**University of Connecticut**
**OpenCommons@UConn**

8-24-2012

# Reducing Knowledge Overconfidence by Reducing the Threat of Knowledge Cue Utilization

Christopher Neil Burrows

*University of Connecticut - Storrs,* christopher.burrows@uconn.edu

Reducing Knowledge Overconfidence by

Reducing the Threat of Knowledge Cue Utilization

Christopher Neil Burrows

B.A., Southern Methodist University, 2007

A Thesis

Submitted in Partial Fulfillment of the

Requirements for the Degree of

Master of Arts (or Science)

at the

University of Connecticut

**2012**

# Master of Arts Thesis

# Reducing Knowledge Overconfidence by

# Reducing the Threat of Knowledge Cue Utilization

Presented by

Christopher Neil Burrows, B.A.

Major Advisor _____

Hart Blanton

Associate Advisor _____

Kerry L. Marsh

Associate Advisor _____

Vicki J. Magley

University of Connecticut

**2012**

Abstract

Overconfident judgments are common. We are often more confident about things than we should be, and this may lead us to make maladaptive decisions. Debiasing confidence by cuing people in to how confident they should be could help people make better choices. However, people may be unwilling to accept debiasing information if doing so implies their own ignorance. This study examined whether self-affirmation can buffer people against threats to self-image, helping people to accept debiasing cues. I hypothesized that combining a cue with self-affirmation would lead to enhanced debiasing over cues or self-affirmation alone. In order to investigate this hypothesis, first a pilot study was used to create veridical cues for implementation in the main experimental design. The experiment used a memory task in which participants were asked to rate their confidence for future recall. The experiment used a 2 (cue) x 2 (self-affirmation) experimental design. No evidence for the effect of self-affirmation, or the interaction between cue and self-affirmation was found. Consequently, the experimental hypothesis was not supported. Future research may seek to investigate whether a stronger self-affirmation manipulation, or whether using a task with greater inherent self-threat, changes these results.

**Reducing Knowledge Overconfidence by**

**Reducing the Threat of Knowledge Cue Utilization**

Evidence suggests that people are often more confident than knowledgeable (Fischhoff, Slovic, & Lichtenstein, 1977). Such knowledge overconfidence can have consequential effects, if it leads to maladaptive choices. Consider as an example a college student who incorrectly concludes that she has mastered material for her next exam. As a result of this misperception, she decides to stop studying too soon and goes on to perform poorly on her exam. Her decision to stop studying early could be interpreted as evidence that she had too much confidence in her knowledge. Ideally, her appraisals of her own knowledge would be more closely linked to her mastery of the material, such that she only stops studying when she knows enough to meet her goals on the exam. This example highlights the potential value in trying to "debias" confidence judgments, so that individuals make more adaptive choices.

One obvious first step in debiasing overconfidence might be to present some kind of information suggesting that a correction is necessary. With the college student, for instance, one might present objective information to her suggesting that she does not know what she believes she does know. Evidence suggests, however, that such information might fail to influence the decisions of many students. Pulford and Colman (1997) found, for instance, that overconfidence can persist even after individuals receive immediate feedback revealing their lack of mastery on a knowledge test. An alternative strategy might be to provide the student with less personal information suggesting to her that most students like her perform poorly on the test she will take. Research suggests, however, that people also fail to use external decision aids that might guide knowledge

evaluations and instead persist in basing their judgments on personal intuitions (Sieck &

Arkes, 2005).

It is the combination of overconfidence in knowledge judgments and reluctance to

utilize external cues that might debiase confidence judgments that is the focus of the

current research. My focus is on seeking methods to get individuals to give more weight

to external cues, when they have the potential to debias confident judgments. Such a cue

can be termed a *veridical cue* if it is based on objective information about performance.

The veridical cue should be a good predictor of performance, because it is based on real

data. However, I argue that people often fail to rely on veridical knowledge cues because

there is an ego threat inherent in admitting one's own ignorance. As a result, there may be

value in pairing veridical decision aids with efforts to "buffer" individuals against ego

threats, so that they will utilize the veridical decision aids made available to them.

Evidence for this thesis can be found in the cognitive dissonance literature, to which I

now turn.

*Cognitive Dissonance Perspectives on Overconfidence and Debiasing*

The Cognitive Dissonance literature has long focused on people's motive to

maintain consonance between their cognitions and behavior (Festinger, 1957). Over time,

this motive has come to be described as a motive to maintain a desirable self-image, in

light of one's own actions (Aronson, 1969; Steele, 1988). The classic example is of

forced compliance, as in Festinger and Carlsmith 1959, wherein participants were asked

to lie about a boring task they had completed, and tell another student it was fun. The

self-image of the participant as an honest person was dissonant with the statement made

about the task, and some participants resolved this by changing their attitude toward the task to be consonant with their statement.

Work on post-decisional dissonance presents a mechanism by which Dissonance influences judgment confidence. The post-decisional dissonance model predicts that, when a decision is made, dissonance should occur to the extent that there are lingering doubts about the choice made (See Brehm, 1956). Knox and Inkster (1968) demonstrated this principle with gamblers betting on horse races. They showed that confidence in the success of a bet was higher after the bet had been placed than before. Confidence was presumed to have been driven up to reduce gamblers' post decisional dissonance. In this instance, gamblers may have wanted to avoid the thought that they had wasted their money, and this motivation may have resulted in elevated confidence ratings.

Knox and Inkster (1968) focused on confidence for a past event, but people's confidence about future task performance may also be subject to the effects of motivated change. Consider a student who has an upcoming exam. Although the student's confidence in exam performance should be related to task relevant knowledge, which should in turn predict actual performance, there may be motivated influences as well. In this example, a strong personal investment in academic achievement may bias confidence, independent of exam preparation. This can result in a student feeling confident about an upcoming exam, but performing below their own expectations.

A cognitive dissonance account of precisely this type of knowledge overconfidence was tested in Blanton, Pelham, DeHart, and Carvallo (2001). They found a relationship between the importance students place on an exam question and the confidence they had in their performance. This relationship held, even after controlling

for actual knowledge.   The tendency to feel confident, independent of accuracy, suggests that participants were basing their judgments on something other than a cold analysis of their knowledge.  Moreover, Blanton et al. (2001) found that confidence was least tied to accuracy when ego threat was high.  This was demonstrated through an experimental manipulation.  Prior to the confidence ratings, participants either were or were not reminded of an option to drop the lowest exam.  Blanton et al. found that confidence was highest (independent of accuracy) when there was no such reminder.  This pattern suggests that one reason knowledge confidence operated independent of actual knowledge was that it helped diminish the threat of doing poorly, when there was greater ego involvement in a strong exam performance.

This finding suggests a possibility of debiasing confidence ratings. Insofar as dissonance threats represent threats to self-image, one might reduce or remove these threats as a possible avenue for debiasing. One way to tackle threats to self-image is to use self-affirmations, consistent with Steele's (1988) model of cognitive dissonance. A self-affirmation can take the form of a task in which people are encouraged to affirm important self-values, which act as buffers against possible threats to self image encountered after the affirmation.

Evidence from Blanton et al. (2001) suggests that self-affirmation may debias overconfidence. Their study demonstrated that confidence in a judgment task increased in tandem with increased importance placed on the task, independent of actual performance or ability.  In a design echoing earlier work on post-decisional dissonance, they found that participants in a taste test were more confident about their ability to discriminate between Coke and Pepsi if they initially indicated a strong preference for one flavor over

the other. This effect occurred independent of actual taste-discrimination ability, suggesting that there was ego involvement. Specifically, individuals who felt they preferred one cola over the other were motivated to believe that they could discriminate between the colas. They were thus more confident, independent of their accuracy. More important for debiasing, however, when a self affirmation task was given before confidence was measured, this relationship was diminished. Among affirmed participants, the stated preference in a cola did not exert influence on confidence after controlling for accuracy.  This study appears to show that the mere act of affirming an individual will remove the influence of motivation on confidence judgments.

*A Study Limitation*

Blanton et al.'s (2001) interpretation of their finding was that an affirmation can be sufficient to debias confidence judgments. However, there may have been a key artifact in their study that was helping them attain their effect. In their original design, an external cue was presented to all participants. All participants were given information regarding the average performance of students on the taste-discrimination task. This cue was intended to reduce noise in participants' confidence estimates but it might have played a more vital role. An alternative interpretation of Blanton et al.'s (2001) results is that self-affirmation influenced the acceptance of this external cue for participants in this condition. In this interpretation, an affirmation will not influence confidence by itself. If this interpretation is accurate, when cues needed to calibrate levels of confidence are not present, self-affirmation should have little to no effect on overconfidence.

*Cognitive Attempts to Debias Overconfidence*

Research to date has not tested if affirmations help individuals make better use of a debiasing cue. To the contrary, the approaches taken towards overconfidence debiasing have often focused on cognitive mechanisms, with little concern for motivation (see Lichtenstein, Fischhoff, and Phillips, 1981). As a result, the potential ego threat inherent in accepting debiasing information presented by an experimenter has not been considered. Instead of examining the motivations people may have for maintaining overconfidence, research has centered on cognitive biases such as people's difficulty in understanding probability (Lichtenstein & Fischhoff, 1980) and biased information processing (Koriat, Lichtenstein, & Fischhoff, 1980). Debiasing often takes the form of training individuals to specifically correct these cognitive biases, or, as recently shown in studies by Zimmerman and Kelley (2010), simply allowing participants to repeat tasks so that they can learn from their previous performances.

The debiasing method I am proposing is straightforward compared to many cognitive attempts, as it rests on simply providing participants with veridical information that can guide confidence estimates. I argue that such an approach may be threatening, however, and there is some evidence to support this view. A study by Sieck and Arkes (2005) showed that baserate information can increase the correspondence of the relationship between performance and confidence, but cues in this experiment were effective only when presented to participants using an enhanced method. This method involved an elaborate three-step procedure designed to drive home the risk of overconfidence. First researchers showed participants how their past performance compared to their earlier confidence ratings. Participants were then asked to decide if

their prior performance suggested they were originally overconfident or underconfident. Finally, participants reported how much their original confidence ratings should change for them to match their later performance. I am proposing that this level of involvement is necessary, in part, because participants are motivated not to use the cue provided by these researchers, but I am predicting that my participants will use a veridical cue without such heavy involvement, if an affirmation is presented.

*Summary and Study Overview*

I propose a motivated model of confidence, similar to Blanton et al. (2001), but with the caveat that people do require information from which to base their confidence. Rather than claiming that self-affirmation has an effect on confidence by itself, this model posits that self-affirmation has influence over the degree to which veridical information is accepted and used as a basis for confidence. In a task where people likely already have inflated confidence judgments regarding their future performance, external cues that suggest future performance does not match expected performance might be threatening to self-esteem. Because of this threat to self-esteem, people may benefit from being affirmed before they are able to incorporate these cues into their confidence judgments.

Additionally, this model posits the influence of self-affirmation on cue acceptance, not on confidence directly. Consequently, self-affirmation debiasing should work on a cognitive bias, not just overconfidence where strong motives for confidence inflation are present. The two studies described by Blanton et al. (2001) used tasks where ego-involvement was clearly indicated. Therefore, replication of these tasks would not enable a rigorous test of this model.

The task I chose to use was a standard judgment of learning (JOL) procedure. Much like in the example of the student taking an exam, participants following this procedure have a study period during which they have to judge the likelihood that they will be able to perform well on a following exam. In the learning period, participants are presented with word pairs that they are to remember and asked to estimate the likelihood they will remember the word pair during a subsequent learning task. To increase the likelihood that participants will overestimate their learning for at least some of the words, I also took advantage of a JOL finding reported in Zimmerman and Kelley (2010). They showed that confidence was inflated by word valence in a JOL task. That is, participants over-estimate the likelihood of remembering word pairs if each word has a negative emotional valence (e.g., sick, ugly). Although this effect probably does not reflect a self-esteem motive, I am suggesting that there may be esteem threat in relying on an accuracy cue, when it informs participants that they will have difficulty remembering such words.

In such a task, self-affirmation's direct effect on confidence should be negligible under the current proposed model. Self-affirmation debiasing of such a bias affords the strongest test of this model. Within the framework of a cognitive bias on JOL, debiasing as a result only of a self-affirmation manipulation is not expected. However, presentation of external cues may debias confidence alone. More importantly, the model of self-affirmation as an influence on cue acceptance predicts that the greatest debiasing should be shown when self-affirmation and cues are delivered together, rather than separate.

*Overview*

The study was broken into two phases of development. First, participants were recruited for an initial pilot study, and then, after the pilot study was completed, more

participants were recruited for the main experiment. This final experiment included between-subjects manipulations of external cue presentation and self-affirmation. A pilot study was required to generate veridical cues for participants to use in the main experiment. To this end, participants in the pilot study took part in a simplified version of our main experiment's design.

<div align="center">*Pilot Study*</div>

*Participants*

Participants in the pilot study were 22 college students recruited from the University of Connecticut Psychology Participant Pool for participation in a "memory study."

*Materials*

As in Zimmerman and Kelley (2010), all words used as stimuli in this experiment were taken from Affective Norms for English Words (ANEW; Bradley & Lang, 1999). In total, 84 words taken directly from Zimmerman and Kelley (2010) were used in the pilot study. Another 14 words were taken from ANEW. These pairs were broken in two series of 42 word pairs, where each series consisted of 35 word pairs from Zimmerman and Kelley (2010), and either 7 more words pairs from the same study, or 7 new words pairs as taken from ANEW.

To record participants' JOL, or their confidence for correctly recalling words, an eleven-point scale was used with labels that indicated percentage chance that a word would be remembered. Points on the scale were marked as 0%, 10%, etc up to 100%. On this scale, 0% was described as representing no confidence in a word being correctly recalled, whereas 100% represented total confidence in a word being correctly recalled.

*Procedure*

The JOL task consisted of a simple paired associates learning task. During this first phase of the task, participants were instructed to try to memorize 42 word pairs shown onscreen for 5 seconds. One word in the pair was always marked as the cue word, and the other as the target word. The cue word was explained to be the word used in the recall phase to prompt recall of the target word.

During the testing phase, the second part of the memory task, only one half of the pair, the cue word, was presented. All 42 cue words were presented in the same order of as they had been shown in the study phase. The participant was required to recall and enter the target word by keyboard into a text box on their computer screen. Participants were allowed to proceed through the testing phase at their own pace. During this second phase of the JOL task, all words typed into the computer by participants were stored as text. A recall score was determined by comparing words recalled by participants to the actual target word requested, with a match being marked as correct, and any other answer marked as incorrect. A match was defined as either the target word, including misspelled target words, or a singular or plural form of the target word.

*Results and Discussion*

Participants' recall of words in the pilot study offered data on which to provide meaningful cues in the main experiment. In addition, it provided some insight into how word actual recall related to estimated recall. Mean recall for positive words was 42.86%, for neutral words 28.57% and for negative words 27.47%. Mean JOL for positive words was 51.33%, for neutral words 44.91% and for negative words 48.91% (see Table 1.). These results replicate Zimmerman et al. (2010) in that mean JOL for negative and

positive words were relatively high, but negative words were recalled correctly much less often than positive words. This indicates that, although predicted recall was higher than actual recall for all three types of words on average, the largest discrepancy was for negative words.

The advantage of including a mix of neutral and positive words in the Main Study that follows is that there is variability in the degree of correspondence between actual and estimated recall.  By including all three types, participants thus encountered a mix of word pairs where they might do relatively well by trusting their private estimation (e.g., neutral pairs) but other word pairs (e.g., negative pairs) where they should not go with their gut but instead trust the cue.  To maximize this form of variability, seven negative word pairs deemed too easy, by virtue of being recalled correctly more often than other similar words in pilot testing, were not used in the main study. All other words were used as stimuli in the main experiments. Percentage recall from the pilot study was rounded up to the nearest ten, and these were used as the cues presented to participants in the main experiment (see Table 2.).

## Main Study

### *Participants*

Participants were recruited from the University of Connecticut Psychology Participant Pool from for participation in a "memory study." In total, 146 college students were recruited.

### *Materials*

Two memory tasks were used in this experiment. The main memory task used words that were taken from the pilot study that were found to have both high and low

discrepancy between JOL estimation and recall. In addition, a practice task was added to give participants familiarity with the task to reduce noise in responding. The practice task used 84 words taken from ANEW (Bradley & Lang, 1999). These words were chosen because their affect rating was close to 4.5 on a 9 point scale, making them relatively neutral. A separate list of 84 words were used for the second task.

*Procedure*

The main experiment was split into three phases, first a practice memory task, then a debiasing procedure, and then a second experiment memory task. All memory tasks were identical to the one described above in the pilot study, except for minor changes as noted. These words acted as a practice block for participants to familiarize themselves with the memory task procedure and the reporting of JOLs.[1]

After the practice memory task was completed, participants were given their assigned debiasing procedure, based on a 2 (Affirmation or No-Affirmation) X 2 (External Cue or No-Cue) factorial design. Participants were assigned to complete either a self-affirmation task or a US state and city naming task. The self-affirmation task was presented as a thought-listing questionnaire. Using the same task as in Blanton et al. (2001), participants were first asked to write a description of a value, talent, relationship, or identity that was both important to them and made them feel proud. To provide a double dose of the affirmation, in order to maximize the effects of the affirmation lasted the entire course of the experiment, participants were asked to write a second time about something they were good at and did better than average. The two writing tasks,

---

[1] The advantage of adding a practice task is that it gave participants greater familiarity with the task so that they would not be distracted learning a new task after the affirmation procedure.

completed back-to-back, allowed participants to affirm whatever aspect of themselves they chose to describe.

A U.S. state and world city naming task was used as a control condition to the affirmation task. Completing this task put participants through a similar, but non affirming task to self-affirmed participants. This task, like the self-affirmation task, consisted of two parts; a state naming task and a city naming task. The U.S. state naming task consisted of participants naming the first 30 U.S. states they could think of. In case participants were unable to name 30 states, and to avoid any threat to self that may occur in such a circumstance, little emphasis was placed on participants actually meeting this goal within the time allowed for the task. To prevent this task providing a self-esteem boost by a participant naming all the US states, participants were limited to naming 30 states and were not able to list all 50 states. After naming states, the participants were asked to name 30 cities located anywhere in the world. Similar instructions to the state naming task were provided, whereby it was made clear that participants did not have to name 30 cities if they were unable, but should try to enter as many as they could.

The final part of this study was a second memory task. This task was the same as the first, but using the words adapted from the pilot study. Presentation was fixed so that words did not repeat by valence, to avoid clumping of words of the same valence that might occur if the order of presentation was completely randomized.

As participants made their JOL response in the final memory task, half of all participants were provided with feedback on average participant performance for that specific word pair (taken from the Pilot study). Participants were presented with feedback

on the scale used to provide JOLs, but were free to ignore feedback and provide a different JOL than the one suggested.

<div align="center">Results</div>

*Analysis of the External Cue*

For external cues to be of use in debiasing overconfidence, those cues must be related to actual recall performance. To explore this issue, I examined the estimated recall and actual recall of participants in only the control condition, as they reveal how individuals performed before any experimental manipulations were applied. For participants in this condition, a single residual score was computed, by subtracting predicted recall from actual recall. These results of this analysis are shown in Table 3. This analysis revealed that, on average, actual word recall was 3.24% lower than the external cue ($SD$ = 12.16). This suggests that overall, the feedback was reasonably accurate. There was variability in the utility of the cues, however. Specifically, two words had performance differences larger than two standard deviations from zero (Beach-Circus; Betray Terrible), indicating that cues were less well linked to actual performance. Because one might question if the cues for these words provided a veridical basis for judging later performance, these words were both removed from further analyses.

Overall residual provides information on only one aspect of the relationship of cue to recall. There also is value in examining if the estimated recall from a given cue covaries with actual recall. After removing the two outliers, a simple correlation of the data in Table 3 revealed that external cue and true recall were correlated $r(38)$ = .72, p < 0.001, indicating a strong positive relationship between the external cue given and recall.

This speaks to the potential validity of the cues, suggesting that participants should estimate a greater likelihood of correct recall for word pairs when the cue suggests recall is likely and a lower likelihood of correct recall for word pairs when the cue suggests recall is unlikely.

Finally, regression analyses were performed to examine how the metric in the external cue mapped on to the likelihood of correctly recalling a word.  Specifically, recall was regressed on cue. If external cues are to be of use to participants for debiasing, it is not sufficient for the correlation coefficient to be strong (as indicated above), but the metric of the two scores must line up.  Consider the regression line linking the external cue to recall.  This linear relationship can be represented by

$$Y = B(X) + a$$

where Y is recall as predicted by variable X, the external cue, with a slope of B and an intercept of a. If the cue provides completely veridical prediction of recall, then this equation should have B = 1 and a = 0. This ideal relationship is represented by the dotted identity line in Figure 1.  In contrast, the data produced a linear relationship with a slope of B = 0.73, and an intercept of 6.82, the solid line in Figure 1.

The results of the above regression show that the external cue was related to true recall at close to a one-to-one relationship. As a cue, this relationship is desirable, as the cue was designed to be a good anchor for participants to use when determining the confidence they should have in task performance. However, for the cue to improve prediction, it also needs to be more strongly linked to recall than a participant's intuitive confidence. To examine this issue, I next explored the relationship between recall and JOL for each word pair.

The correlation of recall and average word confidence was quite strong, $r(38) =$ .52, p = 0.001. This suggests that participants did have insight into which word pairs were relatively easy and which were relatively difficult. However, as before, a simple correlation coefficient conveys only some information. To examine the metric relations between recall and prediction, a regression of recall on to average word confidence was performed. This revealed a slope of B = 2.48, and an intercept of -64.63. This slope is further from the ideal slope of 1 than the slope for the external cue. Moreover, the intercept is further away from the ideal of 0 than the intercept produced generated by the external cue. This comparison indicates that the cue should be a better predictor of recall than intuitive confidence, and so participants would increase correspondence between confidence and recall to the extent that they are given the external cue and make use of it.

Particularly visible in figure 2 is a marked restriction in range for confidence ratings. Confidence ratings had a standard deviation of only 3.09%, with a range of between 33.08% and 43.59%, across all words. This compared to a standard deviation of 14.84%, with a range of between 30.26% and 58.97%, for recall. This suggests that participants in the control condition were unable to track the range of difficulty of words when reporting confidence. Words that were difficult to remember on average, and had correspondingly low scores on recall, tended to be rated at levels of confidence similar to most other, easier words.

Further, the regression of confidence on recall, with a large negative intercept, and steep slope, showed that average word confidence was often higher than suggested by recall. Twenty-six of the total 40 words fall below the ideal relationship posited in the above paragraphs. This suggests that, for a majority of words, participants made

judgments that would be labeled as "overconfident" by some conventions (e.g.,

Fischhoff, Slovic, & Lichtenstein, 1977).[2]

*Utilization of Cue*

These analyses provide strong evidence that participants' recall to JOL calibration

may be improved if they do make use of the external cue. This was not possible in the

control condition, because the cue was not provided. My theory predicts, however, that

provision of a cue will only improve the relationship between prediction and performance

when an affirmation is provided. However, my theory also predicts that provision of an

affirmation alone should in no way influence the relationship between confidence and

recall.

To test these predictions, I ran a multilevel regression analysis to determine if the

relationship between confidence for a given word pair and recall for a given word pair

changed, as a function of experimental condition. In this analysis, word pairs were

nested within participants, wherein confidence was regressed on to recall. The logic to

this analysis is that stated confidence for a given word pair should predict later recall for

that same word pair, to the extent that individuals are able to correctly detect which

words they can predict and which they cannot. However, if my theory holds, there

should be a three-way interaction between recall, affirmation, and the provision of

feedback, such that later recall is more strongly predictive of confidence estimates when

both an external cue and an affirmation are provided.

The findings from the multi-level model proposed above are found in Table 4.

The intercept for this model, 36.51%, represents average word confidence for words that

---

[2] As in the pilot, JOL over-estimation was greatest in the negative-negative word pairs but this was not a focus of the analyses that follow.

were not recalled successfully by participants in the control condition. If participants had

perfect knowledge of word difficulty, their confidence for these words might be expected

to be 0%. In contrast, an intercept of 36.51 indicates that participants in this experiment

tended to be overconfident for words that they were not able to successfully recall.

The effect of recall on confidence was significant, p < 0.001 , and can be

interpreted as participants tending to be 6.17% more confident for successfully recalled

words. This means, as expected, that participants were able to adjust their confidence

inline with word difficulty, even without being presented with a word cue or self-

affirmation manipulation. However, a modest increase of 6.17% confidence for correctly

recalled words could be argued to represent only limited insight on the part of the

participant into their ability to successfully memorize a word.

The effect of the between subjects variables on confidence for unsuccessfully

recalled words was also examined. None of these predictors were significant in this

model. This indicates a lack of support for the thesis, as well as a failure to replicate the

findings in Blanton et al. (2001).  This suggests that participants did not alter their

confidence due to the presence of an affirmation, that they were not influenced by the

presence or absence of a veridical cue and they did not overcome resistance to utilizing

the cue when presented with an affirmation.  Although this indicates a failure to support

the primary or alternative hypotheses, some nonsignificant trends are suggestive.

Recall one debiasing effect on confidence would be shown by confident

participants showing a willingness to use feedback (regardless of whether or not there

was an affirmatison).  There was a nonsignificant trend ($p = .14$) for participants to be

less confident when veridical cues were presented. This effect reduces average

confidence by 6.08%; reducing the intecept term to 30.43% for incorrectly recalled words. This reduction moves confidence toward the ideal of 0%, and thus improves the calibration between recall and confidence.

In addition, there was a nonsignificant trend for recall and the presence of a veridical cue to interact ($p = .13$). This effect means that average confidence rose by an additional 3.06% for correctly recalled words, when the cue was present. This increase in confidence for correctly recalled words is smaller than the increase seen when the cue is not present, but an increase in confidence is still in the desired direction for increasing calibration of confidence and recall.[3]

Other than examining whether the combination of the cue and affirmation condition resulted in debiasing, it is also possible to examine whether affirmation enhanced the correspondence of cue to confidence. That is, an analysis can investigate whether participants were able to take the cue given and use that as an anchor for their confidence, and whether affirmation influenced the extent to which this occurred. For this analysis, it makes sense to only examine participants who were presented with a cue, as only these participants had the opportunity to use the cue. Using a similar model to the above multi-level model, confidence was regressed on recall at level one, with the percent value of the cue given and affirmation as level two predictors (see Table 5.).

The intercept for this model was 26.06% confidence, indicating that participants were confident of getting words correct 26.06% of the time on average, controlling for the effect of recall, cue and affirmation. Interestingly, in this model the relationship between recall and confidence was not significant. However, the main effect of cue was

---

[3] Removing affirmation from the model , word cue had a significant main effect on confidence, $p < .05$, and word cue significantly interacted with the presence of a word cue, $p < .05$.

significant, and an interaction between cue and recall approached significance. The nature of this effect was that participants gave more weight to a cue for words they later correctly remembered ($B = .11$) than for words they later failed to remember ($B = .15$). No other interactions in this model were significant.

*Exploratory Analyses*

Another way to investigate the influence of threat on cue acceptance is to examine whether particularly threatening cues were of less use in debiasing overconfidence. One way to approach this issue is to use practice trials as a measure for average confidence and recall for each participant. These scores of confidence and recall are independent of debiasing procedures used during experiment trials. Using practice recall and confidence, it is possible to determine the relation of participants' recall to their confidence. It is then possible to investigate how this relationship tended to influence confidence in the experiment. If participants tend to predict a better performance than they delivered at practice, then they may have more difficulty accepting cues instructing them to lower confidence. That is, participants with low practice recall but high practice confidence might be assumed to be participants with motivated high confidence. Such participants were originally predicted to be the most resistant to cue based debiasing. If cue threat had any impact on results, debiasing due to cue would be attenuated to the extent that a participant was over confident at practice.

Using a between subjects regression approach, aggregate scores for confidence and recall were created for each participant for both their practice and experiment trials. Average confidence for experiment trials was regressed on experiment trial recall, practice confidence, practice recall, and a binary variable for cue condition. If practice

overconfidence had a an effect on cue debiasing then a three-way interaction between practice recall, practice confidence and cue condition should be significant in this model.

The three-way interaction between practice recall, practice confidence and cue condition interaction was found to be significant, $p = 0.02$, controlling for the effect of experiment recall. The pattern was not as predicted but suggestive of a potentially interesting psychological dynamic (see Table 6.).

The top panel of Figure 3 shows estimated experiment confidence of participants with low practice recall (8.53% accuracy). This reveals little more than a main effect of confidence, with those higher in confidence being more accurate than those low in confidence. In contrast, the bottom panel of Figure 3 shows that among participants with high practice recall (48.64%) there was an interaction pattern, such that the cue lowered the confidence of those originally high in confidence and raised the confidence of those originally low.  This pattern suggests that the feedback was being utilized more among those who were originally high in recall in the practice trials.

## Discussion

The central hypothesis of this study was that self-affirmation would enhance the confidence debiasing effect of a veridical cue. However, none of the analyses conducted offered support for this hypothesis. A multi-level model was used to investigate the effect of self-affirmation together with cue condition on the participant level variables. In this model, confidence was regressed on recall, the participant level variable, and also on self-affirmation and cue condition, which were group level variables. An interaction between self-affirmation and cue condition was predicted by the hypothesis that self-affirmation

enhanced the effect of cue. However, the predicted interaction between self-affirmation and cue condition was not found.

Part of the logic of this study was that debiasing via self-affirmation alone, as postulated in Blanton et al. (2001), should not occur. As a consequence, it was predicted that self-affirmation would not have a main effect on confidence. Consistent with this prediction, no analyses found evidence for the influence of self-affirmation on debiasing.

Another hypothesis, based on previous research on feedback and cues, was that presentation of veridical cues should have a debiasing effect on confidence, even without being combined with self-affirmation. In the multi-level model referred to above, evidence for the effect of cue would be shown by a significant interaction of cue condition with recall, or a main effect of cue condition. Neither of these effects were significant in the multi-level model tested. However, the effect of cue on confidence showed some indication of trending toward significance.[4]

All analyses of confidence compared to recall showed that calibration of recall to confidence was generally poor and biased in the direction of participants giving a greater likelihood of recall than their percent likelihood of recall. In contrast, preliminary analyses of cue showed that the cues used in the study were highly correlated with actual recall, and that cues for a given word tended to be better indicators of recall for that word than participants' own JOLs. Consequently, a participant using cues alone for their JOL would likely have achieved better JOL to recall calibration than would be observed among participants relying more on naïve confidence estimates. However, participants

---

[4] As noted in the results, removing affirmation from the multi-level model led cue condition to become a significant predictor of confidence for both incorrectly and correctly recalled worlds. This effect was clearly fragile, however, and can only be taken as a suggestion that feedback had some influence on participant confidence.

that were given cues did not solely rely on the cue to guide their confidence: when a cue was presented to participants, cue and confidence were not perfectly correlated, and in fact, although there was a strong association between cue and confidence, there was evidence that cues tended not to be fully accepted by participants. This suggests that cues could be used more by participants, but this study did not demonstrate a mechanism that allowed that to happen.

The aim of this study was to demonstrate that self-affirmation could act as buffer against threatening information, or the veridical cues, and thus would offer a method by which cues could be accepted and used by individuals rather than rejected and ignored. Unfortunately, no evidence was found for self-affirmation aiding debiasing in any way. Self-affirmation was not a significant predictor in either of two multi-level models analyzed, indicating that self-affirmation had no effect on debiasing. In other words, self-affirmation did not directly influence the relationship between recall and confidence. In this respect, the study did not replicate the findings of Blanton et al. (2001). The failure to replicate the direct debiasing of self-affirmation does not strengthen the competing hypothesis, that self-affirmation interacts with cue presentation, as no evidence was found to back either the old or the new hypotheses.

*Reasons for Self-Affirmation Null Findings*

One possible reason that self-affirmation did not have an effect on confidence is that the self-affirmation manipulation itself failed, such that participants in the study were not sufficiently affirmed by the manipulation used. Unfortunately, due to fears about manipulation checks influencing self-affirmation tasks (Steele & Liu, 1983), a manipulation check could not be incorporated in the experiment design. As a result, it is

not possible to do much more than speculate about whether the self-affirmation task worked or not. However, there was initial concern when designing this experiment that self-affirmation may "run out" over the course of a very long memory task. In an effort to boost affirmation and potentially lengthen its effect, two self-affirmation manipulations were used instead of just one. Although this precaution was taken in designing the experiment, the self-affirmation may still have weakened over time and not lasted over the course of the entire memory task. However, some exploratory analyses not reported above argued against this.  Separate regression models were run examining either only the first third of words seen by a participant, or the first three words. These analyses also failed to reveal an effect of self-affirmation. This suggests that either the self-affirmation did not enhance self-regard, or it did it had no influence on confidence debiasing.

Other exploratory analyses suggested that cues were rejected by participants, but not due to threat. Instead, analyses suggested cue debiasing tended to be highest among participants with high levels of recall during the practice task, raising the confidence of those initially very low and lowering the confidence of those initially high. This finding suggests that cues were not rejected by all participants. Unfortunately, the specific pattern argued against there being ego threat in the task developed, for which the affirmation given might help debias. For instance, it is arguable that highly confident participants should be the participants threatened most by cue debiasing, because those participants would have to give up their high confidence in order to accept debiasing cues. In fact, analyses showed that highly confident participants that were also highly accurate during practice trials tended to have little problem accepting cues, but any participant with low

practice trial recall, including participants that were particularly overconfident, tended to ignore cues, perhaps because they were not engaged in the task.

Importantly, these analyses showed no indication that cues were threatening to participants in the way theorized above. Any suggestion that participants did not find cues to be threatening is problematic for the hypotheses of this study because, if cues were not threatening to participants then there would be no way for self-affirmation to influence debiasing.

*Further Research*

Future research could attempt to discover whether introducing a stronger self-affirmation manipulation, or more threatening cues, changes the results of the current study. Such changes might determine, for instance, whether self-affirmation can interact with cue to enable greater debiasing, or whether the theory put forward by Blanton et al. (2001) , whereby debiasing functions as a main effect of self-affirmation, is correct. A starting point would be to use a task with greater ego-involvement on the part of participants, to ensure that cues are threatening to participants. Using an IQ test as the task would ensure that most participants are engaged in performing well on the task, and that cues would be particularly threatening.

If a debiasing technique can be established using self-affirmation, another potential step would be to move toward a practical application. Arguably, debiasing confidence for word-pair memorization task has little real-world utility, whereas debiasing for some test performances (e.g., IQ tests, or academic standardized test) would suggest some practical use for this kind technique. An older study by Fischoff et al. (1977) demonstrated that overconfidence when gambling can lead to real monetary

losses. That kind of real world impact could also be added to such research methodology used in this study. Recalling the example of the overconfident student, who skips studying when she really needed to buckle down, perhaps some kind of procedure could be developed to examine such decisions to continue or stop learning. Ideally, development of a self-affirmation debiasing technique would help to debias overconfident students, and that debiased students would then reap the benefits of extra study, leading to enhanced learning of important material and better exam performance.

References

Aronson, E. (1969). The theory of cognitive dissonance: A current perspective. In L. Berkowitz (Ed.). *Advances in Experimental Social Psychology, Vol. 4* (pp. 1–34). New York, NY: Academic Press.

Blanton, H., Pelham, B. W., DeHart, T., & Carvallo, M. (2001). Overconfidence as dissonance reduction. *Journal of Experimental Social Psychology, 37*(5), 373-385.

Bradley, M. M., & Lang, P. J. (1999). Affective norms for English words (ANEW): Instruction manual and affective ratings. Technical report C-1, The Center for Research in Psychophysiology, University of Florida.

Brehm, J. W. (1956). Postdecision changes in the desirability of alternatives. *The Journal of Abnormal and Social Psychology, 52*(3), 384-389.

Festinger, L. (1957). A theory of cognitive dissonance. Stanford University Press.

Festinger, L., & Carlsmith, J. M. (1959). Cognitive consequences of forced compliance. *The Journal of Abnormal and Social Psychology,58*(2), 203-210.

Fischhoff, B., Slovic, P., & Lichtenstein, S. (1977). Knowing with certainty: The appropriateness of extreme confidence. *Journal of Experimental Psychology: Human Perception and Performance, 3*(4), 552-564.

Knox, R. E., & Inkster, J. A. (1968). Postdecision dissonance at post time. *Journal of Personality and Social Psychology, 8*(4, Pt.1), 319-323.

Koriat, A., Lichtenstein, S., & Fischhoff, B. (1980). Reasons for confidence. *Journal of Experimental Psychology: Human Learning and Memory, 6*(2), 107-118.

Lichtenstein, S., & Fischhoff, B. (1980). Training for calibration. *Organizational Behavior & Human Performance, 26*(2), 149-171.

Lichtenstein, S., Fischhoff, B., & Phillips, L. D. (1982). Calibration of Probabilities: The State of the Art to 1980. In D. Kahneman, A. Tversky (Eds.), *Judgement under Uncertainty: Heuristics and Biases* (pp. 306-334). Cambridge, UK: Cambridge University Press.

Pulford, B. D., & Colman, A. M. (1997). Overconfidence: Feedback and item difficulty effects. *Personality and Individual Differences,23*(1), 125-133.

Sieck, W. R., & Arkes, H. R. (2005). The Recalcitrance of Overconfidence and its Contribution to Decision Aid Neglect. *Journal of Behavioral Decision Making, 18*(1), 29-53.

Steele, C. M. (1988). The psychology of self-affirmation: Sustaining the integrity of the self. In L. Berkowitz, L. Berkowitz (Eds.) ,*Advances in Experimental Social Psychology, Vol. 21: Social psychological studies of the self: Perspectives and programs* (pp. 261-302). San Diego, CA: Academic Press.

Steele, C. M., & Liu, T. J. (1983). Dissonance processes as self-affirmation. *Journal of*

*Personality and Social Psychology, 45*(1), 5-19.

Zimmerman, C. A., & Kelley, C. M. (2010). "I'll remember this!" Effects of emotionality

on memory predictions versus memory performance. *Journal of Memory and*

*Language, 62*(3), 240-253.

Table 1

*Average pilot study recall and JOL for valenced word pairs.*

|  | Positive Words | Neutral Words | Negative Words |
|---|---|---|---|
| Mean Recall (%) | 42.86 | 28.57 | 27.47 |
| SD Recall (%) | 8.38 | 12.32 | 18.42 |
| Mean JOL (%) | 51.33 | 44.91 | 48.91 |
| SD JOL (%) | 7.54 | 5.39 | 8.38 |

*Note.* JOL = Judgment of Learning

Table 2

*Pilot study recall and JOL for each word pair.*

| Cue Word, Target Word | JOL (%) | Recall (%) |
|---|---|---|
| ADVICE, STATUE | 54.35 | 19.23 |
| BANNER, GENDER | 47.19 | 23.08 |
| CANNON, RAIN | 50.66 | 53.85 |
| CONTEXT, BASKET | 46.53 | 19.23 |
| DETAIL, ERRAND | 47.74 | 19.23 |
| ELEVATOR, HIGHWAY | 52.22 | 57.69 |
| KETTLE, FABRIC | 44.65 | 23.08 |
| KEY, WHISTLE | 45.83 | 26.92 |
| KNOT, PASSAGE | 38.78 | 23.08 |
| MANNER, SHADOW | 36.38 | 23.08 |
| SALUTE, CORK | 37.78 | 23.08 |
| TOOLS, HABIT | 42.50 | 30.77 |
| TRUNK, SPHERE | 41.15 | 23.08 |
| VIOLIN, AVENUE | 43.02 | 34.62 |
| BATH, CHAMPION | 40.87 | 42.31 |
| BEACH, CIRCUS | 55.67 | 61.54 |
| CAKE, PRIEST | 48.04 | 38.46 |
| CROWN, MOVIE | 63.37 | 42.31 |
| EXERCISE, DIAMOND | 46.88 | 38.46 |
| GARDEN, TALENT | 48.96 | 42.31 |
| HUG, TREASURE | 59.90 | 38.46 |
| JOKE, LEADER | 60.56 | 34.62 |
| LUST, DIPLOMA | 36.60 | 38.46 |
| PERFUME, BUNNY | 45.38 | 50 |
| RADIO, RESCUE | 53.30 | 26.92 |

| | | |
|---|---|---|
| REUNION, CASH | 51.78 | 50 |
| TOY, FAME | 55.07 | 46.15 |
| ZEST, WEDDING | 52.31 | 50 |
| ABUSE, RABIES | 46.11 | 15.38 |
| AFRAID, MUTILATE | 43.33 | 15.38 |
| BEES, BOMB* | 57.78 | 69.23 |
| BETRAY, TERRIBLE | 60.56 | 30.77 |
| DEBT, BULLETS* | 32.22 | 23.08 |
| DISLOYAL, SUFFOCATE | 52.22 | 23.08 |
| DUMP, CANCER* | 47.78 | 61.54 |
| HATRED, TORTURE | 57.78 | 23.08 |
| KNIFE, FEVER* | 42.22 | 30.77 |
| PRISON, BURN* | 59.44 | 38.46 |
| SLAUGHTER, TRAGEDY | 66.11 | 15.38 |
| STAIN, DIVORCE* | 52.78 | 76.92 |
| TERRORIST, ULCER* | 48.33 | 38.46 |
| TOBACCO, ASSAULT* | 42.78 | 38.46 |
| ACCIDENT, SLUM* | 49.38 | 30.77 |
| CANCER, REJECTED | 50.63 | 30.77 |
| COFFIN, STRESS* | 46.13 | 38.46 |
| DROWN, SLAVE* | 50.00 | 38.46 |
| INFECTION, RAGE | 33.63 | 0 |
| KILLER, DISASTER* | 55.63 | 23.08 |
| MOLD, TRAITOR* | 37.88 | 23.08 |
| MURDERER, DISTRESSED | 45.63 | 15.38 |
| NIGHTMARE, POLLUTE | 38.75 | 7.69 |
| POVERTY, SNAKE | 38.75 | 7.69 |
| RAPE, BANKRUPT | 51.88 | 7.69 |
| ROBBER, POISON | 56.25 | 23.08 |
| TERRIFIED, POISON* | 57.50 | 15.38 |
| VICTIM, THORN | 48.13 | 7.69 |
| Mean | 48.52 | 31.59 |
| SD | 7.76 | 16.20 |

*Note.* Word pairs marked * not used after the pilot study

Table 3

*Main study recall, cue and residual for each word pair.*

| Cue Word, Target Word | Recall (%) | Cue (%) | Residual |
|---|---|---|---|
| BEACH, CIRCUS | 38.46 | 70 | -31.54 |
| BETRAY, TERRIBLE | 15.38 | 40 | -24.62 |
| GARDEN, TALENT | 28.21 | 50 | -21.79 |

| | | | |
|---|---|---|---|
| MANNER, SHADOW | 10.26 | 30 | -19.74 |
| TOOLS, HABIT | 20.51 | 40 | -19.49 |
| DISLOYAL, SUFFOCATE | 46.15 | 40 | -17.18 |
| AFRAID, MUTILATE | 5.13 | 20 | -14.87 |
| SLAUGHTER, TRAGEDY | 5.13 | 20 | -14.87 |
| TERRIFIED, POISON | 5.13 | 20 | -14.87 |
| BATH, CHAMPION | 35.9 | 50 | -14.1 |
| HATRED, TORTURE | 12.82 | 30 | -12.05 |
| SALUTE, CORK | 20.51 | 30 | -9.49 |
| HUG, TREASURE | 30.77 | 40 | -9.23 |
| CANNON, RAIN | 51.28 | 60 | -8.72 |
| KETTLE, FABRIC | 23.08 | 30 | -6.92 |
| VIOLIN, AVENUE | 33.33 | 40 | -6.67 |
| PERFUME, BUNNY | 43.59 | 50 | -6.41 |
| ADVICE, STATUE | 15.38 | 20 | -4.62 |
| CONTEXT, BASKET | 15.38 | 20 | -4.62 |
| DETAIL, ERRAND | 15.38 | 20 | -4.62 |
| KNOT, PASSAGE | 25.64 | 30 | -4.36 |
| CAKE, PRIEST | 35.9 | 40 | -4.1 |
| EXERCISE, DIAMOND | 38.46 | 40 | -1.54 |
| CROWN, MOVIE | 48.72 | 50 | -1.28 |
| REUNION, CASH | 48.72 | 50 | -1.28 |
| ELEVATOR, HIGHWAY | 58.97 | 60 | -1.03 |
| ABUSE, RABIES | 20.51 | 20 | 0.51 |
| TOY, FAME | 51.28 | 50 | 1.28 |
| ZEST, WEDDING | 51.28 | 50 | 1.28 |
| MURDERER, DISTRESSED | 17.95 | 30 | 3.08 |
| LUST, DIPLOMA | 43.59 | 40 | 3.59 |
| CANCER, REJECTED | 46.15 | 40 | 6.15 |
| BANNER, GENDER | 38.46 | 30 | 8.46 |
| KEY, WHISTLE | 38.46 | 30 | 8.46 |
| TRUNK, SPHERE | 38.46 | 30 | 8.46 |
| JOKE, LEADER | 48.72 | 40 | 8.72 |
| NIGHTMARE, POLLUTE | 23.08 | 20 | 13.08 |
| POVERTY, SNAKE | 23.08 | 10 | 13.08 |
| RAPE, BANKRUPT | 25.64 | 10 | 15.64 |
| VICTIM, THORN | 25.64 | 10 | 15.64 |
| INFECTION, RAGE | 17.94 | 0 | 17.95 |
| RADIO, RESCUE | 48.72 | 30 | 18.72 |
| Average | 30.1 | 33.33 | -3.24 |
| SD | 14.71 | 15.57 | 12.16 |

*Note.* Word pairs sorted from most negative to most positive residual. Top two word pairs removed from analyses due to large residuals.

Table 4

*Multi-level model of confidence regressed on recall, at the participant level, and cue condition and affirmation between subject variables.*

| Predictor | Estimate | Standard Error | Degrees Of Freedom | t-Value | Significance |
|---|---|---|---|---|---|
| Intercept | 36.51 | 2.8 | 142.75 | 13.03 | < 0.001*** |
| Recall (0, 1) | 6.17 | 1.39 | 124.64 | 4.45 | < 0.001*** |
| Cue Condition (0, 1) | -6.08 | 4.05 | 143.17 | -1.5 | 0.14 |
| Affirmation Condition (0, 1) | 0.70 | 4.05 | 143 | 0.17 | 0.86 |
| Recall by Cue | 3.06 | 2.02 | 125.51 | 1.51 | 0.13 |
| Recall by Affirmation | 1.22 | 1.99 | 128.68 | 0.61 | 0.54 |
| Affirmation by Cue | -0.88 | 5.8 | 143.2 | -0.15 | 0.88 |
| Recall by Affirmation by Cue | -0.28 | 2.85 | 127.23 | -0.1 | 0.92 |

*Note.* $*p < .05. **p < .01. ***p < .001.$

Table 5.

*Multi-level model of confidence regressed on recall, at the participant level, and cue value and affirmation between subject variables.*

| Predictor | Estimate | Standard Error | Degrees Of Freedom | t-Value | Significance |
|---|---|---|---|---|---|
| Intercept | 26.06 | 2.73 | 96.72 | 9.54 | < 0.001*** |
| Recall (0, 1) | 3.76 | 2.56 | 422.42 | 1.47 | 0.14 |
| Cue Value | 0.15 | 0.04 | 271.8 | 4.16 | < 0.001*** |
| Affirmation Condition (0, 1) | 0.88 | 3.89 | 96.02 | 0.28 | 0.82 |
| Recall by Cue Value | 0.11 | 0.06 | 274.5 | 1.88 | 0.06 |
| Recall by Affirmation | 4.13 | 3.66 | 461.53 | 1.13 | 0.26 |
| Affirmation by Cue Value | -0.04 | 0.05 | 271.5 | -0.73 | 0.47 |
| Recall by Affirmation by Cue Value | -0.76 | 0.09 | 274.4 | -0.9 | 0.37 |

*Note.* $*p < .05. **p < .01. ***p < .001.$

Table 6.

*Multiple regression aggregate between subject values of experiment confidence on experiment recall, practice confidence, practice recall and cue condition.*

| Predictor | B | Standard Error | Beta | t-Value | Significance |
|---|---|---|---|---|---|
| Intercept | -0.01 | 0.06 | | -0.23 | 0.82 |
| Experiment Recall | 0.23 | 0.06 | 0.27 | 3.96 | < 0.001*** |
| Practice Confidence | 0.67 | 0.15 | 0.66 | 4.58 | < 0.001*** |
| Practice Recall | -0.16 | 0.23 | -0.17 | -0.70 | 0.48 |
| Cue Condition (0, 1) | 0.00 | 0.08 | 0.00 | -0.02 | 0.99 |
| Practice Confidence by Practice Recall | 0.57 | 0.43 | 0.40 | 1.33 | 0.18 |
| Practice Confidence by Cue Condition | 0.02 | 0.18 | 0.03 | 0.12 | 0.90 |
| Practice Recall by Cue Condition | 0.43 | 0.27 | 0.49 | 1.60 | 0.11 |
| Practice Confidence by Practice Recall by Cue Condition | -1.30 | 0.53 | -0.82 | -2.45 | 0.02* |

*Note.* $*p < .05$. $**p < .01$. $***p < .001$.

Figure 1.

Figure 2.

**Low Practice Recall**



**High Practice Recall**
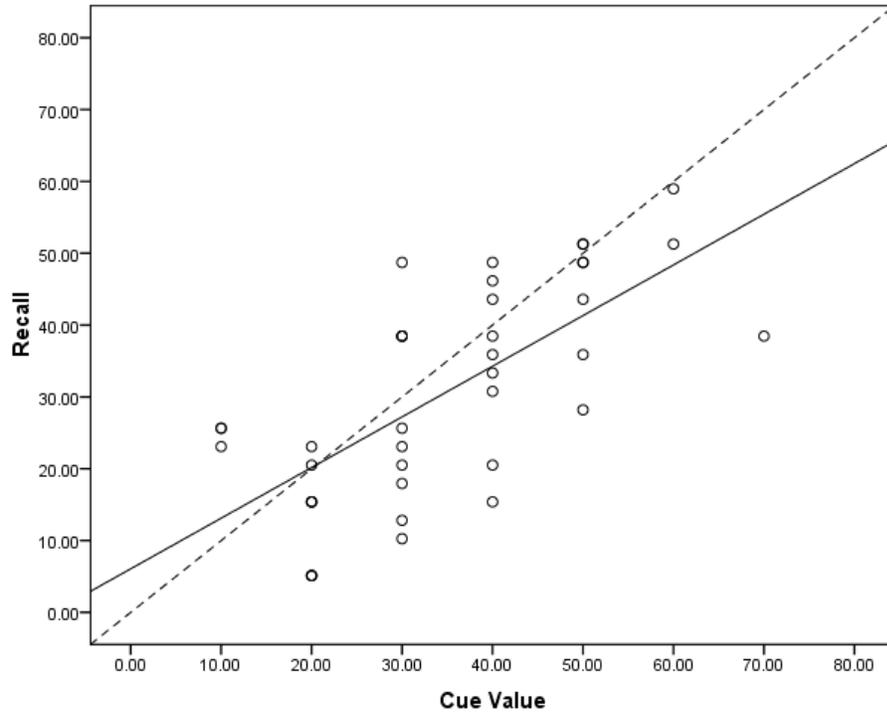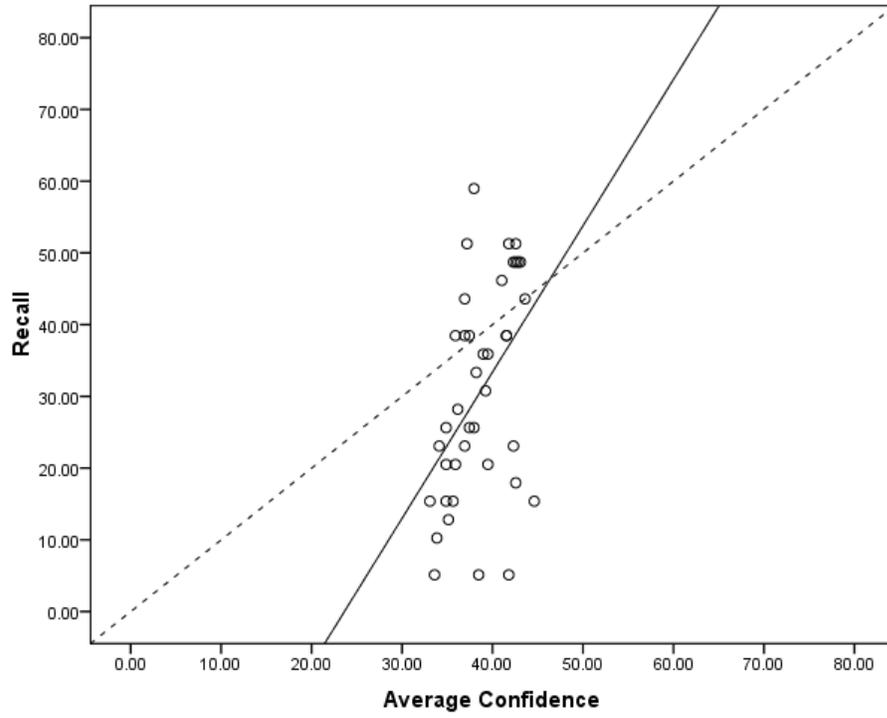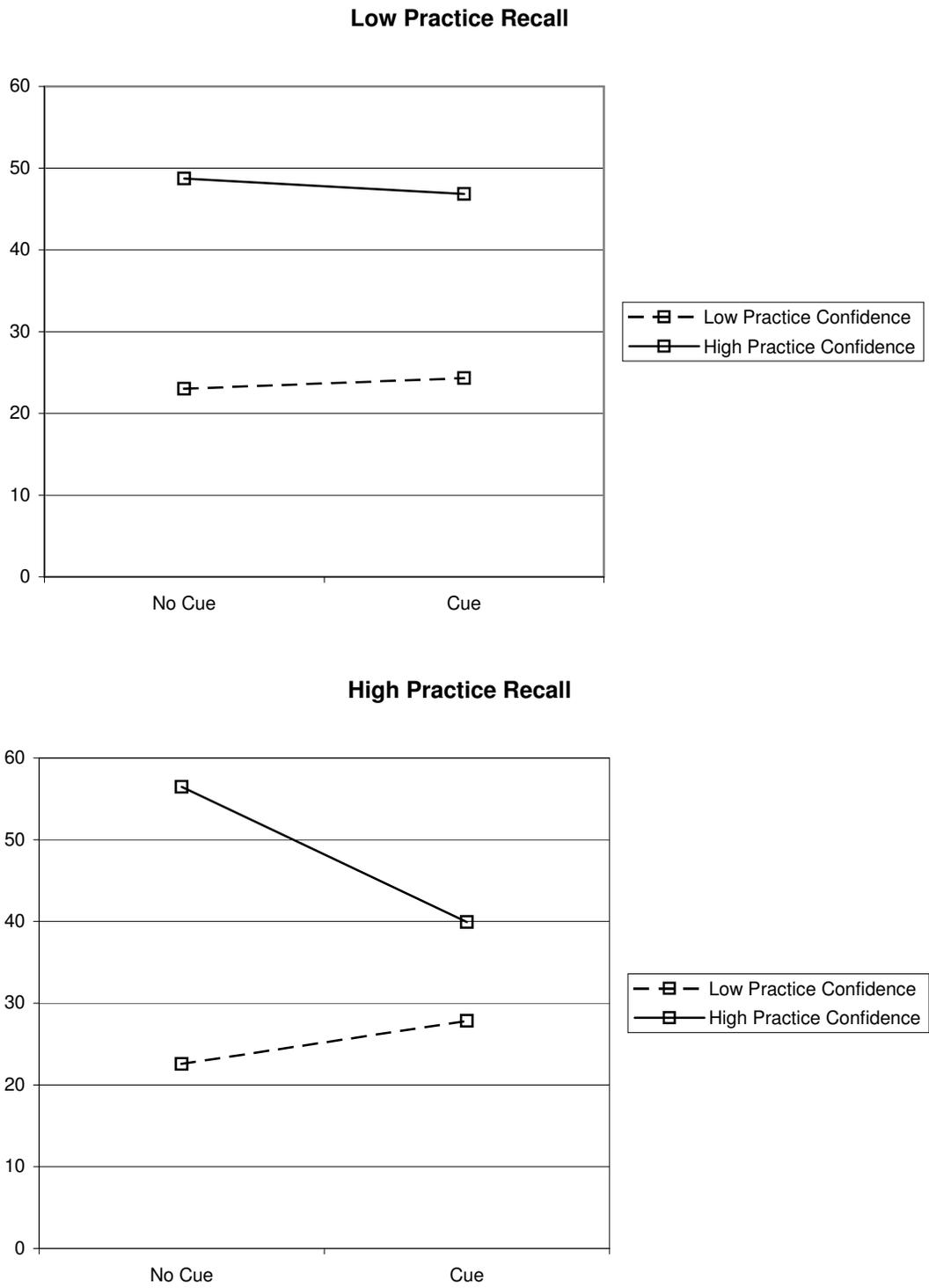


Figure 3.

Figure Captions

Figure 1. Average word recall regressed on word cue value for control participants.

Figure 2. Average word recall regressed on word Judgment of Learning confidence for

control participants.

Figure 3. Effect of the three-way interaction of cue condition, average practice task

confidence, and average practice task recall on average experiment task confidence,

controlling for average experiment task accuracy.