

## Response to Reviews

Is there a Criterion in Criterial Learning?  
Insights from Studying Feedback Delays

*Journal of Experimental Psychology: Learning, Memory, and Cognition*  
Manuscript number: XLM-2024-0436

Matthew J. Crossley, Benjamin O. Pelzer, F. Gregory Ashby

We would like to thank the editor and reviewers for their excellent suggestions and thoughtful comments. We worked very hard to respond to all of these. As a result, we are confident that the revised manuscript is significantly stronger than our original submission. Please note that we changed the title of the manuscript to more directly convey the main point of our paper.

Our specific response to each reviewer comment is as follows.

### Reviewer 1

We were pleased that Reviewer 1 “really enjoyed reading the manuscript” and “highly recommend(ed) this paper to be published after addressing a number of mostly minor concerns.”

#### Reviewer comment 1

1. Exploration of model space: When reading the description of the models it was not entirely clear to me whether the three presented models explore the space of possible models to a sufficient degree or whether there might be some models missing from consideration that could change the overall message.

#### Author response

This is an excellent and important point. In response, we greatly expanded the set of models we explored, and we also now classify these models more systematically. The previous version examined 3 models, whereas the revision examines 12. These 12 models are partitioned into two classes: one that assumes a criterion is learned and one that assumes SR learning with no criterion. The criterial-learning class includes 8 different models that were created by factorially combining three different binary assumptions: perceptual drift (yes or no), criterial drift (yes or no), and delay-sensitive updating (yes or no). The SR learning class

includes four models created by factorially combining two binary assumptions: perceptual drift (yes or no) and delay-sensitive feedback (yes or no).

## **Reviewer comment 2**

2. Simulation results figures: There were some aspects of the results figures for the simulation that I found difficult to understand:

- (a) Even after reading the full paper and thinking about it for a while, I do not understand what panels B are supposed to show. Does this only show the subset of the parameter space in which one of the two relevant patterns was observed? If not, shouldn't there be data points everywhere?
- (b) I do not understand why in Figure 1, for example, panel C shows some cases in which the predicted qualitative pattern is either "control impaired" or "other", but there is no corresponding bar in panel A.
- (c) The bars in Figure 2A are smaller than in Figures 1A and 3A. Why?

## **Author response**

With the new approach to model exploration, we have heavily revised the simulation results figures. We hope that the new figures are clearer and address the reviewer's concerns.

## **Reviewer comment 3**

3. Calling the simulation Experiment 1: I found it extremely confusing that the simulation using parameter-space partitioning was called "Experiment 1". If I read "Experiment" I expect to see some empirical data obtained from participants and not solely a simulation. I suggest renaming this section to something clearer. For example "Simulation Study" or "Parameter-Space Partitioning Across Models".

## **Author response**

After thinking about this a bit, we agree that referring to the computational modeling analyses as Experiment 1 was a mistake. First, we agree with the reviewer that this nomenclature is confusing, but we also came to realize that the equations in this section might dissuade empirically oriented readers from reading about our empirical results. As a result, the revised manuscript renames the section describing the computational modeling "New Computational Models of Criterial Learning" and this more technical section has been moved to later in the manuscript – after the two empirical experiments are described.

## **Reviewer comment 4**

4. Outliers in Exp. 2:

- While I think the argument for removing outliers in Exp. 2 is fine, the actual criterion defining an outlier is not actually given (which feels a bit ironic in a paper about criteria).

### Author response

Good point. We now say “for each participant, we used the best-fitting two component model to compute the likelihood that their performance was best described by each component. Participants whose data were estimated to most likely belong to the low-performance distribution were classified as outliers and excluded from further analysis.”

### Reviewer comment 5

4. Outliers in Exp. 2:

- How do the results look with these outliers included? Does the difference between condition vanish? Adding a footnote or so should be enough.
- The fact that some participants did not even learn a single criterion feels very surprising to me. Is this common in such experiments? Please add a sentence embedding this finding within your general experience running such experiments.

### Author response

????

### Reviewer comment 6

Minor points:

- Does Eq. 13 show results for trial  $n$  or  $n + 1$ ? The text and equation are inconsistent.

### Author response

We corrected this typo.

### Reviewer comment 7

- The description regarding the randomisation of Experiment 2 on page 12 (paragraph before “Procedure”) seems to be inconsistent with Figure 4. Is the dimension picked at random for each trial or always the same for certain problem numbers (as implied by Figure 4)?

### Author response

We agree that the description of the design of this experiment (now Experiment 1) needed improvement. We revised the paragraph Reviewer 1 refers to in this comment as follows.

“Each participant practiced each of 14 one-dimensional category-learning tasks, or problems, until they responded correctly on 9 of the 10 previous trials, at which point the problem changed. The 14 different category structures are described in Figure 1. The relevant dimension was bar thickness in problems 1 – 7 and bar angle in problems 8 – 14. Each problem included stimuli in two distinct clusters that varied over a restricted range of the

relevant dimension. Critically, the optimal criterion value varied from problem-to-problem with respect to its position within the stimulus range. For example, for some problems the optimal criterion was below the midpoint of the range and for other problems it was above the midpoint.”

### **Reviewer comment 8**

- In figure 4, the top says Problems “1 thru 13”, but only shows 7 problems.

### **Author response**

We corrected this typo.

## **Reviewer 2**

We were pleased that Reviewer 2 believed our “manuscript reports a sophisticated approach to an important problem” and that the reviewer “especially like(d) the use of mathematical models.”

### **Reviewer comment 1**

I was concerned that the models seem to have multiple varying features. It might be more useful to isolate specific processes while holding other model elements constant.

### **Author response**

This is the same as Reviewer 1’s first comment. As described there, we strongly agree with this concern, and as a result, we greatly expanded the number of models we examined, and we reclassified these according to the fundamental assumption of whether or not they assumed a criterion value is specifically learned and stored trial-by-trial in memory.

### **Reviewer comment 2**

I think it makes more sense to have a prediction simulation section and Experiments 1 and 2 rather than to consider the simulation the first experiment.

### **Author response**

This is the same as Reviewer 1’s third comment. As noted there, we agree and as a result, we now refer to the two empirical experiments as Experiment 1 and 2, and the computational modeling section has been renamed to “New Computational Models of Criterial Learning” and moved to after the two experiments.

### **Reviewer comment 3**

The data in Figure 6 can likely be fit much more closely by an ExGaussian distribution than either the 1 or 2 component Gaussian mixtures that are shown.

#### **Author response**

We thank the reviewer for this suggestion. We now estimate both a single-component truncated ExGaussian and a two-component truncated ExGaussian mixture via MLE with proper normalization on the tasks bounded support [9, 512], and compared them via AIC, BIC, and a parametric bootstrap likelihood-ratio test (LRT). The results from moving from the mixture of Gaussians used in our original submission to the truncated ExGaussian mixture used in the revision does not change the qualitative pattern of our results in any way. However, we agree that the ExGaussian and ExGaussian mixture better capture the skewed and bounded distribution of our trials-to-criterion data so we have adopted it in the revision.

### **Reviewer comment 4**

Replication of the criterial learning results from the currently-labeled Experiment 2 seems like a good idea given the fairly low sample sizes per condition.

#### **Author response**

We agree with the reviewer that a replication would be valuable, but we believe the data and the expanded modeling analysis currently presented represent a significant contribution to the field.

### **Reviewer comment 5**

Did stimuli remain on the screen across the feedback delay? I don't think so. If not, what happens if the stimulus reappears at the time of feedback? I could imagine this would eliminate or dramatically reduce the negative effect of feedback delay. If so, is that consistent with the conclusions? That is, would the models that predict a deleterious effect of feedback delay also predict that it should be eliminated with "reminder" stimuli at the time of feedback?

#### **Author response**

As shown in Figures 2 and 5, the stimulus was presented up until a response was made. At this time the stimulus was removed and replaced by a noise mask until feedback was delivered. The question of how our results might have changed if the stimulus were to reappear at the time of feedback is both interesting and insightful. In the very first study that investigated the effects of feedback delays on category learning (Maddox, Ashby, & Bohil, 2003, JEP:LM&C) we found that a short 2.5 s feedback delay impaired information-integration category learning, but not rule-based learning – but only if a noise mask was presented during the time between response and feedback. The published 2003 article doesn't explore this

issue, but we hypothesized at the time that in the absence of the mask, participants might be able to use eidetic imagery to keep the relevant synapses active during the delay period. As a result, we agree with the reviewer’s prediction that re-presenting the stimulus at the time of feedback would eliminate the impaired learning that we observed. This interesting point seems tangential to our main point, however, and as a result, we did not address this point in the revision.

## **Reviewer comment 6**

For the currently-labeled Experiment 3, reporting fail-to-reject results for t-tests is not a valid method for establishing that feedback delay has no effect. This is especially a concern given the low participant count per condition and the fact that empirical effects go in the “right” direction (lower performance for delayed feedback).

### **Author response**

This is a completely valid point. As in almost all psychological experiments, our Experiment 2 lacks the statistical power to make strong inferences about the null results we obtained there. We now explicitly acknowledge this point in paragraph 2 of the section entitled “Discussion of Experiments 1 and 2.” Also, rather than try to make a statistical argument in support of our null results, we now use a converging operations approach by appealing to the literature. The new relevant paragraph is as follows.

“Whereas the results of Experiment 1 seem clear, the Experiment 2 results are weaker since they suggest a null result. Although insufficient power makes it impossible to rule out the possibility that Experiment 2 missed some small effect of feedback delay, the null results found there are consistent with a wide variety of other evidence. In particular, a variety of other studies have reported that feedback delays of up to 10 s do not significantly slow one-dimensional rule-based learning (Dunn et al., 2012; Ell et al., 2009; Maddox et al., 2003; Maddox & Ing, 2005; Worthy et al., 2013). In addition, many studies have reported evidence that one-dimensional rule-based learning recruits working memory and executive attention, but not procedural learning (for reviews, see, e.g., Ashby & Valentin, 2017; Ashby et al., 2020). So the results of Experiment 2 are consistent with all of these studies. On the other hand, several of these feedback-delay studies did not use binary-valued stimulus dimensions, so some criterial learning was required. If so, then why did they all report no effect of feedback delay? It is important to note that criterial learning was not the focus of any of these studies, and as a result, in all of these previous studies the optimal criterion was always set exactly midway between the category prototypes, which makes criterial learning trivial. For example, under these conditions, criterial learning might not even require feedback. The unsupervised category-learning experiments reported by Ashby et al. (1999) provide strong support for this because all of their rule-based participants learned the correct criterion (which was midway between the category means), even though the task was completely unsupervised.”

## **Reviewer comment 7**

I recommend running a bigger replication study that combines both tasks to demonstrate the claimed interaction (feedback delay impairs performance only for the criterion learning task). The replication study should be preregistered with a sample size justification based on the interaction effect defined by the across-experiment comparison in the current data. Adding a “stimulus reminder at feedback” condition to this new study might be helpful in clarifying the mechanisms producing the delay effect.

### **Author response**

We agree that a replication/extension study would strengthen our manuscript – as it would for almost all empirical papers. We chose not to pursue this suggestion, however. Even so, we believe we did significantly strengthen our manuscript by greatly expanding the class of computational models that we investigated, and by better situating our results within the existing literature.

## **Action Editor**

### **AE comment 1**

The final paragraph of the introduction summarizes the experiments. This seems more appropriate in the abstract, given that no information about the experiments are presented in the abstract.

### **Author response**

Thanks for this suggestion. The revised abstract now has more detail about the experiments.

### **AE comment 2**

Equation 1 needs some explanation, as this implies that the decision  $R_n$  is made immediately on stimulus presentation. In other words,  $R_n$  is not an outcome of some time-extended cognitive process. Am I misunderstanding equation 1?

### **Author response**

This is a valid and insightful point. Equation 1 does assume that the response is made immediately upon stimulus presentation. Because we are not trying to account for response time, we decided to not explicitly model the time-extended cognitive processes that govern evidence accumulation. We now state this explicitly in the revision. Specifically, immediately after equation 1, we added the following text:

“Note that this equation indicates that the response is made immediately upon stimulus presentation. We therefore are not explicitly modeling the time-extended cognitive processes that govern evidence accumulation and the corresponding response time.”

### **AE comment 3**

What was the reason behind the choice to have equation 2 only be sensitive to negative feedback?

#### **Author response**

The rationale behind the choice to have equation 2 only be sensitive to negative feedback is now stated explicitly in the first two sentences of the paragraph that includes equation 2. Those sentences read as follows.

“If positive feedback is received, then the value of the criterion remains unchanged for the next trial (except for drift – see equation 3). The rationale here is that if the response is correct, then the observer has effectively gained zero information about how their criterion should be modified.”

### **AE comment 4**

How is “optimal” defined in the current case, given that response time does not factor in the update rule?

#### **Author response**

By “optimal”, we mean the value of the criterion that maximizes accuracy. We now say this explicitly in the revision. This is now explained in the last sentence of the paragraph that includes equation 2, which reads as follows.

“This updating rule will gradually converge on the optimal criterion value (i.e., the value of the criterion that maximizes accuracy).”

### **AE comment 5**

The use of parameter space partitioning is sensible given the many differences among the three models. However, PSP operates on non-stochastic models. In addition, the way you have approached PSP is not how the original PSP is implemented (i.e., using MCMC). This should be mentioned explicitly to avoid confusion.

#### **Author response**

Thanks for raising this important point. We have revised the manuscript to make this point clear by adding the following text to the beginning of the third paragraph of the section entitled “Simulation Results.”

“Classical PSP is defined for non-stochastic models and is often implemented by exploring the parameter space via a Markov chain Monte Carlo algorithm. In contrast, we used brute force search over a grid of plausible parameter values. Because our models are stochastic, at each sampled set of parameter values we generated multiple simulations and classified the results into one of a small set of qualitative outcome patterns. ”



## **AE comment 6**

The three-dimensional plots in figures 1B, 2B, and 3B do not come out right on 2D. You could improve these by having 2D projections on the planes. Better still, three boxplots would be more appropriate as there is no additional information by using 3D plots.

### **Author response**

This was a good suggestion. We revised the two figures (8 and 9) that describe the simulation results accordingly – i.e., both figures are 2D and include boxplots.

## **AE comment 7**

Although not a focus in this manuscript, I wondered whether the models are able to account for finer-grained data patterns such as complete RT distributions and speed-accuracy tradeoff functions.

### **Author response**

As described in our response to AE comment 2, the models make no attempt to account for response time (RT). The models could all be generalized to make RT predictions, but any such generalization would require adding extra assumptions about the time course of the decision-making process. The focus of our research was on *what* is updated trial-by-trial (i.e., a criterion vs SR associations), rather than on the dynamics of the updating process. For this reason, we chose a theoretically minimal landscape. Specifically, we were concerned that adding any extra assumptions that were not directly related to our primary goal would weaken the inferences possible from our model simulations. For example, it would be difficult to rule out the possibility that any qualitative result we established was a consequence of one of these added assumptions, rather than because of a difference in what was being updated on each trial. Even so, we agree that extending the models proposed here to make RT predictions could be a fruitful avenue for future work.

## **AE comment 8**

The experiments are numbered 2 and 3, instead of 1 and 2. The models were following the methods of experiment 1, which is missing. Experiment 1 is also mentioned in the discussion separately with experiment 2. It looks like experiment 1 has been deleted from this manuscript, but the model simulations and discussion still refer to it.

### **Author response**

We thank the reviewer for catching this confusion. No study was deleted from the manuscript. Rather, whether or not the computational modeling section should be considered an “experiment” was a point of oscillation when writing. This comment, along with comment 3 of Reviewer 1 and comment 2 of Reviewer 2 helped us see that referring to the modeling section as Experiment 1 was a mistake. As explained in our responses to both of those comments,

the revision refers to the two empirical experiments as Experiments 1 and 2. These are now presented first, and then following by a modeling section entitled “New Computational Models of Criterial Learning.”

## **AE comment 9**

In figure 6, it is clear that the participants form two subcohorts. However, the use of a Gaussian seems odd given that there can not be any negative values for trials-to-criterion. Given that the authors model this frequency histogram to exclude participants, more appropriate distributions, such as Poisson should be used.

### **Author response**

This is a good point that is similar to comment 3 of Reviewer 2. Rather than use a mixture model based on the Poisson distribution, we took the advice of Reviewer 2 and used the ExGaussian as our base distribution. Specifically, as described in our response to comment 3 of Reviewer 2, we now fit the trials-to-criterion data collapsed across conditions and participants with both a single-component truncated ExGaussian and a two-component truncated ExGaussian mixture via MLE with proper normalization on the tasks bounded support [9, 512].

## **AE comment 10**

Given the experimental results, what is the added benefit of the three models? Experiment 2 is the critical experiment showing that 3.5-0.5 leads to more trials-to-criterion than 0.5-3.5 or 0.5-0.5. The models do not provide any detailed understanding of mechanisms and they are not fitted to actual data. The general message from the experiments is that any updating rule should be sensitive to feedback delay. The models do not go beyond this.

### **Author response**

The experiments show that feedback delay impairs criterial learning, but they do not speak to the question of underlying mechanism. The modeling is critical to this question. Specifically, the computational analysis attempts to establish constraints on underlying mechanisms that are compatible with the empirical results. We believe that PSP is a better approach, given this goal, than the more traditional approach of fitting the models to the observed data. This is now explained in the following new paragraph (paragraph 2 of the “Simulation Results” section):

“Our approach was to generate predictions for each of the three experimental conditions across a wide range of parameter values. We then used parameter-space partitioning (PSP) to evaluate the performance of each model (Pitt et al., 2006). Traditional computational modeling identifies parameter values that allow the model to provide the best possible fit to the observed data. This traditional approach therefore treats the single observed data set being fit as a gold standard. Our goal instead was to investigate the range of possible predictions of each model. We were interested in questions such as: can a model predict any

possible outcome or is it constrained to always predict, for example, that feedback delays impair learning? PSP was designed to answer such questions.”

Furthermore, we believe this PSP analysis produced some surprising and important results. First, it showed that SR-learning accounts are compatible with our results – even if they assume that the amount of synaptic updating that occurs trial-by-trial is unaffected by feedback delay. Second, the PSP analysis showed that, in contrast, stored-criterion accounts of Experiments 1 and 2 must assume that feedback delays impair the criterion-updating process – a result that greatly constrains the type of memory that could be used to store and update the criterion. For example, it rules out working memory because many previous results show that working memory is immune to short 3 s delays. We believe these inferences would be impossible without our computational modeling approach. In this way, the models contribute the most important piece: they pin down the process-level mechanisms consistent with the experiments.