

Drawing conclusions from choice response time models: A tutorial using the linear ballistic accumulator

Chris Donkin^{a,*}, Scott Brown^b, Andrew Heathcote^b

^a Department of Psychological and Brain Sciences, Indiana University, United States

^b School of Psychology, The University of Newcastle, Australia

ARTICLE INFO

Article history:

Received 30 April 2009

Received in revised form

23 September 2010

Available online 2 November 2010

Keywords:

Response time models

Tutorial

Linear ballistic accumulator model

Evidence accumulator models

ABSTRACT

Cognitive models of choice and response times can lead to deeper insights into the processes underlying decisions than standard analyses of accuracy and response time data. The application of these models, however, has historically been reserved for the authors of the models, and their associates. Recently, choice response time models have become more accessible through the release of user-friendly software for estimating their parameters. The aim of this tutorial is to provide guidance about the process of using these parameter estimates and associated model fits to make conclusions about experimental data. We use an application of one response time model, the linear ballistic accumulator, as an example to demonstrate the steps required to select an appropriate parametric characterization of a data set. We also discuss how to evaluate the quality of the agreement between model and data, including guidelines for presenting model predictions for group-level data.

© 2010 Elsevier Inc. All rights reserved.

1. Introduction

Evidence accumulation models of choice response time (RT) are increasingly used to examine the psychological processes underlying rapid decisions. The central assumption of these models is that the decision maker accumulates evidence for potential choices and makes a decision once the evidence reaches a threshold amount. The predicted time to make a response is the time taken to accumulate evidence, plus “non-decision time”, which is the time for other necessary processes, such as stimulus encoding and response execution. The parameters of evidence accumulation models quantify different aspects of the decision process, such as the rate of evidence accumulation, response caution (the amount of evidence required for a response) and response bias (different caution for different responses). Variations among experimental conditions in these parameters, and in non-decision time, can provide insights into latent psychological processes beyond those available from traditional approaches, such as independent analyses of accuracy and mean RT.

Theories based on the idea of evidence accumulation have been successfully applied to many different paradigms, including: simple perceptual decisions (Usher & McClelland, 2001), visual short-term memory (Smith & Ratcliff, 2009), absolute identification (Brown, Marley, Donkin, & Heathcote, 2008), lexical decision

(Ratcliff, Gomez, & McKoon, 2004; Wagenmakers, Ratcliff, Gomez, & McKoon, 2008), the link between depression and anxiety (White, Ratcliff, Vasey, & McKoon, 2009, 2010), and the neural correlates of behavioral measures (Farrell, Ratcliff, Cherian, & Segraves, 2006; Forstmann et al., 2008; Ho, Brown, & Serences, 2009). Many different evidence accumulation models have been proposed, including Ratcliff’s diffusion model (Ratcliff, 1978), the Poisson counter model (Pike, 1966; Van Zandt, Colonius, & Proctor, 2000), the accumulator model (Smith & Vickers, 1988), the leaky competing accumulator model (Usher & McClelland, 2001), and ballistic accumulator models (Brown & Heathcote, 2005, 2008). We will focus on the recently proposed linear ballistic accumulator (LBA) model because it is mathematically simple, and because it was the model used by the authors of the data set we use as an example in this tutorial (Forstmann et al., 2008). Although our focus here is on the LBA model, the techniques we illustrate for model selection and evaluation are applicable to all evidence accumulation models.

Applying an RT model to data involves – at minimum – estimating parameters from data. Brown and Heathcote (2008) and Donkin, Averell, Brown, and Heathcote (2009) provide computational routines for LBA parameter estimation. Similarly, Vandekerckhove and Tuerlinckx (2007) provide methods and advice for estimating the parameters of Ratcliff’s (1978) diffusion model (see also Tuerlinckx, 2004; Tuerlinckx, Maris, Ratcliff, & De Boeck, 2001; Vandekerckhove & Tuerlinckx, 2008). More generally, Myung (2003) and Van Zandt (2000) provide excellent tutorials on how to estimate parameters for psychological models. However, when using a choice RT model, it is not a trivial step to go from

* Corresponding author.

E-mail address: christopher.donkin@gmail.com (C. Donkin).

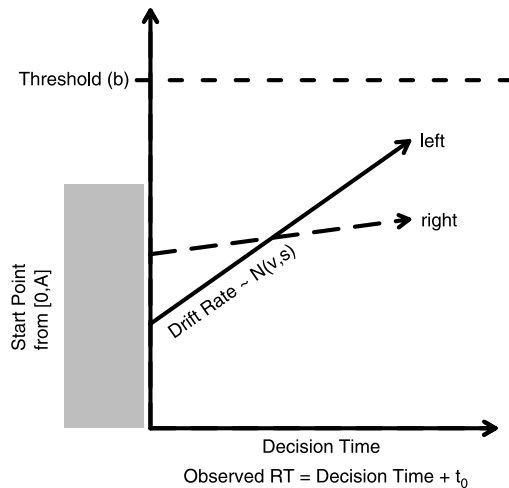


Fig. 1. A typical LBA decision for the task in Forstmann et al. (2008). In the illustrated trial, a left-moving stimulus has been presented and so drift rates for the left and right accumulators have been sampled normal distributions with means v and $1 - v$, respectively, and a common standard deviation s .

estimating parameters to drawing psychologically meaningful conclusions. The aim of the current tutorial is to bridge the gap between estimating parameters and interpreting data. We present a step-by-step analysis of data from a simple perceptual task (Forstmann et al., 2008) to illustrate this process. The aim of this tutorial is to provide new users with guidance about conventions and assumptions that are often not reported, or only briefly reported, in applications of choice RT models.

2. The linear ballistic accumulator

Fig. 1 illustrates decision processing in a pair of LBA units. Suppose that the figure represents a single trial in Forstmann et al.'s (2008) experiment, in which participants must choose whether a cloud of dots appears to be moving to the left or to the right, requiring a “left” or “right” response, respectively. Presentation of the stimulus causes evidence to accumulate for both the “left” and “right” responses separately, as indicated by the two lines (one solid and one dotted) in Fig. 1. The vertical axis of the figure represents the amount of evidence that has been accumulated, and the horizontal axis shows how much decision time has passed. The amount of evidence in each accumulator increases linearly with time. The choice made corresponds to the accumulator whose evidence total reaches the response threshold first. Decision time corresponds to the time taken for that accumulator to reach threshold. The predicted RT is the sum of decision time and non-decision time, quantified by parameter t_0 .

The slopes of the lines in Fig. 1 indicate the rates at which evidence is accumulated for each response, and are usually referred to as the drift rates. If the physical stimulus favors a “left” response, the drift rate for the “left” response accumulator will usually be larger than that for the “right” response accumulator. Drift rates are assumed to be set by physical stimulus properties and by the demands of the task. For example, in Forstmann et al.'s (2008) task, a correct “left” decision is made easier by making the displayed dots drift more steadily to the left. This would provide more evidence that “left” was the correct response, and so the drift rate for that response would increase. Drift rates are also assumed to be modulated by sensory and attentional processing, and the overall efficiency of the cognitive system. For example, Schmiedek, Oberauer, Wilhelm, Süß, and Wittmann (2007) found larger drift rates for participants with higher working memory capacity and fluid intelligence. In the LBA there is one drift rate for each

accumulator, corresponding in this application to “left” and “right” responses. The relative size of drift rate parameters describes differences in task performance between different conditions or groups. Although not explicitly illustrated in Fig. 1, drift rates in the LBA are assumed to vary randomly and independently between accumulators from trial-to-trial according to a normal distribution with mean v and standard deviation s , reflecting trial-to-trial fluctuations in factors such as attention.

The amount of evidence in each accumulator before the beginning of the decision process also varies from trial-to-trial. The starting evidence for each accumulator is assumed to follow a uniform distribution whose minimum value is set (without loss of generality) at zero evidence for all accumulators, and whose upper value is determined by a parameter A . Hence, the average amount (across trials) of evidence in each accumulator before accumulation begins is $\frac{A}{2}$. The response threshold is quantified by parameter b , represented by the horizontal dotted line in Fig. 1. The value of b is constrained to be greater than A so that a response cannot be made without accumulating some evidence. The difference $(b - \frac{A}{2})$ provides a measure of average “response caution”, as it is the average amount of evidence that must be accumulated to trigger a response. In Fig. 1, both accumulators have the same b and A parameters so the same amount of evidence is required, on average, before either response is made. Participants can choose to favor one particular response (i.e., a response bias), by setting a smaller value of b for the corresponding accumulator. Such response bias leads to a speed-accuracy trade-off, as the preferred response is made more quickly, but it is also made more often when incorrect, reducing accuracy. Bias towards a response by a particular accumulator can also be caused by increasing its A parameter, but changes in the A parameter are not usually assumed to be under the participant's control.

The time taken for each accumulator to reach threshold on any given trial is the distance between the response threshold and the start point of activation, divided by the rate of evidence accumulation. The observed decision time on any given trial, however, is the time for the fastest accumulator to reach threshold. The formula for the distribution across trials of the time taken for the fastest accumulator to reach threshold is given by Brown and Heathcote (2008). This formula makes it possible to estimate the model's parameters from data.

2.1. Example LBA application

Choice RT models are most appropriate for paradigms requiring simple and rapid decisions.¹ Forstmann et al.'s (2008) participants made simple decisions with average RTs around one second, so the paradigm is appropriate. Their experiment investigated the neural correlates of the trade-off between speed and accuracy, by testing predictions from a neurophysiological theory of how response caution is implemented by sub-cortical decision circuits. They presented participants with a cloud of 120 moving dots, of which 60% moved coherently to either the left or right of the screen, while the remaining 40% moved in random directions. Participants were asked in which direction (either “left” or “right”) the cloud appeared to move. Several seconds before the decision stimulus participants were given one of three cues, indicating whether they should try to make a very accurate response (accuracy emphasis), or a very fast response (speed emphasis), or try to balance accuracy and speed (neutral emphasis). Twenty participants each completed 280 trials per emphasis condition; other methodological details can be found in the original article.

¹ Although similar models have been extended to more complicated judgments (e.g., Bussemeyer & Townsend, 1992, 1993).

The manipulation of response caution had the expected effect. On average, participants were faster under speed emphasis ($\bar{RT} = 429$ ms) than under neutral emphasis ($\bar{RT} = 515$ ms) or accuracy emphasis ($\bar{RT} = 555$ ms). The faster responses came at the cost of lower accuracy: in the speed condition 77% of responses were correct, whereas in the neutral and accuracy conditions 86% and 87% of responses, respectively, were correct. This data pattern – trading accuracy for speed – is consistent with the effects of manipulating response caution in a choice response model (i.e., moving the response threshold higher and lower). However, it is also possible that participants were doing something more complicated. For example, non-decision processes (t_0) might also have been faster under speed emphasis, or the quality of information (drift rate, v) might have been greater under accuracy emphasis. Forstmann et al. (2008) examined these possibilities by comparing the fit of the LBA model using a range of different parameter constraints. This analysis allowed them to infer which cognitive processes were influenced by the experimental manipulation. In the next section we address in detail the problem of selecting the best set of parameter constraints. First, however, we briefly review parameter estimation for choice RT models and some other assumed knowledge.

3. Fitting the model

3.1. Parameter estimation

A choice RT model, like any quantitative theory, is defined by numerical parameters, and changing these parameters changes the model's predictions about RT and accuracy. For example, increasing the response threshold parameter increases accuracy and both slows and increases the variability of predicted RTs. Increasing the drift rate also increases accuracy, but it has the opposite effect on RT, reducing both its mean and variability. Non-decision time affects mean RT, but has no effect on RT variability or on accuracy.

The initial aim of fitting a model is to find parameter values which yield model predictions that adequately match the observed data. The degree of match is quantified by an objective function, which takes into account both decision accuracy and the distribution of RTs for each type of response. Automated search algorithms are used to find the best-fitting parameter values – those that optimize the objective function – given a particular set of parameter constraints.

There is a vast amount of literature, including several tutorials, dealing with the choice of objective function and optimization algorithm (e.g., Heathcote, Brown, & Mewhort, 2002; Myung, 2003; Ratcliff & Tuerlinckx, 2002; Van Zandt, 2000). For the purpose of this tutorial we assume that the reader has a reasonable grasp of the issues associated with parameter estimation. Our starting assumption is that the reader is capable of finding the best-fitting values of a set of parameters, defined by a particular set of parametric constraints. The Appendix contains a general review of issues related to fitting and parameter estimation for choice RT models, with a particular emphasis on the model and fitting methods used here. Readers less familiar with the LBA model might also benefit from reviewing the software and methods described by Donkin et al. (2009).

The purpose of this tutorial is to go beyond finding the best estimates for free parameters by describing how to find the best set of free parameters to estimate. This issue is critical since choice RT models are often used to draw inferences about which cognitive processes are influenced by experimental manipulations. These conclusions are drawn by determining which parameters of the model systematically vary across a set of conditions produced by experimental manipulations. In most applications of choice RT

models, the authors present only the best parameterization of a model—the smallest set of parameters needed to vary across experimental conditions in order to account for the data. Little discussion, however, is generally given to the many assumptions and decisions which yield these parameters. The aim of this tutorial is to go step-by-step through the process of identifying a best set of parameters. To illustrate, we analyze Forstmann et al.'s (2008) data, discussing many of the issues regarding the selection of which parameters *can* and *should* change across experimental conditions.

3.2. Which parameters change across conditions?

Forstmann et al.'s (2008) experiment had three emphasis conditions (speed, neutral and accuracy), and in each of these conditions there were two types of stimuli (coherent motion to the left or to the right). One of the central tasks of cognitive modeling for these data is to investigate which aspects of cognitive processing were influenced by the experimental factors. In model terms, we want to know which parameters changed across each condition.

3.2.1. A priori assumptions

To begin, we first decide which parameters potentially *could* change. The LBA model has five parameters that determine behavior in any condition (b, A, s, t_0, v). When there are two choices there are two accumulators, one for each decision (e.g., one corresponding to the response “left” and one to the response “right”). This means that there could be up to 10 parameters that vary for each particular combination of stimulus and emphasis conditions, for a total of 60 parameters. Fortunately, this type of freedom, though possible, is not usually required, because sensible *a priori* constraints can be placed on parameters across conditions. We elaborate these constraints by considering three factors in succession: “left” vs. “right” responses; left-moving vs. right-moving stimuli; and decision caution conditions (speed, neutral or accuracy emphasis).

The two possible responses (“left” and “right”, corresponding to the two accumulators) should share many parameters. Usually, t_0 can be fixed at the same value for both because, for example, in most cases it is reasonable to assume that the time to execute each response is the same. This assumption is plausible for the present data, but may break down in unusual paradigms, such as when one response is harder to produce than another. In contrast, the evidence threshold parameter (b) and the starting point distribution parameter (A) might reasonably vary between responses—because participants might be biased toward one response over the other. For example, if participants are biased to respond “left” rather than “right”, this can be reflected in a smaller value of b and/or a larger value of A for the “left” accumulator than the “right” accumulator.

Response biases may or may not occur, depending on individual differences between participants. However, when choice accuracy is above chance, every participant ought to demonstrate a difference in drift rates between “left” and “right” accumulators as a function of which response is correct for a given stimulus. For example, on trials where the stimulus drifts to the right, the mean drift rate should be higher for the accumulator corresponding to the “right” response than for the “left” response. Often, greater simplification can be obtained by fixing the mean drift rate for the incorrect response at one minus the mean drift rate for the correct response. This restriction has commonly been applied because it also satisfies a scaling property applying to all choice RT models, which requires at least one parameter to be fixed in order to obtain unique estimates of the remaining parameters. However, applying this restriction to drift rates across *all* conditions provides greater constraint than necessary to satisfy the scaling property, and can

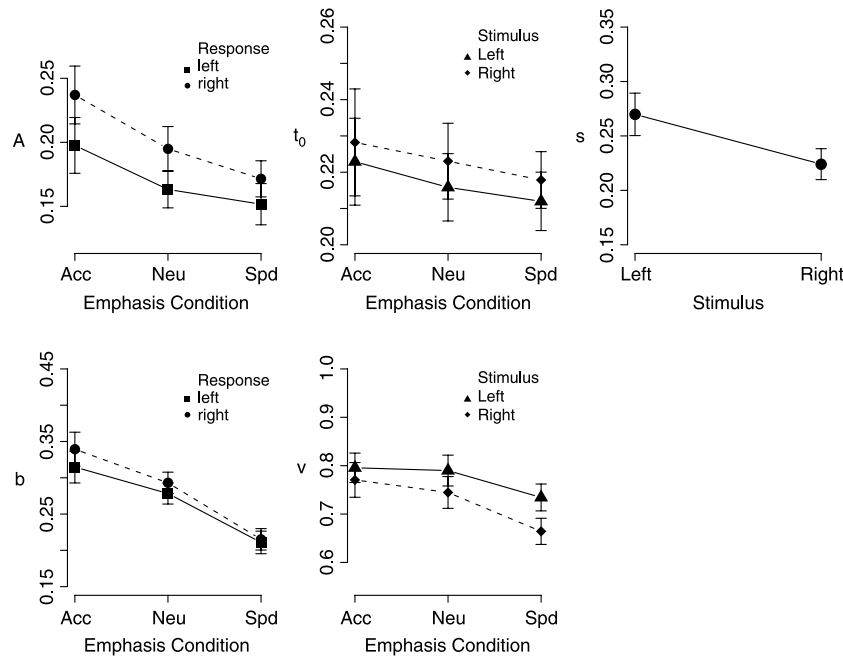


Fig. 2. Parameter estimates averaged over participants across emphasis conditions, responses and stimuli. Error bars are ± 1 standard error.

result in poor fits. We advise careful consideration about whether this restriction is justifiable on theoretical grounds, before applying it (for further discussion see Donkin, Brown, & Heathcote, 2009).²

Finally, though it is possible that between-trial variability in the drift rate, s , can differ across responses, it has been fixed to the same value in all applications of the LBA to date, and the same is true, to our knowledge, of analogous parameters in other evidence accumulation models. We follow this convention here, but note that this is an additional, and untested, assumption.

In summary, the only parameters that we allow to take on different values for “left” than “right” response options are b and A . We will further assume that v is constrained to sum to one across “left” and “right” response accumulators within any condition. This is a reasonable assumption for Forstmann et al.’s (2008) experiment, corresponding to the idea that increased evidence for one response (e.g., more dots moving left) necessarily implies decreased evidence for the other response (e.g., fewer dots moving right).

Next, consider which parameters could vary between left- and right-moving stimuli. In Forstmann et al.’s (2008) experiment, left- and right-moving stimuli were randomly ordered over trials. It is typically assumed that changing response threshold settings is a relatively slow process, so threshold parameters (b and A) do not depend on the stimulus presented for the current decision (Ratcliff, 1978). The general version of this principle is that b and A are kept fixed across conditions whenever the participant is unable to predict which of those conditions will occur next. Since the other parameters (v , s and t_0) are assumed to be influenced by stimulus properties, they should be free to vary across stimulus types. For example, left-moving stimuli might provide more salient motion cues than right-moving stimuli, which should be reflected in a higher left-moving stimulus drift rate.

Finally, consider which parameters might vary between speed, neutral and accuracy emphasis conditions. Following convention,

between-trial variability in drift rate (s) is usually fixed across experimental conditions, particularly those not stimulus-based—although this assumption is not strictly necessary. All other parameters (b , A , v and t_0) could feasibly be influenced by response emphasis. Indeed, this was the central question for Forstmann et al.’s (2008) data analysis—which cognitive processes (i.e., parameters) were influenced by the response caution manipulation?

Together, this relatively liberal set of constraints, based on conventions and theoretical plausibility, reduces the number of free parameters from 60 to 26. To keep notation compact, we subscript parameters differing between left-moving and right-moving stimuli with “left” and “right”, and use “L” and “R” subscripts to indicate parameters corresponding to evidence accumulators for “left” and “right” responses. With this notation, the 26 free parameters are: s_{left} , s_{right} and three sets of b_L , b_R , A_L , A_R , v_{left} , v_{right} , $t_{0\text{left}}$ and $t_{0\text{right}}$, one for each emphasis condition.

3.2.2. Which parameters need to change to fit the data?

We next assess which of the 26 free parameters are actually needed to fit the data. There are two ways this has been approached in the literature. The first is to fit the model to each participant’s data with all 26 free parameters, and then examine how parameter estimates differ across conditions. To demonstrate this approach, we fit each individual’s data from Forstmann et al.’s (2008) experiment using maximum likelihood estimation (MLE) and a SIMPLEX search algorithm (see the Appendix for computational details). As discussed in the Appendix, obtaining good parameter estimates for such a complex model (26 free parameters) is not easy—a drawback of this first approach.

Fig. 2 suggests some general ideas about which parameters need to vary across conditions. For example, on average both the drift rate (v) and its standard deviation across trials (s) were larger for left-moving than right-moving stimuli. However, the much smaller corresponding difference (relative to the standard error bars) in the non-decision time plot suggests that t_0 probably did not change between stimulus types. Similarly, the evidence threshold (b) and start point variability (A) parameters did not change much between “left” and “right” responses. These two parameters, however, changed substantially between the three response caution conditions (left to right across the plots). In

² An even greater restriction is often applied: the mean drift rate for correct responses is constrained to be less than one. This ensures that the mean drift rate for incorrect responses is greater than zero, but again the appropriateness of this further restriction must be examined.

contrast, non-decision time and drift rate showed much smaller changes between emphasis conditions.

Differences between average parameter estimates can be tested for statistical reliability using a repeated measures analysis of variance (ANOVA). The results of these tests help to decide which parameters were affected by which manipulations. However, such tests bear only on the question of whether the population means for the parameters vary between conditions. It is possible that there is no difference in population means between conditions and yet each individual differs systematically between conditions in a way that cancels out on average. In this case fixing parameters to be the same over conditions in all individual participant fits may distort the selection of the best set of parameter constraints (i.e., model selection). A full solution to this dilemma requires a “random effects” approach, which produces explicit estimates of population means and variability for each parameter type by fitting the model to all data from a group of participants simultaneously (Averell & Heathcote, *in press*; Lee, *in press*; Morey, *in press*; Pratte & Rouder, *in press*). Vandekerckhove, Tuerlinckx, and Lee (*in press*) develop this approach for Ratcliff’s diffusion and we are presently doing the same for the LBA (see Donkin et al., 2009). However, random effects models impose a greater computational burden, and so in this tutorial we focus on methods based on fitting each participant’s data separately.

Given the limitations of this initial individual participant free-fitting method, we recommend it to be augmented with a second method based on sequential model building. This method also uses individual analysis, but it is still sensitive to the need to allow for individual differences. The key to this approach is to fit many different versions of the model, beginning with the simplest version (identical parameters for all conditions; only five free parameters in our example) and ending with the most complex (with 26 free parameters in our example). This approach can be computationally demanding because there might be very many intermediate models to analyze. The intermediate models are formed by considering all factorial combinations of parameter constraints. For example, after estimating the simplest model, one might next estimate a model where drift rate was free to vary between left-moving and right-moving stimuli. After that, both of those first two models would be used to start parameter searches for even more complex models, perhaps with the boundary separation parameter free to vary across speed/accuracy emphasis conditions. This process continues through to the most complex model.

After estimating the parameters of all the intermediate models, one model is selected that best satisfies the trade-off between simplicity and goodness-of-fit. When each intermediate model is nested within a more complex version, as is assumed here, each model’s goodness-of-fit must necessarily improve with the number of estimated parameters. However, when the improvement is small it may be due to “over-fitting”, where the extra parameters serve only to account for unsystematic variation in the data. Such over-fitting is undesirable because it leads to a model that predicts new data poorly, and may produce theoretically misleading patterns of parameter estimates (for additional details see the following special issues on model selection: Myung, Forster, & Browne, 2000; Wagenmakers & Waldorp, 2006).

One method of identifying over-fitting is to choose the model with the smallest AIC (Akaike Information Criterion, Akaike, 1974) or BIC (the Bayesian Information Criterion, Schwarz, 1978). Parameter estimation using maximum likelihood is most appropriate for this purpose as both information criteria are calculated by adding a penalty to minus twice the log-likelihood of the model. The penalty quantifies model complexity based on the number of estimated model parameters, k ($2k$ for AIC and $\log(N) \times k$ for BIC, where N is the number of data points).

We use BIC in our application, as the AIC prefers overly complex models in large samples, although we acknowledge that this preference is debatable. Both methods are limited because they do not take account of differences in functional form complexity (Pitt & Myung, 2002) between models: that is, both statistics treat all parameters equally in terms of model complexity, but this may not be true. For example, even when two models have the same number of parameters one model may have more flexibility in fitting data due to differences in the way that the models restrict interactions amongst parameters. Model selection methods that address this issue, such as the Deviance Information Criterion (DIC, Spiegelhalter, Best, Carlin, & van der Linde, 2002) require Bayesian estimation, which we do not address here (see Donkin et al., 2009, for details of Bayesian LBA estimation).

The ideal data driven version of this second, nested model, approach requires fitting of all factorial combinations of restrictions on the 26 parameters. However, that is not always feasible because it can require estimating parameters for many thousands of models. If the computational load is too great, one can make an initial simplification by fixing parameters whenever the free estimates (from Fig. 2) strongly suggest that those parameters do not change across conditions. For example, non-decision time does not appear to be influenced much by either different stimulus or emphasis conditions, suggesting that just one t_0 estimate will do for all six conditions. Similarly, response threshold and start point variability appear not to vary across emphasis conditions but not across responses. Based on these observations we can narrow down the options to constrain the most complicated model to one with 15 free parameters: b , A , v_{left} and v_{right} varying across the three emphasis conditions, s_{right} and s_{left} , and t_0 fixed across all conditions. It is important at this point to remember the limitations, discussed above, of making inferences about the population parameters based on average parameter estimates from the minimally constrained model. These limitations mean that the shortcut method we used to move from the 26-parameter model to the 15-parameter model should only be employed when the computational burden associated with exhaustively estimating the intermediate models is too great.

We also note that the interpretation of Fig. 2 and the subsequent choices about which parameters should be fixed have a subjective element. For example, we chose to fix A across responses but let it vary across emphasis conditions.³ It is certainly arguable from the upper left plot in Fig. 2 that other ways of constraining the A parameter are plausible. A more formal method would be to perform ANOVA on the parameter estimates, but we are cautious about recommending the blind application of this approach given its inherent limitations in terms of providing positive evidence in favor of a null difference. Instead, when it is not clear whether a parameter might vary across conditions, it is best to use the methods outlined in the next section to further investigate.

3.3. Model selection example

The 15 free parameters that may account for our data are generated from five ways that parameters vary across experimental conditions; b , A , and v vary over emphasis condition, and v and s also vary across left- vs. right-moving stimuli. We will call these variations the five “features” of the model. One way to determine which features are required by the data is to fit all 32 possible combinations of models made up of these features, i.e. one model with all five features, the five models with four features, the

³ Further investigation (not reported here, using the methods outlined below, confirmed that A should be fixed across responses.

Table 1

The most complex, and the BIC-best model: which parameters varied across responses, stimuli, and emphasis conditions; also the number of parameters (k), and total BIC.

Model	Factor					k	Σ BIC	
	Responses	Stimuli	Emphasis					
Most complex	–	v	s	b	A	v	15	–15 226
Best intermediate	–	v	–	b	–	–	8	–15 653

10 models with three features, etc., and select the best model using BIC values. This method is by far the most comprehensive, and indeed could be extended to all nine parameter variation features that were earlier deemed plausible. Such an approach using all nine features would require 512 models to be fit separately to each participant's data. Forstmann et al.'s (2008) model selection analysis was based on an exhaustive evaluation of set of models of this size, although with slightly different features.⁴ It is a matter of judgment for the individual researcher to balance the computational burden of exhaustive search using all features against the subjectivity of identifying short-cuts to reduce the number of features.

We estimated the parameters for all 32 possible combinations of models made up of the five features. Models were fit individually to each of the 20 participants, and for each participant a set of parameters and a corresponding BIC were obtained. We use BIC summed over participants, which we will call "total BIC", to describe group-level results. The results of the total BIC analysis for the most complex model and for the intermediate model which yielded the smallest total BIC value are reported in Table 1. The most complex model, with all five features, has a much larger BIC value (–15 226) than the best intermediate model (–15 653) indicating that this intermediate model provides a much better compromise between goodness-of-fit and model complexity.⁵ The best intermediate model has only two features, and eight parameters. Averaged across participants, those parameters were: $v_{\text{left}} = 0.72 \text{ s}^{-1}$ and $v_{\text{right}} = 0.67 \text{ s}^{-1}$, $b_{\text{acc}} = 0.29$, $b_{\text{neu}} = 0.27$, $b_{\text{speed}} = 0.17$, $s = 0.22 \text{ s}^{-1}$, $A = 0.15$ and $t_0 = 0.11 \text{ s}$.

The drift rate estimates for the eight-parameter model suggest that the participants studied by Forstmann et al. (2008) were able to extract information from left-moving stimuli around 8% faster than for right-moving stimuli. The eight-parameter model also implies that the response caution manipulation affected only one cognitive process: the amount of evidence required before responding. Relative to the neutral condition, participants set evidence thresholds 7% higher under accuracy emphasis and 37% lower under speed emphasis. Note that the final model chosen in the present analysis is very similar to that selected in the original paper, except that the model selected here equated between-trial drift rate variability for left- and right-moving stimuli (this possibility was not examined in the original analysis). Our results are thus consistent with the major conclusion of the original paper, that the emphasis instructions selectively affect the response threshold.

4. Evaluating and presenting model fit

Another important model selection criterion is the descriptive adequacy of the model, which can be assessed graphically. A model is

inadequate if it fails to describe theoretically important patterns in the data. Similarly, if the parameter estimates vary across conditions in ways that make no psychological sense, the model is suspect in terms of its theoretical adequacy. The average parameter estimates for the best intermediate model given in the last section appear to be adequate on the latter grounds, as did the corresponding parameter estimates for all individuals. In this section we describe how to graphically check model adequacy.

The match between model and data should be assessed for each individual participant. However, the final communication of results almost always requires a summary of the grouped data. Such averages can fail to represent the individual participants, depending on how they are constructed. As an extreme example, suppose that an experiment had just two participants, one who responded very quickly and another who responded very slowly. In this case, an "average" histogram formed by pooling participant data could be bimodal, and so not be representative of either individual.⁶ Because of this issue it is often better to first calculate statistics which summarize the RT distribution and then average those. Regardless of the method used, one should always check how well averaged data matches the individual participants.

The agreement between model and data is usually assessed by plotting together predicted and observed statistics that summarize RT distributions and response probabilities. Histograms depicting the observed RT distribution are often overlaid with the predicted probability density function (PDF) from the model, to assess model fit. Such plots are simple to interpret, but do not always highlight the shortcomings of the model. Cumulative probability plots (e.g. Forstmann et al., 2008), or quantile probability (QP) plots (e.g. Ratcliff & Smith, 2004), are more complicated to produce and read, but can better illustrate differences between the model and data. Group QP or cumulative probability plots, which are obtained by averaging quantiles for each individual, also have the advantage that they tend to be more representative of individual results (e.g., such averages do not suffer from the bimodality problem that can occur with histograms). To represent the model predictions using group plots, one calculates the model's predicted quantiles for each individual and averages these together in the same way as the data. This means that we apply the same averaging process to create summary information for model predictions as for the data, and so both summaries are equally subject to any distorting effects of averaging.

Fig. 3 summarizes Forstmann et al.'s (2008) data and the corresponding LBA model fits using group QP plots. QP plots are an efficient way of displaying the important information from a set of choice RT data—the horizontal axis displays response probability (accuracy) information and the vertical axis displays information about the RT distribution. There are six QP plots in Fig. 3, with each plot representing average results from a single experimental condition. Each plot contains two sets of vertically aligned points, illustrating the RT distributions for correct and incorrect responses from one experimental condition. The horizontal position of a set of vertically aligned points represents the proportion of responses making up that RT distribution. For example, in the top left panel of Fig. 3 the observed quantiles (solid squares) sit above 0.89 on the horizontal axis, indicating that on average 89% of responses were correct in that condition (left-moving stimuli under accuracy emphasis). Note also that this implies that 11% (i.e., 100%–89%) of responses were incorrect. Hence, the quantiles for these errors are displayed at 0.11 on the horizontal axis. In general, points to the left and right of 0.5 on a QP plot indicate incorrect and correct

⁴ Sequential model selection techniques, such as the forward, backward and stepwise methods commonly used in linear regression model selection, provide an alternative method of reducing computational cost. Such techniques could be applied to selection amongst choice RT models, either based on likelihood ratio tests, or based on BIC (Hoeting, Madigan, Raftery, & Volinsky, 1999), but as in linear regression they are not guaranteed to find the best model.

⁵ See Wagenmakers and Farrell (2004) for formal methods of comparing BIC differences that can be employed when results are not so clear cut.

⁶ Even though data grouped this way will not necessarily look like any individual's data, a similarly grouped graph of the model predictions still provides a valid assessment of model adequacy.

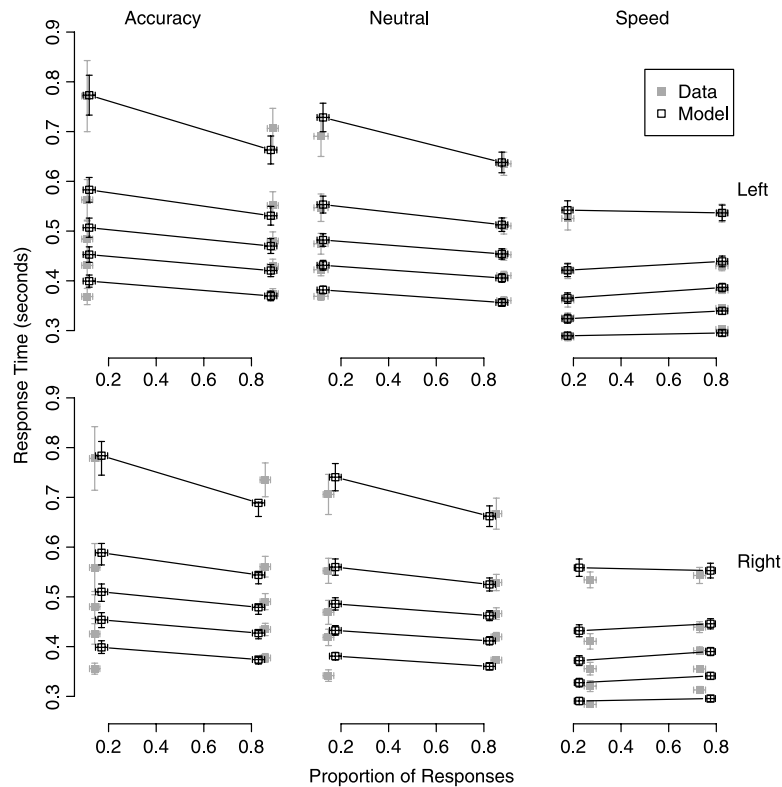


Fig. 3. A quantile probability plot for the data from Forstmann et al. (2008). Observed and predicted quantiles are represented by solid and open symbols, respectively. Responses to left-moving and right-moving stimuli are represented in the top and bottom rows, respectively. Accuracy, neutral and speed emphasis conditions are shown in the left, center and right columns, respectively. Error bars show standard errors across participants for both data and model predictions.

responses, respectively. The vertical positions of the points are determined by five quantile estimates (0.1, 0.3, 0.5, 0.7 and 0.9). For example, the 0.1 quantile estimate corresponds to the value below which 10% of the RT values in the distribution fall. Taken together, the five quantile values summarize the RT distribution. For example, the first filled square above 0.89 in the top left panel shows the 0.1 quantile for correct responses in that condition, the next square above this shows the 0.3 quantile, and so on. The unfilled squares provide the same information, but for the distributions predicted by the LBA model rather than for the observed data.

Note that QP plots can be constructed with any desired set of quantiles, such as with deciles or semi-deciles. Using more than five quantile estimates will provide a more detailed description of the RT distributions, but can also make the plots difficult to read. Similarly, results for more than one condition can be given in the same graph. This often works well when the conditions differ sufficiently in accuracy. For example, we could have given results for accuracy and speed conditions in the same panel. In contrast, the neutral and accuracy conditions are quite similar in accuracy, so providing results for these two conditions in the one plot made the QP plot hard to interpret.

Fig. 3 reveals the following general patterns: Responses in the speed emphasis conditions (right column of Fig. 3) are faster, as indicated by their lower position on the vertical axis, than in accuracy and neutral conditions (left and center columns, respectively). Responses for left-moving stimuli (top row) are more accurate than for right-moving stimuli (bottom row). This shows up in the figure because quantiles in the top row sit at more extreme horizontal positions than those in the bottom row. For example, quantiles for left stimuli in the speed condition are positioned at 0.18 and 0.82 on the horizontal axis, while the same quantiles for right stimuli sit above 0.27 and 0.73, indicating more incorrect and fewer correct responses for right-moving than

left-moving stimuli. The addition of lines joining corresponding quantiles for correct and incorrect responses in Fig. 3 highlights a theoretically important issue, the relative speed of correct and incorrect responses for different emphasis conditions. In these data, incorrect responses are generally faster than correct responses in the speed emphasis condition, but this difference is reversed in the accuracy emphasis condition. The model does a good job of accounting for these patterns. However, the QP plot also reveals shortcomings of the model, with the most evident being a tendency to predict too many incorrect responses for right-moving stimuli. If this failing were considered to be practically or theoretically important selection of a more complex model that addresses this issue might be warranted.

Producing a QP plot requires calculation of the 0.1, 0.3, 0.5, 0.7, and 0.9 quantiles for observed and predicted RT distributions. Quantile estimates from the observed data can be calculated using functions available in most statistical software (for more details see Heathcote et al., 2002; Van Zandt, 2000). Quantiles were calculated for each individual participant and then averaged together to create the observed quantiles in Fig. 3. Note that it is important to check that the summary information to be presented in a QP plot is representative of individuals. In these data, more than 80% of individual quantile estimates were within 50 ms of their respective average values, suggesting that our averages were representative of individual RT distributions.

Calculating the quantile values predicted by the model is a little more difficult. There are two standard approaches: either using a search algorithm to invert the cumulative distribution function (CDF) of the model, or via simulation. To generate the predictions shown in Fig. 3 we used the conceptually simpler, but computationally more expensive, simulation method (see Appendix for details on the search method). To calculate predicted quantiles via simulation we took each individual's best-fitting parameters and used them to sample one million data points in

each condition. Note that LBA simulation is computationally very cheap (see Donkin et al., 2009); this is not necessarily the case for other choice RT models (e.g., see Brown, Ratcliff, & Smith, 2006, for details of simulating the Ratcliff diffusion model). The simulated data followed the exact same design as the empirical data—i.e., three emphasis conditions and two stimulus conditions, where only drift rate changes over stimulus conditions and only response threshold changes over emphasis conditions. Finally, we calculated quantiles from these simulated data, and averaged across participants in the same way as for the observed data.

Rather than plotting the predicted quantiles averaged over individuals, Ratcliff and colleagues suggest fitting the model to the average observed quantiles to create model predictions for QP plots (e.g., Ratcliff, 2002; Ratcliff et al., 2004). This approach can appear to indicate a better fit than the method we describe here, since model predictions will be based on parameters which optimize that fit. However, one risk with this approach is that the newly estimated parameters may not be representative of the parameters of any individual. For this reason, Ratcliff and colleagues always assess how closely these new parameters match the average individual participant parameters. Further, a quantile-based objective function must be used for estimating parameters from the average observed quantiles; MLE cannot be used (see Appendix).

5. Discussion

In recent years, the once-difficult process of obtaining the parameters for choice RT models has been made much easier by the provision of software that automates the parameter estimation process. Our aim was to build on these developments by providing advice about, and a detailed example of, the many extra steps involved in moving from simple parameter estimation to a more meaningful analysis. We focused on an application of the LBA model (Brown & Heathcote, 2008) to data reported by Forstmann et al. (2008). We illustrated the canonical problems in such modeling, by first describing how there are – potentially – 60 free parameters even for Forstmann et al.'s quite simple experiment. We then illustrated how this number can be reduced to 26 free parameters in our example by *a priori* considerations. Exploratory analysis of the 26-parameter model identified several parameters that, on average, did not change substantially across experimental conditions. This led to an even simpler model with 15 parameters. Finally, we exhaustively fit 32 versions of the 15-parameter model and selected a final model with 8 free parameters that provided the best trade-off between goodness-of-fit and model complexity.

We also showed how to check the descriptive adequacy of the final model using QP plots. The selected model provided a good fit that captured theoretically important features of the data, and is consistent with Forstmann et al.'s (2008) conclusion that a manipulation of response emphasis selectively influenced the amount of evidence required for a decision. The same conclusion has also been made based on applications of the Ratcliff diffusion model to data from similar paradigms (e.g., Ratcliff & Rouder, 1998; see Donkin, Brown, Heathcote, & Wagenmakers, *in press* for a detailed analysis of the relationship between the parameters of these two models).

The model selection process we described relies on fixing some parameters across different conditions. In a between-subjects manipulation, different conditions are populated by different people, meaning that certain parameters would have to be fixed across participants—this requires modeling random effects. Most often for choice RT analyses, this problem is handled by estimating model parameters separately for individual subjects, then using standard null hypothesis significance testing (NHST) to determine which parameters vary across the between-subjects conditions.

As an example, imagine that Forstmann et al. (2008) had tested both an older and a younger group of participants and thus had an additional between-subject factor. Standard practice would be to fit each individual from both the younger and older groups with the model we previously selected, giving observed distributions of each parameter for younger and older participants. NHST inferential tests could then be used to determine whether the average of certain parameters differ between younger and older groups. For example, an independent sample *t*-test could be used to determine whether the average non-decision time parameter is different for older and younger participants. This approach has been used to identify the effects of aging on decision processes (e.g., Ratcliff, Thapar, Gomez, & McKoon, 2004; Ratcliff, Thapar, & McKoon, 2007).

The method of using NHST to determine differences between subjects carries with it all of the usual drawbacks. These may be particularly problematic for the application of choice RT models because the typical question is whether some parameter does *not* change across conditions. For example, Ratcliff and colleagues often find that old and young participants do not differ significantly in their drift rate parameters. It is difficult to know if this lack of significance is due to power or whether drift rate is truly equal for old and young participants. Bayesian hypothesis tests such as the Savage–Dickey test (Wagenmakers, Lodewyckx, Kuriyal, & Grasman, 2010) allow for direct assessment of the truth of the null hypothesis but require the posterior distribution of model parameters. Donkin et al. (2009) and Vandekerckhove et al. (*in press*) offer software for producing these distributions, but at present the approach is limited by its computational cost.

Random effects modeling, in which parametric distributions are used to describe subject-level parameters, provides another way to investigate differences across both between-subject conditions. For example, hyper-parameters describing the distribution of individual parameter estimates within younger and older groups can be estimated and compared. However, most researchers still use the approaches presented in this tutorial because they are many orders of magnitude faster than Bayesian Markov chain Monte Carlo methods typically required for random effects estimation. Further, we believe that it is prudent to preface random effects modeling with individual fitting. For example, Pinheiro and Bates (2000) recommend ‘An “inside-out” model building approach ... starting with individual fits ... to decide on the random effects structure’ (p. 133). Consequently, we expect the individual analysis approach of this tutorial to remain relevant. Further, the central issue that this tutorial has addressed – choosing a set of model constraints from among many possible sets – applies equally to all approaches.

Appendix. Additional details

A.1. The objective function

A widely used objective function is the likelihood of a model with parameters θ given data $\mathbf{x} : L(\theta|\mathbf{x})$. Finding parameter estimates by optimizing this objective function is called maximum likelihood estimation (MLE, see Myung, 2003, for a tutorial on MLE methods). For given parameters, a choice RT model defines a joint density over response (i.e., the choice) and the response time. This density function is used to define the likelihood function, so that MLE naturally takes into account both accuracy and RT information (see Brown & Heathcote, 2008, for details of the LBA model's density functions).

MLE is a default choice in many areas of statistics because it is unbiased for large samples, and because no other method is more efficient, as long as certain regularity conditions on the model are satisfied. However, choice RT models do not usually satisfy these

conditions because they predict distributions whose support is determined by an estimated parameter, t_0 (Heathcote, Brown, & Cousineau, 2004). This can cause maximum likelihood methods to spuriously estimate t_0 as equal to the minimum RT in a data sample, with concomitant mis-estimation of the other parameters. Although it is important to be aware of this problem, it can usually be avoided by censoring implausibly fast RT data (e.g., responses faster than 200 ms, which are likely the result of anticipation). Slow outliers, due to processes such as distraction, are more problematic as they are harder to detect than fast outliers. Heathcote et al. (2002) showed that – even when estimating simple and regular RT models – an estimation method based on data quantiles (quantile maximum probability estimation, QMPE; see also Heathcote & Brown, 2004) could be more efficient and less biased than MLE in small samples.

Appropriately selected quantiles can summarize an RT distribution, and more quantiles lead to a more accurate summary. There are several objective functions that use quantiles to summarize the observed RT distributions, and compare these against model predictions. Besides QMPE, these functions include the Kolmogorov–Smirnov statistic (Voss, Rothermund, & Voss, 2004; Voss & Voss, 2007), χ^2 (Ratcliff, 2002), and weighted least squares error (Ratcliff & Tuerlinckx, 2002). In computational terms, the quantile-based objective functions require evaluation of the model's cumulative distribution function (CDF); this is different from MLE which requires evaluation of the model's PDF. The difference means that quantile-based methods are especially useful for models that have easy-to-use algorithms to calculate their CDF, but not PDF (such as Ratcliff's diffusion model).

With the exception of the Kolmogorov–Smirnov based approach, choice RT modelers have mostly used a coarse set of five quantiles: 0.1, 0.3, 0.5, 0.7, and 0.9 (Ratcliff & Tuerlinckx, 2002), which we will describe as the “standard” quantile set. Summarizing the observed RT distributions with such a coarse set has the advantage that fitting is only weakly influenced by fast and slow outliers. For example, even if 2% of the data were from a fast-guessing contaminant process, these would all fall below the 10% quantile estimate, which would thus be only mildly affected. Note, however, that there is no necessity for a coarse set of quantiles between the smallest and largest values to gain this advantage in robustness, and that Heathcote and Brown (2004) found that the advantages of QMPE in small samples only emerged with fine-grained quantile sets.

Forstmann et al. (2008) analyzed their data using QMPE. Above, we reported analyses of the same data using MLE, and we found that the two approaches agreed closely in this data set. In general, we have found that, as long as sensible precautions are enforced pertaining to outliers, and good starting points are used, MLE performs very well for the LBA model. MLE can also enjoy a substantial advantage over quantile-based methods in terms of computational speed. A further advantage is that MLE produces maximized likelihood values, which can be useful for performing the model selection analyses discussed in the main body of the paper.

A.2. Finding optimal parameters

A variety of optimization methods are available, but in our experience these algorithms differ mostly in speed and numerical stability rather than in their ability to find the best set of parameters. Here we use the SIMPLEX algorithm (Nelder & Mead, 1965), which is the most commonly used search algorithm for choice RT fitting because of its ease of use. More computationally efficient optimization algorithms, which require analytic derivatives of the objective function, are not generally used because the required derivatives are not easily

available. Other algorithms operate by numerically estimating derivatives, which can improve efficiency, but also decrease numerical stability. We recommend that users explore these different optimization approaches and then use the method, or combination of methods, that is both fast and stable in their application.

All parameter search algorithms, such as SIMPLEX, need a set of parameters to begin their search, and the consequences of a poor set of starting parameters can be dire—the parameter search can get stuck on an estimate that matches the data better than all nearby parameter sets, but which is much worse than estimates further away (a “local optimum”). Hence, the search should always begin with parameters that produce model predictions that are reasonably close to the data. Identifying such starting values is a difficult problem in itself, with no general solution. One easy way to generate start points is to use parameters which have been reliably shown to produce reasonable predictions for similar choice RT data sets. For example, Matzke and Wagenmakers (2009) provide the average parameter values for the diffusion model based on a large number of fits of the model to data, and Donkin et al. (in press) provide equivalent values for the LBA. A number of searches may then be run from a range of start points obtained by randomly perturbing the initial start point.

A second method is to use heuristics that obtain start points based on the data to be fit. We have found the following set of heuristics useful for the LBA model. We first set the drift rate distribution parameters: standard deviation, $s = 0.3$, and mean, $v = \frac{1}{2} + \Phi^{-1}(p)$, where p is the probability of making the response for this accumulator, and Φ is a normal CDF with mean 0 and standard deviation $s\sqrt{2}$. For t_0 we use 9/10 of the value of the minimum RT from the data. For the maximum of the uniform start point distribution, A , we take twice the inter-quartile range of the RT distribution, and finally we set the response threshold, b , at $1.25 \times A$. The heuristic values are calculated separately for each experimental condition, and then averaged over conditions for parameters that are constrained to be equal across conditions. Again a number of searches may then be run from a range of start points obtained by randomly perturbing the heuristically obtained start point.

Even when a good start point is obtained it is also possible that the search algorithm may terminate search prematurely. When using the SIMPLEX algorithm, this can sometimes be avoided by performing repeated searches, with each new search using the best estimate found by the last fit as its starting point. This method can be effective because each new search typically starts with a large simplex that explores parameter values relatively distant from those explored in the final stages of the previous search.⁷

Problems with poor parameter estimation become more severe as more free parameters are estimated. In terms of the methods we describe above, this means that parameter estimation can be very difficult for the most complex models, which might unfairly disadvantage those models. Fortunately, the nested model building approach that we advocate provides a natural solution to this problem—start points for more complex models can be generated from parameter estimates for simpler (nested) models. For example, consider data from a 2×3 factorial design. Suppose the two factors, A and B , are both assumed to affect drift rate. We might first fit a simple model with just one drift rate for all

⁷ The performance of the SIMPLEX algorithm degrades markedly when the number of parameters being optimized is large. We have generally found adequate performance with up to 20–30 parameters, at least with repeated fitting, and when a large number of iterations is employed—about 500 times the number of free parameters. Beyond 30 parameters, specialized search algorithms designed for dealing with high dimensional search spaces may be required.

six conditions, starting from parameters obtained by averaging heuristic estimates based on data from each condition. Suppose this results in a best-fitting drift rate of 1. We could then fit a model in which drift rate varied across the two levels of factor A, say A1 and A2, using 1 as the start point for both levels. Suppose the best-fitting parameters turned out to be 0.5 and 1.5 for A1 and A2, respectively. The two best-fitting parameters can then serve to create start points for the full factorial model in which drift rate varies over both factors A and B (i.e., the start points for the three levels of factor B in A1 would be 0.5, and the start points for the three levels of A2 would be 1.5). Similarly, we could also fit an intermediate model in which drift rate varies over only the levels of B, then use that to provide a second set of start points for the full factorial model. Although we have found this method to be very effective with complex models, there is still no guarantee that the best parameter estimates will be found. Hence, it is important to check the quality of fits graphically, as we described in the main body of this paper, and to try different starting values in order to see if improvements can be obtained for any poor quality fits.

As well as checking goodness-of-fit graphically, the parameter estimates themselves should also be checked for *a priori* plausibility. Plausibility can be judged relative to typical parameter ranges (e.g., for the LBA see Donkin et al., in press). When fitting data from a number of participants, consistency across participants can also be assessed. It sometimes happens that these checks reveal wildly large or small estimates of a particular parameter. This usually occurs when the value of that parameter receives little constraint from the data (i.e., large changes in its value lead to small changes in the objective function). This lack of constraint can be seen by making a graph of the objective function for a range of values of the suspect parameter around its estimated value while keeping the remaining parameter values fixed at their estimated values; sometimes called a “profile plot”. A flat profile plot indicates a poorly constrained parameter.

For the LBA, poor constraint is most commonly associated with the drift rate parameter for incorrect responses. This is particularly the case where observed accuracy is high, as there are then very few incorrect responses to constrain the estimate of this parameter. This problem did not occur with our example data, as Forstmann et al.'s (2008) participants made several dozen incorrect responses even in the most accurate condition. However, in paradigms where accuracy is very high, prohibitively many trials may be required to obtain enough error responses. Ludwig, Farrell, Ellis, and Gilchrist (2009) describe a method of circumventing this problem, using the LBA, that relies only on an estimate of the proportion of error responses rather than using error RT.

An alternative approach, useful for addressing under-constraint for any type of parameter, is to constrain parameters to be the same across conditions based on theoretical considerations. For example, the manipulation of response caution used in Forstmann et al.'s (2008) experiment is commonly assumed to not effect drift rate parameters. If this constraint is enforced by assuming the same value of error drift rate across a range of condition estimates of this parameter will be constrained as long as there are sufficiently many incorrect responses in total across all conditions.

A.3. Maximum likelihood estimation

In this section, we describe in detail how maximum likelihood estimation was carried out for the example data. First, consider data just from one emphasis condition in Forstmann et al. (2008). Let us assume that only drift rate differs between left and right responses. We could fit an LBA model for these data with five parameters to estimate: $(b, A, v_{\text{left}}, v_{\text{right}}, S, t_0)$. The v_{left} and v_{right} parameters represent mean drift rates for correct responses to left and right stimuli, respectively, and in the following we refer

to them generically as v_c . We fix the mean drift rates for error responses at $v_e = 1 - v_c$ for both left and right stimuli.

The likelihood function for the LBA model (see Brown & Heathcote, 2008) is relatively easy to compute as it is specified in terms of basic functions and the integral of a normal distribution, which has fast and accurate numerical approximations. Brown and Heathcote also provide computer code that evaluates the likelihood function. These routines take in a set of parameter values, and a response time, say t , and return the probability density that the accumulator corresponding to the first response has reached threshold before any other, and at time t . To get the likelihood for a *correct* response then the drift rate for the first accumulator is set at v_c and for the second accumulator at v_e (to get the likelihood for an *incorrect* response, one simply swaps the values of the drift rates to be v_e and v_c , respectively).

To construct the maximum likelihood objective function, we evaluate the likelihood function at each and every observed RT value, and multiply them together—this gives the likelihood of parameter set θ given the entire data set. We use every RT value because we want the set of parameters $\theta = (b, A, v_{\text{left}}, v_{\text{right}}, S, t_0)$ that is most likely given the four RT distributions under consideration: correct and error responses for left and right stimuli. For every RT value in each of these distributions, we take the following steps:

1. Identify the appropriate drift rate (v_{left} or v_{right}) depending on the stimulus presented on the given trial (left or right).
2. If the response associated with this RT was correct, set v_1 to the drift rate identified from Step #1. If the response was incorrect, set v_1 to one minus the drift rate from Step #1.
3. Set $v_2 = 1 - v_1$.
4. Subtract the parameter t_0 from the observed RT, as the likelihood equations given by Brown and Heathcote (2008) provide the likelihood for the *decision* time, which is $RT - t_0$.
5. Using the equation for the PDF (Equation 3 in Brown & Heathcote, 2008), and the drift rates from Steps #2 and #3, and the parameters from above, evaluate the likelihood function for this observation.

Once this operation is performed for every observation, the likelihood function can be obtained by simply multiplying together all the likelihood values (from Step #5) for all the data. However, it is usual to instead take the logarithm of all likelihoods, and then add these log-likelihoods together, to improve numerical stability. Optimizing the summed log-likelihoods is equivalent to optimizing the product of the likelihoods, as these two quantities are monotonically related.

A.4. A fast method for producing predicted quantiles

Predicted quantiles can be obtained more quickly than through the simulation method described in the main text by evaluation of the inverse of the CDF of a model. The CDF, $F(t|\theta)$, of the model gives the proportion of responses made before time t , given parameters θ . To find predicted quantile values, we require the inverse of the CDF—the proportion of responses made before time t . For choice RT models, however, we are interested in the conditional CDF for one of the possible responses, which does not reach a probability of 1 as t increases. For example, if a model predicted that only 65% of responses were accurate in a given

task, the conditional CDF for correct responses would only reach a probability of 0.65 at t grows large. To evaluate the predicted 0.1, 0.3, 0.5, 0.7 and 0.9 quantiles for the *correct* RT distribution, we must identify those values of t for which the conditional CDF for correct responses, $F(t|\theta)$ equals $0.1p$, $0.3p$, $0.5p$, $0.7p$ and $0.9p$, where p is the predicted response accuracy ($p = 0.65$ in our example). When calculating quantiles for the incorrect response distribution, the value p is replaced by $1 - p$, which is the probability of an incorrect response ($1 - p = 0.35$ for errors in our example).

For example, consider the 0.1 quantile of a correct RT distribution. The LBA predicts that 0.1 of this RT distribution is reached at $F^{-1}(0.1p)$, where p is the predicted proportion of correct responses. In other words the predicted 0.1 quantile value, say t , is the inverse conditional CDF for correct responses evaluated at $0.1p$. To get the predicted quantile value we first need to calculate the predicted value of p , which is done by evaluating the CDF at ∞ : $p = F(\infty)$. The inverse of the CDF does not have a closed-form expression that can be easily evaluated. Instead, we employ a numerical solution. We are attempting to solve $F^{-1}(0.1p) = t$, which is equivalent to $F(t) = 0.1p$. We are left, therefore, with the expression, $F(t) - 0.1p = 0$, which we can now solve using a standard root finding algorithm. The value of t returned by this algorithm, plus t_0 , is exactly the 0.1 quantile prediction for correct responses. We can repeat this process for error responses by replacing all instances of p with $(1 - p)$.

The following steps summarize the calculation of a given quantile corresponding to probability q for correct and incorrect responses:

1. Do steps #1-#3 specified above for maximum likelihood estimation to get the appropriate sets of parameters for correct (θ_c) and incorrect (θ_e) responses.
2. Obtain the predicted proportion of correct responses (p) by evaluating the CDF at infinity, given the parameters for correct responses that were chosen in Step #1.
3. Use a root finding algorithm to get the correct and incorrect quantiles. These correspond to the value of t for which $F(t|\theta_c) = qp$ and $F(t|\theta_e) = q(1 - p)$, respectively.

The QP plot graphs both data and model quantiles for correct and incorrect responses in each experimental condition, as explained in the main body of this paper. Note that code for calculating predicted quantiles using the LBA is available in Donkin et al. (2009).

References

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, *19*, 716–723.
- Averell, L., & Heathcote, A. (2009). The form of the forgetting curve and the fate of memories. *Journal of Mathematical Psychology* (in press).
- Brown, S., & Heathcote, A. (2005). A ballistic model of choice response time. *Psychological Review*, *112*, 117–128.
- Brown, S., & Heathcote, A. J. (2008). The simplest complete model of choice reaction time: linear ballistic accumulation. *Cognitive Psychology*, *57*, 153–178.
- Brown, S., Marley, A. A. J., Donkin, C., & Heathcote, A. (2008). An integrated model of choices and response times in absolute identification. *Psychological Review*, *115*, 396–425.
- Brown, S., Ratcliff, R., & Smith, P. (2006). Evaluating methods for approximating stochastic differential equations. *Journal of Mathematical Psychology*, *50*, 402–410.
- Busemeyer, J. R., & Townsend, J. T. (1992). Fundamental derivations from decision field theory. *Mathematical Social Sciences*, *23*, 255–282.
- Busemeyer, J. R., & Townsend, J. T. (1993). Decision field theory: a dynamic-cognitive approach to decision making. *Psychological Review*, *100*, 432–459.
- Donkin, C., Averell, L., Brown, S., & Heathcote, A. (2009). Getting more from accuracy and response time data: methods for fitting the linear ballistic accumulator. *Behavior Research Methods*, *41*, 1095–1110.
- Donkin, C., Brown, S., & Heathcote, A. (2009). The over-constraint of response time models: rethinking the scaling problem. *Psychonomic Bulletin & Review*, *16*, 1129–1135.
- Donkin, C., Brown, S., Heathcote, A., & Wagenmakers, E. J. (in press). Diffusion versus linear ballistic accumulation: different models for response time, same conclusions about psychological mechanisms? *Psychonomic Bulletin & Review*.
- Farrell, S., Ratcliff, R., Cherian, A., & Segraves, M. (2006). Modeling unidimensional categorization in monkeys. *Learning and Behavior*, *34*, 86–101.
- Forstmann, B. U., Dutilh, G., Brown, S., Neumann, J., von Cramon, D. Y., & Ridderinkhof, K. R. (2008). The striatum facilitates decision-making under time pressure. *Proceedings of the National Academy of Sciences*, *105*, 17538–17542.
- Heathcote, A., & Brown, S. D. (2004). Reply to speckman and roudser: a theoretical basis for QML. *Psychonomic Bulletin & Review*, *11*, 577–578.
- Heathcote, A., Brown, S. D., & Cousineau, D. (2004). QMPE: estimating lognormal, wald and weibull RT distributions with a parameter dependent lower bound. *Behavior Research Methods, Instruments, & Computers*, *36*, 277–290.
- Heathcote, A., Brown, S. D., & Mewhort, D. J. K. (2002). Quantile maximum likelihood estimation of response time distributions. *Psychonomic Bulletin & Review*, *9*, 394–401.
- Ho, T., Brown, S., & Serences, J. (2009). Domain general mechanisms of perceptual decision making in human cortex. *Journal of Neuroscience*, *29*, 8675–8687.
- Hoeting, J. A., Madigan, D., Raftery, A. E., & Volinsky, C. T. (1999). Bayesian model averaging: a tutorial. *Statistical Science*, *14*, 382–417.
- Lee, M. D. How cognitive modeling can benefit from hierarchical Bayesian models. *Journal of Mathematical Psychology* (in press).
- Ludwig, C. J., Farell, S., Ellis, L. A., & Gilchrist, I. D. (2009). The mechanism underlying inhibition of saccadic return. *Cognitive Psychology*, *59*, 180–202.
- Matzke, D., & Wagenmakers, E. J. (2009). Psychological interpretation of ex-Gaussian and shifted wald parameters: a diffusion model analysis. *Psychonomic Bulletin & Review*, *16*, 798–817.
- Morey, R. D. A Bayesian hierarchical model for the measurement of working memory capacity. *Journal of Mathematical Psychology* (in press).
- Myung, I. J. (2003). Tutorial on maximum likelihood estimation. *Journal of Mathematical Psychology*, *47*, 90–100.
- Myung, I. J., Forster, M. R., & Browne, M. W. (2000). Model selection [Special issue]. *Journal of Mathematical Psychology*, *44*(1–2).
- Nelder, J. A., & Mead, R. (1965). A simplex algorithm for function minimization. *Computer Journal*, *7*, 308–313.
- Pike, A. R. (1966). Stochastic models of choice behaviour: response probabilities and latencies of finite Markov chain systems. *British Journal of Mathematical and Statistical Psychology*, *21*, 161–182.
- Pinheiro, J. C., & Bates, D. M. (2000). *Mixed-effects models in s and s-plus*. New York: Springer.
- Pitt, M., & Myung, I. J. (2002). When a good fit can be bad. *Trends in Cognitive Sciences*, *6*, 421–425.
- Pratte, M., & Rouder, J. Hierarchical single- and dual-process models of recognition memory. *Journal of Mathematical Psychology* (in press).
- Ratcliff, R. (1978). A theory of memory retrieval. *Psychological Review*, *85*, 59–108.
- Ratcliff, R. (2002). A diffusion model account of response time and accuracy in a brightness discrimination task: fitting real data and failing to fit fake but plausible data. *Psychonomic Bulletin & Review*, *9*, 278–291.
- Ratcliff, R., Gomez, P., & McKoon, G. (2004). Diffusion model account of lexical decision. *Psychological Review*, *111*, 159–182.
- Ratcliff, R., & Rouder, J. N. (1998). Modeling response times for two-choice decisions. *Psychological Science*, *9*, 347–356.
- Ratcliff, R., & Smith, P. L. (2004). A comparison of sequential sampling models for two-choice reaction time. *Psychological Review*, *111*, 333–367.
- Ratcliff, R., Thapar, A., Gomez, P., & McKoon, G. (2004). A diffusion model analysis of the effects of aging in the lexical-decision task. *Psychology and Aging*, *19*, 278–289.
- Ratcliff, R., Thapar, A., & McKoon, G. (2007). Application of the diffusion model to two-choice tasks for adults 75–90 years old. *Psychology and Aging*, *22*, 56–66.
- Ratcliff, R., & Tuerlinckx, F. (2002). Estimating parameters of the diffusion model: approaches to dealing with contaminant reaction times and parameter variability. *Psychonomic Bulletin & Review*, *9*, 438–481.
- Schmiedek, F., Oberauer, K., Wilhelm, O., Süß, H.-M., & Wittmann, W. W. (2007). Individual differences in components of reaction time distributions and their relations to working memory and intelligence. *Journal of Experimental Psychology: General*, *136*, 414–429.
- Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, *6*, 461–464.
- Smith, P. L., & Ratcliff, R. (2009). An integrated theory of attention and decision making in visual signal detection. *Psychological Review*, *116*, 283–317.
- Smith, P. L., & Vickers, D. (1988). The accumulator model of two-choice discrimination. *Journal of Mathematical Psychology*, *32*, 135–168.
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P., & van der Linde, A. (2002). Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society. Series B*, *64*, 583–639.
- Tuerlinckx, F. (2004). The efficient computation of the cumulative distribution and probability density functions in the diffusion model. *Behavior Research Methods, Instruments, & Computers*, *36*, 702–716.

- Tuerlinckx, F., Maris, E., Ratcliff, R., & De Boeck, P. (2001). A comparison of four methods for simulating the diffusion process. *Behavior Research Methods, Instruments, & Computers*, 33, 443–456.
- Usher, M., & McClelland, J. L. (2001). On the time course of perceptual choice: the leaky competing accumulator model. *Psychological Review*, 108, 550–592.
- Vandekerckhove, J., & Tuerlinckx, F. (2007). Fitting the Ratcliff diffusion model to experimental data. *Psychonomic Bulletin & Review*, 14, 1011–1026.
- Vandekerckhove, J., & Tuerlinckx, F. (2008). Diffusion model analysis with MATLAB: a DMAT primer. *Behavior Research Methods*, 40, 61–72.
- Vandekerckhove, J., Tuerlinckx, F., & Lee, M. D. (in press). Hierarchical Bayesian diffusion models for two-choice response times. *Psychological Methods*.
- Van Zandt, T. (2000). How to fit a response time distribution. *Psychonomic Bulletin & Review*, 7, 424–465.
- Van Zandt, T., Colonius, H., & Proctor, R. W. (2000). A comparison of two response time models applied to perceptual matching. *Psychonomic Bulletin & Review*, 7, 208–256.
- Voss, A., Rothermund, K., & Voss, J. (2004). Interpreting the parameters of the diffusion model: an empirical validation. *Memory & Cognition*, 32, 1206–1220.
- Voss, A., & Voss, J. (2007). Fast-dm: a free program for efficient diffusion model analysis. *Behavior Research Methods*, 39, 767–775.
- Wagenmakers, E.-J., & Farrell, S. (2004). AIC model selection using Akaike weights. *Psychonomic Bulletin & Review*, 11, 192–196.
- Wagenmakers, E.-J., Lodewyckx, T., Kuriyal, H., & Grasman, R. (2010). Bayesian hypothesis testing for psychologists: a tutorial on the savage-dickey method. *Cognitive Psychology*, 60, 158–189.
- Wagenmakers, E.-J., Ratcliff, R., Gomez, P., & McKoon, G. (2008). A diffusion model account of criterion shifts in the lexical decision task. *Journal of Memory and Language*, 58, 140–159.
- Wagenmakers, E.-J., & Waldorp, L. (2006). Model selection: theoretical developments and applications [Special issue]. *Journal of Mathematical Psychology*, 50(2).
- White, C., Ratcliff, R., Vasey, M., & McKoon, G. (2009). Dysphoria and memory for emotional material: a diffusion model analysis. *Cognition and Emotion*, 23, 181–205.
- White, C., Ratcliff, R., Vasey, M., & McKoon, G. (2010). Using diffusion models to understand clinical disorders. *Journal of Mathematical Psychology*, 54, 39–52.