

Equal Evidence Perceptual Tasks Suggest a Key Role for Interactive Competition in  
Decision-Making

Ryan P. Kirkpatrick <sup>1</sup>, Brandon M. Turner <sup>2</sup>, Per B. Sederberg <sup>1\*</sup>

<sup>1</sup> Department of Psychology, University of Virginia, USA, <sup>2</sup> Department of Psychology, The  
Ohio State University, USA, \* Corresponding Author, Per B. Sederberg, 434-924-5725  
(phone), pbs5u@virginia.edu (email)

Author Note

Portions of this work were presented in partial fulfillment of a master of arts degree for the first author at the University of Virginia in 2019. The ideas and results discussed here have been presented at conferences by the first author since 2017.

## Abstract

The dynamics of decision-making have been widely studied over the past several decades through the lens of an overarching theory called sequential sampling theory (SST). Within SST, choices are represented as accumulators, each of which races toward a decision boundary by drawing stochastic samples of evidence through time. Although progress has been made in understanding how decisions are made within the SST framework, considerable debate centers on whether the accumulators exhibit dependency during the evidence accumulation process; namely, whether accumulators are independent, fully dependent, or partially dependent. To evaluate which type of dependency is the most plausible representation of human decision-making, we applied a novel twist on two classic perceptual tasks; namely, in addition to the classic paradigm (i.e., the unequal-evidence conditions), we used stimuli that provided different magnitudes of equal-evidence (i.e., the equal-evidence conditions). In equal-evidence conditions, response times systematically decreased with increases in the magnitude of evidence, whereas in unequal-evidence conditions, response times systematically increased as the difference in evidence between the two alternatives decreased. We designed a spectrum of models that ranged from independent accumulation to fully dependent accumulation, while also examining the effects of within-trial and between-trial variability. We then fit the set of models to our two experiments and found that models instantiating the principles of partial dependency provided the best fit to the data. Our results further suggest that mechanisms inducing partial dependency, such as lateral inhibition, are beneficial for understanding complex decision-making dynamics, even when the task is relatively simple.

*Keywords:* perceptual decision-making, sequential sampling models, Bayesian inference, leaky competing accumulator model, response time and accuracy

## Equal Evidence Perceptual Tasks Suggest a Key Role for Interactive Competition in Decision-Making

### **Introduction**

For decades, decision-making researchers have proposed various concepts detailing how the state of the evidence evolves throughout the decision-making process, yet a general consensus concerning the nature of the decision-making process remains elusive (Carland, Thura, & Cisek, 2015; Jones & Dzhafarov, 2014; Ratcliff, 2006; Ratcliff & Smith, 2004; Teodorescu & Usher, 2013). Because decision models play such an important role in enhancing our understanding of individual differences in cognitive dynamics, uncertainty about the general architecture of evidence accumulation has produced a growing tension. We argue that one explanation for our lack of consensus is the sufficiency criterion in the model development process. The field has evolved a set of benchmarks that all models must pass in order to be considered a reasonable description of the decision-making process. For example, a standard paradigm is to present stimuli with varying levels of support for one of multiple (e.g., two) alternatives (Erlick, 1961; Lee & Janke, 1964; Ratcliff, 2006; Ratcliff & Rouder, 1998; Swensson, 1972). Another benchmark involves the data that result from changes in the task instructions, such as emphasizing the speed or accuracy of the decision (Ratcliff & McKoon, 2008; Strayer & Kramer, 1994; Vickers, Burt, Smith, & Brown, 1985; Vickers & Smith, 1989; Vuckovic, Kwantes, Humphreys, & Neal, 2014; Wagenmakers, Ratcliff, Gomez, & McKoon, 2008; Wickelgren, 1977). Somewhat paradoxically, due to our consensus about experimental benchmarks, all serious theoretical contenders have been optimized to pass these benchmarks with ease, ultimately creating a theoretical stalemate until other experimental benchmarks are employed (Brown & Heathcote, 2008; Ratcliff, 1978; Ratcliff & Rouder, 1998; Tsetsos, Usher, & McClelland, 2011; Usher & McClelland, 2001).

In this article, we present a simple, yet nonstandard manipulation nested within the standard benchmarks across two different experiments. Specifically, the tasks include conditions with equal-evidence, variable total sums of evidence within the stimulus, and

varying differences in supporting evidence between the options. This stimulus design affords a variety of conditions and patterns that must be captured simultaneously, and thus provides strong constraints on extant theoretical treatments of how stimulus evidence maps onto decision variables (i.e., choice and response time).

Although the experiments reported below are clearly valuable as a benchmark for subsequent theoretical developments, they are ultimately a means to an end. In this article, our goal is to test and evaluate the relative capabilities of different evidence accumulation architectures. We focus our analyses on three types of dependency that may exist among choice alternatives: fully dependent, partially dependent, and independent. To provide a rigorous evaluation, we designed a set of 12 models, each of which instantiate various forms of dependency among alternatives. The data alone provide strong evidence against fully dependent and independent architectures, leaving only partially dependent architectures as a reasonable explanation for our data. Although we also provide full evaluations of within- and between-trial variability, ultimately these analyses reveal that such variation cannot compensate for the architectural deficits of fully dependent and independent accumulator models.

The outline of this article is as follows. First, we review extant theoretical accounts of evidence accumulation, emphasizing their differences with respect to assumptions about dependency among choice alternatives. Second, we review various experimental paradigms that motivated our specific design. Third, we discuss previous model comparison efforts, noting the distribution of models tested as well as the experimental data used to evaluate the models. Fourth, we provide some “theoretical predictions” from the three classes of model architecture (i.e., fully dependent, independent, and partially dependent) for the types of manipulations used in our experiments. Although these evaluations are inconclusive because they do not explore the full range of possible forms of models within each class, we hope this section orients the reader to the notion that the architecture alone mandates specific predictions for the patterns in behavioral data that can be expected from the experimental

designs we use in this article. Fifth, we describe our two experiments and the pattern of key variables in succession. Sixth, we provide two detailed analyses. In the first analysis, we compare the full class of models by fitting each of them to the two experiments. Evidence for each model is evaluated, and a discussion about consensus is provided. In the second analysis, we further explore the role of within- and between-trial variability in accounting for these data, because combinations of these mechanisms have been provided as explanations for similar patterns of results. The results of our analyses indicate that partially dependent accumulator models, and in particular the Leaky Competing Accumulator model, are the preferred models for fitting to these kinds of data. Subsequently, we close with a discussion concerning previous modeling efforts involving the Leaky Competing Accumulator model and speculate about other models that could potentially fit our data.

### **Extant Theories about Evidence Accumulation**

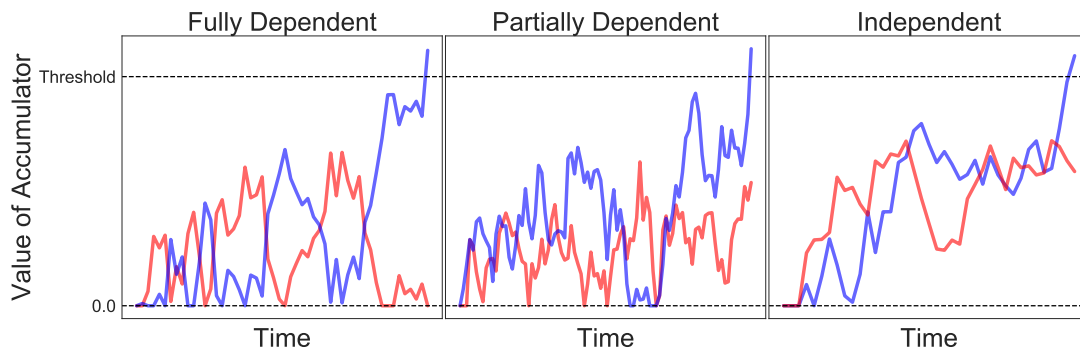
The process of decision-making has been mechanistically specified via the overarching theory of sequential sampling (Busemeyer & Townsend, 1993; Cisek, Puskas, & El-Murr, 2009; Krajbich & Rangel, 2011; Ratcliff, 1978; Usher & McClelland, 2001; Wang, 2002). Models within this framework assume that each choice alternative is represented as an accumulator, and these accumulators race toward a decision threshold by integrating noisy evidence for their corresponding alternatives. The integration of evidence with respect to time allows evidence for a particular alternative to accumulate, where each response alternative will accumulate at a different rate based on assumptions about how the physical stimulus maps onto the psychological perception. Once one of the alternatives reaches a prespecified amount of evidence (i.e., a threshold), a decision is made corresponding to the winning accumulator. The latency between stimulus presentation and an accumulator reaching a threshold is the decision time, but due to elicitation details such as visual encoding and motor response, a nondecision time parameter is often used to (linearly) shift the decision time to an explicit prediction about the response time from an experimental task.

One way to differentiate among many sequential sampling models is by degree of linkage between the accumulators in the model. By this approach, there are three different classes of accumulator models: fully dependent, partially dependent, and independent. Figure 1 illustrates how the evidence accumulation process transpires for each class of model. The model that was used to create the first panel in Figure 1 was the fixed FFI model, the model for the second panel was the LCA model, and the model for the third panel was the Race model with an LCA architecture. These models will be described in more detail in subsequent sections. In fully dependent models, evidence for one option is also evidence against the other option. As illustrated in the left panel of Figure 1, full dependency among alternatives implies perfect anti-correlation in the evidence accumulation process (Ratcliff, 1978, 2006; Ratcliff & McKoon, 2008; Ratcliff & Rouder, 1998; Turner, Sederberg, & McClelland, 2016).

In partially dependent accumulator models, each option is represented by a separate accumulator, but, at each timestep, the value of each accumulator is affected by the input into the other accumulators or the value of the other accumulators (Shadlen & Newsome, 2001; Usher & McClelland, 2001). Typically, in these models, stronger evidence in support of one option creates greater suppression of the evidence supporting the other options. As illustrated in the middle panel of Figure 1, partially dependent models assume each piece of evidence affects every option to some extent, but the drive to each accumulator is not perfectly anticorrelated as it is in fully dependent accumulator models. The inhibition within the evidence accumulation process was inspired by, and is often likened to, global and local inhibitory dynamics in the brain (Bogacz, Brown, Moehlis, Holmes, & Cohen, 2006; Shadlen & Newsome, 2001; Usher & McClelland, 2001; van Ravenzwaaij, van der Maas, & Wagenmakers, 2012; Wang, 2002).

Independent accumulator models accumulate information pertaining to the corresponding alternative with zero consultation of other accumulation occurring in the decision process (Brown & Heathcote, 2008; LaBerge, 1962; Merkle & Van Zandt, 2006; Pike,

1971; Rouder, Province, Morey, Gomez, & Heathcote, 2015; Vickers, 1970). Importantly, as illustrated in the right panel of Figure 1, the evidence accumulation process of one accumulator does not explicitly depend on the state of the other accumulator.



*Figure 1. Illustration of the evidence accumulation process for each class of accumulator model.* Each subfigure illustrates how each accumulator inhibits or does not inhibit the other accumulator during the evidence accumulation process for a two-choice decision for each class of accumulator model. The blue line represents the accumulator that eventually crosses the threshold while the red line represents the accumulator that fails to cross the threshold before the blue accumulator. Consistent with the models in subsequent sections, in this illustration, a lower bound is in place for all models to prevent the accumulator values from dropping below zero. The upper dotted black line represents the threshold. See the Model Analysis 1 section for details of the models used to generate these example simulations.

In general, sequential sampling models have successfully captured a variety of decision-making data, which includes experiments presenting participants with difficult decisions where the differences in quality between the competing options are small (potentially the most constraining for the models). Some studies have used random dot kinematograms where the direction of coherent dot movement switches mid-trial (Holmes, Trueblood, & Heathcote, 2016; Tsetsos, Gao, McClelland, & Usher, 2012; Winkel, Keuken, van Maanen, Wagenmakers, & Forstmann, 2014). Other studies have included conditions in their random dot motion (RDM) task where the motion coherence is equal to zero (no

coherent movement in any direction) (Heekeren, Marrett, Ruff, Bandettini, & Ungerleider, 2006; Palmer, Huk, & Shadlen, 2005).

Studies have also examined perceptual decision tasks with multiple alternatives (Krajbich & Rangel, 2011; Niwa & Ditterich, 2008; Tsetsos et al., 2011; Usher & McClelland, 2004). Several recent studies have examined decisions between choices with equal-evidence supporting each choice and multiple levels of total evidence (Krajbich, Armel, & Rangel, 2010; Pais et al., 2013; Pirrone, Azab, Hayden, Stafford, & Marshall, 2017; Smith & Krajbich, 2018). In the present study, we fit multiple representatives of each class of accumulator model to similar decisions as many of those listed above to determine which model is best able to capture all of these data patterns simultaneously.

Some previous studies have compared the three classes of accumulator models directly, but with conflicting results. Ratcliff and Smith (2004) fit candidates from all three classes of accumulator models to signal detection and lexical decision data. These authors found that the fully dependent candidate fit both the signal detection and lexical decision data better than the independent model and both the fully dependent and independent candidates fit the lexical decision data better than the partially dependent model. However, Teodorescu and Usher (2013) found the partially dependent model fit a fluctuating brightness discrimination task better than the independent model, and Teodorescu et al. (2015) found the partially dependent model fit a different fluctuating brightness discrimination task better than the fully dependent model. As we outline below, we believe our work resolves some of this conflict concerning which of the extant models provides the best account of perceptual decision-making.

Although previous studies have examined the fit of sequential sampling models to many of the patterns observed in our data, there is no study to our knowledge that has attempted to fit equal-evidence conditions, zero evidence conditions, and unequal-evidence conditions simultaneously. Thus, we contribute a 2-choice RDM task where we manipulate the proportion of coherent motion in the left direction and the proportion of coherent motion



in the right direction within the same trial. Throughout the article, we will refer to the proportion of coherent dot motion as simply “coherence”. To ensure that our results were not related to perceptual effects that can contaminate responses made during RDM tasks (Anstis, 1980; Pilly & Seitz, 2009), we also created a different perceptual task where participants perform contrast judgments on grating stimuli. Both tasks have four equal-evidence conditions with varying amounts of total evidence and 12 unequal-evidence conditions with differing amounts of disparity in the evidence supporting the choices. We hypothesized that, by including multiple equal-evidence conditions in our task, the changes in response time at differing levels of evidence would be challenging for some accumulator models to fit simultaneously with the unequal conditions. In the following section, we provide theoretical predictions of the models to identify which qualitative patterns of results from our task could potentially be challenging for a given theory of decision-making.

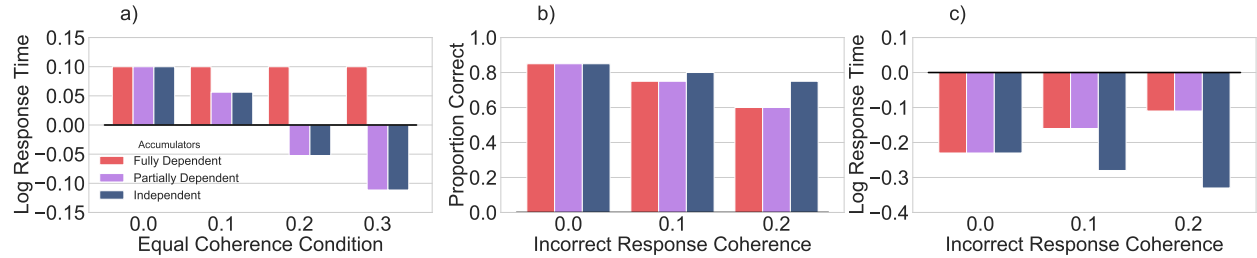
### **Theoretical predictions**

Before discussing our theoretical predictions, we explain here why our response distributions throughout the article are logarithmically transformed and then summarized by calculating the mean of the distribution. Although we appreciate the rich history emphasizing the shape of the response time distribution in discriminating among models, and, in fact, all our model fits are based on the full set of responses, in the subsequent figures we report only the mean of the log-transformed response times because these statistics alone are enough to discriminate among the models, given the design of our experiments. In general, response time distributions have long tails that can shift the calculation of the mean rightward even if the shift is only based on a few slow responses rather than the slower responses making up a significant density of the distribution. Transforming the response time distribution using the logarithmic function before calculating the mean normalizes the distribution and shifts the tail back toward the highest density region. The logarithmic transformation subsequently makes comparisons of means to be

unbiased by the slow responses occupying the tail of the distribution.

We illustrate in Figure 2 the qualitative predictions of the three classes of models for the patterns of response times and accuracies observed in our RDM task. To create this theoretical predictions figure, we assumed that all accumulator architectures made the same prediction at the baseline 0.0 condition. Our task induces a wide array of model predictions because it includes conditions with equal-evidence supporting each option at various levels of coherence and is designed such that comparisons can be made between conditions at different total levels of evidence. Although all models predict chance accuracy for each equal-coherence condition (not shown), the models make divergent predictions with regard to response time. As illustrated in Figure 2a, a fully dependent accumulator architecture would make the same response time prediction regardless of the coherence level in the equal-coherence conditions. This is because evidence in each direction is perfectly anti-correlated, such that evidence supporting one option is perfect evidence against the other option. In this case, the evidence does not drive the accumulation process because there is no evidence supporting either option over the alternative, and the decision boundary is only reached due to noise in the model architecture (Turner, Gao, Koenig, Palfy, & L. McClelland, 2017). If we assume the noise is fixed across conditions, the fully dependent architecture will predict the same response times for each equal-coherence condition. To address this limitation, some researchers have proposed mechanisms that adjust the overall level of noise based on the magnitude of the inputs (Ratcliff, Voskuilen, & Teodorescu, 2018; Teodorescu et al., 2015). In Model Analysis 2, we investigate the utility of this approach with several models that examine the effect of both between-trial and within-trial variability on the model fits.

As shown in Figure 2a for the equal-coherence conditions, independent and partially dependent accumulator theories make a prediction in the equal-coherence conditions that is distinct from the fully dependent models. Independent accumulator theories and partially dependent accumulator theories would predict faster response times as the coherence in one



*Figure 2. Theoretical predictions of each accumulator architecture.* Each subfigure illustrates the qualitative theoretical predictions of each accumulator architecture in a subset of trials in an RDM task where the coherence is manipulated in both directions in the same trial. In all subfigures, we assume the accumulator architectures make equivalent predictions at the baseline 0.0 coherence condition. Thus, the predictions are meant to illustrate the model behavior as a function of coherence level, not mean performance level when averaged across coherence. a) Log response time predictions of each accumulator architecture in the equal-coherence conditions. b) Proportion correct predictions of each accumulator architecture in conditions where the coherence in exactly one direction is equal to 0.3. c) Log response time predictions of each accumulator architecture in conditions where the coherence in exactly one direction is equal to 0.3. These response times are associated with a correct response.

direction increases because the evidence becomes stronger as the coherence increases (Teodorescu & Usher, 2013). In an independent accumulator architecture, the strength of one accumulator does not directly affect the other accumulator, so the response times will become faster as the evidence increases. In a partially dependent accumulator architecture, the added drive in the system is reflected in the rate of accumulation that, unlike in the fully dependent accumulator theories, is not cancelled out because each accumulator only partially inhibits the other accumulators.

Figure 2b shows the predicted proportion correct for each of the three classes of accumulator theories for conditions where one coherence is equal to 0.3 and the other coherence is equal to some value less than 0.3. As the difference in coherence between each

direction becomes smaller, fully dependent and partially dependent theories predict less accurate responses due to stronger inhibition by the accumulator representing the incorrect response on the accumulator representing the correct response. Independent accumulator theories also predict less accurate responses as the difference in coherence becomes smaller, however the severity of this predicted decrease in accuracy will be smaller than that predicted by both classes of dependent accumulator theories. Because the accumulators in these theories are independent, there is no inhibition preventing the accumulator representing the correct response from crossing the threshold the majority of the time. In these theories, the decrease in proportion correct can be attributed to the accumulator representing the incorrect response crossing the threshold more often at higher levels of evidence.

In Figure 2c, predicted response times for correct responses are illustrated for conditions where one coherence is equal to 0.3 and the other coherence is equal to some value less than 0.3. Both partially and fully dependent accumulator theories predict slower response times as the disparity between the coherence in the two directions decreases (i.e., as the decision becomes more difficult). For example, if the coherence in the left direction is 0.0 and is 0.3 in the right direction, there is less inhibition on the correct choice (0.3) than there is when the coherence in the left direction is 0.2. Less inhibition allows the accumulator representing the correct response to cross the decision threshold faster than that accumulator could if it faced more inhibition.

Independent accumulator theories predict equal or faster response times as the disparity between the coherence in the two directions decreases. For independent accumulator theories, the predicted response time distribution for a specific choice involves a calculation of the probability of making that choice, times the probability that one has not made the alternative choice at each time point  $t$ . To see how these predictions play out in our example, suppose we are trying to compute the probability of choosing the rightward response when the coherence for right is 0.3 and the coherence for the left is 0.1. In this situation, because the difference in coherence between the two options is relatively large, it will be highly likely

that the 0.3 choice will be made relative to the 0.1 condition. Essentially, because there is little chance that the leftward accumulator will win, it becomes inconsequential in predicting the response time distributions regardless of choice. However, if the evidence for the leftward response increases to say 0.2, the probability of making a leftward choice has substantially increased. Now when making predictions for the probability of any response at time  $t$ , we must seriously consider the probability that a leftward choice will be made. For small  $t$ , the probability of making a response at that time increases, and this increase for small  $t$  decreases the probability of making a response at larger  $t$  because it becomes generally more likely that, if a response is made, the response is made at an earlier time point. The result is that although a comparison between say 0.3 and 0.1 is more accurate, it may actually be slower than a comparison between 0.3 and 0.2. This phenomenon is known as statistical facilitation (Luce, 1986; Raab, 1962; Townsend & Nozawa, 1995), where, when one coherence is 0.3, the accumulator representing the direction with a coherence of 0.2 crosses the threshold more often than the accumulator representing the direction with a coherence of 0.1. This process removes winning trials from the 0.3 coherence, resulting in a slower response time distribution for the correct response of 0.3 when the other coherence is 0.2 than when the other coherence is 0.1. It is important to note that, although we are predicting the direction of effect, we are making no claim with regard to whether that effect will be significant, which depends on other factors, such as the number of trials and participants.

We test these theoretical predictions in our analyses of the data from our two experiments. Our first analysis examines the accuracy and response time patterns observed in our experiments. The second analysis examines the fits of models representing each of our three classes of accumulator dependency (independent, partially dependent, and fully dependent) to our two data sets. The representative models we selected are: the partially dependent Feed-Forward Inhibition model (FFI; Shadlen & Newsome, 2001), the partially dependent Leaky Competing Accumulator model (LCA; Usher & McClelland, 2001), and the independent Linear Ballistic Accumulator model (LBA; Brown & Heathcote, 2008) as well as

two independent accumulator variants of the LCA model and one additional fully dependent implementation of the FFI model. The fully dependent implementation of the FFI model has a similar accumulation process to the popular Diffusion Decision Model (DDM; Ratcliff, 1978, 2006; Ratcliff & McKoon, 2008; Ratcliff & Rouder, 1998; Turner et al., 2016). In our third analysis, we examine the impact of additional sources of variability on the model fits to both the RDM and the grating data.

## Experiment 1

Our first experiment used an RDM task to examine the effects of unequal-coherence and equal-coherence stimuli on decision-making variables. To the best of our knowledge, performance in an RDM task involving both types of coherence conditions has never been examined. Although various combinations of each type of coherence have been investigated, as we will show below, together they provide a unique and powerful test of extant theories of evidence accumulation.

### Participants

16 undergraduate students attending the Ohio State University participated in the study as a requirement for an introductory psychology course. 10 of the participants were male, and the participants averaged 18.88 years of age ( $SD=1.218$ ). The study protocol was approved by the Institutional Review Board for Human Subjects at the Ohio State University.

### Stimuli and apparatus

The experiment was written and displayed using the State Machine Interface Library for Experiments (<https://github.com/compmem/smile>). The stimuli were generated on Debian Linux operating systems with Nvidia graphics cards. The stimulus was composed of 100 dots contained in a circle with a 200 pixel-length radius. The dots were 3 pixels by 3 pixels. The lifespan for each dot was randomly chosen from between 0.25 and 1.25 seconds.

Each dot would appear at a random location within the stimulus window and moved in its predetermined direction until its lifespan expired or it left the radius of the circle, at which point a new dot would be generated.

## **Procedure**

Participants were instructed to quickly and accurately choose the direction (left or right) with the most coherent dot movement. Both the proportion of coherent dot movement to the left and to the right were manipulated in the same stimulus, while the remaining dots moved in random directions. The proportion of dot coherence in either direction could be either 0.0, 0.1, 0.2, or 0.3. Thus, as shown in Figure 3b, there were 16 different experimental conditions the participants could encounter, including four different conditions where the proportion of dot coherence in the left and right directions was equivalent.

Before stimulus presentation, a fixation cross was presented on screen for a random duration between 0.75 and 1.25 seconds. In each trial, the stimulus remained on the screen until the participant indicated which direction had more coherent dot movement by pressing either the “D” key for the left direction or the “K” key for the right direction. Immediately following the response, feedback was presented in the center of the computer screen. The feedback was presented for 1 second and consisted of a green check mark symbol for a correct response, a red “X” symbol for an incorrect response, or the expression “Too Fast!” if the participant answered before 100 milliseconds had elapsed after stimulus onset. In the equal-coherence trials, the green checkmark and the red “X” symbols were presented randomly, such that “D” was the correct response in exactly half of these trials. Each participant completed 8 blocks with 60 trials each for a total of 480 trials.

## **Behavioral Results**

Responses faster than 0.2 seconds and slower than 5 seconds were excluded from data analysis (<3.2% of all data). In all figures, error bars represent Loftus and Masson corrected 95% confidence intervals (Loftus & Masson, 1994). Figure 4a shows how the coherence

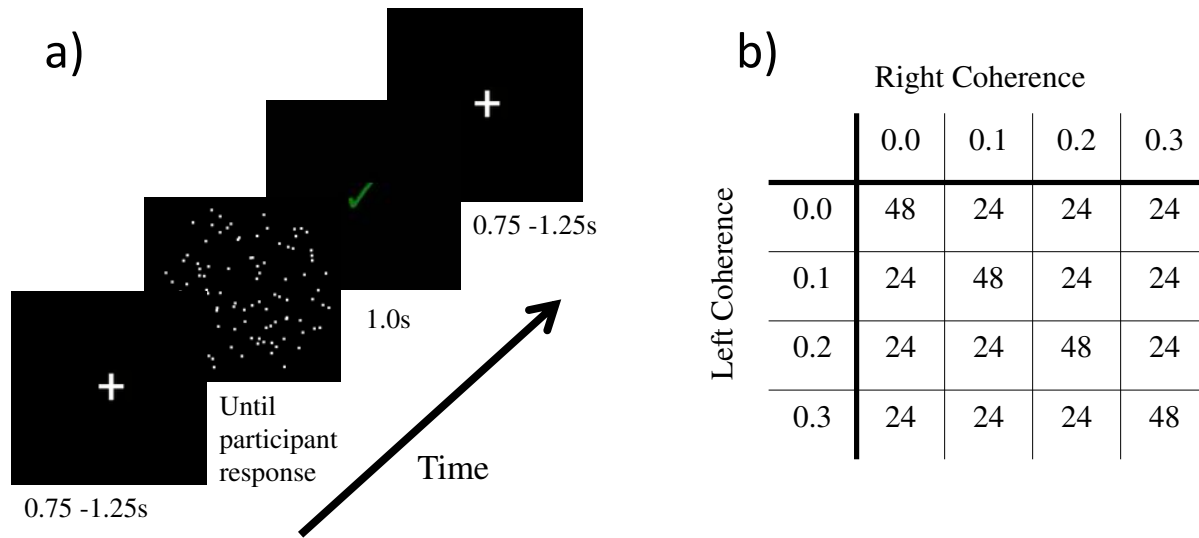


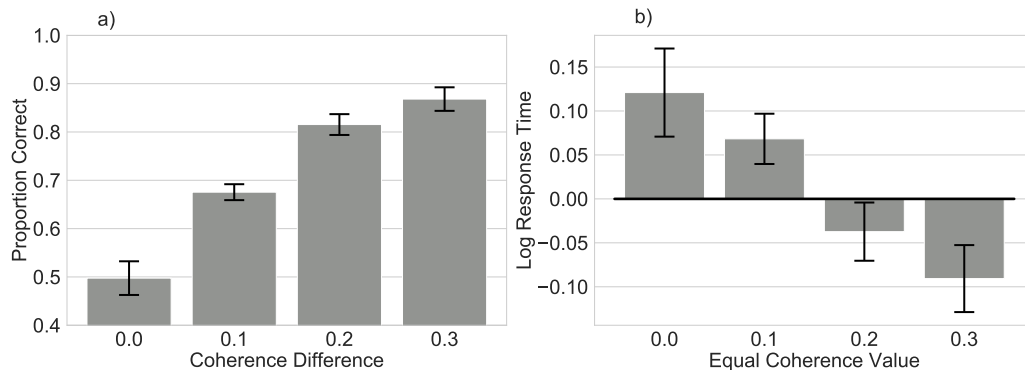
Figure 3. Illustration of the experiment and table of the experimental conditions. a)

Depiction of one trial of the experiment with the length of presentation below each image. b)

Number of trials for each experimental condition. The main diagonal shows the number of trials for each of the equal-coherence conditions.



manipulation affects accuracy. We tested this finding with a mixed effects linear regression model that predicted accuracy from coherence difference with a random intercept for each participant. As the difference in coherence between the left and right directions increased, participants responded with greater accuracy,  $B=1.296$ , 95% CI [1.132, 1.458],  $p<0.001$ . Participant accuracy was at chance overall when all of the equal-coherence data were combined.



*Figure 4. Proportion correct as a function of coherence difference and response times in each equal-coherence condition. a) All of the participant data were organized by difference in coherence between the left and right directions. First, proportion correct was calculated within participant, and then the mean was calculated between participants. b) From only the four equal-coherence conditions, response times across all participants were log transformed, then the mean was calculated between participants. In both subfigures, the error bars represent Loftus and Masson (1994) corrected 95% confidence intervals.*

Figure 4b shows average log response time as a function of equal-coherence condition. As the coherence in both directions increased, the average log response time decreased. We tested this finding with a mixed effects linear regression model that predicted log response time from equal-coherence condition with a random intercept for each participant. We found response time decreased as the coherence in both directions increased,  $B=-0.763$ , 95% CI [-0.955, -0.571],  $p<0.001$ . To rephrase this result, participants responded more quickly on average when they had more evidence to inform their decision, even though the evidence

supporting each direction was equivalent. This pattern of results supports the predictions of independent and partially dependent architectures, but contradicts the predictions of fully dependent architectures with noise that is value-independent.

Figure 5 shows proportion correct and response time across all participants as a function of incorrect coherence condition. Here, proportion correct and response time are shown for only the conditions where exactly one direction had a coherence of 0.3. We chose to highlight these three unequal-coherence conditions to illustrate how proportion correct and response times change as the difference in evidence supporting each option becomes smaller. Figure 5a shows proportion correct across all participants as a function of incorrect coherence condition. There is a clear decrease in participant accuracy as the difference in coherence between the two directions becomes smaller. We confirmed this result with a mixed effects linear regression model that predicted accuracy from coherence condition where exactly one coherence was 0.3. The model indicated that accuracy decreased as the difference between the incorrect coherence and the coherence of 0.3 decreased,  $B=-1.078$ , 95% CI  $[-1.330, -0.826]$ ,  $p<0.001$ . The relatively pronounced difference in proportion correct between the 0.0 and 0.2 conditions provides support for the fully dependent and partially dependent accumulator theories more than it supports the independent accumulator theories.

Figure 5b shows average log response time for the correct response across all participants as a function of incorrect coherence condition. As the difference between the incorrect coherence and the coherence of 0.3 decreased, average log response time increased. We confirmed the validity of this result with a mixed effects linear regression model that predicted log response time from coherence condition where exactly one coherence was 0.3. The model indicated that response time increased as the difference between the incorrect coherence and the coherence of 0.3 decreased,  $B=0.613$ , 95% CI  $[0.352, 0.874]$ ,  $p<0.001$ . This pattern of results supports the predictions of both partially dependent and fully dependent accumulator theories, but opposes the predictions of independent accumulator theories.

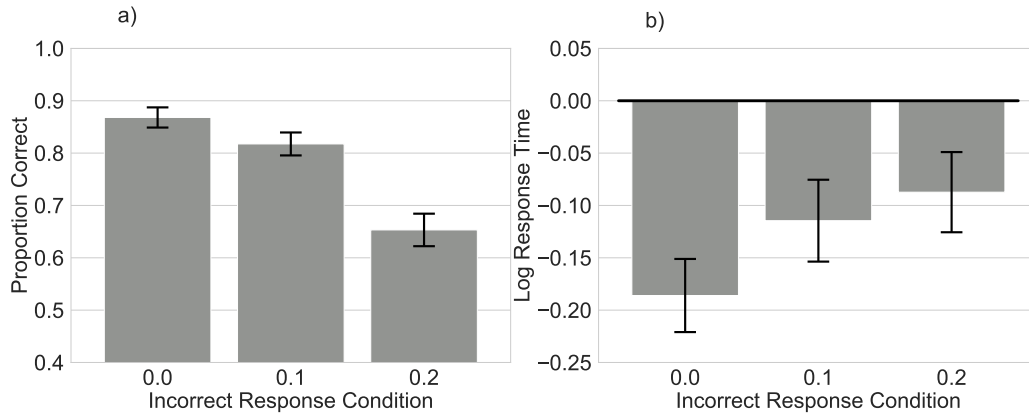


Figure 5. Response times and proportion correct in conditions where exactly one coherence equals 0.3. a) From only the conditions where exactly one direction had a coherence of 0.3, proportion correct was calculated between participants. b) From only the conditions where exactly one direction had a coherence of 0.3, response times associated with the correct response across all participants were log transformed. Then the mean was calculated between participants. In both subfigures, the error bars represent Loftus and Masson (1994) corrected 95% confidence intervals.

## Discussion

We have shown how a simple perceptual decision-making task can potentially distinguish between different theories of decision-making. In equal-coherence conditions, participants responded more quickly as the total level of evidence increased, a finding that contradicts the predictions of fully dependent accumulator theories. Fully dependent theories predict equivalent response times regardless of the total level of evidence. Independent and partially dependent theories correctly predict the pattern of response times observed in the equal-coherence conditions. However, independent theories predict slightly faster response times in unequal-coherence conditions as the difference between the two options is reduced, whereas we observed significantly slower response times in our results. Only partially-dependent accumulator theories correctly predict all the behavioral patterns. To ensure that the results in our RDM task were not caused by perceptual effects that can

affect responses in these kinds of tasks (Anstis, 1980; Pilly & Seitz, 2009), we also collected data in a similar task that used grating stimuli instead of RDM stimuli.

## Experiment 2

In this study, we set out to replicate our findings Experiment 1 using similar conditions, but with a different perceptual decision-making task. We decided this replication was necessary because results from some RDM tasks are affected by perceptual contaminants outside of the processes of interest (Anstis, 1980; Pilly & Seitz, 2009). These studies discuss how motion in an RDM task can be incorrectly classified by our visual systems, resulting in evidence supporting an option that actually should have no support. In this study, we presented our participants with two grating stimuli on each screen. Each grating stimulus had either the same or a different contrast than its counterpart which allowed us to create a similar set of manipulations as in Experiment 1.

### Participants

23 undergraduate students attending the University of Virginia participated in the study as a requirement for the Psychology program. 14 of the participants were female, and the participants averaged 18.696 years of age ( $SD=0.748$ ). The study protocol was approved by the Institutional Review Board for Social and Behavioral Research at the University of Virginia.

### Stimuli and apparatus

The experiment was written and displayed using the State Machine Interface Library for Experiments (<https://github.com/compmem/smile>). The stimuli were generated on Windows operating systems with Nvidia graphics cards. Each stimulus was composed of two sinusoidal gratings separated by 60 pixels with a separate, Gaussian envelope obscuring each grating. The key parameter manipulated in each stimulus was how similar the stimulus is to the experiment background, a parameter we will call the “contrast” for the remainder of this

article. As we have defined contrast, a grating with a larger contrast is easier to discriminate from the experiment background and a grating with a smaller contrast is more difficult to discriminate from the experiment background. Each grating had a 150 pixel-length radius with a 180 degree orientation. The phase shift of the sine wave controlling the grating was 0 cycles and the frequency of the sine wave was 20 cycles per pixel. The standard deviation of the Gaussian envelope was 7.5 pixels.

### **Procedure**

Participants were instructed to quickly and accurately choose the most clear grating (left or right). The contrast of each grating could be either 0.40, 0.43, 0.46, or 0.49. Thus, as shown in Figure 6b, there were 16 different experimental conditions the participants could encounter, including four different conditions where the contrast of the left grating was equal to the contrast of the right grating.

Before stimulus presentation, a fixation cross was presented on screen for a random length between 0.75 and 1.25 seconds. In each trial, the stimulus remained on the screen until the participant indicated which grating stimulus was the most clear by pressing either the “D” key for the left grating or the “K” key for the right grating. Immediately following the response, feedback was presented in the center of the computer screen. The feedback was presented for 1 second and consisted of a green check mark symbol for a correct response, a red “X” symbol for an incorrect response, or the expression “Too Fast!” if the participant answered before 100 milliseconds had elapsed after stimulus onset. In the equal-contrast trials, the green checkmark and the red “X” symbols were presented randomly, such that “D” was the correct response in exactly half of these trials. Each participant completed 8 blocks with 60 trials each for a total of 480 trials.

### **Results**

Responses faster than 0.2 seconds and slower than 5 seconds were excluded from data analysis (<2.6% of all data). Figure 7a shows how the contrast manipulation affects

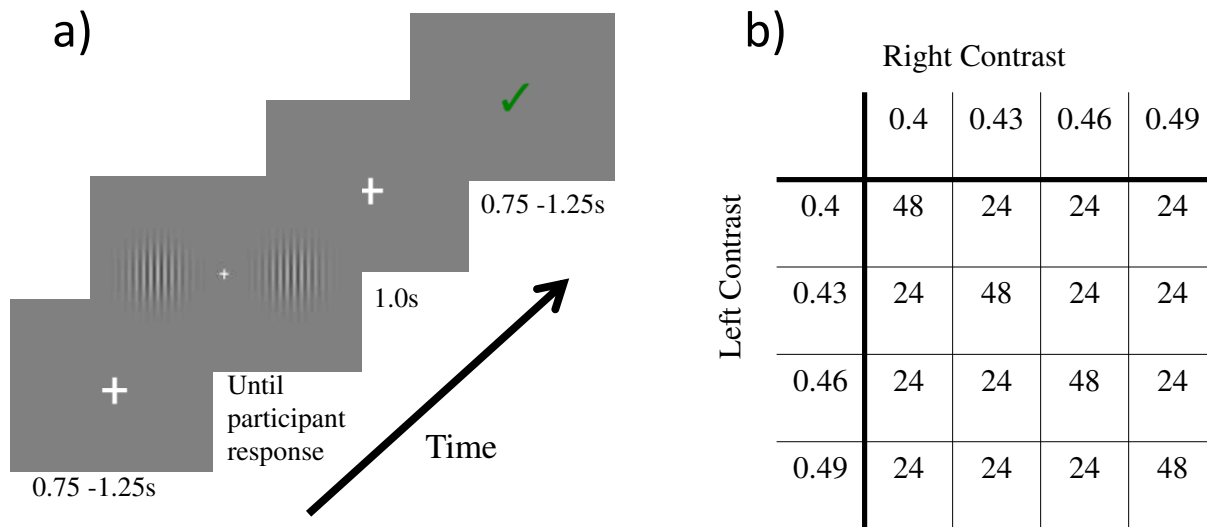


Figure 6. Illustration of the grating experiment and table of the experimental conditions. a) Depiction of one trial of the experiment with the length of presentation below each image. b) Number of trials for each experimental condition. The main diagonal shows the number of trials for each of the equal-contrast conditions.

accuracy. We tested this finding with a mixed effects linear regression model that predicted accuracy from contrast difference with a random intercept for each participant. Consistent with Experiment 1, as the difference in contrast between the left and right gratings increased, participants responded with greater accuracy,  $B=4.731$ , 95% CI [4.272, 5.191],  $p<0.001$ . Participant accuracy was at chance overall when all of the equal-contrast data were combined.

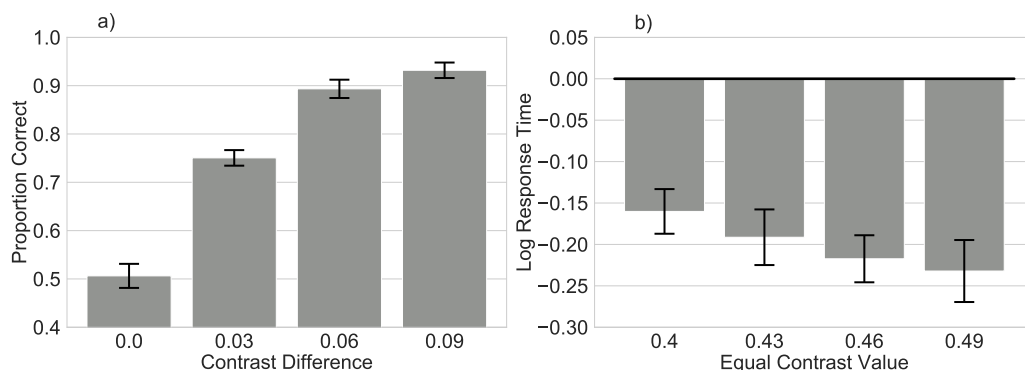


Figure 7. Proportion correct as a function of contrast difference and response times in each equal-contrast condition. a) All of the participant data were organized by difference in contrast between the left and right direction. First, proportion correct was calculated within participant, and then the mean was calculated between participants. b) From only the four equal-contrast conditions, response times across all participants were log transformed. Then the mean was calculated between participants. In both subfigures, the error bars represent Loftus and Masson (1994) corrected 95% confidence intervals.

Figure 7b shows average log response time as a function of equal-contrast condition. Consistent with Experiment 1, as the contrast between each grating stimulus and the background increased, the average log response time decreased. We tested this finding with a mixed effects linear regression model that predicted log response time from the equal-contrast condition with a random intercept for each participant. We found response time decreased as the contrast between each grating stimulus and the background increased,  $B=-0.825$ , 95% CI [-1.277, -0.372],  $p<0.001$ .

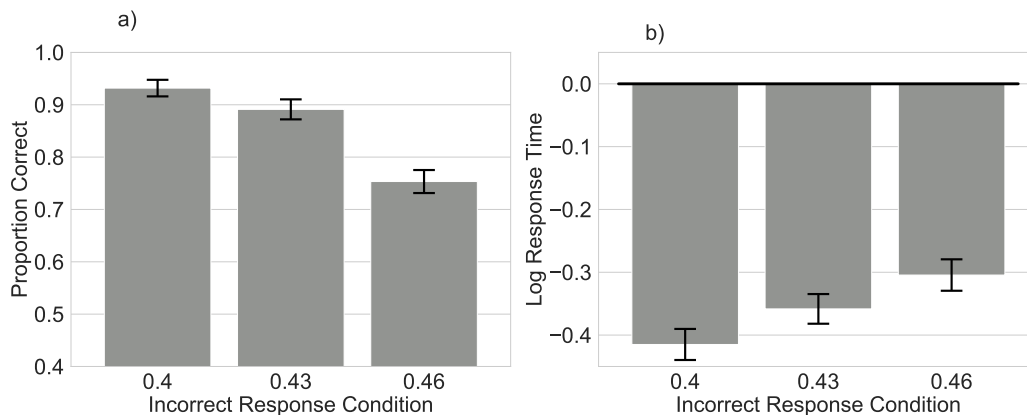


Figure 8. Response times and proportion correct in conditions where exactly one contrast equals 0.49. a) From only the conditions where exactly one grating had the highest contrast of 0.49, proportion correct was calculated between participants. b) From only the conditions where exactly one grating had the highest contrast of 0.49, response times associated with the correct response across all participants were log transformed. Then the mean was calculated between participants. In both subfigures, the error bars represent Loftus and Masson (1994) corrected 95% confidence intervals.

Figure 8a shows the observed proportion correct across all participants as a function of incorrect contrast condition. There is a clear decrease in participant accuracy as the difference in contrast between the two stimuli becomes smaller which is consistent with the results from Experiment 1. We confirmed this result with a mixed effects linear regression model that predicted accuracy from contrast condition where exactly one contrast was the highest contrast of 0.49. The model indicated that accuracy decreased as the difference between the incorrect contrast and the highest contrast decreased,  $B=-2.975$ , 95% CI  $[-3.515, -2.435]$ ,  $p<0.001$ .

Figure 8b shows average log response time for the correct response across all participants as a function of incorrect contrast condition. As the difference between the incorrect contrast and the highest contrast decreased, average log response time increased, which is consistent with Experiment 1. We confirmed the presence of this pattern with a



mixed effects linear regression model that predicted log response time from contrast condition where exactly one contrast was the highest contrast of 0.49. The model indicated that response time increased as the difference between the incorrect contrast and the highest contrast decreased,  $B=1.842$ , 95% CI [1.246, 2.438],  $p<0.001$ .

## Discussion

Experiment 2 replicated the key findings of Experiment 1, providing evidence that the results from Experiment 1 cannot be fully explained by perceptual effects that could have contaminated the RDM stimuli. In Experiment 2, as in Experiment 1, we observed faster response times in the equal-contrast conditions as the contrast between the stimuli and the background increased. In the unequal-contrast conditions, we observed slower response times and reduced accuracy as the difference in contrast between the two gratings decreased which was consistent with Experiment 1. Taken together, our results from Experiments 1 and 2 indicate this task could provide sufficient constraint to discriminate between the classes of accumulator theories, with the greatest support for the partially-dependent accumulator theories, provided those classes do generally predict the patterns of behavior we illustrated in the introduction. However, as we noted above, these theoretical predictions are based on the general architecture of the models and with a single parameter setting; that is, they do not address the many types of predictions the models could produce with different parameter settings, known as model flexibility. To examine whether any of the models are flexible enough to still capture patterns in our data, in the following sections, we fit each of the models using simulation-based Bayesian methods, which are known to balance model fit with model flexibility (Lee, 2008; Palestro, Sederberg, Osth, Van Zandt, & Turner, 2018; Turner, Sederberg, Brown, & Steyvers, 2013; Turner et al., 2016; B. M. Turner & Van Zandt, 2018).

### Model Analysis 1

Although our behavioral analyses provide strong evidence in support of partially dependent accumulator models, the prior analyses did not provide a thorough quantitative

assessment of specific models. Our prior analyses also did not explore the influence of variability in theories of decision-making, so we explored these influences in Model Analysis 2. Thus, we fit several well-studied models that implement either independent, fully dependent, or partially dependent accumulators in the evidence accumulation process to the data from Experiments 1 and 2. We fit each of the FFI model (Shadlen & Newsome, 2001), the LCA model (Usher & McClelland, 2001), and the LBA model (Brown & Heathcote, 2008) to the full choice-response time distributions for each coherence condition using a simulation-based, Bayesian approach. In addition, we fit an additional fully dependent variant of the FFI model and two additional independent variants of the LCA model to the data, giving rise to 6 total models spanning the three classes of accumulator models: independent, partially dependent, and fully dependent.

### The Feed-Forward Inhibition (FFI) Model

We chose the FFI model because it has a parameter that can be fixed to create a fully dependent model that resembles the DDM (Turner et al., 2016) or can be left free to create a partially dependent model. Because only one parameter is changed in the model between the two variants, we can directly compare how the fully dependent and partially dependent versions capture the essential trends in our data.

The FFI model assumes the inhibition on each accumulator is based on the average stimulus input to the other alternatives, such that

$$dx_i = \left(\rho_i - \frac{\nu}{C-1} \sum_{j \neq i} \rho_j\right) \frac{dt}{\tau} + \xi \sqrt{\frac{dt}{\tau}}, \quad (1)$$

$$\xi \sim N(0, \eta), \quad (2)$$

$$x_i \rightarrow \max(x_i, 0), \quad (3)$$

where  $\nu$  is the FFI parameter,  $\rho_i$  represents the rate of evidence accumulation for the  $i$ th alternative,  $dx_i$  represents the change in the value of the accumulator  $x_i$  at each timestep,  $\xi \sim N(0, \eta)$  represents the within-trial variability ( $\xi$  is recalculated at each timestep), and  $C$

represents the number of choice alternatives. We fixed  $dt = 0.01$ ,  $\tau = 0.1$ , and  $\eta = 1$ . Each accumulator  $x_i$  is initialized to 0. At each timestep, each accumulator value changes until one accumulator has a value greater than or equal to the threshold  $\alpha$ . At which point, the value of the threshold-surpassing accumulator plus non-decision time  $t_0$  (a parameter representing perceptual and motor response) represents the response time for choice  $c$  made by the participant. Note, during this accumulation process, a lower boundary (Bogacz, Usher, Zhang, & McClelland, 2007; Diederich, 1995) is in place such that if the accumulator becomes negative, it is reset to zero.

In fitting the FFI model to the data, for each of the four  $\rho_i$  parameters, we specified a prior of a truncated normal distribution with a mean of 2.5, standard deviation of 5, lower bound of 0, and upper bound of 10. For  $\nu$ , we specified a prior of a normal distribution with a mean of 0 and standard deviation of 1.4 altered by an inverse logit transform. For  $\alpha$ , the prior specified was a truncated normal distribution with a mean of 2.5, standard deviation of 10, lower bound of 0, and upper bound of 30. We specified a prior for  $t_0$  of a uniform distribution with a lower bound of 0 and upper bound of the minimum observed response (unique for each participant).

We also fit a variant of the FFI model to the data where we fixed  $\nu$  to 1 (Turner et al., 2016). For clarity, we will call this variant of the model the fixed FFI model and the variant with  $\nu$  free the free FFI model. By our definitions, the fixed FFI model is a fully dependent model and the free FFI model is a partially dependent model. In the two-alternative case, the fixed FFI model behaves similarly to the DDM in that the evidence accumulation process is completely anticorrelated. The only discrepancy between the fixed FFI and the DDM occurs when one accumulator is pushed downward to zero evidence. In the DDM, the losing accumulator would continue to decrease (i.e., become negative), whereas in the fixed FFI, the losing accumulator would just be perpetually reset to zero. This is a subtle difference that does not affect the model predictions given that losing accumulators tend to continue losing in our paradigm, but it is worth noting that the fixed FFI is not a perfect

analogue to the DDM. We specified the same priors for the fixed FFI model as the free FFI with the exception of  $\nu$ .

### The Leaky Competing Accumulator (LCA) Model

We included the LCA model in our study because it makes a different prediction from the FFI model about how partial inhibition during the evidence accumulation process is implemented. Whereas the FFI model explains that inhibition occurs via input competition, the LCA model explains that inhibition occurs laterally or based on the total value of the accumulator at each timestep. By comparing the fit of each model to the same dataset, we can determine which mechanism provides a more satisfactory explanation of the observed patterns. Because the leak-only LCA model is simply the LCA model without lateral inhibition, we will be able to further assess the importance of lateral inhibition to the model by comparing the fits of the two models. Furthermore, we can evaluate the importance of the passive decay of evidence by comparing the leak-only LCA model and the race LCA model, because that parameter is the only difference between the two models.

The LCA model is a partially dependent accumulator model where the stimulus input to one accumulator does not directly inhibit the values of the other accumulators. Instead, at each timestep, the total value of the other accumulators inhibits each accumulator. The model accounts for the passive decay of old information with a leakage parameter. The following differential equation describes the accumulation process in the LCA model:

$$dx_i = (\rho_i - \kappa x_i - \beta \sum_{j \neq i} x_j) \frac{dt}{\tau} + \xi \sqrt{\frac{dt}{\tau}}, \quad (4)$$

$$\xi \sim N(0, \eta), \quad (5)$$

$$x_i \rightarrow \max(x_i, 0), \quad (6)$$

where  $\kappa$  is the leakage parameter,  $\beta$  is the lateral inhibition parameter,  $\rho_i$  represents the rate of evidence accumulation for the  $i$ th alternative,  $dx_i$  represents the change in the value of accumulator  $x_i$  at each timestep, and  $\xi \sim N(0, \eta)$  represents the within-trial variability ( $\xi$  is

recalculated at each timestep). We again fixed  $dt = 0.01$ ,  $\tau = 0.1$ , and  $\eta = 1$ . In the LCA model, each accumulator  $x_i$  is initialized to 0. At each timestep, each accumulator value changes until one accumulator has a value greater than or equal to the threshold  $\alpha$ . At which point, the value of the threshold-surpassing accumulator plus non-decision time  $t_0$  represents the response time for choice  $i$  made by the participant. Again, during this accumulation process, a lower reflecting boundary is in place such that the value of the accumulator cannot become negative. In fitting the LCA model to the data, we specified the same priors for each of the  $\rho$ ,  $\alpha$ , and  $t_0$  parameters as for FFI. For  $\beta$  and  $k$ , we specified a prior of a normal distribution with a mean of 0 and standard deviation of 1.4 passed through an inverse logit transform.

We fit two additional variants of the LCA model to the data: the leak-only LCA and the race LCA models. In the leak-only LCA model, we fixed  $\beta$  to 0. We fit the leak-only LCA to the data to illustrate the importance of lateral inhibition to the LCA model (Purcell et al., 2010). In the race LCA model, we fixed  $\beta$  and  $\kappa$  to 0. This change converts the LCA model into a racing diffusion model with a lower bound where the accumulators are independent of each other (Bogacz et al., 2007). Both the leak-only LCA model and the race LCA model are independent accumulator models because they both lack the lateral inhibition parameter.

**The Linear Ballistic Accumulator (LBA) Model.** Unlike the other models we examined, the LBA model represents the evidence accumulation process linearly rather than with noisy samples of evidence. The LBA model also considers variability in the starting point and variability in the drift rates as parameters in the evidence accumulation process, so we can compare the fits of this model to the other models and determine if these considerations are advantageous for fitting to these data. Note, we do not consider the LBA model to be a true between-trial variability model because variability in the drift rate and variability in the starting point are required for this model to generate response time distributions, because there is no within-trial variability present in the LBA model as there is in LCA and FFI models.

The LBA model represents each accumulator independently and assumes the path of each accumulator towards the threshold is linear as represented in the following equations:

$$s_i \sim U(0, z), \quad (7)$$

$$d_i \sim N(\rho_i, \eta), \quad (8)$$

$$x_i = \frac{\alpha - s_i}{d_i}, \quad (9)$$

where  $s_i$  and  $d_i$  are, for the  $i$ th alternative, the starting point and the drive in the accumulation process respectively,  $\alpha$  is the threshold, and  $\eta$  is fixed to 1. The LBA explains variability in the data by drawing uniform random values for the starting point of each alternative and by drawing normally-distributed random values for the drive in the accumulation process of each alternative. The decision time  $x_i$  is calculated from the above equation and added to non-decision time  $t_0$  to generate the response time. In fitting the LBA model to the data, we specified the same priors for each of the  $\rho$ ,  $\alpha$ , and  $t_0$  parameters as for FFI and LCA. For  $z$ , the prior specified was a truncated normal distribution with a mean of 2.5, standard deviation of 10, lower bound of 0, and upper bound of 30.

Table 1

*Summary of free parameters in the examined models and the priors for those parameters.*

*The logistic operation represents the inverse logit transform.*

Category	Parameter	Description	Prior
All models	$\alpha$	Decision threshold	N(2.5, 10, 0, 30)
	$t_0$	Non-decision time	U(0, min_rt)
	$\rho_i$	Drift rate for choice $i$ (4 separate parameters)	N(2.5, 5, 0, 10)
LCA	$\kappa$	Decay of information over time	Logistic(N(0, 1.4))
	$\beta$	Strength of lateral inhibition	Logistic(N(0, 1.4))
FFI	$\nu$	Feed-forward inhibition	Logistic(N(0, 1.4))
LBA	$z$	Starting point variability	N(2.5, 10, 0, 30)

### Details of the model-fitting process

As with the behavioral analysis, responses faster than 0.2 seconds and slower than 5 seconds were excluded from the model fitting process (<3.2% of all data). We fit each model individually to each participant, and we applied Differential Evolution Markov Chain Monte Carlo (DE-MCMC) via the RunDEMC library (<https://github.com/compmem/RunDEMC>) to sample from the posterior distribution (Turner et al., 2013). For the DE-MCMC procedure, we initialized  $10k$  parallel chains, where  $k$  is the number of free parameters per participant and simulated with the procedure for 400 iterations of burn-in followed by 1000 samples from the posterior per chain. To evaluate the quality of each proposal, we used the probability density approximation (PDA) method to approximate the likelihood function for each model (Turner & Sederberg, 2014). We chose this procedure because it efficiently generates informative proposals by harnessing information about the structure of the posterior distribution. For PDA, we simulated each parameter proposal 50,000 times to generate an approximate density function which we then used to calculate the log likelihood. The prior distributions for each parameter that are necessary to perform Bayesian inference are listed in Table 1. To compare the models and their variants, we used the Bayesian Predictive Information Criterion (BPIC) (Ando, 2007). BPIC is more stable and penalizes for complexity more than metrics such as the DIC and has a greater scope than other model selection criteria (Ando, 2007). BPIC is designed such that models with smaller BPIC values are preferred over models with larger BPIC values.

For each of the model variants, we mapped each of the four coherences (0, 0.1, 0.2, 0.3) onto their own “drift rate” parameter, denoted as  $\rho$  in the previous model descriptions. When ordered by unique coherence pairing, all 480 of our experimental trials can be grouped into ten partitions of data. We then simulated each model variant for each of these ten partitions using the appropriate pair of these  $\rho$  parameters for the given partition (i.e.  $\rho_{0.0}$  and  $\rho_{0.1}$  for the partition where coherence was 0.0 in one direction and 0.1 in the other). Consequently, when iterating toward the best-fitting sets of parameter values, each model

must find the four  $\rho$  parameter values that provide the best account of ten conditions given the architecture of that model. Each of the ten conditions corresponds to a different response time distribution, which means the selection of a given  $\rho$  parameter may result in a perfect fit to one condition, but at the cost of a substantially poorer fit to another condition.

## Results

Table 2 and Table 3 show the mean best-fitting parameter values for each model fit to the data from Experiment 1 and Experiment 2. Shown in Figure 9 are the BPIC values calculated for each model variant, mean-centered for each participant. For the model fits to the data from Experiment 1, the LCA model had the lowest BPIC value for 13 participants, the LBA model had the lowest value for 2 participants, and the fixed FFI model had the lowest value for 1 participant. For the model fits to the data for Experiment 2, the LCA model had the lowest BPIC value for 16 participants, the LBA model had the lowest value for 5 participants, and the fixed FFI model had the lowest value for 2 participants. Clearly, out of these model variants, the LCA model provides the best fit to these perceptual decision-making datasets. In all participants, the LCA model fit the data better than the race and leak-only model variants. This suggests the addition of the lateral inhibition parameter to the LCA model architecture is indeed necessary to provide a good fit to this dataset.

Table 2

*Mean best fitting parameter values calculated between participants for each model fit to the data from Experiment 1. The standard deviation is given in parentheses.*

<i>Model</i>	$\rho_{0.0}$	$\rho_{0.1}$	$\rho_{0.2}$	$\rho_{0.3}$	$\kappa$	$\beta$	$\nu$	$z$	$\alpha$	$t_0$
Fixed	2.916	3.109	3.297	3.493	-	-	-	-	6.122	0.101
FFI	(1.407)	(1.413)	(1.395)	(1.373)					(1.269)	(0.114)
Free	0.036	0.201	0.372	0.57	-	-	0.832	-	5.929	0.093
FFI	(0.045)	(0.109)	(0.184)	(0.207)			(0.298)		(1.481)	(0.113)



<i>Model</i>	$\rho_{0.0}$	$\rho_{0.1}$	$\rho_{0.2}$	$\rho_{0.3}$	$\kappa$	$\beta$	$\nu$	$z$	$\alpha$	$t_0$
LCA	3.167	3.363	3.547	3.723	0.449	0.543	-	-	7.351	0.115
	(1.583)	(1.608)	(1.635)	(1.617)	(0.234)	(0.214)			(3.184)	(0.073)
Leak	1.800	1.915	2.034	2.177	0.624	-	-	-	4.971	0.151
Only	(1.340)	(1.352)	(1.393)	(1.398)	(0.178)				(2.091)	(0.084)
LCA										
Race	0.019	0.078	0.181	0.296	-	-	-	-	4.574	0.096
	(0.032)	(0.074)	(0.116)	(0.144)				(0.98)	(0.103)	
LBA	1.024	1.422	1.821	2.204	-	-	-	3.112	3.885	0.127
	(0.625)	(0.660)	(0.684)	(0.736)				(3.617)	(3.639)	(0.101)

Table 3

*Mean best fitting parameter values calculated between participants for each model fit to the data from Experiment 2. The standard deviation is given in parentheses.*

<i>Model</i>	$\rho_{0.0}$	$\rho_{0.1}$	$\rho_{0.2}$	$\rho_{0.3}$	$\kappa$	$\beta$	$\nu$	$z$	$\alpha$	$t_0$
Fixed	2.139	2.416	2.705	2.998	-	-	-	-	5.011	0.171
FFI	(1.495)	(1.478)	(1.463)	(1.492)					(0.97)	(0.098)
Free	1.751	2.029	2.306	2.608	-	-	0.997	-	5.056	0.165
FFI	(1.501)	(1.487)	(1.511)	(1.507)			(0.008)		(0.976)	(0.099)
LCA	4.921	5.214	5.462	5.688	0.386	0.575	-	-	9.319	0.123
	(1.404)	(1.433)	(1.457)	(1.472)	(0.186)	(0.179)			(3.126)	(0.103)
Leak	2.417	2.555	2.684	2.821	0.661	-	-	-	5.227	0.167
Only	(1.698)	(1.691)	(1.702)	(1.720)	(0.198)				(2.205)	(0.129)
LCA										
Race	0.041	0.117	0.236	0.346	-	-	-	-	3.738	0.173
	(0.061)	(0.127)	(0.132)	(0.149)					(0.668)	(0.098)

<i>Model</i>	$\rho_{0.0}$	$\rho_{0.1}$	$\rho_{0.2}$	$\rho_{0.3}$	$\kappa$	$\beta$	$\nu$	$z$	$\alpha$	$t_0$
LBA	0.887	1.262	1.600	1.922	-	-	-	1.022	1.785	0.13
	(0.671)	(0.657)	(0.575)	(0.615)				(0.995)	(1.016)	(0.119)

A comparison can be made between the FFI model with  $\nu$  fixed to 1.0 (which is similar to the DDM) and the FFI model with  $\nu$  free. The fixed FFI model had a lower BPIC value than the free FFI model in 14 participants from Experiment 1 and a lower BPIC value than the free FFI model in 22 participants from Experiment 2. Since the calculation of BPIC penalizes the inclusion of extra parameters, it would appear the benefit of freeing this parameter in the FFI model fails to overcome the complexity costs. The LBA model fit better in 10 participants than a theoretically similar race model with an LCA architecture from Experiment 1 and fit better in 21 participants from Experiment 2. Unlike the race model, the LBA model takes into account starting point variability and relies upon its between-trial variability mechanism to generate a distribution of response times. In this instance, it appears that one or both of these mechanisms have provided an advantage when fitting to these data.

Figure 10 shows all the examined dependent accumulator models simulate nearly the same increase in participant accuracy with increasing coherence difference as we observed in Experiment 1. For Experiment 2, of the dependent accumulator models, the LCA model missed the most on the 0.03 and 0.06 contrast difference conditions, despite simulating essentially the same proportion correct as observed in the 0.0 and 0.09 contrast difference conditions. Figure 10 shows, in the non-zero coherence difference conditions from both Experiment 1 and Experiment 2, the independent accumulator models consistently predict a lower proportion correct than we observed. The independent models have too few mechanisms to account for all of the effects in these data, so, throughout the course of the model fitting process, parameters are selected that can fit as much data as possible. In this instance, it seems the best-fitting parameter values were chosen to fit the equal-coherence

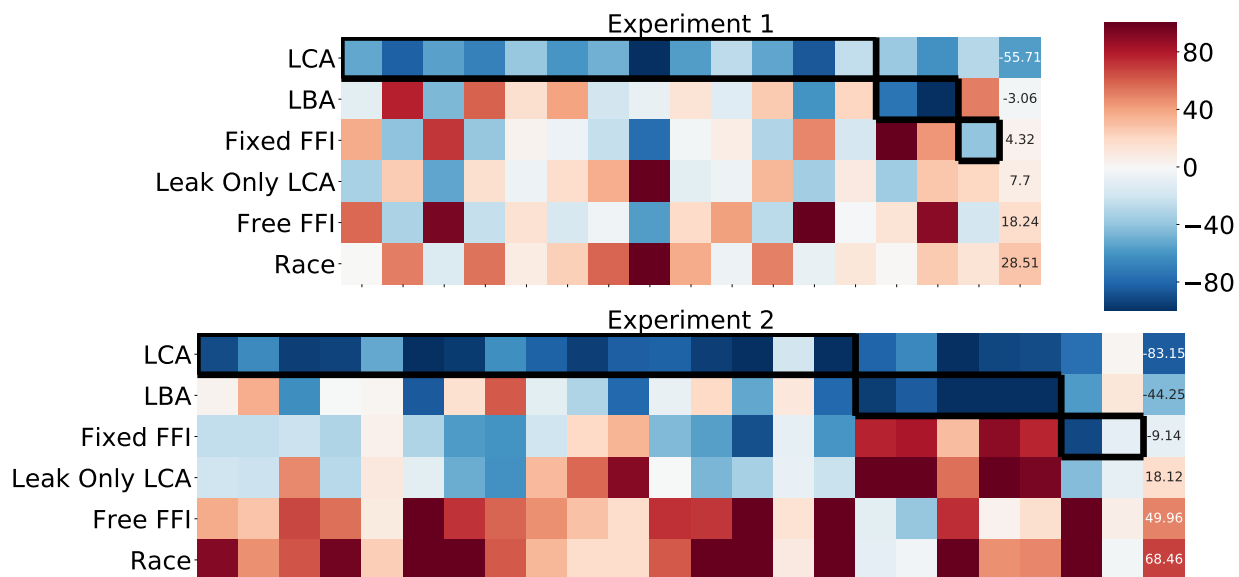


Figure 9. Heatmap showing mean-centered BPIC values for each model and for each participant. This figure is organized by participant and model variant with each column representing a single participant and each row represents a single model variant. Each square in the figure represents the mean-centered BPIC for that particular model with the mean calculated across the 6 model variants for one participant. Cooler colors represent lower (preferred) BPIC values and warmer colors represent higher BPIC values. The squares outlined by the black line represent the model variants with the lowest BPIC value of the 6 variants. The final column represents the mean BPIC value for the model calculated across participants. The LCA model had the lowest BPIC value for 13 of the 16 participants in Experiment 1 and 16 of the 23 participants in Experiment 2.

data more because most observations lie within that subset. This resulted in a poorer fit to the remainder of the proportion correct data.

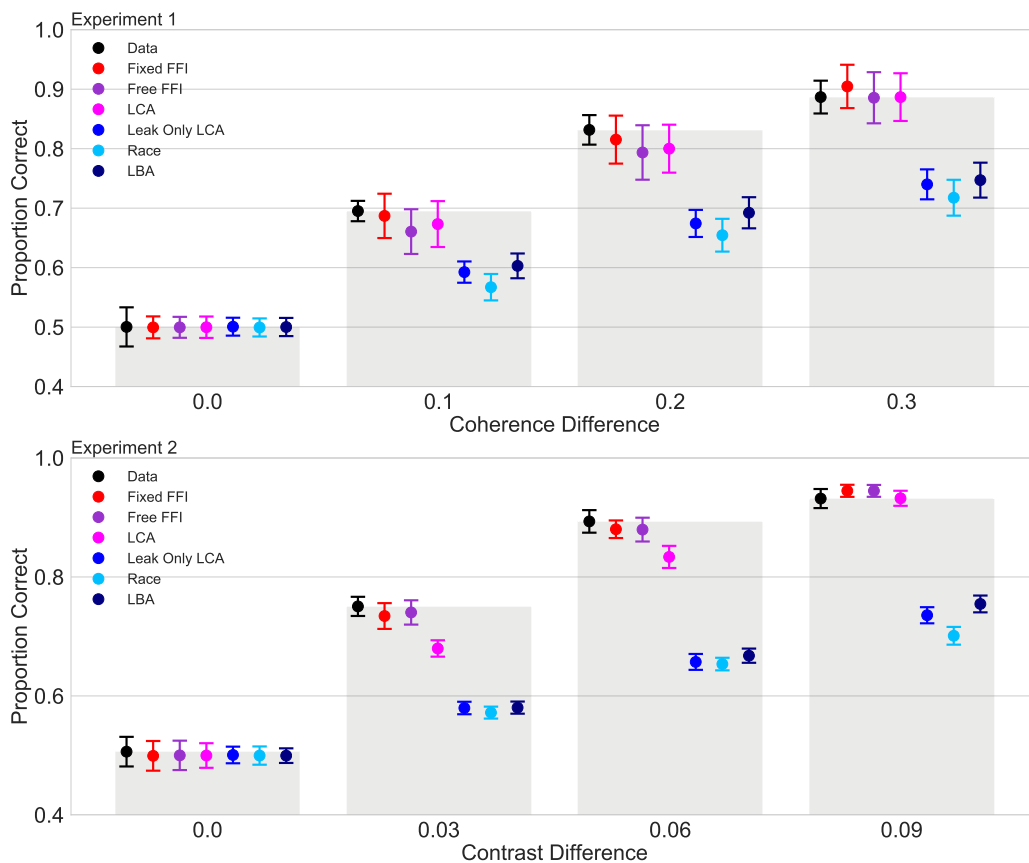


Figure 10. Simulated proportion correct as a function of coherence difference or contrast difference. The gray bars represent the same observed mean proportion correct data from Figures 4a and 7a. Using the unique best-fitting parameter values of each participant, each accumulator model was simulated within participant to generate the proportion correct for each coherence grouping for the data from Experiment 1 and for each contrast grouping for the data from Experiment 2. The mean was then calculated between the participant simulations for each model, just as in the actual data. The error bars represent Loftus and Masson (1994) corrected 95% confidence intervals.

Figure 11 shows the model fits to response times from the equal-evidence conditions from Experiment 1 and Experiment 2. Out of the dependent accumulator models, the LCA model provided the best fit to the subset of response times shown in Figure 11. The LCA model closely matched the response times in each condition of Experiment 1. For Experiment 2, the LCA model simulated progressively faster response times as the contrast increased but missed on both the 0.40 and 0.49 conditions. The Free FFI model captured the general differences in response times reasonably well but had more difficulty capturing the response times of the 0.2 coherence and 0.3 coherence conditions from Experiment 1.

For this subset of response times from Experiment 1, the Fixed FFI model produced essentially the same mean log response time for each condition. In the fixed FFI model, evidence in each direction is nearly perfectly anticorrelated. So, when evidence in both directions is equal, the drift is equal to 0 and only noise is driving the accumulation process. Thus, given a sufficient number of simulations, the fixed FFI model will predict essentially the same mean log response time for each equal-coherence condition, regardless of the level of coherence. For Experiment 2, the fixed FFI model again produced essentially the same mean log response time for each equal-contrast condition. The free FFI model also produced the same mean log response time for each equal-contrast condition. As shown in Table 3, this is because, to achieve the best fit to this dataset, the best fitting value for the input competition parameter in the free FFI model was close to 1.0. When this parameter is set to 1.0, the free FFI model is equivalent to the fixed FFI model.

One question to ask is why the best fitting input competition parameter ended up being 1.0 at the end of the model fitting process for the free FFI model in this case. The fits to Experiment 1 showed that the free FFI model can capture both the progressively decreasing response times in the equal evidence conditions and the progressively increasing response times in this subset of unequal evidence conditions, but the free FFI model was unable to simulate both patterns of response times for Experiment 2. As can be observed in Table 1, we did not provide a strong constraint on the input competition parameter with our

choice of prior distribution. Thus, over the course of the model fitting process, the algorithm determined that an input competition parameter near one provided the best fit to the most data.

One explanation for why the algorithm estimated a value of one on average for this parameter to fit Experiment 2 but estimated the parameter to be .832 on average to fit Experiment 1 could be the differences in response times between the two tasks. In the unequal evidence conditions, response times increased as the difference in evidence decreased in Experiment 1 more than the response times increased in Experiment 2. In other words, in the unequal evidence conditions, the effect of evidence difference on response time was more pronounced in Experiment 1 than in Experiment 2. Similarly, in the equal evidence conditions, response times decreased more as the total evidence present in the stimulus increased in Experiment 1 than the response times decreased in Experiment 2. In the equal evidence conditions, the effect of total evidence on response time was more pronounced in Experiment 1 than Experiment 2. To account for these differences between the conditions, the algorithm found that 1.0 was the value that could best account for the multiple patterns in the dataset rather than a value less than 1.0 which would have created less input competition.

In general, the independent accumulator models produced faster response times as the evidence increased in the equal-evidence conditions, but, in many conditions, produced slower or faster response times than we observed. In Experiment 1, the independent models closely matched the response times in the 0.1 equal-coherence condition but produced slower response times than we observed in the 0.0 condition, slightly faster response times than we observed in the 0.2 condition, and much faster response times than we observed in the 0.3 condition. In Experiment 2, the independent models matched the 0.43 condition but greatly missed the observed response times of the other three conditions.

In the conditions where exactly one coherence was equal to 0.3 in Experiment 1, the two partially dependent accumulator models provided a good fit to the mean log response

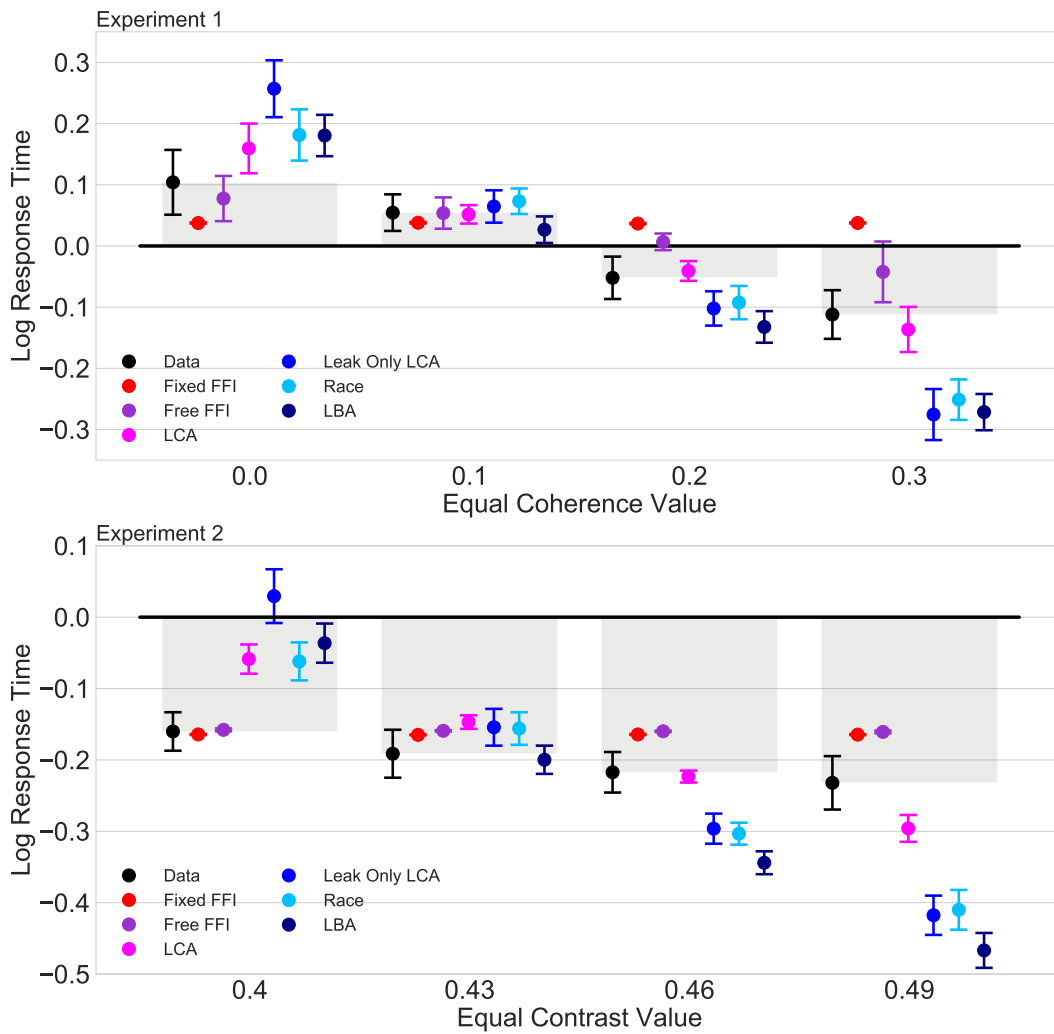


Figure 11. Simulated log response time as a function of equal-coherence or equal-contrast condition. The gray bars represent the same observed mean log response time data from Figures 4b and 7b. Using the unique best-fitting parameter values of the participant, each dependent and independent accumulator model was simulated within participant to generate response time distributions for each equal-coherence or equal-contrast grouping. Then mean log response time was calculated between participants. The error bars represent Loftus and Masson (1994) corrected 95% confidence intervals.

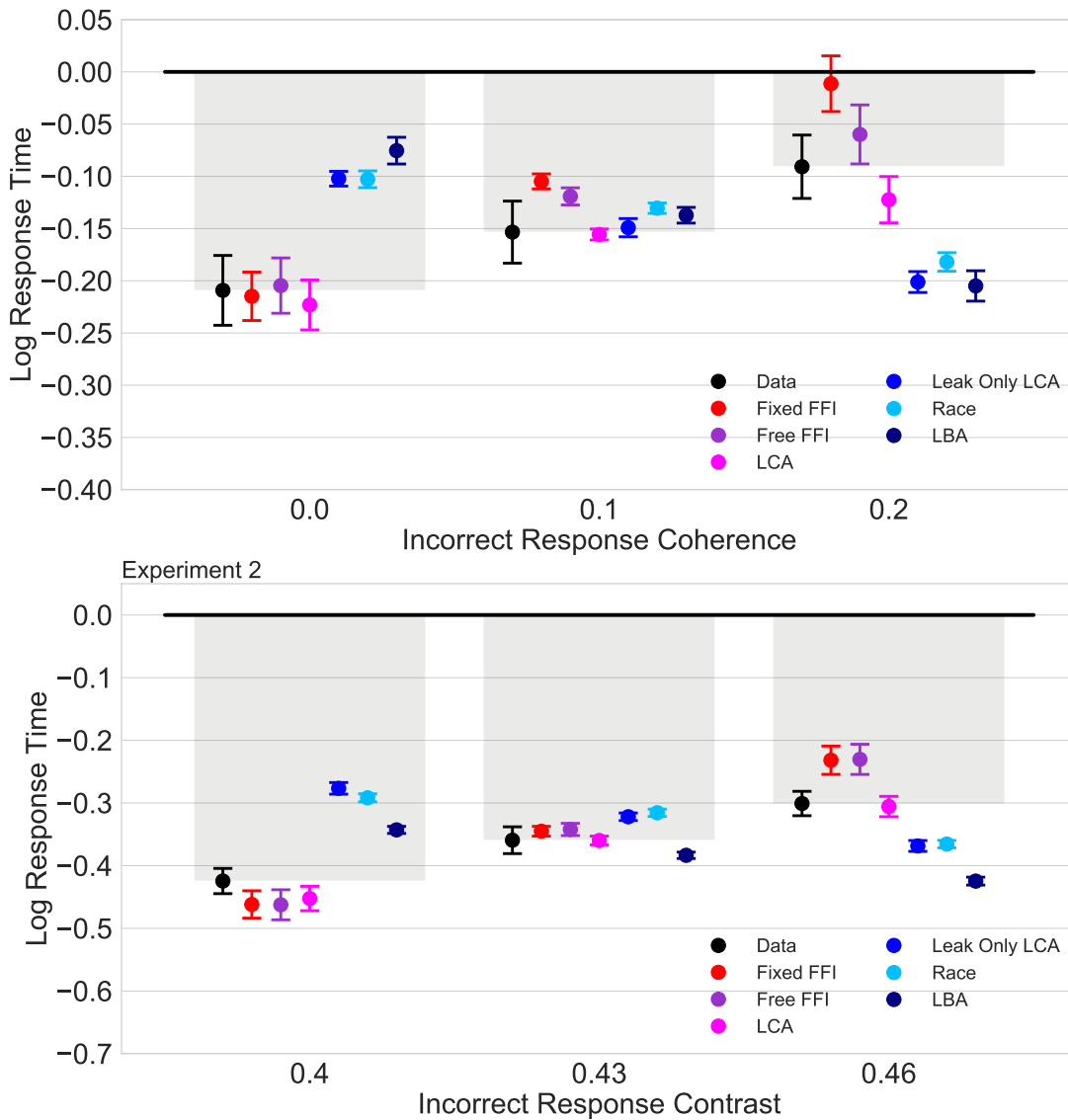


Figure 12. Simulated log response time as a function of coherence or contrast condition. The gray bars represent the same observed mean log response time data from Figures 5b and 8b. Only the conditions where exactly one direction had a coherence of 0.3 in Experiment 1 are displayed. In Experiment 2, only the conditions where exactly one grating had the highest contrast of 0.49 are displayed. For the Experiment 1 results, the response times in this figure are those paired with the correct response of 0.3. For the Experiment 2 results, the response times in this figure are those paired with the correct response of 0.49. Using the unique best-fitting parameter values of the participant, each dependent and independent



times displayed in Figure 12 and to the proportion correct data displayed in Figure 12. The fully dependent fixed FFI model provided a good fit to the proportion correct data observed in these conditions but generated faster response times than were observed in two of the conditions displayed in Figure 12. All of the independent models showed the opposite pattern of results for the response times than we observed in these conditions. The mean log response times simulated by the independent models decreased as incorrect response coherence approached 0.3. In other words, the independent models produced faster response times as the trials became objectively more difficult for the participants. In addition, the independent models produced much lower proportion correct than we observed in each of these three conditions.

In the conditions where exactly one contrast was the highest contrast of 0.49 in Experiment 2, each of the dependent accumulator models simulated slower response times for correct responses as the contrast difference decreased. Of the dependent accumulator models, the LCA model appears to have matched the observed response times the best. Crucially, the independent accumulator models simulated faster response times as the contrast difference decreased. In general, the dependent accumulator models correctly produced the observed subset of response times and the independent accumulator models failed to produce the observed response times. Of the six models we examined, the LCA model simulated the closest proportion correct to the observed data. The other dependent accumulator models simulated the same proportion correct for each condition, resulting in a substantially higher proportion correct than observed in two of these conditions. In each of these conditions, the independent accumulator models simulated a lower proportion correct than observed.

## Discussion

Overall, the partially dependent accumulator models, and, in particular, the LCA model, consistently provided the best account of the data from Experiment 1 and Experiment 2. These models correctly simulated the faster response times with greater

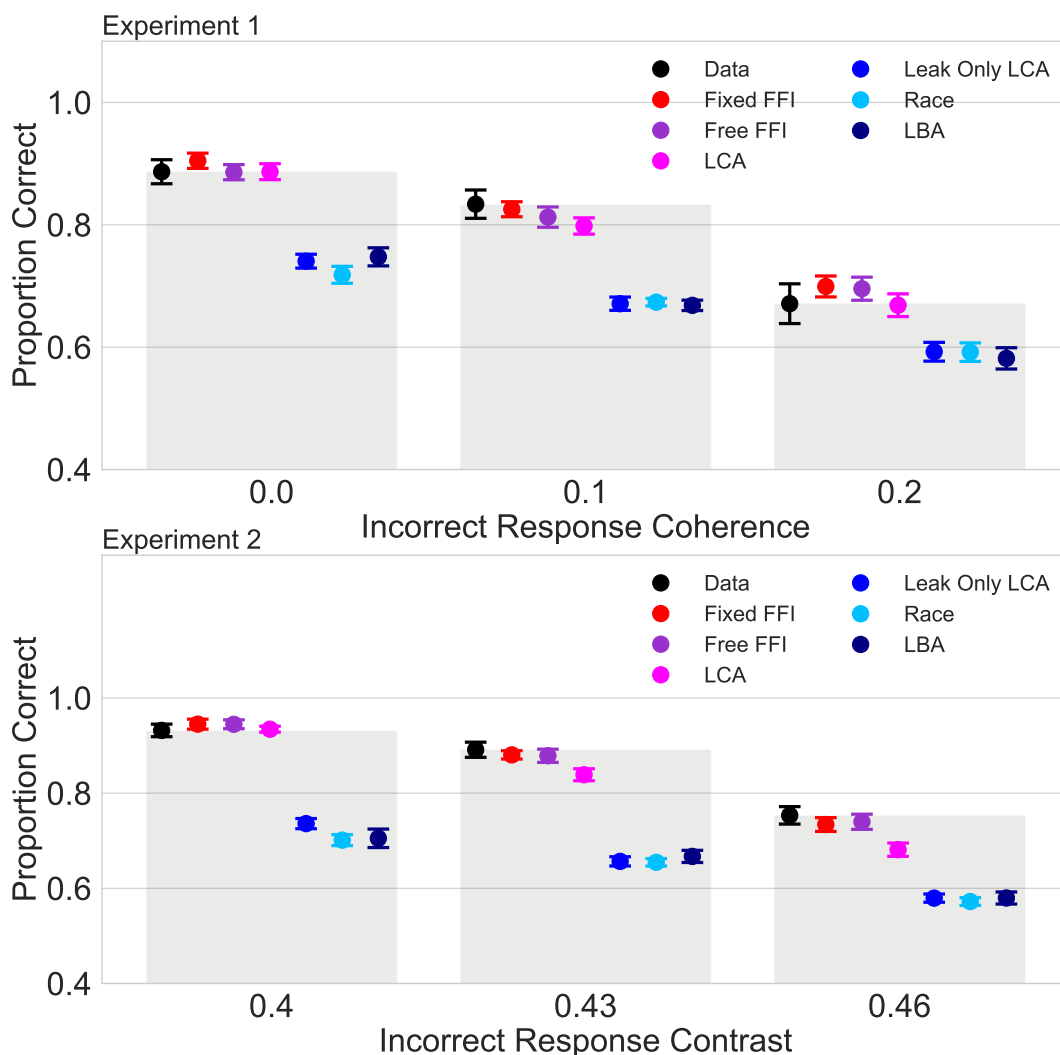


Figure 13. Simulated proportion correct as a function of coherence or contrast condition.

The gray bars represent the same observed proportion correct data from Figures 5a and 8a. Only the conditions where exactly one direction had a coherence of 0.3 in Experiment 1 are displayed. In Experiment 2, only the conditions where exactly one grating had the highest contrast of 0.49 are displayed. Using the unique best-fitting parameter values of the participant, each dependent and independent accumulator model was simulated within participant to calculate the proportion correct for each coherence or contrast grouping. Then mean proportion correct was calculated between participants. The error bars represent Loftus and Masson (1994) corrected 95% confidence intervals.

evidence in the equal-evidence conditions and slower response times as the difference between the evidence in the unequal-evidence conditions decreased. In contrast, while the fixed FFI model featuring fully dependent accumulators correctly simulated slower response times as the difference in evidence increased in the unequal-evidence conditions, it failed to simulate the faster response times observed as the evidence increased in the equal-evidence conditions. Furthermore, the three independent accumulator models we examined failed to simulate the slowing of response times for the correct choice of the 0.3 coherence when the competing coherence increases from 0.0 to 0.2. The independent accumulator models also predicted less accurate responses than we observed in each of the three unequal-evidence conditions we examined here. Thus, the models with partially dependent accumulators are the preferred models for fitting to our data.

Our tasks are similar to many standard perceptual decision-making tasks that have historically supported the models we examined, but combines all of these conditions together to create a challenging set of patterns of results for extant models of decision-making to fit. A priori, we expected the independent accumulator models to provide a good fit to the equal-evidence data and provide a poor fit to the unequal-evidence data because the accumulator for the incorrect response does not directly inhibit the accumulator for the correct response. As predicted, because the independent accumulator models have no inhibitory components by definition, the models could not generate the slower response times as the decision becomes more difficult in unequal-evidence conditions. The independent accumulator models did generate faster response times as the coherence increased in the equal-evidence conditions, but, unexpectedly, the models generated slower response times than were observed in the 0.0 equal-coherence condition and faster response times than were observed in the 0.3 equal-coherence condition from Experiment 1. The independent accumulator models also, contrary to our predictions, generated slower response times than were observed in the 0.4 equal-contrast condition and faster response times than were observed in the 0.46 and 0.49 equal-contrast conditions from Experiment 2. We suspect the

multiple constraints imposed by our task design caused these misses in some conditions where the theory would suggest a more accurate fit.

Fully dependent accumulator models have difficulty matching the response times observed in equal-evidence conditions. In a two-choice decision, this is the case because fully dependent models represent the drift rate towards the decision threshold as the difference between the two drift rates. Because the drift rates for both alternatives have the same value, the difference between them is zero. With a drift rate of zero, the evidence accumulator crosses the decision threshold solely based on random noise in the signal, which, on average, is equivalent across conditions.

Partially dependent accumulator models have the mechanisms to account for the choice–response time patterns observed in both the equal-evidence and unequal-evidence conditions of our experiments. In these models, for the equal-evidence conditions, the drift rates for the accumulators representing each direction will be greater with more evidence. Now the difference between these drift rates is zero, but, since the accumulators only partially inhibit each other, the accumulators will exceed the decision threshold more quickly with more evidence and produce the patterns observed in our data. Partially dependent accumulator models also capture the unequal-evidence data because the degree of inhibition on each accumulator will increase as the coherence difference decreases. As the decision becomes more difficult, response times are slower and accuracy approaches chance. Both features are captured well in partially dependent accumulator models.

Model Analysis 1 provided strong evidence supporting partially dependent accumulator models (and, in particular, the LCA model) over both fully dependent and independent accumulator models. Recent work has examined how the mechanisms of within and between-trial variability can be implemented in decision-making models to capture magnitude effects in simple perceptual decision-making tasks similar to those presented in this paper (Ratcliff et al., 2018; Teodorescu et al., 2015). We wanted to determine if adding these additional sources of variability into our examined models would change the fits and, in

particular, give the fully dependent fixed FFI model the flexibility to provide a good fit to all of the patterns observed in our data.

## Model Analysis 2

As mentioned in the introduction, one approach to accounting for the decrease in response time as evidence increases in the equal-evidence conditions is to add sources of variability that are correlated with the evidence levels. In a similar manner to our experiments, both Teodorescu et al. (2015) and Ratcliff et al. (2018) have examined the fits of the LCA model and the DDM to data exhibiting magnitude effects. Teodorescu and colleagues (2015) found that an implementation of the DDM where processing noise in the model increases with increased input intensity fit 3 out of 6 participants worse than the LCA model in the choice–response time data from a two-choice fluctuating brightness discrimination task, thereby not favoring either model. However, Ratcliff and colleagues (2018) found that a similar implementation of the DDM provided a better fit in more participants than the LCA model to their fluctuating brightness discrimination task data.

We believe that one reason that these two studies published conflicting results is that the specific models compared differed by too many mechanisms for the key mechanism of the LCA model, lateral inhibition, to be fairly compared to the inhibitory process of the DDM. To uncover which model framework (partially dependent or fully dependent accumulation) and which additional sources of variability (between-trial or within-trial) provide the best account of these data, we decided to fit the relevant models from this line of research to our two experiments. We believe that we created a fair comparison of the extant mechanisms of the models by giving both the fixed FFI and LCA models the same additional sources of variability. Our work sits in contrast to this previous research where the LCA model was given no additional sources of variability whereas the DDM was given an additional source of either between-trial or within-trial variability. To implement the models in this way, we modified the LCA and the fixed FFI models to represent the drift rate as a linear function of

the difference in magnitude of the stimulus strengths and added parameters and mechanisms to link within-trial and between-trial variability to the coherence of the input. Our modifications are similar to the modifications made in the previous work that examined the fits of these models to data where the overall evidence present in the stimulus was manipulated between trial while maintaining a constant difference in evidence between the choice alternatives (Ratcliff et al., 2018; Teodorescu et al., 2015). We will now explain each of model variants in this study in turn.

### No Additional Variability Models

We examined both the fixed FFI model and the LCA model with no additional sources of variability than the models from Model Analysis 1, but with the change to how the drift rate in each model is calculated. In the remainder of the paper, we will call these models the NV fixed FFI model and the NV LCA model respectively. We fit these models to our data to determine if representing the rate of accumulation as a function of the difference in the stimulus values provides an acceptable (or even improved) fit than assigning a separate drift rate for each coherence or contrast value. One advantage of these models is that they have two fewer parameters than their model counterparts in Model Analysis 1. First, we standardized the stimulus values to be on a 0 to 1 scale (a separate 0 to 1 scale for the RDM task and for the grating task). Then, using the same 10 conditions that we used in Model Analysis 1, we calculated the differences between the stimulus values and scaled that difference by the parameter  $\mu$ . Then our drift rate for each condition becomes the scaled difference value added to our single  $\rho$  parameter, which is still a free parameter in this model. The  $\rho$  parameter is necessary to allow the drift rate to be greater than 0 in the conditions where the evidence supporting each direction is equivalent. The following equations describe the NV fixed FFI model for 2-alternative forced-choice tasks:

$$m_1 = \mu(I_1 - I_2)\frac{dt}{\tau} + \xi_1\sqrt{\frac{dt}{\tau}}, \quad (10)$$

$$m_2 = \mu(I_2 - I_1)\frac{dt}{\tau} + \xi_2\sqrt{\frac{dt}{\tau}}, \quad (11)$$

$$dx_1 = \rho\frac{dt}{\tau} + m_1 - m_2, \quad (12)$$

$$dx_2 = \rho\frac{dt}{\tau} + m_2 - m_1, \quad (13)$$

$$\xi_1 \sim N(0, \eta), \quad (14)$$

$$\xi_2 \sim N(0, \eta), \quad (15)$$

where  $I_1$  and  $I_2$  represent the normalized stimulus value for stimuli 1 and 2 respectively,  $\mu$  modulates the effect of the difference in supporting evidence upon the rate of evidence accumulation,  $\rho$  represents the base rate of evidence accumulation,  $dx_1$  and  $dx_2$  represent the change in the value of the accumulators  $x_1$  and  $x_2$  respectively at each timestep, and  $\xi_1 \sim N(0, \eta)$  and  $\xi_2 \sim N(0, \eta)$  represent the within-trial variability for the drive to accumulators  $x_1$  and  $x_2$  respectively (each  $\xi$  is recalculated at each timestep). The following equations describe the NV LCA model for 2-alternative forced-choice tasks:

$$m_1 = \mu(I_1 - I_2)\frac{dt}{\tau} + \xi_1\sqrt{\frac{dt}{\tau}}, \quad (16)$$

$$m_2 = \mu(I_2 - I_1)\frac{dt}{\tau} + \xi_2\sqrt{\frac{dt}{\tau}}, \quad (17)$$

$$dx_1 = m_1 + (\rho - \kappa x_1 - \beta x_2)\frac{dt}{\tau}, \quad (18)$$

$$dx_2 = m_2 + (\rho - \kappa x_2 - \beta x_1)\frac{dt}{\tau}, \quad (19)$$

$$\xi_1 \sim N(0, \eta), \quad (20)$$

$$\xi_2 \sim N(0, \eta), \quad (21)$$

where  $I_1$ ,  $I_2$ ,  $\mu$ ,  $dx_1$ ,  $dx_2$ ,  $\rho$ ,  $\xi_1$ , and  $\xi_2$  represent the same processes as in the NV fixed FFI model,  $\kappa$  represents the passive decay of evidence, and  $\beta$  represents the amount of lateral inhibition applied by one accumulator upon the other accumulator.

### Between-Trial Variability Models

Between-trial variability in drift rate is a powerful mechanism that has been used to fit fully dependent accumulator models to a variety of choice–response time data (Ratcliff, 2006; Ratcliff & McKoon, 2008; Ratcliff & Rouder, 1998; Ratcliff et al., 2018). Ratcliff et al. (2018) have proposed to linearly scale the between-trial variability parameter as a function of the evidence present in the stimulus as a mechanism to allow fully dependent accumulator models to account for the magnitude effect. To test the fidelity of their proposal, we included the linearly scaling mechanism in both the LCA and fixed FFI models, and compared them to the standard LCA and fixed FFI models, respectively. For the remainder of this paper, we will abbreviate the models that linearly scale the between-trial variability as the BTV LCA model and the BTV fixed FFI model, respectively. For each of these models, we scale and calculate the differences in stimulus strength in the same manner as the NV models. The key difference between the BTV models and the NV models is how the drift rate is calculated. For each simulation, the drift rate for options 1 and 2 are drawn from a normal distribution where the mean of the distribution varies with the difference in stimulus strength between the two options and the standard deviation of the distribution varies with the strength of options 1 and 2 respectively. Thus, the drift rate for a choice depends on both the evidence supporting that choice, and the difference in evidence for that choice and the other choices. The following equations describe the BTV fixed FFI model for 2-alternative forced-choice tasks:

$$m_1 = \mu(I_1 - I_2), \quad (22)$$

$$m_2 = \mu(I_2 - I_1), \quad (23)$$



$$s_1 = \sigma I_1, \quad (24)$$

$$s_2 = \sigma I_2, \quad (25)$$

$$d_1 \sim N(m_1, s_1), \quad (26)$$

$$d_2 \sim N(m_2, s_2), \quad (27)$$

$$dr_1 = d_1 \frac{dt}{\tau} + \xi_1 \sqrt{\frac{dt}{\tau}}, \quad (28)$$

$$dr_2 = d_2 \frac{dt}{\tau} + \xi_2 \sqrt{\frac{dt}{\tau}}, \quad (29)$$

$$dx_1 = \rho \frac{dt}{\tau} + dr_1 - dr_2, \quad (30)$$

$$dx_2 = \rho \frac{dt}{\tau} + dr_2 - dr_1, \quad (31)$$

$$\xi_1 \sim N(0, \eta), \quad (32)$$

$$\xi_2 \sim N(0, \eta), \quad (33)$$

where  $I_1$ ,  $I_2$ ,  $\mu$ ,  $dx_1$ ,  $dx_2$ ,  $\rho$ ,  $\xi_1$  and  $\xi_2$  represent the same processes as in the NV models,  $\sigma$  modulates the effect of stimulus strength on the standard deviation of the normal distribution from which drift rate is calculated,  $d_1$  and  $d_2$  represent the drift rates for options 1 and 2 respectively and  $dr_1$  and  $dr_2$  represent the drive to the  $x_1$  and  $x_2$  accumulators respectively. The following equations describe the BTV LCA model for 2-alternative forced-choice tasks:

$$m_1 = \mu(I_1 - I_2), \quad (34)$$

$$m_2 = \mu(I_2 - I_1), \quad (35)$$

$$s_1 = \sigma I_1, \quad (36)$$

$$s_2 = \sigma I_2, \quad (37)$$

$$d_1 \sim N(m_1, s_1), \quad (38)$$

$$d_2 \sim N(m_2, s_2), \quad (39)$$

$$dr_1 = d_1 \frac{dt}{\tau} + \xi_1 \sqrt{\frac{dt}{\tau}}, \quad (40)$$

$$dr_2 = d_2 \frac{dt}{\tau} + \xi_2 \sqrt{\frac{dt}{\tau}}, \quad (41)$$

$$dx_1 = dr_1 + (\rho - \kappa x_1 - \beta x_2) \frac{dt}{\tau}, \quad (42)$$

$$dx_2 = dr_2 + (\rho - \kappa x_2 - \beta x_1) \frac{dt}{\tau}, \quad (43)$$

$$\xi_1 \sim N(0, \eta), \quad (44)$$

$$\xi_2 \sim N(0, \eta), \quad (45)$$

where  $\kappa$  and  $\beta$  represent the same leak and lateral inhibition processes as in the NV LCA model and the remaining parameters represent the same processes as in the BTV fixed FFI model.

### Within Trial Variability Models

Researchers have also accounted for magnitude effects in simple perceptual decision-making tasks by including an additional source of within-trial variability in their model that increases as a function of the stimulus strength (Teodorescu et al., 2015).

Within-trial variability has been an important mechanism in several recent studies because

activation-dependent, multiplicative noise introduces magnitude sensitivity at the level of the input and is subsequently independent of the main decision mechanism (Louie, Khaw, & Glimcher, 2013; Teodorescu et al., 2015). For our purposes, the within-trial variability mechanism could allow the fixed FFI model to simulate faster response times as the evidence increases in the equal-evidence conditions and overcome the strong input competition in the model, which led to response time distributions that did not match the observed data in Model Analysis 1. In our implementation, we modulated within-trial variability through the parameter  $v$ . As with the BTV models above, we added the within-trial variability mechanism to both the LCA model and the fixed FFI model. We will refer to these models the WTV LCA model and the WTV fixed FFI model for the remainder of the paper. The WTV fixed FFI model is described by the following equations for 2-alternative forced-choice tasks:

$$m_1 = \mu(I_1 - I_2) \frac{dt}{\tau} + \xi_1 \sqrt{\frac{dt}{\tau}}, \quad (46)$$

$$m_2 = \mu(I_2 - I_1) \frac{dt}{\tau} + \xi_2 \sqrt{\frac{dt}{\tau}}, \quad (47)$$

$$dx_1 = \rho \frac{dt}{\tau} + m_1 - m_2, \quad (48)$$

$$dx_2 = \rho \frac{dt}{\tau} + m_2 - m_1, \quad (49)$$

$$\xi_1 \sim N(0, \eta + vI_1), \quad (50)$$

$$\xi_2 \sim N(0, \eta + vI_2), \quad (51)$$

where  $I_1$ ,  $I_2$ ,  $\mu$ ,  $dx_1$ ,  $dx_2$  and  $\rho$  represent the same processes as in the NV models  $\xi_1$  and  $xi_2$  represent the within-trial variability for stimuli 1 and 2 respectively, and  $v$  is the sensitivity of within-trial variability to the stimulus strength. The following equations describe the WTV LCA model for 2-alternative forced-choice tasks:

$$m_1 = \mu(I_1 - I_2) \frac{dt}{\tau} + \xi_1 \sqrt{\frac{dt}{\tau}}, \quad (52)$$

$$m_2 = \mu(I_2 - I_1) \frac{dt}{\tau} + \xi_2 \sqrt{\frac{dt}{\tau}}, \quad (53)$$

$$dx_1 = m_1 + (\rho - \kappa x_1 - \beta x_2) \frac{dt}{\tau}, \quad (54)$$

$$dx_2 = m_2 + (\rho - \kappa x_2 - \beta x_1) \frac{dt}{\tau}, \quad (55)$$

$$\xi_1 \sim N(0, \eta + vI_1), \quad (56)$$

$$\xi_2 \sim N(0, \eta + vI_2), \quad (57)$$

where  $\kappa$  and  $\beta$  represent the same leak and lateral inhibition processes as in the NV LCA model and the remaining parameters represent the same processes as in the WTV fixed FFI model.

The prior distributions for each parameter in each of the six models are listed in Table 4. We fixed  $dt = 0.01$ ,  $\tau = 0.1$ , and  $\eta = 1$ . As in Model Analysis 1, during the evidence accumulation process, the value of the accumulator in each of the six model variants in this study is prevented from falling below 0 via the following equations

$$x_1 \rightarrow \max(x_1, 0), \quad (58)$$

$$x_2 \rightarrow \max(x_2, 0), \quad (59)$$

Table 4

*Summary of free parameters in the examined models and the priors for those parameters.  $C$  represents the half cauchy distribution, and the logistic operation represents the inverse logit transform.*

Category	Parameter	Description	Prior
All models	$\alpha$	Decision threshold	N(2.5, 10, 0, 30)
	$t_0$	Non-decision time	U(0, min_rt)
	$\rho$	Base level of drift rate	N(2.5, 5, 0, 10)

Category	Parameter	Description	Prior
LCA	$\mu$	Effect of stimulus difference on drift rate	$N(0.75, 2, 0, 10)$
	$\kappa$	Decay of information over time	$\text{Logistic}(N(0, 1.4))$
	$\beta$	Strength of lateral inhibition	$\text{Logistic}(N(0, 1.4))$
BTV	$\sigma$	Sensitivity of BTV to stimulus strength	$C(0, 5.0)$
WTV	$v$	Sensitivity of WTV to stimulus strength	$\gamma(4, 1)$

## Results

Table 5 and Table 6 show the mean best-fitting parameter values for each model fit to the data from Experiment 1 and Experiment 2. Figure 14 shows the mean-centered BPIC values for each model and participant from Experiment 1 and 2. In Experiment 1, the BTV LCA had the lowest BPIC value in 5 participants, the NV LCA had the lowest BPIC value in 9 participants, the NV fixed FFI had the lowest BPIC value in 1 participant, and the WTV fixed FFI had the lowest BPIC in 1 participant. Overall, the LCA models had the lowest BPIC values for 14 of the 16 participants. In Experiment 2, the BTV LCA had the lowest BPIC value in 11 participants, the NV LCA had the lowest BPIC value in 10 participants, the WTV LCA had the lowest BPIC value in 1 participant, and the NV fixed FFI had the lowest BPIC value in 1 participant. Overall, the LCA models had the lowest BPIC values for 22 of the 23 participants, providing more evidence that the partially dependent LCA model is the preferred model for our datasets over the fully dependent fixed FFI model, even when additional sources of variability are considered.

To confirm the representations for the models of Model Analysis 2 provide a better fit than their counterparts in Model Analysis 1, we compared the NV fixed FFI model with the Model Analysis 1 fixed FFI model and the NV LCA model with the Model Analysis 1 LCA model. The NV fixed FFI model had lower BPIC values than the fixed FFI model of Model Analysis 1 for all 16 of the participants from Experiment 1 and for all 23 of the participants

from Experiment 2. The NV LCA model had lower BPIC values than the LCA model of Model Analysis 1 for all 16 of the participants from Experiment 1 and for 21 of the 23 participants from Experiment 2. Taken together, these results provide evidence that these NV models fit the data from both experiments better than the models of Model Analysis 1.

In Appendix C, we examined six additional model variants. In one set of model variants, we fit three additional free FFI models: an NV free FFI model, a BTV free FFI model, and a WTV free FFI model. We also fit additional NV LCA, BTV LCA, and WTV LCA models where the drift rate is not a function of the stimulus difference and the drift rate for each accumulator is simply a function of the stimulus value that accumulator represents. We call these LCA models the LCA ND models where ND means “no difference”. When comparing the models by BPIC value, we found that the NV LCA, BTV LCA, and WTV LCA models each fit more participants better than the free FFI or LCA ND model with the same additional sources of variability. Furthermore, the NV LCA ND, BTV LCA ND, and WTV LCA ND models each fit better than the fixed FFI and free FFI models with the same additional sources of variability.

Table 5

*Mean best fitting parameter values calculated between participants for each model fit to the data from Experiment 1. The standard deviation is given in parentheses.*

<i>Model</i>	$\rho$	$\mu$	$\sigma$	$v$	$\kappa$	$\beta$	$\alpha$	$t_0$
BTV	3.161	0.351	0.319	-	0.466	0.522	6.506	0.149
LCA	(2.406)	(0.192)	(0.283)		(0.134)	(0.167)	(3.642)	(0.087)
WTV	3.344	0.282	-	0.856	0.47	0.517	6.711	0.164
LCA	(2.766)	(0.161)		(0.422)	(0.142)	(0.169)	(3.951)	(0.095)
NV	3.12	0.301	-	-	0.477	0.542	6.414	0.156
LCA	(2.482)	(0.138)			(0.164)	(0.125)	(3.616)	(0.074)
BTV	0.011	0.381	0.619	-	-	-	6.814	0.102



<i>Model</i>	$\rho$	$\mu$	$\sigma$	$v$	$\kappa$	$\beta$	$\alpha$	$t_0$
NV	0.003	0.419	-	-	-	-	5.087	0.164
Fixed	(0.004)	(0.129)					(0.973)	(0.097)
FFI								

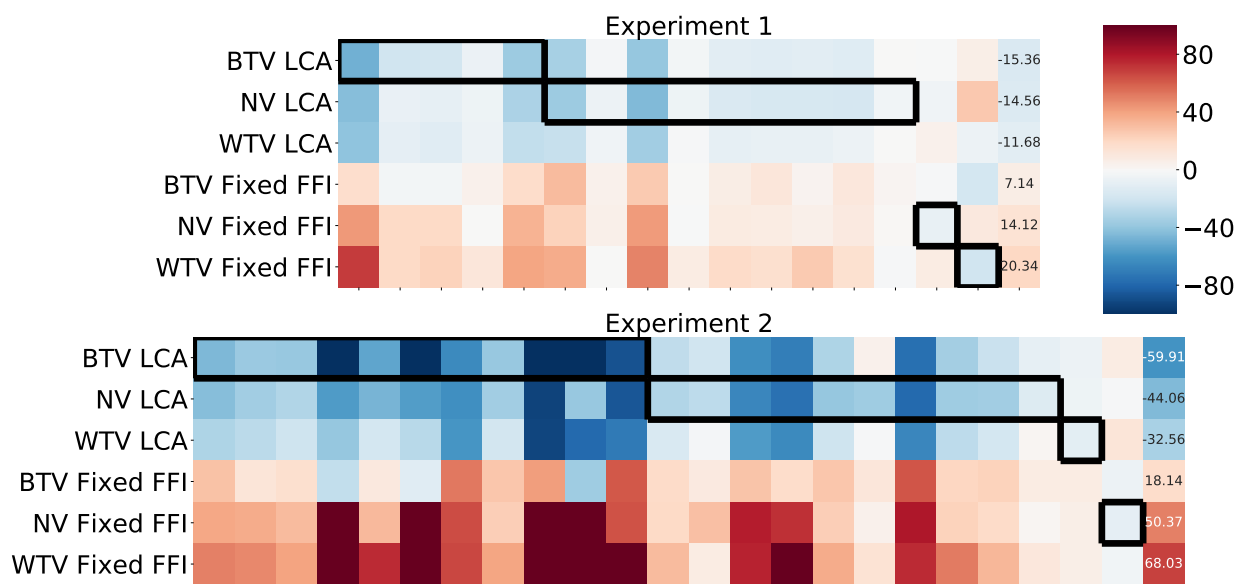


Figure 14. Heatmap showing mean-centered BPIC values for each model and for each participant. This figure is organized by participant and model variant with each column representing a single participant and each row representing a single model variant. Each square in the figure represents the mean-centered BPIC for that particular model with the mean calculated across the 6 model variants for the one participant. Cooler colors represent lower (preferred) BPIC values and warmer colors represent higher BPIC values. The squares outlined by the black line represent the model variants with the lowest BPIC value of the 6 variants. The final column represents the mean BPIC value for the model calculated across participants. The model with the LCA architecture had the numerically lowest BPIC value for 14 of the 16 participants in Experiment 1 and for 22 of the 23 participants in Experiment 2.



Figure 15 shows the fit of each model to the proportion correct in the RDM data from Experiment 1 and grating data from Experiment 2. For Experiment 1, the WTV LCA, NV LCA, and NV fixed FFI models closely match the observed proportion correct in every coherence difference condition. The BTV LCA and WTV fixed FFI models slightly underestimate the observed proportion correct in the 0.1 and 0.2 coherence difference conditions but simulate the correct proportion correct in the other two conditions for Experiment 1. The BTV fixed FFI model underestimates the proportion correct the more than the other 5 model variants, but this model correctly simulates increasing accuracy as coherence difference increases. For Experiment 2, all models provide similar fits to these data with the BTV fixed FFI having a slightly worse fit than the other models to the 0.1 and 0.2 conditions.

Figure 16 shows the fit of each model to the equal-evidence response time data from Experiments 1 and 2. Each model struggled to fit to some aspect of the data from Experiment 1 with the BTV LCA model providing the closest fit overall. As expected, for both Experiments 1 and 2, both the NV LCA model and the NV Fixed FFI model could not simulate the decrease in response time as the evidence increased. This is because the difference in stimulus strength is zero regardless of equal-evidence condition, and there is no other mechanism besides the difference that can adjust the response times as the stimulus magnitude changes. A stimulus strength of zero will result in the same response times irrespective of the amount of evidence present in the stimulus.

For Experiment 1, the BTV fixed FFI and WTV fixed FFI models simulate different response times than we observed in the equal-coherence data because the amount of inhibition upon each accumulator is high in spite of the additional sources of variability. The additional sources of variability allow these fixed FFI models to simulate faster response times as the amount of evidence present in the stimulus increases, but the response times are faster than we observe in the 0.3 condition and slower than we observe in the 0.0 condition. The BTV LCA model can fit these data because each accumulator only partially inhibits the

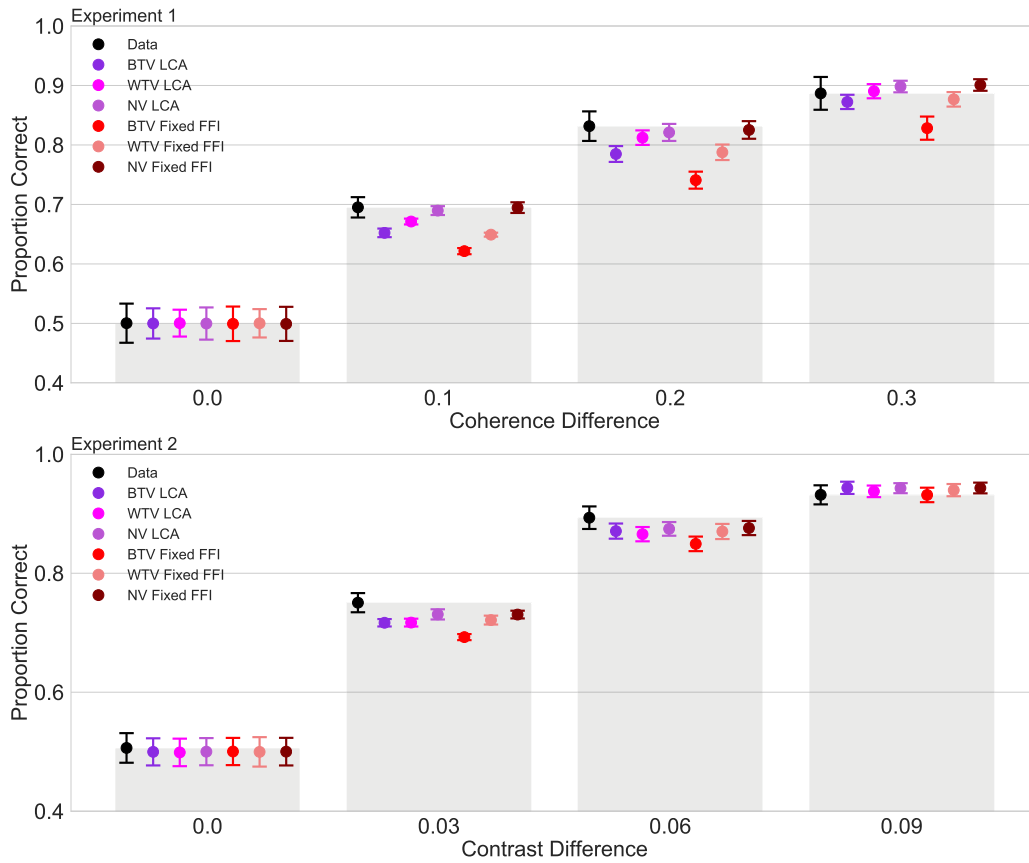


Figure 15. Simulated proportion correct as a function of coherence difference or contrast difference. The gray bars represent the same observed mean proportion correct data from Figures 4a and 7a. Using the unique best-fitting parameter values of the participant, each dependent and independent accumulator model was simulated within participant to generate the proportion correct for each coherence grouping for the data from Experiment 1 and for each contrast grouping for the data from Experiment 2. The mean was then calculated between the participant simulations for each model. The error bars represent Loftus and Masson (1994) corrected 95% confidence intervals.

other accumulator in the model which gives the model more flexibility in fitting each participant. For Experiment 2, the BTV and WTV LCA models fit the equal-contrast data well in general, but the BTV and WTV fixed FFI models simulate response times that are too slow in the 0.4 and 0.43 conditions. It appears that the mechanisms in these models are too insensitive to simulate the slope of the observed decrease in response time as a function of contrast condition.

Figure 16 shows the fit of each model to the response times associated with the correct response in conditions where exactly one coherence is equal to 0.3 in Experiment 1 and where exactly one contrast is equal to the highest contrast in Experiment 2. For Experiment 1, the BTV LCA model closely matches each of these three conditions. The two NV models match the 0.0 condition but simulate slower response times than observed in the other two conditions. The WTV models and the BTV Fixed FFI model simulate faster response times than observed in all three conditions. For Experiment 2, all models simulate faster response times than observed in the 0.4 condition. For the other two conditions, the response times simulated by the BTV models are the closest to the observed response times. The response times simulated by the NV models are slower than the observed response times in the most difficult of these three conditions, the 0.46 condition.

Figure 18 shows the fit of each model to the proportion correct data in the conditions where exactly one coherence is equal to 0.3 in Experiment 1 and where exactly one contrast is equal to the highest contrast in Experiment 2. The NV models and the WTV LCA model provide the best fit to this subset of data from Experiment 1. The BTV LCA and WTV fixed FFI simulate a slightly lower proportion correct than observed in the 0.1 condition. The BTV fixed FFI model simulates a lower proportion correct than observed in each of the three conditions. The NV models provide the best fit to this subset of data from Experiment 2. The two WTV models and the BTV LCA model simulate a slightly lower proportion correct than observed in the 0.46 condition. The BTV fixed FFI model simulates a lower proportion correct than observed in both the 0.43 and 0.46 conditions.

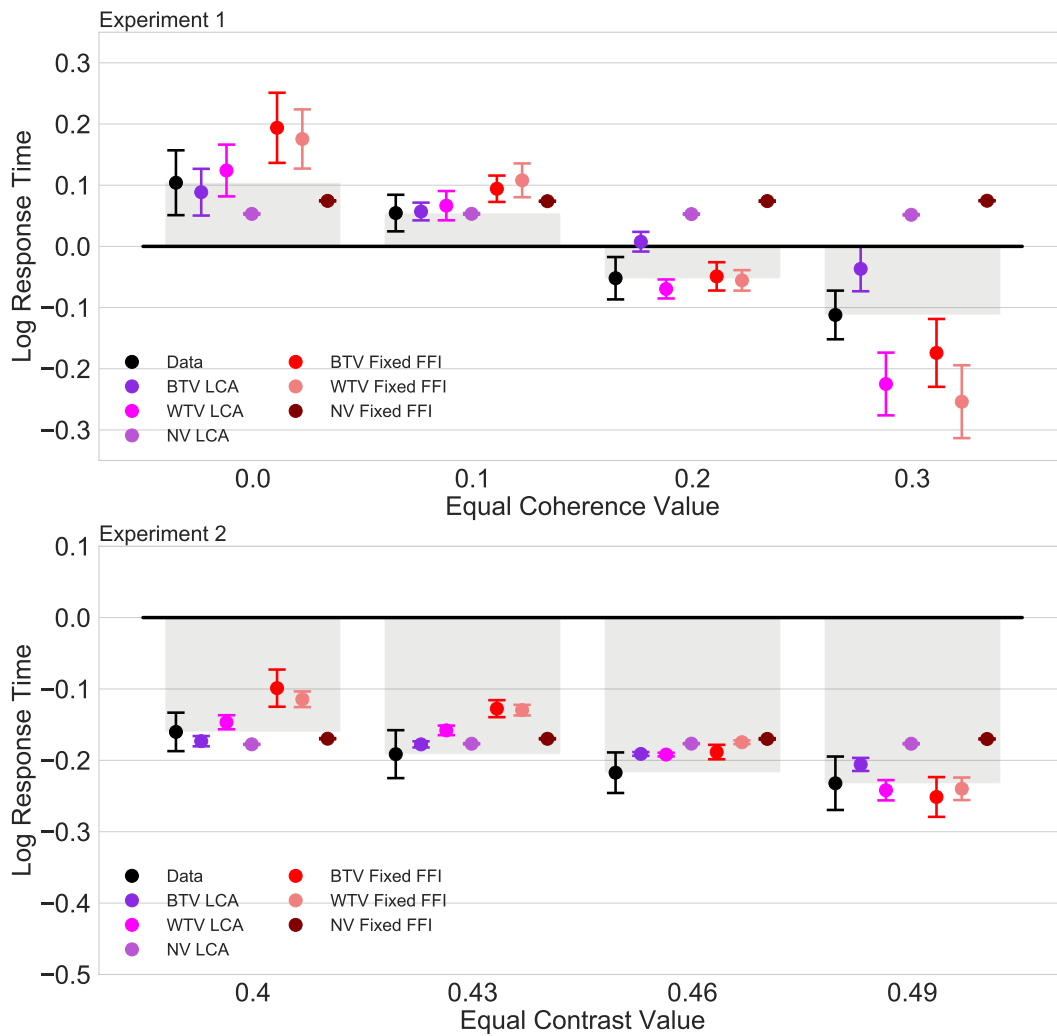


Figure 16. Simulated log response time as a function of equal-coherence or equal-contrast condition. The gray bars represent the same observed mean log response time data from Figures 4b and 7b. These response times were log transformed and then the mean was calculated between participants. Using the unique best-fitting parameter values of the participant, each model was simulated within participant to generate response time distributions for each equal-evidence grouping. Then mean log response time was calculated between participants. The error bars represent Loftus and Masson (1994) corrected 95% confidence intervals.

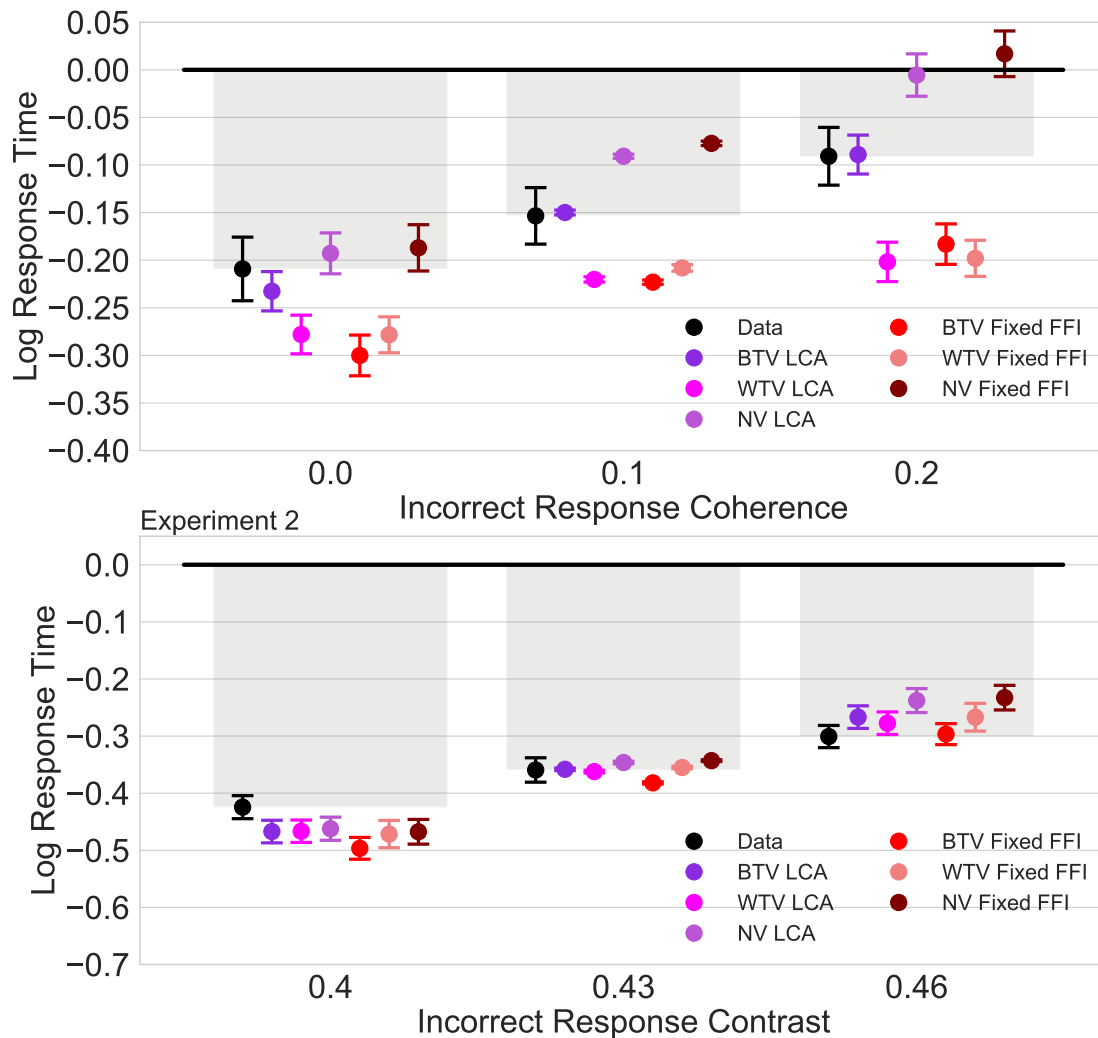


Figure 17. Simulated log response time as a function of coherence or contrast condition. The gray bars represent the same observed mean log response time data from Figures 5b and 8b. For Experiment 1, only the conditions where exactly one direction had a coherence of 0.3 are displayed. For Experiment 2, only the conditions where exactly one grating had the highest contrast of 0.49 are displayed. For the Experiment 1 results, the response times in this figure are those paired with the correct response of 0.3. For the Experiment 2 results, the response times in this figure are those paired with the correct response of 0.49. Using the unique best-fitting parameter values of the participant, each dependent and independent accumulator model was simulated within participant to generate response time distributions for each coherence or contrast grouping. Then mean log response time was calculated

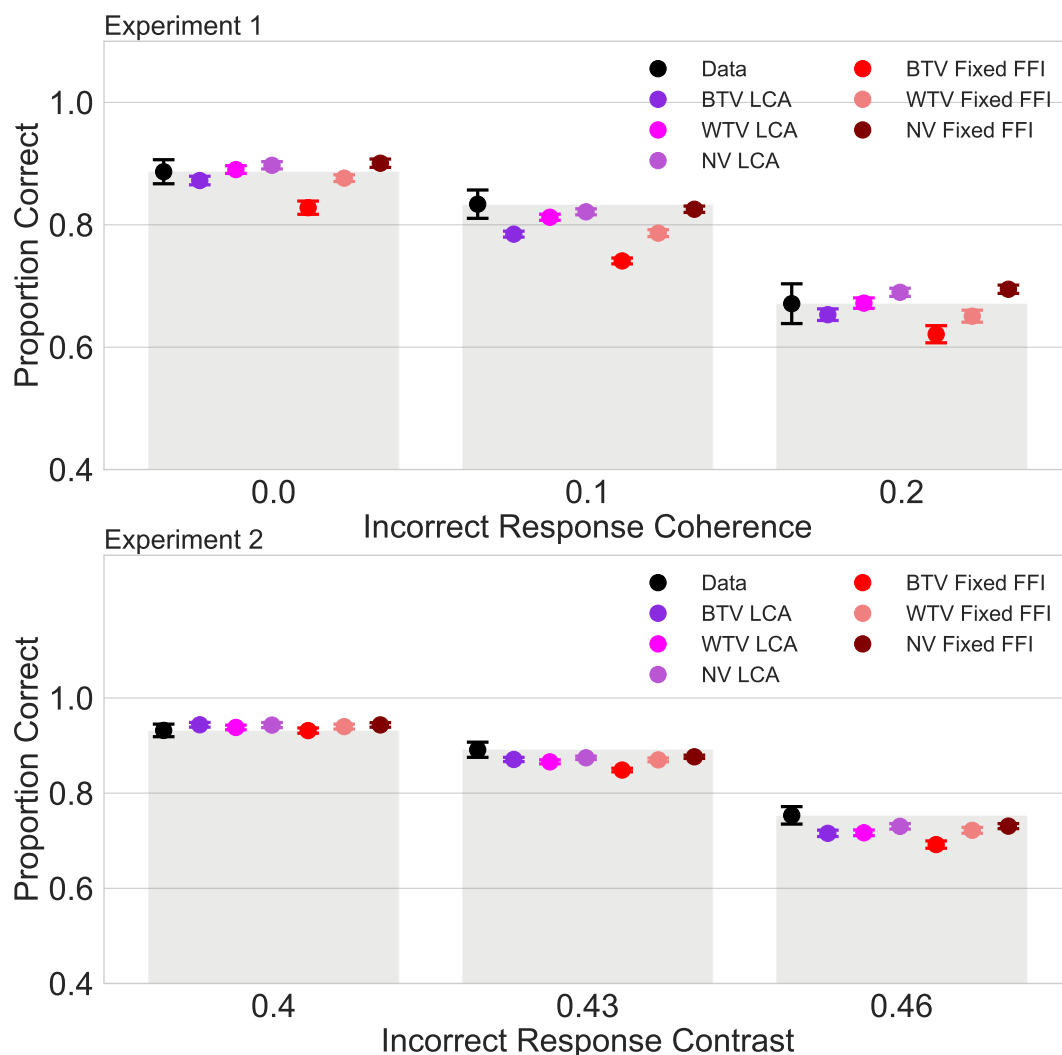


Figure 18. Simulated proportion correct as a function of coherence or contrast condition.

The gray bars represent the same observed proportion correct data from Figures 5a and 8a. In Experiment 1, only the conditions where exactly one direction had a coherence of 0.3 are displayed. In Experiment 2, only the conditions where exactly one grating had the highest contrast of 0.49 are displayed. Using the unique best-fitting parameter values of the participant, each dependent and independent accumulator model was simulated within participant to calculate the proportion correct for each coherence or contrast grouping. Then mean proportion correct was calculated between participants. The error bars represent Loftus and Masson corrected 95% confidence intervals.

## Discussion

We fit both the LCA and fixed FFI models with additional sources of variability that have been proposed by other researchers to the RDM data from Experiment 1 and to the grating data from Experiment 2. Across the two experiments, an LCA model variant had the lowest BPIC value in 36 of the 39 participants. Both how the stimulus was represented in the model and how the sources of variability were added to each model were identical for the LCA and for the fixed FFI model. This suggests that even when between-trial variability gave rise to an improved fit, lateral inhibition is a key mechanism for attaining the best fit to our datasets instead of the fully dependent mechanism of the fixed FFI model. In general, the qualitative fits of all models to most conditions was good, but the BTV and WTV fixed FFI models simulated faster response times than were observed in the highest evidence equal-coherence condition of Experiment 1 and slower response times than were observed in the lowest evidence equal-contrast condition of Experiment 2.

## General Discussion

For over half a century, the principles underlying sequential sampling theory (SST) have proven themselves as an invaluable guide for developing specific models of decision-making. The most successful extant models of choice–response time all inherit the basic architecture of SST, but make different assumptions about the type of dependency among accumulators and the role of trial-to-trial variability. The set of modifications have all proven useful under different considerations such as dynamic, time-varying information (Tsetsos et al., 2011; Usher & McClelland, 2001), statistical fluctuations from unobservable attentional sources (Franco-Watkins & Johnson, 2011; Krajbich & Rangel, 2011; Mittner et al., 2014; Turner et al., 2017; Turner, van Maanen, & Forstmann, 2015), and mathematical tractability (Brown & Heathcote, 2008; Trueblood & Heathcote, 2014; B. M. Turner et al., 2018). Given an environment where an assortment of different theoretical and practical pressures on model development have been applied, it is predictable that a range of models

have evolved that optimize for specific objectives, while still capturing key regularities in empirical data.

Each of the leading, extant models of choice–response time instantiate a particular set of mechanisms, and these combinations have proven useful in capturing patterns of behavioral data across a wide array of experiments. However, by stepping back from the particular constellations of model mechanisms, in this article, we have defined a set of three classes that can be used to group possible types of conceivable models (i.e., models that have and have not yet been realized): fully dependent, partially dependent, and independent accumulation. Although many mechanisms can be mixed to produce new model variants within these classes, the architectural constraints provide strong guidelines on how those models can ultimately behave. When the patterns of predictions defined by architecture are compared against experimental data involving exhaustive and systematic configurations of perceptual evidence embedded within the stimulus itself, strong tests emerge that can be used to provide theoretical support for types of architectures, and hence, classes of models.

A novel contribution of our task was the factorial manipulation of stimulus “coherence” for one of two options. That is, each stimulus was a mixture of evidence for alternative one (e.g., “left” response) and alternative two (e.g., “right” response), and the degree of evidence was manipulated such that stimuli could have equal amounts of evidence for both alternatives, zero evidence for both alternatives (i.e., a special case of equal evidence), or could have preferential evidence for one alternative. We used this basic factorial schematic within the classic paradigm of random dot motion, as well as contrast discrimination of grating stimuli. Both experiments exhibited strong consistency in terms of the patterns of data that emerged, providing evidence that the design itself reveals insight into decision-making processes, and the tasks are neither subject to idiosyncratic features of the stimuli nor the apparatus.

We then defined a set of models each inheriting a specific type of architectural constraint: full dependence, partial dependence, or independence. A set of 12 model variants



were examined across two sets of analyses, ultimately revealing that while various mechanisms could be used to improve fits to data within a class, the architecture of dependency imposed on the class was strong enough to provide consensus about which type of dependency best captured choice–response time data across both experiments. These analyses revealed that partial dependency was the strongest contributor to model performance, evaluated as the total number of individuals best accounted for by each model. Specifically, the mechanisms of lateral inhibition and leakage as specified in the LCA model showed clear promise in explaining human decision-making across the tasks investigated here.

Given that our results show such strong consistency for partially dependent accumulation, one may wonder why other researchers with similar curiosities either failed to find consistency or found results that favored an alternative architecture. Similar to our study, recent work has examined the fits of the LCA model and the DDM to data where the overall evidence present in the stimulus changes, but the difference in evidence between the two options is consistent (Ratcliff et al., 2018; Teodorescu et al., 2015). Teodorescu and colleagues (2015) found that an implementation of the DDM where processing noise in the model increases with increased input intensity fit 3 out of 6 participants worse than the LCA model in the choice–response time data from a two-choice fluctuating brightness discrimination task. Ratcliff and colleagues (2018) found that a similar implementation of the DDM provided a better fit in more participants than the LCA model to their fluctuating brightness discrimination task data. Contrary to these two studies, in our study, the LCA models fit to 36 of the 39 participants better than the fixed FFI model similar to the DDM implementation of the previous two experiments. This suggests that allowing either within- or between-trial variability to be a function of the magnitude of the input is insufficient for the DDM to fit these data better than the LCA model.

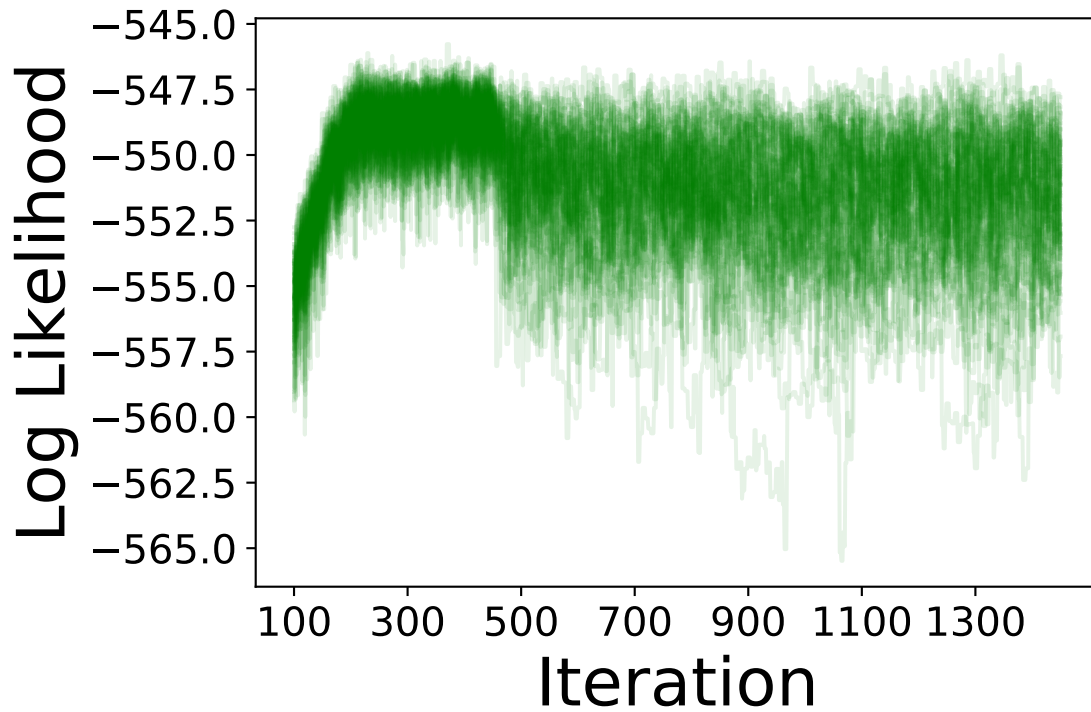
Within the category of partially dependent accumulator models, we found the LCA model to fit the data better than the free FFI model. In Model Analysis 1, the LCA model had the lower BPIC value than the free FFI model for all 39 participants across the two

experiments. In Appendix C, we show the fits of the free FFI models with the same additional sources of variability as the models in Model Analysis 2. When both the LCA and free FFI models had the same sources of variability, the LCA model fit more participants better than the free FFI model. This is evidence that the interactive competition mechanisms of the LCA provide a better fit to these data than the input competition mechanisms of the free FFI model. Competition between the stimulus values themselves alone is not sufficient to account for the patterns observed in our dataset and lateral inhibition is necessary to provide the best explanation of our data.

One explanation for the discrepancy between the results of our study and the results of previous studies could be the specific models compared differed by too many mechanisms for the key mechanism of the LCA model, lateral inhibition, to be fairly compared to the full inhibition process of the DDM model. Particularly relevant to this paper, Teodorescu and colleagues (2015) compared models similar to the WTV LCA and WTV fixed FFI models, but did not consider the between-trial variability versions of the LCA and fixed FFI models, our BTV LCA and BTV FFI, respectively. Ratcliff and colleagues (2018) compared models similar to the BTV and WTV fixed FFI models and the NV LCA model, but did not include models similar to the BTV or WTV LCA models. In our study, the BTV LCA model quantitatively fit the data better than the NV LCA model indicating that the between-trial variability mechanism specifically improved the fit of the BTV model to our datasets over that of the NV LCA model. Furthermore, the only difference between the BTV LCA and BTV fixed FFI models were the passive decay of evidence and lateral inhibition processes of the LCA model against the fully dependent input competition process of the fixed FFI model. Because the BTV LCA model fit our datasets better than both the BTV and WTV fixed FFI model, it is possible that the BTV LCA model would have fit the data from Ratcliff and colleagues (2018) better than the models they tested that were similar to the BTV and WTV fixed FFI models. However, it is interesting that even the NV LCA model outperformed the BTV and WTV variants of the fixed FFI models used here.

Given its prior success in capturing a wide range of decision-making data, it is perhaps unsurprising that the LCA model fit our data well. The LCA model has been successfully fit to data from binary perceptual judgment tasks (Usher & McClelland, 2001), numerosity judgment tasks (Usher & McClelland, 2001), value-based multi-attribute choice tasks (Bogacz et al., 2007; Turner et al., 2018; Usher & McClelland, 2004), fluctuating brightness discrimination tasks (Teodorescu et al., 2015; Teodorescu & Usher, 2013; Tsetsos et al., 2011) and other random dot motion tasks (Turner et al., 2016). Previous research has identified lateral inhibition, the key mechanism of the LCA model, in the brain during the decision-making process. Researchers have identified lateral inhibition in the inferior temporal cortex of nonhuman primates in visual search tasks in electrophysiological studies (Chelazzi, Miller, Duncan, & Desimone, 1993; Desimone, 1998; Reynolds, Chelazzi, & Desimone, 1996), in nonhuman primates in economic decision-making tasks in an electrophysiological study (Padoa-Schioppa, 2013), in humans in a flanker task in a study of event-related potentials (Gratton, Coles, Sirevaag, Eriksen, & Donchin, 1988), and in humans in an intertemporal choice task which found correlations between the lateral inhibition parameters in the model with the BOLD response in the dorsomedial frontal cortex, right and left dorsolateral prefrontal cortex, and right posterior parietal cortex (Turner et al., 2018).

Despite the success of the LCA model, one reason researchers have been hesitant to fit this model to their data may be because it does not have an analytic likelihood function, which means it requires simulations that are potentially computationally intensive (Trueblood & Heathcote, 2014; Turner et al., 2018; Turner & Sederberg, 2014). Even though the LCA model does not have an analytic likelihood function, the LCA parameters have been determined to be recoverable if a DE-MCMC sampler is used with the PDA method to fit to either a sufficiently constrained or sufficiently large dataset (Miletić, Turner, Forstmann, & van Maanen, 2017; Turner, 2019). We believe our fitting procedure and dataset satisfy these criteria. In Figure 19, we show the DE-MCMC chains for the LCA



*Figure 19. DE-MCMC chains for the LCA model fit to the data of a sample participant.*

This figure shows each of the 80 chains from the DE-MCMC and PDA algorithm, which was implemented to fit the LCA model to a sample participant of Experiment 1 in Model Analysis 1. For the clarity of the figure, the first 100 iterations of the process are not shown. At iteration 400, we switched over from burnin mode to sampling mode. After switching over from sampling mode, we waited for 200 additional samples before considering the chains as being a true sample of the posterior distribution.

model fit to the data of one of the participants from Experiment 1. From this figure, it can be observed that proper mixing is occurring such that we can be confident that the posterior distributions will be well-estimated. In Appendix A, we show a successful recovery of the parameters of the BTV LCA model. Each of the generating BTV LCA model parameters was within the respective posterior distribution. To overcome the computational intensity of fitting the LCA model to data, our model fitting code was written such that it could be parallelized on a graphics processing unit (GPU). As a result, it only took 1 to 2 hours of compute time for each participant to produce full posterior distributions of the model parameters. Hence, we believe that computational concerns are no longer reasonable justifications for excluding the LCA from serious investigations of human decision-making.

Another possible reason for the discrepancy between our results and other researchers is due strictly to methodological concerns. In nearly all previous applications, quantiles have been used to reduce the full set of data down to a set of 10 summary statistics. Models are then fit to these 10 summary statistics by minimizing the distance between the model predictions and the observed data. Although reducing the full choice–response time distribution down to a set of summary statistics clearly reduces the computational burden of fitting a model to data, the benefit also comes with an important cost. Namely, reducing the data to summary statistics will necessarily reduce the informativeness of said data unless the summary statistics are sufficient for the model parameters (Palestro et al., 2018; Turner & Sederberg, 2014; Turner & Van Zandt, 2012; B. M. Turner & Van Zandt, 2018). However, for simulation-based models such as the ones used here, demonstrating whether or not summary statistics are sufficient is difficult, if not impossible. To provide some indication of whether or not statistics are sufficient for models of choice–response time, one can compare estimated posterior distributions under two conditions: fitting the model using quantiles, and fitting the model using the full set of data. If the two posteriors closely align, then the set of summary statistics could be declared as jointly sufficient for the parameters of the model under consideration (Molloy, Galdo, Bahg, Liu, & Turner, 2019). Turner and

Sederberg (2014) demonstrated using the LBA model (which has a tractable likelihood function) that quantiles are not jointly sufficient for the LBA parameters, and thus do not convey the same information as the full choice–response time distribution when fitting the model to data. Given the consistency between the parameters of the DDM and LBA (Donkin, Brown, Heathcote, & Wagenmakers, 2011; Rodriguez, Turner, & McClure, 2014), we believe that quantiles are also insufficient for conveying the full granularity of information in the data to the parameters of other models such as the DDM and LCA, which will clearly affect conclusions about model parameters and relative model fits to data. Our study fit 12 models sampled from the three architectures of accumulator dependency to the full response time distributions for each participant individually. Both the level of constraint provided by our tasks and the methods we applied to fit the models to our data provide what we consider to be an unbiased and comprehensive assessment of the strengths and weaknesses of the examined accumulator models, focused on the theoretical question of the nature of choice dependency.

There are models we did not examine in this study that may have the mechanisms necessary to fit this dataset. The Advantage LBA model of van Ravenzwaaij, Brown, Marley, & Heathcote (2019) is structured such that both the sum of the stimulus values and the difference between stimulus values are preprocessing steps of the evidence accumulation process. It is possible that the sum of the stimulus values would allow the model to simulate the faster response times observed as the total evidence present in the stimulus increases (including in the equal evidence conditions) and the difference between the stimulus values would allow the model to simulate the slower response times observed as the difference in evidence supporting each side of the stimulus decreases. The full normalization model of Louie et al. (2013) is structured such that it continuously transitions between a relative ratio model and an independent race model. Because this model has both dependent and independent accumulation processes, it may be sufficiently flexible to capture the patterns of response times observed in both the equal-evidence and unequal-evidence conditions. A

differential relativity account may also be able to capture all of the patterns observed in our dataset (Moreno-Bote, 2010; Teodorescu et al., 2015; Zylberberg, Barttfeld, & Sigman, 2012). Differential relativity models are designed such that one parameter can manipulate the degree of dependence of the accumulators, allowing the model to flexibly switch between independent and dependent processes using input competition in a similar manner to the FFI model. The main distinction between the differential relativity models and the FFI model is that the FFI model has a neurally inspired lower boundary preventing the accumulator from activating below 0. Thus, it is reasonable to assume the two models will make similar predictions for this dataset. Although these hybrid models with both independent and dependent accumulation styles could potentially capture our data, ultimately they are nested within the architecture of partial dependency. One approach to establishing the level of dependency would be to fit a hybrid model to our data and report the estimates of the parameter that allows the hybrid model to transition from fully dependent to independent. However, another approach, used here, is to define sets of models according to the extreme classes that hybrid models transition between. Evidence for the type of dependency can then be assessed by examining relative model fits, rather than parameter estimates. Both approaches are appropriate, but as our goal was to evaluate whether or not the extreme positions of independence and full dependence were valid, we chose model evaluation metrics over parameter estimates within hybrid models (but see Osth & Dennis (2015) for an example).

## Conclusions

Sequential sampling theory has proven to be a plausible framework for investigating the dynamics of decision-making. Since its inception, many models have been developed that provide exquisite accounts of behavioral data for a handful of empirical benchmarks of decision-making. Although these models make use of a variety of different mechanisms, at their core is an explicit architectural assumption about the type of dependency that

describes the accumulation of evidence. Here we have investigated the plausibility of three types of dependency – fully dependent, partially dependent, and independent – using a combination of a fully factorial experimental design and a consortium of models, followed by model evaluation techniques. In the end, our results provide remarkable consensus that the type of dependency among choice alternatives is neither independent nor fully dependent. Instead, partial dependency, implemented via lateral inhibition between the accumulators, best accounts for human decision-making, and the degree of dependency is an important source of variation among individuals.

### **Appendix A: LCA Parameter Recovery**

In this section, we provide evidence that the LCA parameters can be recovered when the LCA model is fit to data collected from our experimental paradigms. We selected participant 9 from Experiment 1 and used the best fitting parameters obtained from the fit of the BTV LCA to choice–response time data participant 9 produced to generate simulated data. Consistent with the amount of trials participants completed in Experiment 1, we simulated 48 trials for each of the 10 conditions for a total of 480 trials of simulated data. Using the same DE-MCMC and PDA methods as Model Analyses 1 and 2, we fit the BTV LCA to the simulated data. Figure 20 shows the results of this BTV LCA parameter recovery. The posterior distributions of every parameter from the fit of BTV LCA to the simulated data include the true, generating parameters. This is initial evidence that the parameters from the BTV LCA model can be recovered when the model is fit to data from our tasks.

### **Appendix B: Cross Validation of the LCA model**

To assess whether the LCA model variants from Model Analysis 2 could be considered overparameterized as compared to the LCA model from Model Analysis 1, we performed a cross validation analysis. If the models from Model Analysis 2 had overfit the data, this cross validation analysis would show that the 4 drift-rate (DR) LCA (the model from Model Analysis 1) fit the data better than the other models (contrary to what we found in our



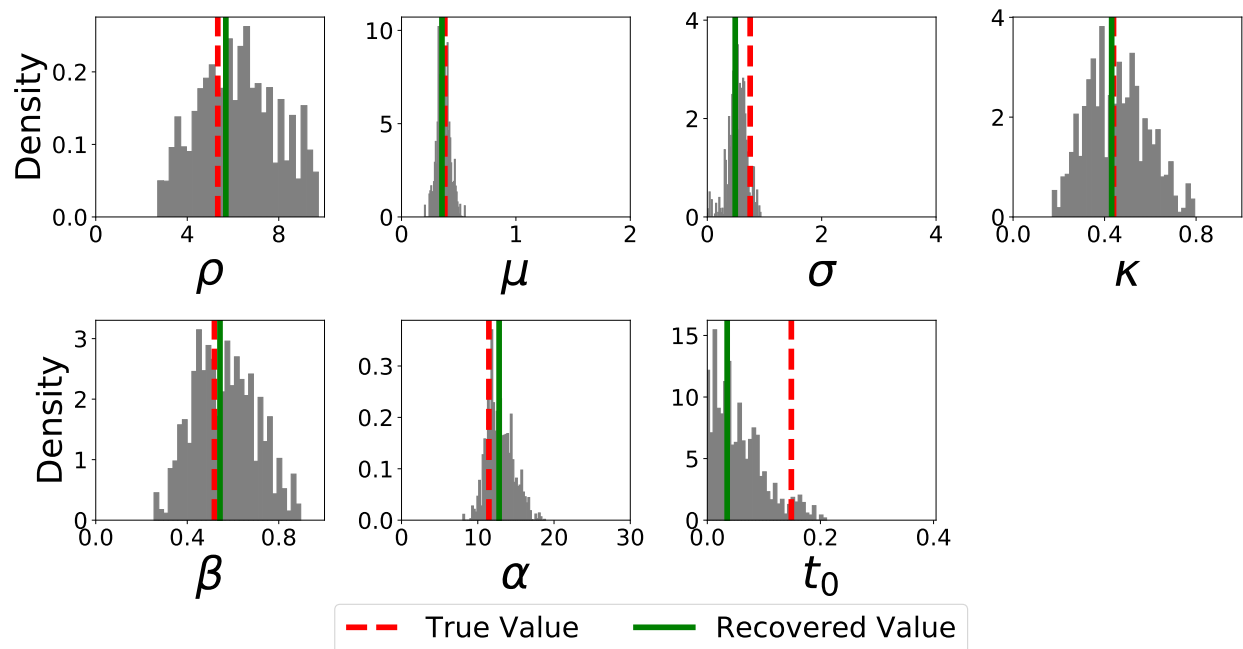


Figure 20. Recovery of BTV LCA parameter values. We fit the BTV LCA model to simulated data generated by the same model to determine if the parameters of the BTV LCA model are recoverable. The red lines represent the parameter values that we used to generate the simulated data, the gray bars represent the posterior distributions for each parameter, and the green lines represent the best fitting parameter of the fits to the simulated data. The lower and upper bounds on each x-axis are the same values as the lower and upper bounds of the prior distributions we specified for each parameter except for  $\mu$  and  $\sigma$  where we restricted the bounds on the x-axis to a narrower range for greater visibility of the results.

initial analyses). Conversely, if the models from Model Analysis 2 had not overfit the data, the analysis would show the same models that fit the best in the initial fitting procedure to fit the best in the cross validation analysis.

For this analysis, we chose two participants from Experiment 1, which we will call participant A and participant B here. In our initial analyses, we found that the BTV LCA model had the lowest BPIC value of the four LCA models for participant A and the WTV LCA had the lowest BPIC value for participant B. So the BTV LCA model fit participant A

better than the other three models and the WTV LCA model fit participant B better than the other three models. For each participant, we specified 10 pairs of training and testing datasets where the training set was all of the data from Experiment 1 for that participant save one condition (i.e. the 0.0-0.2 coherence condition) and the testing set consisted of the data from the one condition that was not included in the training data. Using the same fitting methods as Model Analyses 1 and 2, we then fit each of our four LCA models (BTV LCA, WTV LCA, NV LCA, and 4 DR LCA) to each of the 20 training datasets (10 for each participant) to obtain full posterior distributions for each model and dataset. Next, we sampled from each posterior distribution and resimulated the respective model using those samples from the testing data conditions. Using the simulated data, we calculated the log likelihood value for each posterior sample using the kernel density estimation procedure and combined these samples together to create a distribution of log likelihood values.

The BPIC values in each subfigure of Figures 21 and 22 are calculated from the distributions of log likelihood values plus the log prior for each testing dataset with only one condition removed. So, for example, to generate the blue BTV LCA point in the subfigure with the heading 0.0, 0.0 for Participant A, we sampled from the posterior distribution generated when we fit to the training data and then resimulated the BTV LCA model 50,000 times using 0.0 as for both coherence values for each of those posterior samples. For each of those 50,000 simulations, we created a probability distribution using a kernel density estimate with Silverman's rule of thumb, took the log of each value in the probability distribution, and then summed the log probability values together to get the log likelihood. We then added the log likelihood to the log prior for that posterior sample. This was repeated for each of our posterior samples to form a distribution of weights which was used to calculate the BPIC

We calculated the BPIC for each of the four models for each participant and testing data set. The results of our analysis are displayed in Figures 21 and 22. For the datasets of participant A where data from one condition were the testing data, the NV LCA had the lowest BPIC value for 7 of the 10 conditions, the 4 DR LCA had the lowest BPIC value for 2

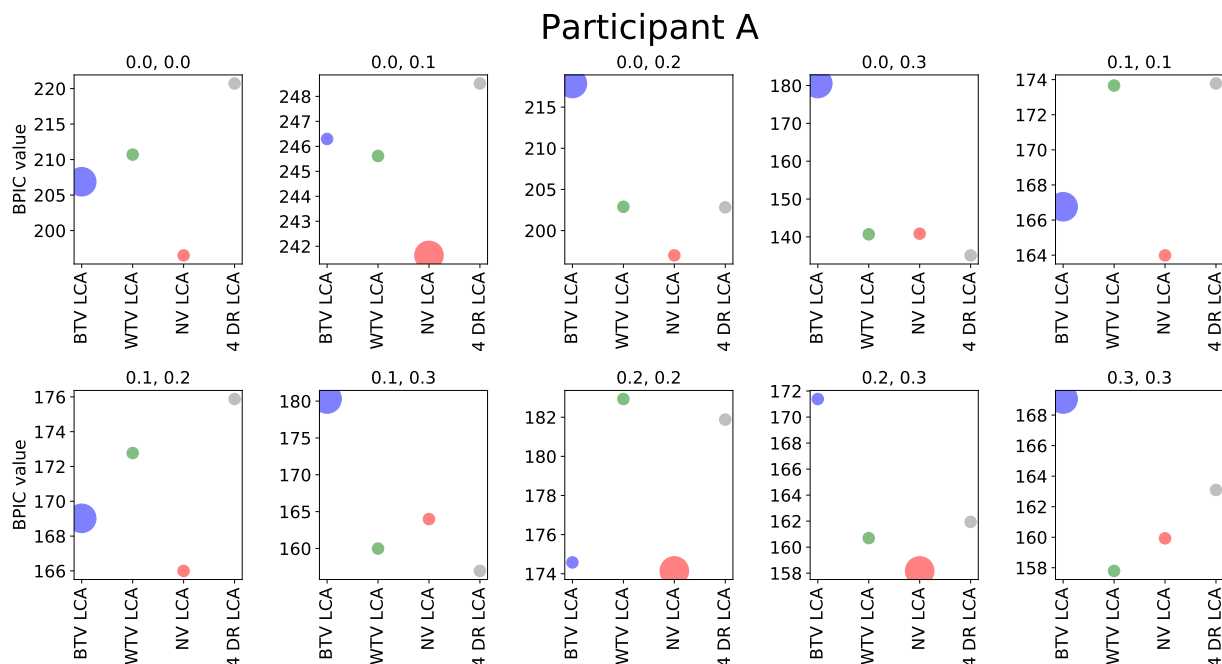


Figure 21. Results of the cross validation of the LCA model for participant A. The details of the procedure used to generate these results are found in the text of this section. The heading of each subfigure indicates the condition whose trials formed the testing set. Each dot represents the BPIC value for one model for one set of testing data. The larger dot in each figure indicates which model had the lowest BPIC score when fit to the training data.

of the 10 conditions, and WTV LCA had the lowest BPIC for 1 condition. For the datasets of participant B where data from one condition were the testing data, the WTV LCA and NV LCA had the lowest BPIC value for 3 conditions while the BTV LCA and 4 DR LCA had the lowest BPIC value for 2 conditions. We also calculated the mean BPIC across the 10 single condition testing data sets for each model. For participant A, NV LCA had a mean BPIC of 176.2, WTV LCA had a mean BPIC of 180.7, 4 DR LCA had a mean BPIC of 182.0, and BTV LCA had a mean BPIC of 188.2. For participant B, 4 DR LCA had a mean BPIC of 158.1, WTV LCA had a mean BPIC of 158.8, BTV LCA had a mean BPIC of 160.7, and NV LCA had a mean BPIC of 166.3.

None of the four models systematically had the poorest fit for any subset of training

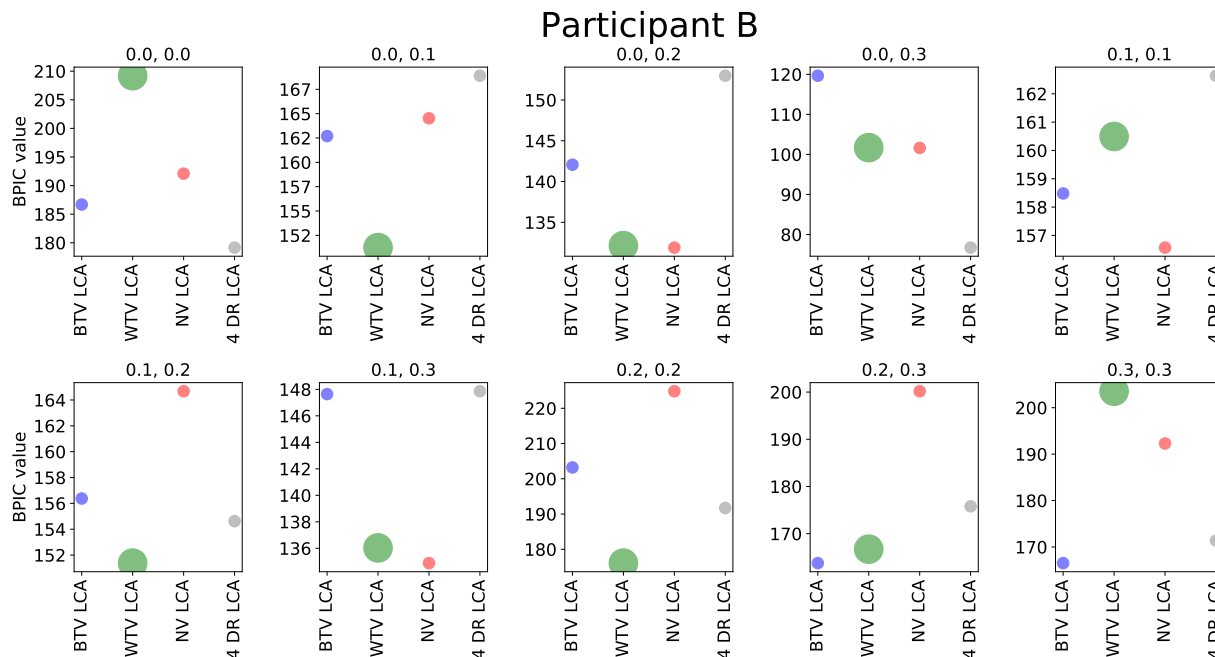


Figure 22. Results of the cross validation of the LCA model for participant B. The details of the procedure used to generate these results are found in the text of this section. The heading of each subfigure indicates the condition whose trials formed the testing set. Each dot represents the BPIC value for one model for one set of testing data. The larger dot in each figure indicates which model had the lowest BPIC score when fit to the training data. This speaks to the validity of our approach in that the extra mechanisms given to the WTV LCA and BTV LCA models provide additional explanatory power to the model without the cost of overfitting to non-psychological processes.

**Appendix C: Analysis of the free FFI model with additional sources of variability and the LCA models without the stimulus difference in the drift rate computation**

Model Analysis 2 showed that the LCA model benefited from additional sources of variability which suggests that other partially dependent models could potentially fit our data better than the LCA model if given additional sources of variability. We fit one such

model, the free FFI model from Model Analysis 1, with the drift rate representation and additional sources of variability from Model Analysis 2. We briefly describe these three additional free FFI model variants in the following equations.

The following equations describe the NV free FFI model for 2-alternative forced-choice tasks:

$$m_1 = \mu(I_1 - I_2) \frac{dt}{\tau} + \xi_1 \sqrt{\frac{dt}{\tau}}, \quad (60)$$

$$m_2 = \mu(I_2 - I_1) \frac{dt}{\tau} + \xi_2 \sqrt{\frac{dt}{\tau}}, \quad (61)$$

$$dx_1 = \rho \frac{dt}{\tau} + m_1 - \nu m_2, \quad (62)$$

$$dx_2 = \rho \frac{dt}{\tau} + m_2 - \nu m_1, \quad (63)$$

$$\xi_1 \sim N(0, \eta), \quad (64)$$

$$\xi_2 \sim N(0, \eta), \quad (65)$$

where  $\nu$  is the FFI parameter and the other parameters represent the same processes as the NV fixed FFI model.

The following equations describe the BTV free FFI model for 2-alternative forced-choice tasks:

$$m_1 = \mu(I_1 - I_2), \quad (66)$$

$$m_2 = \mu(I_2 - I_1), \quad (67)$$

$$s_1 = \sigma I_1, \quad (68)$$

$$s_2 = \sigma I_2, \quad (69)$$

$$d_1 \sim N(m_1, s_1), \quad (70)$$

$$d_2 \sim N(m_2, s_2), \quad (71)$$

$$dr_1 = d_1 \frac{dt}{\tau} + \xi_1 \sqrt{\frac{dt}{\tau}}, \quad (72)$$

$$dr_2 = d_2 \frac{dt}{\tau} + \xi_2 \sqrt{\frac{dt}{\tau}}, \quad (73)$$

$$dx_1 = \rho \frac{dt}{\tau} + dr_1 - \nu dr_2, \quad (74)$$

$$dx_2 = \rho \frac{dt}{\tau} + dr_2 - \nu dr_1, \quad (75)$$

$$\xi_1 \sim N(0, \eta), \quad (76)$$

$$\xi_2 \sim N(0, \eta), \quad (77)$$

where  $\nu$  is the FFI parameter and the other parameters represent the same processes as the BTV fixed FFI model.

The WTV free FFI model is described by the following equations for 2-alternative forced-choice tasks:

$$m_1 = \mu(I_1 - I_2) \frac{dt}{\tau} + \xi_1 \sqrt{\frac{dt}{\tau}}, \quad (78)$$

$$m_2 = \mu(I_2 - I_1) \frac{dt}{\tau} + \xi_2 \sqrt{\frac{dt}{\tau}}, \quad (79)$$

$$dx_1 = \rho \frac{dt}{\tau} + m_1 - \nu m_2, \quad (80)$$

$$dx_2 = \rho \frac{dt}{\tau} + m_2 - \nu m_1, \quad (81)$$

$$\xi_1 \sim N(0, \eta + \nu I_1), \quad (82)$$

$$\xi_2 \sim N(0, \eta + \nu I_2), \quad (83)$$

where  $\nu$  is the FFI parameter and the other parameters represent the same processes as the WTV fixed FFI model.

To make the comparison as fair as possible, in Model Analysis 2, we originally compared the LCA model to the fixed FFI with the same representation of drift rate as the difference between the normalized stimulus values. Model Analysis 1 provided evidence that this mechanism was likely unnecessary for the LCA model to fit our data well. Thus, in addition, we fit each LCA variant from Model Analysis 2 where the drift rate for each accumulator was a function of only the stimulus value represented by that accumulator instead of the difference between the stimulus values. We briefly describe these “no stimulus difference” (ND) LCA variants in the following equations.

The following equations describe the NV LCA ND model for 2-alternative forced-choice tasks:

$$m_1 = \mu I_1 \frac{dt}{\tau} + \xi_1 \sqrt{\frac{dt}{\tau}}, \quad (84)$$

$$m_2 = \mu I_2 \frac{dt}{\tau} + \xi_2 \sqrt{\frac{dt}{\tau}}, \quad (85)$$

$$dx_1 = m_1 + (\rho - \kappa x_1 - \beta x_2) \frac{dt}{\tau}, \quad (86)$$

$$dx_2 = m_2 + (\rho - \kappa x_2 - \beta x_1) \frac{dt}{\tau}, \quad (87)$$

$$\xi_1 \sim N(0, \eta), \quad (88)$$

$$\xi_2 \sim N(0, \eta), \quad (89)$$

where all parameters represent the same processes as the NV LCA model.

The following equations describe the BTV LCA ND model for 2-alternative forced-choice tasks:

$$m_1 = \mu I_1, \quad (90)$$

$$m_2 = \mu I_2, \quad (91)$$

$$s_1 = \sigma I_1, \quad (92)$$

$$s_2 = \sigma I_2, \quad (93)$$

$$d_1 \sim N(m_1, s_1), \quad (94)$$

$$d_2 \sim N(m_2, s_2), \quad (95)$$

$$dr_1 = d_1 \frac{dt}{\tau} + \xi_1 \sqrt{\frac{dt}{\tau}}, \quad (96)$$

$$dr_2 = d_2 \frac{dt}{\tau} + \xi_2 \sqrt{\frac{dt}{\tau}}, \quad (97)$$

$$dx_1 = dr_1 + (\rho - \kappa x_1 - \beta x_2) \frac{dt}{\tau}, \quad (98)$$

$$dx_2 = dr_2 + (\rho - \kappa x_2 - \beta x_1) \frac{dt}{\tau}, \quad (99)$$

$$\xi_1 \sim N(0, \eta), \quad (100)$$

$$\xi_2 \sim N(0, \eta), \quad (101)$$

where all parameters represent the same processes as the BTV LCA model.

The following equations describe the WTV LCA ND model for 2-alternative forced-choice tasks:

$$m_1 = \mu I_1 \frac{dt}{\tau} + \xi_1 \sqrt{\frac{dt}{\tau}}, \quad (102)$$

$$m_2 = \mu I_2 \frac{dt}{\tau} + \xi_2 \sqrt{\frac{dt}{\tau}}, \quad (103)$$

$$dx_1 = m_1 + (\rho - \kappa x_1 - \beta x_2) \frac{dt}{\tau}, \quad (104)$$

$$dx_2 = m_2 + (\rho - \kappa x_2 - \beta x_1) \frac{dt}{\tau}, \quad (105)$$

$$\xi_1 \sim N(0, \eta + v I_1), \quad (106)$$





<i>Model</i>	$\rho$	$\mu$	$\sigma$	$v$	$\nu$	$\kappa$	$\beta$	$\alpha$	$t_0$
WTV	0.029	0.294	-	0.855	0.144	-	-	4.573	0.123
Free	(0.071)	(0.187)		(0.409)	0.11			(1.088)	(0.123)
FFI									
NV	0.022	0.329	-	-	0.152	-	-	4.355	0.119
Free	(0.071)	(0.174)			(0.154)			(1.032)	(0.117)
FFI									

Table 8

*Mean best fitting parameter values calculated between participants for each model fit to the data from Experiment 2. The standard deviation is given in parentheses.*

<i>Model</i>	$\rho$	$\mu$	$\sigma$	$v$	$\nu$	$\kappa$	$\beta$	$\alpha$	$t_0$
BTV	6.146	0.797	0.088	-	-	0.499	0.688	9.098	0.147
LCA	(1.952)	(0.243)	(0.103)			(0.22)	(0.194)	(2.708)	(0.108)
ND									
WTV	6.456	0.782	-	0.278	-	0.443	0.64	10.26	0.129
LCA	(1.782)	(0.23)		(0.082)		(0.18)	(0.178)	(3.327)	(0.105)
ND									
NV	5.678	0.781	-	-	-	0.42	0.619	9.74	0.132
LCA	(1.849)	(0.218)				(0.21)	(0.195)	(3.412)	(0.102)
ND									
BTV	0.094	0.58	0.28	-	0.104	-	-	3.547	0.212
Free	(0.138)	(0.219)	(0.153)		(0.063)			(0.606)	(0.102)
FFI									
WTV	0.077	0.509	-	0.41	0.105	-	-	3.643	0.199
Free	(0.127)	(0.173)		(0.126)	0.087			(0.614)	(0.095)

<i>Model</i>	$\rho$	$\mu$	$\sigma$	$\nu$	$\nu$	$\kappa$	$\beta$	$\alpha$	$t_0$
FFI									
NV	0.074	0.505	-	-	0.093	-	-	3.531	0.196
Free	(0.131)	(0.155)			(0.066)			(0.617)	(0.099)
FFI									

The means and standard deviations of the best fitting parameter values calculated between participants for each model are given in Tables 7 and 8 for Experiments 1 and 2, respectively. Figure 23 displays the mean-centered BPIC values for each of the 6 model variants introduced in Model Analysis 2 and the 6 model variants introduced in this Appendix for each participant from Experiment 1 and 2. During the model fitting process, the WTV free FFI model had extremely poor mixing for one participant from Experiment 1 due to fast outlier response times between 0.2 and 0.3 seconds. The poor mixing hindered our ability to use BPIC values to discriminate between the models even though the best fitting parameter values suggested a reasonable fit for the WTV free FFI model. Thus, we re-fit all 12 models for this one participant for illustration purposes where (after first removing all response times less than 0.2 seconds and greater than 5.0 seconds) we iteratively removed trials where the response times were outside 3 standard deviations of the mean (~3% of all trials). This outlier correction resulted in significantly improved mixing for the WTV free FFI model for this one participant and subsequently removed the distortion from the BPIC figure.

In Experiment 1, the BTV LCA had the lowest BPIC value in 3 participants, the NV LCA had the lowest BPIC value in 7 participants, the NV LCA ND had the lowest BPIC value in 3 participants, the BTV LCA ND had the lowest BPIC value in 1 participant, the NV free FFI had the lowest BPIC value in 1 participant, and the WTV fixed FFI had the lowest BPIC in 1 participant. Overall, the LCA models had the lowest BPIC values for 14 of the 16 participants. In Experiment 2, the BTV LCA had the lowest BPIC value in 8

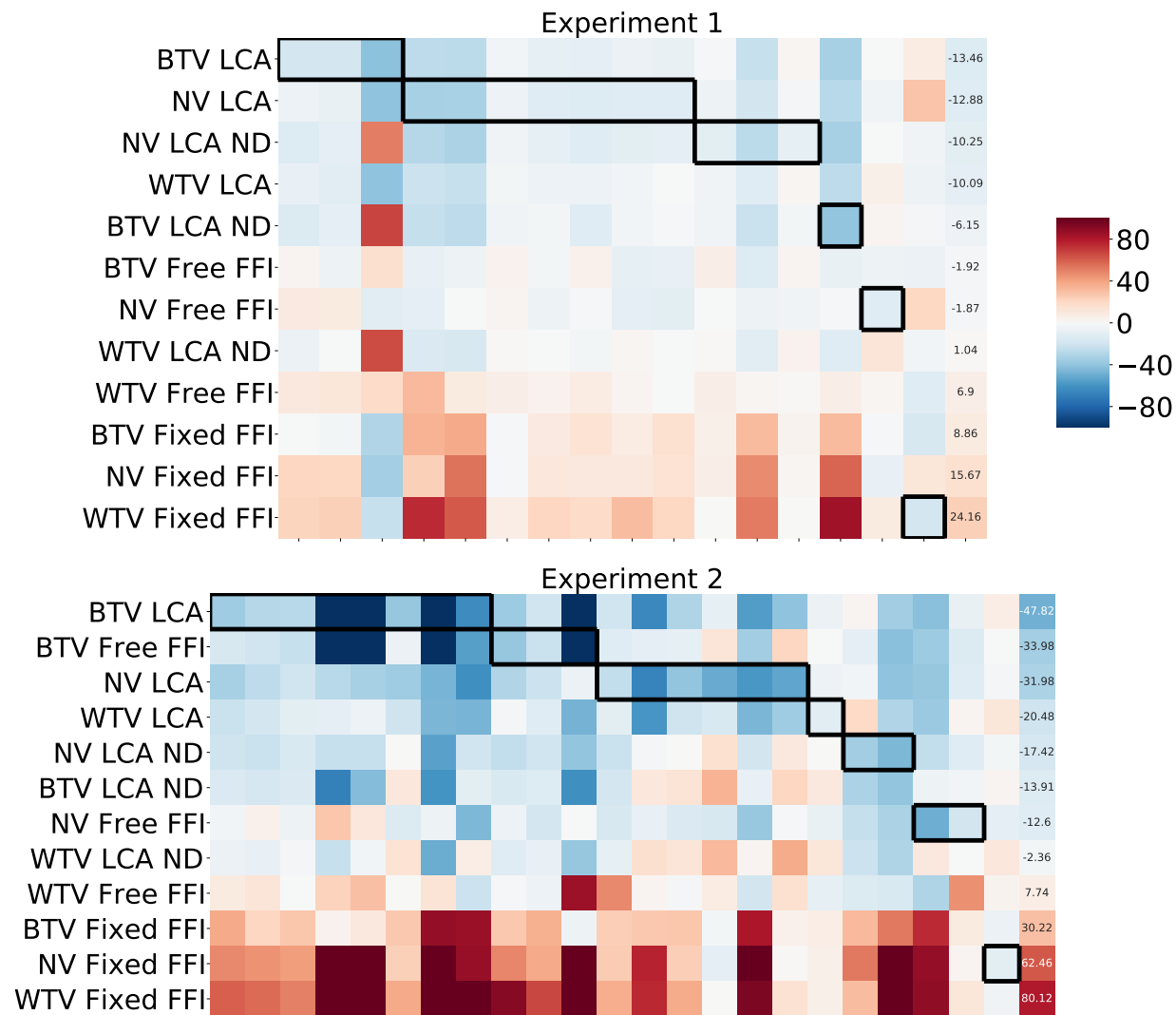


Figure 23. Heatmap showing mean-centered BPIC values for each model and for each participant. This figure is organized by participant and model variants with each column representing a single participant and each row representing a single model variant. Each square in the figure represents the mean-centered BPIC for that particular model with the mean calculated across the 12 model variants for the one participant. Cooler colors represent lower (preferred) BPIC values and warmer colors represent higher BPIC values. The squares outlined by the black line represent the model variants with the lowest BPIC value of the 12 variants. The final column represents the mean BPIC value for the model calculated across participants. The model with the LCA architecture had the numerically lowest BPIC value for 14 of the 16 participants in Experiment 1 and for 17 of the 23 participants in Experiment 2.

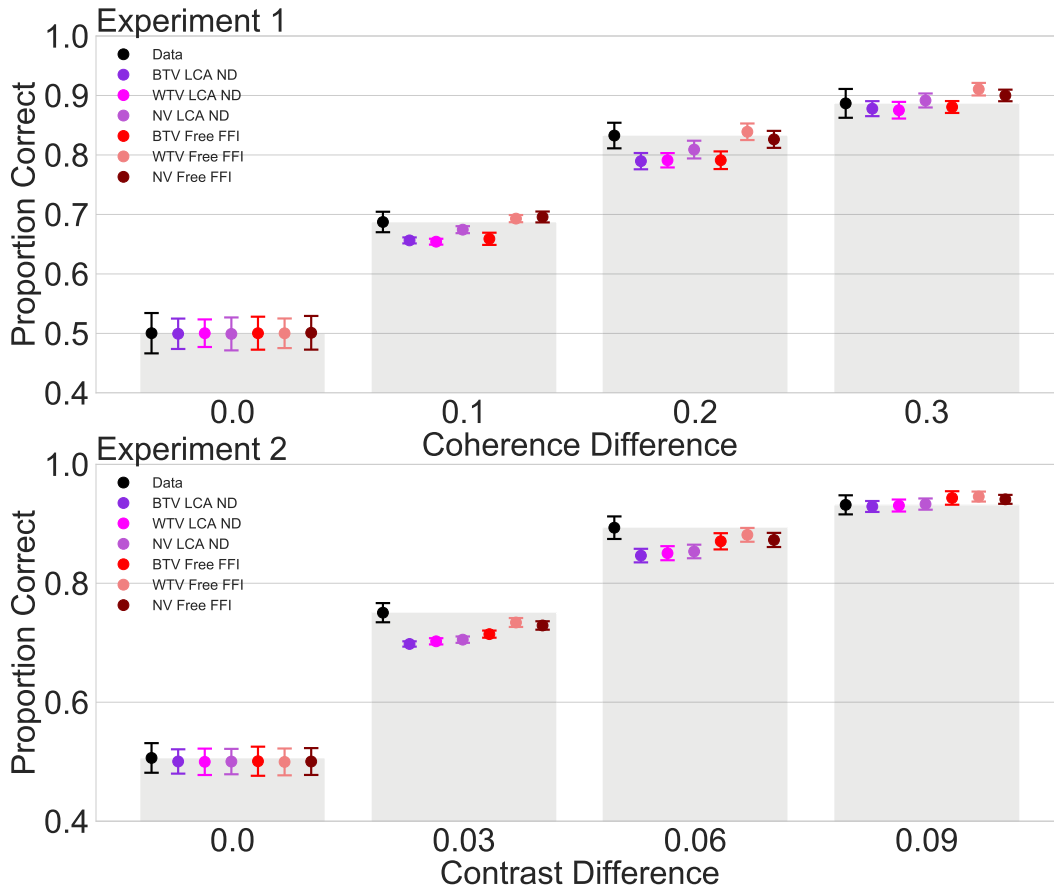


Figure 24. Simulated proportion correct as a function of coherence difference or contrast difference. The gray bars represent the same observed mean proportion correct data from Figures 4a and 7a. Using the unique best-fitting parameter values of the participant, each dependent and independent accumulator model was simulated within-participant to generate the proportion correct for each coherence grouping for the data from Experiment 1 and for each contrast grouping for the data from Experiment 2. The mean was then calculated between the participant simulations for each model. The error bars represent Loftus and Masson (1994) corrected 95% confidence intervals.

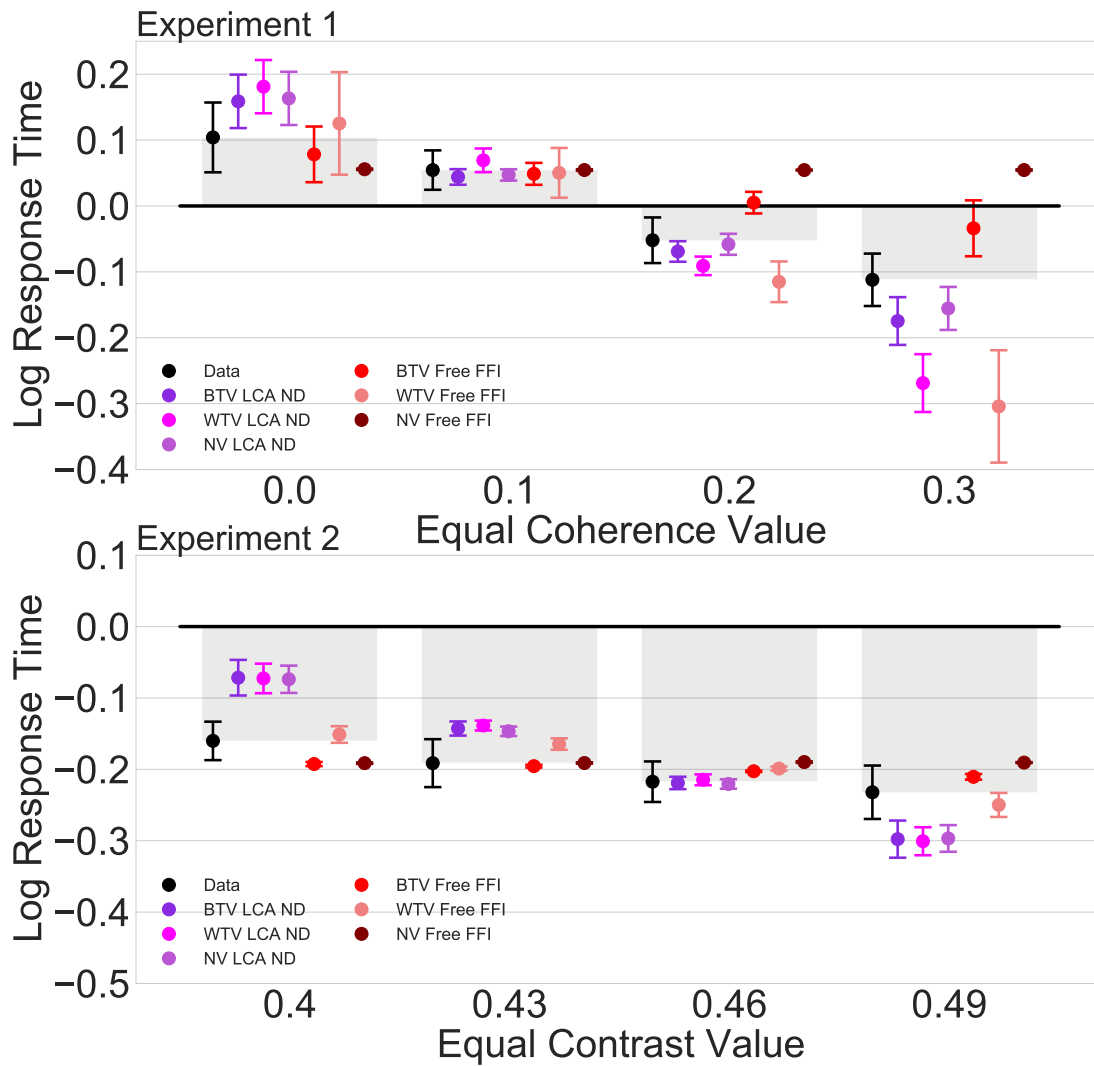


Figure 25. Simulated log response time as a function of equal-coherence or equal-contrast condition. The gray bars represent the same observed mean log response time data from Figures 4b and 7b. These response times were log transformed and then the mean was calculated between participants. Using the unique best-fitting parameter values of the participant, each model was simulated within participant to generate response time distributions for each equal-evidence grouping. Then mean log response time was calculated between participants. The error bars represent Loftus and Masson (1994) corrected 95% confidence intervals.

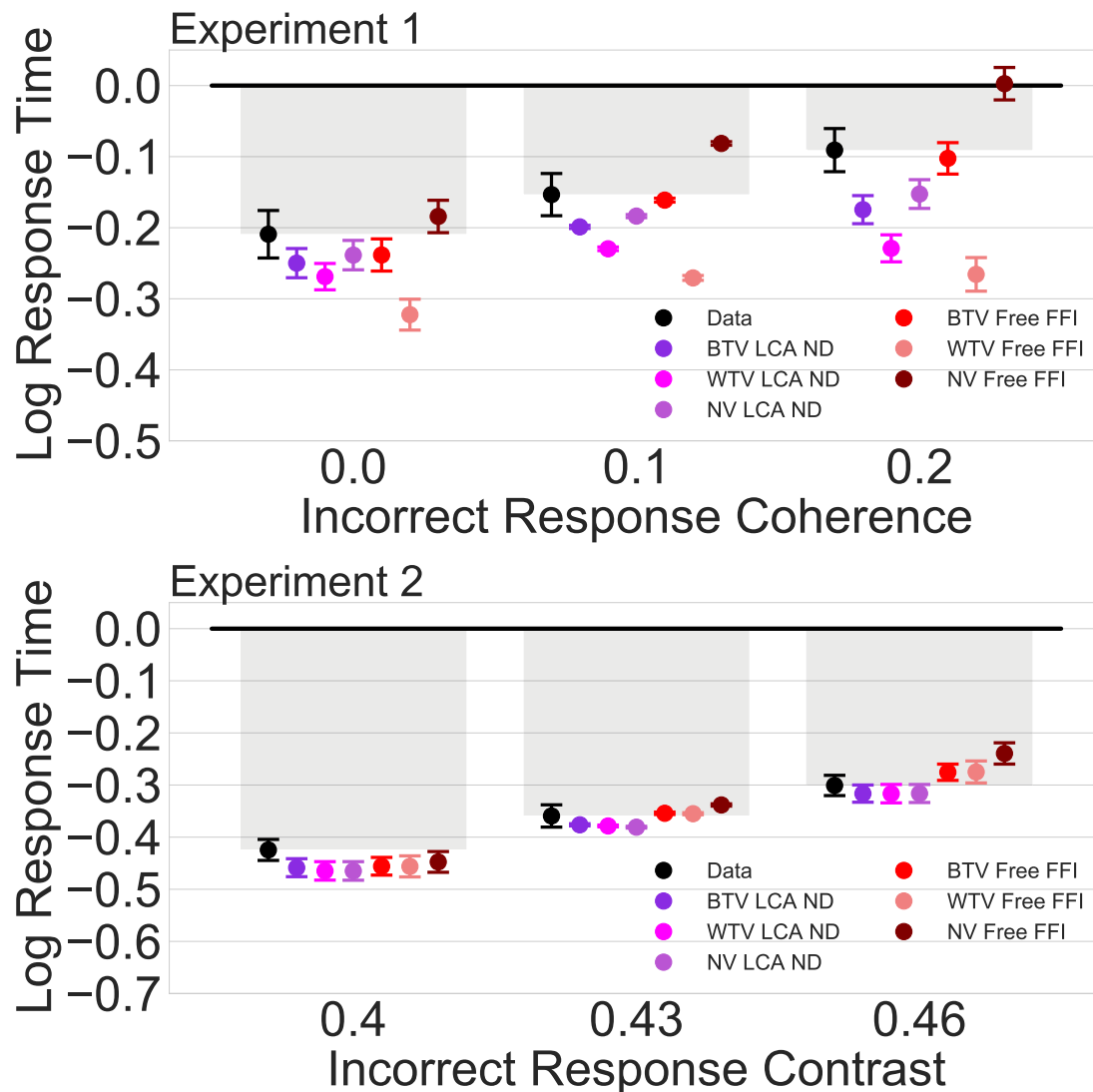


Figure 26. Simulated log response time as a function of coherence or contrast condition. The gray bars represent the same observed mean log response time data from Figures 5b and 8b. For Experiment 1, only the conditions where exactly one direction had a coherence of 0.3 are displayed. For Experiment 2, only the conditions where exactly one grating had the highest contrast of 0.49 are displayed. For the Experiment 1 results, the response times in this figure are those paired with the correct response of 0.3. For the Experiment 2 results, the response times in this figure are those paired with the correct response of 0.49. Using the unique best-fitting parameter values of the participant, each dependent and independent accumulator model was simulated within participant to generate response time distributions for each coherence or contrast grouping. Then mean log response time was calculated

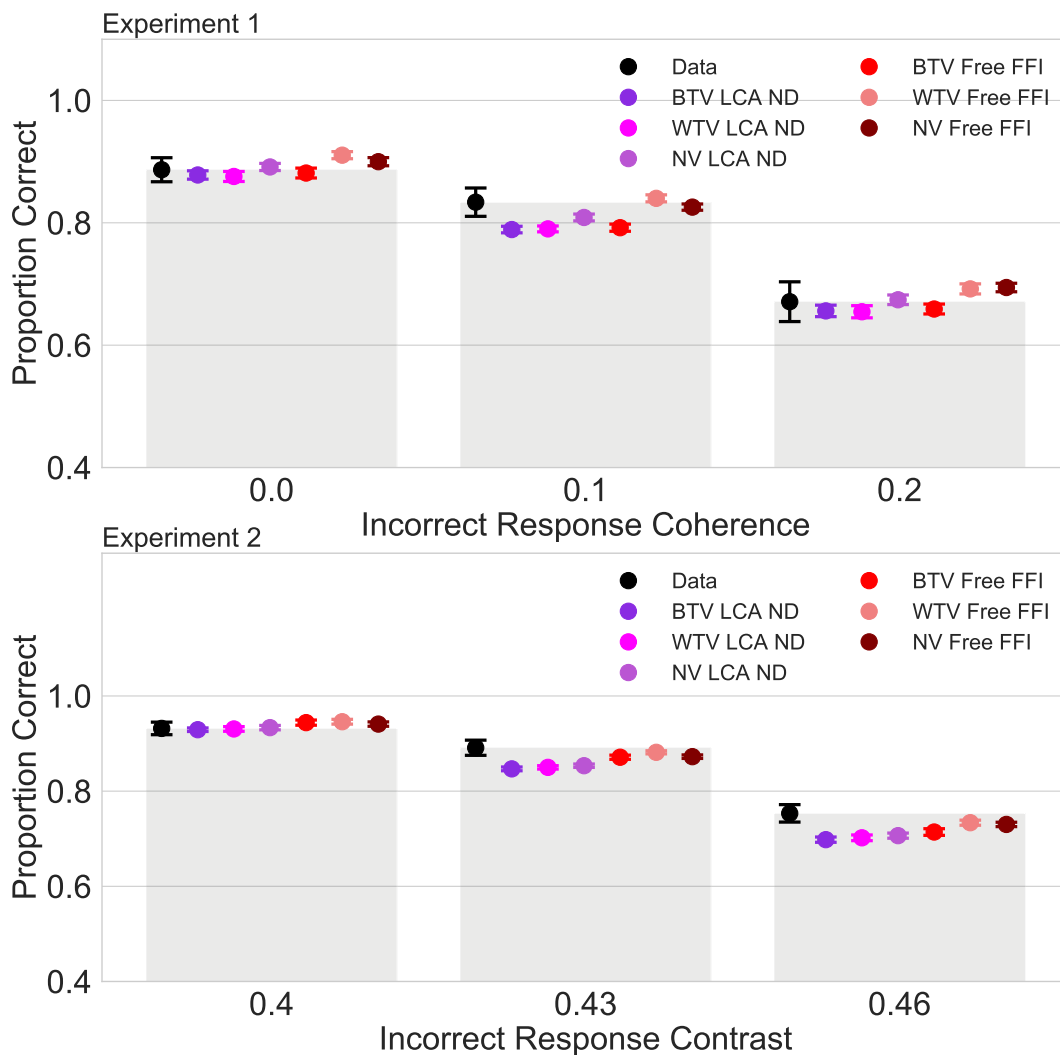


Figure 27. Simulated proportion correct as a function of coherence or contrast condition.

The gray bars represent the same observed proportion correct data from Figures 5a and 8a. In Experiment 1, only the conditions where exactly one direction had a coherence of 0.3 are displayed. In Experiment 2, only the conditions where exactly one grating had the highest contrast of 0.49 are displayed. Using the unique best-fitting parameter values of the participant, each dependent and independent accumulator model was simulated within participant to calculate the proportion correct for each coherence or contrast grouping. Then mean proportion correct was calculated between participants. The error bars represent Loftus and Masson corrected 95% confidence intervals.



participants, the BTV free FFI had the lowest BPIC value in 3 participants, the NV LCA had the lowest BPIC value in 6 participants, the WTV LCA had the lowest BPIC value in 1 participant, the NV LCA ND had the lowest BPIC value in 2 participants, the NV free FFI had the lowest BPIC value in 2 participants, and the NV fixed FFI had the lowest BPIC value in 1 participant. Overall, the LCA models had the lowest BPIC values for 17 of the 23 participants, providing more evidence that the interactive competition LCA model is the preferred model for our datasets over the input competition fixed and free FFI models, even when additional sources of variability are considered.

To assess whether the LCA models from Model Analysis 2 would be improved if the stimulus difference component of the drift rate calculation was removed, we compared the LCA models introduced in Model Analysis 2 to the LCA models introduced in this Appendix. Across the two experiments, the BTV LCA model fit 30 of 39 participants better than the BTV LCA ND model. The NV LCA model fit 26 of 39 participants better than the NV LCA ND model. Similarly, the WTV LCA fit 32 of 39 participants better than the WTV LCA ND model. Each LCA model from Model Analysis 2 provided the best fit to more participants than their counterpart introduced in this Appendix which indicates that calculating drift rate as a function of the difference in the stimulus values is an important mechanism for fitting these data.

To provide further evidence that interactive competition is a better mechanism than input competition for fitting these data, we also compared the free FFI models with additional sources of variability to the LCA models with additional sources of variability from Model Analysis 2. Across the two experiments, the BTV LCA model fit 29 of 39 participants better than the BTV free FFI model. The NV LCA model fit 31 of 39 participants better than the NV free FFI model. The WTV LCA model fit 33 of 39 participants better than the WTV free FFI model. Each LCA model from Model Analysis 2 provided the best fit to more participants than their counterpart free FFI model which indicates interactive competition is the preferred mechanism to input competition for fitting these data.

Figure 24 shows the fits of the free FFI models and the LCA ND models to the proportion correct data as a function of the difference in stimulus value for both experiments. All models generally match the proportion correct data well with the largest misses for all models appearing in the 0.03 contrast difference condition of Experiment 2.

Figure 25 shows the fits of the free FFI models and the LCA ND models to the response times of the equal evidence conditions for both experiments. For Experiment 1, the LCA ND models generally fit the data well. The biggest miss of the models is the WTV LCA ND model for the 0.3 coherence condition. The WTV and BTV free FFI model fit the 0.0 and 0.1 conditions well but miss on the 0.2 and 0.3 conditions. The NV free FFI model simulates the same response time for each condition for Experiment 1 and the same response time for each condition for Experiment 2. For Experiment 2, the LCA ND models fit the data similarly to each other. All three models miss greatly on the 0.4 and 0.49 conditions, miss slightly on the 0.43 condition and fit the 0.46 condition well. The WTV free FFI model provides the best visual fit to this subset of data of these 6 models and the BTV free FFI model fits all conditions except the 0.4 condition well.

Figure 26 shows the fits of the free FFI models and the LCA ND models to the response times of a subset of the unequal evidence conditions for both experiments. For Experiment 1, the LCA ND models generally predict faster response times than observed with the WTV LCA ND predicting the fastest response times. The BTV free FFI model provides the best fit to these data, the WTV free FFI predicts even faster response times than the response times predicted by the LCA ND models, and the NV free FFI predicts slower response times than observed in two of the conditions. For Experiment 2, all models generally provide a good fit to the data with the fit of the free FFI models being worse than the LCA ND models in the 0.46 contrast condition.

Figure 27 shows the fits of the free FFI models and the LCA ND models to the proportion correct data of a subset of the unequal evidence conditions for both experiments. For Experiment 1, all models generally fit the data well with the LCA ND models fitting the

data worse in the 0.1 coherence condition and the free FFI models fitting the data worse in the 0.2 coherence condition. For Experiment 2, all models generally fit the data well with all models underestimating the proportion correct in the 0.46 contrast condition. Taken together, the quantitative BPIC results and the visualizations of the model fits indicate that the LCA ND models generally fit the data better than the free FFI models but neither set of models fit the data as well as the LCA models from Model Analysis 2.

### References

- Ando, T. (2007). Bayesian predictive information criterion for the evaluation of hierarchical Bayesian and empirical Bayes models. *Biometrika*, *94*(2), 443–458.
- Anstis, S. M. (1980). The perception of apparent movement. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, *290*(1038), 153–168.
- Bogacz, R., Brown, E., Moehlis, J., Holmes, P., & Cohen, J. D. (2006). The physics of optimal decision making: A formal analysis of models of performance in two-alternative forced-choice tasks. *Psychological Review*, *113*(4), 700–765.
- Bogacz, R., Usher, M., Zhang, J., & McClelland, J. L. (2007). Extending a biologically inspired model of choice: Multi-alternatives, nonlinearity and value-based multidimensional choice. *Philosophical Transactions of the Royal Society B: Biological Sciences*, *362*(1485), 1655–1670.
- Brown, S. D., & Heathcote, A. (2008). The simplest complete model of choice response time: Linear ballistic accumulation. *Cognitive Psychology*, *57*(3), 153–178.
- Busemeyer, J. R., & Townsend, J. T. (1993). Decision field theory: A dynamic-cognitive approach to decision making in an uncertain environment. *Psychological Review*, *100*(3), 432.
- Carland, M. A., Thura, D., & Cisek, P. (2015). The urgency-gating model can explain the effects of early evidence. *Psychonomic Bulletin & Review*, *22*(6), 1830–1838.
- Chelazzi, L., Miller, E. K., Duncan, J., & Desimone, R. (1993). A neural basis for

visual search in inferior temporal cortex. *Nature*, *363*(6427), 345–347.

Cisek, P., Puskas, G. A., & El-Murr, S. (2009). Decisions in Changing Conditions: The Urgency-Gating Model. *Journal of Neuroscience*, *29*(37), 11560–11571.

Desimone, R. (1998). Visual attention mediated by biased competition in extrastriate visual cortex. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, *353*(1373), 1245–1255.

Diederich, A. (1995). Intersensory facilitation of reaction time: Evaluation of counter and diffusion coactivation models. *Journal of Mathematical Psychology*, *39*(2), 197–215.

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.

Erlick, D. E. (1961). Judgments of the relative frequency of a sequential series of two events. *Journal of Experimental Psychology*, *62*(2), 105–112.

Franco-Watkins, A. M., & Johnson, J. G. (2011). Decision moving window: Using interactive eye tracking to examine decision processes. *Behavior Research Methods*, *43*(3), 853–863.

Gratton, G., Coles, M. H. C., Sirevaag, E. J., Eriksen, C. J., & Donchin, E. (1988). Pre- and poststimulus activation of response channels: A psychophysiological analysis. *Journal of Experimental Psychology: Human Perception and Performance*, *14*(3), 331–344.

Heekeren, H. R., Marrett, S., Ruff, D. A., Bandettini, P. A., & Ungerleider, L. G. (2006). Involvement of human left dorsolateral prefrontal cortex in perceptual decision making is independent of response modality. *Proceedings of the National Academy of Sciences*, *103*(26), 10023–10028.

Holmes, W. R., Trueblood, J. S., & Heathcote, A. (2016). A new framework for modeling decisions about changing information: The Piecewise Linear Ballistic Accumulator model. *Cognitive Psychology*, *85*, 1–29.

Jones, M., & Dzhafarov, E. N. (2014). Unfalsifiability and mutual translatability of

major modeling schemes for choice reaction time. *Psychological Review*, *121*(1), 1–32.

Krajovich, I., Armel, C., & Rangel, A. (2010). Visual fixations and the computation and comparison of value in simple choice. *Nature Neuroscience*, *13*(10), 1292–1298.

Krajovich, I., & Rangel, A. (2011). Multialternative drift-diffusion model predicts the relationship between visual fixations and choice in value-based decisions. *Proceedings of the National Academy of Sciences*, *108*(33), 13852–13857.

LaBerge, D. (1962). A recruitment theory of simple behavior. *Psychometrika*, *27*(4), 375–396.

Lee, M. D. (2008). Three case studies in the Bayesian analysis of cognitive models. *Psychonomic Bulletin & Review*, *15*(1), 1–15.

Lee, W., & Janke, M. (1964). Categorizing externally distributed stimulus samples for three continua. *Journal of Experimental Psychology*, *68*(4), 376–382.

Loftus, G. R., & Masson, M. E. (1994). Using confidence intervals in within-subject designs. *Psychonomic Bulletin & Review*, *1*(4), 476–490.

Louie, K., Khaw, M. W., & Glimcher, P. W. (2013). Normalization is a general neural mechanism for context-dependent decision making. *Proceedings of the National Academy of Sciences*, *110*(15), 6139–6144.

Luce, R. D. (1986). Response times: Their role in inferring elementary mental organization.

Merkle, E. C., & Van Zandt, T. (2006). An application of the Poisson race model to confidence calibration. *Journal of Experimental Psychology: General*, *135*(3), 391–408.

Miletić, S., Turner, B. M., Forstmann, B. U., & van Maanen, L. (2017). Parameter recovery for the Leaky Competing Accumulator model. *Journal of Mathematical Psychology*, *76*, 25–50.

Mittner, M., Boekel, W., Tucker, A. M., Turner, B. M., Heathcote, A., & Forstmann, B. U. (2014). When the Brain Takes a Break: A Model-Based Analysis of Mind Wandering. *The Journal of Neuroscience*, *34*(49), 16286–16295.

Molloy, M. F., Galdo, M., Bahg, G., Liu, Q., & Turner, B. M. (2019). What's in a response time?: On the importance of response time measures in constraining models of context effects. *Decision*, *6*(2), 171–200.

Moreno-Bote, R. (2010). Decision Confidence and Uncertainty in Diffusion Models with Partially Correlated Neuronal Integrators. *Neural Computation*, *22*(7), 1786–1811.

Niwa, M., & Ditterich, J. (2008). Perceptual Decisions between Multiple Directions of Visual Motion. *Journal of Neuroscience*, *28*(17), 4435–4445.

Osth, A. F., & Dennis, S. (2015). Sources of interference in item and associative recognition memory. *Psychological Review*, *122*(2), 260–311.

Padoa-Schioppa, C. (2013). Neuronal Origins of Choice Variability in Economic Decisions. *Neuron*, *80*(5), 1322–1336.

Pais, D., Hogan, P. M., Schlegel, T., Franks, N. R., Leonard, N. E., & Marshall, J. A. R. (2013). A Mechanism for Value-Sensitive Decision-Making. *PLoS ONE*, *8*(9), e73216.

Palestro, J. J., Sederberg, P. B., Osth, A. F., Van Zandt, T., & Turner, B. M. (2018). Likelihood-free bayesian inference in cognitive science.

Palmer, J., Huk, A. C., & Shadlen, M. N. (2005). The effect of stimulus strength on the speed and accuracy of a perceptual decision. *Journal of Vision*, *5*(5), 376–404.

Pike, A. R. (1971). The latencies of correct and incorrect responses in discrimination and detection tasks: Their interpretation in terms of a model based on simple counting. *Attention, Perception, & Psychophysics*, *9*(6), 455–460.

Pilly, P. K., & Seitz, A. R. (2009). What a difference a parameter makes: A psychophysical comparison of random dot motion algorithms. *Vision Research*, *49*(13), 1599–1612.

Pirrone, A., Azab, H., Hayden, B. Y., Stafford, T., & Marshall, J. A. R. (2017). Evidence for the SpeedValue Trade-Off: Human and Monkey Decision Making Is Magnitude Sensitive. *Decision*, 129–142.

Purcell, B. A., Heitz, R. P., Cohen, J. Y., Schall, J. D., Logan, G. D., & Palmeri, T. J.

(2010). Neurally constrained modeling of perceptual decision making. *Psychological Review*, *117*(4), 1113–1143.

Raab, D. H. (1962). Statistical facilitation of simple reaction times. *Transactions of the New York Academy of Sciences*, *24*, 574–590.

Ratcliff, R. (1978). A theory of memory retrieval. *Psychological Review*, *85*(2), 59–108.

Ratcliff, R. (2006). Modeling response signal and response time data73. *Cognitive Psychology*, *53*(3), 195–237.

Ratcliff, R., & McKoon, G. (2008). The diffusion decision model: Theory and data for two-choice decision tasks. *Neural Computation*, *20*(4), 873–922.

Ratcliff, R., & Rouder, J. N. (1998). Modeling response times for two-choice decisions. *Psychological Science*, *9*(5), 347–356.

Ratcliff, R., & Smith, P. L. (2004). A Comparison of Sequential Sampling Models for Two-Choice Reaction Time. *Psychological Review*, *111*(2), 333–367.

Ratcliff, R., Voskuilen, C., & Teodorescu, A. (2018). Modeling 2-alternative forced-choice tasks: Accounting for both magnitude and difference effects. *Cognitive Psychology*, *103*, 1–22.

Reynolds, J. H., Chelazzi, L., & Desimone, R. (1996). Competitive Mechanisms Subserve Attention in Macaque Areas V2 and V4. *Journal of Cognitive Neuroscience*, *8*(4), 311–327.

Rodriguez, C. A., Turner, B. M., & McClure, S. M. (2014). Intertemporal Choice as Discounted Value Accumulation. *PLoS ONE*, *9*(2), e90138.

Rouder, J. N., Province, J. M., Morey, R. D., Gomez, P., & Heathcote, A. (2015). The lognormal race: A cognitive-process model of choice and latency with desirable psychometric properties. *Psychometrika*, *80*(2), 491–513.

Shadlen, M. N., & Newsome, W. T. (2001). Neural Basis of a Perceptual Decision in the Parietal Cortex (Area LIP) of the Rhesus Monkey. *Journal of Neurophysiology*, *86*(4), 1916–1936.

- Smith, S. M., & Krajbich, I. (2018). Gaze Amplifies Value in Decision Making. *Psychological Science*, 1–13.
- Strayer, D. L., & Kramer, A. F. (1994). Strategies and automaticity: 1. Basic findings and conceptual framework. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 20(2), 318–341.
- Swensson, R. G. (1972). The elusive tradeoff: Speed vs accuracy in visual discrimination tasks. *Perception & Psychophysics*, 12(1), 16–32.
- Teodorescu, A. R., Moran, R., & Usher, M. (2015). Absolutely relative or relatively absolute: Violations of value invariance in human decision making. *Psychonomic Bulletin & Review*, 23(1), 22–38.
- Teodorescu, A. R., & Usher, M. (2013). Disentangling decision models: From independence to competition. *Psychological Review*, 120(1), 1–38.
- Townsend, James T., & Nozawa, G. (1995). Spatio-temporal properties of elementary perception: An investigation of parallel, serial, and coactive theories. *Journal of Mathematical Psychology*, 39, 321–359.
- Trueblood, B., Jennifer S, & Heathcote, A. (2014). The Multiattribute Linear Ballistic Accumulator Model of Context Effects in Multialternative Choice. *Psychological Review*, 121(2), 179–205.
- Tsetsos, K., Gao, J., McClelland, J. L., & Usher, M. (2012). Using Time-Varying Evidence to Test Models of Decision Dynamics: Bounded Diffusion vs. The Leaky Competing Accumulator Model. *Frontiers in Neuroscience*, 6.
- Tsetsos, K., Usher, M., & McClelland, J. L. (2011). Testing Multi-Alternative Decision Models with Non-Stationary Evidence. *Frontiers in Neuroscience*, 5.
- Turner, B. M. (2019). Toward a common representational framework for adaptation. *Psychological Review*.
- Turner, B. M., Gao, J., Koenig, S., Palfy, D., & L. McClelland, J. (2017). The dynamics of multimodal integration: The averaging diffusion model. *Psychonomic Bulletin*



*Review*, 24(6), 1819–1843.

Turner, B. M., Schley, D. R., Muller, C., & Tsetsos, K. (2018). Competing theories of multialternative, multiattribute preferential choice. *Psychological Review*, 125(3), 329–362.

Turner, B. M., & Sederberg, P. B. (2014). A generalized, likelihood-free method for posterior estimation. *Psychonomic Bulletin & Review*, 21(2), 227–250.

Turner, B. M., Sederberg, P. B., Brown, S. D., & Steyvers, M. (2013). A method for efficiently sampling from distributions with correlated dimensions. *Psychological Methods*, 18(3), 368–384.

Turner, B. M., Sederberg, P. B., & McClelland, J. L. (2016). Bayesian analysis of simulation-based models. *Journal of Mathematical Psychology*, 72, 191–199.

Turner, B. M., van Maanen, L., & Forstmann, B. U. (2015). Informing cognitive abstractions through neuroimaging: The neural drift diffusion model. *Psychological Review*, 122(2), 312–336.

Turner, B. M., & Van Zandt, T. (2012). A tutorial on approximate Bayesian computation. *Journal of Mathematical Psychology*, 56(2), 69–85.

Turner, B. M., & Van Zandt, T. (2018). Approximating Bayesian Inference through Model Simulation. *Trends in Cognitive Sciences*, 22(9), 826–840.

Usher, M., & McClelland, J. L. (2001). The time course of perceptual choice: The leaky, competing accumulator model. *Psychological Review*, 108(3), 550–592.

Usher, M., & McClelland, J. L. (2004). Loss Aversion and Inhibition in Dynamical Models of Multialternative Choice. *Psychological Review*, 111(3), 757–769.

van Ravenzwaaij, D., Brown, S. D., Marley, A. A. J., & Heathcote, A. (2019). Accumulating advantages: A new conceptualization of rapid multiple choice. *Psychological Review*, 127(2), 186–215.

van Ravenzwaaij, D., van der Maas, H. L. J., & Wagenmakers, E.-J. (2012). Optimal decision making in neural inhibition models. *Psychological Review*, 119(1), 201–215.

Vickers, D. (1970). Evidence for an Accumulator Model of Psychophysical

Discrimination. *Ergonomics*, *13*(1), 37–58.

Vickers, D., Burt, J., Smith, P., & Brown, M. (1985). Experimental paradigms emphasising state or process limitations: I effects on speed-accuracy tradeoffs. *Acta Psychologica*, *59*(2), 129–161.

Vickers, D., & Smith, P. L. (1989). Modeling evidence accumulation with partial loss in expanded judgment. *Journal of Experimental Psychology: Human Perception and Performance*, *15*(4), 797–815.

Vuckovic, A., Kwantes, P. J., Humphreys, M., & Neal, A. (2014). A sequential sampling account of response bias and speed-accuracy tradeoffs in a conflict detection task. *Journal of Experimental Psychology: Applied*, *20*(1), 55–68.

Wagenmakers, E.-J., Ratcliff, R., Gomez, P., & McKoon, G. (2008). A diffusion model account of criterion shifts in the lexical decision task. *Journal of Memory and Language*, *58*(1), 140–159.

Wang, X.-J. (2002). Probabilistic decision making by slow reverberation in cortical circuits. *Neuron*, *36*(5), 955–968.

Wickelgren, W. A. (1977). Speed-accuracy tradeoff and information processing dynamics. *Acta Psychologica*, *41*(1), 105–112.

Winkel, J., Keuken, M. C., van Maanen, L., Wagenmakers, E.-J., & Forstmann, B. U. (2014). Early evidence affects later decisions: Why evidence accumulation is required to explain response time data. *Psychonomic Bulletin & Review*, 774–784.

Zylberberg, A., Barttfeld, P., & Sigman, M. (2012). The construction of confidence in a perceptual decision. *Frontiers in Integrative Neuroscience*, *6*.