# Supplemental Material

Angus Reynolds

December 2019

# 1   Likelihood Functions

The likelihoods are defined with respect to LBA pdf and CDF functions provided in Brown and Heathcote (2008) or in the rtdists package. https://CRAN.R-project.org/package=rtdists. The typical LBA model has thresholds parameters $A$ and $b$, where $A$ is the height of the uniform start point noise and $b$ is the response threshold height. The MTR adds a parameter $d$, which is the intermediate threshold. We can consider two versions of this model. One where d is bounded between $A$ and $b$, and one where it is just bounded above by $b$. The LBA also has drift rate parameters v and sv; the mean drift rate and the standard deviation of drift rate. Lastly, there is non-decision time $t_{er}$, which is subtracted from the response time to give the estimated decision time. Since we are working with the LBA, we will defining the MTR likelihood functions with respect to the LBA density function $f$, its cumulative function $F$ or the survival function $S = 1 - F$. For example $f(A, b, v, sv, dt)$ is the density of an LBA function finishing in time dt given threshold parameters $A$ and $b$ and drift rate parameters $v$ and $sv$. For brevity we will sometimes drop the $dt$ term and write $f(A, b, v, sv)$. Given that for a specific stimulus there are possible definitive responses i=1,2 made in time $dt = RT - t_{er}$, the likelihood for a definitive response i=1 is:

$$L(dt, i = 1) = f(A_1, b_1, v_1, sv_1, dt) \times S(A_2, d_2, v_2, sv_2, dt) \tag{1}$$

This is the instantaneous density of accumulator 1 finishing in time dt, multiplied by the probability that accumulator 2 has not reached even the lower $d_2$ threshold in time $dt$. Note that in many instances $A_1$ and $A_2$, will be the same.

The likelihood of a don't know response is the sum of two LBA likelihoods:

$$L(dt, DK) = f(A_1, b_1, v_1, sv_1) \times (F(A_2, d_2, v_2, sv_2) - F(A_2, b_2, v_2, sv_2)) +$$
$$f(A_2, b_2, v_2, sv_2) \times (F(A_1, d_1, v_1, sv_1) - F(A_1, b_1, v_1, sv_1)) \tag{2}$$

Here for each sum we have the density of the winner, multiplied by a the probability that the losing accumulator is between the $d$ and $b$ threshold.

If $d < A$ the definitive response likelihood changes to:

$$L(dt, i = 1) = f(A_1, b_1, v_1, sv_1, dt) \times \frac{d_2}{A - 2} \times S(A_2, d_2, v_2, sv_2, dt) \tag{3}$$

Here $\frac{d_2}{A_2}$ is the probability that the uniform start point for the second accumulator was bellow $d_2$ when sampling started, so that a definitive response was possible.

The don't know likelihood changes to :

$$\begin{aligned}
L(dt, DK) = f(A_1, b_1, v_1, sv_1) \times (&\frac{d_2}{A - 2} \times (F(A_2, d_2, v_2, sv_2) - F(A_2, b_2, v_2, sv_2)) + \\
&(1 - \frac{d_2}{A_2}) \times S(A_2 - d_2, b_2 - d_2, v_2, sv_2)) + \\
f(A_2, b_2, v_2, sv_2) \times (&\frac{d_2}{A - 2} \times (F(A_1, d_1, v_1, sv_1) - F(A_1, b_1, v_1, sv_1)) + \\
&(1 - \frac{d_1}{A_1}) \times S(A_1 - d_1, b_1 - d_1, v_1, sv_1))
\end{aligned} \tag{4}$$

# 2 Experiment 1

## 2.1 Linear mixed effects models

Definitive Responses

A strong response bias was evident favouring left responses, with the pattern of effects across accuracy and RT for left vs. right targets being exactly as predicted by a lower threshold for left than right responses in an evidence-accumulation model. Overall accuracy was much greater when the target was on the left (82.1%) than the right (60.4%), $\chi^2(1) = 350$, p $<$.001. This was accompanied by faster responses to left (1.1s) than right (1.15s) targets, $\chi^2(1) = 27$, p $<$.001. The latter effect interacted strongly with response accuracy, with correct responses to left targets being 0.09s faster than errors, whereas correct responses to right targets were 0.07s slower than errors, $\chi^2(1) = 77$, p $<$.001. Left responses (i.e., correct responses for left targets and error responses for right targets) are faster because it takes like less time to accrue sufficient evidence to reach the lower left-accumulator threshold. Accuracy is higher for left targets because the higher right threshold causes a bias to make left responses. The speed-accuracy trade-off manipulation was clearly successful, with faster but less accurate responses in the speed-emphasis condition (1.04s, 70.5%) than the accuracy-emphasis condition (1.21s, 74.6%), $\chi^2(1) = 11.7$, p $<$.001 and $\chi^2(1)$ = 706, p $<$.001, respectively. For mean RT there was a large interaction of this effect with the error-cost manipulation; the speedup was 0.22s in the high error-cost condition, but only 0.13s in the low error-cost condition. Surprisingly, few effects involving the similarity manipulation were

2

significant. The interaction between similarity and error cost was the largest, $\chi^2(1) = 3.7$, p = .054, due to a trend for an advantage in accuracy for the high over low error-cost conditions to be larger for low-similarity pairs (8.0%) than high-similarity pairs (3.6%). In terms of mean RT there was a just-significant interaction with both speed vs. accuracy and error cost, $\chi^2(1) = 4.0$, p = .046, whereby responding was slower for higher- than lower-similarity pairs in the low error-cost speed condition (by 0.015s) and high-error cost accuracy condition (by 0.011s), but negligibly slower in the low error-cost accuracy condition (by 0.001s) and faster in the high error-cost speed condition (by 0.019s). There was also a stronger interaction of similarity and speed vs. accuracy manipulations with response, $\chi^2(1) = 8.35$, p = .004. This reflects under speed emphasis left responses were quicker than right responses by a larger margin for higher- than lower-similarity pairs (0.09s vs. 0.06s) but the reverse was true in the accuracy condition (0.08s vs. 0.1s).

Don't-know Responses

Don't-know responses were on average fairly common, but less so for left targets (38%) than right targets (41.3%), $\chi^2(1) = 12$, p $<$.001. The latter difference was almost entirely in the speed condition (37.5% vs. 43.3%) with little difference in the accuracy condition (38.4% vs. 39.4%), $\chi^2(1) = 5.6$, p = .018. There was also an interaction between error cost and similarity; there were more don't-know responses under high than lower error cost for lower-similarity pairs (by 7.7%), but that increase was smaller for higher-similarity pairs (by 2.4%), $\chi^2(1) = 6.7$, p = .01. Overall, mean RT of don't-know responses (1.13s) were slower than left responses (1.1s), $\chi^2(1) = 5.1$, p = .02, but faster than right responses (1.18s), $\chi^2(1) = 39.3$, p $<$.001. The effect of the speed vs. accuracy manipulation on don't-know responses (0.2s) was larger than on left responses (0.17s), $\chi^2(1) = 41.5$, p $<$.001, and also than on right responses (0.18s), $\chi^2(1) = 25.6$, p $<$.001.

## 2.2 MTR Model Description

As described in the main paper, choice thresholds and intermediate thresholds varied by speed/accuracy emphasis and the response (left or right), resulting in 8 threshold parameters. There were 8 mean rate parameters, half for the matching accumulator and half for the mismatching accumulator, with different values for lower- and higher-similarity pairs in the speed and accuracy conditions. There was a single non-decision time estimate. One value of the rate standard deviation was estimated, for the matching accumulator, with the value for the mismatching accumulator fixed at 1 in order to make the model identifiable (Donkin, Brown & Heathcote, 2009). Finally, one value of start-point variability was estimated for all conditions, so in total 19 parameters were estimated per participant. In addition to this model, we ran numerous other models with fewer, greater and sometimes the equivalent number of parameters. Ultimately, we chose on this model as it provided the best subjective fit to the patterns in the data, while also attaining a good score on the Deviance Information Criterion (DIC). We find that with MTR models, DIC tends to prefer complex models.

Later we show that when these models are too complex, parameter recovery suffers, since there are numerous ways the model can account for the same pattern in data, something that DIC does not penalize for. A 23 parameter model that let word-frequency effect the placement of the d threshold provided improved DIC over any 19 parameter models we tested. However, this requires the more tenuous assumption that subjects judge the relative similarity of the two stimulus and then adjust their confidence levels during the decision. Similarity is not known before the trial begins. To properly test if subjects adjust threshold in this manner, a design where the similarity of the two stimulus is provided immediately before each trial so that it would be more plausible that subjects could adjust thresholds in response. Additional model fits can be found at osf.io/6h4qe/

Population means for each parameter were given Gaussian priors that were truncated below at zero for all but the mean rate parameters and truncated above at 1 for non-decision time parameters. Start-point noise and threshold priors had a mean of 1, as did mean rates for mismatching accumulators and the matching rate standard deviation. Matching accumulators had a prior mean of two, and non-decision time parameters a prior mean of 1s. All mean rate parameter population priors had a standard deviation of 4, thresholds and start-point noise a standard deviation of 4, and all other parameters a standard deviation of 1. Estimation was carried out using the DMC software using the methods described in Heathcote et al. (2019). Briefly, the procedure was to fit subjects individually first to give good starting points for the hierarchical modelling. Then three sampling functions are used; the first automatically runs samples until there are no stuck chains (chains whose likelihood is far from all the other chains) in the sampling, then a function to run the hierarchical chains until they are converged, and finally a function that runs extra samples until a desired number of posterior samples is achieved. This resulted in 60 chains each with 250 samples per chain after thinning out 9 out of 10 samples (i.e. requiring 2500 samples per chain) used to form the final posterior sample. Convergence was confirmed visually and through low values of Brooks and Gelman's (1998) multivariate $R^2$ (1.01 at the hyper level and an average of 1.04 over participants for both low and high error-cost fits).

## 2.3 Extra Model Parameters.

As shown in Figure S1, non-decision time was faster in the low than high error cost condition (p <.001), by about 0.08s, whereas start-point noise was larger in the low than high error cost condition. The rate standard deviation for the matching accumulator was estimated at 0.86 for both low and high error cost conditions, clearly less than the fixed rate standard deviation of 1 for the mismatching accumulator (ps <.001).
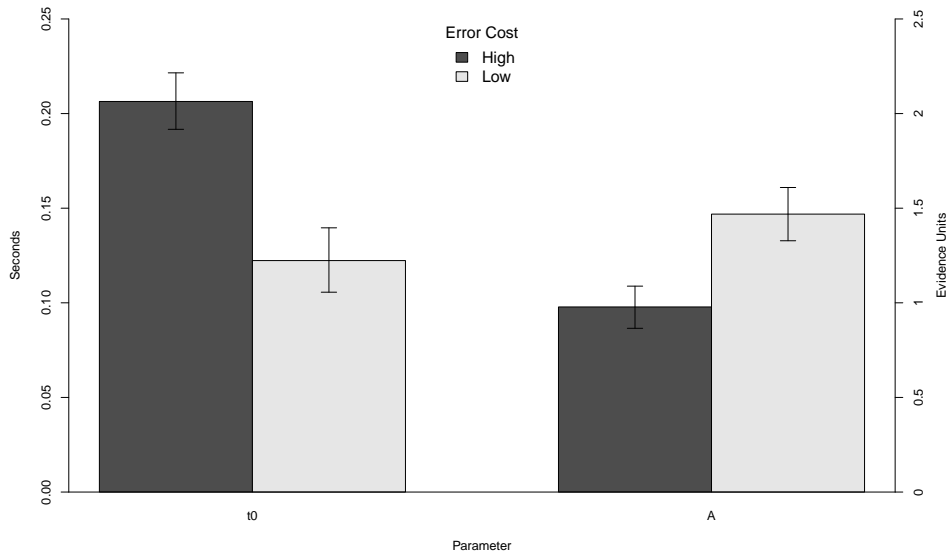
Figure S1: Experiment 1 non-decision time and start-point noise estimates with 95% credible intervals for high and low error-cost conditions.

## 2.4 Global Fit

Figures S2 and S3 show the model average fit to each condition. The error bars denote the 95% credible intervals for the model predictions. The global fit graphs show the improved accuracy on left responses, with little change to the probability of don't know responses. With respect to RT, the model consistently predicts that don't know responses are slower than definitive responses, which is typically the pattern in the data with some exceptions. The low error cost speed condition trials have faster don't know responses, particularly for the $10^{th}$ percentile. Figures S4 shows the averaged model PDF for each response, collapsing over difficulty. This figure highlights how the model predicts very little change to don't know use by stimulus, but does account for the change in definitive responses. Figures S5 and S6 shows the CDFs of the data and model fit to experiment 1.
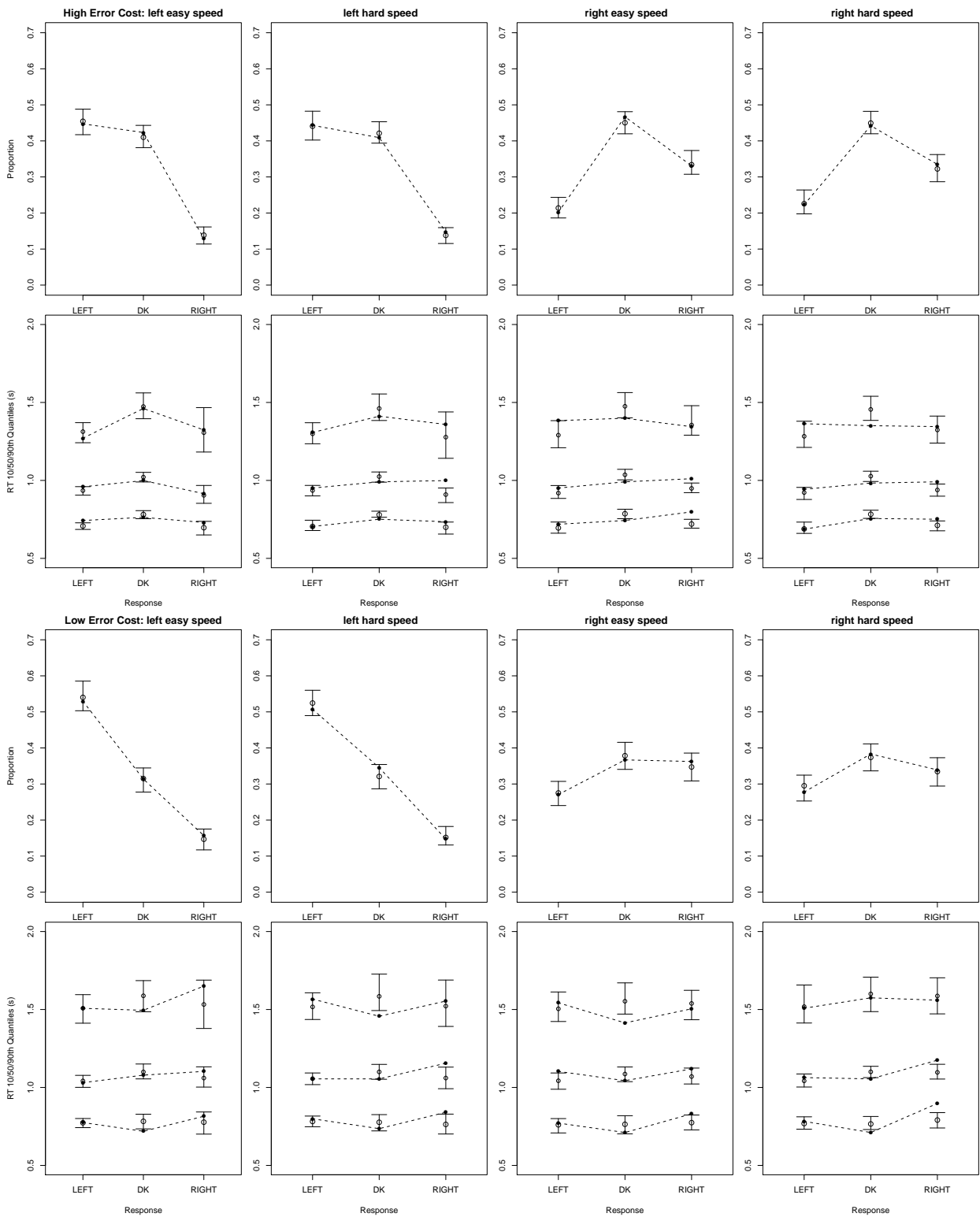
Figure S2: Average model fit to Experiment 1, speed emphasis trials. Response proportions and 10, 50 and 90$^{th}$ percentiles. The arrows show the 95% credible interval of the model predictions.
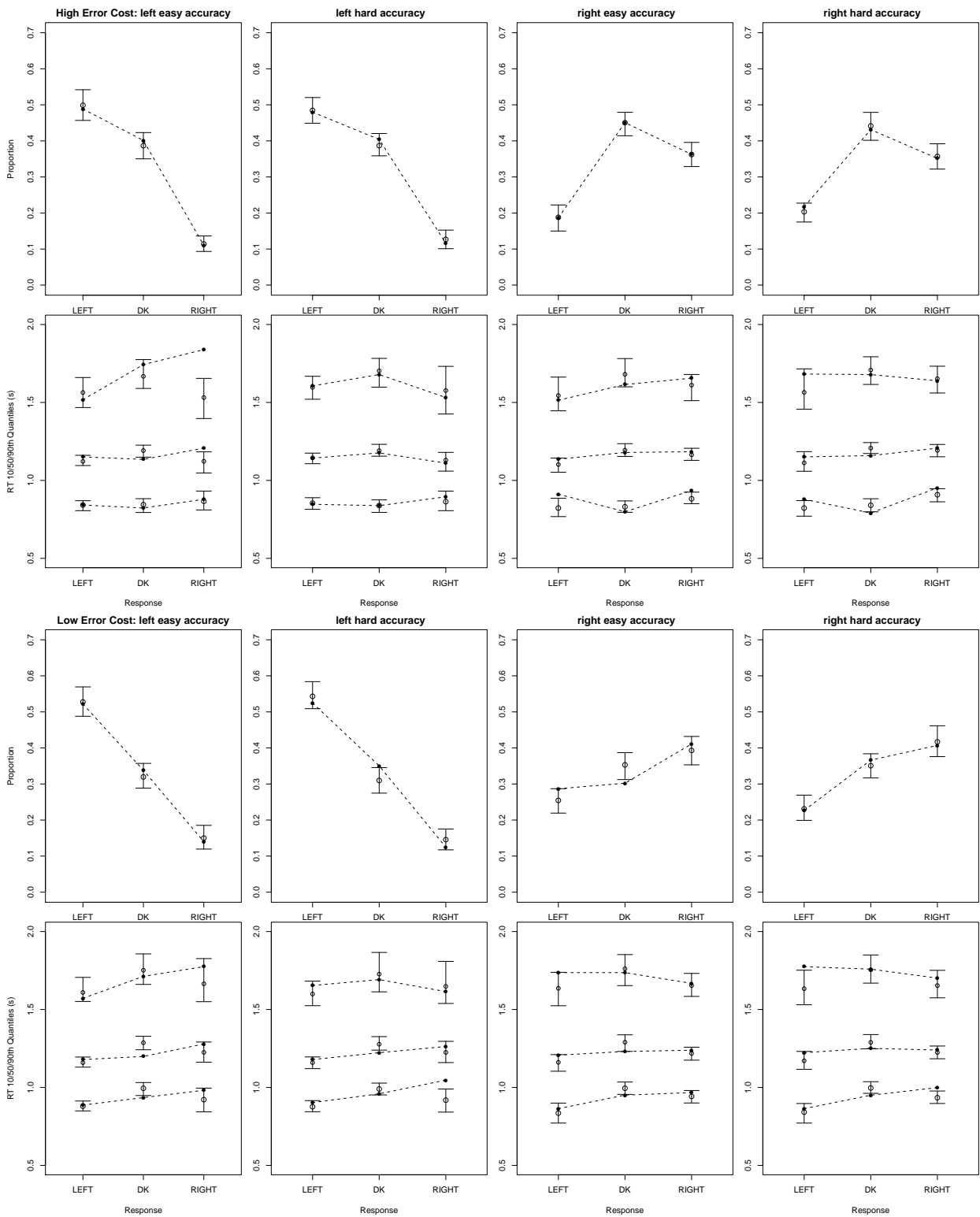
Figure S3: Average model fit to Experiment 1, accuracy emphasis trials. Response proportions and 10, 50 and 90$^{th}$ percentiles. The arrows show the 95% credible interval of the model predictions.
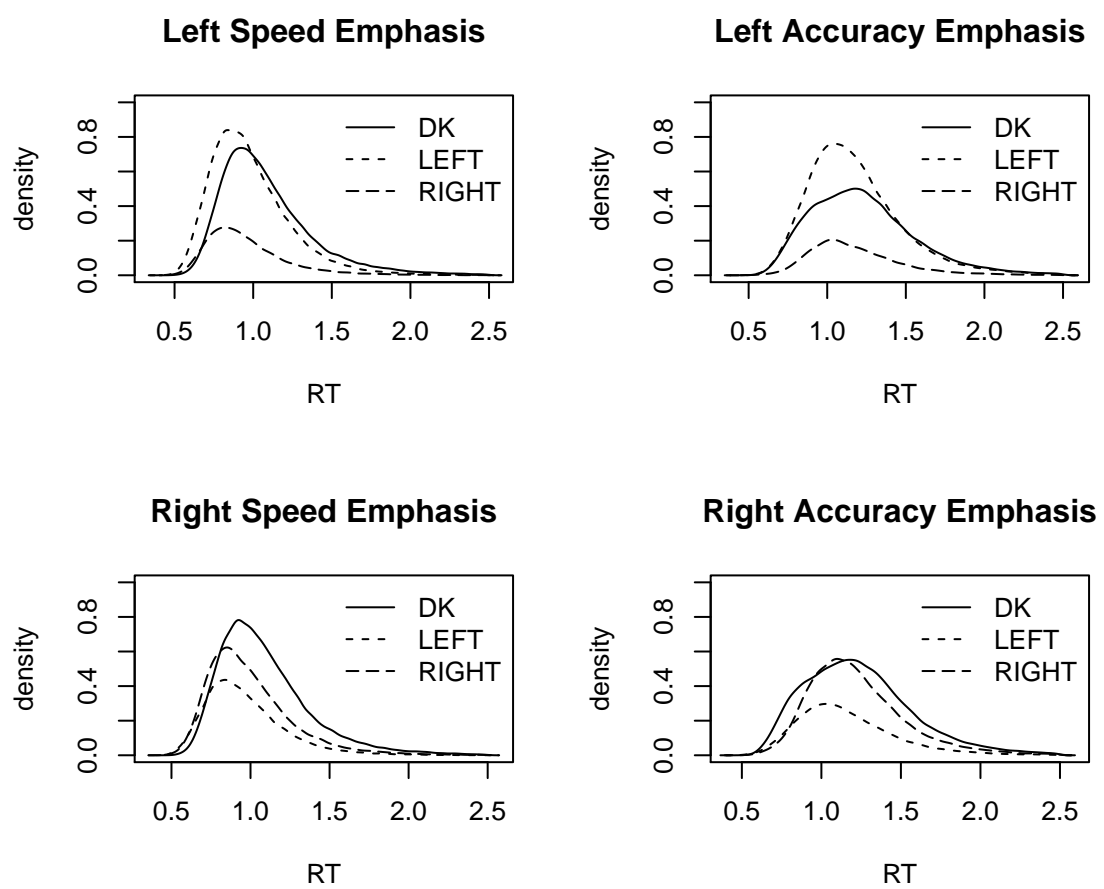
Figure S4: MTR Probability Density Function, averaging over all High Error cost subjects and collapsing accross difficulty for experiment 1.
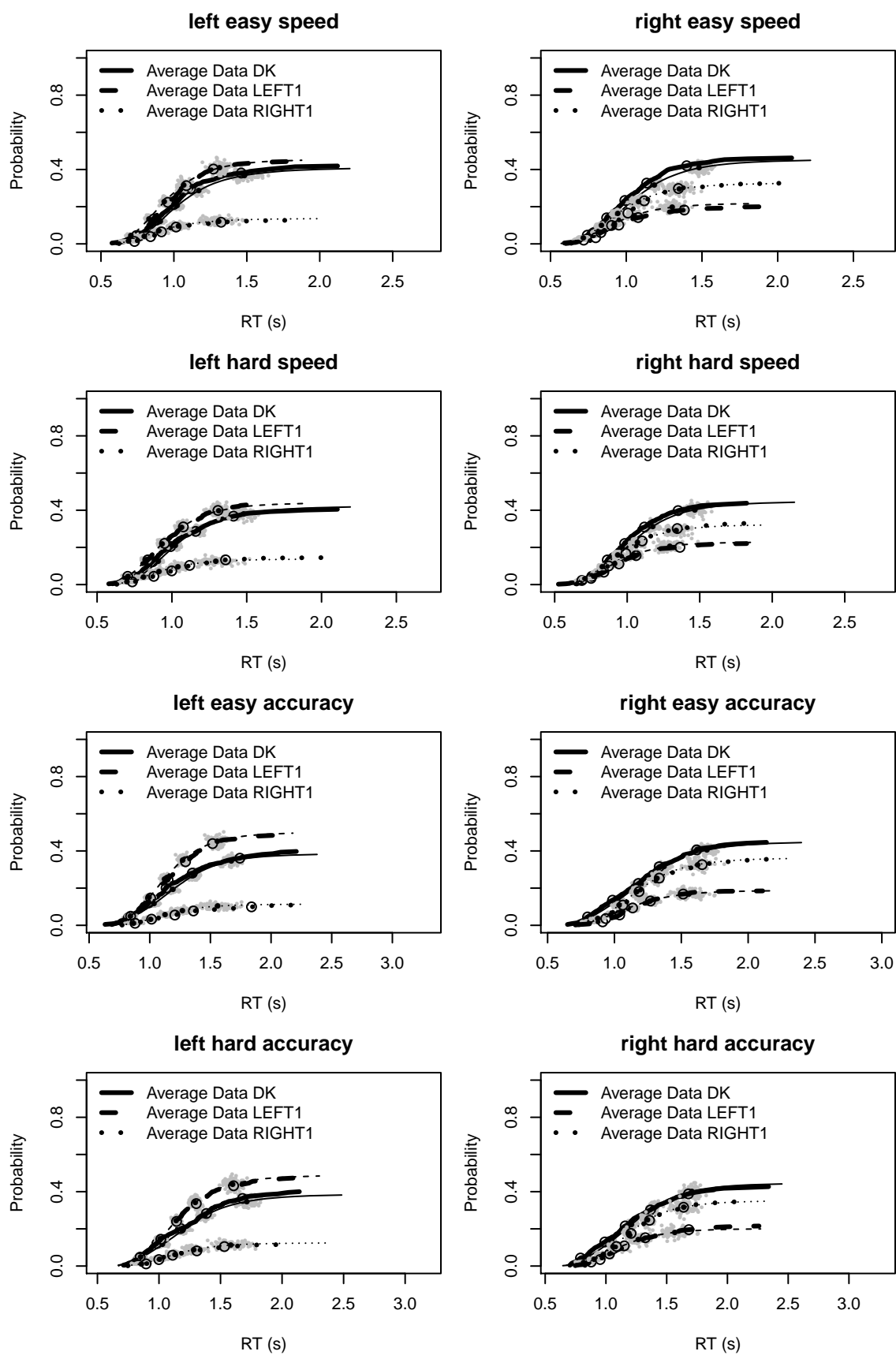
Figure S5: MTR CDF, averaging over all high error cost subjects for experiment 1.
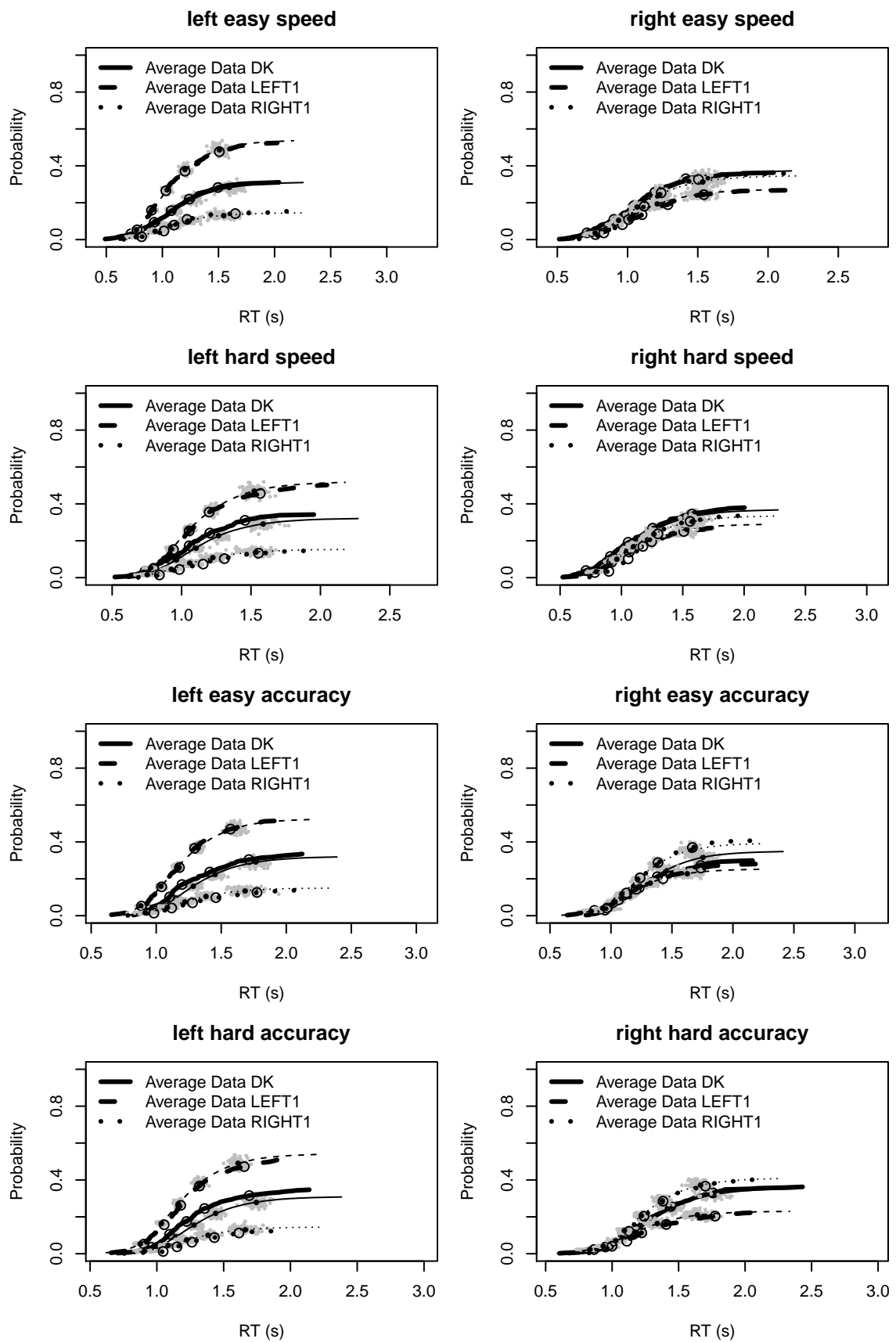
Figure S6: MTR CDF, averaging over all low error cost subjects for experiment 1.

# 3 Experiment 2

## 3.1 Linear mixed effects models.

Accuracy was lower under speed (66.5%) than accuracy (69.4%) emphasis, $\chi^2(1) = 8.2$, p = 0.004. In contrast to the first experiment there was a strong effect of similarity, with high similarity lures barely above chance (52.9%), and low similarity lures (66.8%) a little less accurate than targets (74%), $\chi^2(2) = 328$, p < .001. Accuracy was higher when the subject was on a high score (69.5%) than a low score (65.5%), $\chi^2(1) = 5.07$, p < .024. There was also a significant interaction between the stimulus and error cost effects, where for higher similarity lures accuracy was 1.6% worse for lower than higher error cost, but this reversed for lower similarity lures (low more accurate than high by 1.6%) and even more so for targets (a 4.6% advantage), $\chi^2(2) = 6.2$, p = .045.

Mean RT was substantially faster under speed emphasis (0.796s) than accuracy emphasis (1.04s), $\chi^2(1) = 3090$, p < .001. It was also faster for targets (0.882s) than either lower similarity lures (0.945s) or higher similarity lures (0.941s), $\chi^2(2) = 203$, p < .001. These two effects interacted, $\chi^2(2) = 13.4$, p = .001, with speed faster than accuracy emphasis by 0.215s for targets and 0.271s and 0.293s for lower and higher similarity lures respectively(Figure 11). While the effect of score on mean RT was significant, $\chi^2(1) = 24.5$, p < .001, low score responses were on average only 18ms faster.

There were strong two-way interactions between the speed of correct vs. error responses and both stimulus, $\chi^2(2) = 662$, p < .001, and speed vs. accuracy emphasis, $\chi^2(1) = 17.2$, p < .001, which combined in a just significant interaction between all three factors, $\chi^2(2) = 8.7$, p = .013. For targets, errors were slower than corrects not only in the accuracy condition (by 0.212s) but also to a lesser degree in the speed condition (by 0.095s). For lures, errors were faster than corrects, with the difference for lower similarity lures being larger in the speed condition (0.126s) than the accuracy condition (0.081s) but the reverse being true for higher similarity lures (0.103s vs. 0.147s).

Don't-know responses were more common for lures (28% and 29% for lower and higher similarity respectively) than targets (17.4%), $\chi^2(2) = 216$, p < .001. They were also more common when error-cost was high than low, $\chi^2(1) = 13.4$, p < .001, and this effect interacted with speed vs. accuracy, $\chi^2(1) = 4.7$, p = .03, due to a greater increase under speed emphasis (12%) than under speed emphasis (8.8%). Don't Know responses were more common when the subject had collected 2 or more points (23.9%) than when they had 0 or 1 points (20.4 %) , $\chi^2(1) = 95.4$, p < .001. This effect interacted with error cost too, $\chi^2(1) = 4.7$, p = .03 with the difference between low and high score don't know probabilities 14.2% in the high error cost condition and 9.3% in the low error cost condition. This is intuitive as in the low error cost condition to maximise your score you would ignore the don't know response always, whereas in the high error cost condition
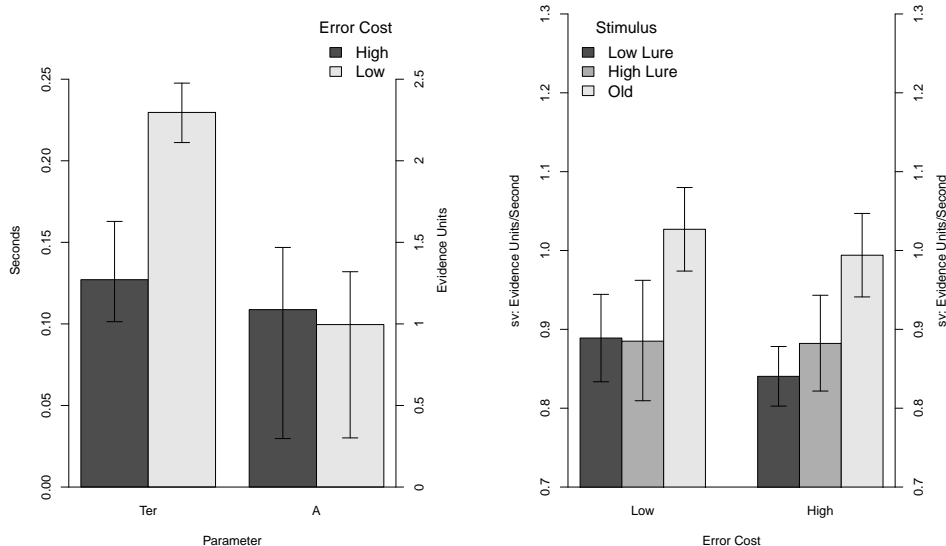
Figure S7: Experiment 2 MTR parameter estimates-19 parameter model. Left hand panel shows $t_er$ and $A$. The right hand panel shows the matching *sv* estimates. Mismatching *sv* was fixed to 1.

it is advantageous to adjust the don't know thresholds as your score increases.

Mean RT for don't-know responses (1.00) was significantly slower than target responses (0.855s) but did not differ significantly from lure responses (0.993s), $\chi^2(2) = 942$, p $<$.001. Don't-know responses were much slower under accuracy than speed emphasis (by 0.325s), which was larger than the effect on both new responses (0.274s), $\chi^2(1) = 45.7$ and that of old responses (0.221s), $\chi^2(2) = 36.6$, p $<$.001.

Error cost also strongly interacted with response, with don't-know responses 0.035s faster for high than low error cost, but new and old responses faster for low than high error cost (by 0.050s and 0.067s, respectively), $\chi^2(2) = 40.8$, p $<$.001. The error cost effect on don't-know responses was modulated by emphasis, being virtually absent under accuracy emphasis (.003s vs. 0.037s under speed emphasis), whereas there was little effect on new (speed: 0.052s, accuracy: 0.047s) or target (speed: 0.067s, accuracy: 0.054s) responses, $\chi^2(2) = 8.8$, p = .012.

## 3.2 Extra Model Parameters

Figure S7 shows the parameter estimates for non decision time, start point noise and the matching drift rate standard deviation. Non decision time was slower in low error cost conditions (p<.001). Standard deviations were less for lure stimulus trials than target stimulus trials (p¡.001).

## 3.3 Probability Density Functions

## 3.4 Global Fit

Figures S9 and S10 show the data and model averaged CDFs of the 19 parameter models for experiment 2. Figure S8 shows the aggregated PDF to the high error cost subjects. Figure
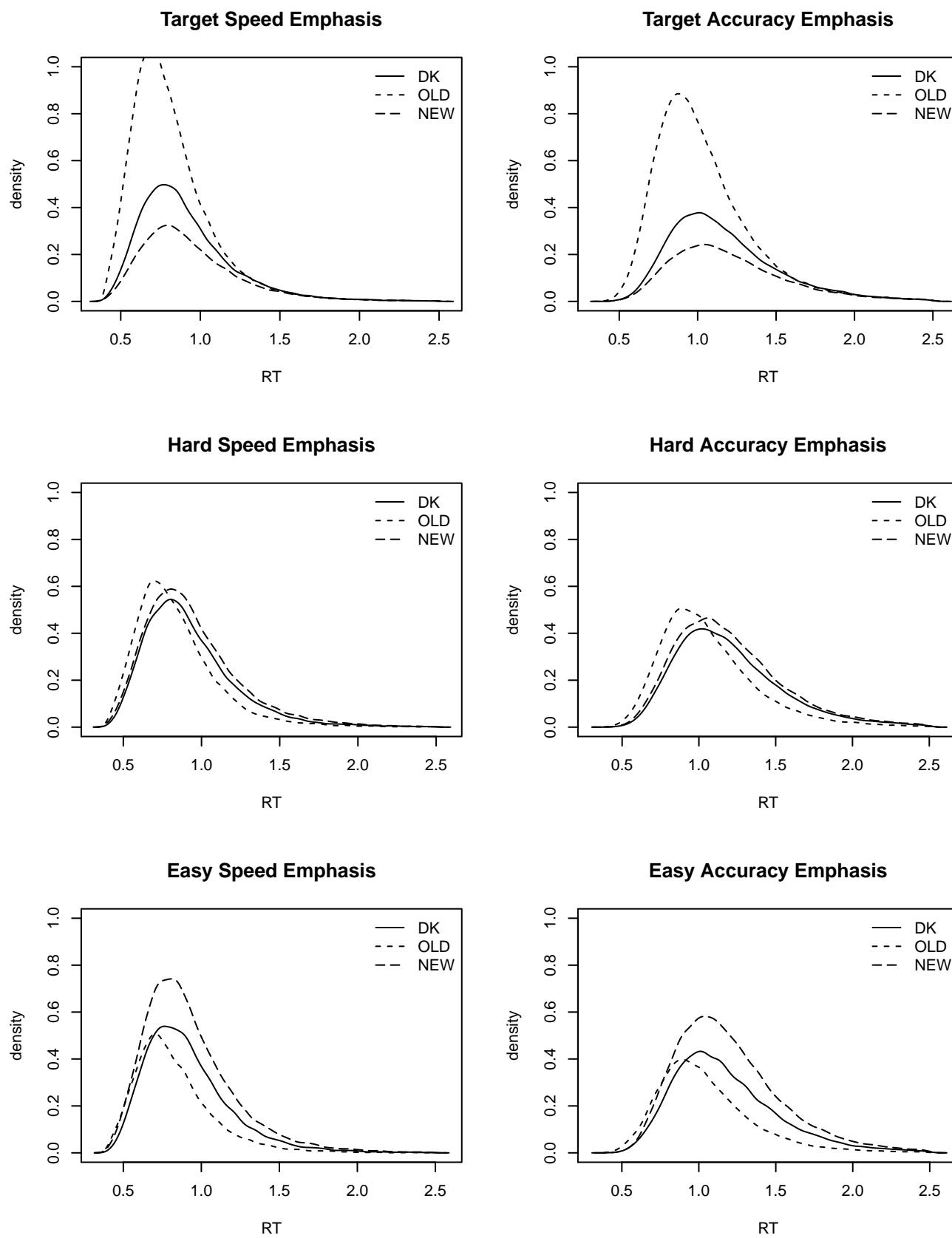
Figure S8: MTR Probability Density Function, averaging over all High Error cost subjects for experiment 2.
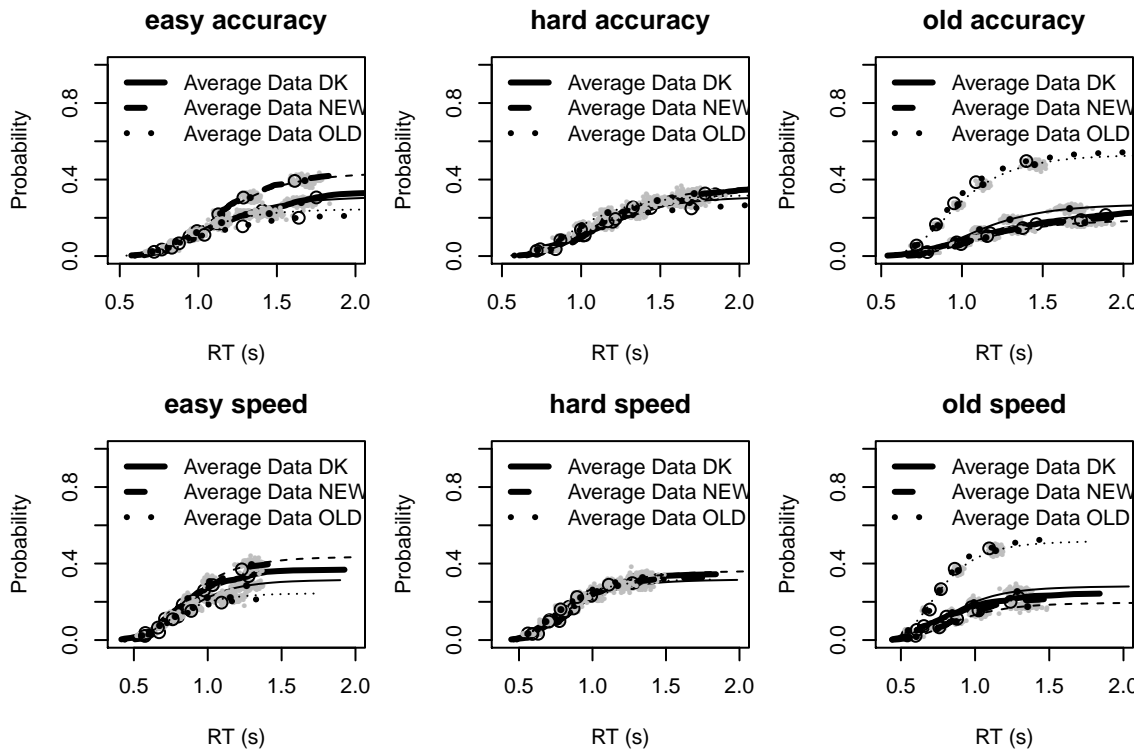
Figure S9: MTR Cumulative Density Function, averaging over all High Error cost subjects for experiment 2.
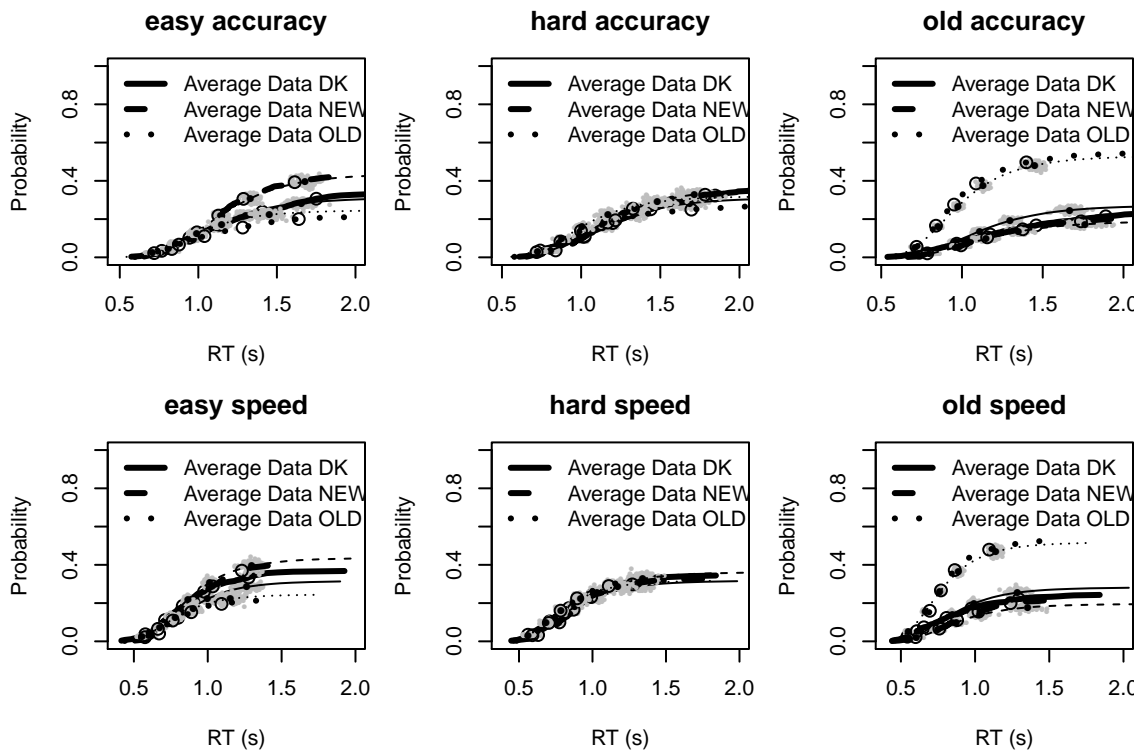


Figure S10: MTR Cumulative Density Function, averaging over all High Error cost subjects for experiment 2.

S11,S12,S13,S14 show the same fits with the model and data response proportions and the 10th, 50th and 90th RT percentiles with 95% credible intervals of the model predictions. Overall the model does quite a good job fitting the data, though there are some problems with the 90th percentile in particular in some cells.

## 3.5 Additional Model Fit

In addition to the model fits presented in the paper, we ran numerous other MTR models with differing levels of model complexity. Fits can be found at osf.io/6h4qe/. Here we present a model that while marginally better than the paper model by its DIC score, it has 10 more parameters. There was not a large amount of data collected for each subject, so we are less confident that there is enough data to reliable distinguish between a speed/accuracy effect of threshold and drift rate. In the main paper we opted for a simpler model that only identified a difference on thresholds. Here we present the more complex model fit that includes the manipulation of speed emphasis on drift rate parameters.

Figure S15 and S16 show the association between DK and the probability of a don't-know response. As in Experiment 1 there was a strong linear relationship with a slope close to two for both high error cost (slope = 1.9, r2 = .99, p <.001) and low error cost (slope = 1.83, r2 = .97, p <.001). As shown in Figures 12c and 12d, the same type of strong linear relationship held between the difference in don't-know probability between high and low trials for both high error cost (slope = 2.2, r2 = .96, p <.001) and low error cost (slope = 2.1, r2 = .87, p <.001) conditions.

Figure S17 shows that non decision time was smaller in the higher error cost condition and that the start point noise was slightly smaller on average in the low error cost condition. The drift rate variance on target stimulus trials was larger than easy or hard lure stimulus trials (p <.001).

Figure S19 shows that DK estimates in the high error-cost condition were larger for high than low score trials (ps <.001, except for the target accumulator under speed emphasis, p = .049). This was also the case for low error cost in the speed condition for the lure accumulator, and in the accuracy condition for the target accumulator (ps <.001), whereas in the accuracy condition for the lure accumulator the opposite trend was evident (p = .012) and for the target accumulator in the speed condition DK was near floor for both low and high score trials.
In the high error-cost condition DK was larger for the accuracy than speed condition for the target accumulator and vice versa for the lure accumulator (ps <.001). With low error cost DK for accuracy emphasis was generally greater than for speed emphasis (ps <.025) except for the lure accumulator in high score trials where this reversed (p = .021). In the low error-cost condition DK was consistently greater for the lure than target accumulator (ps <.001). This was also the case

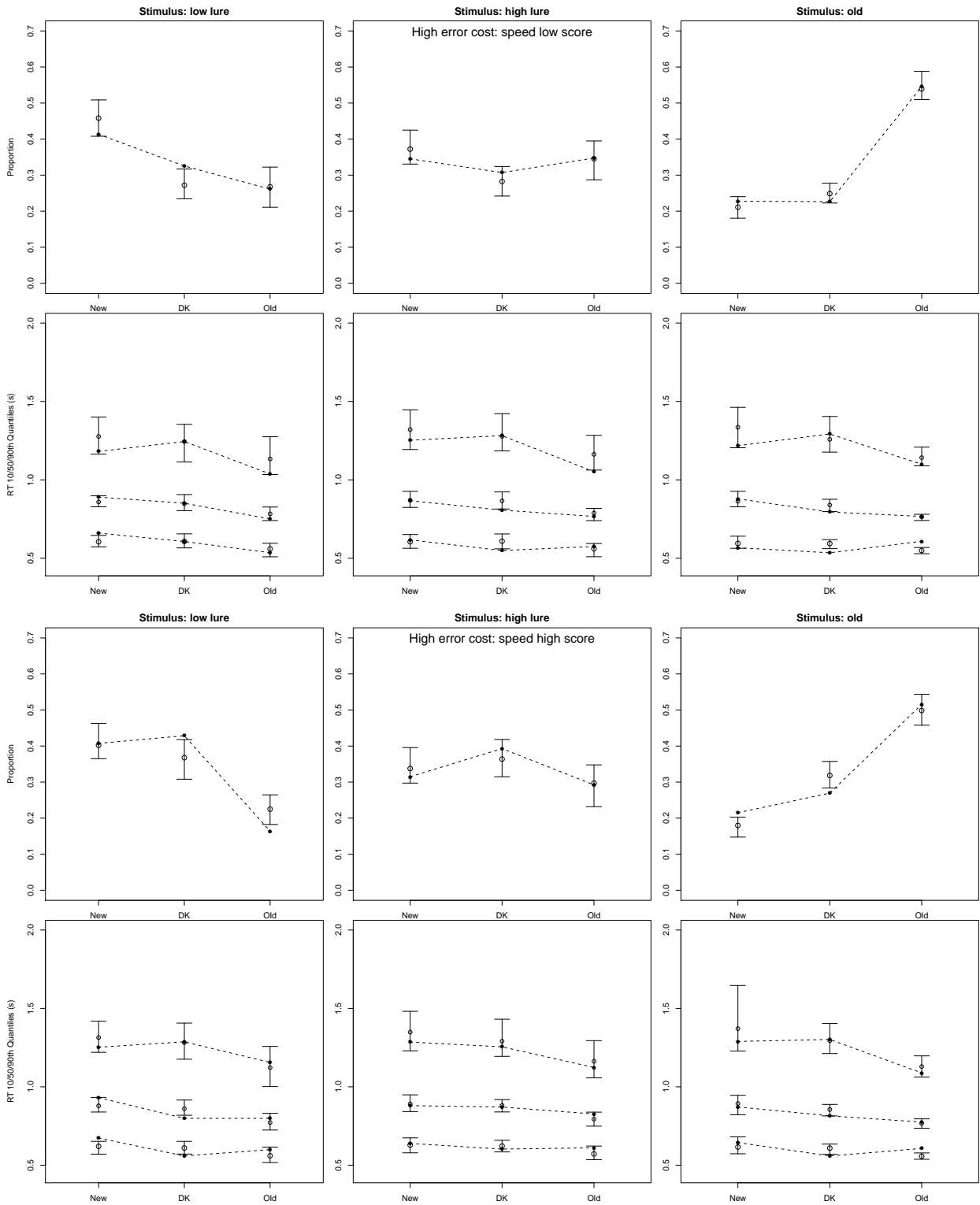Figure S11: Average model fit to Experiment 2, speed emphasis, high error cost, 19 parameter paper model. Response proportions and 10, 50 and $90^{th}$ percentiles. The arrows show the 95% credible interval of the model predictions.

Figure S12: Average model fit to Experiment 2, accuracy emphasis, high error cost, 19 parameter paper model. Response proportions and 10, 50 and $90^{th}$ percentiles. The arrows show the 95% credible interval of the model predictions.

Figure S13: Average model fit to Experiment 2, speed emphasis, low error cost, 19 parameter paper model. Response proportions and 10, 50 and $90^{th}$ percentiles. The arrows show the 95% credible interval of the model predictions.

Figure S14: Average model fit to Experiment 2, accuracy emphasis, low error cost, 19 parameter paper model. Response proportions and 10, 50 and $90^{th}$ percentiles. The arrows show the 95% credible interval of the model predictions.
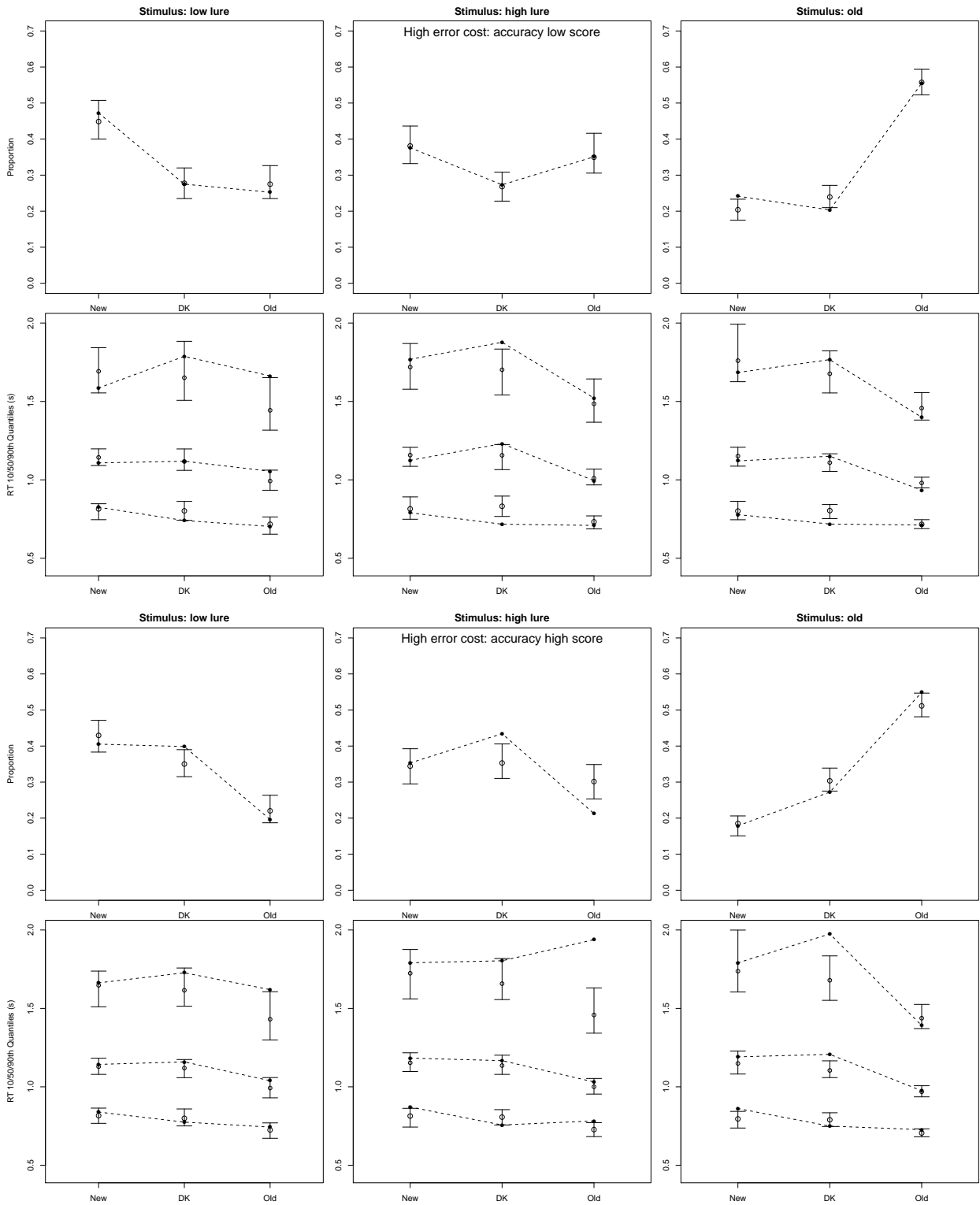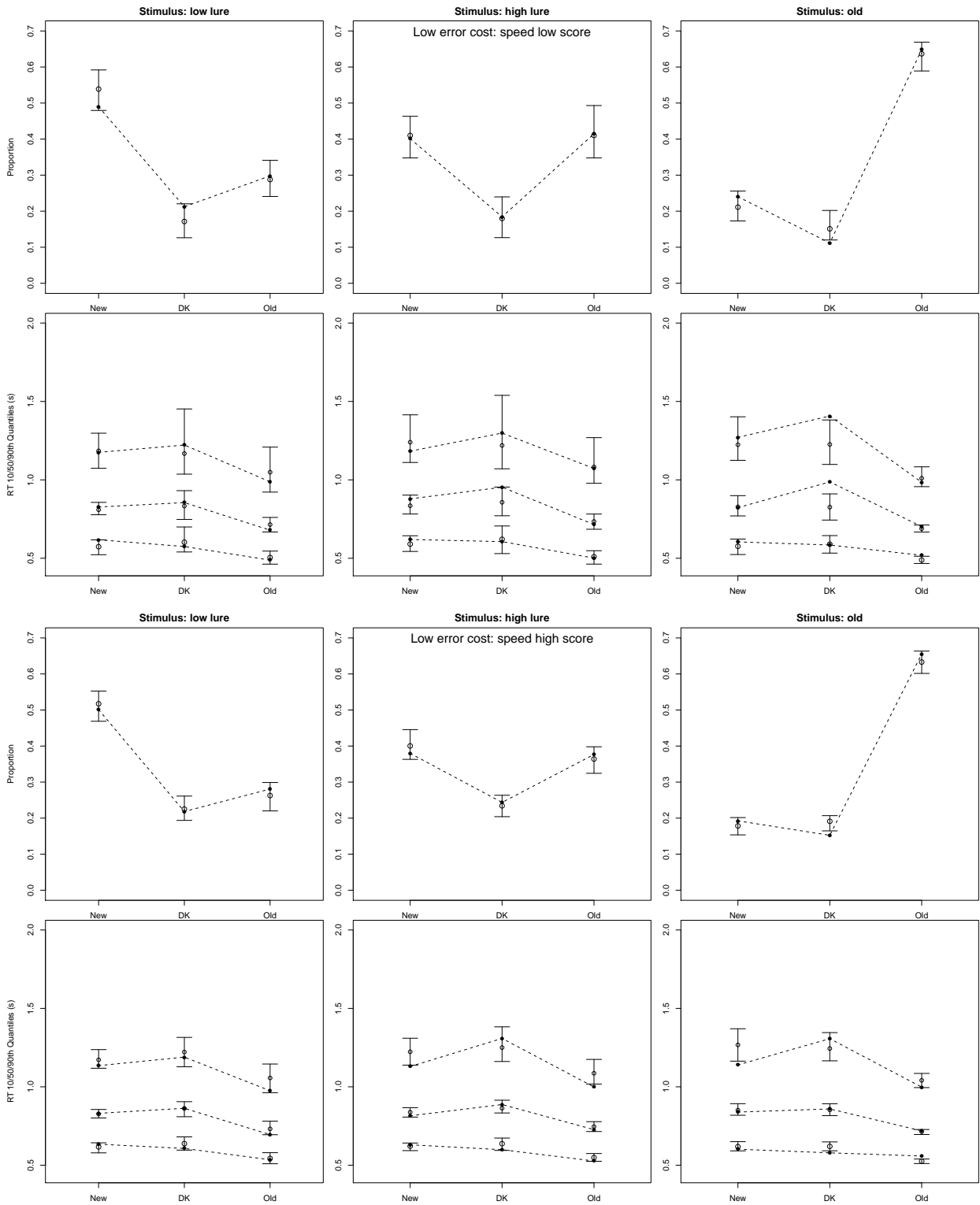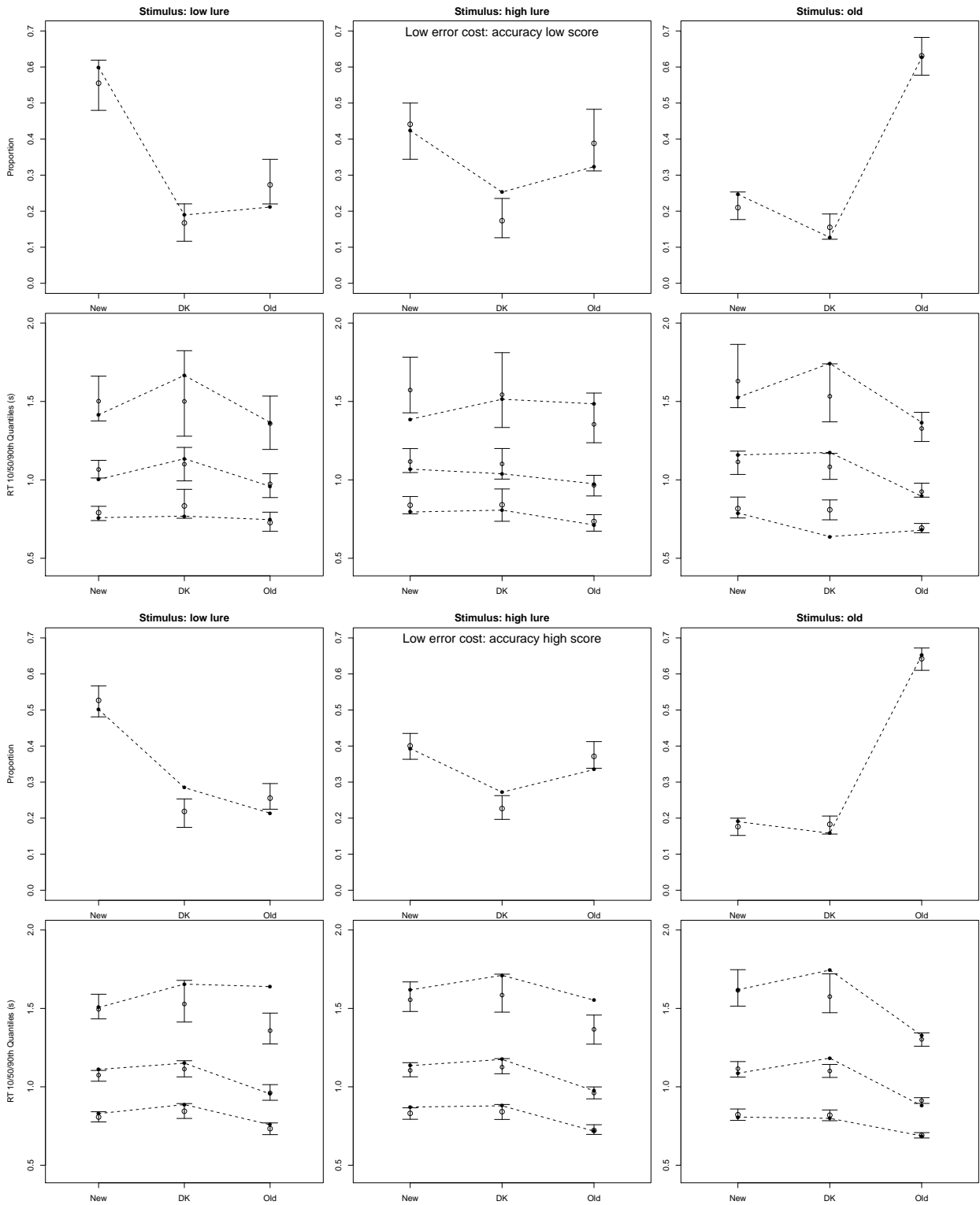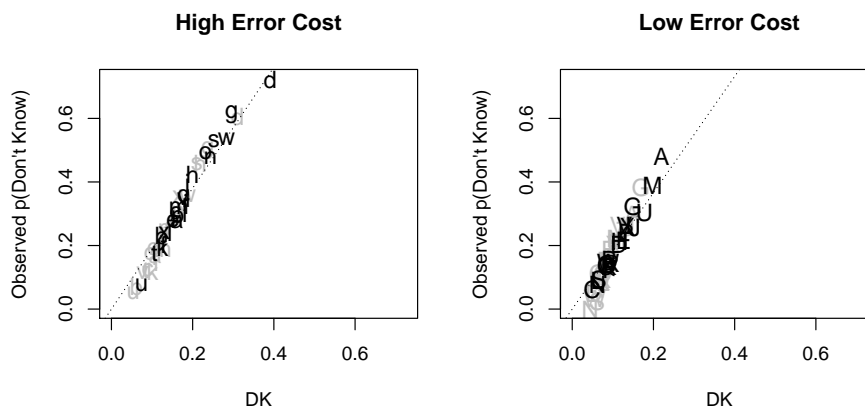
Figure S15: Experiment 2 MTR- 29 parameter model. Observed probability of a don't-know response for low-score score (grey letters) and high-score (black letters) trials as a function of the average of DK estimates over lure and target accumulators. Letters a ... x correspond to participants in the high error-cost condition and A ... X to participants in the low error-cost condition.

for high error cost under speed emphasis (ps <.001) but that reversed under accuracy emphasis (ps <.001).

Figure S18 shows that choice threshold estimates were greater under accuracy than speed emphasis (ps <.001), and also greater for the lure than target accumulator (ps <.001).

Figure S20shows that in most cases the match-accumulator rate was greater than the mismatch-accumulator rate (ps <.001) but there were exceptions. For hard lures under speed emphasis the two were almost equivalent in the low error-cost condition (p = .35) and in the high error-cost condition mismatch was greater than match (p <.001). For targets under accuracy emphasis there was also support for reversal (p = .015). The match-mismatch difference was always larger for easy than hard lures (ps <.001). It was also always larger under accuracy emphasis than speed emphasis, both for low error cost (by 0.52, p <.001) and high error cost (by 0.25, p <.001). Under accuracy emphasis rates favoured the lure accumulator over the target accumulator (i.e., larger rates for the matching accumulator for lure stimuli and the mismatching accumulator for target stimuli), for both high error costs (by 0.91, p <.001) and low error costs (by 0.36, p <.001). In contrast, under speed emphasis they favoured the target accumulator (i.e., larger rates for the mismatching accumulator for lure stimuli and the matching accumulator for target stimuli), for both high error costs (by 0.42, p <.001) and low error costs (by 0.16, p <.001).

Figures S11 - S14 show the global fit to the 29 parameter model to experiment 2. By a visual inspection, the 29 parameter model only provides a minor improvement over the 19 parameter model, with some reduction in the width of error bars and slightly fewer instances of a 95% interval missing the data on the 90th RT percentile.

Figure S16: Experiment 2 MTR- 29 parameter model. Observed probability of the different between don't-know responses for high and low score trials as a function of the difference in DK estimates between old and new accumulators. Letters a ... x correspond to participants in the high error-cost condition and A ... X to participants in the low error-cost condition.



Figure S17: Experiment 2 MTR parameter estimates-29 parameter model. Left hand panel shows $t_e r$ and $A$. The right hand panel shows the matching $sv$ estimates. Mismatching $sv$ was fixed to 1.



Figure S18: Experiment 2 MTR parameter estimates-29 parameter model. The left hand panel shows the $B$ distance values for the high error cost subjects. The right hand panel shows the low error cost subjects.

Figure S19: Experiment 2 MTR parameter estimates-29 parameter model. Top row of figures shows the $1 - D/b$ values for the high error cost subjects. The bottom row shows the low error cost subjects.

Figure S20: Experiment 2 MTR parameter estimates-29 parameter model. Top row of figures shows the mean drift rate *v* values for the high error cost subjects. The bottom row shows the low error cost subjects.

Figure S21: Average model fit to Experiment 2, speed emphasis, high error cost, 29 parameter paper model. Response proportions and 10, 50 and $90^{th}$ percentiles. The arrows show the 95% credible interval of the model predictions.

Figure S22: Average model fit to Experiment 2, accuracy emphasis, high error cost, 29 parameter paper model. Response proportions and 10, 50 and $90^{th}$ percentiles. The arrows show the 95% credible interval of the model predictions.

Figure S23: Average model fit to Experiment 2, speed emphasis, low error cost, 29 parameter paper model. Response proportions and 10, 50 and $90^{th}$ percentiles. The arrows show the 95% credible interval of the model predictions.

Figure S24: Average model fit to Experiment 2, accuracy emphasis, low error cost, 29 parameter paper model. Response proportions and 10, 50 and $90^{th}$ percentiles. The arrows show the 95% credible interval of the model predictions.

# 4   Simulation Studies

## 4.1   2AFC

Due to the presence of a stimulus bias in the 2AFC task, there are many potential models that could potentially describe the bias. Within the standard LBA framework, the bias towards left or right responses could be due to a threshold bias, a drift rate bias or both. Additionally, since the don't know winner is unobserved, we wanted to check which of these models can parameter recover. For the simulation studies we ignored the effect of difficulty, as we simply wanted to test how complex a mode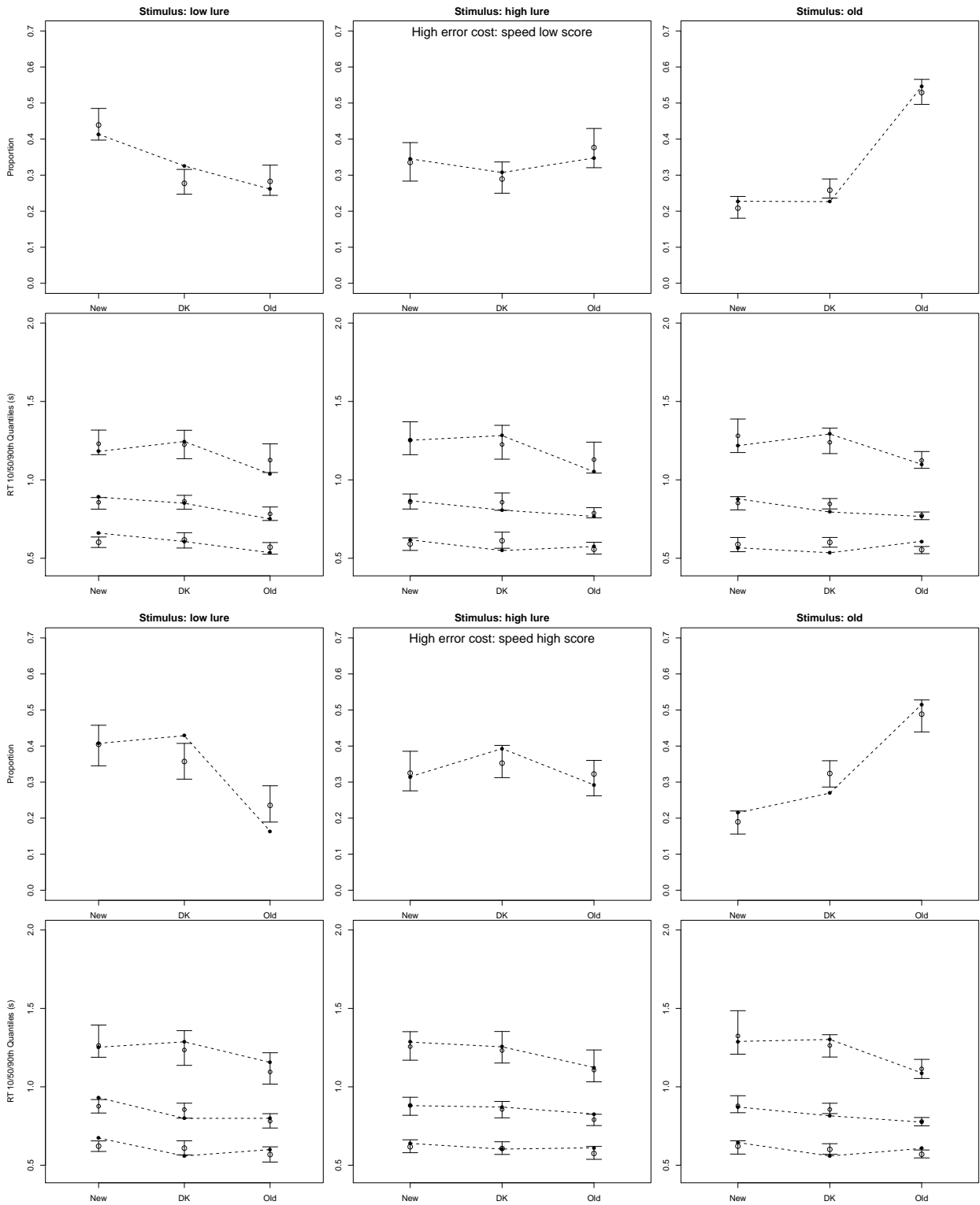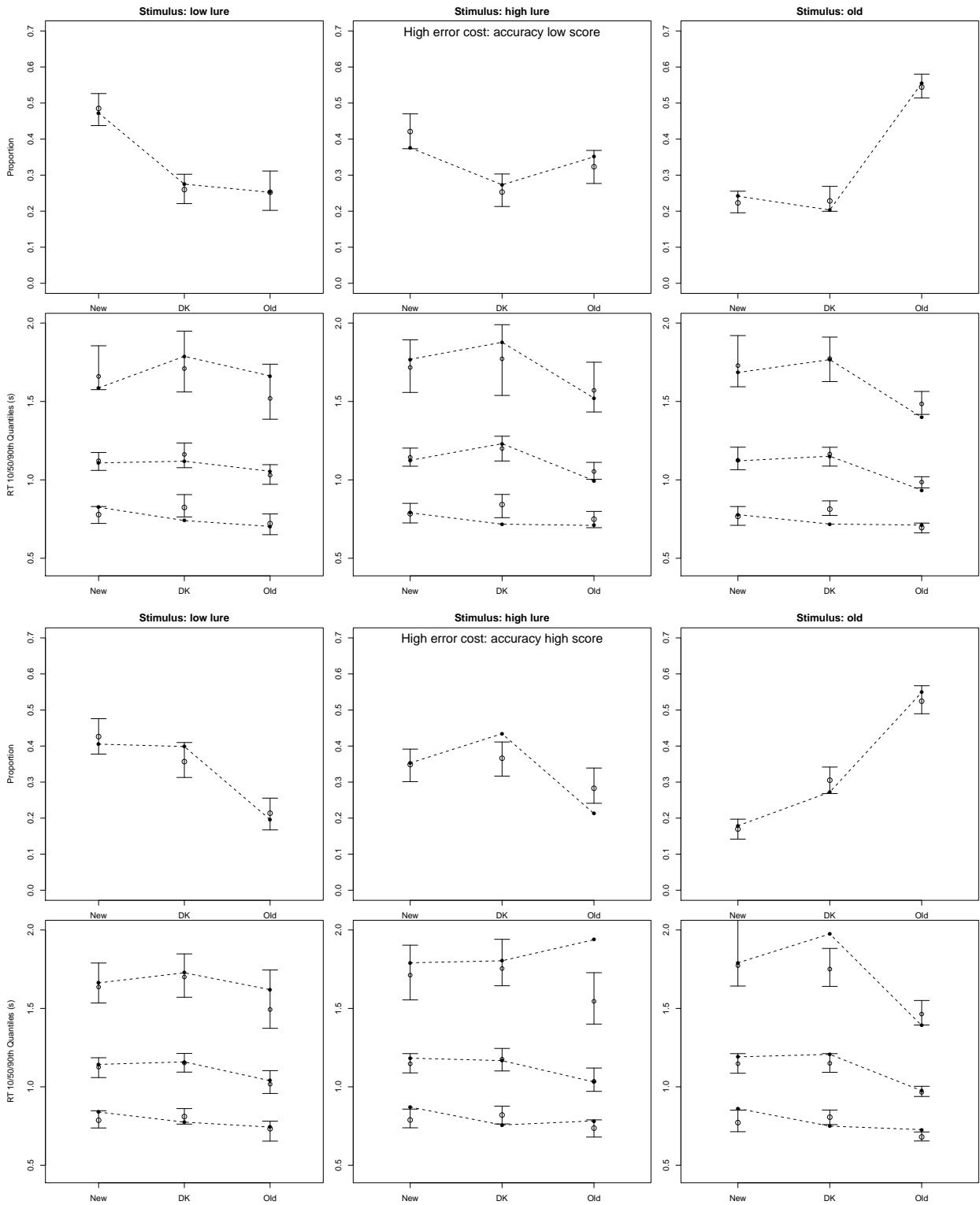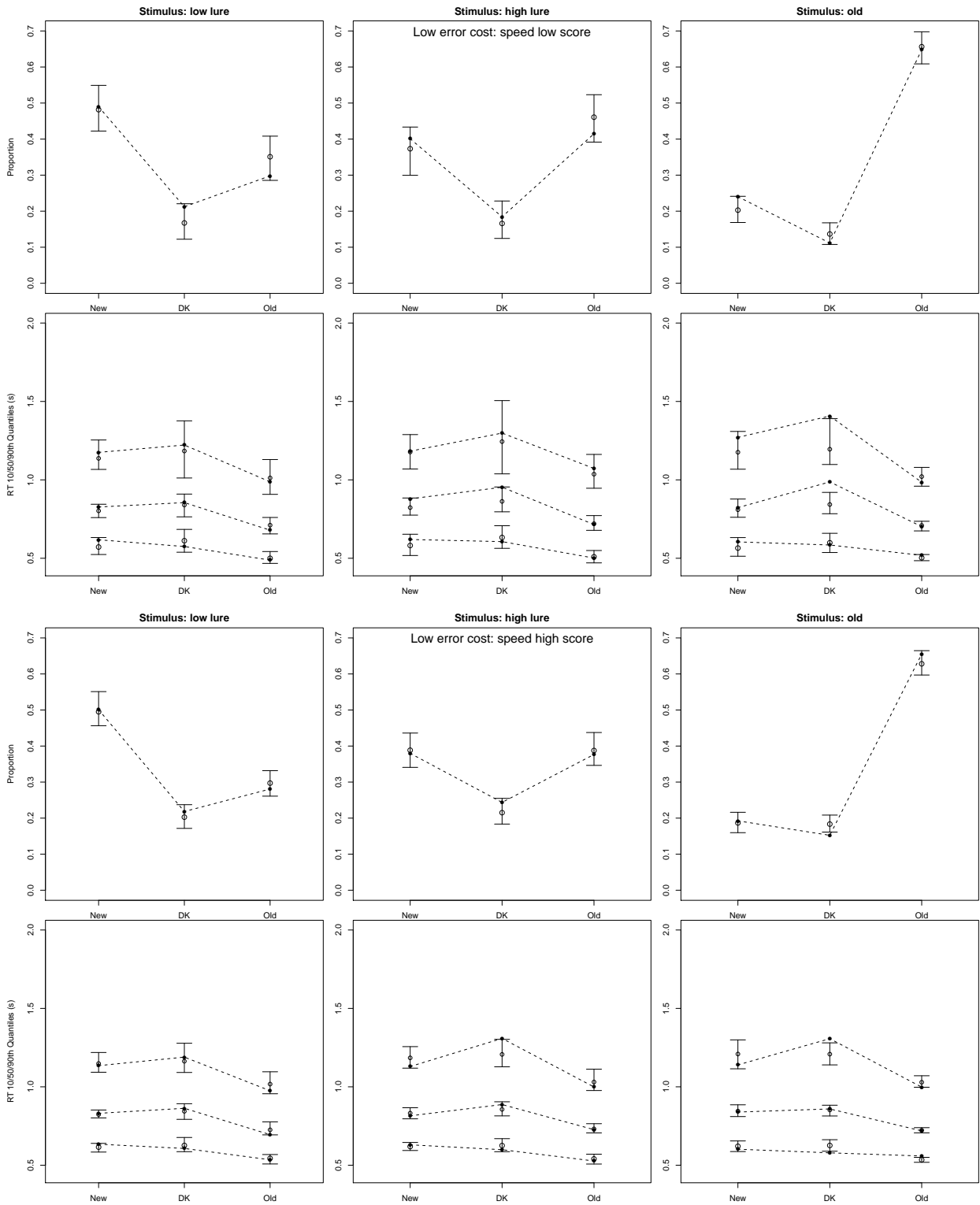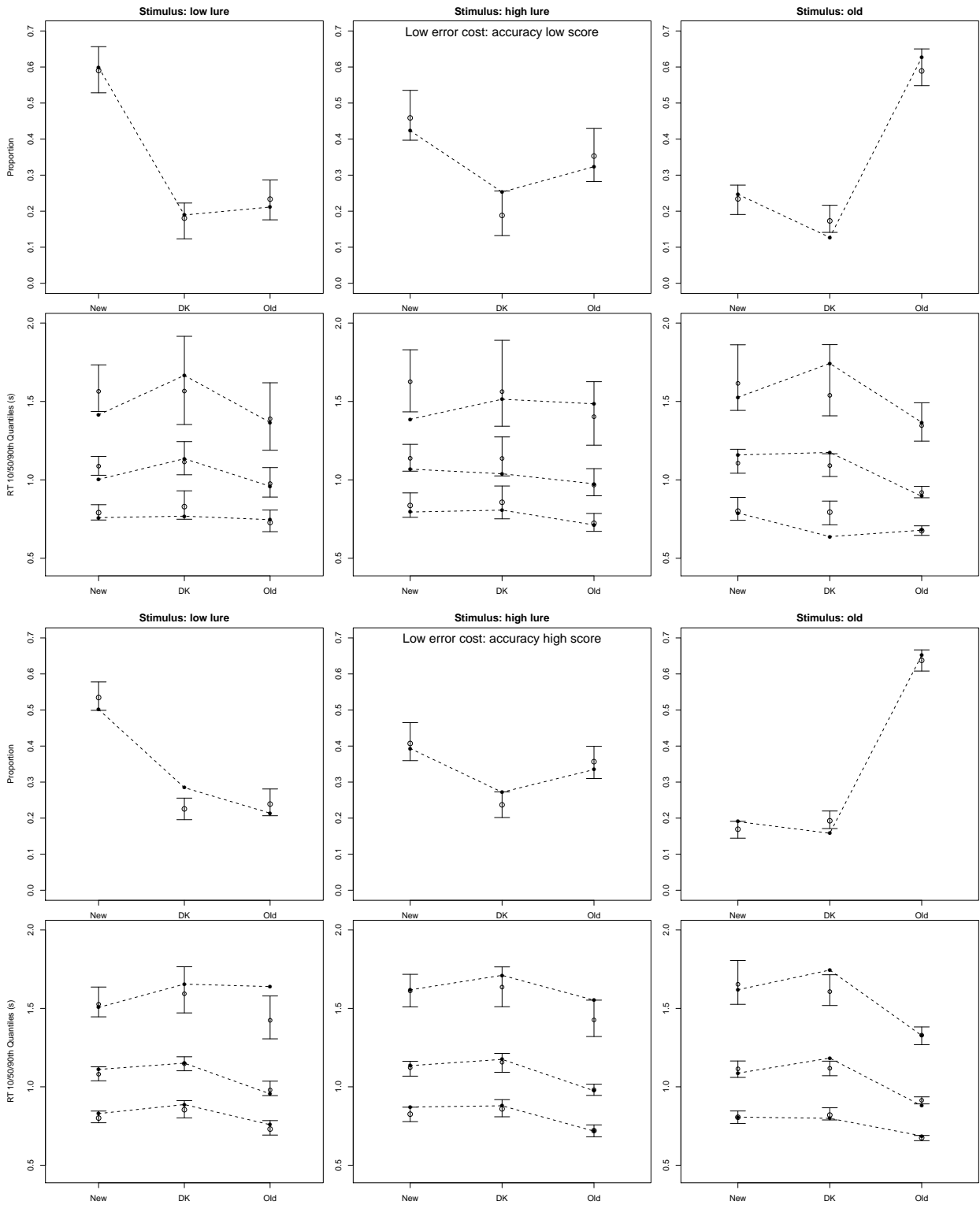l with just speed/accuracy manipulation and a large bias could be recovered, both of which are manipulations that can plausible effect both threshold and drift rate parameters. We assume that since difficulty can only plausibly manipulate drift rate parameters, that there would be no issue with adding difficulty manipulations to the drift rates. Parameters: The model parameterisation defined external parameters $A$, $B$ and $qD$. $A$ is the maximum value of start point noise. $B$ is the distance between $A$ and the threshold height $b$ (an internal parameter). The formula for $d$ is $d = plogis(qD) * b$, where plogis is the logistic distribution function. Hence $qD$ is the quantile function of the relative value of $d$ with respect to $b$. Note that for values of qD above 4 or bellow -4, increases in $qD$ cause minor shifts to the value of $d$. Finally note that we can constrain $qD$ more than $B$, and this will allow $d$ to vary over the same conditions that $B$ varies over, but the relative value of $d/b$ will be constrained.

The first simulation studied took a complex model to see if it could parameter recover. This model had a single start point noise parameter, but let response and speed emphasis effect the B and d thresholds. Stimulus, speed emphasis and response effected drift rates, creating 8 drift rate parameters. Finally the variance of the matching drift rate and the non-decision time were estimated. In total there were 19 parameters to estimate.

This model was applied to six sets of simulated trials. In the simulations, Don't Know thresholds were either fixed to be equal, biased towards one response or biased towards the opposite response. This was then repeated for drift rates either equal across stimulus or varying across stimulus. The 20 parameter model only successfully parameter recovered from simulations when the relative don't know threshold height was equal across response or when the drift rates were equal across stimulus. Note that some estimates show large ranges for qD parameters. These are q-logistic transformed proportions, so a difference from 4 to 8 means an effective proportion change from 0.98 to 0.99.

Figure S1 confirms that when the when the simulated value of $d/b$ is equal across each accumulator, that the model does parameter recover. Figures S2 and S3 shows that when there is a large

|  | 2.5% Estimate | True | 97.5% Estimate |
|---|---|---|---|
| A | 0.839 | 0.996 | 1.145 |
| B.speed.LEFT | 2.069 | 2.186 | 2.505 |
| B.accuracy.LEFT | 2.780 | 2.928 | 3.286 |
| B.speed.RIGHT | 2.071 | 2.200 | 2.520 |
| B.accuracy.RIGHT | 2.717 | 2.840 | 3.177 |
| qD.speed.LEFT | 1.190 | 1.309 | 1.465 |
| qD.accuracy.LEFT | 1.296 | 1.406 | 1.545 |
| qD.speed.RIGHT | 1.201 | 1.309 | 1.466 |
| qD.accuracy.RIGHT | 1.285 | 1.406 | 1.530 |
| v.left.speed.true | 3.078 | 3.130 | 3.288 |
| v.right.speed.true | 2.899 | 2.944 | 3.113 |
| v.left.accuracy.true | 3.321 | 3.368 | 3.527 |
| v.right.accuracy.true | 2.951 | 2.994 | 3.136 |
| v.left.speed.false | 2.340 | 2.422 | 2.602 |
| v.right.speed.false | 2.587 | 2.645 | 2.834 |
| v.left.accuracy.false | 2.271 | 2.340 | 2.502 |
| v.right.accuracy.false | 2.673 | 2.733 | 2.918 |
| sv.true | 0.821 | 0.829 | 0.841 |
| $t_{er}$ | 0.152 | 0.198 | 0.219 |

Table S1: MTR simulation study. 19 parameter model without difficulty. Simulated Drift rates differ by stimulus, simulated relative $d$ thresholds constant accross response.

|  | 2.5% Estimate | True | 97.5% Estimate |
|---|---|---|---|
| A | 0.008 | 1.341 | 0.442* |
| B.speed.LEFT | 2.194 | 3.121 | 2.521* |
| B.accuracy.LEFT | 3.194 | 4.071 | 3.578* |
| B.speed.RIGHT | 3.329 | 2.529 | 3.730* |
| B.accuracy.RIGHT | 4.021 | 3.091 | 4.438* |
| qD.speed.LEFT | 0.095 | 8.022 | 0.266* |
| qD.accuracy.LEFT | 0.523 | 7.630 | 0.696* |
| qD.speed.RIGHT | 2.571 | 0.453 | 3.015* |
| qD.accuracy.RIGHT | 2.235 | 0.592 | 2.695* |
| v.left.speed.true | 2.407 | 4.409 | 2.533* |
| v.right.speed.true | 3.844 | 2.771 | 3.996* |
| v.left.accuracy.true | 3.044 | 4.603 | 3.209* |
| v.right.accuracy.true | 3.689 | 2.731 | 3.865* |
| v.left.speed.false | 3.379 | 1.939 | 3.542* |
| v.right.speed.false | 1.811 | 4.120 | 1.956* |
| v.left.accuracy.false | 3.105 | 1.773 | 3.313* |
| v.right.accuracy.false | 2.320 | 4.171 | 2.513* |
| sv.true | 0.871 | 0.945 | 0.887* |
| $t_{er}$ | 0.112 | 0.151 | 0.143* |

Table S2: MTR simulation study. 19 parameter model without difficulty. Simulated Drift rates differ by stimulus, simulated larger $d$ on the left response accumulator.

|  | 2.5% Estimate | True | 97.5% Estimate |
|---|---|---|---|
| A | 1.653 | 1.341 | 1.879* |
| B.speed.LEFT | 3.305 | 3.121 | 3.869* |
| B.accuracy.LEFT | 4.315 | 4.071 | 4.928* |
| B.speed.RIGHT | 1.346 | 2.529 | 1.571* |
| B.accuracy.RIGHT | 1.464 | 3.091 | 1.700* |
| qD.speed.LEFT | 1.463 | 0.453 | 1.894* |
| qD.accuracy.LEFT | 1.719 | 0.592 | 2.189* |
| qD.speed.RIGHT | 0.776 | 8.022 | 0.993* |
| qD.accuracy.RIGHT | 0.723 | 7.630 | 0.921* |
| v.left.speed.true | 5.120 | 4.409 | 5.338* |
| v.right.speed.true | 0.465 | 2.771 | 0.862* |
| v.left.accuracy.true | 5.255 | 4.603 | 5.466* |
| v.right.accuracy.true | 0.251 | 2.731 | 0.658* |
| v.left.speed.false | -1.010 | 1.939 | -0.407* |
| v.right.speed.false | 4.986 | 4.120 | 5.215* |
| v.left.accuracy.false | -1.772 | 1.773 | -1.045* |
| v.right.accuracy.false | 5.008 | 4.171 | 5.235* |
| sv.true | 0.998 | 0.945 | 1.022* |
| $t_{er}$ | 0.113 | 0.151 | 0.178 |

Table S3: MTR simulation study. 19 parameter model without difficulty. Simulated Drift rates differ by stimulus, simulated larger $d$ on the right response accumulator.

bias in $d$ thresholds, that the models does not parameter recover with none of the true values within the 95% credible interval. This means that parameter recovery for this model is highly unstable and should not be trusted.

Next we test if a simpler model can parameter recover when the relative threshold heights are constant across accumulator. This still means that the $d$ values will be different if there is a $b$ threshold bias towards one response, but the value of $d/b$ will be constant.

The 17 parameter model drops the difference in $qD$ by response. Figure S4 confirms that this model can parameter recover, even when there is both a threshold bias and drift rate bias to account for.

Finally, we wanted to confirm that it was possible to recover parameters without requiring that the don't know thresholds be equal across response, and to fix drift rates by stimulus instead. Drift rates were fixed to be equivalent by stimulus but we added back the complexity by introducing a difficulty manipulation to drift rates. This is the 20 parameter model that was used in the main paper. Figure S5 confirms that this model does parameter recover. 5 out of 19 parameters were outside the expected region, but in each case the bias was negligible in magnitude, unlike the results in S2 and S3.

## 4.2   Experiment 2

A parameter study from the most complex 29 parameter model was run. We found that with asymptotic data, the model was well recovered. This is despite stimulus varying by both stimulus and response and thresholds varying by response too which caused problems in the 2AFC experiment.

|  | 2.5% Estimate | True | 97.5% Estimate |
|---|---|---|---|
| A | 0.923 | 0.996 | 1.168 |
| B.speed.LEFT | 1.896 | 2.186 | 2.260 |
| B.accuracy.LEFT | 2.643 | 2.928 | 3.062 |
| B.speed.RIGHT | 1.938 | 2.200 | 2.300 |
| B.accuracy.RIGHT | 2.549 | 2.840 | 2.939 |
| qD.speed | 1.276 | 1.309 | 1.316 |
| qD.accuracy | 1.389 | 1.406 | 1.429 |
| v.left.speed.true | 3.019 | 3.130 | 3.138 |
| v.right.speed.true | 2.882 | 2.944 | 2.999 |
| v.left.accuracy.true | 3.311 | 3.368 | 3.438 |
| v.right.accuracy.true | 2.911 | 2.994 | 3.035 |
| v.left.speed.false | 2.350 | 2.422 | 2.476 |
| v.right.speed.false | 2.527 | 2.645 | 2.658 |
| v.left.accuracy.false | 2.239 | 2.340 | 2.374 |
| v.right.accuracy.false | 2.656 | 2.733 | 2.797 |
| sv.true | 0.825 | 0.829 | 0.843 |
| $t_{er}$ | 0.184 | 0.198 | 0.243 |

Table S4: MTR simulation study. 17 parameter model without difficulty. Simulated Drift rates differ by stimulus, thresholds vary by response.

|  | 2.5% Estimate | True | 97.5% Estimate |
|---|---|---|---|
| A | 1.349 | 1.541 | 1.523* |
| B.speed.LEFT | 1.767 | 1.857 | 2.021 |
| B.accuracy.LEFT | 2.466 | 2.530 | 2.745 |
| B.speed.RIGHT | 1.943 | 2.031 | 2.203 |
| B.accuracy.RIGHT | 2.634 | 2.695 | 2.916 |
| qD.speed.LEFT | 1.558 | 1.581 | 1.634 |
| qD.accuracy.LEFT | 1.378 | 1.404 | 1.438 |
| qD.speed.RIGHT | 1.833 | 1.938 | 1.914* |
| qD.accuracy.RIGHT | 1.969 | 2.044 | 2.063 |
| v.easy.speed.true | 3.241 | 3.370 | 3.349* |
| v.hard.speed.true | 3.237 | 3.359 | 3.345* |
| v.easy.accuracy.true | 3.446 | 3.507 | 3.552 |
| v.hard.accuracy.true | 3.401 | 3.450 | 3.504 |
| v.easy.speed.false | 2.614 | 2.740 | 2.729* |
| v.hard.speed.false | 2.706 | 2.805 | 2.821 |
| v.easy.accuracy.false | 2.634 | 2.715 | 2.751 |
| v.hard.accuracy.false | 2.640 | 2.701 | 2.750 |
| sv.true | 0.904 | 0.913 | 0.925 |
| $t_{er}$ | 0.266 | 0.301 | 0.311 |

Table S5: MTR simulation study. 19 parameter model with difficulty effect but not stimulus effect on drift rates. Simulated Drift rates differ by stimulus, thresholds vary by response.

|  | 2.5% Estimate | True | 97.5% Estimate |
|---|---|---|---|
| A | 1.768 | 1.799 | 1.896 |
| B.accuracy.NEW | 2.259 | 2.410 | 2.418 |
| B.speed.NEW | 1.605 | 1.685 | 1.745 |
| B.accuracy.OLD | 1.551 | 1.622 | 1.694 |
| B.speed.OLD | 1.400 | 1.506 | 1.543 |
| qD.low.accuracy.NEW | 2.117 | 2.454 | 2.504 |
| qD.high.accuracy.NEW | 1.726 | 1.922 | 1.980 |
| qD.low.speed.NEW | 1.464 | 1.478 | 1.537 |
| qD.high.speed.NEW | 1.069 | 1.101 | 1.124 |
| qD.low.accuracy.OLD | 1.562 | 1.609 | 1.803 |
| qD.high.accuracy.OLD | 1.090 | 1.104 | 1.268 |
| qD.low.speed.OLD | 4.534 | 6.375 | 9.642 |
| qD.high.speed.OLD | 4.911 | 6.234 | 9.064 |
| v.easy.accuracy.NEW | 3.104 | 3.221 | 3.243 |
| v.hard.accuracy.NEW | 2.906 | 3.027 | 3.042 |
| v.old.accuracy.NEW | 2.532 | 2.679 | 2.691 |
| v.easy.speed.NEW | 3.155 | 3.189 | 3.291 |
| v.hard.speed.NEW | 2.920 | 2.973 | 3.054 |
| v.old.speed.NEW | 2.526 | 2.552 | 2.658 |
| v.easy.accuracy.OLD | 1.713 | 1.747 | 1.915 |
| v.hard.accuracy.OLD | 1.870 | 1.905 | 2.053 |
| v.old.accuracy.OLD | 2.516 | 2.553 | 2.693 |
| v.easy.speed.OLD | 2.946 | 3.062 | 3.098 |
| v.hard.speed.OLD | 3.191 | 3.265 | 3.339 |
| v.old.speed.OLD | 3.532 | 3.630 | 3.686 |
| sv.easy.true | 0.851 | 0.883 | 0.889 |
| sv.hard.true | 0.897 | 0.903 | 0.938 |
| sv.old.true | 1.029 | 1.028 | 1.067 |
| $t_{er}$ | 0.196 | 0.203 | 0.221 |

Table S6: MTR simulation study. 29 parameter model from experiment 2. Simulated Drift rates differ by stimulus, thresholds vary by response.

Speculation: because there are three types of stimuli and still only two responses the model is less able to flip parameters around.

## 4.3   LBA vs MTR parameter recover

A don't know experiment gives slightly more information about responses given we know when uncertain decisions are made (and thus that both responses must have been above the intermediate threshold when the response was made), but less information is given in that we don't know which of the two responses finished first. To compare estimation properties we simulated don't know data using the design from Experiment 1 and then fit the data with an MTR, and a binary LBA model with the don't know responses identified. We found the models recovered the LBA parameters were recovered with approximately equal bias and uncertainty by both models, suggesting that there is no difference in terms of simple estimation of parameters between binary and don't-know models.

|  | Bias | | SD | |
| --- | --- | --- | --- | --- |
|  | LBA | MTR | LBA | MTR |
| A | 0.29 | 0.32 | 0.24 | 0.26 |
| B.speed.LEFT | -0.31 | -0.27 | 0.30 | 0.32 |
| B.accuracy.LEFT | -0.32 | -0.29 | 0.33 | 0.35 |
| B.speed.RIGHT | -0.30 | -0.27 | 0.31 | 0.32 |
| B.accuracy.RIGHT | -0.34 | -0.29 | 0.34 | 0.36 |
| v.easy.speed.true | -0.03 | 0.00 | 0.14 | 0.16 |
| v.hard.speed.true | -0.02 | 0.01 | 0.14 | 0.16 |
| v.easy.accuracy.true | -0.02 | 0.01 | 0.14 | 0.16 |
| v.hard.accuracy.true | 0.00 | 0.04 | 0.14 | 0.16 |
| v.easy.speed.false | -0.04 | 0.00 | 0.15 | 0.16 |
| v.hard.speed.false | -0.05 | -0.01 | 0.15 | 0.16 |
| v.easy.accuracy.false | -0.04 | 0.00 | 0.15 | 0.16 |
| v.hard.accuracy.false | -0.03 | 0.01 | 0.15 | 0.16 |
| sv.true | 0.00 | 0.00 | 0.03 | 0.03 |
| $t_{er}$ | 0.05 | 0.05 | 0.05 | 0.05 |

Table S7: LBA and don't know MTR model parameter recovery bias and posterior standard deviations.