., Refs 1,15) and how they have been 97  
42 applied to particular aspects of cognition (e.g., 98  
43 Refs 2,16,17), or how cognitive models can be 99  
44 applied to neural data (e.g., Ref 18). This article first 100  
45 describes how parameter estimation, model predic- 101  
46 tion, and model selection have been traditionally 102  
47 carried out in the cognitive modeling literature.<sup>1,11</sup> The 103  
48 bulk of the article describes how *Bayesian statistics* 104  
49 can provide an alternative, coherent, and principled 105  
50 approach to these elements of modeling. 106

51 To be clear, Bayesian principles have made 107  
52 inroads into cognitive science and cognitive modeling 108  
53 in two different ways: One way is assuming that the 109  
54 mind and brain are inherently Bayesian, that human 110

learning and cognition follow the principles of Bayes- 57  
ian probabilistic inference (e.g., Refs 19–21 but see 58  
Ref 22). Such Bayes-in-the-head models are often 59  
referred to as *Bayesian cognitive models*. Instead, 60  
here we are using Bayesian statistics not as a princi- 61  
ple to explain cognition, but as a tool to evaluate 62  
models of cognition. Bayesian statistics can be used 63  
to evaluate cognitive models, whether those cognitive 64  
models are themselves Bayesian or not. In fact, while 65  
there are indeed many successful Bayesian cognitive 66  
models, the majority of cognitive models are non- 67  
Bayesian, in the sense that learning and cognition are 68  
not governed by Bayesian probabilistic inference. We 69  
can test quantitatively the adequacy of non-Bayesian 70  
models using Bayesian statistics. 71

72 The ideal intended audience for this article is 72  
73 someone who is familiar with cognitive models and 73  
74 how those models are traditionally fitted and evalu- 74  
75 ated, and wants to understand how Bayesian statis- 75  
76 tics might be used as a tool to do that fitting and 76  
77 evaluating. This could be because they want to use 77  
78 Bayesian statistics in their own work or merely want 78  
79 to better understand how other researchers are using 79  
80 Bayesian statistics in published articles. Some famil- 80  
81 iarity with Bayes rule would be beneficial, but we 81  
82 have tried to craft this article without requiring much 82  
83 background knowledge of Bayesian statistics. 83  
84 Because this article reflects the intersection of cog- 84  
85 nitive modeling and Bayesian statistics, we recommend 85  
86 readers new to both of these topics to read some of 86  
87 the more introductory textbooks and articles we ref- 87  
88 erence (e.g., Refs 1,11) for cognitive modeling, and 88  
89 (e.g., Refs 23–25) for Bayesian statistics. 89  
90

## 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105 106 107 108 109 110

## TRADITIONAL APPROACHES TO EVALUATE COGNITIVE MODELS

Traditionally, parameter estimation, model predic- 95  
tion, and model selection in cognitive modeling have 96  
been carried out using various optimization and sta- 97  
tistical techniques, each with their own advantages 98  
and disadvantages.<sup>1,11</sup> 99

100 Estimating best-fitting parameters of a cognitive 100  
101 model involves finding the values of model parame- 101  
102 ters that minimize or maximize some objective func- 102  
103 tion measuring how much the model predictions 103  
104 mismatch the observed data. The earliest published 104  
105 modeling work often minimized sum of squared error 105  
106 or root mean squared error (RMSE), or maximized 106  
107 correlation ( $r$ ) between model and data or percent 107  
108 variance accounted for by the model; while these 108  
109 objective functions may be fine for linear statistical 109  
110 and mechanistic models, they are often inappropriate 110

1 for nonlinear models, which is a property of many  
2 cognitive models.

3 A better approach, used in both nonlinear sta-  
4 tistical modeling and cognitive modeling, is maxi-  
5 mum likelihood estimation (MLE<sup>26</sup>). Formally, and  
6 abstractly, the goal is to find the parameter vector,  $\theta$ ,  
7 that maximizes the likelihood of the data,  $D$ , given  
8 the parameters, or maximizing  $p(D|\theta)$ ; practically  
9 speaking, you will often see work minimizing the  
10 negative of the log likelihood (minimizing the nega-  
11 tive only because many optimization routines are set  
12 up to minimize by default, and taking the log because  
13 calculating likelihoods often involve numerous multi-  
14 plications of small numbers that would lead to  
15 numeric underflow without taking logs). Of course,  
16 MLE assumes that this likelihood exists or can be  
17 approximated, which we will assume to be the case  
18 throughout the rest of this article.

19 Numerous methods exist for optimizing param-  
20 eters given some objective function. For the majority  
21 of cognitive models, direct application of calculus for  
22 optimization is so unwieldy as to render it largely  
23 impractical. For instance, a model with  $n$  parameters  
24 requires  $n$  partial derivatives, a solution to a system  
25 of  $n$  (often nonlinear) equations, and a series of tests  
26 to rule out local minima. Use of optimization algo-  
27 rithms that require the first or second derivatives can  
28 be somewhat less cumbersome, but they quickly  
29 become tedious to implement, computationally  
30 expensive, and error-prone as dimensionality  
31 increases (but see Ref 27).

32 So instead, optimization is often done using  
33 techniques such as hill-climbing, such as the well-  
34 known Simplex method;<sup>28</sup> Simplex is the default  
35 optimization algorithm in Matlab and Python when  
36 derivative are unspecified or unavailable. Given a  
37 starting point for the hill climb (often many different  
38 starting points are used), maximizing likelihood or  
39 minimizing RMSE yields a point estimate for the  
40 best-fitting parameter vector. This omits any measure  
41 of uncertainty we might want to know in the param-  
42 eter estimate; for example, if we obtain a maximum  
43 likelihood estimate of 1.03 for a given parameter, is  
44 there any chance that the parameter could be 1.04 or  
45 1.02 instead? To obtain estimates of uncertainty,  
46 other procedure must be carried out such as paramet-  
47 ric or nonparametric bootstrap sampling,<sup>29</sup> which  
48 gives an estimate of a confidence interval around a  
49 maximum likelihood point estimate.

50 Best-fitting parameters generated by MLE or  
51 other methods are often then used to generate model  
52 predictions. Especially for models with stochastic ele-  
53 ments, this may involve simulating the model hun-  
54 dreds or thousands of times, aggregating the

57 predictions in some way, and then comparing predic-  
58 tions against the observed data. Typical of a tradi-  
59 tional approach, and a potential disadvantage as  
60 well, is that these predictions are based on point esti-  
61 mates alone, ignoring any uncertainty of the parame-  
62 ter estimates.

63 Models are compared quantitatively according  
64 to how well they fit the observed data. But a mere  
65 comparison of quantitative fit, declaring the ‘winner’  
66 as the model with the largest likelihood or with the  
67 smallest RMSE, is nonsensical because doing so does  
68 not take into account the relative complexity or flexi-  
69 bility of the competing models. A more general  
70 model is mathematically guaranteed to fit at least as  
71 well, if not better, than a special case of that general  
72 model. The question is not whether the special case  
73 fits worse—it generally will—but whether it fits sig-  
74 nificantly worse than the more general model. Also  
75 in the case where one model is not a special case of  
76 another model, a model with more free parameters  
77 or with more flexibility in terms of the range of pre-  
78 dictions it can make will undoubtedly fit better than  
79 a model with fewer parameters and greater restric-  
80 tions. The question is not whether the more complex  
81 model fits better—it often will—but whether it fits  
82 better even when that model is appropriately penal-  
83 ized for its greater complexity. Traditional methods  
84 that penalize likelihood measures of fit based on  
85 complexity include the Bayesian Information Crite-  
86 rion (BIC;<sup>30</sup>) and the Akaike Information Criterion  
87 (AIC;<sup>31</sup>). While these methods are easy to  
88 implement—they simply involve adding a penalty  
89 term based on the number of free parameters—they  
90 largely ignore other important aspects of model com-  
91 plexity such as the functional form of the model and  
92 the size of space of possible predictions (e.g.,<sup>10</sup>).

## 93 INTRODUCING THE BAYESIAN 94 STATISTICAL APPROACH 95

96 In this article, we outline a cognitive modeling  
97 approach to parameter estimation, model prediction,  
98 and model selection that uses *Bayesian Statistics*. The  
99 Bayesian approach answers some of the limitations  
100 of the traditional approaches outlined above. Interest  
101 in Bayes has exploded over the past decade or more.  
102 As noted earlier, we do not discuss ideas of how  
103 human cognition might be based on Bayesian princi-  
104 ples.<sup>19,20,32</sup> We will also not discuss general aspects  
105 of Bayesian data analysis,<sup>23,33</sup> although there are cer-  
106 tainly parallels to what we discuss here. Rather, we  
107 will discuss how the Bayesian approach can be  
108 applied to cognitive models in order to perform  
109  
110

parameter estimation, model prediction, and model selection in a manner that is arguably both logically consistent and principled.<sup>34</sup> Therefore, this article is aimed at readers who are familiar with cognitive modeling, but are less familiar with how the Bayesian statistical approach can be applied to parameter estimation, prediction, or model selection in the context of cognitive modeling.

To state the obvious, Bayesian analysis is based on Bayes' rule. For example, we can use Bayes' rule to compute the full posterior probability distribution of the parameters given the data,  $p(\theta|D)$ :

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)} \quad (1)$$

The posterior probability of the parameters given the data is the product of the likelihood,  $p(D|\theta)$ , and the prior probability of the parameters  $p(\theta)$ , normalized by the probability of the data,  $p(D)$ . The likelihood of the data given the parameters is defined by the cognitive model in exactly the same way it would be specified for MLE. The prior distribution reflects a subjective belief about parameter values before observing the data.<sup>35</sup> These beliefs could be vague in the form of flat or relatively flat priors, they could be based on previous data used to estimate the parameters, or they could be based on values of the parameters that are meaningful by some theoretical, objective, or subjective criterion. Because Bayes' rule yields a posterior joint probability distribution of parameter values, it provides naturally a measure of parameter uncertainty, and hence an elegant solution to model prediction that takes into account that uncertainty. The Bayesian approach also allows a means for model comparison and model selection that takes into account model complexity in a natural, coherent, and comprehensive manner.

While traditional techniques of parameter estimation require a search of parameter space to minimize or maximize some objective function, Bayes' rule tells us directly the probability of the parameter values for a model given the observed data—at least in principle. As we will see, despite its apparent simplicity, even for models with only modest complexity, Bayes' rule cannot be solved analytically and requires computational estimation. One reason for the relatively recent explosion of interest in an idea first suggested in the 1700s,<sup>36</sup> and characterized and formalized mathematically many decades ago, has been the joint development of computational techniques and the availability of powerful computer hardware to make Bayesian analysis tractable. As

such, interest in the Bayesian approach to evaluating cognitive models has grown significantly over the past few years. Some applications include evaluating variants of signal detection theory (e.g., Refs 37,38), multinomial processing trees (e.g., Ref 39), individual differences (e.g., Refs 7,40–42), decision making (e.g., Refs 43,44), multidimensional scaling (e.g., Refs 37,42), choice response time (e.g., Refs 45–48), memory (e.g., Refs 49–52), and joint modeling of neural and behavioral data (e.g., Refs 6,53).

What follows is a combination of a review and tutorial of Bayesian approaches to evaluating cognitive models, with an added aim of pointing the reader to emerging new developments. We structure the remainder of the article around the three key elements of cognitive modeling outlined earlier—parameter estimation, model prediction, and model selection. For each, we describe the underlying concepts and mathematics, followed by the computational techniques used in practice. Throughout, we use a cognitive model called the *Linear Ballistic Accumulator* model (LBA;<sup>54</sup>) as a running example. We chose the LBA because it is a general model of choice response time that can be applied to a wide variety of tasks; we also recently published a companion paper<sup>45</sup> that outlines how to implement the LBA in a Bayesian statistical language called Stan.<sup>27</sup>

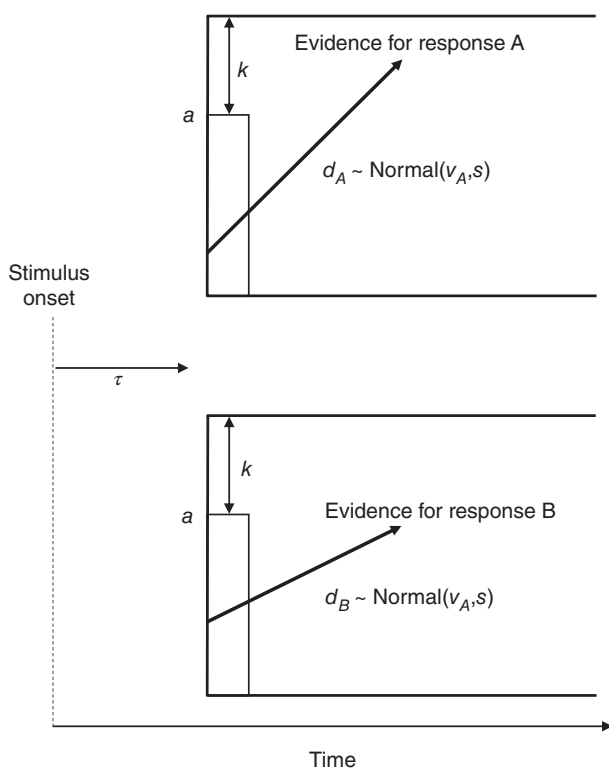
There are many excellent reviews and textbooks on Bayesian approaches to evaluating cognitive models and Bayesian statistics more generally. See Lee and Wagenmakers<sup>55</sup> for a practical introduction to Bayesian approaches to cognitive modeling, Rouder and Lu<sup>38</sup> for a more mathematically rich tutorial, Lee<sup>37,56</sup> for illustrative examples of how the Bayesian framework can be applied to cognitive models, and Shiffrin et al.<sup>57</sup> for a tutorial on Bayesian model selection. For textbooks on Bayesian statistics in general, see, for example, Kruschke<sup>23</sup> for the uninitiated and Gelman et al.<sup>25</sup> for the more heroic; Jackman<sup>24</sup> and Lynch<sup>58</sup> are texts specifically geared toward social scientists.

### Example: The LBA Model

LBA belongs to a class of models known as *sequential sampling models*<sup>2</sup> that describes an evidence accumulation process of decision making over time and can be used to predict both response probabilities and response times. LBA assumes that after a stimulus is presented, its representation is perceptually encoded and compared to some form of knowledge that will drive a decision. The time to complete this encoding process is given by the parameter  $\tau$  (this parameter also includes motor execution time to



produce an overt response). After encoding, evidence begins to accrue in independent *accumulators* that each correspond to one response alternative,  $i$ . Figure 1 shows an example of LBA with two accumulators, one for Response A and one for Response B. The rate at which evidence accumulates for response  $i$  is given by its corresponding drift rate,  $d_i$ . Drift rates are assumed to vary across trials, and are sampled from a normal distribution with mean  $v_i$  and standard deviation  $s$ . The starting point of evidence accumulation is also assumed to vary over trials, sampled from a uniform distribution,  $U(0, a)$ , where  $a$  is the maximum (note that in some articles that use LBA, the maximum is denoted by an



**FIGURE 1** | The Linear Ballistic Accumulator (LBA) model is an example of a formal cognitive model that predicts response probabilities and distributions of response times. LBA can be used to decompose response time and accuracy into core cognitive parameters: evidence accumulation, response caution, and perceptual encoding. LBA assumes that after the stimulus is perceptually encoded after time  $\tau$ , evidence toward each response alternative,  $i$ , accumulates with drift rate  $d_i$ . Drift rates across trials are sampled from normal distributions with mean  $v_i$  and standard deviation  $s$ . In our examples, we constrain mean drift rate for the Response B accumulator to be 1 minus mean drift rate for the Response A accumulator. The starting point of the evidence accumulation process on each is sampled from a uniform distribution between 0 and  $a$ . The response is determined by the first accumulator to reach threshold  $a + k$ .

uppercase  $A$ ; here we chose to instead use the lowercase  $a$  to ensure no confusion with the response alternative  $A$ ). The evidence accumulation terminates and a response is made when the first accumulator reaches its threshold  $a + k$ , where  $k$  is the *relative threshold*.

Compared to other accumulator models that assume noisy accumulation or lateral inhibition (e.g., Ref 59), LBA assumes a linear rise to threshold, which significantly simplifies its mathematical formulation. Brown and Heathcote<sup>54</sup> derived the likelihood function for the LBA. If  $D^A$  is a vector of observed (Data) response times for Response A, and  $D^B$  is a vector of observed (Data) response times for Response B, the likelihood for the combined observed data vector,  $D$ , is the product of the two likelihoods:

$$p(D|\theta) = p(D^A|\theta)p(D^B|\theta). \quad (2)$$

For the full mathematical description of the likelihoods we refer the reader to Brown and Heathcote.<sup>54</sup>

### Bayesian Parameter Estimation

One key difference between traditional and Bayesian approaches, is that Bayesian statistics treats data as well as unknown parameters as random variables. This allows us to write a joint distribution of the data,  $D$ , and the parameter(s),  $\theta$ :

$$p(D, \theta). \quad (3)$$

We can re-express the joint distribution via the definition of conditional probability:

$$p(D, \theta) = p(D|\theta)p(\theta) = p(\theta|D)p(D). \quad (4)$$

Rearranging, we obtain Bayes' rule:

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)} = \frac{p(D|\theta)p(\theta)}{\int p(D|\theta')p(\theta')d\theta'}, \quad (5)$$

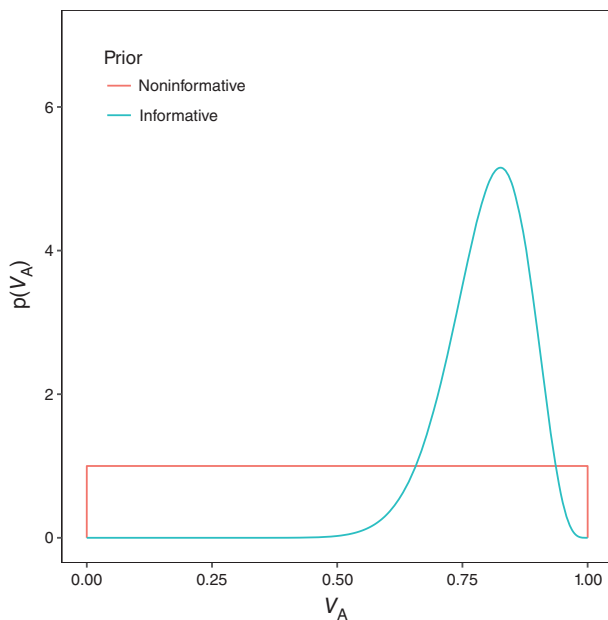
where  $p(\theta|D)$  is the posterior distribution,  $p(D|\theta)$  is the likelihood function (the same as that found in MLE),  $p(\theta)$  is the prior distribution, and  $p(D)$  is a normalization constant, which ensures the posterior integrates to 1, and is referred to as the *marginal likelihood*, or sometimes as the *evidence*. Here, the parameter vector  $\theta'$  has a prime within the integral to make clear that  $\theta'$  is different from  $\theta$ ; for simplicity, from now on, we will drop the superscript on parameters appearing within the integral.

Note that this form of Bayes' rule can be thought of as being implicitly conditioned on a given model. For example,  $p(\theta|\mathbf{D})$  under the LBA model will be different than  $p(\theta|\mathbf{D})$  under the diffusion model. To make the model we are working with explicit, sometimes Bayes' rule will be written including a model notation ( $\mathcal{M}$ ) explicitly:

$$p(\theta|\mathbf{D}, \mathcal{M}) = \frac{p(\mathbf{D}|\theta, \mathcal{M})p(\theta|\mathcal{M})}{p(\mathbf{D}|\mathcal{M})}. \quad (6)$$

This more complete formulation will be important later when we discuss model selection. For now, we will omit the explicit model notation.

Unlike traditional parameter estimation approaches that result in point estimates, the Bayesian approach results in a full posterior probability distribution of the parameters. Summary statistics like the mean, mode,<sup>a</sup> and standard deviation of a parameter can be computed, as can the correlations



**FIGURE 2** | Prior distributions represent our subjective *a priori* beliefs about parameter values. This figure illustrates two different prior distributions chosen for  $v_A$  in the LBA model, which in one common instantiation is constrained to fall between 0 and 1. Both priors are beta distributions with different shapes (controlled by the rate parameter  $\alpha$  and the shape parameter  $\beta$  of the beta). When  $\alpha = \beta$ , the beta distribution is mathematically equivalent to the uniform distribution, depicted in red. This prior is *noninformative* because it represents the belief that all values of  $v_A$  are equally likely. Depicted in blue is a beta distribution with  $\alpha = 5$  and  $\beta = 20$  and represents the prior belief that relatively large values of  $v_A$  are more likely than relatively small values. This prior is *informative* because we have concentrated a large amount of mass over a relatively small range of parameter values.

between parameters to detect potential parameter trade-offs. One statistic often used to summarize the amount of uncertainty in a Bayesian parameter estimate is called the *95% Highest Density Interval* (HDI); this is the shortest interval that contains 95% of the mass of the posterior distribution.<sup>60</sup> Smaller HDIs indicate a posterior distribution with concentrated mass over a smaller amount of parameter space (and thus higher probability) than larger HDIs. Although many loosely equate the HDI with confidence intervals in traditional methods,<sup>61,62</sup> Bayesian HDIs are not the same thing.<sup>53,55</sup>

As noted earlier, in addition to specifying the likelihood of the data given the parameters,  $p(\mathbf{D}|\theta)$ , the Bayesian approach also requires us to specify priors,  $p(\theta)$ . Broadly speaking, there are two types of priors: *informative* and *noninformative*. An informative prior represents a strong *a priori* belief about the parameter values. These prior beliefs might represent expert knowledge about the parameters, previous fits of the model to other data, biological or other constraints on possible parameter values, or theoretical information regarding allowable parameter values. A noninformative prior represents a similar degree of belief across all possible parameter values (or a very wide range of parameter values). For example, if a parameter can only exist within a bounded range (not including positive or negative infinity), a useful noninformative prior might be a uniform distribution over that bounded range.

How to set priors is hotly debated in the Bayesian literature (e.g., Ref 63). For some, any informative prior is seen as too subjective, and has been part of a general critique of Bayesian statistics.<sup>64</sup> For others, an informative prior is viewed as an integral part of the model, forcing the theorist to formalize assumptions about model parameters.<sup>65,66</sup> In the case of cognitive psychology, we often have too little prior information to set informative priors in a way that would be universally accepted, so we often use noninformative priors. There is a large body of research devoted to developing noninformative priors which can be used as (at least arguably) a reasonable default.<sup>67–69</sup> Figure 2 shows examples of informative and noninformative priors based on the beta distributions.

The last part of Bayes' rule is its denominator (Eq. (6)),  $p(\mathbf{D})$ . This marginal likelihood is obtained by integrating the product of the likelihood and prior over the entire parameter space. When, as usually is the case for a cognitive model,  $\theta$  is a vector of parameters, this will be a multivariate integral. For example, the simplest version of the LBA posterior distribution (fixing  $s$  for identifiability) expands to:

$$p(v_1, v_2, a, \tau, b | \mathbf{D}) = \frac{p(\mathbf{D} | v_1, v_2, a, \tau, k) p(v_1, v_2, a, \tau, k)}{\int_{v_1}^b \int_{v_2}^b \int_a^b \int_{\tau}^b p(\mathbf{D} | v_1, v_2, a, \tau, k) p(v_1, v_2, a, \tau, k) dv_1 dv_2 da d\tau dk} \quad (7)$$

This is LBA applied to a single subject and a single experimental condition. Multiple conditions manipulating difficulty, bias, and speed–accuracy trade-off will significantly increase the number of model parameters, and hence significantly increase the complexity of the integral. Fitting a hierarchical model with many subjects simultaneously will further increase the number of model parameters and the complexity of the integral by orders of magnitude.

Multivariate integrals of complex nonlinear functions (such as likelihoods defined by cognitive models multiplied by the priors on parameters) nearly always defy any closed form solution. While standard numeric integration techniques such as quadrature exist, they can only be applied to numeric integrals over a handful of variables. These techniques are impractical for models with dozens or hundreds of parameters, which are often the case with hierarchical cognitive models; in fact it is not hard to reach a level of complexity that would require longer than the lifetime of the universe to solve the resulting multivariate integral using standard numerical methods.

Thankfully, modern computer hardware permits the use of *Monte Carlo techniques* (e.g., Ref 70). These allow one to draw random samples,  $\theta_i$ , directly from the posterior,  $p(\theta | \mathbf{D})$ , without having to explicitly solve an intractable integral; from these samples, summaries, and inferences about the underlying parameter distribution are possible.

## Computing the Posterior: Markov Chain Monte Carlo

*Markov Chain Monte Carlo* (MCMC) methods<sup>55,70–72</sup> can efficiently sample from high-dimensional probability density functions. Using random numbers to solve (hard) problems in generally are referred to as *Monte Carlo* techniques. With a traditional random number generator—think `rand()` or `randn()` in Matlab—sequential random samples (random numbers) are statistically independent of one another.<sup>b</sup> In MCMC, sequential samples are not statistically independent, but depend on the previous sample. Such dependency makes the process *Markov*. The sequence of such random numbers forms the chain in *Markov Chain*.

Perhaps the first MCMC method was the *Metropolis algorithm*.<sup>73</sup> In its simplest form, for the case of a single parameter, it begins by picking a

random initial value of the chain,  $\theta_0$ . Each step of the chain represents the next (potential) random sample from the probability density function. On each step  $i$  of the chain, a proposed random sample,  $\theta_*$ , is generated by adding random noise (often from a zero-centered normal distribution),  $\epsilon_i$ , to the previous random sample,  $\theta_{i-1}$ . The proposed sample,  $\theta_*$ , is always accepted if it has a higher probability density than that of the previous sample,  $\theta_{i-1}$ . If the proposed sample has a lower probability density than the previous sample, then the proposed sampled is accepted probabilistically, with an acceptance probability equal to the ratio of the probability density of the proposed sample versus the probability density of the previous sample.

The Metropolis algorithm is completely generic and can be applied to any probability density function. For the case of sampling from posterior probabilities in a Bayesian analysis, we can formalize the probability of acceptance of the next random sample in the chain as:

$$p(\text{accept}) = \min\left(1, \frac{p(\theta_* | \mathbf{D})}{p(\theta_{i-1} | \mathbf{D})}\right). \quad (8)$$

If the sample is accepted then  $\theta_i = \theta_*$ , otherwise  $\theta_i = \theta_{i-1}$ . The chain of  $\theta_i$  values represents random samples drawn from  $p(\theta | \mathbf{D})$ . Those random samples from the posterior can be used to calculate quantities like the mean, MAP, or HDI of parameter  $\theta$ . As noted earlier, being a chain of samples, one random number is not independent of the previous random number, unlike standard, non-MCMC, random number generators; such autocorrelation is not always an issue in practice, but techniques like *thinning*, whereby only every 10th or 50th or 100th sample in the chain are kept as true samples, are sometimes used.

At first blush, it appears as if we have done nothing to make the problem any more tractable. After all, calculating the posterior distributions still requires calculating an integral in the denominator of Bayes' rule. But note that the denominator is the same whether calculating  $p(\theta_* | \mathbf{D})$  or calculating  $p(\theta_{i-1} | \mathbf{D})$ . Those two denominators cancel each other in the ratio. Therefore, the acceptance formula simplifies to:

$$p(\text{accept}) = \min\left(1, \frac{p(\mathbf{D} | \theta_*) p(\theta_*)}{p(\mathbf{D} | \theta_{i-1}) p(\theta_{i-1})}\right), \quad (9)$$

with the integrals eliminated entirely. This is form of the Metropolis acceptance ratio most commonly seen

1 in the literature. Here we have illustrated the algo-  
 2 rithm assuming a single model parameter, but the  
 3 method can be extended to posteriors on many  
 4 parameters. This makes an intractable problem tract-  
 5 able, albeit at the computational cost of calculating  
 6 long chains of sampled random numbers (or sampled  
 7 vectors of random numbers in the case of multidimensional posteriors).

9 The Metropolis algorithm was later generalized  
 10 to arbitrary proposal distributions, including those  
 11 that are asymmetric, in the *Metropolis–Hastings*  
 12 *algorithm*;<sup>74</sup> asymmetric proposal distributions are  
 13 more efficient, for example, in the case of parameters  
 14 defined within bounded regions. MCMC was further  
 15 extended with the development of the *Gibbs*  
 16 *sampler*,<sup>75,76</sup> a special case of Metropolis–Hastings  
 17 where the acceptance probability is always 1, where  
 18 proposals are drawn from the full conditional distri-  
 19 butions for each parameter one at a time.<sup>70</sup> In  
 20 MCMC, it is valid to take Gibbs steps on some  
 21 parameters and Metropolis steps on others as in the  
 22 *Metropolis-within-Gibbs Sampler*.<sup>77</sup> While some  
 23 modelers program their own MCMC algorithms by  
 24 hand, there are a variety of software toolboxes that  
 25 function as black box inference engines performing  
 26 ‘automatic Bayesian inference’ via built-in  
 27 Metropolis–Hastings, Gibbs, and other samplers.  
 28 These include *WinBUGS*,<sup>78</sup> *JAGS*,<sup>79</sup> and *Stan*<sup>27</sup> with  
 29 its advanced MCMC algorithm based on *Hamiltonian Monte Carlo*.<sup>80</sup> These toolboxes only require  
 31 the user to specify the model (statistical or cognitive)  
 32 in a probabilistic programming language and then let  
 33 an automated inference engine generate samples from  
 34 the posterior distribution.

35 These toolkits have many built-in probability distri-  
 36 butions that make programming Bayesian statistical  
 37 models fairly straightforward. However, cognitive  
 38 models, such as LBA, require specialized likelihood  
 39 functions that are not pre-packaged with any toolbox.  
 40 For many, if not most, cognitive models, it is necessary  
 41 to implement custom probability distributions. *Win-*  
 42 *BUGS* and *JAGS* allow this, but require relatively low-  
 43 level programming in C++<sup>81</sup>; *Stan* allows this within  
 44 the same *Stan* programming language directly.<sup>45</sup>

## 46 LBA Example

48 Here we illustrate component-wise Metropolis with  
 49 the LBA model, in a simple example assuming data  
 50 from a single subject in a single condition. The  
 51 MCMC chain begins by initializing model param-  
 52 eters:  $v_A^0, v_B^0, s^0, \tau^0, b^0$ , with  $a$  fixed at 1 for identifiability.<sup>82</sup> On each step  $i$  of a chain of length  $N$ ,  $i = \{1, \dots, N\}$ , we do a Metropolis step for each parameter of

57 the model in random order. For each parameter, ran-  
 58 dom noise (here normally distributed with mean zero  
 59 and a small standard deviation of .05) is added to  
 60 the previous sample to produce a proposed sample.  
 61 The variance of the proposal distribution was found  
 62 after some experimentation by running the chain for  
 63 a short time and observing whether the chain was  
 64 efficiently exploring the parameter space.

65 To give a concrete example of how the algo-  
 66 rithm works, suppose we update  $v_A$  first:

$$67 v_A^* = v_A^{j-1} + \epsilon_{v_A}^i \quad (10)$$

70 The probability of accepting  $v_A^*$  is given by:

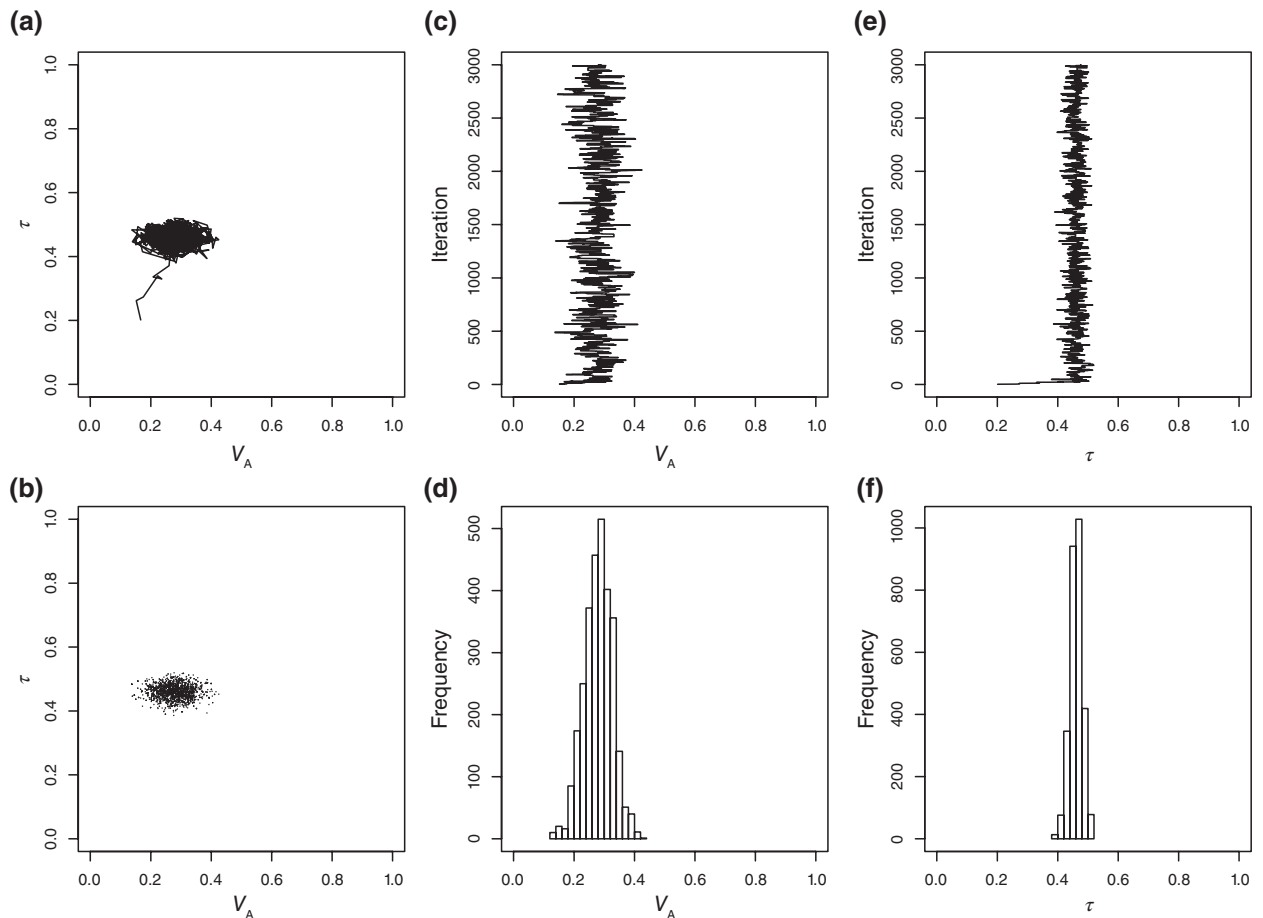
$$71 p(\text{accept}) = \min \left( 1, \frac{p(D|v_A^*, \theta') p(v_A^*)}{p(D|v_A^{j-1}, \theta') p(v_A^{j-1})} \right), \quad (11)$$

72 where  $\theta'$  contains all the most recently updated  
 73 parameters except for  $v_A$  (not shown in the prior cal-  
 74 culation because these terms cancel out due to an  
 75 assumption of independent priors for different  
 76 parameters); in this case,  $\theta' = (v_B^0, s^0, \tau^0, b^0)$ . If the  
 77 proposal is accepted, we set  $v_A^j = v_A^*$ , otherwise  
 78  $v_A^j = v_A^{j-1}$ . We repeat for all parameters of the model  
 79 to constitute a single MCMC step. We continue sam-  
 80 pling, creating an MCMC of length  $N$  or until the  
 81 algorithm converges on the full posterior distribution  
 82 by some criterion (for best practices on diagnosing  
 83 MCMC chain convergence, e.g., Ref 25).

84 Figure 3 shows an example of an MCMC chain  
 85 and the resulting samples for two of the five LBA  
 86 parameters; while we plot only two variables, all five  
 87 parameters are sampled. Panel (a) shows that the  
 88 chain stepped toward a higher value of  $\tau$  and a  
 89 slightly lower value of  $v_A$ . Panel (b) shows the distri-  
 90 bution of samples that were obtained after letting the  
 91 chain run for 3000 iterations and discarding the first  
 92 100 samples as *burn-in*; ‘burn-in’ represents samples  
 93 at the beginning of the MCMC chain that may not  
 94 be representative of the posterior distribution because  
 95 they are sampled from an extremely low density  
 96 region and simply reflect the initialization point for a  
 97 chain. Panel (c) shows the individual chain for  $v_A$   
 98 over the course of the iterations and these samples  
 99 are plotted as a histogram in Panel (d); panels (e) and  
 100 (f) show the same for  $\tau$ .

101 These Bayesian approaches require a mathe-  
 102 matically specified likelihood function,  $p(D|\theta)$ , for  
 103 the cognitive model. Unfortunately, predictions from  
 104 many interesting and important models in cognitive  
 105 psychology are based on computer simulation, they





**FIGURE 3** | (a) The path of the Markov chain for  $\tau$  and  $v_A$ . The chain begins in low density region around  $\tau = .2$  and  $v_A = .2$  and quickly moves to a higher density region as per the Metropolis acceptance probability ratio. (b) (below (a)) The resulting samples drawn from the joint posterior distribution of  $\tau$  and  $v_A$ , excluding the first 100 samples as burn-in. (c) The path the chain took over the marginal distribution for  $v_A$  at each iteration of the algorithm. The resulting marginal distribution is plotted below in (d). The path of the chain over the marginal distribution of  $\tau$  and the resulting samples are shown in (e) and (f), respectively.

are not defined by an explicit closed-form likelihood function. Traditionally, such models are fitted by simulating the model thousands of times for a given set of parameters, computing the discrepancy between model predictions and observed data, adjusting the parameters to minimize or maximize some objective function by hill-climbing or some other optimization technique. Fortunately, new methods are being proposed to allow Bayesian approaches to be applied to simulation-based cognitive models without any explicitly specified likelihood function (e.g., Refs 83–87).

### BAYESIAN MODEL PREDICTION

Once a cognitive model’s parameters have been estimated, a common next step is to generate model predictions and compare those qualitatively and quantitatively with observed data. In a traditional

modeling approach, one obtains a point estimate of the parameters,  $\hat{\theta}$ , that maximize likelihood or minimize RMSE, and then uses that point estimate to generate model predictions.

Bayesian prediction is different because Bayesian analysis produces a full joint posterior distribution of parameter values,  $p(\theta|D)$ . So cognitive model predictions should be based on that full parameter distribution, not a point estimate. The probability of a prediction,  $\hat{D}$ , is conditionalized on the model parameters,  $\theta$ , which in turn are conditionalized on the observed data,  $D$ . By the law of total probability, we can characterize model prediction as:

$$p(\hat{D}|D) = \int p(\hat{D}|\theta)p(\theta|D)d\theta. \tag{12}$$

This is known as the *posterior predictive distribution*. Consider an extreme case where the posterior,  $p(\theta|$

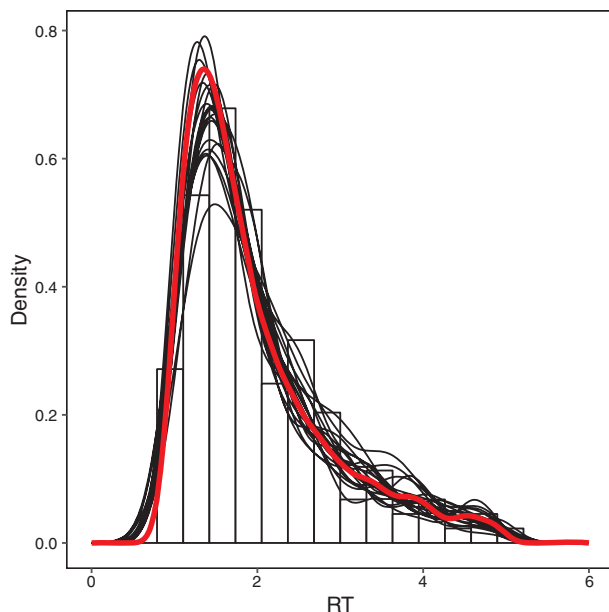
$D$ ), had all of its mass at a single point rather than a distribution; after observing the data, only a single vector of parameters,  $\hat{\theta}$ , is possible. In that case, the posterior predictive would reduce to the likelihood,  $p(\hat{D}|\hat{\theta})$ . For all other cases, the posterior predictive distribution is effectively a weighted average of the likelihood,  $p(\hat{D}|\theta)$ , with weights determined by the posterior,  $p(\theta|D)$ .

## COMPUTING THE POSTERIOR PREDICTIVE

On the surface, calculating the posterior predictive in Eq. (13) looks daunting. After all, it requires solving a multivariate integral, which could well be composed of hundreds or thousands of variables for a complex hierarchical model. But it turns out this is a fairly straightforward because a commonly used Monte Carlo integration technique can be applied. Since this technique might not be familiar to all readers, we provide a brief introduction.

Consider first the definition of expected value:

$$E[x] = \int xp(x)dx \approx \frac{1}{N} \sum_{i=1}^N x_i, \quad (13)$$



**FIGURE 4** | Shows the predictive distribution for each of several posterior samples (black lines) and the overall posterior predictive distribution (red line) plotted against the response time distribution simulated from the LBA (histogram bars).

where  $p(x)$  is the probability of  $x$ . This can be generalized to the expected value of a function  $g$  applied to  $x$  as:

$$E[g(x)] = \int g(x)p(x)dx \approx \frac{1}{N} \sum_{i=1}^N g(x_i). \quad (14)$$

The *expected value* is, as the name suggests, the long-run, expected theoretical mean for random variables  $x$  having probability  $p(x)$ , whereas the *average* is the empirical mean of observed data points  $x_i$  that have been sampled from  $p(x)$ .

Monte Carlo integration turns these formulae on their heads. Imagine instead of trying to estimate the theoretical mean or the empirical mean per the definition of expected value, you are trying to solve an integral. Suppose you need to solve an integral that has the functional form of Eqs. (13) or (14), where  $p(x)$  happens to be in the form of a probability density function from which samples can be generated. This integral can be approximated by drawing samples  $x_i$  from  $p(x)$  and averaging all the  $x_i$  (Eq. (13)) or averaging all the  $g(x_i)$  (Eq. (14)). By this Monte Carlo method, integrals that could be difficult or impossible to solve analytically are approximated computationally by sampling lots of random numbers.

In the present case,

$$\int p(\hat{D}|\theta)p(\theta|D)d\theta \approx \frac{1}{N} \sum_{i=1}^N p(\hat{D}|\theta_i) \quad (15)$$

where  $\theta_i$  is randomly sampled from  $p(\theta|D)$ . This is a form of Eq. (14). How do we generate those random samples? Well,  $p(\theta|D)$  is produced by the MCMC procedures described in the previous section, which in fact instantiates a random number generator for  $p(\theta|D)$ .

To generate predictions, we can simply simulate data from the model given each posterior sample,  $\theta_i$ , of the Markov chain and take the average. Figure 4 shows the histogram of response time data to which the LBA model was fit. Black lines plot a subset of the distributions,  $p(\hat{D}|\theta_i)$ , and the red line plots the posterior predictive distribution,  $p(\hat{D}|D)$ . The variability across  $p(\hat{D}|\theta_i)$  is reflective of the uncertainty associated with the parameters in the model. The posterior predictive distribution takes into account all of this uncertainty by averaging over all possible  $p(\hat{D}|\theta_i)$ . The posterior predictive distribution can

then be plotted over the relative frequency distribution given by the observed data to perform a *posterior predictive check*. If the predictive distribution aligns with the empirical distribution, then we can make a qualitative judgment as to whether the model adequately accounts for the data.

It is often the case that there are visual discrepancies between the observed data and the model predictions. Sometimes these discrepancies may be large enough to warrant rejection of the model. What is ‘large enough’? Sometimes it may be visually obvious when to reject the model and start over and other times it is not. One way to quantify these discrepancies is to construct a so-called discrepancy function between the observed and predicted data, denoted  $T(\hat{D}, D)$ . The discrepancy function measures how different the predicted values are from the observed. Using the same discrepancy function, one then measures the difference between the predictions,  $\hat{D}$ , and replications of the predictions  $\hat{D}_{rep}$  with  $T(\hat{D}, \hat{D}_{rep})$ .

The probability that  $T(\hat{D}, \hat{D}_{rep})$  is less than  $T(\hat{D}, D)$  is referred to as the Bayesian  $p$  value.<sup>88,89</sup>

The outcome of the Bayesian  $p$  value heavily relies on the choice of the discrepancy function, which is largely an arbitrary choice.<sup>90</sup> While this is problematic for the Bayesian  $p$  value, it might serve as a more objective alternative to a visual posterior predictive check.

## BAYESIAN MODEL SELECTION

Often, there are several competing cognitive models that the researcher would like to contrast, picking the model that ‘best’ explains the observed data. In a sense, while in parameter estimation we are interested in the probabilities of parameter values given data, in model selection we are interested in probabilities of models given data. In both the traditional and Bayesian approaches, to determine the ‘best’ fitting model, a trade-off must be made between overall goodness-of-fit and model complexity.<sup>91</sup>

As noted earlier, one traditional approach to model selection among non-nested models involves computing the maximum likelihood and then penalizing the model based on its number of parameters. The AIC<sup>31</sup> and the BIC<sup>30</sup> are both based on this type of rule. The BIC is given by:

$$BIC = -2 \ln p(\mathbf{D} | \hat{\theta}) + k \ln(n), \quad (16)$$

where  $k$  is the number of parameters in the model and  $n$  is the number of data points (the AIC penalty term is similarly additive, but only involves  $2k$ ). These approaches are computationally simple to implement, but they ignore both parameter uncertainty and the functional forms of the models (e.g., Ref 10).

By contrast, the Bayesian framework promises to provide a principled—if computationally more challenging—approach to select among competing models while taking into account both parameter uncertainty and the functional form of the model. Bayesian model selection can be thought of as weighing the evidence provided by the data in favor of alternative models.<sup>67</sup> In Bayesian terms, we are interested in the probability of model  $k$ ,  $\mathcal{M}_k$ , given data,  $\mathbf{D}$ . That probability can be found by a simple application of Bayes’ rule:

$$p(\mathcal{M}_k | \mathbf{D}) = \frac{p(\mathbf{D} | \mathcal{M}_k) p(\mathcal{M}_k)}{\sum_{j=1}^M p(\mathbf{D} | \mathcal{M}_j) p(\mathcal{M}_j)}. \quad (17)$$

where  $p(\mathcal{M}_k)$  is the prior probability of model  $k$ ,  $p(\mathbf{D} | \mathcal{M}_k)$  is the *marginal likelihood* for model  $k$ , and the sum in the denominator is a normalizing constant over all  $M$  possible models under consideration.

In the case of comparing two models,  $\mathcal{M}_1$  and  $\mathcal{M}_2$ , we consider the ratio of the posterior probabilities:

$$\frac{p(\mathcal{M}_1 | \mathbf{D})}{p(\mathcal{M}_2 | \mathbf{D})} = \frac{p(\mathbf{D} | \mathcal{M}_1) p(\mathcal{M}_1)}{\sum_{j=1}^M p(\mathbf{D} | \mathcal{M}_j) p(\mathcal{M}_j)} / \frac{p(\mathbf{D} | \mathcal{M}_2) p(\mathcal{M}_2)}{\sum_{j=1}^M p(\mathbf{D} | \mathcal{M}_j) p(\mathcal{M}_j)}. \quad (18)$$

The normalizing constants in the denominators drop out and we can rewrite the posterior odds as:

$$\frac{p(\mathcal{M}_1 | \mathbf{D})}{p(\mathcal{M}_2 | \mathbf{D})} = \frac{p(\mathbf{D} | \mathcal{M}_1) p(\mathcal{M}_1)}{p(\mathbf{D} | \mathcal{M}_2) p(\mathcal{M}_2)}. \quad (19)$$

The transformation from prior to posterior odds is determined by the ratio of the marginal likelihoods for each model. This transformation is the weight of the evidence provided by the data and is called the *Bayes Factor*,<sup>10,92–94</sup> denoted  $B_{12}$ :

$$B_{12} = \frac{p(\mathbf{D} | \mathcal{M}_1)}{p(\mathbf{D} | \mathcal{M}_2)}. \quad (20)$$

Note that the Bayes factor does not depend on the prior odds of the models. Arguably, this is convenient because there might be disagreement among theorists

as to the prior probability of the alternative models. If it is principled, one could set the prior odds to give each model equal footing.<sup>95–97</sup> In that case, the  $p(\mathcal{M}_k)$  terms in Eq. (19) cancel out, and the posterior odds equal the Bayes factor. As a rule of thumb, a  $B_{12}$  greater than 3 is generally considered to be positive evidence for Model 1, while a  $B_{12}$  greater than 10 is generally considered to be strong evidence for Model 1 (e.g., Ref 93). Of course, reciprocals of these Bayes factors provide corresponding levels of evidence for Model 2; indeed, one of the great strengths of Bayesian statistics in general is that it lets us evaluate both sides a comparison, not just one.

While simple to write out, how do we calculate  $p(\mathbf{D}|\mathcal{M}_k)$ ? Recall from earlier that we can rewrite Bayes' rule for parameter estimation in a form that makes the assumed model,  $\mathcal{M}$ , explicit:

$$p(\boldsymbol{\theta}|\mathbf{D},\mathcal{M}) = \frac{p(\mathbf{D}|\boldsymbol{\theta},\mathcal{M})p(\boldsymbol{\theta}|\mathcal{M})}{p(\mathbf{D}|\mathcal{M})}. \quad (21)$$

Now it is the marginal likelihood in the denominator,  $p(\mathbf{D}|\mathcal{M})$ , that is the focus of attention. But recall that we needed the development of MCMC techniques because the marginal likelihood is nearly always impossible to solve analytically and very difficult to estimate computationally. The marginal likelihoods,  $p(\mathbf{D}|\mathcal{M}_k)$ , integrate over the entire parameter space. While computationally challenging to estimate, the marginal likelihood takes into account both uncertainty associated with the parameters and the functional form of the model, providing a natural and principled penalty for model complexity.

## Computing the Bayes Factor

As previously shown, the marginal likelihood is obtained by integrating over the parameters:

$$p(\mathbf{D}|\mathcal{M}) = \int p(\mathbf{D}|\boldsymbol{\theta},\mathcal{M})p(\boldsymbol{\theta}|\mathcal{M})d\boldsymbol{\theta}. \quad (22)$$

So the computation of the Bayes factor can be written as:

$$B_{12} = \frac{\int p(\mathbf{D}|\boldsymbol{\theta},\mathcal{M}_1)p(\boldsymbol{\theta}|\mathcal{M}_1)d\boldsymbol{\theta}}{\int p(\mathbf{D}|\boldsymbol{\theta},\mathcal{M}_2)p(\boldsymbol{\theta}|\mathcal{M}_2)d\boldsymbol{\theta}} \quad (23)$$

These integrals are intractable for complex models with many parameters; standard numerical methods for solving integrals may work for models with relatively few parameters, but are not scalable to complex nonlinear models with many parameters. What about the Monte Carlo method we discussed earlier

for Bayesian prediction? Well, if we sample random numbers from the prior, it is *theoretically* possible to estimate marginal likelihoods thus:

$$\int p(\mathbf{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} \approx \frac{1}{N} \sum_{i=1}^N p(\mathbf{D}|\boldsymbol{\theta}_i), \quad (24)$$

where  $\boldsymbol{\theta}_i$  represent random samples from the prior. Unfortunately, estimating marginal likelihood this way turns out to be so highly inefficient as to be rendered impractical; this is in part because priors are often relatively flat, covering jointly a huge expanse of possible parameter space.<sup>93</sup>

That said, recent advancements in specialized computer hardware, referred to as graphical processing units (GPUs), have begun to make this task more feasible for some models.<sup>98</sup> GPUs were originally designed for the efficient control of computer graphics. Unlike CPUs, which process a single instruction at a time, GPUs implement massively parallel architectures and are now applied to range of scientific computing problems (e.g., Ref 99). The Monte Carlo marginal likelihood estimator in Eq. (24) is a perfect example of such a problem, in which the sampling procedure and likelihood computation can both be parallelized on a GPU. Evans and Brown,<sup>98</sup> using the LBA as an example, showed that a GPU can produce a marginal likelihood estimate containing 100,000,000 samples in mere minutes. For a typical CPU (around 2 GHz, circa 2017), that same computation might take a couple of days. Although the GPU method is promising, it is unknown whether it reasonably scales with increases in the dimensionality of the model. Evans and Brown showed the marginal likelihood estimates were highly variable for moderately sized models with multiple subjects. The GPU method also requires specialized hardware, which might be prohibitively expensive for some. However, there are other methods available that provide ways of estimating the Bayes factor with standard CPUs (for a tutorial, see Ref 100).

## Nested Model Comparison

Model comparison can be broadly classified into *nested model comparison* and *non-nested model comparison*. In nested model comparison, the two models being compared have the same parameters, but restriction is placed on the parameters of one of the models. For example, we might want to compare a model in which a particular parameter  $\theta$  is equal to some specific value  $\theta_1$  (e.g., in a model where  $\theta$  acts additively, the specific value might be  $\theta_1 = 0$ , and in

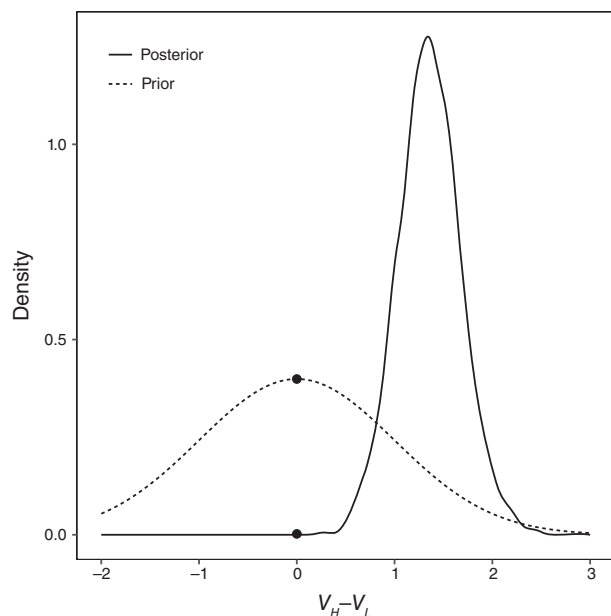


a model where  $\theta$  acts multiplicatively, the specific value might be  $\theta_1 = 1$ ) versus a model in which the parameter  $\theta$  is allowed to vary freely. For this nested model situation, we can estimate the Bayes factor by simply computing the ratio of the density of the posterior and prior at  $\theta_1$ . This is known as the *Savage–Dickey Ratio Test*.<sup>101</sup>

Consider two models, one nested within the other, where a restriction is placed on one of the parameters in one of the models,  $\theta = \theta_1$ . We refer to the restricted model as  $\mathcal{M}_R$  and the full model as  $\mathcal{M}_F$ . It can be shown (e.g., Ref 102) that the Bayes factor is the ratio of the posterior density at  $\theta_1$  and the prior density at  $\theta_1$ :

$$B_{RF} = \frac{p(D|\mathcal{M}_R)}{p(D|\mathcal{M}_F)} = \frac{p(\theta = \theta_1|D)}{p(\theta = \theta_1)}. \quad (25)$$

This entails that the Bayes factor can be computed by sampling from the posterior as is done via MCMC and determining the density at the desired point, here  $\theta_1$ . Values of  $B_{RF}$  greater than 1 indicate evidence for the restricted model, while values less than 1 indicates evidence for the full model, with strength of evidence given by the size of those ratios.



**FIGURE 5** | Graphical depiction of the Savage-Dickey ratio test. The dotted line is the prior placed on the effect size and solid line is posterior. The black dots represent the height of the prior and posterior when the effect size is 0. The ratio of these heights is the Bayes factor, the weight of the evidence. The height of the posterior at zero is 178 times less than the prior at zero, indicating the data have decreased our belief in the effect size being zero by a factor of 178.

As an example, consider a perceptual decision task with two conditions in which the discriminability of the target stimulus was manipulated. The hypothesis is that drift rate should vary across conditions. Let  $v_L$  and  $v_H$  denote the drift rate for correct responses for the condition in which discriminability is low and high, respectively. To test the hypothesis, we must determine whether the difference between  $v_L$  and  $v_H$  is equal to zero. One way to test this via the Savage–Dickey ratio would be to determine the following quantity:

$$B_{RF} = \frac{p(v_H - v_L = 0|D)}{p(v_H - v_L = 0)}. \quad (26)$$

Values greater than 1 would indicate evidence for the restricted model in which  $v_L$  equals  $v_H$ , while evidence less than 1 indicates evidence for the full model in which  $v_L$  and  $v_H$  are free to vary. Visually, the numerator is the height of the marginal posterior distribution  $p(v_H - v_L = 0|D)$  when  $v_L$  equals  $v_H$  and the denominator is the height of the prior when  $v_L$  equals  $v_H$ . Figure 5 shows a graphical example of what a Savage–Dickey test might look like after obtaining MCMC samples. The dotted line shows the density of the prior distribution  $p(v_H - v_L)$ . The black points show the height of the posterior and prior when  $v_H - v_L = 0$ . The ratio of these heights is equal to  $B_{RF}$ . In this example, the height of the marginal posterior at zero is  $p(v_H - v_L = 0|D) \approx .002$ . For the prior, the density is roughly  $p(v_H - v_L) = .40$ . Thus,  $B_{RF} \approx \frac{.002}{.40} = .005$ . This indicates that the data are roughly 200 times (i.e.,  $1/.005$ ) more likely under the full model than they are under the restricted model. There are several packages in R and Python that can compute the height of the estimated density from MCMC samples. In this example, we use a nonparametric density estimator suggested by Wagenmakers et al.<sup>101</sup>

## Non-nested Model Comparison

With non-nested model comparison, the models are completely different with different parameters, so computational short-cuts like Savage–Dickey are not available. For example, the LBA and the Leaky Competing Accumulator model<sup>59</sup> are both used to model choice response times, but make different assumptions and have different sets of parameters. These two models are therefore non-nested and if we wanted to compare their performance using the Bayes factor we would not be able to use the Savage–Dickey test.

1 However, there are some general computational  
2 methods that have been devised to compute the marginal  
3 likelihoods for any model (for a review, see Ref  
4 103). Most of these methods are quite complex and we  
5 will not detail them here, and only provide pointers for  
6 interested readers. One class of methods computes an  
7 estimate of the marginal likelihoods via Monte Carlo  
8 sampling. These include importance sampling,<sup>104</sup> reciprocal  
9 importance sampling,<sup>105</sup> annealed importance  
10 sampling,<sup>106</sup> bridge sampling,<sup>107</sup> Chib's method,<sup>108,109</sup>  
11 nested sampling,<sup>110</sup> and thermodynamic integration.  
12 <sup>111,112</sup> The thermodynamic approach has received  
13 significant attention in fields like biology<sup>112</sup> and ecology  
14 (e.g., Ref 113) in part because it is a general method  
15 that can be applied to any model with little modification  
16 to existing model code.

17 Another class of methods is called *transdimensional*  
18 *MCMC*, in which the competing models are  
19 placed within one 'supermodel.' On each step of the  
20 algorithm, a model index variable indicates one of  
21 the two models. The ratio of the proportion of times  
22 each model is visited equals the Bayes factor. Examples  
23 of transdimensional MCMC algorithms include  
24 *reversible-jump MCMC*<sup>114</sup> and the *product space*  
25 *method*.<sup>115,116</sup>

26 Lastly, there are information criterion  
27 approaches that are similar to BIC and AIC, but  
28 takes into account the uncertainty in the parameter  
29 estimates by considering the entirety of the MCMC  
30 sample. Examples of such information criteria  
31 include the *Bayesian Predictive Information Criterion*  
32 (<sup>117</sup>), the *Widely Applicable Bayesian Information*  
33 *Criterion* (WBIC;<sup>118,119</sup>), the *Widely Applicable*  
34 *Information Criterion* (WAIC;<sup>120,121</sup>), and the *Deviance*  
35 *Information Criterion* (DIC;<sup>122</sup>).<sup>c</sup>

### 36 37 38 Sensitivity to the Prior

39 Although the Bayes factor has gained significant traction  
40 in model selection in psychology,<sup>10,57,94,123</sup> and  
41 has been applied to model selection problems in a  
42 variety of domains (e.g., Refs 41,124,125), one  
43 potentially contentious issue is the Bayes factor's sensitivity  
44 to the priors on  $\theta$ .<sup>93,96,97</sup> This sensitivity contrasts  
45 with Bayesian parameter estimation, where any  
46 influence of priors is largely overwhelmed by the likelihood  
47 given a sufficient amount of data.

48 In Bayesian model selection, the average predictive  
49 performance of the model over the entire parameter  
50 space is assessed. Any prior that assigns low  
51 weight to high likelihood areas or high weight to low  
52 likelihood areas will penalize the model. And informative  
53 priors that end up in line with the likelihood  
54 will reward a model more so than vague priors or

57 one with informative priors that are not in line with  
58 the likelihood. Unlike parameter estimation, effects  
59 of the prior are not significantly diminished by having  
60 large amounts of data. Some argue that subjective  
61 prior belief about parameters should not significantly  
62 affect model selection.

63 In response, others have argued that priors  
64 placed on model parameters are a vital component of  
65 the theory, and therefore model selection should be  
66 sensitive to the priors.<sup>65,97,126</sup> Theoretical development  
67 could reflect different theoretical assumptions  
68 of the priors on  $\theta$  and model selection could be performed  
69 over these different instantiations. For example,  
70 Vanpaemel and Lee<sup>66</sup> formalized different  
71 assumptions about optimal dimensional attention  
72 weight parameters in the generalized context model  
73 of categorization using different priors (GCM<sup>127</sup>);  
74 the Bayes factor was then used select among the  
75 alternatives. The Bayes factor's sensitivity to the prior  
76 can be seen as advantageous when testing different  
77 theoretical assumptions.

78 Other times, we do not have strong theoretical  
79 assumptions that we can instantiate in priors. For  
80 example, if we are just beginning to develop a new  
81 model, we usually do not know what sort of parameter  
82 values we should expect, or do not have a theory  
83 of how those parameters might be set. In these cases,  
84 we are more concerned about the robustness of the  
85 inference made under a particular prior. To test the  
86 robustness of inference, we can conduct a *sensitivity*  
87 *analysis*<sup>10,93,128</sup> in which the Bayes factor is computed  
88 over a range of different priors. If the Bayes  
89 factor is qualitatively consistent across different prior  
90 settings, then we know that our inferences are robust  
91 under different assumptions about the prior.

### 92 93 94 CONCLUSION AND SUMMARY

95 Formal cognitive models describe psychological  
96 mechanisms in terms of mathematical structures. We  
97 contrasted traditional and Bayesian approaches to  
98 cognitive modeling, focusing on issues of parameter  
99 estimation, model prediction, and model selection.  
100 The goal of parameter estimation is to determine  
101 those parameters which provide the best fit of a cognitive  
102 model to the data. While traditional  
103 approaches, such as maximum likelihood, treat  
104 parameters as point estimates, the Bayesian approach  
105 quantifies uncertainty about the parameter estimates  
106 in terms of complete probability distributions. The  
107 goal of model prediction is to predict new data based  
108 on what has been observed. While traditional  
109 approaches use point estimate values to predict new  
110

1 data, the Bayesian approach takes into account the  
 2 uncertainty in the entire posterior distribution to gener-  
 3 ate model predictions. Lastly, the goal of model  
 4 selection is to select the cognitive model which best  
 5 explains the data. Traditional approaches use infor-  
 6 mation theoretic measures, such as AIC and BIC, that  
 7 do not take into account parameter uncertainty and  
 8 function form. The Bayesian approach, based on the  
 9 Bayes factor, takes into account the uncertainty over  
 10 the entire parameter space and balances complexity  
 11 versus fit to accomplish model selection.

12 An important topic we did not have space to  
 13 discuss is Bayesian hierarchical approaches. The  
 14 Bayesian approaches we discussed here could be  
 15 applied to individual participants  $i$ , separately obtain-  
 16 ing posterior distributions  $\theta_i|D$  for each participant.  
 17 Although this approach is useful for estimating  
 18 parameters at the individual level, we are often inter-  
 19 ested in both individual-level and group-level perfor-  
 20 mance. We can model both via a conceptually  
 21 straightforward extension to the Bayesian approach  
 22 that treats unknown parameters at the group level as  
 23 random variables as well as those that describe each  
 24 individual participant (e.g.,<sup>25</sup>).

25 By simultaneously estimating both group and  
 26 individual-level parameters, such hierarchical Bayes-  
 27 ian methods largely solve the problem of  
 28 aggregation,<sup>129</sup> which has been a key issue in cogni-  
 29 tive modeling for decades (e.g., Ref 130). When fit-  
 30 ting a model to data, often data will be aggregated  
 31 over trials or subjects depending on whether the  
 32 interest lies in group-level or individual-level conclu-  
 33 sions. However, the conclusions that one might draw  
 34 from a model fitted to aggregated data must be  
 35 drawn carefully because models can behave differ-  
 36 ently when fit to group and individual data.<sup>8,131</sup>  
 37 Bayesian hierarchical methods have provided a solu-  
 38 tion to problems associated with aggregation in a  
 39 wide range of different areas including recognition  
 40 memory,<sup>52</sup> multidimensional scaling,<sup>132,133</sup> and cate-  
 41 gory learning.<sup>7</sup> In addition, hierarchical Bayesian  
 42 approaches offer various avenues for cognitive  
 43 modeling including modeling multiple tasks within a  
 44 single model, assigning subjects to latent classes, and

57 modeling individual differences (for reviews, see Refs  
 58 8,56,129).

59 While Bayesian approaches to cognitive model-  
 60 ing clearly provide advantages over many traditional  
 61 approaches, one thing that should be clear from this  
 62 review is that Bayesian approaches carry the cost of  
 63 being more computationally intensive. We discussed  
 64 how MCMC algorithms have allowed Bayesian  
 65 models with hundreds or thousands of parameters—  
 66 especially in the case of hierarchical models—to be  
 67 fitted to data. The Bayesian approach provides a  
 68 coherent way to update beliefs in light of data and  
 69 offers an extremely flexible framework to fit  
 70 individual- and group-level parameters not only in  
 71 theory but also in practice.

## NOTES

<sup>a</sup> The mode of the posterior is often called the *maximum a posteriori* (MAP) estimate in Bayesian analysis.

<sup>b</sup> Technically, of course, the pseudo-random number generators used in nearly all programming environments are based on a completely deterministic algorithm; they produce a sequence of random samples that cannot be distinguished statistically from those produced by a true random process.

<sup>c</sup> While the WBIC requires samples from the posterior raised to the power of  $1/\ln(n)$ , where  $n$  is the number of data points, the WAIC and DIC only require samples from the posterior and are, therefore, easily computed. The DIC is known to have problems penalizing the model for complexity, as it is known to sometimes yield a negative estimate for the number of effective parameters in the model. The WAIC does not suffer from this and has many other advantages.<sup>120</sup> The WAIC is given by:

$$\text{WAIC} = \sum_{i=1}^n \ln \left( \frac{1}{S} \sum_{s=1}^S p(D_i | \theta_s) \right) - \sum_{i=1}^n V_{s=1}^S (\ln p(D_i | \theta_s)),$$

96 where  $n$  is the number of data points,  $S$  is the number of  
 97 posterior samples, and  $V_{s=1}^S (\ln p(D_i | \theta_s))$  is the sample vari-  
 98 ance of the log likelihood of the data point,  $D_i$ , under all  
 99 the posterior samples. The first term acts as the goodness-  
 100 of-fit measure, while the second term penalizes the model  
 101 for its complexity.

## ACKNOWLEDGMENTS

50 This work was supported by NSF SBE-1257098, NSF SBE-1640681, NEI R01 EY021833, the Temporal  
 51 Dynamics of Learning Center (NSF SMA-1041755), the Vanderbilt Vision Research Center (NEI P30-  
 52 EY008126), a Discovery Grant from Vanderbilt University, and a training grant from the NEI (T32-  
 53 EY07135). The authors thank Brandon Turner and Gabriel Tillman for their helpful comments and criticisms.  
 54

## REFERENCES

1. Lewandowsky S, Farrel S. *Computational Modeling in Cognition: Principles and Practice*. Thousand Oaks, CA: Sage Publications; 2011.
2. Ratcliff R, Smith PL. A comparison of sequential sampling models for two-choice reaction time. *Psychol Rev* 2004, 111:333–367.
3. Nosofsky RM. Exemplar-based approach to relating categorization, identification, and recognition. In: Ashby G, ed. *Multidimensional Models of Perception and Cognition*. Hillsdale, NJ: Lawrence Erlbaum; 1992, 363–394.
4. Richler JJ, Palmeri TJ. Visual category learning. *WIREs Cogn Sci* 2014, 5:75–94.
5. Polyn SM, Norman KA, Kahana MJ. A context maintenance and retrieval model of organizational processes in free recall. *Psychol Rev* 2009, 116:129–156.
6. Turner BM, Forstmann BU, Love BC, Palmeri TJ. Approaches to analysis in model-based cognitive neuroscience. *J Math Psychol* 2017, 76:65–79 <https://doi.org/10.1016/j.jmp.2016.01.001>.
7. Bartlema A, Lee MD, Wetzels R, Vanpaemel W. A Bayesian hierarchical mixture approach to individual differences: case studies in selective attention and representation in category learning. *J Math Psychol* 2014, 59:132–150.
8. Shen J, Palmeri T. Modelling individual difference in visual categorization. *Vis Cogn* 2016, 24:260–283 <https://doi.org/10.1080/13506285.2016.1236053>.
9. Wiecki TV, Poland J, Frank MJ. Model-based cognitive neuroscience approaches to computational psychiatry: Clustering and classification. *Clin Psychol Sci* 2015, 3:378–399.
10. Myung IJ, Pitt MA. Applying Occam's razor in modeling cognition: a Bayesian approach. *Psychon Bull Rev* 1997, 4:79–95.
11. Busemeyer JR, Diederich A. *Cognitive Modeling*. Thousand Oaks, CA: Sage Publications; 2010.
12. Newell A. *You Can't Play 20 Questions with Nature and Win: Projective Comments on the Papers of this Symposium*. New York: Academic Press; 1973.
13. Allport DA. Critical notice: The state of cognitive psychology. *Q J Exp Psychol* 1975, 27:141–152.
14. Shiffrin RM, Nobel PA. The art of model development and testing. *Behav Res Methods Instrum Comput* 1997, 29:6–14.
15. Polk TA, Seifert CM. *Cognitive Modeling*. Cambridge, MA: MIT Press; 2002.
16. Clark SE, Gronlund SD. Global matching models of recognition memory: how the models match the data. *Psychon Bull Rev* 1996, 3:37–60 <https://doi.org/10.3758/BF03210740>.
17. Pothos EM, Wills AJ. *Formal Approaches to Categorization*. Cambridge, UK: Cambridge University Press; 2011.
18. Palmeri TJ, Love BC, Turner BM. Model-based cognitive neuroscience. *J Math Psychol* 2017, 76:59–64 <https://doi.org/10.1016/j.jmp.2016.10.010>.
19. Jacobs RA, Kruschke JK. Bayesian learning theory applied to human cognition. *WIREs Cogn Sci* 2011, 2:8–21.
20. Griffiths T, Kemp C, Tenenbaum J. Bayesian models of cognition. In: Sun R, ed. *The Cambridge Handbook of Expertise and Expert Performance*. New York: Cambridge University Press; 2008, 59–100.
21. Anderson JR. *The Adaptive Character of Thought*. Hillsdale NJ: Lawrence Erlbaum; 1990.
22. Jones M, Love BC. Bayesian fundamentalism or enlightenment? On the explanatory status and theoretical contributions of Bayesian models of cognition. *Behav Brain Sci* 2011, 34:169–231.
23. Kruschke JK. *Doing Bayesian Data Analysis: A Tutorial with R, JAGS, and Stan*. 2nd ed. Burlington, MA: Academic Press; 2014.
24. Jackman S. *Bayesian Analysis for the Social Sciences*. John Wiley & Sons; 2009.
25. Gelman A, Carlin JB, Stern HS, Rubin DB. *Bayesian Data Analysis*. Boca Raton, FL: Chapman & Hall/CRC; 2014.
26. Myung IJ. Tutorial on maximum likelihood estimation. *J Math Psychol* 2003, 47:90–100.
27. Carpenter B, Gelman A, Hoffman M, Lee D, Goodrich B, Betancourt M, et al. Stan: a probabilistic programming language. *J Stat Softw* 2016, 76:1–32.
28. Nelder JA, Mead R. A simplex method for function minimization. *Comput J* 1964, 7:308–313 <https://doi.org/10.1093/comjnl/7.4.308>.
29. Wagenmakers E-J, Ratcliff R, Gomez R, Iverson GJ. Assessing model mimicry using the parametric bootstrap. *J Math Psychol* 2004, 48:28–50.
30. Schwarz G. Estimating the dimension of a model. *Ann Stat* 1978, 6:461–464 <https://doi.org/10.1214/aos/1176344136>.
31. Akaike H. A new look at the statistical model identification. *IEEE Trans Autom Control* 1974, 19:716–723 <https://doi.org/10.1109/TAC.1974.1100705>.
32. Chater N, Oaksford M, Hahn U, Heit E. Bayesian models of cognition. *WIREs Cogn Sci* 2010, 1:811–823.
33. Kruschke JK. Bayesian data analysis. *WIREs Cogn Sci* 2010, 1:658–676.
34. Bernardo JM, Smith AFM. *Bayesian Theory*. New York: Wiley; 1994.



- 1 35. de Finetti B. Foresight: its logical laws in subjective  
2 sources. In: Kyburg H, Smokler H, eds. *Studies in*  
3 *Subjective Probability*. New York: Wiley; 1964,  
4 93–158.
- 5 36. Bayes T. A letter from the late reverend Mr. Thomas  
6 Bayes, FRS to John Canton, MA and FRS. *Philos*  
7 *Trans* 1763, 53:269–271.
- 8 37. Lee MD. Three case studies in the Bayesian analysis  
9 of cognitive models. *Psychon Bull Rev* 2008,  
10 15:1–15.
- 11 38. Rouder JN, Lu J. An introduction to Bayesian hierar-  
12 chical models with an application in the theory of sig-  
13 nal detection. *Psychon Bull Rev* 2005, 12:573–604.
- 14 39. Matzke D, Dolan CV, Batchelder WH,  
15 Wagenmakers E-J. Bayesian estimation of multinomial  
16 processing tree models with heterogeneity in partic-  
17 ipants and items. *Psychometrika* 2015,  
18 80:205–235.
- 19 40. Katahira K. How hierarchical models improve point  
20 estimates of model parameters at the individual level.  
21 *J Math Psychol* 2016, 73:37–58 [https://doi.org/10.](https://doi.org/10.1016/j.jmp.2016.03.007)  
22 [1016/j.jmp.2016.03.007](https://doi.org/10.1016/j.jmp.2016.03.007).
- 23 41. Lee MD, Webb MR. Modeling individual differences  
24 in cognition. *Psychon Bull Rev* 2005, 12:605–621.
- 25 42. Okada K, Lee MD. A Bayesian approach to modeling  
26 group and individual differences in multidimensional  
27 scaling. *J Math Psychol* 2016, 70:35–44.
- 28 43. Merkle EC, Smithson M, Verkuilen J. Hierarchical  
29 models of simple mechanisms underlying confidence  
30 in decision making. *J Math Psychol* 2011, 55:57–67  
31 <https://doi.org/10.1016/j.jmp.2010.08.011>.
- 32 44. Ravenzwaaj D, Moore CP, Lee MD, Newell BR. A  
33 hierarchical Bayesian modeling approach to searching  
34 and stopping in multi-attribute judgment. *Cogn Sci*  
35 2014, 38:1384–1405.
- 36 45. Annis J, Miller BJ, Palmeri TJ. Bayesian inference  
37 with Stan: a tutorial on adding custom distributions.  
38 *Behav Res Methods* 2016.
- 39 46. Ratcliff R, Childers R. Individual differences and fit-  
40 ting methods for the two-choice diffusion model of  
41 decision making. *Decision* 2015, 2:237–279.
- 42 47. Vandekerckhove J, Tuerlinckx F, Lee MD. Hierarchi-  
43 cal diffusion models for two-choice response times.  
44 *Psychol Methods* 2011, 16:44–62.
- 45 48. Wiecki TV, Sofer I, Frank MJ. HDDM: hierarchical  
46 bayesian estimation of the drift-diffusion model in  
47 python. *Front Neuroinform* 2013, 7:14.
- 48 49. Annis J, Lenes JG, Westfall HA, Criss AH,  
49 Malmberg KJ. The list-length effect does not discrimi-  
50 nate between models of recognition memory. *J Mem*  
51 *Lang* 2015, 85:27–41.
- 52 50. Dennis S, Lee MD, Kinnell A. Bayesian analysis of  
53 recognition memory: the case of the list-length effect.  
54 *J Mem Lang* 2008, 59:361–376.
51. Morey RD. A Bayesian hierarchical model for the  
57 measurement of working memory capacity. *J Math*  
58 *Psychol* 2011, 55:8–24 [https://doi.org/10.1016/j.jmp.](https://doi.org/10.1016/j.jmp.2010.08.008)  
59 [2010.08.008](https://doi.org/10.1016/j.jmp.2010.08.008). 60
52. Pratte MS, Rouder JN. Hierarchical single- and dual-  
61 process models of recognition memory. *J Math Psy-*  
62 *chol* 2011, 55:36–46 [https://doi.org/10.1016/j.jmp.](https://doi.org/10.1016/j.jmp.2010.08.007)  
63 [2010.08.007](https://doi.org/10.1016/j.jmp.2010.08.007). 64
53. Turner BM, Forstmann BU, Wagenmakers E-J,  
65 Brown SD, Sederberg PB, Steyvers M. A Bayesian  
66 framework for simultaneously modeling neural and  
67 behavioral data. *Neuroimage* 2013, 72:193–206  
68 <https://doi.org/10.1016/j.neuroimage.2013.01.048>. 69
54. Brown SD, Heathcote A. The simplest complete  
70 model of choice response time: linear ballistic accu-  
71 mulation. *Cogn Psychol* 2008, 57:153–178. 72
55. Lee MD, Wagenmakers E-J. *Bayesian Cognitive*  
73 *Modeling: A Practical Course*. Cambridge, UK: Cam-  
74 bridge University Press; 2014. 75
56. Lee MD. How cognitive modeling can benefit from  
76 hierarchical Bayesian models. *J Math Psychol*  
77 2011, 55:1–7. 78
57. Shiffrin RM, Lee MD, Kim W, Wagenmakers E-J. A  
79 survey of model evaluation approaches with a tutorial  
80 on hierarchical bayesian methods. *Cogn Sci* 2008,  
81 32:1248–1284. 82
58. Lynch SM. *Introduction to Applied Bayesian Statis-*  
83 *tics and Estimation for Social Scientists*. New York:  
84 Springer Science & Business Media; 2007. 85
59. Usher M, McClelland JL. The time course of percep-  
86 tual choice: the leaky, competing accumulator model.  
87 *Psychol Rev* 2001, 108:550–592. 88
60. Hyndman RJ. Computing and graphing highest den-  
89 sity regions. *Am Stat* 1996, 50:120–126. 90
61. Jaynes ET, Kempthorne O. Confidence intervals vs  
91 Bayesian intervals. In: Harper WL, Hooker C, eds.  
92 *Foundations of Probability Theory, Statistical Infer-*  
93 *ence, and Statistical Theories of Science*. Dordrecht,  
94 The Netherlands: Springer; 1976, 175–257. 95
62. Morey RD, Hoekstra R, Rouder JN, Lee MD,  
96 Wagenmakers E-J. The fallacy of placing confidence  
97 in confidence intervals. *Psychon Bull Rev* 2016,  
98 23:103–123 [https://doi.org/10.3758/s13423-015-](https://doi.org/10.3758/s13423-015-0947-8)  
99 [0947-8](https://doi.org/10.3758/s13423-015-0947-8). 100
63. Lindley D. That wretched prior. *Significance* 2004,  
101 1:85–87. 102
64. Efron B. Why isn't everyone a Bayesian? *Am Stat.*  
103 1986, 40:1–5. 104
65. Lee MD, Vanpaemel W. Determining informative  
105 priors for cognitive models. *Psychon Bull Rev* 2017. 106
66. Vanpaemel W, Lee MD. Using priors to formalize  
107 theory: optimal attention and the generalized context  
108 model. *Psychon Bull Rev* 2012, 19:1047–1056  
109 <https://doi.org/10.3758/s13423-012-0300-4>. 110

67. Jeffreys H. *Theory of Probability*. 3rd ed. Oxford, UK: Oxford University Press; 1961.
68. Kass RE, Wasserman L. The selection of prior distributions by formal rules. *J Am Stat Assoc* 1996, 91:1343–1370.
69. Spiegelhalter DJ, Smith AFM. Exact and approximate posterior moments for a normal location. *J R Stat Soc Ser B Stat Methodol*. 1982, 54:793–804.
70. Brooks S, Gelman A, Jones G, Meng X-L. *Handbook of Markov Chain Monte Carlo*. Boca Raton, FL: Chapman & Hall/CRC Press; 2011.
71. Green PJ, Łatuszyński K, Pereyra M, Robert CP. Bayesian computation: a summary of the current state, and samples backwards and forwards. *Stat Comput* 2015, 25:835–862.
72. van Ravenzwaaij D, Cassey P, Brown SD. A simple introduction to Markov chain Monte Carlo sampling. *Psychon Bull Rev* 2016 <https://doi.org/10.3758/s13423-016-1015-8>.
73. Metropolis N, Rosenbluth AW, Rosenbluth MN, Teller AH, Teller E. Equation of state calculations by fast computing machines. *J Chem Phys* 1953, 21:1087–1092 <https://doi.org/10.1063/1.1699114>.
74. Hastings WK. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* 1970, 57:97–109.
75. Gelfand A, Smith A. Sampling-based approaches to calculating marginal densities. *J Am Stat Assoc* 1990, 85:398–409.
76. Geman S, Geman D. Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Trans Pattern Anal Mach Intell* 1984, PAMI-6:721–741.
77. Gilks WR, Best NG, Tan KKC. Adaptive rejection metropolis sampling within Gibbs sampling. *J R Stat Soc Ser C Appl Stat* 1995, 44:455–472.
78. Lunn DJ, Thomas A, Best N, Spiegelhalter D. WinBUGS – a Bayesian modelling framework: concepts, structure, and extensibility. *Stat Comput*. 2000, 10:325–337.
79. Plummer M. JAGS: a program for analysis of Bayesian graphical models using Gibbs sampling. In: *Proceedings of 3rd International Workshop on Distributed Statistical Computing (DSC 2003)*, 2003, 20–2 [doi:https://doi.org/10.1.1.13.3406](https://doi.org/10.1.1.13.3406).
80. Neal RM. MCMC using Hamiltonian dynamics. In: Brooks S, Gelman A, Jones G, Meng X-L, eds. *Handbook of Markov Chain Monte Carlo*. Boca Raton, FL: Chapman & Hall/CRC Press; 2011, 113–162.
81. Wabersich D, Vandekerckhove J. Extending JAGS. A tutorial on adding custom distributions to JAGS (with a diffusion model example). *Behav Res Methods* 2014, 46:15–28 <https://doi.org/10.3758/s13428-013-0369-3>.
82. Donkin C, Brown S, Heathcote A. Drawing conclusions from choice response time models: a tutorial using the linear ballistic accumulator. *J Math Psychol* 2011, 55:140–151.
83. Holmes WR. A practical guide to the probability density approximation (PDA) with improved implementation and error characterization. *J Math Psychol* 2015, 68–69:13–24.
84. Turner BM, Sederberg PB. Approximate Bayesian computation with differential evolution. *J Math Psychol* 2012, 56:375–385 <https://doi.org/10.1016/j.jmp.2012.06.004>.
85. Turner BM, Sederberg PB, McClelland JL. Bayesian analysis of simulation-based models. *J Math Psychol* 2014, 72:191–199 <https://doi.org/10.1016/j.jmp.2014.10.001>.
86. Turner BM, Van Zandt T. A tutorial on approximate Bayesian computation. *J Math Psychol* 2012, 56:69–85 <https://doi.org/10.1016/j.jmp.2012.02.005>.
87. Turner BM, Van Zandt T. Hierarchical approximate Bayesian computation. *Psychometrika* 2013, 79:185–209.
88. Gelman A. Two simple examples for understanding posterior p-values whose distributions are far from uniform. *Electron J Stat* 2013, 7:2595–2602.
89. Zhang JL. Comparative investigation of three Bayesian p values. *Comput Stat Data Anal* 2014, 79:277–291 <https://doi.org/10.1016/j.csda.2014.05.012>.
90. Kruschke JK. Posterior predictive checks can and should be Bayesian: comment on Gelman and Shalizi, “philosophy and the practice of Bayesian statistics”. *Br J Math Stat Psychol* 2013, 66:45–56.
91. Myung IJ. The importance of complexity in model selection. *J Math Psychol* 2000, 44:190–204 <https://doi.org/10.1006/jmps.1999.1283>.
92. Good IJ. The interface between statistics and philosophy of science. *Stat Sci* 1988, 3:386–412 <https://doi.org/10.1214/ss/1177013604>.
93. Kass RE, Raftery AE. Bayes factors. *J Am Stat Assoc* 1995, 90:773–795.
94. Shiffrin RM, Chandramouli SH, Grunwald PD. Bayes factors, relations to minimum description length, and overlapping model classes. *J Math Psychol*. 2015, 72:56–77 <https://doi.org/10.1016/j.jmp.2015.11.002>.
95. Kruschke JK. Bayesian assessment of null values via parameter estimation and model comparison. *Perspect Psychol Sci* 2011, 6:299–312.
96. Liu CC, Aitkin M. Bayes factors: prior sensitivity and model generalizability. *J Math Psychol* 2008, 52:362–375 <https://doi.org/10.1016/j.jmp.2008.03.002>.
97. Vanpaemel W. Prior sensitivity in theory testing: an apologia for the Bayes factor. *J Math Psychol*

- 2010, 54:491–498 <https://doi.org/10.1016/j.jmp.2010.07.003>.
98. Evans NJ, Brown SD. Bayes factors for the Linear Ballistic Accumulator Model of Decision-Making. *Behav Res Methods* 2017 Advance online publication, <https://doi.org/10.3758/s13428-017-0887-5>.
99. Luebke D. CUDA: Scalable parallel programming for high-performance scientific computing. In: *5th IEEE International Symposium on Biomedical Imaging: From Nano to Macro, 2008 ISBI 2008*, 2008, 836–838.
100. Annis J, Evans NJ, Miller BJ, Palmeri TJ. Thermodynamic integration and steppingstone sampling methods for estimating Bayes factors: a tutorial for psychologists
101. Wagenmakers E-J, Lodewyckx T, Kuriyal H, Grasman R. Bayesian hypothesis testing for psychologists: a tutorial on the savage-dickey method. *Cogn Psychol* 2010, 60:158–189 <https://doi.org/10.1016/j.cogpsych.2009.12.001>.
102. Morey RD, Rouder JN, Pratte MS, Speckman PL. Using MCMC chain outputs to efficiently estimate Bayes factors. *J Math Psychol* 2011, 55:368–378.
103. Friel N, Wyse J. Estimating the evidence - a review. *Stat Neerl* 2012, 66:288–308.
104. Geweke BYJ. Bayesian inference in iconometric models using monte carlo integration. *Econometrica* 1989, 57:1317–1339.
105. Newton M, Raftery A. Approximate Bayesian inference with the weighted likelihood bootstrap. *J R Stat Soc Ser B Stat Methodol* 1994, 56:3–48.
106. Neal RM. Annealed importance sampling. *Stat Comput* 2001, 11:125–139.
107. Meng X-L, Wong HW. Simulating ratios of normalizing constants via a simple identity: a theoretical exploration. *Stat Sin* 1996, 6:831–860.
108. Chib S, Jeliazkov I. Marginal likelihood from the metropolis-Hastings output. *J Am Stat Assoc* 2001, 96:270–281 <https://doi.org/10.2307/2291521>.
109. Chib S. Marginal likelihood from the Gibbs output. *J Am Stat Assoc* 1995, 90:1313–1321 <https://doi.org/10.2307/2291521>.
110. Skilling J. Nested sampling for Bayesian computations. *Bayesian Anal* 2006, 1:833–860 <https://doi.org/10.1214/06-BA127>.
111. Friel N, Pettitt AN. Marginal likelihood estimation via power posteriors. *J R Stat Soc Ser B Stat Methodol*. 2008, 70:589–607.
112. Lartillot N, Philippe H. Computing bayes factors using thermodynamic integration. *Syst Biol* 2006, 55:195–207 <https://doi.org/10.1080/10635150500433722>.
113. Liu P, Elshall AS, Ye M, Beerli P, Zeng X, Lu D, Tao Y. Evaluating marginal likelihood with thermodynamic integration method and comparison with several other numerical methods. *Water Resour Res* 2016, 52:734–758.
114. Green PJ. Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika* 1995, 82:711–732.
115. Carlin BP, Chib S. Bayesian model choice via Markov chain Monte Carlo methods. *J R Stat Soc Ser B (Stat Methodol)* 1995, 57:473–484.
116. Lodewyckx T, Kim W, Lee MD, Tuerlinckx F, Kuppens P, Wagenmakers E-J. A tutorial on Bayes factor estimation with the product space method. *J Math Psychol*. 2011, 55:331–347.
117. Ando T. Bayesian predictive information criterion for the evaluation of hierarchical Bayesian and empirical Bayes models. *Biometrika* 2017, 94:443–458.
118. Friel N, Mckeone JP, Oates CJ, Pettitt AN. Investigation of the widely applicable Bayesian information criterion. *Stat Comput* 2016, 27:833–844.
119. Watanabe S. A widely applicable Bayesian information criterion. *Mach Learn Res* 2013, 14:867–897.
120. Gelman A, Hwang J, Vehtari A. Understanding predictive information criteria for Bayesian models. *Stat Comput* 2014, 24:997–1016.
121. Watanabe S. Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. *J Mach Learn Res* 2010, 11:3571–3594.
122. Spiegelhalter DJ, Best NG, Carlin BP, van der Linde A. Bayesian measures of model complexity and fit. *J R Stat Soc Ser B Stat Methodol* 2002, 64:583–639.
123. Lee MD, Wagenmakers E-J. Bayesian statistical inference in psychology: comment on Trafimow (2003). *Psychol Rev*. 2005, 112:662–668.
124. Lee MD. Determining the dimensionality of multidimensional scaling representations for cognitive modeling. *J Math Psychol* 2001, 45:149–166 <https://doi.org/10.1006/jmps.1999.1300>.
125. Rouder JN, Lu J, Speckman P, Sun D. A hierarchical model for estimating response time distributions. *Psychon Bull Rev* 2005, 12:195–223.
126. Vanpaemel W. Constructing informative model priors using hierarchical methods. *J Math Psychol* 2011, 55:106–117 <https://doi.org/10.1016/j.jmp.2010.08.005>.
127. Nosofsky RM. Attention, similarity, and the identification-categorization relationship. *J Exp Psychol Gen* 1986, 115:39–57.
128. Sinharay S, Stern HS. On the sensitivity of Bayes factors to the prior distributions. *Am Stat* 2002, 56:196–201.
129. Rouder JN, Morey RD, Pratte MS. Hierarchical Bayesian models. In: Batchelder WH, Colonius H, Dzhamfarov E, Myung JI, eds. *New Handbook of Mathematical Psychology*. Measurement and

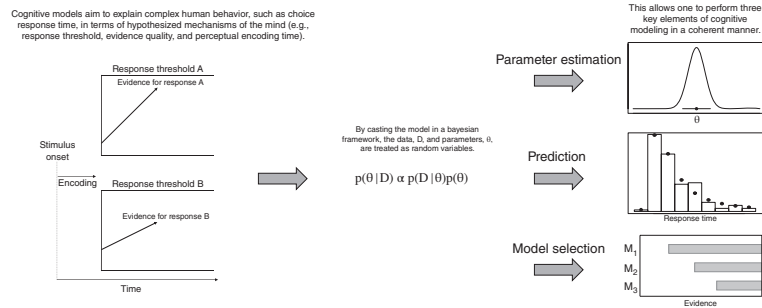
1	Methodology, vol. 1. London, UK: Cambridge University Press; 2013.	
2		
3	130. Estes WK. The problem of inference from curves based on group data. <i>Psychol Bull</i> 1956, 53:134–140.	57
4		58
5		59
6	131. Cohen AL, Sanborn AN, Shiffrin RM. Model evaluation using grouped or individual data. <i>Psychon Bull Rev</i> 2008, 15:692–712.	60
7		61
8		62
9		63
10		64
11		65
12		66
13		67
14		68
15		69
16		70
17		71
18		72
19		73
20		74
21		75
22		76
23		77
24		78
25		79
26		80
27		81
28		82
29		83
30		84
31		85
32		86
33		87
34		88
35		89
36		90
37		91
38		92
39		93
40		94
41		95
42		96
43		97
44		98
45		99
46		100
47		101
48		102
49		103
50		104
51		105
52		106
53		107
54		108
		109
		110



Graphical abstract

Bayesian statistical approaches to evaluate cognitive models

Jeffrey Annis<sup>1</sup>, Thomas J. Palmeri<sup>1</sup>



Cognitive models aim to explain complex human behavior such as choice response time, in terms of hypothesized mechanisms of the mind (e.g., response threshold, evidence quality, and perceptual encoding time).