



Hidden Markov Models of Evidence Accumulation in Speeded Decision Tasks

Šimon Kucharský¹ · N.-Han Tran² · Karel Veldkamp¹ · Maartje Raijmakers^{1,3} · Ingmar Visser^{1,4}

Accepted: 12 August 2021
© The Author(s) 2021

Abstract

Speeded decision tasks are usually modeled within the evidence accumulation framework, enabling inferences on latent cognitive parameters, and capturing dependencies between the observed response times and accuracy. An example is the speed-accuracy trade-off, where people sacrifice speed for accuracy (or vice versa). Different views on this phenomenon lead to the idea that participants may not be able to control this trade-off on a continuum, but rather switch between distinct states (Dutilh et al., *Cognitive Science* 35(2):211–250, 2010). Hidden Markov models are used to account for switching between distinct states. However, combining evidence accumulation models with a hidden Markov structure is a challenging problem, as evidence accumulation models typically come with identification and computational issues that make them challenging on their own. Thus, an integration of hidden Markov models with evidence accumulation models has still remained elusive, even though such models would allow researchers to capture potential dependencies between response times and accuracy within the states, while concomitantly capturing different behavioral modes during cognitive processing. This article presents a model that uses an evidence accumulation model as part of a hidden Markov structure. This model is considered as a proof of principle that evidence accumulation models can be combined with Markov switching models. As such, the article considers a very simple case of a simplified Linear Ballistic Accumulation. An extensive simulation study was conducted to validate the model's implementation according to principles of robust Bayesian workflow. Example reanalysis of data from Dutilh et al. (*Cognitive Science* 35(2):211–250, 2010) demonstrates the application of the new model. The article concludes with limitations and future extensions or alternatives to the model and its application.

Keywords Evidence accumulation · Speeded decision · Speed-accuracy trade-off · Response times · Hidden Markov models · Phase transition

Introduction

Evidence accumulation models (EAMs) have become widely popular for explaining the generative process

of response times and response accuracy in elementary cognitive tasks (Evans & Wagenmakers, 2019). The strength of EAMs is their ability to accurately describe the speed-accuracy trade-off in speeded decision paradigms. The speed-accuracy trade-off is the conundrum that typically occurs when participants are instructed to make faster decisions, thereby increasing their proportion of errors (Bogacz et al., 2010; Wickelgren, 1977; Luce, 1991). The trade-off implies that in some situations, people can be slow and accurate, whereas fast and inaccurate in other situations. The dependency between response times and responses generally frustrates interpretation of response time and accuracy at face value. EAMs aim to capture and explain this dependency between response times and accuracy, and enable inference on the latent cognitive constructs and a mechanistic explanation of the observed response time and accuracy. Thus, such analyses often enable us to tell, for example, whether slowing down is caused by increased

✉ Šimon Kucharský
s.kucharsky@uva.nl

¹ Department of Psychology, Faculty of Social and Behavioural Sciences, University of Amsterdam, Amsterdam, The Netherlands

² Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany

³ Department of Educational Studies and Learn!, Faculty of Behavioral and Movement Sciences, Free University Amsterdam, Amsterdam, The Netherlands

⁴ Amsterdam Brain & Cognition (ABC), University of Amsterdam, Amsterdam, The Netherlands

response caution, increased difficulty or decreased ability of the respondent (van der Maas et al., 2011; Evans & Wagenmakers, 2019).

The traditional view of the speed-accuracy trade-off is that of a continuous function. That is, people are able to control their responses on the entire continuum from “slow and accurate” to “fast and inaccurate”. This is an intrinsic assumption of EAMs which makes it possible to manipulate parameters associated with “response caution” to make more or less accurate (on average) decisions by slower or faster (on average) responding. Under such a view, it is in principle possible to hold average accuracy to any value between a chance performance and a maximum possible accuracy (often near 100%), by adjusting how fast one needs to be.

An opposing view is that of a “discontinuity” hypothesis (Dutilh et al., 2010), which states that people are not able to trade accuracy for response time on a continuous function, but rather switch between different stable states. The discontinuity hypothesis in speeded decision-making is strongly associated with thinking about two particular response modes: a stimulus controlled mode and a guessing mode (Ollman, 1966). Under the stimulus controlled mode, one is maximizing response accuracy while sacrificing speed; whereas under the guessing mode, choices are made at random for the sake of responding relatively fast. Hence, there are two modes of behavior under discontinuity hypothesis. Such dual behavioral modes are present in many models of cognitive processing (e.g., dual processing theory Evans 2008).

The discontinuity hypothesis has an increasing relevance in the speeded decision paradigm because it is able to explain specific observed relationships between decision outcomes and reaction times that standard EAMs cannot account for (Dutilh et al., 2010; van Maanen et al., 2016; Molenaar et al., 2016). One of the most elaborate theoretical and empirical investigations of the “discontinuity” hypothesis is the phase transition model for the speed-accuracy trade-off (Dutilh et al., 2010), which added several more predictions regarding the dynamics of switching between the controlled and guessing state. These phenomena can be modeled using hidden Markov models (HMM, Visser et al., 2009; Visser, 2011). Dutilh et al. (2010) used HMMs to model their data such that response time and accuracy are independent conditional on the state. Specifically, the model assumed that the responses are generated from a categorical distribution and response times from the log-normal distribution, independently of each other. Thus, the speed-accuracy trade-off is described only by assuming one slow and accurate state, and one fast and inaccurate state. However, at least under the controlled state, evidence accumulation presumably takes place to generate the responses,

and so can lead to continuous speed-accuracy trade-off typical for EAMs, although within a smaller range than assumed under the continuous hypothesis. Thus, inference on the latent cognitive constructs given by the EAM might be the preferred option, but is neglected under the current HMM implementations of the phase transition model. Combining EAM with HMM would thus result in a model that is discontinuous on the larger scale (between state speed-accuracy trade-off), and continuous on the smaller scale (within state speed-accuracy trade-off), representing a third theoretical possibility beyond purely continuous and purely discontinuous models (Dutilh et al., 2010).

Fitting an HMM combined with an EAM would enable researchers to test specific predictions coming from the phase transition model as well as utilizing the strength of the EAM framework to account for the continuous speed-accuracy trade-off within the states. The ability of EAMs to infer the latent cognitive constructs liberates researchers from defining the states solely in terms of their behavioral outcomes. For instance, instead of describing the controlled state on the observed behavioral outcomes only (i.e., “slow and accurate”), EAMs allows researchers to form a mechanistic explanation of the observed behavioral outcomes using the latent cognitive constructs (i.e., “high response caution and high drift rate”). Further, capturing residual dependency between the observable variables conditionally on the latent states could improve performance of an HMM in terms of classification accuracy.

However, fitting EAMs can be a challenging endeavor, especially for more complicated models that allow for various sources of within and between trial variability, which often exhibit strong mimicry between different parameters, and as such belong to the category of “sloppy models” (Apgar et al., 2010; Gutenkunst et al., 2007). More complicated models, such as leaky competitor models, are not analytically tractable, and subject to highly specific simulation-based fitting methods (Evans, 2019). Thus, combining EAMs with HMMs, which themselves come with several computational (e.g., evaluation of the likelihood of the whole data sequence, Visser, 2011) and practical (e.g., label switching, Spezia, 2009) challenges, is highly demanding. The only successful applications of HMMs in these tasks is in combination with models that cannot capture possible residual dependencies, usually log-normal models or shifted Wald models for response times (Dutilh et al., 2010; Molenaar et al., 2016; Timmers, 2019). Yet, even the supposedly simplest complete model of response times and accuracy — the Linear Ballistic Accumulation model (LBA; Brown & Heathcote, 2008) — has proven to be difficult to combine with an HMM structure or even as a simple independent mixture

(Veldkamp, 2020); this may not come as a surprise considering the general identifiability issues of the standard LBA model (Evans, 2020).

Given the potential of complex cognitive models to suffer from computational issues, it is important to present evidence that the model implementation is correct and that the procedure used to fit the model on realistic data (in terms of plausible values but also size) indeed succeeds in recovering the information that is used for inferences. The importance of validating models in terms of practical applicability is ever more increasing with the growing heterogeneity of approaches for fitting complex models, as well as modern approaches to build custom models tailored to specific purposes.

This need is taken seriously in this article which implements and validates a simple (constrained) version of the LBA model as part of an HMM. This model makes it possible to capture the discontinuity of the speed-accuracy trade-off by the HMM part, while concomitantly striving to capture the residual dependency between speed and accuracy within the states. Further, the model retains the fundamental inferential advantages of an EAM framework, but is analytically tractable and stable enough to be used with standard, state-of-the-art, modeling tools. To our knowledge, this is the first working combination of an HMM and an EAM, and serves as a proof of concept.

The structure of this article is as follows. First, the model is described in conceptual terms to explain the core assumptions and mechanics. Second, a simulation study summarizes all steps that were followed when building and validating the model in accordance with a robust Bayesian workflow (Schad et al., 2019; Talts et al., 2018; Lee et al. 2019). The model validation is followed with an empirical example to demonstrate the full inferential power of the model on experimental data. The article concludes with discussion and future potential directions towards improving the model.

Model

The general architecture of the model for response times and choices that we adopt here is the same as for the Linear Ballistic Accumulator (LBA; Brown & Heathcote, 2008). In the standard LBA, each response option is associated with its own evidence accumulator. Each accumulator rises linearly towards a threshold from a randomly drawn starting point, with its own specific drift rate, drawn from some distribution (commonly a normal distribution that is truncated at zero). The first accumulator that reaches its decision threshold triggers the corresponding response. Figure 1 explains the basic mechanics of typical LBA model.

Although the LBA became a popular choice for analyzing response times and accuracy, more recently evidence has surfaced suggesting practical identifiability issues of the standard LBA model — especially when trying to quantify differences in parameters such as decision boundary or drift rates between experimental conditions (Evans, 2020). Given that HMMs can be viewed as way to quantify differences between “conditions” (states) which themselves need to be inferred from the data, (lack of) identifiability of the standard LBA in combination with HMMs is a concern (especially in the upper bound of the starting point Veldkamp, 2020; Timmers, 2019).

However, there exists a number of potential remedies to solve the identifiability issue of the standard LBA. These remedies involve constraining the LBA model in some way while retaining as much flexibility of the model as possible to account for different patterns in the data, and to still allow inferences on the most fundamental parts of the evidence accumulation decision process (e.g., speed of accumulation, response caution). For example, a relatively well established set of constraints is to ensure that the average drift rates across accumulators are equal to some constant value (e.g., a scaling value of 1; Evans, 2020; Donkin et al., 2011; Visser & Poessé, 2017). Such constraints may be accompanied by implementing equality constraints on parameters such as the upper bound of the starting point or the standard deviation

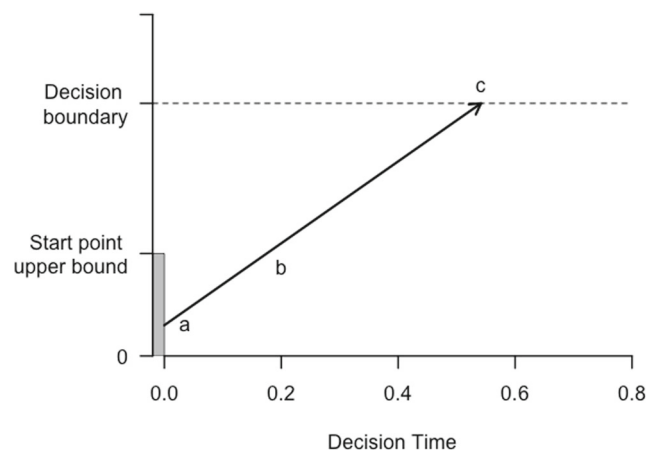


Fig. 1 Linear Ballistic Accumulator (Brown & Heathcote, 2008). Each response outcome has an independent accumulator. For simplicity, the plot shows only one accumulator. (a) Starting point for each accumulator is generated from Uniform distribution between zero and the upper bound of the starting point. (b) Accumulator is launched from the starting point and with a drift rate that is generated from normal distribution with a mean drift and standard deviation of drift rate. (c) Decision is made based on which accumulator hits the decision boundary first. Final response time is the sum of the decision time (the time it took the first accumulator reach the boundary) and a non-decision time (a fixed time for encoding the stimuli and motoric response)

of the drift rates. In the context of different conditions, even more stringent (equality) constraints are possible, such as equating parameters (such as drift rate for the “error” response) across conditions (Evans, 2020).

This article aims to provide a proof of concept that EAMs and HMMs can be combined into a single model. The present application simplifies the LBA model to a bare minimum and acts as a sanity check — in case even very minimalist EAM models cannot be employed as part of a HMM model, there is little reason to expect that more complex, complete and computationally demanding models of decision-making will be more successful.

The bare minimum, simple instance of LBA is achieved in this article by setting several constraints on the parameters. For practical reasons, we will refer to this model as sLBA, a short for “simplified Linear Ballistic Accumulator”. Most significantly, the model implemented in this article fixes all starting points at zero, effectively removing the variability of the starting point. As commonly done in the LBA, we constrain the drift rates to sum to unity. In addition to that, the drift rates are assumed to have equal standard deviations across accumulators. Full details on the model, its likelihood and identifiability are described in Appendix A, additional helpful derivations can be found in Nakahara et al. (2006). Figure 2 explains the model in additional detail.

The simplification achieved by removing the variability of the starting point makes the model coarsely similar to the LATER model (Linear Approach to Threshold with Ergodic Rate, Carpenter, 1981; Noorani and Carpenter, 2016), with the difference that the current model explicitly evaluates the likelihood of observing the first accumulator that reached the threshold according to the general race equations (see Heathcote & Love, 2012), and contains additional parameters (such as non-decision time). Therefore, it enables researchers to model accuracy in addition to response times as opposed to the LATER model (see Ratcliff, 2001 for critique of LATER for inability to do so).

The constraints employed in this application greatly reduce the complexity compared to the standard LBA model. Specifically, our model for responses and response times on a two-choice task contains the following parameters: the average drift rate for the correct (v_1) and incorrect (v_2) responses, the standard deviation of the drift rates (σ), the decision threshold (α), and the non-decision time (τ). The latter three parameters are equal for both accumulators.

The purpose of simplifying the LBA model is to employ it as a distribution of response times and responses in an HMM. Specifically, the current model assumes two latent states: A “controlled” state ($s = 1$) and a “guessing” state ($s = 2$). These states evolve according to a Markov chain, which is characterized by the initial (π_1 and π_2) and transition state probabilities ρ_{ij} , where the first index i

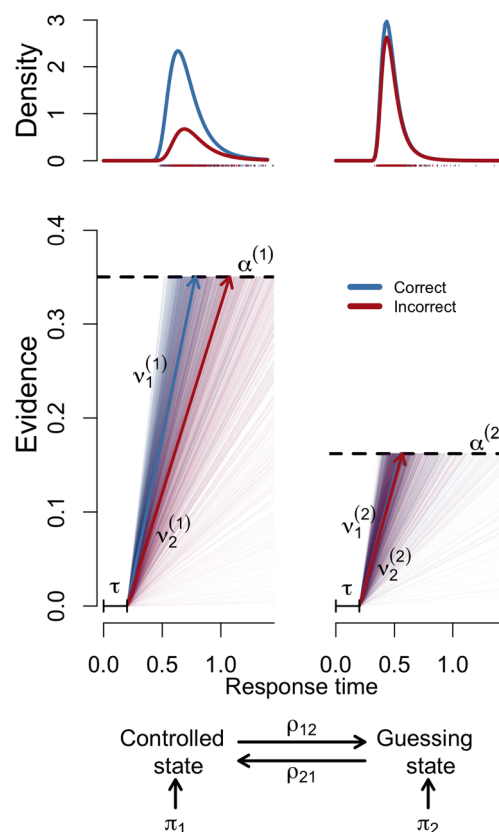


Fig. 2 HMM combined with sLBA. Bottom panel: Latent controlled and guessing states evolve as a Markov chain, with initial state probabilities π_1 and π_2 , and transition probabilities ρ_{12} and ρ_{21} . Middle panel: Non-decision time τ shifts the response times. Correct and incorrect responses launch an accumulator (starting at 0), with a drift rate drawn from a truncated normal distribution with mean drift rate v and a standard deviation σ . The plot shows the average drift rates as thick arrows, and realizations of the random process as thin lines to represent the randomness of the process. Accumulator that reaches the decision boundary α first launches corresponding response. Average drift rates and decision boundary can differ between the states. Top panel: Under the controlled state (left), the expected response times are larger than under the guessing state (right), but the accuracy is higher (i.e., the decision boundary is reached by the correct accumulator more often)

corresponds to the outgoing state and j corresponds to the incoming state: For example, ρ_{12} is the probability that the participants switch from the controlled state to the guessing state.

Traditionally, these states would be equipped by their own distribution of response times and responses, possessing their own parameters. That is, we could use the LBA model for each latent state of the HMM, and estimate the drift rate for the correct responses for the first state $v_1^{(1)}$, second state $v_1^{(2)}$, and similarly for all of the parameters. However, we further reduce the complexity of the model by equating some parameters between states. Specifically, we assume that the difference between the guessing state and the controlled state is evoked by differences between

average drift rates and decision thresholds. The rest of the parameters are held equal across the states. Thus, equality constraints $\sigma^{(1)} = \sigma^{(2)}$ and $\tau^{(1)} = \tau^{(2)}$ are used to further simplify the model.

Additionally, there are some notable considerations regarding the controlled and guessing states, which will later help setting priors and preventing label switching. Specifically, the controlled state has higher average drift rate for the correct response than the guessing state ($v_1^{(1)} > v_1^{(2)}$), and consequently $v_2^{(1)} < v_2^{(2)}$ due to the sum-to-one constraint of the drift rates, see Appendix A) at the expense of having higher decision threshold ($\alpha^{(1)} > \alpha^{(2)}$). Further, if the second state truly is guessing, the drift rates under this state should be roughly the same: $v_1^{(2)} \approx v_2^{(2)} \approx 0.5$.

Implementation

We implemented the HMM and LBA model in a probabilistic modeling language Stan (Carpenter et al. 2017); specifically, v2.24.0 release candidate of CmdStan (<https://github.com/stan-dev/cmdstan/releases/tag/v2.24.0-rc1>, Stan Development Team, 2020). In this version of Stan, several new functions were introduced that implement the forward algorithm for calculating the log-likelihood of the data sequence, while marginalizing out the latent state parameters (for easy introduction, see Visser, 2011), which makes estimating HMM models in Stan much easier, computationally cheaper, and less error-prone than before (which required manual coding of the forward algorithm). The sLBA distribution of response times and responses was custom coded in the Stan language. We executed CmdStan from the statistical computing language R (R Core Team, 2020) using the R package `cmdstanr` (Gabry & Češnovar, 2020). The code is available at https://github.com/Kucharssim/hmm_slba.

Label Switching

Finite mixture models and Hidden Markov models share the characteristic that the likelihood of the models is typically invariant to the permutation of the latent state labels (Jasra et al., 2005; Spezia, 2009). This means that fitting the model can result in different estimates, depending on towards which state configurations the fitting procedure leads to. In the current context of guessing and controlled state, it is not possible on the basis of the model likelihood alone to state whether component 1 should be controlled or guessing state and vice versa — both options lead to the same likelihood value. There are several perspectives on dealing with potential label switching, perspectives that differ in terms of what types of applications and inferential paradigms one follows. For example, in maximum likelihood paradigm, label switching is not a severe problem as the analyst can

simply relabel the states after the model has been fitted, based on how the parameter estimates can be interpreted. In Bayesian framework (especially with MCMC), the problem is more complicated as the label switching can manifest in different ways, and can also depend on the sampler (and its settings) one uses to obtain the estimates of the entire posterior distribution. Common remedies of label switching are, for example, (1) change the model so that emission distributions under each state are uniquely identified, (2) establish parameter inequalities which leads to identifying the labels, (3) use of informative priors that lead to better identification of the a priori constraints, (4) some form of state relabeling of the posterior samples, among others. Usually, various remedies are combined together as the solutions do not work in generality for all possible mixture problems and applications.

In the current application, we heavily rely on approach 3), whereby specifying informative priors leads to soft identification of the state labels, i.e., associating slow and accurate responding with a (controlled) state 1 and fast and inaccurate responding with a (guessing) state 2. However, it is important that even which informative priors, one is only increasing the a priori probability of some state configuration, but does not render other configurations impossible. In fact, the other state configurations are still *valid* modes of the joint posterior space, albeit less plausible according to the prior specification. In some applications (estimating marginal likelihood in order to conduct model comparison; Frühwirth-Schnatter, 2004), it is actually desirable to make sure that the sampler is switching between state labeling freely, to ensure that the MCMC sampling efficiently explores the joint posterior in its entirety. In purely estimation settings (which is the case of this article which is not concerned by model comparison), one does not need to ensure that all valid modes of the posteriors are explored efficiently, as long as the main mode is explored well, which, among others, entails checking whether the labels did *not* switch, either within- or between- the MCMC chains.

Simulation Study

In order to investigate the quality of inferences we draw from the model, a simulation study was conducted. Specifically, we conducted the simulation in accordance with a principled Bayesian workflow (Schad et al., 2019). The simulation study consists of (1) prior predictive checks to identify priors that reflect our domain specific knowledge, (2) a computational faithfulness check to test correct posterior distribution approximation, and (3) model sensitivity analysis to investigate how well the estimated posterior mean of parameter matches the true

data generating value, and the amount of updating (i.e., how much are the parameters informed by the data). Additionally, as is the case in classical model validation simulation, we report standard parameter recovery results, including coverage probabilities of credible intervals.

Prior Predictives

Choosing prior distributions is an integral part of the Bayesian model-building process because the prior should reflect theoretical assumptions and cumulative knowledge about the parameter space as well as aid model convergence (Vanpaemel, 2011; Gershman, 2016). Ideally, the priors should be informed and constrained by a large collection of previous studies (Tran et al., 2020; van Zwet & Gelman, 2021) to yield more efficient sampling and plausible estimates. In the current study, we selected prior distributions to constrain parameter values to reasonable regions of the parameter space (e.g., non-decision time must be positive, therefore we used an exponential distribution) and to nudge the model towards convergence. Concomitantly, our prior distributions were informed by the large collection of literature on evidence accumulation models applied to lexical and perceptual decision tasks (Tran et al., 2020). Interested readers who want to apply our models to different experimental tasks or non-standard populations might want to consult the corpus of literature specific to the application to adjust the prior distributions.

To place priors that reflect our expectations about data from the tasks to which the model will be applied, we conducted prior predictive simulations. In particular, we first set out to generate 1,000 data sets each of 200 trials, which is generally a lower bar for running speeded decision tasks. Then, the following expectations of the generated data are defined, specified in terms of summary statistics across the 200 observations per data set. Throughout, response times are measured and reported in seconds. In case response times are measured in different units, the priors should be re-scaled appropriately.

Latent State Distribution

First, we expect that the number of trials participants spend in one or another state will be relatively even, and that it is very rare that participants would complete all 200 trials in a single state. The evenness is achieved by composing a symmetric initial state probabilities vector π and a symmetric transition matrix $P = \begin{bmatrix} \rho_1 & \\ & \rho_2 \end{bmatrix}$. Further, we assume that the states are relatively sticky, therefore there will be a tendency to stay in the current state rather than switching to another state. Specifically, the average run length is expected to be approximately between 5 and 10,

and that in at least 50% of the simulations the proportion of the trials under the controlled state ranges between 30 and 70%.

We chose the following priors

$$\pi \sim \text{Dirichlet}(5, 5)$$

$$\rho_1 \sim \text{Dirichet}(8, 2)$$

$$\rho_2 \sim \text{Dirichet}(2, 8).$$

The initial state probabilities are assigned a symmetric Dirichlet prior. The hyperparameters slightly favor probabilities closer to 0.5. Usually, the initial state probabilities are not the focus of inference as they depend mostly on just the first trial. Thus, slightly informative priors were chosen to help the model to converge. For the transition probabilities, Dirichlet priors that favor “sticky” states were chosen. Specifically, the mean probability of staying under the current state is 0.8. There is still considerable uncertainty about how sticky the two states are: 90% of the prior mass for the probability of persisting in the current state lies between 0.63 and 0.94.

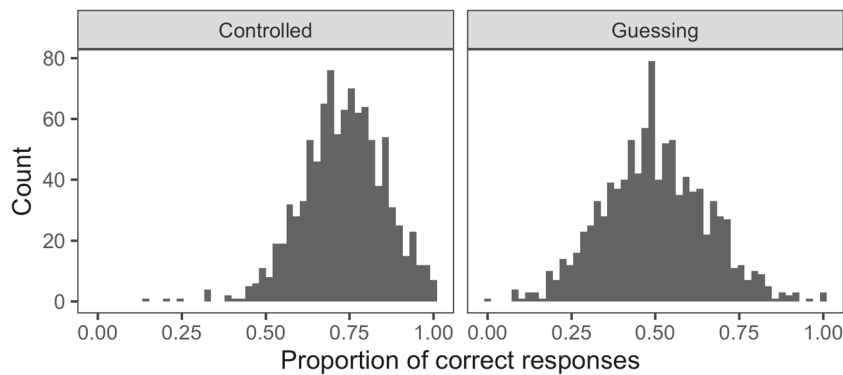
The results of the prior predictive simulation showed that the median of the average run length is 6.25 (IQR [4.35, 9.524]). The distribution of the average run length is positively skewed. Although it could be expected in many experiments that run lengths could be higher, the priors would have to be much more informative (pushing the probability of staying in a current state closer to one) than the current settings. However, that would give only a very narrow range of the values used for validating the models. Therefore, the current setting of the prior is a compromise between prior expectations about the data and the need to validate the model on a wider range of parameter values. Regarding the percentage of trials in the controlled state, the distribution over the 1,000 simulations had a median of 0.51 (IQR [0.35, 0.67]).

Response and Response Time Distributions

We expect that the distributions of the responses will be the following. Under the controlled state, the proportion of correct responses is well above chance; we assume that under the controlled state, there is almost zero probability that a person would have accuracy smaller than 50%, and that it is possible to achieve relatively high accuracy on average ($\approx 75\%$). Under the guessing state, we assume that the average accuracy is exactly 50%.

For the distributions of the response times, we have the following expectations. First, the response times under the controlled state are on average slower than responses under the guessing state. Second, the responses under the guessing state are relatively rapid: responses in simple perceptual decision tasks can be faster than 1 s on average. Third, the

Fig. 3 Prior predictive distribution of the response accuracy (proportion of correct answers)



majority of response times does not exceed 5 s (Tran et al., 2020).

Based on these considerations and prior predictive simulations, the following prior specification for the LBA parameters were identified as suitable:

$$\begin{aligned} \nu^{(1)} &\sim \text{Dirichlet}(14, 6) \\ \nu^{(2)} &\sim \text{Dirichlet}(10, 10) \\ \alpha^{(1)} &\sim \text{Gaussian}(0.5, 0.1)_{(0, \infty)} \\ \alpha^{(2)} &\sim \text{Gaussian}(0.25, 0.05)_{(0, \infty)} \\ \sigma &\sim \text{Gaussian}(0.4, 0.1)_{(0, \infty)} \\ \tau &\sim \text{Exponential}(5) \end{aligned}$$

Figure 3 and Table 1 summarize the prior predictive distribution of the accuracy (proportion of correct answers) under the two states separately. As desired, the accuracy under the controlled state is well above chance, whereas under the guessing state it clusters around 50%. There is considerable variability under both states, leaving the possibility for the model to learn from the data.

Figure 4 and Table 2 summarize the prior predictive distributions of the average response times for correct and incorrect responses under the two states separately. As desired, the average response times are slower under the controlled state than under the guessing state. The majority of the average response times under the guessing state are below 1 s, whereas under the controlled state cluster around 1 s. There are no large differences between response times for correct and incorrect responses under the two states separately, although the average response times for incorrect responses under the controlled state show higher variance

than for the correct responses. However, this phenomenon might be caused by the fact that there are more correct responses than incorrect responses under the guessing state, resulting in higher standard errors for the averages of the incorrect responses.

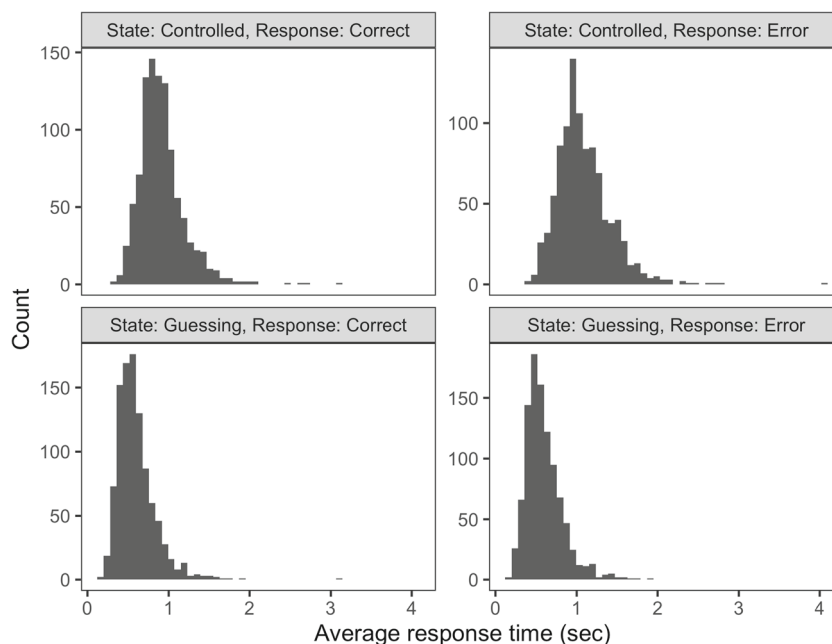
The prior distributions specified above may seem extremely informative, introducing “subjective” bias to the analysis. However, we believe the prior distributions are justified by our prior predictive simulations and based on cumulative characterizations of psychological processes underlying a lexical decision and a perceptual decision task of EAMs (Tran et al., 2020). Further, prior distributions may be also regarded as constraining the parameter space to plausible values (Tran et al., 2020; Vanpaemel, 2011; Kennedy et al., 2019), similarly as a traditional statistician would decide on ranges of parameters for a simulation study. In the current study, the prior distributions actually cover slightly more volume of the parameter space than is typical in simulation studies of similar type (e.g., Donkin et al., 2011; Visser & Poessé, 2017). Lastly, priors on the parameters in both states (e.g., $\alpha^{(1)}$ and $\alpha^{(2)}$) are used to primarily separate the latent states from each other, and associate the first state with the controlled state (and conversely the second state with the guessing state). Using informed priors in such occasions prevents label switching problems, and gently nudges the model towards convergence.¹ However, the prior specification does not ensure that the labels do not switch at all. When fitting the models, we performed additional checks using the posterior samples to check whether the labels indeed converged to the modes of the posteriors we intended.

Table 1 Descriptives of the prior predictive distribution of the response accuracy (proportion of correct answers)

State	Mean	SD	Quantile				
			2.5%	25%	50%	75%	97.5%
Controlled	0.73	0.12	0.48	0.65	0.73	0.81	0.96
Guessing	0.50	0.16	0.21	0.39	0.50	0.60	0.81

¹There are other techniques to identify states and prevent label switching (Jasra et al., 2005). For example, a common approach is to put an order constraint on the model parameters, for example, $\alpha^{(1)} < \alpha^{(2)}$, by using a transformation $\alpha_2 := \alpha_1 + \exp(\theta)$. Such a “hard” order restriction is effective in dealing with label switching, but makes it harder to reason about the prior specification. Further, “hard” order restrictions can hinder computing normalizing constants, in case one is eager to quantify the marginal likelihood (evidence) of the model (Frühwirth-Schnatter, 2004; 2019).

Fig. 4 Prior predictive distribution of the average response times



Computational Faithfulness

There are many ways in which model implementation can fail, especially in case of Bayesian models requiring MCMC. Possible problems might arise due to error in specification of the likelihood (or just insufficiently robust implementation), the use of difficult parameterizations, or a simple coding error. Another problem may arise when the model combined with the priors and the data result in a very complex parameter space for the MCMC algorithm to navigate, which may lead to inefficient exploration of the target posterior distribution. Such issues can lead to biased estimates, underestimating the uncertainty of parameters, or simply wrong inferences.

For the endless possibilities in which model implementation can fail, there was a lot of recent advancement in techniques that aim to check for *computational faithfulness* of a model — in the context of the Bayesian framework, this means testing whether the proposed MCMC procedure yields valid approximations of the posterior distributions (Schad et al., 2019). One established technique is

Simulation-based calibration (SBC, Talts et al., 2018). As the model that we propose in this article is definitely suspect for computational problems, we use SBC to check our model implementation (although it could be argued that such checks should be done by default for non-standard models at least). Since these checks are not yet the standard in cognitive modeling literature (Schad et al., 2019), we briefly summarize the rationale behind SBC here, although the interested reader should refer to excellent articles by Talts et al. (2018) and Schad et al. (2019).

To check whether the method used for approximating the posterior distribution $\pi(\theta|\tilde{y})$ is correct, the following steps can be done: (1) draw from the prior distribution $\tilde{\theta} \sim \pi(\theta)$, (2) draw a data set from the model using the generated values of the parameters, $\tilde{y} \sim \pi(\tilde{y}|\tilde{\theta})$, and (3) fit the model on the generated data to obtain the posterior distribution $\pi(\theta|\tilde{y})$. The draws from such an obtained distribution, across many repeated replications of this procedure, should give back the prior distribution of the parameters $\pi(\theta)$. In short, SBC builds on the fact that (Talts et al., 2018)

$$\pi(\theta) = \int \int \pi(\theta|\tilde{y})\pi(\tilde{y}|\tilde{\theta})\pi(\tilde{\theta})d\tilde{y}d\tilde{\theta}, \tag{1}$$

which means that we can recover analytically the prior distribution on model parameters $\pi(\theta)$ by averaging the posterior distribution $\pi(\theta|\tilde{y})$ weighted by the prior predictive distribution $\int \pi(\tilde{y}|\tilde{\theta})\pi(\tilde{\theta})d\tilde{\theta}$. In order to check whether the prior distribution is indeed recovered, for each repetition, we compare the draw from the prior (that generated the data) to the samples from the posterior, and count the posterior samples that are smaller than the draw from the prior. If these two distributions are the same, every

Table 2 Descriptives of the prior predictive distribution of the average response times

State	Response	Mean	SD	Quantile				
				2.5%	25%	50%	75%	97.5%
Controlled	Correct	0.92	0.28	0.49	0.73	0.87	1.03	1.57
Controlled	Error	1.09	0.34	0.59	0.87	1.03	1.26	1.82
Guessing	Correct	0.60	0.24	0.28	0.44	0.55	0.70	1.19
Guessing	Error	0.60	0.23	0.27	0.44	0.55	0.70	1.18

rank (i.e., the count of posterior samples that are smaller than the generating parameter value) would be equally likely — yielding an approximately uniformly distributed rank statistic (Talts et al., 2018).

Using the already created ensemble of 1,000 prior predictive data sets in “Prior Predictives”, each of the data sets was fitted using Hamiltonian Monte Carlo supplied by Stan (Carpenter et al. 2017). Due to computational constraints (typical run of a model averages roughly about 500 sampling iterations per minute on Apple’s MacBook Air edition 2017), each model ran only with one chain for 500 warmup and 1,000 sampling iterations. Starting points were generated by drawing independent samples from the priors. In case the model label switched, the model was reran (at maximum five times). Model switching was detected by comparing the true (generative) states to the estimated states (identified using modal assignment based on mean state probabilities using the forward-backward algorithm). This resulted in non-label switching MCMC samples for 945 data sets out of the total 1,000. Since

only 783 repetitions achieved acceptable values of the (split-half) Gelman-Rubin \hat{R} statistic (Gelman & Rubin, 1992) between 0.99 and 1.01 for all of the parameters, we selected several data sets at random from non-converged cases and refitted them with 4 chains, 1,000 warmup and 1,000 sampling iterations. The new model fits had good \hat{R} for all parameters, suggesting that the unsatisfactory convergence diagnostics were a consequence of the small number of MCMC iterations during the simulation. We excluded from the results only the repetitions that label switched, but kept those that did not yield satisfactory convergence diagnostics. Because the SBC rank statistic is sensitive to potential autocorrelation of the chain, the posterior samples were thinned by a factor of 50 — leading to the rank statistic ranging between 0 and 20.

Figure 5 shows the histogram of the SBC rank statistic for each of the parameter separately. Figure 6 shows the difference between the cumulative distribution and the theoretical cumulative distribution of a uniformly distributed variable (Talts et al., 2018).

Fig. 5 Simulation-based calibration: Histogram of the rank statistic. The dashed lines correspond to the lower and upper limits of the 95% interval under the null hypothesis that the rank statistic is uniformly distributed

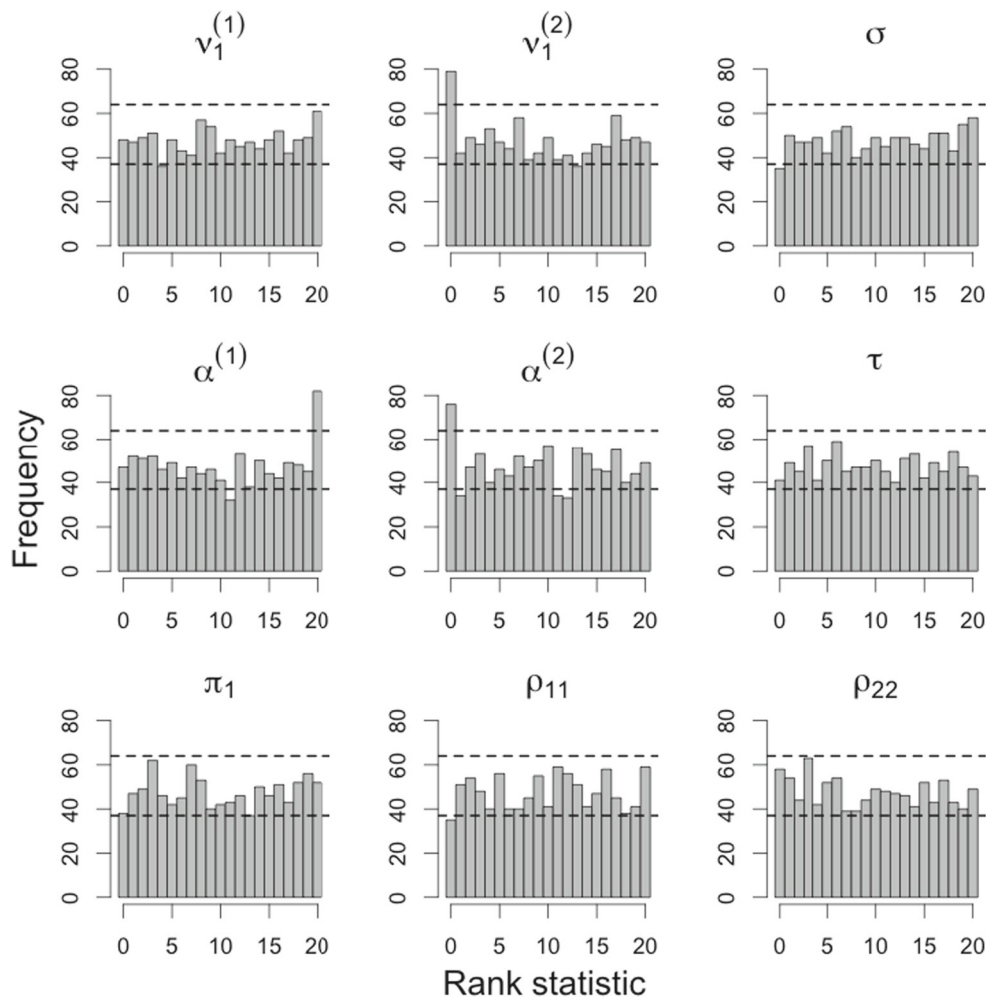
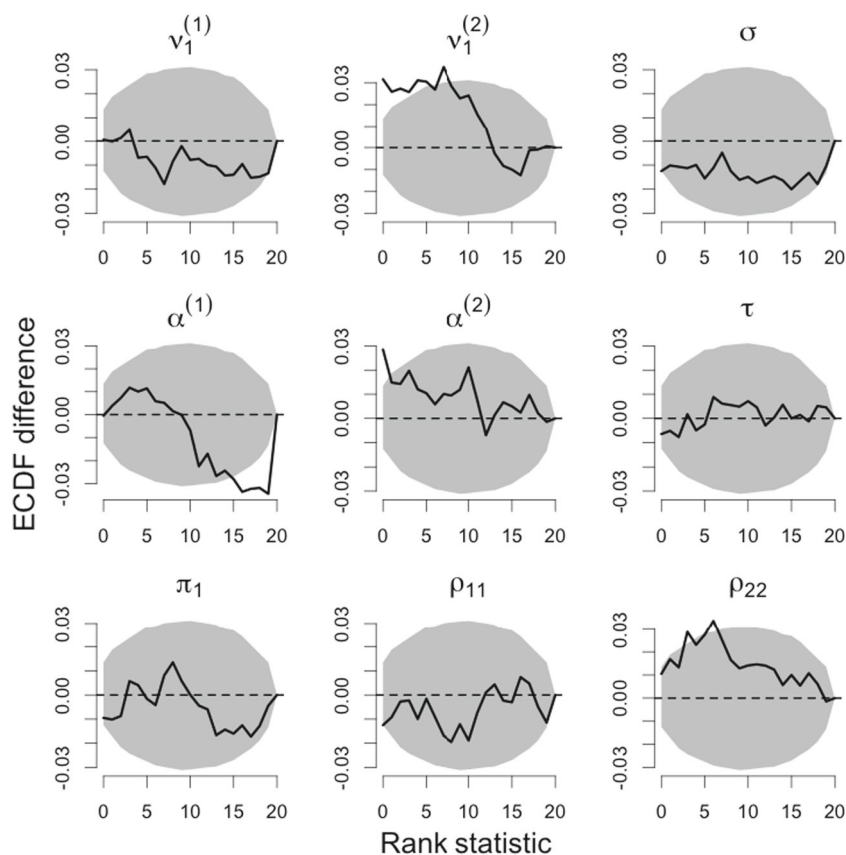


Fig. 6 Simulation-based calibration: ECDF of the rank statistic minus the ECDF of a uniformly distributed variable. The shaded area corresponds to the 95% interval under the null hypothesis that the rank statistic is uniformly distributed



The results show that none of the parameters exhibits typical patterns present in case that the posterior approximation is under-dispersed or over-dispersed compared to the true posterior (which would manifest as a U or \cap shape of the rank distribution; Talts et al., 2018). Further, the distribution of rank statistics for most of the parameters seem consistent with a uniform distribution, suggesting that the posterior approximation is very close to the true posterior. However, three parameters seem potentially problematic: the rank statistic for $\alpha^{(1)}$, $\alpha^{(2)}$, and $v_1^{(2)}$ show an excess of frequencies at 20 and 0, respectively, suggesting that $\alpha^{(1)}$ approximation could be underestimating the true posterior, whereas $\alpha^{(2)}$ and $v_1^{(2)}$ approximations could be overestimating the true posterior. However, this observation could also arise if the thinning was not efficient to reduce the auto-correlation of the chain (autocorrelation can result in excess of ranks at the edge of the distribution; Talts et al., 2018). Additionally Fig. 6 reveals that the rank distribution for ρ_{22} also potentially deviates from the uniform distribution. However, this deviance is not associated with any typical problem in posterior approximations, lacking a meaningful interpretation.

SBC gave us assurance that our model is capable of approximating the posterior distribution for most of the parameters. Three potentially problematic parameters remain, although the deviance from the expected results it

small. Potential explanations for these deviances could be the constraints to resolve label switching (which could cause the truncation of the parameters for one state near values for the same parameter from the other state), or unsuccessful reduction of the auto correlations of the MCMC chains (which could be solved by running the procedure for more iterations and use higher thinning.)

Model Sensitivity

Next, the goal was to investigate for each parameter, (1) how well the posterior mean matches the true data generating value of the parameter, and (2) how much uncertainty is removed when updating the prior to the posterior. This is useful to investigate the bias-variance trade-off for each parameter, and to adjust our expectations regarding how much we can learn about parameters, given a data set of a specified size (in this simulation, number of trials = 200).

To answer (1), posterior z-scores for each parameter are defined as:

$$z = \frac{\mu_{\text{posterior}} - \tilde{\theta}}{\sigma_{\text{posterior}}}, \tag{2}$$

that is, the difference between the posterior mean and the true parameter value is divided by the posterior standard deviation. The posterior z-scores tell us how far the posterior

expectation is from the true value, relative to the posterior uncertainty. The distribution of the posterior z -scores should have a mean close to 0 (if not, the posterior expectation is a biased estimator).

To answer (2), posterior contraction for each parameter is defined as:

$$\text{contraction} = 1 - \frac{\sigma_{\text{posterior}}^2}{\sigma_{\text{prior}}^2}. \tag{3}$$

If the posterior contraction approaches one, the variance of the posterior is negligible compared to the variance of the prior, indicating that the model learned a lot about the parameter of interest. Conversely, if the posterior contraction is close to zero, there is not much information in the data about the parameter, resulting in the inability to reduce the prior uncertainty.

These two variables are plotted against each other in a scatter plot, which provides useful diagnostic insights (Schad et al., 2019). Specifically, for each parameter, and each simulation which did not label switch, the posterior z -scores and posterior contraction are plotted on the y -axis and x -axis, respectively. Figure 7 shows the diagnostic plot for the nine parameters with equal axes between them to enable comparison between parameters.

All of the parameters cluster around z -scores of 0 (dashed horizontal line), suggesting that neither of the

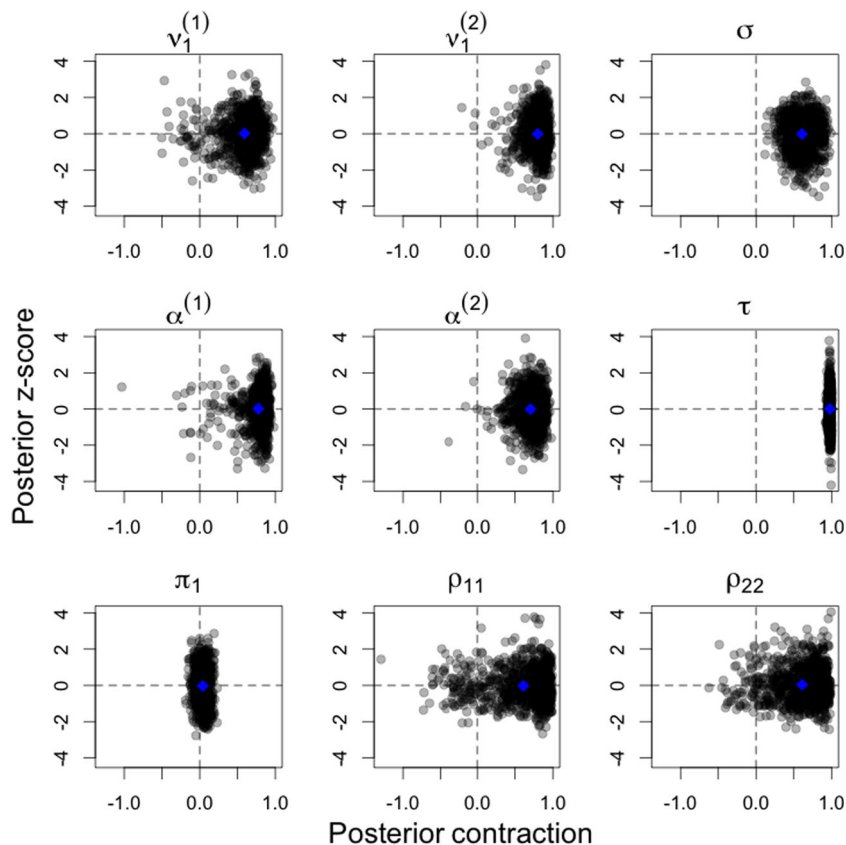
parameters exhibits systematic bias. However, there are large differences between parameters in terms of posterior contraction. The most contraction is present for the non-decision time τ , followed by the rest of the LBA parameters. We could expect that the contraction would increase with the number of trials. The worst results concern the initial state probability π_1 : The posterior contraction basically stays at zero. However, this is expected as the initial state probability is affected mostly by just the first trial, and as such, there is not much information in the data about it. Increasing the number of trials would not help to identify this parameter, only repeated experiments would.

In general, the sensitivity analyses suggest that the amount of learning about the parameters of interest could be satisfactory given the typical experimental designs (our simulation was based on 200 trials per experiment, whereas typical decision tasks experiments could count multiples of that number), especially for the LBA parameters.

Parameter Recovery and Coverage Probability

Traditional simulation studies aim to validate statistical models and assess the quality of a point estimator of a given parameter of interest. Additionally, such simulations are accompanied by assessment procedures. This section adheres to this tradition: for each of the parameters (that are

Fig. 7 Model sensitivity plot for all nine parameters. Blue diamond shapes depict the means of the distributions



not a linear combination of others) we report the standard “parameter recovery” results.

The simulation was done using two estimation techniques: the maximum a posteriori (MAP) estimation, and the posterior expectation (i.e., the mean of the posterior distribution). MAP is useful in situations where researcher needs to obtain estimates quickly, and does not need to express the uncertainty in the estimates. As the rest of the article focuses on full Bayesian inference, MAP results are presented only in the [Supplementary Information](#). Pearson’s correlation coefficient between the estimated parameter value and its true values serves as a rough indicator of parameter recovery. High correlations indicate that the model is able to pick up variation in the parameter. Additionally, scatter plots visualizing the relationship between the true and estimated parameter values show the precise relationship between the true and estimated values of the parameters.

We also investigate the coverage performance of the central credible intervals. For each parameter, the frequency with which 50% and 80% central credible intervals contain the true data generating value was recorded. The confidence levels are relatively low compared to traditionally reported values, because we have only 1,000 MCMC samples per parameter due to computational constraints, which results in

low precision in the tails of the posterior distributions (i.e., the tail effective sample size was generally too low).

Posterior Expectation

Figure 8 shows the scatter plot between the true (x-axis) and estimated (y-axis) values (i.e., means of the posteriors) for the nine free parameters in the model: the drift for the correct choice under the controlled state ($v_1^{(1)}$), the drift for the correct choice under the guessing state ($v_1^{(2)}$), the standard deviation of drifts (σ), the decision boundary under the controlled ($\alpha^{(1)}$) and guessing ($\alpha^{(2)}$) state, the non-decision time (τ), the initial probability of the controlled state (π_1), the probability of dwelling in the controlled (ρ_{11}) and the guessing (ρ_{22}) state. The correlations for the LBA parameters range from high ($r = 0.77$ for $v_1^{(1)}$) to nearly perfect ($r = 0.99$ for τ) and the points lie close to the identity line, suggesting good recovery of the LBA parameters. An exception is the parameter σ , which shows a pattern of underestimating the true values, if the true value is relatively high.

As for the parameters characterizing the evolution of the latent states, the recovery of the initial state probability is sub optimal ($r = 0.22$). This is expected, as there is not much information in the data about this parameter (it

Fig. 8 Parameter recovery using posterior expectation. Correlation plots between the true values (x-axis) and the estimated values (y-axis). The slope line shows the identity function

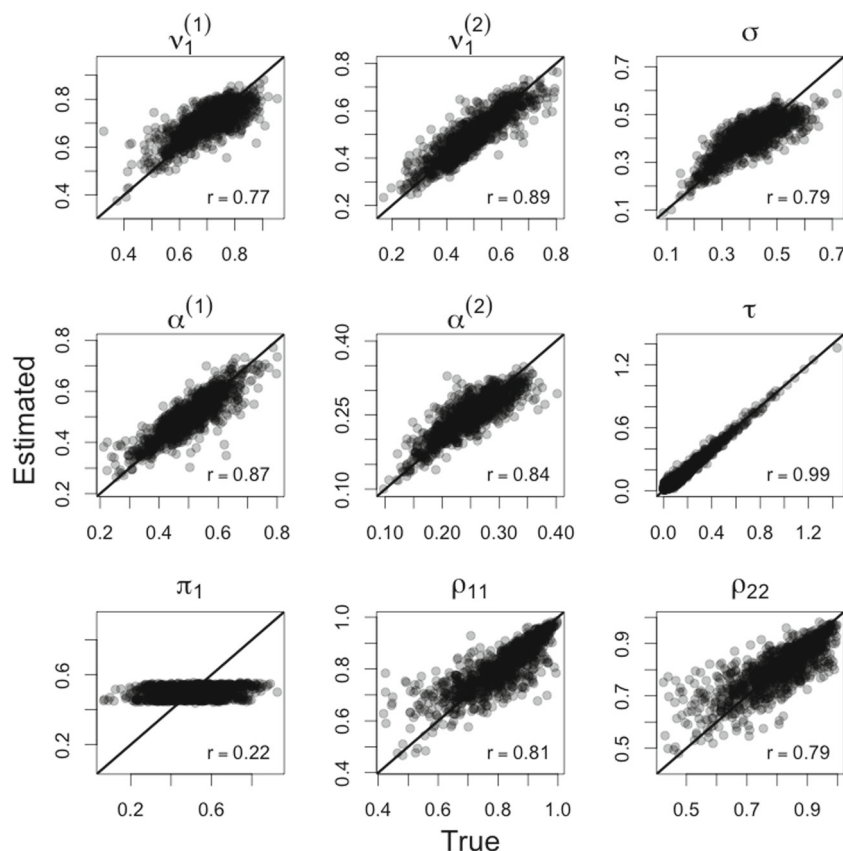


Table 3 The relative frequency with which 50% and 80% credible interval contained the true parameter value. The numbers in the brackets correspond to the 95% Jeffreys credible interval for binomial proportion (Brown et al., 2001)

	50% CI coverage	80% CI coverage
$v_1^{(1)}$	0.52 [0.49, 0.55]	0.79 [0.76, 0.82]
$v_1^{(2)}$	0.48 [0.45, 0.51]	0.79 [0.76, 0.82]
σ	0.51 [0.48, 0.54]	0.82 [0.80, 0.85]
$\alpha^{(1)}$	0.49 [0.45, 0.52]	0.78 [0.76, 0.81]
$\alpha^{(2)}$	0.51 [0.48, 0.54]	0.81 [0.79, 0.84]
τ	0.50 [0.47, 0.53]	0.81 [0.79, 0.84]
π_1	0.49 [0.45, 0.52]	0.80 [0.78, 0.83]
ρ_{11}	0.52 [0.49, 0.56]	0.83 [0.81, 0.86]
ρ_{22}	0.51 [0.48, 0.54]	0.80 [0.77, 0.82]

mostly depends on the state of the first trial), and so it is highly dependent on the prior. This parameter is not to be interpreted, however, unless the model is fitted on repeated trial sequences (so that there are more “first” trial observations). The recovery of the two “dwelling” probabilities are satisfactory.

Coverage of the Credible Intervals

Using the MCMC samples, we computed the 50% and 80% central credible intervals for each parameter under each fitted model (that did not label switch), and checked whether the true value of the parameter lies within that interval. Table 3 shows that the relative frequencies with which the CIs cover the true value is very close to the nominal value of the confidence level. Thus, we did not observe that the credible intervals would be poorly calibrated with respect to their frequentist properties. It is important to keep in mind, though, that this is not a proof of well calibrated CIs in general (e.g., for all possible parameter values and all confidence levels).

Conclusion

We followed general recommendations for a principled Bayesian workflow for building and validating bespoke cognitive models (Schad et al., 2019; Tran et al., 2020; Kennedy et al., 2019). Knowledge about data typical in two-choice speeded decision tasks was used to define the prior distributions on the model parameters. The MCMC procedure yielded accurate approximations of the posterior distributions using simulation-based calibration. SBC further yielded good results except for three parameters for which slight bias could have potentially occurred. Model sensitivity analysis revealed that the model is able to learn about

the parameters of interest while not introducing substantial systematic bias to the estimates. The standard parameter recovery resulted in acceptable results. Further, the 50% and 80% credible intervals had coverage probabilities at their nominal levels. Results of the simulation study hence suggest that further work on improving the model is not absolutely necessary before applying it to real data.

Example: Dutilh et al. (2010) Study

This section demonstrates the use of our model on a real data set from an experiment reported by Dutilh et al. (2010). In this experiment, 11 participants took part in a lexical decision task (participants A–C in Experiment 1a and participants D–G in Experiment 1bL) and perceptual decision task (participants H–K in Experiment 1bV). Despite the fact that the experiments are based on a different modality, the analysis stayed the same as the data have the same structure regarding the application of the HMM. Specifically, participants were asked to give answers on a two-choice task with varying degrees of pay-off for response time and response accuracy: the sum of the pay-off was a given constant, but the difference between them varied, thus leading to trials preferring accuracy (high reward for getting the answer correctly) to trials preferring speed (high reward for responding fast). Dutilh et al. (2010) originally fitted a two state HMMs where the emission distribution for the response times was assumed log-normal, and the distribution for the responses a categorical (i.e., assuming independence of response times and accuracy after conditioning on the state). Here, the EAM HMM model is applied to each of the participants separately, and the model fit is assessed using posterior predictives.

Method

We fitted each participants’ data using the model described in “Model” and priors developed in “Prior Predictives”. Specifically, for each participant, we ran eight MCMC chains with a 1,000 warmup and 1,000 sampling iterations using Stan (Carpenter et al. 2017), with the tuning parameter δ_{adapt} increased to 0.9. Starting points were randomly generated from the prior. Some initial values yielded likelihoods that were too low, leading to failure of the chain initialization. If seven out of the eight chains failed to initialize, the model was reran. If at least two chains managed to run, we inspected the Gelman-Rubin potential scale reduction factor \hat{R} (Gelman & Rubin, 1992), traceplots of the MCMC chains, and parameter estimates, to detect possible label switching. Label switching was identified if $\bar{\alpha}^{(1)} < \bar{\alpha}^{(2)}$ or $\bar{v}_1^{(1)} < \bar{v}_2^{(2)}$ (the two conditions coincided in 100% of the cases). If label switching occurred,

we reran the eight chains. Once we were able to run at least two chains without label switching, we proceeded to fit data from another participant.

Results

Model fit for two participants needed to be run three times and for one participant five times due to seven chains failing to initialize. Further, models needed to be rerun twice for one participant and three times for four participants due to between chain label switching. The final fits for two participants ended with two valid chains, for six participants with three valid chains, and for three participants with four valid chains. Therefore, the number of posterior samples used for inference ranged between 2,000 and 4,000. None of the models yielded divergent transitions. All \hat{R} statistics range between 0.99 and 1.01, and traceplots of the MCMC chains show typical caterpillar shape without a visible drift. Thus, the final model fits do not exhibit convergence issues.

For each participant, we performed several fit diagnostics, to assess whether (and how) the model misfits the data. In the interest of brevity, results for only the first participant from each of the sub-experiments are shown (i.e.,

participant A, participant D, and participant H). The rest of the results can be found online at https://github.com/Kucharssim/hmm_slba/tree/master/figures.

First, we simulated the posterior predictives for response times and accuracy and plotted them against the observed data. Figure 9 shows the posterior predictive distribution for the response times summarized as 80% and 50% quantiles of the posterior predictive distribution for each trial (light red and dark red, respectively), and the median of the posterior predictive distribution (red line). The black line shows the observed response times at a particular trial. Figure 10 shows the posterior predictive distribution for the responses. Specifically, the red line shows the predicted probability of a correct response for a particular trial, whereas the black dots points the observed responses. For ease of the visual comparison, the observed responses were smoothed by calculating their moving average with a window of 10 trials, which is shown as a black line.

In general, the posterior predictives capture the observed data well. Specifically, the model is able to replicate the bi-modality of the response times and captures the runs of trials with predominantly correct responses relatively well. The model also seems to capture correctly that the

Fig. 9 Posterior predictives for the response times for three participants. Only the first 300 trials are shown

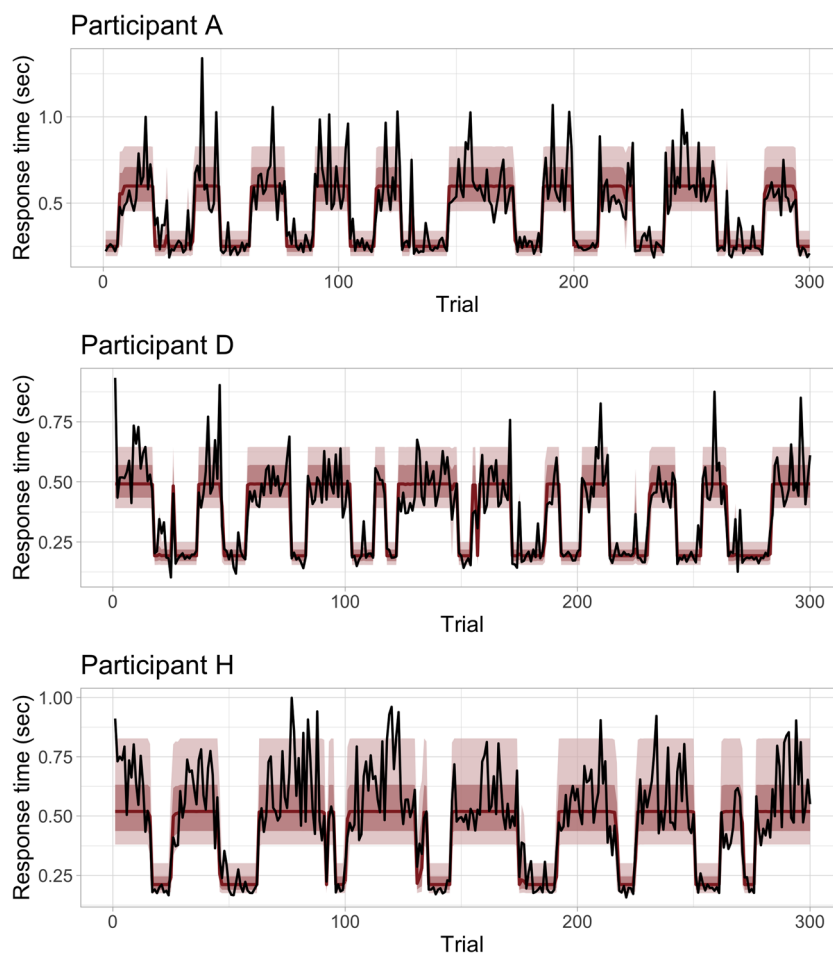
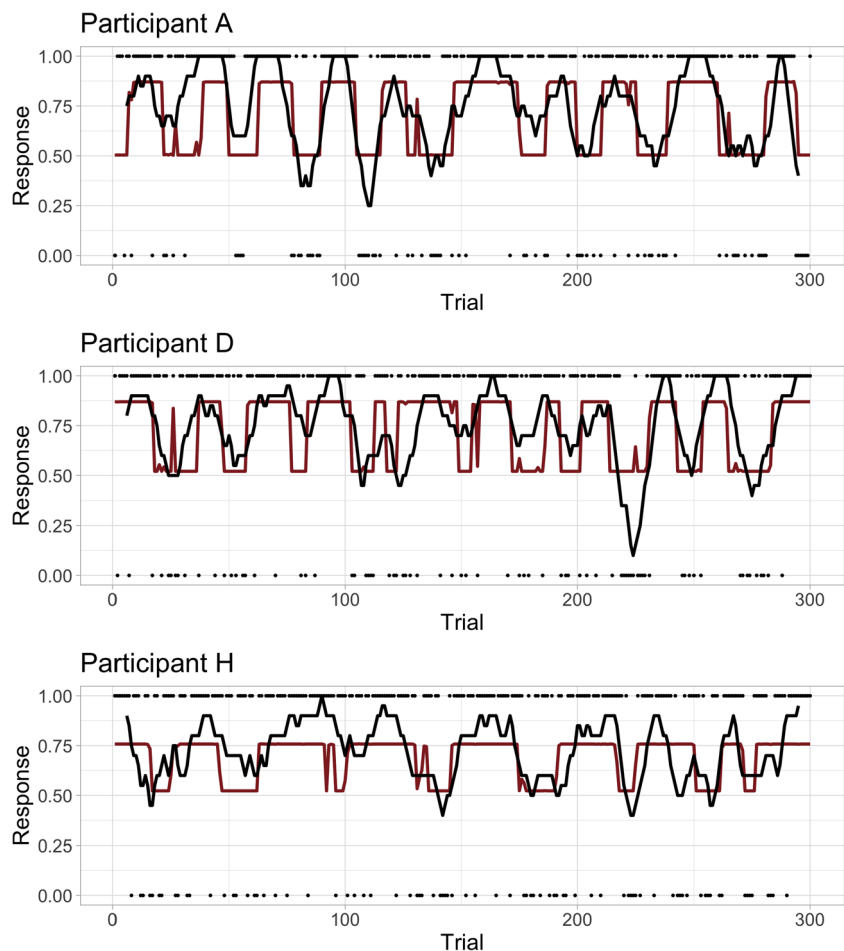


Fig. 10 Posterior predictives for the responses for three participants. Only the first 300 trials are shown



response times under the guessing (fast) state have smaller variance than under the controlled state. However, for some participants, there seem to be many outliers (i.e., slow responses) that are not predicted by the model, suggesting that the model of the response times has perhaps tails that are too thin.

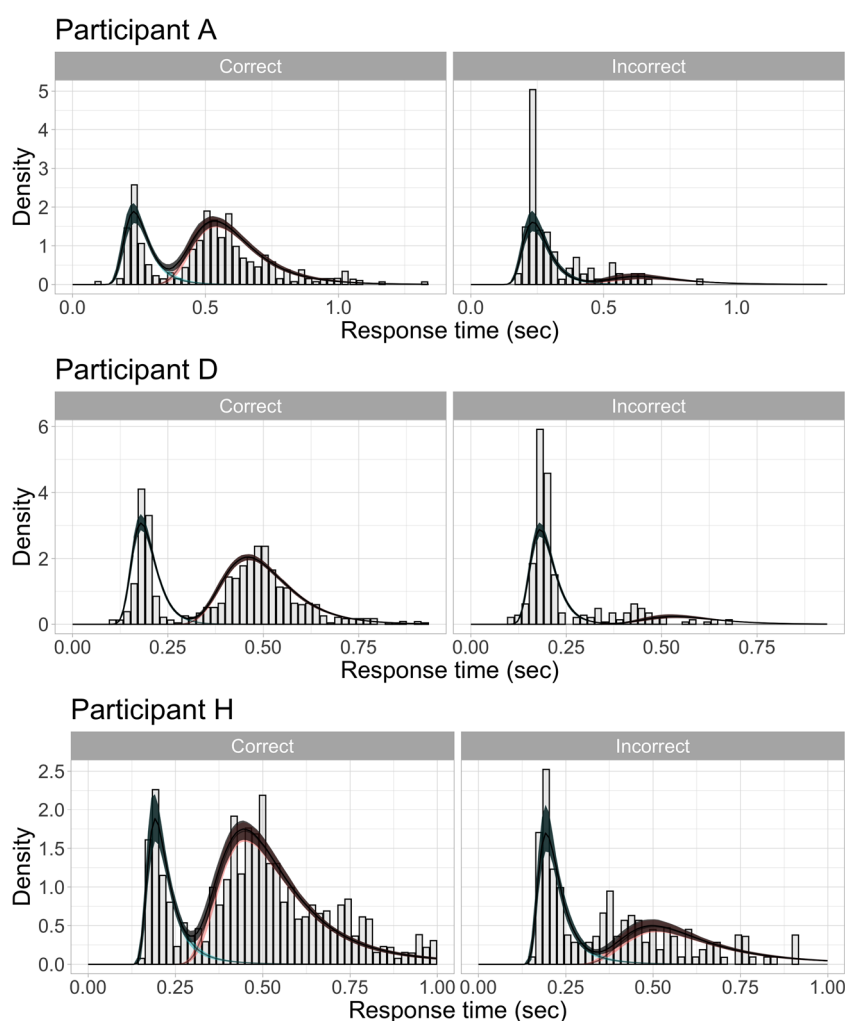
We also assessed how well the model predicts the response time distributions for correct and incorrect responses. Figure 11 shows the observed response times of the correct and incorrect responses as histograms, overlaid with the predicted density of the response times — shown as a black line and 90% CI band. Further, the blue and red lines show the densities under the guessing and controlled state, respectively. Figure 12 shows the observed and predicted cumulative distribution functions conditioned on the state and response.

The distribution plots show good model fits, as the bimodality of the response times is captured correctly, as well as the proportions of correct and incorrect answers under the states. However, for some participants, there are clear signs of a slight misfit. For example, the predicted distribution of the response times of incorrect answers under the controlled state is shifted slightly to the right compared

to the empirical distribution (this shift is the most visible for participant H). Further, there is a general tendency of the model to overestimate the variance of the response times under the guessing state, which might be a consequence of equating the standard deviation of the drift rate (σ) across all accumulators and states. Another possibility would be to enable bias, by setting different decision boundaries for each of the accumulators. These alterations to the model would increase its flexibility and should be validated using simulations - therefore, such additions should be the focus of future projects. In general, the tendency of the model to imply slightly slower incorrect responses than the data suggests, could be also caused by the fact that the number of incorrect responses under the controlled state is low, generally about 10% of the trials (see Fig. 12). It is possible that the likelihood is then dominated by the distribution of the responses under the guessing state, thus favoring a better fit towards them.

Parameter estimates for each participant are attached in Appendix B. Posterior contraction for all participants was close to one for most of the parameters, indicating that there occurred substantial updating of the priors through

Fig. 11 Observed and predicted response times distribution of correct and incorrect responses



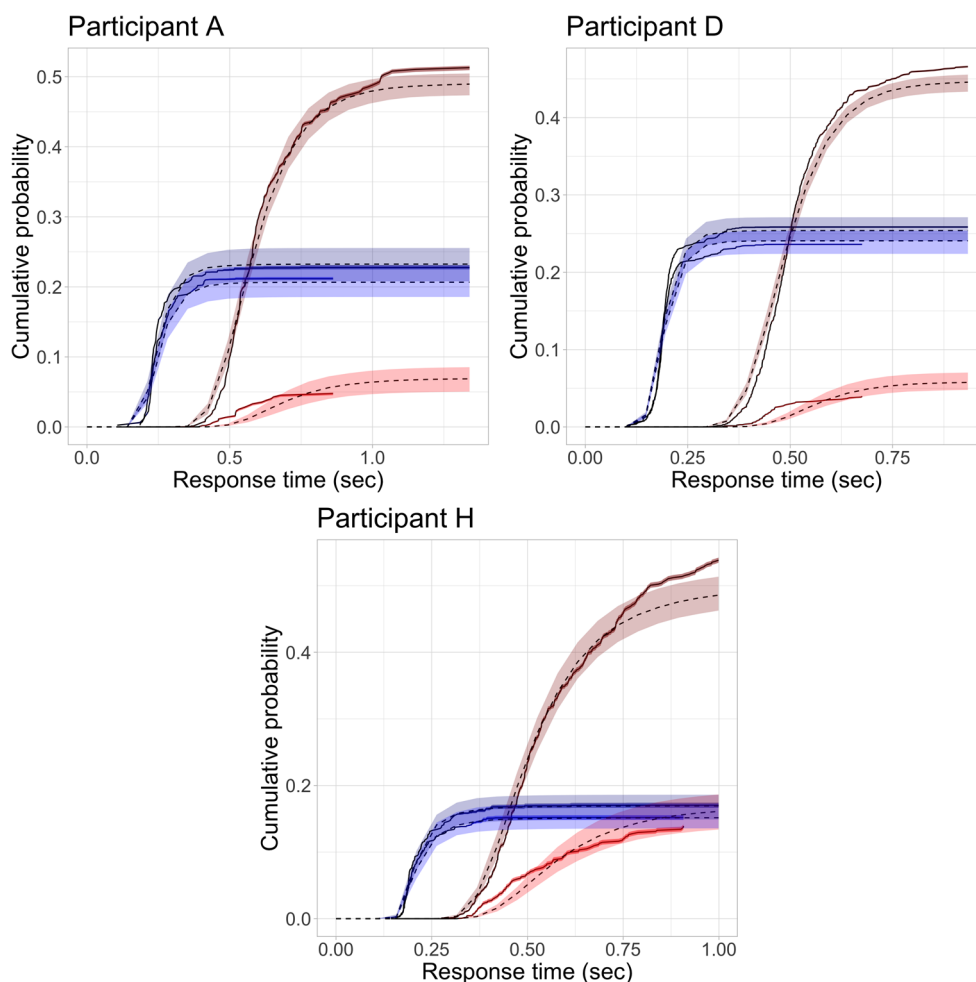
the observed data, in line with the simulation results which showed strong updating of priors despite relatively modest number of trials ($n = 200$) in the simulations. An exception was the parameter π_1 which does not update much, a result that was expected following the simulation results as well. Although there seems to be variability between participants' parameter estimates, there are common patterns that to some degree apply to all participants. Generally, the states of the HMMs are sticky, with a probability of remaining in the current state at about 90% of the trials for both of the states. This percentage is (likely) dependent on the experimental design of Dutilh et al. (2010) who varied the pay-off balance in a structured way depending on the participant's actions, and should not be interpreted as a general tendency of people to stick in the current state to exactly this extent.

As for the parameters that were held fixed across states and accumulators, the non-decision time τ is negligible for the majority of participants; the longest non-decision time occurred for participant B with about 0.11 s (110 ms), with some participants as short as about 0.01 s (10 ms).

Non-decision time is largely informed by the fastest responses in the data (i.e., the shortest response time gives the upper bound of the parameter). It is possible that loosening up equality constraint between the states would reveal that non-decision time is larger under the controlled state than under the guessing state, representing additional encoding time and executing a motoric response after a decision is made; which could also slightly improve the model fit especially regarding the relatively more variable response times under the controlled state. Relatively surprising were the values of the standard deviation of the drift rates σ , with posterior means ranging between 0.13 and 0.27 — quite smaller than specified by the priors ($\sigma \sim \text{Gaussian}(0.4, 0.1)_{(0, \infty)}$) — suggesting that the variability of the response times is smaller than implied by the prior. Future studies should pay specific attention to variability of the response times in prior predictive simulations.

Shorter response times in the actual data compared to the prior predictive expectations resulted also in a relative mismatch between the prior settings for the decision

Fig. 12 Observed and predicted cumulative distribution conditioned on the state (blue = guessing, red = controlled) and response (dark = correct, light = incorrect)



boundaries under the two states. Specifically, the posterior means of the decision boundary under the controlled state ranged between 0.24 and 0.37 (whereas the prior was set $\alpha^{(1)} \sim \text{Gaussian}(0.5, 0.1)_{(0, \infty)}$). The posterior means of the decision boundary under the guessing state was as low as between 0.08 and 0.18 (prior $\alpha^{(2)} \sim \text{Gaussian}(0.25, 0.05)_{(0, \infty)}$).

As expected, the average drift rate of the correct response under the guessing state is usually very close to 0.5, implying 50% accuracy. Under the controlled state, the posterior mean of the average drift rate of the correct response ranged between 0.58 and 0.65. This is slightly smaller than the prior expectation (which on average expects about 0.7), although it still leads to relatively high accuracy (at minimum 75%, and leading to accuracy as high as 90%) due to the small standard deviations of the drift rates.

Thanks to the fact that our model is an EAM model, it is possible to inspect the pattern of the discontinuous speed-accuracy trade-off within and between participants in terms of the latent cognitive parameters that control speed of the evidence accumulation (v) and the response caution (α). Figure 13 shows this between state trade-off and reveals striking similarity between participants.

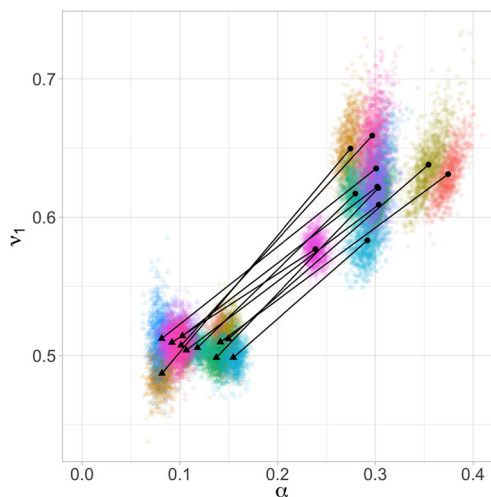


Fig. 13 Speed-accuracy trade-off for all participants in the Dutilh et al. (2010) data set. Black dots show the posterior mean of each participants' decision boundary ($\alpha^{(1)}$) and drift rate for the correct response ($v_1^{(1)}$) under the controlled state, triangles the same but under the guessing state. Lines connect the posterior means for separate participants. Colored points show the samples from the joint posterior distributions

General Conclusion and Discussion

This article presented a robust implementation of a model that combines an EAM with an HMM structure. To our knowledge, this is the first successful implementation combining both structures in one model. The model was built to capture the two state hypothesis following from the phase transition model of the speed-accuracy trade-off (Dutilh et al., 2010) — that there is a guessing and a controlled state between which participants switch. This hypothesis can be represented by an HMM structure. Compared to previous HMM applications on speeded decision tasks, our model uses an EAM framework for the joint distributions of the responses and response times, and thus enables inference on latent cognitive parameters, such as response caution or drift rate (Evans & Wagenmakers, 2019).

The model was validated using extensive simulations and by applying it to real data. The simulations suggested that the model implementation was robust and did not show pathological behavior. Further, the model achieved good parameter recovery and coverage probabilities of the credible intervals. In the empirical example, the model was fitted to eleven participants who partook in the Dutilh et al. (2010) study. The results demonstrate that the model shows a good fit to the data and is able to capture most of the patterns in the data. However, the model also showed a slight systematic misfit because the predicted error responses under the controlled state were slower than that of the data (a typical example of a phenomenon known as fast errors; Tillman and Evans, 2020). The results suggested quite strong consistency between participants in terms of the speed-accuracy trade-off — suggesting that the inaccessibility region (i.e., a region of speed of accumulation and response caution which “cannot be accessed”, resulting in switching between two discrete states) predicted by the phase transition model could be qualitatively similar across participants (see Fig. 13).

We used a full Bayesian framework in this article, and with it comes the perks of defining the prior distributions on the parameters. Setting well behaved priors is important in any Bayesian application as they define the subset of the parameter space that generates data that are expected in a particular application of the model. Because the EAMs can cover a lot of heterogeneous experimental paradigms (with heterogeneous scales of the data), it is important to decide on priors in respect to the specific application of the model, preferably after consulting related research literature, careful reasoning about the experimental design and the particular parameterization of the model. The empirical analysis pointed to some discrepancies between empirical parameter estimates and their priors that highlight misalignment between the priors and the data. Ideally, such

discrepancies would be minimized to avoid a prior-data conflict (possibly leading to problems with estimation, Box 1980; Evans & Moshonov 2006). In our application, the discrepancy between the priors and the data arose mainly because we a priori expected longer and more variable response times than was the case in the Dutilh et al. (2010) study. For the purpose of model validation through extensive simulation, such discrepancy is not a critical problem as the simulation covered cases with potentially more variability and outliers (which usually cause problems in fitting), thus exposing the model to a robustness test.

It is important to reiterate that the priors in this model also serve another purpose: to solve the label switching problem. As is commonly the case in HMMs, the current model is identified only up to the permutation of the state labels. The priors in this article were used to nudge the model towards one specific permutation — to associate the first state with the controlled response, and the second state with the guessing response. Such use of the priors was possible because we specifically assumed the controlled and guessing state, and followed the implications from the theory about them (Dutilh et al., 2010). In case the expectation regarding the state identity is more vague (e.g., when expecting only that the distributions might be multimodal), such use of priors becomes much more problematic on both the conceptual and practical level. On the other hand, some prior specifications could have been even more informative in the current application. For instance, under the controlled state, the drift rate for the correct response should be higher than the drift rate for the incorrect response as the other alternative would imply that the respondent’s performance is below chance level.

Despite our efforts to solve label switching using informative priors, the issue of switching labels still persists, albeit to a lesser degree than without informative priors. Specifically, the use of soft order constraints (by specifying prior distributions that heighten prior probability of a specific state configuration) does not ensure that the labels do not switch at all. To this end, we were forced to perform additional checks of label switching to ensure that the model converged to the solution we preferred, and refitting the model if it did not. Virtually the same estimation results would have been obtained if traditional order constraints were used, by effectively truncating the parameter space to the region which corresponds to the appropriate state interpretation, although in case one would want to perform model comparison using marginal likelihoods, the decision of whether or not to use order restriction would make a difference. Implementing order restrictions would also make it harder to reason about theoretically justified prior specification. For the sake of simplicity, this article did not focus on developing such approach, as its focus was to demonstrate the possibility of combining EAMs with

HMMs at least in estimation context. Developing proper ways how to identify the model using order constraints, set reasonable priors, and compute marginal likelihoods would be additional ways how to take the current modeling framework towards more general applications.

One of the future applications would be to actually put the continuous and discontinuous debate under a test. In this article, we presented a model that assumes both discontinuous, between state trade-off, and continuous, within state trade-off inherent to the EAM. Utilizing Bayesian framework makes it naturally attractive to use marginal likelihoods to compare simple EAMs, HMM combined with an EAM, and a HMM that assume local independence of response times and accuracy, to assess which of the hypotheses are supported by the data. Although methods for estimating marginal likelihoods for EAMs are available (Evans & Brown, 2018; Gronau et al., 2019), the HMM extensions will lead to further problems, as estimating marginal likelihoods for finite mixture models and HMMs is a notoriously difficult problem (Frühwirth-Schnatter, 2004). Nevertheless, combining clever constraints (so as to prevent label switching) and development of principled priors would enable the use of efficient techniques for estimating marginal likelihoods such as bridge sampling and its extensions (Meng & Wong, 1996; Gronau et al. 2017; Gronau et al., 2019; Gronau et al., 2017), which are now becoming more available than ever. Of course, multi-model inference would also benefit from simulation-based calibration approaches build on similar principles as that of single model inference shown in this article (Schad et al., 2021).

An alternative to identifying the HMMs using the priors is to assume functionally different emission distributions under the states. For example, as Dutilh et al. (2010) point out, it is questionable to assume that guessing requires evidence to make a response. Therefore, using an EAM to represent the guessing state probably leads to model misspecification, as under guessing there is no evidence accumulation (about the correct response). Such misspecification could be fixed, for example, by assuming that the response time of guessing is just a simple response time (Luce, 1991), and model it appropriately by a single accumulator independent of the response (which would be a categorical variable with proportion of correct answer fixed at 0.5). In the context of the phase transition model, such an assumption could further improve the model.

Additional advantages of utilizing Bayesian inference and implementation in Stan is the relative ease with which the model could be extended from single-participant model to multiple-participants model and let the individual parameters be estimated in a hierarchical structure. Hierarchical models have the advantage that they can improve individual estimates by pooling information across the sample. Such approach would also improve the amount of information

used for estimating the prior probability of the starting state, which is poorly identified in the single-participant model.

In this article, we used a minimal linear ballistic model to ensure computational stability of the model. However, such a model can hardly be considered adequate for characterizing all phenomena of the speeded decision paradigm, and the current results already revealed some ways in which the current model misfits the data. Thus, it is desirable to find ways how to extend or improve the current model, while ensuring that the quality of inferences and implementation does not decline. One alternative to improve the current model is to use the full LBA model where the variability of the starting point is not fixed at zero (Brown & Heathcote, 2008). Another would be to build on a different evidence accumulation mechanism (such as replacing the ballistic accumulation with sequential sampling models) — for example, the Diffusion Decision model (DDM, Ratcliff & McKoon, 2008) or the Racing diffusion model (Tillman et al., 2020). Regardless of which framework will be in the end more successful in combination with a HMM, we believe it is important to start with a minimal existing model that captures the most crude phenomena from the speeded decision framework, and expand from there. In the case of a DDM, that would be to start with the simplest four parameter model because it can be implemented in a fast and robust way (Wabersich & Vandekerckhove, 2014; Navarro & Fuss, 2009) and generally focus on the most important sources of variability at first (Tillman et al., 2020). Then — provided that model validations are satisfactory — it is possible to add more parameters. In each stage of the model building, it is important to stick to the model validation procedures, some of which were demonstrated in the current article.

Further development and additions to the model should probably also be combined with simplifications. Such simplifications, as for example, simplifying the distribution under the guessing state (as discussed above) can provide more computational stability and provide degrees of freedom to extend the model under the controlled state.

The current model provides a proof of principle of a combination of an EAM with an HMM, and as such can lead to further interesting applications and extensions, as it opens new possibilities regarding modeling continuous and discontinuous patterns of response times and accuracy in a single modeling framework. Although the current article focused solely on speeded decision tasks, questions about the continuous and discontinuous relations between response times and accuracy is ubiquitous in higher cognitive applications as well, including study of more complex cognitive tasks and development of strategies used to solve these tasks (van der Maas & Jansen, 2003; Raijmakers et al., 2014; Hofman et al., 2018). An interesting feature of higher level cognitive tasks that might be relevant

to explore using the current framework is the emergence of more efficient strategies, that lead to qualitatively better response accuracy as well as shorter response times. Such strategies have been described in many applications, such as multiplication tasks (Hofman et al., 2018), Mastermind game (Gierasimczuk et al., 2013; Kucharský et al., 2020), or Progressive matrices tasks (Vigneau et al., 2006; Laurence et al., 2018). Combination of HMM with EAM in this context would enable uncovering different relations between response times and accuracy depending on whether we look within or between strategies — it is possible to imagine that an efficient strategy would be faster and more accurate than less efficient strategy, but within those strategies separately, we will see the traditional speed-accuracy trade-off whereby increasing response caution increases accuracy at the cost of speed, which would be captured by the EAM part of the model.

Appendix A. Derivation of the simplified LBA model

Here, we provide the derivation of the likelihood function for the simplified LBA model. We assume that each choice option is associated with an accumulator of evidence. These accumulators are independent of each other and the first accumulator that reaches its decision threshold launches the decision associated with it. This leads to general race equations (Heathcote & Love, 2012), the probability density of observing response a with the reaction time rt comprises of the probability density that an accumulator associated with response a hits the threshold at time rt times the probability that none of the other accumulators has hit the threshold at an earlier time point:

$$sLBA(rt, a|v, \sigma, \alpha, \tau) = f(rt|v_a, \sigma_a, \alpha_a, \tau_a) \times \prod_{k \neq a} [1 - F(rt|v_k, \sigma_k, \alpha_k, \tau_k)], \quad (A.1)$$

with v_a the mean drift rate, σ_a the standard deviation of drift rate, α_a the decision boundary, and τ_a the non-decision time for the accumulator a .

The density of the passage time for each accumulator $f(rt)$ is specified as follows:

$$\begin{aligned} rt &= \tau + t \\ t &= \frac{\alpha}{\delta} \\ \delta &\sim \text{Gaussian}(v, \sigma)_{(0, \infty)}. \end{aligned} \quad (A.2)$$

We assume that the passage time is a sum of the non-decision time and the decision time t , where the decision time is a result of a linear rise of evidence towards a decision threshold α , at a drift rate δ drawn randomly from a Gaussian distribution with mean v and standard deviation

σ , truncated at 0 on the lower bound. The truncation is assumed because we do not allow for the possibility of a non-response (i.e., that all drifts in a particular trial are negative, thus never cross the decision threshold). We do not assume any randomness in the parameters τ, α, v and σ , hence, the only missing piece in deriving $f(rt)$ is the change of variables $rt = \tau + \alpha/\delta$.

First, we derive the density of the latent drift (δ), which is defined as a truncated normal distribution for $\delta \geq 0$ and zero otherwise:

$$g(\delta|v, \sigma) = \frac{1}{\sigma} \times \frac{\phi\left(\frac{\delta-v}{\sigma}\right)}{1 - \Phi\left(\frac{-v}{\sigma}\right)}, \quad (A.3)$$

where $\phi(\cdot)$ is the pdf and $\Phi(\cdot)$ the cdf of the standard normal distribution, respectively.

Next, we determine the density of the variable t , which arises as a scaled reciprocal truncated normal variable for $t \geq 0$ and zero otherwise (see also Nakahara et al. 2006):

$$h(t|v, \sigma, \alpha) = \frac{\alpha}{t^2} \times g\left(\frac{\alpha}{t}|v, \sigma\right) \quad (A.4)$$

Finally, to obtain the density of the passage time rt , we shift the distribution of the decision time t by τ , which results in the following pdf:

$$f(rt|v, \sigma, \alpha, \tau) = h(rt-\tau|v, \sigma, \alpha) = \frac{\alpha}{(rt-\tau)^2} \times g\left(\frac{\alpha}{rt-\tau}|v, \sigma\right), \quad (A.5)$$

for $rt > \tau$ and zero otherwise.

The cumulative probability function of the passage times, $F(rt|v, \sigma, \alpha, \tau)$, is relatively easier to compute, by realizing that the only source of randomness in this model is the distribution of the latent drift δ . Thus,

$$\begin{aligned} P(rt \leq X) &= P(\delta \leq Y) \\ Y &= \frac{\alpha}{X - \tau}, \end{aligned} \quad (A.6)$$

which leads to

$$F(rt|v, \sigma, \alpha, \tau) = G\left(\frac{\alpha}{rt-\tau}|v, \sigma\right), \quad (A.7)$$

where $G(\cdot|v, \sigma)$ is the cdf of a normal distribution truncated at zero.

Identifiability and a Minimal Model

If we had only response time data without choices (e.g., from a single choice response time task), the entire likelihood would be given by the distribution of the passage times for a single accumulator $f(rt|v, \sigma, \alpha, \tau)$. Such distribution is a ballistic analogue to the shifted Wald distribution (otherwise known as inverse Gaussian distribution) of response times (Chhikara & Folks, 1988;

Anders et al., 2016), and would similarly require fixing one of the parameters ν , σ , or α to achieve identifiability.

Once we have multiple choice tasks, it is possible to estimate more parameters per accumulator, as is the case for the LBA (Brown & Heathcote, 2008). However, some identifiability constraints still need to be put in place. In this paper, we use the following set of identifiability constraints:

$$\sum_i \nu_i = 1,$$

$$1 \geq \nu_i \geq 0.$$

That is, we use the sum-to-one constraint common for the LBA model (Brown & Heathcote, 2008; Visser & Poessé, 2017), and make it even slightly more severe by assuming that no average drift rate can be negative. The second, additional constraint is convenient for Bayesian implementation as it allows using Dirichlet priors on the drifts.

The simplified LBA model can be achieved by additionally assuming that the non-decision time is equal between the accumulators — usually EAM models assume that non-decision time is by definition the time spend on processes that are not related to the decision — such as encoding and executing motoric responses (Evans & Wagenmakers, 2019). Further, we may equate σ and α between the accumulators. The minimal model for a two-choice task would then contain five parameters: $\theta = (\nu_1, \nu_2, \sigma, \alpha, \tau)$, of which four of them are “free” (ν_1 and ν_2 are collinear due to the sum-to-one constraint). In general, the simplified LBA model would have $K + 3$ parameters (of which $K + 2$ are free), where K is the number of response options (accumulators).

Appendix B. Parameter estimates of the Dutilh et al. (2010) data

Table 4 Descriptives of the posterior draws for Participant A from Dutilh et al. (2010)

Parameter	Mean	Median	SD	Quantile		\hat{R}	ESS		Contraction
				5%	95%		Bulk	Tail	
$\nu_1^{(1)}$	0.63	0.63	0.02	0.61	0.66	1.001	3319	2939	0.975
$\nu_1^{(2)}$	0.51	0.51	0.01	0.49	0.53	1.000	4090	2860	0.992
$\alpha^{(1)}$	0.37	0.37	0.01	0.36	0.39	1.003	2540	2191	0.991
$\alpha^{(2)}$	0.14	0.14	0.00	0.13	0.15	1.002	2069	2483	0.992
σ	0.16	0.16	0.01	0.15	0.18	1.000	2250	2672	0.991
τ	0.01	0.01	0.01	0.00	0.02	1.003	1602	1690	0.999
π_1	0.46	0.46	0.15	0.22	0.70	1.001	4497	2559	0.051
ρ_{11}	0.92	0.92	0.02	0.88	0.95	1.001	4483	2909	0.973
ρ_{22}	0.89	0.90	0.02	0.85	0.93	1.002	3901	2381	0.960

Table 5 Descriptives of the posterior draws for Participant B from Dutilh et al. (2010)

Parameter	Mean	Median	SD	Quantile		\hat{R}	ESS		Contraction
				5%	95%		Bulk	Tail	
$\nu_1^{(1)}$	0.65	0.65	0.02	0.62	0.68	1.004	1837	1704	0.960
$\nu_1^{(2)}$	0.49	0.49	0.01	0.47	0.51	1.000	3065	2167	0.990
$\alpha^{(1)}$	0.27	0.27	0.01	0.26	0.29	1.000	1979	1934	0.994
$\alpha^{(2)}$	0.08	0.08	0.01	0.07	0.09	1.005	1168	1061	0.978
σ	0.18	0.18	0.02	0.16	0.21	1.003	1271	1370	0.975
τ	0.11	0.11	0.01	0.08	0.13	1.005	1127	1063	0.996
π_1	0.45	0.45	0.14	0.22	0.70	1.001	3029	1997	0.070
ρ_{11}	0.90	0.90	0.02	0.87	0.93	1.001	3038	1897	0.978
ρ_{22}	0.84	0.84	0.03	0.80	0.89	1.001	3049	2364	0.946

Table 6 Descriptives of the posterior draws for Participant C from Dutilh et al. (2010)

Parameter	Mean	Median	SD	Quantile		\hat{R}	ESS		Contraction
				5%	95%		Bulk	Tail	
$v_1^{(1)}$	0.64	0.64	0.02	0.61	0.68	1.001	2190	1837	0.953
$v_1^{(2)}$	0.51	0.51	0.01	0.49	0.53	1.000	2883	2091	0.987
$\alpha^{(1)}$	0.35	0.35	0.01	0.34	0.37	1.002	1985	1831	0.986
$\alpha^{(2)}$	0.15	0.15	0.01	0.14	0.16	1.001	1693	1564	0.984
σ	0.17	0.17	0.01	0.15	0.19	1.001	2022	1734	0.984
τ	0.01	0.01	0.01	0.00	0.03	1.002	1358	1622	0.998
π_1	0.46	0.46	0.14	0.23	0.69	1.001	3171	2226	0.120
ρ_{11}	0.91	0.92	0.02	0.88	0.94	1.000	3279	1883	0.968
ρ_{22}	0.87	0.88	0.03	0.82	0.92	1.002	2925	2082	0.937

Table 7 Descriptives of the posterior draws for Participant D from Dutilh et al. (2010)

Parameter	Mean	Median	SD	Quantile		\hat{R}	ESS		Contraction
				5%	95%		Bulk	Tail	
$v_1^{(1)}$	0.61	0.61	0.01	0.60	0.62	1.000	2911	2213	0.994
$v_1^{(2)}$	0.50	0.50	0.00	0.50	0.51	1.004	3268	1746	0.998
$\alpha^{(1)}$	0.30	0.30	0.00	0.30	0.31	1.000	2889	1793	0.998
$\alpha^{(2)}$	0.11	0.11	0.00	0.10	0.11	1.001	1391	1591	0.999
σ	0.13	0.13	0.00	0.12	0.14	1.000	2095	2116	0.998
τ	0.00	0.00	0.00	0.00	0.01	1.001	1131	1488	1.000
π_1	0.54	0.54	0.15	0.29	0.78	1.000	3930	2281	0.027
ρ_{11}	0.90	0.90	0.01	0.88	0.92	1.000	3998	2251	0.987
ρ_{22}	0.90	0.90	0.01	0.88	0.92	1.000	3513	1906	0.987

Table 8 Descriptives of the posterior draws for Participant E from Dutilh et al. (2010)

Parameter	Mean	Median	SD	Quantile		\hat{R}	ESS		Contraction
				5%	95%		Bulk	Tail	
$v_1^{(1)}$	0.62	0.62	0.01	0.60	0.65	1.001	2303	2036	0.978
$v_1^{(2)}$	0.50	0.50	0.01	0.49	0.51	1.000	2858	2045	0.996
$\alpha^{(1)}$	0.30	0.30	0.01	0.29	0.32	1.000	1530	1782	0.994
$\alpha^{(2)}$	0.14	0.14	0.01	0.12	0.15	1.001	957	1785	0.983
σ	0.15	0.14	0.01	0.13	0.16	1.001	1458	1674	0.990
τ	0.02	0.01	0.01	0.00	0.04	1.001	899	987	0.997
π_1	0.46	0.45	0.14	0.23	0.70	1.000	2769	1768	0.079
ρ_{11}	0.85	0.85	0.02	0.80	0.88	1.002	2862	1848	0.959
ρ_{22}	0.85	0.85	0.02	0.81	0.89	1.000	2668	1749	0.957

Table 9 Descriptives of the posterior draws for Participant F from Dutilh et al. (2010)

Parameter	Mean	Median	SD	Quantile		\hat{R}	ESS		Contraction
				5%	95%		Bulk	Tail	
$v_1^{(1)}$	0.62	0.62	0.01	0.60	0.64	1.002	1999	2295	0.984
$v_1^{(2)}$	0.51	0.51	0.01	0.50	0.51	1.001	3617	2235	0.998
$\alpha^{(1)}$	0.28	0.28	0.01	0.27	0.29	1.003	1206	1413	0.994
$\alpha^{(2)}$	0.12	0.12	0.01	0.11	0.13	1.004	893	803	0.975
σ	0.16	0.16	0.01	0.14	0.18	1.003	1023	974	0.987
τ	0.05	0.05	0.01	0.02	0.07	1.004	874	815	0.995
π_1	0.45	0.45	0.14	0.23	0.70	1.004	2860	1943	0.102
ρ_{11}	0.91	0.91	0.01	0.88	0.93	1.002	2486	1753	0.986
ρ_{22}	0.91	0.91	0.01	0.89	0.93	1.001	2647	1798	0.988

Table 10 Descriptives of the posterior draws for Participant G from Dutilh et al. (2010)

Parameter	Mean	Median	SD	Quantile		\hat{R}	ESS		Contraction
				5%	95%		Bulk	Tail	
$v_1^{(1)}$	0.58	0.58	0.01	0.56	0.61	1.001	3076	2370	0.977
$v_1^{(2)}$	0.50	0.50	0.01	0.48	0.51	1.000	2903	2069	0.993
$\alpha^{(1)}$	0.29	0.29	0.01	0.28	0.31	1.001	1334	1996	0.990
$\alpha^{(2)}$	0.15	0.16	0.01	0.14	0.17	1.001	1109	1461	0.978
σ	0.17	0.17	0.01	0.15	0.19	1.000	2029	1885	0.986
τ	0.02	0.01	0.01	0.00	0.04	1.001	1049	1069	0.997
π_1	0.46	0.46	0.14	0.23	0.69	1.001	2437	1881	0.114
ρ_{11}	0.89	0.89	0.03	0.84	0.93	1.000	2661	2149	0.953
ρ_{22}	0.88	0.89	0.03	0.84	0.93	1.001	2175	2087	0.946

Table 11 Descriptives of the posterior draws for Participant H from Dutilh et al. (2010)

Parameter	Mean	Median	SD	Quantile		\hat{R}	ESS		Contraction
				5%	95%		Bulk	Tail	
$v_1^{(1)}$	0.64	0.63	0.02	0.61	0.67	1.000	2787	2656	0.959
$v_1^{(2)}$	0.51	0.51	0.02	0.48	0.54	1.001	3673	2653	0.978
$\alpha^{(1)}$	0.30	0.30	0.01	0.29	0.32	1.000	2982	2630	0.991
$\alpha^{(2)}$	0.08	0.08	0.01	0.07	0.09	1.002	1922	1484	0.977
σ	0.27	0.27	0.02	0.23	0.31	1.001	2001	1784	0.938
τ	0.09	0.09	0.01	0.06	0.10	1.002	1825	1520	0.997
π_1	0.55	0.55	0.14	0.30	0.77	1.002	4024	2054	0.067
ρ_{11}	0.94	0.94	0.01	0.92	0.96	1.002	3678	2633	0.989
ρ_{22}	0.88	0.88	0.02	0.84	0.92	1.003	3683	2612	0.958

Table 12 Descriptives of the posterior draws for Participant I from Dutilh et al. (2010)

Parameter	Mean	Median	SD	Quantile		\hat{R}	ESS		Contraction
				5%	95%		Bulk	Tail	
$v_1^{(1)}$	0.62	0.62	0.02	0.60	0.65	1.001	1423	1521	0.968
$v_1^{(2)}$	0.51	0.51	0.01	0.50	0.53	1.000	2217	1289	0.991
$\alpha^{(1)}$	0.30	0.30	0.01	0.29	0.32	1.000	1851	1275	0.993
$\alpha^{(2)}$	0.10	0.10	0.01	0.09	0.12	1.001	899	934	0.978
σ	0.26	0.25	0.02	0.22	0.30	1.001	1074	1177	0.944
τ	0.06	0.06	0.01	0.04	0.08	1.001	854	789	0.996
π_1	0.55	0.55	0.15	0.30	0.80	1.000	2496	1225	-0.001
ρ_{11}	0.91	0.91	0.01	0.89	0.93	1.001	2211	1267	0.986
ρ_{22}	0.90	0.90	0.02	0.88	0.93	1.000	2047	1255	0.982

Table 13 Descriptives of the posterior draws for Participant J from Dutilh et al. (2010)

Parameter	Mean	Median	SD	Quantile		\hat{R}	ESS		Contraction
				5%	95%		Bulk	Tail	
$v_1^{(1)}$	0.58	0.58	0.01	0.56	0.59	1.000	4004	3489	0.992
$v_1^{(2)}$	0.51	0.51	0.01	0.50	0.52	1.002	4176	2785	0.995
$\alpha^{(1)}$	0.24	0.24	0.01	0.23	0.25	1.001	2186	2528	0.996
$\alpha^{(2)}$	0.09	0.09	0.01	0.08	0.10	1.001	1731	1602	0.984
σ	0.18	0.18	0.01	0.16	0.20	1.002	2166	2552	0.988
τ	0.05	0.06	0.01	0.04	0.07	1.002	1674	1606	0.997
π_1	0.45	0.45	0.14	0.22	0.69	1.001	4561	2501	0.103
ρ_{11}	0.94	0.94	0.01	0.92	0.96	1.000	3888	2076	0.991
ρ_{22}	0.89	0.89	0.02	0.86	0.92	1.002	4567	2907	0.977

Table 14 Descriptives of the posterior draws for Participant K from Dutilh et al. (2010)

Parameter	Mean	Median	SD	Quantile		\hat{R}	ESS		Contraction
				5%	95%		Bulk	Tail	
$v_1^{(1)}$	0.66	0.66	0.02	0.63	0.69	1.001	1492	1412	0.953
$v_1^{(2)}$	0.51	0.51	0.01	0.49	0.53	1.000	1757	1341	0.990
$\alpha^{(1)}$	0.30	0.30	0.01	0.28	0.31	1.000	1590	1497	0.992
$\alpha^{(2)}$	0.10	0.10	0.01	0.09	0.11	1.002	769	778	0.985
σ	0.21	0.21	0.02	0.19	0.24	1.000	944	1039	0.973
τ	0.04	0.05	0.01	0.03	0.06	1.002	708	725	0.997
π_1	0.46	0.46	0.15	0.22	0.70	1.000	2083	1334	0.040
ρ_{11}	0.91	0.92	0.02	0.88	0.94	1.002	1898	1417	0.977
ρ_{22}	0.92	0.92	0.02	0.89	0.94	1.000	2218	1371	0.978

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1007/s42113-021-00115-0>.

Author Contribution Ingmar Visser provided the concept of the article. Karel Veldkamp, Šimon Kucharský, and Ingmar Visser conducted initial feasibility study that provided insights into the issues associated with this topic. Šimon Kucharský and N.-Han Tran developed the model presented in this article and drafted the initial manuscript. Šimon Kucharský implemented the model, conducted the simulation study and analyzed the data. N.-Han Tran checked the correctness and reproducibility of the code. All authors contributed to the final version of the manuscript.

Funding Šimon Kucharský was supported by the NWO (Nederlandse Organisatie voor Wetenschappelijk Onderzoek) grant no. 406.10.559.

Code Availability The code and data used in this article are publicly available at https://github.com/Kucharssim/hmm_slba.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Anders, R., Alario, F. X., & van Maanen, L. (2016). The shifted Wald distribution for response time data analysis. *Psychological Methods*, 21(3), 309–327.
- Apgar, J. F., Witmer, D. K., White, F. M., & Tidor, B. (2010). Sloppy models, parameter uncertainty, and the role of experimental design. *Molecular BioSystems*, 6(10), 1890–1900.
- Bogacz, R., Wagenmakers, E.-J., Forstmann, B. U., & Nieuwenhuis, S. (2010). The neural basis of the speed–accuracy tradeoff. *Trends in neurosciences*, 33(1), 10–16.
- Box, G. E. (1980). Sampling and Bayes' inference in scientific modelling and robustness. *Journal of the Royal Statistical Society: Series A (General)*, 143(4), 383–404.
- Brown, L. D., Cai, T. T., & DasGupta, A. (2001). Interval estimation for a binomial proportion. *Statistical Science*, 16(2), 101–117.
- Brown, S. D., & Heathcote, A. (2008). The simplest complete model of choice response time: Linear ballistic accumulation. *Cognitive Psychology*, 57(3), 153–178.
- Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., et al. (2017). Stan: A probabilistic programming language. *Journal of Statistical Software*, 76(1), 1–32.
- Carpenter, R. (1981). Oculomotor procrastination. In D.F. Fisher, R.A. Monty, & J.W. Senders (Eds.) *Eye Movements: Cognition and Visual Perception*. Hillsdale: Lawrence Erlbaum Associates.
- Chhikara, R., & Folks, L. J. (1988). The inverse Gaussian distribution: theory, methodology, and applications. CRC Press.
- Donkin, C., Brown, S., Heathcote, A., & Wagenmakers, E.-J. (2011). Diffusion versus linear ballistic accumulation: different models but the same conclusions about psychological processes?. *Psychonomic Bulletin & Review*, 18(1), 61–69.
- Dutilh, G., Wagenmakers, E.-J., Visser, I., & van der Maas, H.L. (2010). A phase transition model for the speed-accuracy trade-off in response time experiments. *Cognitive Science*, 35(2), 211–250.
- Evans, J. (2008). Dual-processing accounts of reasoning, judgment, and social cognition. *Annual Review of Psychology*, 59, 255–278.
- Evans, M., & Moshonov, H. (2006). Checking for prior-data conflict. *Bayesian Analysis*, 1(4), 893–914.
- Evans, N. (2019). A method, framework, and tutorial for efficiently simulating models of decision-making. *Behavior Research Methods*, 51(5), 2390–2404.
- Evans, N. (2020). Same model, different conclusions: An identifiability issue in the linear ballistic accumulator model of decision-making. PsyArXiv.
- Evans, N., & Brown, S. D. (2018). Bayes factors for the linear ballistic accumulator model of decision-making. *Behavior Research Methods*, 50(2), 589–603.
- Evans, N., & Wagenmakers, E.-J. (2019). Evidence accumulation models: Current limitations and future directions. *The Quantitative Methods for Psychology*, 16(2), 73–90.
- Frühwirth-Schnatter, S. (2004). Estimating marginal likelihoods for mixture and Markov switching models using bridge sampling techniques. *The Econometrics Journal*, 7(1), 143–167.
- Frühwirth-Schnatter, S. (2019). Keeping the balance-ridge sampling for marginal likelihood estimation in finite mixture, mixture of experts and Markov mixture models. *Brazilian Journal of Probability and Statistics*, 33(4), 706–733.
- Gabry, J., & Češnovar, R. (2020). cmdstanr: R Interface to 'CmdStan' [Computer software manual]. Retrieved from <https://CRAN.R-project.org/package=cmdstanr> (R package version 2.19.3).
- Gelman, A., & Rubin, D. B. (1992). Inference from iterative simulation using multiple sequences. *Statistical Science*, 7(4), 457–472.
- Gershman, S. J. (2016). Empirical priors for reinforcement learning models. *Journal of Mathematical Psychology*, 71, 1–6. <https://doi.org/10.1016/j.jmp.2016.01.006>.
- Gierasimczuk, N., van der Maas, H. L., & Raijmakers, M.E. (2013). An analytic tableaux model for deductive mastermind empirically tested with a massively used online learning system. *Journal of Logic, Language and Information*, 22(3), 297–314.
- Gronau, Q. F., Heathcote, A., & Matzke, D. (2019). Computing Bayes factors for evidence-accumulation models using Warp-III bridge sampling. *Behavior Research Methods*, 52, 1–20.
- Gronau, Q. F., Sarafoglou, A., Matzke, D., Ly, A., Boehm, U., Marsman, M., et al. (2017). A tutorial on bridge sampling. *Journal of Mathematical Psychology*, 81, 80–97.
- Gronau, Q. F., Singmann, H., & Wagenmakers, E.-J. (2017). bridge-sampling: An r package for estimating normalizing constants. arXiv:1710.08162.
- Gronau, Q. F., Wagenmakers, E.-J., Heck, D. W., & Matzke, D. (2019). A simple method for comparing complex models: Bayesian model comparison for hierarchical multinomial processing tree models using warp-iii bridge sampling. *Psychometrika*, 84(1), 261–284.
- Gutenkunst, R. N., Waterfall, J. J., Casey, F. P., Brown, K. S., Myers, C. R., & Sethna, J.P. (2007). Universally sloppy parameter sensitivities in systems biology models. *PLOS Computational Biology*, 3(10), e189.
- Heathcote, A., & Love, J. (2012). Linear deterministic accumulator models of simple choice. *Frontiers in Psychology*, 3, 292.
- Hofman, A. D., Visser, I., Jansen, B. R., Marsman, M., & van der Maas, H.L. (2018). Fast and slow strategies in multiplication. *Learning and Individual Differences*, 68, 30–40.
- Jasra, A., Holmes, C. C., & Stephens, D.A. (2005). Markov chain monte carlo methods and the label switching problem in bayesian mixture modeling. *Statistical Science*, 20(1), 50–67.

- Kennedy, L., Simpson, D., & Gelman, A. (2019). The experiment is just as important as the likelihood in understanding the prior: a cautionary note on robust cognitive modeling. *Computational Brain & Behavior*, 2(3), 210–217. <https://doi.org/10.1007/s42113-019-00051-0>.
- Kucharský, Š., Visser, I., Truțescu, G.-O., Laurence, P. G., Zaharieva, M., & Raijmakers, M.E (2020). Cognitive strategies revealed by clustering eye movement transitions. *Journal of Eye Movement Research*, 13(1), 1–20.
- Laurence, P. G., Mecca, T. P., Serpa, A., Martin, R., & Macedo, E.C (2018). Eye movements and cognitive strategy in a fluid intelligence test: Item type analysis. *Frontiers in Psychology*, 9, 380.
- Lee, M. D., Criss, A. H., Devezer, B., Donkin, C., Etz, A., Leite, F. P., et al. (2019). Robust modeling in cognitive science. *Computational Brain & Behavior*, 2(3-4), 141–153.
- Luce, R. D. (1991). *Response times: Their role in inferring elementary mental organization*, 2nd edn. Oxford: Oxford University Press.
- Meng, X.iao.-L.i., & Wong, W. H. (1996). Simulating ratios of normalizing constants via a simple identity: a theoretical exploration. *Statistica Sinica*, 6(4), 831–860.
- Molenaar, D., Oberski, D., Vermunt, J., & De Boeck, P. (2016). Hidden Markov item response theory models for responses and response times. *Multivariate Behavioral Research*, 51(5), 606–626.
- Nakahara, H., Nakamura, K., & Hikosaka, O. (2006). Extended LATER model can account for trial-by-trial variability of both pre-and post-processes. *Neural Networks*, 19(8), 1027–1046.
- Navarro, D. J., & Fuss, I. G. (2009). Fast and accurate calculations for first-passage times in Wiener diffusion models. *Journal of Mathematical Psychology*, 53(4), 222–230.
- Noorani, I., & Carpenter, R. H. S. (2016). The LATER model of reaction time and decision. *Neuroscience & Biobehavioral Reviews*, 64, 229–251.
- Ollman, R. (1966). Fast guesses in choice reaction time. *Psychonomic Science*, 6(4), 155–156.
- R Core Team (2020). R: A language and environment for statistical computing. [Computer software manual]. Vienna, Austria. Retrieved. <https://www.R-project.org/>.
- Raijmakers, M. E., Schmittmann, V. D., & Visser, I. (2014). Costs and benefits of automatization in category learning of ill-defined rules. *Cognitive Psychology*, 69, 1–24.
- Ratcliff, R. (2001). Putting noise into neurophysiological models of simple decision making. *Nature Neuroscience*, 4(4), 336–336.
- Ratcliff, R., & McKoon, G. (2008). The diffusion decision model: Theory and data for two-choice decision tasks. *Neural computation*, 20(4), 873–922.
- Schad, D. J., Betancourt, M., & Vasisht, S. (2019). Toward a principled Bayesian workflow in cognitive science. arXiv:1904.12765.
- Schad, D. J., Nicenboim, B., Bürkner, P.-C., Betancourt, M., & Vasisht, S. (2021). Workflow techniques for the robust use of bayes factors. arXiv:2103.08744.
- Spezia, L. (2009). Reversible jump and the label switching problem in hidden Markov models. *Journal of Statistical Planning and Inference*, 139(7), 2305–2315.
- Stan Development Team (2020). CmdStan: the command-line interface to stan. [urlhttps://github.com/stan-dev/cmdstan/releases/tag/v2.24.0-rc1](https://github.com/stan-dev/cmdstan/releases/tag/v2.24.0-rc1) (Version 2.24.0 release candidate 1).
- Talts, S., Betancourt, M., Simpson, D., Vehtari, A., & Gelman, A. (2018). Validating bayesian inference algorithms with simulation-based calibration. arXiv:1804.06788.
- Tillman, G., & Evans, N. J. (2020). Redefining qualitative benchmarks of theories and models: An empirical exploration of fast and slow errors in speeded decision-making. PsyArXiv.
- Tillman, G., Van Zandt, T., & Logan, G.D (2020). Sequential sampling models without random between-trial variability: the racing diffusion model of speeded decision making. *Psychonomic Bulletin & Review*.
- Timmers, B. (2019). Mixture components in response times: A hidden Markov modeling approach for evidence accumulation models. Unpublished master's thesis, University of Amsterdam. Retrieved from <https://osf.io/mjptz/>.
- Tran, N.-H., van Maanen, L., Heathcote, A., & Matzke, D. (2020). Systematic parameter reviews in cognitive modeling: Towards a robust and cumulative characterization of psychological processes in the diffusion decision model. *Frontiers in Psychology*, 11, 1–14.
- van der Maas, H. L., & Jansen, B. R. (2003). What response times tell of children's behavior on the balance scale task. *Journal of Experimental Child Psychology*, 85(2), 141–177.
- van der Maas, H. L., Molenaar, D., Maris, G., Kievit, R. A., & Borsboom, D. (2011). Cognitive psychology meets psychometric theory: On the relation between process models for decision making and latent variable models for individual differences. *Psychological Review*, 118(2), 339.
- van Maanen, L., Couto, J., & Lebreton, M. (2016). Three boundary conditions for computing the fixed-point property in binary mixture data. *PloS one*, 11(11), e0167377.
- van Zwet, E., & Gelman, A. (2021). A proposal for informative default priors scaled by the standard error of estimates. *The American Statistician*, 0(0), 1–9. <https://doi.org/10.1080/00031305.2021.1938225>.
- Vanpaemel, W. (2011). Constructing informative model priors using hierarchical methods. *Journal of Mathematical Psychology*, 55(1), 106–117. <http://www.sciencedirect.com/science/article/pii/S002249610001069>.
- Veldkamp, K. (2020). Fitting mixtures of Linear Ballistic Accumulation models. University of Amsterdam. https://github.com/Kucharssim/hmm_lba (Unpublished internship report).
- Vigneau, F., Caissie, A. F., & Bors, D.A. (2006). Eye-movement analysis demonstrates strategic influences on intelligence. *Intelligence*, 34(3), 261–272.
- Visser, I. (2011). Seven things to remember about hidden Markov models: A tutorial on Markovian models for time series. *Journal of Mathematical Psychology*, 55(6), 403–415.
- Visser, I., & Poessé, R. (2017). Parameter recovery, bias and standard errors in the linear ballistic accumulator model. *British Journal of Mathematical and Statistical Psychology*, 70(2), 280–296.
- Visser, I., Raijmakers, M. E., & van der Maas, H.L. (2009). Hidden Markov models for individual time series. In *Dynamic process methodology in the social and developmental sciences* (pp. 269–289). Berlin: Springer.
- Wabersich, D., & Vandekerckhove, J. (2014). The rwiener package: an r package providing distribution functions for the wiener diffusion model. *R Journal*, 6(1).
- Wickelgren, W. A. (1977). Speed-accuracy tradeoff and information processing dynamics. *Acta psychologica*, 41(1), 67–85.