## Widely-Used Measures of Overconfidence Are Confounded With Ability

Stephen A. Spiller

UCLA Anderson School of Management

November 29, 2024

## **Author Note**

Stephen A. Spiller D https://orcid.org/0000-0001-6951-6046

All simulation and analysis code is available at

https://researchbox.org/1597&PEER\_REVIEW\_passcode=ORRDVP. Data and materials from prior

investigations are available via the cited articles. I have no conflicts of interest to disclose.

Correspondence concerning this article should be addressed to Stephen A. Spiller, UCLA

Anderson School of Management, 110 Westwood Plaza, Los Angeles, CA, 90095. Email:

stephen.spiller@anderson.ucla.edu

#### Abstract

The overconfidence concept is one of the great success stories of psychological research, influencing research in other disciplines as well as discourse in the popular press, business, and public policy. Relative to underconfidence, overconfidence at various tasks is purportedly associated with greater narcissism, lower anxiety regarding those tasks, higher status, greater savings, more planning, and numerous other differences. Yet much of this evidence may merely reflect that there are associations with ability rather than overconfidence. This results from two underappreciated properties of typical measures of overconfidence. First, performance is an imperfect measure of ability; accounting for performance does not sufficiently account for ability. Second, self-evaluations of performance should reflect ability in addition to performance; because performance is ambiguous, people should use prior beliefs about their ability. I show these basic principles imply that commonly-used measures of overconfidence are confounded with ability. I support these analytical results by reexamining previously-published findings. In the first analysis, I find overconfidence predicts subsequent performance, consistent with overconfidence as a signal of ability but inconsistent with overconfidence as a bias. In the second set of analyses, I find the purported association between overconfidence and other proposed constructs can be adequately explained through ability alone. I close with recommendations on approaches to recognize and reduce the extent of the problem. This model serves as a stark reminder: when researchers propose that differences in overconfidence are associated with other behaviors, beliefs, or evaluations, they must account for the possibility that differences in ability provide a sufficient explanation.

Keywords: overconfidence, ability, knowledge, performance, measurement error

Overconfidence is widely considered to be a ubiquitous bias. It is reliably reproduced in academic research, worthy of chapters in popular business books, and labeled as "the most significant of the cognitive biases" by a founder of the heuristics-and-biases research program (Kahneman, 2011). Overconfidence refers to the state in which one's beliefs regarding one's ability on some dimension exceed one's true ability on that dimension.<sup>1</sup> Correlates of measures of overconfidence on specific tasks have been widely found. These include narcissism (Ames and Kammrath 2004; Campbell et al. 2004; John and Robins 1994), savings (Avdeenko et al. 2019), advice-seeking (Kramer 2016), financial planning (Anderson et al. 2017; Parker et al. 2012), reduced language anxiety (MacIntyre et al. 1997), social status (Anderson et al. 2012), social class (Belmi et al. 2020), choice of nonlinear incentives (Larkin and Leider 2012), susceptibility to false news (Lyons et al. 2021), search behavior (Moorman et al. 2004), stock ownership (Ke, 2021), short-term debt (Landier and Thesmar 2008), and many more.

But do associations with measures of overconfidence necessarily imply associations with true overconfidence (i.e., beliefs that exceed ability)? No. In this research, I build on a key insight from Moore and Healy's (2008) model of overconfidence: self-evaluations of performance are sensitive to one's prior beliefs.<sup>2</sup> By modeling the consequences of differences in ability, I show that under reasonable assumptions, widely-used measures of overconfidence are confounded with ability. Reanalysis of published findings indicate this confound can be consequential. I do not argue overconfidence, nor differences in overconfidence, nor correlates of overconfidence measures and candidate correlates, even when overconfidence is merely incidental. Thus, observing a correlation between overconfidence measures and that construct.

<sup>&</sup>lt;sup>1</sup> The literature uses a variety of related terms, including overconfidence, biased self-evaluations or self-assessments, unjustified confidence, inappropriate confidence, subjective knowledge when controlling for objective knowledge, and others. There are also direct links to the literature on the correlates of positive self-views and self-enhancement (Taylor and Brown 1988; Colvin et al. 1995) which I return to at the end of the paper.

<sup>&</sup>lt;sup>2</sup> Moore and Healy (2008) present a model in which prior beliefs about a task's simplicity along with ambiguity in self-assessments of performance lead people to overestimate performance on hard tasks, underestimate performance on easy tasks, overplace performance on easy tasks, and underplace performance on hard tasks. Whereas their model explains these patterns across *tasks*, my model considers differences across *people*, holding task constant.

The key issue is that although researchers intend to account for latent ability, they instead account for observed performance. If people are appropriately sensitive to even partially-calibrated prior beliefs (e.g., they incorporate base rates and use Bayesian-like reasoning), their self-evaluations will reflect ability directly, not just through performance. Because self-evaluations reflect ability and accounting for performance is insufficient to account for ability, the resulting associations between measures of overconfidence and other constructs are systematically biased by ability. Although many reports claim to find evidence that *overconfidence* on a task is associated with various correlates, that evidence could instead indicate that *ability* on a task is associated with those correlates. This confound frequently escapes notice. This may be because in contrast to other uses of noisy covariates, when studying overconfidence, (a) it is not always apparent that performance is a noisy measure of ability rather than itself being the focus of inquiry, and (b) it is not always apparent that evaluations ought to regress to one's prior beliefs. This confound can be particularly pernicious because ability is often considered and explicitly ruled out as an alternative explanation of the results based on how the overconfidence measure is constructed.

I begin by describing a typical paradigm used to measure differences in overconfidence, variations on that theme, and why they each result in a problem. I next present a mathematical model to formalize and quantify this bias. I then examine whether these theoretical predictions hold in real data using previously collected datasets. First, I establish the confound generates correlations where we expect none exist. Using data from Moore and Healy (2008), I find measures of overconfidence predict later performance, consistent with an account in which measures of overconfidence are confounded with ability. Second, I establish the confound could plausibly account for reported findings. Using data from Anderson et al. (2012) and Belmi et al. (2020), I reexamine the relationship between overconfidence and proposed causes and consequences. I find it is possible to explain the observed correlations through the confound with ability, even if there were no overconfidence. I close with recommendations to recognize and ameliorate the problem, even if eliminating the problem may be a nearly unattainable target.

#### Theme and Variations: How Differences in Overconfidence Are Measured

Research on correlates of overconfidence has used an extensive array of measures. I consider cases of *overestimation* of absolute performance and *overplacement* of relative performance for which there is a reality criterion against which to compare (Moore and Healy 2008; I do not address *overprecision*, or excessive certainty in one's knowledge). When using a single measure of performance and a single self-evaluation, there are at least 20 different ways that overconfidence may be assessed; adding cases in which self-evaluations reflect future expectations rather than past performance would nearly double this number. Each of these measures of overestimation and overplacement are confounded with ability on the focal task. Such measures vary in terms of whether the measures assess absolute or relative performance, whether self-evaluations assess performance or ability, whether the self-evaluation measure is in the same metric or a different metric as performance, and whether the measures assess overconfidence by including a control variable, calculating a residual, or calculating a difference score.

## **Base Case**

Begin by considering a study designed to assess overconfidence regarding absolute performance using the residual of a self-evaluation measure in the same metric as performance. Participants complete an ability-based task (e.g., a 10-item financial literacy quiz) and then report their self-evaluation of their own performance (e.g., how many of the 10 items do they think that they got correct). The researcher then regresses self-evaluations of performance as the dependent variable on objective performance as the independent variable. The residual of this regression, reflecting how much higher or lower selfevaluations are than is warranted by objective performance, is used as a measure of overconfidence. The researcher then tests whether those residuals are associated with some other measure (e.g., financial planning) in a new analysis.

## **Residual vs. Control vs. Difference**

There are three approaches researchers may use to combine a measure of performance and a selfevaluation to compute overconfidence: they may use residualized self-evaluations, they may control for performance in a multiple regression analysis, or they may calculate a difference score.<sup>3</sup> The first two are nearly identical. All three are problematic, for related but distinguishable reasons detailed in the model.

Using the *residual* approach, researchers regress self-evaluations on objective performance and keep the residuals as measures of overconfidence. By construction, the residuals have a mean of zero: scores ranging from well-calibrated to overconfident will be indistinguishable from those ranging from underconfident to well-calibrated.<sup>4</sup>

Using the *control* approach, researchers regress the proposed correlate of overconfidence on selfevaluations and control for performance. In this approach, the partial regression coefficient estimate on self-evaluation with performance as a control is *precisely* the same as the coefficient estimate on the residualized estimate (with or without performance as a control). Controlling for performance rather than (or in addition to) residualizing self-evaluation has the benefit of reducing error variance in the analysis of the outcome measure, thereby providing a more precise estimate.

Using the *difference* approach, researchers subtract the measure of objective performance from the self-evaluation of performance and keep the difference as a measure of overconfidence. Researchers will sometimes use a difference score and also control for objective performance. If they do, the coefficient on the difference score is precisely equivalent to that on self-evaluation when controlling for performance, so the estimate is equivalent to that from the residual approach. When researchers control for performance in analysis of a proposed correlate, whether the focal variable is the residual, selfevaluation, or difference, the focal coefficient is precisely the same.

Table A1 in the supplement characterizes a sample of recent articles in terms of their usage of residuals, difference scores (including some that control for performance, making their usage equivalent to residuals), both, or measures that are equivalent to either (e.g., a measure of confidence without accounting for performance). Of 31 coded articles, 28 used a difference score or equivalent measure and

<sup>&</sup>lt;sup>3</sup> Parker and Stone (2014) refer to the residual approach as *unjustified confidence* and the difference score approach as *overconfidence*.

<sup>&</sup>lt;sup>4</sup> As a result, participants with negative residuals cannot be said to exhibit a self-diminishment bias: they may instead merely exhibit less of a self-enhancement bias or exhibit no bias at all (cf. John and Robins 1994, p214).

20 used a residual score or equivalent measure. This indicates each type of measure is commonly used.

## Self-Evaluate Using the Same vs. Different Metric

The self-evaluation may be assessed in the same metric as performance or in a different metric. Above, both performance (on a 10-item quiz) and self-evaluation (out of 10 items) are in the same metric. Alternatively, researchers may assess self-evaluations of performance in a different metric (e.g., a 1-7 Likert scale). If the self-evaluation is in a different metric, overconfidence may be assessed using the residual or covariate method, but not difference scores (even if the variables have been standardized).

#### Self-Evaluate Performance vs. Ability

Participants may be asked to evaluate their performance or their ability. The cases above represent self-evaluations of task-specific performance. In other cases, the self-evaluation may instead be an evaluation of ability. For example, after completing a 10-item financial literacy quiz, participants may report how well they performed on a 1 to 7 scale (performance), or they may report how knowledgeable they are about financial matters on a 1 to 7 scale (ability). Researchers residualize this measure of self-evaluated ability on performance (or control for performance in multiple regression) to consider the role of overconfidence. Although typical examples of self-evaluated ability tend to be in a different metric, it could be assessed in the same metric. For example, in principle researchers could inquire about expected performance on a 10-item test drawn from the same test bank in order to assess ability in the same metric, enabling potential use of difference scores. When participants evaluate expected performance on an upcoming task, the evaluation assesses ability.

#### **Absolute vs. Relative Evaluation**

Performance and self-evaluations may be measured in absolute or relative terms. In each case above, the focus is on absolute performance. In Moore and Healy's (2008) parlance, this is overestimation. The same techniques are used when measuring relative performance (i.e., overplacement), such as percentile performance within some specified sample. Self-evaluations of relative standing may be measures of performance on a particular task or measures of evaluations of ability.

## Variations on a Theme

These variations may be assembled in any combination as long as it does not involve taking a difference between two measures in different metrics. Evaluations may also be assessed item-by-item to enable assessments of sensitivity or efficiency (e.g., Burson et al. 2006; Fleming and Lau 2014; Stankov and Crawford 1996). Each of the approaches described above could result in a measure of overconfidence that is confounded with ability. As a result, using any of these measures biases measures of the correlation between overconfidence and other measures. The confound and bias is present whether the residual, covariate, or difference approach is taken, for both overestimation and for overplacement, whether self-evaluations are of performance or of latent ability, and whether they use the same or different metrics.

Such measures of overconfidence are sometimes, but not always, interpreted as measures of stable, general individual differences in overconfidence. For example, in the finance literature, measures of overconfidence on some unrelated tasks are used as correlates of trading activity or stock ownership (e.g., Grinblatt and Keloharju 2009; Ke 2021); this analysis requires both stability over time and consistency across domains to permit the researchers' preferred interpretation. In other cases, these measures are used as measures of possibly-transient, possibly-domain-specific overconfidence. For example, differences in overconfidence on a particular task with a particular partner are correlated with perceived competence and social status as rated by that partner (Anderson et al. 2012); no assumption of stability or consistency is required. Whether there are stable, general individual differences in overconfidence is a topic of ongoing debate (for arguments and evidence against, see e.g., Li et al. 2024; Moore and Dev 2017; Moore and Swift 2011; Moore and Schatz 2017; for evidence in favor, see e.g., Lawson et al. 2023, 2024; Binnedyk and Pennycook 2024). The analysis I present here addresses measured differences in overconfidence, whether they are stable or transient, general or domain-specific.

#### **The Potential Confound**

Before presenting a formal mathematical proof of the problem and simulation results, I provide an informal verbal account. The problem arises from four properties common to the paradigms described above. First, people typically differ in task-relevant ability (whether that ability is mere test-taking ability, general cognitive ability, domain-specific knowledge, current caffeination level, etc.) Second, people typically have at least partial insight into their ability. Third, performance is typically an imperfect measure of ability: it includes some noise and is unlikely to fully and only assess the construct it is intended to measure. Fourth, performance is typically ambiguous to the participant: people often only have an imperfect sense of how they did from the task itself prior to evaluating how they did.

Because performance is ambiguous, principles of Bayesian reasoning require that self-evaluations of performance regress toward people's prior beliefs (Moore and Healy 2008). This leads to regression toward ability on average under the weak assumption that people's beliefs about their own ability are correlated with their true ability. This implies that self-evaluations are a (noisy) weighted average of ability and performance. As a result, observing self-evaluations exceed performance signals that ability likely exceeds performance too. If two quiz-takers who have insight into their own ability each scored a 70%, and one believes she scored an 80% and the other believes he scored a 60%, there is a sound basis one may use to infer that the first test-taker has higher ability than the second.

Whenever performance is a noisy measure of ability, controlling for differences in performance is not sufficient to control for differences in ability (e.g., Birnbaum and Mellers 1979; Cohen et al. 2003; Culpepper & Aguinis 2011; Gillen et al. 2019; Kahneman 1965; Westfall & Yarkoni 2016). Because selfevaluations of performance are regressive toward ability, controlling for performance will result in residual variation in self-evaluation that is attributable to ability. Although the residuals are uncorrelated with performance by construction, they are still correlated with true ability. So measured overconfidence, computed via residuals or by controlling for performance, is confounded with ability. When people selfevaluate ability rather than performance (e.g., when assessing future expectations), the confound is more severe because the measure directly assesses ability rather than merely being contaminated by it.

The version above applies when self-evaluations are residualized or the analysis controls for performance. A related variant applies when researchers use difference scores. If the measure does not fully and only measure what it is believed to measure, performance will exhibit regression to the mean. People who are very high in ability will perform moderately highly, and people who are very low in ability will perform moderately poorly. As a result, the difference measure will also be confounded with ability. Consider again the 10-item quiz designed to measure financial literacy. Unbeknownst to researchers or participants, four of the items inadvertently assess trust in institutions instead. A financially-literate but average-trusting participant expected to get 9 answers correct, actually got 7 correct (5 of the 6 financial literacy questions and 2 of the 4 trust questions), and, due to the inherent ambiguity, reported that they got 8 correct. A less-literate but more-trusting participant expected to get 5 answers correct, actually got 7 correct (3 of the 6 financial literacy questions and all 4 of the trust questions), and due to the inherent ambiguity, reported that they got 7 correct (3 of the 6 financial literacy questions and all 4 of the trust questions), and due to the inherent ambiguity, reported that they got 6 correct. In this example, the apparent overconfidence of the first participant and underconfidence of the second participant reflect true differences in financial literacy, not a surplus nor deficit of confidence.<sup>5,6</sup>

Because private information regarding one's own ability affects self-evaluations but not evaluations of others' performance, this same logic leads to the same confound for overplacement.

## A Simulated Example of the Problem

A visual depiction of the problem is given in Figure 1. In this example, participants vary in ability, have perfect self-insight into their own ability, and have ambiguity about their own performance. In evaluating their performance, they rely on a noisy signal of their own performance and beliefs about their own ability. The outcome is correlated with ability, and only correlated with beliefs because their

<sup>&</sup>lt;sup>5</sup> Imperfect sampling and the role of error in the use of discrepancy scores are repeated themes in the overconfidence literature (e.g., Burson et al. 2006; Erev et al. 1994; Gigerenzer et al. 1991; Juslin 1994; Klayman et al. 1999). The focus of these critiques has been imperfect calibration and findings regarding aggregate overconfidence rather than the implications for individual-level measures of overconfidence described here. Concerns about inappropriate inferences regarding true scores when relying on measured scores are an old problem in measurement (e.g., Cochran 1968; Cronbach and Furby 1970; Lord 1956, 1958, 1960; McNemar 1958; Rogosa et al. 1982; Thomson 1924). This has led to an array of possible approaches to attempt to recover unbiased coefficient estimates (e.g., Cronbach and Furby 1970; Culpepper and Aguinis 2011; Kline 2005; Fuller 1987). Why has such a critical concern not been central in recent discussions of measures of overconfidence? A key factor may be that both overconfidence measures are uniquely susceptible to an illusion that only performance matters, not latent ability, because the performance measure is the target of self-evaluation. But because one's own performance is noisy and ambiguous, requiring people to incorporate their prior beliefs, the conclusion that latent ability does not matter is incorrect. <sup>6</sup> It would be inappropriate to attribute the error to the participant for misunderstanding what the quiz measures. The fact that the participant uses prior beliefs about financial literacy rather than a linear combination of financial literacy and trust in institutions should not be interpreted as overconfidence if they rely on the very construct the researchers themselves believe they are measuring. I return to this point toward the end of paper.

beliefs are accurate. There is no overconfidence in this example. Details regarding the model parameters used to simulate these data are given in the figure caption; the model is described in the next section.

#### Figure 1

Visual Depiction of How Measures of Overconfidence Are Confounded With Ability



*Note*. The 25 highest-ability individuals are depicted as filled circles. The 25 lowest-ability individuals are depicted as open circles. The 50 middle-ability individuals are depicted as crossed circles. Solid lines depict 45-degree lines; dashed lines depict best-fit regression lines. The parameters used for this example from the model detailed later are  $\rho = 1$ ,  $\lambda = .5$ ,  $\sigma_{\nu}^2 = 1$ ,  $\alpha = .5$ ,  $\sigma_{\nu}^2 = .25$ ,  $\beta = 1$ ,  $\sigma_{\epsilon}^2 = .25$ .  $\beta$  represents the coefficient on ability predicting outcome, and  $\sigma_{\epsilon}^2$  represents the error variance of the outcome.

Panel A shows the relationship between ability and performance for 100 people. In this example, the performance measure is both noisy and regressive; the 45-degree line is given by the solid line. The 25 highest-ability people are depicted as filled circles and the 25 lowest-ability people are depicted as open circles; the 50 middle-ability people are depicted as crossed circles.

Panel B depicts self-evaluations as a function of performance when self-evaluations are regressive toward ability. The 50 individuals classified as overconfident by the residual score (i.e., the vertical difference between each point and the dashed best-fit line) include 22 of the 25 highest-ability individuals and only 5 of the 25 lowest-ability individuals, and the 50 individuals classified as underconfident by the residual score include 20 of the 25 lowest-ability individuals and only 3 of the 25 highest-ability individuals. Similarly, the 44 individuals classified as overconfident by the difference score (i.e., the vertical difference between each point and the solid 45-degree line) included 17 of the 25 highest-ability individuals and only 6 of the 25 lowest-ability individuals, whereas the 56 individuals classified as underconfident by the difference score included 19 of the 25 lowest-ability individuals and only 8 of the 25 highest-ability individuals. In this example, classifying individuals by overconfidence effectively, albeit imperfectly, classifies them by ability ( $r_{ability,residual} = 0.60$ ,  $r_{ability,difference} = 0.37$ ).

Panels C and D plot the correspondence between the residual measure (C) and difference score (D) and an arbitrary correlate of ability. As is evident in this example, these correlates of ability are positively correlated with both measures of overconfidence, despite the fact that both measures of overconfidence account for performance and (in this example) there is no true overconfidence. As is derived below, the problem in (C) arises from the simulated measurement error in performance whereas the problem in (D) arises from the simulated regression to the mean in performance.

#### Modeling the Bias in Measures of Overconfidence

I next formalize the model described above. A straightforward extension of Moore and Healy's (2008) model of overconfidence permits a focus on differences between people, so I adapt their notation.<sup>7</sup>

<sup>&</sup>lt;sup>7</sup> This relates to a discussion in Healy and Moore (2007) and footnote 2 in Moore and Healy's (2008) in which luck is separated from expectations of ability. The implication for bias in the measure of overconfidence is not addressed.

## **Model Setup**

## Latent Ability and Beliefs

People, denoted by *i*, differ in their task-specific ability or skill:

$$S_i \sim D(0, 1) \tag{1}$$

where D(0, 1) represents any distribution that has been standardized to have a mean of 0 and variance  $\sigma_s^2$  of 1. Depending on the task, ability may represent mere test-taking ability, general cognitive ability, domain-specific knowledge, current caffeination status, etc.

People's beliefs,  $\tilde{S}_i$ , are a function of their own ability:

$$S_i = \rho S_i + \zeta_i \tag{2}$$

where  $0 \le \rho \le 1$  and  $\zeta$  is independently drawn from any distribution with mean 0 and variance  $\sigma_{\zeta}^2 = 1 - \rho^2$ , such that  $\tilde{S}$  has a variance  $\sigma_{\tilde{S}}^2$  of 1 and the correlation between *S* and  $\tilde{S}$  is given by  $\rho$ .<sup>8</sup> Latent overconfidence is then given by  $\tilde{S}_i - S_i$ .

People often can and do have insight into their own ability, suggesting we ought to expect  $\rho > 0$ . The Subjective Numeracy Scale (Fagerlin et al., 2007) was developed to find a way for people to selfreport their own numeracy using a less-burdensome task than a math test. Objective financial literacy shows correspondence with subjective financial literacy (Lusardi & Mitchell 2017). Objective knowledge and subjective knowledge are correlated across a range of domains (Alba & Hutchinson, 2000; Carlson et al., 2009). Across multiple domains, there is good reason to expect people have at least partial insight into their own abilities. Partial but incomplete insight into one's own ability can be modeled as  $0 < \rho < 1$ .

If  $\rho = 1$ , then  $\tilde{S}_i = S_i$  and there is no latent overconfidence. This perfect self-insight condition is of interest not because it is likely to be accurate, but rather because it presents an important null model to

<sup>&</sup>lt;sup>8</sup> This sets the mean level of overconfidence to 0, but this is merely for convenience. As the current discussion only addresses relatively more or less overconfident individuals, one can arbitrarily shift beliefs by changing the mean of  $\zeta$  without any substantive impact on the argument developed here regarding the confound with ability. This also equates the variance of beliefs to the variance of ability. Relaxing this constraint and freeing the constraint on  $\sigma_{\zeta}^2$  would make the correlation between *S* and  $\tilde{S}$  equal to  $\frac{\rho}{\sqrt{\rho^2 + \sigma_{\zeta}^2}}$ , but no other consequential implications for the model.

address: Is there evidence of a correlation with measures of overconfidence even when latent overconfidence does not exist? An answer of "yes" would be troubling. The assumption of perfect selfinsight is assuredly wrong in many—if not all—cases. But if our approach to assessing overconfidence and its correspondence with other constructs finds evidence in its absence, we must rethink that approach.

#### **Observable Performance and Self-Evaluations**

Although ability,  $S_i$ , varies across people, it is not directly observable. Instead, people's performance,  $P_i$ , is assessed on a particular task. Performance reflects a combination of ability and luck:

$$P_i = \lambda S_i + \nu_i \tag{3}$$

where  $0 \le \lambda \le 1$  and luck,  $v_i$ , is independently drawn from any distribution with mean 0 and variance  $\sigma_v^2$ .  $\lambda$  represents performance's loading on ability. A perfect measure that fully and only captures the focal ability (possibly with measurement error) has  $\lambda = 1$ ; an invalid measure (e.g., pure noise or a measure of an unrelated construct) has  $\lambda = 0$ . Consider a researcher measuring individual differences in intelligence using either (a) a test consisting of three of Raven's progressive matrices, or (b) a phrenologist's head measurements. Both measures contain noise, but for Raven's matrices we expect  $\lambda > 0$  (whether or not  $\lambda = 1$ ) whereas for the phrenologist's head measurements we expect  $\lambda = 0$ .

People's self-evaluations,  $\tilde{P}_i$ , are their noisy attempts to evaluate their own performance,  $P_i$ . After feedback, performance may be unambiguous. But prior to feedback, people generally have ambiguity regarding how they performed. Following Moore and Healy (2008), the presence of such uncertainty should lead to self-evaluations that incorporate prior beliefs through Bayesian-like reasoning (whether or not people are proper Bayesian updaters). For Moore and Healy, these prior beliefs represented beliefs about the simplicity of the task. In the current model, these prior beliefs represent beliefs about one's own ability. The key extension is thus the variability in those prior beliefs. The result of this Bayesian-like reasoning is that people ought to evaluate their own performance as lying between their prior beliefs and their performance, plus noise, where the weight on prior beliefs increases with ambiguity:

$$\tilde{P}_i = \alpha \tilde{S}_i + (1 - \alpha) P_i + v_i \tag{4}$$

where  $0 \le \alpha \le 1$  and  $v_i$  is independently drawn from any distribution with mean 0 and variance  $\sigma_v^{2.9} \alpha$ represents the ambiguity of someone assessing their own performance. As ambiguity increases,  $\alpha$ approaches 1, and self-evaluations of performance reflect beliefs about ability to a greater extent. When self-evaluations are measures of ability rather than measures of performance, as they are for expectations,  $\alpha = 1$ , because the measure is only a measure of ability and is not designed to assess performance at all.

## Figure 2

Measurement Model of Relationships Among Ability, Beliefs, Performance, and Self-Evaluations



*Note*. Panel A depicts imperfect but correlated beliefs. Panel B depicts the equivalent model when constraining  $\rho = 1$ . This implies  $S_i = \tilde{S}_i$  given that  $E[\zeta] = 0$  and  $\sigma_{\zeta}^2 = 1 - \rho^2 = 0$ .

<sup>&</sup>lt;sup>9</sup> If people knew their true performance,  $P_i$ , they could simply report it directly.  $P_i$  enters their beliefs but is not used directly because participants receive a noisy signal of their performance. That noise is then folded into  $v_i$ , leaving the signal to enter the equation directly. See Moore and Healy (2008).

The measurement model is depicted in Figure 2. Panel A depicts the general version, in which beliefs are imperfectly correlated with ability. Panel B represents a simplified version of the model implied by constraining  $\rho = 1$ .

This model distinguishes two reasons why  $P_i$  and  $\tilde{P}_i$  may be uncorrelated. First, they will be uncorrelated if  $\alpha = 1$  and  $\rho = 0$ . This leads to the typical interpretation: performance (may) provide a noisy signal of ability, and people can assess neither how they performed nor their true ability. Second, they will also be uncorrelated if  $\alpha = 1$  and  $\lambda = 0$ . Here, they are uncorrelated because performance does not provide a signal of ability, even though people have insight into their true ability. The comparison of these two cases implies observing no correlation between  $P_i$  and  $\tilde{P}_i$  is not sufficient to claim that  $\rho = 0$ .

If  $\rho = 1$ , then  $\tilde{S}_i = S_i$  because  $E[\zeta] = 0$  and  $\sigma_{\zeta}^2 = 0$ . In this case, there is no overconfidence. If

 $\rho < 1$ , there may be differences in overconfidence, but that does not imply there are *relevant* differences in overconfidence. When researchers introduce a candidate correlate and argue for its relation to overconfidence, they typically (implicitly or explicitly) argue for a causal arrow to or from  $\tilde{S}_i$ . For  $\rho > 0$ , beliefs may be correlated with the third variable entirely due to their mutual correlations with ability.

#### **Computing Measures of Overconfidence**

The researcher aims to isolate the role of overestimation by: (a) regressing  $\tilde{P}_i$  on  $P_i$  and keeping the residual, (b) including both self-evaluation  $\tilde{P}_i$  and performance  $P_i$  in a single multiple regression model predicting the candidate correlate, or (c) taking the difference between  $\tilde{P}_i$  and  $P_i$ . I first address the residual and regression approaches together, as they result in equivalent coefficients, and then the difference score.

#### **Residual and Multiple Regression Approaches**

To calculate overconfidence via residuals, researchers regress self-evaluations on performance:<sup>10</sup>

$$\tilde{P}_i = \gamma P_i + \epsilon_i \tag{5}$$

<sup>&</sup>lt;sup>10</sup> Throughout, I exclude intercepts for simplicity. Because my focus is on differences in overconfidence rather than mean levels of overconfidence, intercepts can be accounted for through recentering.

The residuals,  $e_i = \hat{\epsilon}_i$ , are kept as measures of overconfidence. Because performance is noisy and selfevaluations use prior beliefs about ability, the expected errors (and thus the residuals) vary with ability:

$$E[\epsilon|S] = \rho \left(1 - \frac{\lambda^2}{\lambda^2 + \sigma_v^2}\right) \alpha S \tag{6}$$

The derivation is in the Supplement. For reasonable sample sizes such that  $e_i \cong \epsilon_i$ , the residual from regressing self-evaluation on performance is positively confounded with ability if  $\rho\left(1-\frac{\lambda^2}{\lambda^2+\sigma_v^2}\right)\alpha > 0$ . That is, there is a confound if three conditions hold. First, true ability and beliefs about one own ability must be positively correlated ( $\rho > 0$ ). The confound is maximized if people have perfect self-insight and there are no differences in overconfidence ( $\rho = 1$ ). Second, there must be error in

the performance measure that is not attributable to ability  $(\sigma_{\nu}^2 > 0, \text{ making } \left(1 - \frac{\lambda^2}{\lambda^2 + \sigma_{\nu}^2}\right) > 0)$ . The

absence of measurement error is the exception, not the rule, so this condition is likely to be met.<sup>11</sup> Third, self-evaluations must be related to beliefs conditional on performance ( $\alpha > 0$ ). Any application of basic Bayesian logic in the presence of uncertainty about performance will lead to a direct effect of beliefs on self-evaluations, as will measures that include self-evaluations of ability rather than performance, so this condition is likely to be met as well.

Multiple regression can be written as a regression of residuals on residuals. When predicting a candidate correlate, the coefficient on evaluations controlling for performance is the same as the coefficient on residualized evaluations. The multiple regression estimate will typically be more-precise.

Given the broader literature on measurement error in predictors (e.g., Birnbaum and Mellers 1979; Cohen et al. 2003; Culpepper and Aguinis 2011; Gillen et al. 2019; Kahneman 1965; Westfall and Yarkoni 2016), why does the current paradigm deserve special consideration? First, because performance measures are the target of self-evaluations, there may be a mistaken belief that performance itself is the object of inquiry rather than also representing a measure of latent ability. Second, without the extension of

<sup>&</sup>lt;sup>11</sup> Note that under this interpretation, we assume all variance in performance is attributable to either ability or measurement error. In the general discussion I consider when  $\sigma_{\nu}^2$  includes systematic error (e.g., other constructs).

Moore and Healy's (2008) model, it may not be transparent to all researchers that self-evaluations are directly affected by beliefs about ability. Together, these neglected properties grant a false sense of security regarding the impact of measurement error in performance.

## Difference Score Approach

To assess overconfidence via a difference score, one subtracts performance from self-evaluation:

$$\Delta_i = \tilde{P}_i - P_i \tag{7}$$

In expectation, this difference score is also a function of ability:

$$E[\Delta|S] = (\rho - \lambda)\alpha S \tag{8}$$

The derivation is in the Supplement. The difference is positively confounded with ability if  $(\rho - \lambda)\alpha > 0$ . Once again, it is positively confounded if a specifiable set of conditions hold. First, beliefs must be positively correlated with ability ( $\rho > 0$ ). Second, performance must load sufficiently poorly on ability ( $\lambda < \rho$ ). Third, self-evaluations must be related to beliefs conditional on performance ( $\alpha > 0$ ).<sup>12</sup> If people hold accurate beliefs, then in the idealized case in which the measure of performance fully and only measures the construct that researchers and participants think it measures,  $\lambda = 1$  and there is no association between the difference score and ability. (There will also be no association if beliefs are as imperfectly related to ability as performance is, i.e., if  $\rho = \lambda$ .) If self-evaluation only depends on performance and not beliefs, then  $\alpha = 0$  and there is no relationship between the difference and ability.

Unlike the residual measure, and aligning with typical critiques of difference scores, the difference score measure can also be negatively confounded with ability if  $\rho < \lambda$ .<sup>13</sup> This implies that for

<sup>&</sup>lt;sup>12</sup> In this idealized case, measurement error  $(\sigma_v^2)$  is inconsequential. In practice,  $P_i$  is often bounded such that  $\sigma_v^2 > 0$  would likely drive  $\lambda$  down.

<sup>&</sup>lt;sup>13</sup> Prior critiques have addressed difference scores' confound with their component measures: what appears to be a property of the difference may instead reflect a property of one of the components (e.g., Cronbach and Furby 1970; Cohen et al. 2003; Edwards and Parry 1993; Griffin et al. 1999; Johns 1981; Wall and Payne 1973; Zuckerman and Knee 1996). Response Surface Analysis via polynomial regression (e.g., Edwards 1994; Barranti et al. 2017; Humberg, Dufner, et al. 2019; Humberg, Nestler, and Back 2019) and Condition-based Regression Analysis (e.g., Humberg et al. 2018a, 2019) aim to establish alternative conditions to assess whether the active ingredient is a discrepancy or positive self-evaluation. The present concern is a confound of self-evaluations with ability. It is relevant whether one is interested in the discrepancy or positive self-evaluation. In addition to other concerns regarding these regression-based approaches (Krueger et al. 2017; Fiedler 2021; cf. Humberg et al. 2018b, 2022), they do not distinguish between performance and ability, and so are equally susceptible to the concerns I raise here.

certain parameter configurations, there can be cases in which the residual measure of overconfidence is positively correlated with some outcome, the difference score measure is negatively correlated with that same outcome, and each correlation is entirely attributable to the measure's confound with ability.

## **Comparing the Biases and Additional Variations**

As long as there is ambiguity ( $\alpha > 0$ ) and performance loads on ability ( $\lambda > 0$ ), the residual score's bias will be more-positive than the difference score's bias if  $\sigma_{\nu}^2 > \lambda(\rho - \lambda)$ .<sup>14</sup> This will be the case if any of three conditions are met: (a) there is enough error in performance, (b) self-insight is low enough, or (c) the loading of performance on ability is not close enough to  $\frac{\rho}{2}$ .

If researchers use a difference score but control for performance, the covariate-adjusted regression coefficient and statistical test of the difference score is precisely equivalent to that using the multiple regression approach. The intuition is that the coefficient on the difference score is interpreted as "all else constant," and when "all else" includes performance, the only way the difference changes holding performance constant is by the evaluation changing. Thus, because the multiple regression approach leads to the same bias as the residual score approach, using difference scores while controlling for performance has the same bias as the residual score calculation.

For both residual and difference measures, the same confound holds for both overestimation of absolute performance and overplacement of relative performance. Evaluations of one's own performance are a function of one's own idiosyncratic ability and idiosyncratic beliefs; evaluations of others' performances are not. As a result, the conditional expectation of evaluations of others' performances are also not a function of one's own ability. The additional terms in overplacement drop out, leaving a bias in one's relative performance (overplacement) that matches the bias in one's absolute performance (overplacement) that matches the bias in one's absolute performance

<sup>&</sup>lt;sup>14</sup> If  $\alpha = 0$ , neither coefficient is biased; if  $\lambda = 0$ , they each exhibit the same bias.

<sup>&</sup>lt;sup>15</sup> This may appear to be at odds with Moore and Healy (2008) who show reversals between overestimation and overplacement. Because *S* in their model represents a quiz's simplicity, it takes a common value for every individual such that prior beliefs are relevant for both oneself and others. In my model,  $S_i$  represents an individual's ability, so it takes an independent value for every individual such that prior beliefs are not relevant for others' performance.

## Figure 3

Simulation Results of Bias in Residual and Difference Scores as a Function of  $\rho$ ,  $\lambda$ ,  $\alpha$ , and  $\sigma_{\nu}^2$ 



*Note*.  $\sigma_v^2$  is plotted but not apparent as it does not affect the bias. For the residual panel, when  $\sigma_v^2 = \lambda = 0$ , estimates are missing. This is because performance does not vary and so cannot be used as a covariate.

## Simulation Results

Simulations show that these asymptotic results hold in reasonable sample sizes. For all factorial combinations of  $\rho$ ,  $\lambda$ ,  $\alpha$ ,  $\sigma_v^2$ , and  $\sigma_v^2$  taking a value in [0, 0.2, 0.4, 0.6, 0.8, 1.0], I simulated 1,000 samples of 100 observations each. In each sample, ability was drawn from a standard normal distribution, and

error terms were drawn from standard normal distributions scaled by the corresponding variance. Residuals and difference scores were predicted by ability as derived in the model. Average coefficients are depicted in Figure 3. Repeating the simulations using log-normal distributions recovered the same average coefficients, reflecting the principle that the variables need not be normally distributed.

#### **Two Empirical Applications**

Common measures of overestimation and overplacement are confounded with ability in reasonable null models in which overconfidence does not exist or does not matter (i.e., only relates to candidate constructs through ability). The magnitude of the confound depends on the parameter values and the model is undoubtedly a simplification. Does the proposed confound actually matter in practice? Two empirical applications indicate yes.

In the first application, I examine a case in which an outcome measure (performance on a related task) should be related to ability and should not be related to overconfidence. If the outcome is related to measures of overconfidence, this can be explained parsimoniously through the confound with ability. Given strong prior beliefs about true relationships among ability, overconfidence, and the outcome measure, this demonstrates that the confound can hold in practice.

In the second application, I re-examine a related pair of published findings in which the results are interpreted in terms of overconfidence. Using the current model, I find that they are also compatible with plausible parameter values in which there is no consequential overconfidence (or no overconfidence at all). This indicates the confound provides a plausible alternative interpretation of published results. Whereas the first application indicates measures of overconfidence can lead to problematic conclusions, the second application indicates evidence from the literature which has been interpreted in terms of overconfidence may not be sufficient to imply a causal or consequent role of overconfidence, and are consistent with perfectly-calibrated beliefs (i.e., the absence of overconfidence).

#### **Empirical Application I: Overconfidence Predicts Performance**

The analysis above indicates overconfidence measures are biased under plausible conditions. Can this bias affect inferences? The answer depends on typical parameters: how strongly performance loads

on ability, how much measurement error is in performance, whether people exhibit (even imperfect) Bayesian updating, and the degree to which people have self-insight into their own ability. To test this, I first examine a case where there is a measure of performance, people self-evaluate their performance, and there is an outcome measure which is a priori likely related to ability and unrelated to overconfidence.

Using data from Moore and Healy (2008), I consider a case where the relevant outcome is a future measure of performance.<sup>16</sup> Future performance cannot cause past overconfidence, and it is unlikely that past overconfidence causes future performance in a way that is not attenuated as the number of intervening tasks increases. As a result, finding that past residuals or difference scores predict future performance would suggest that past residuals or difference scores are confounded with ability.

## **Overconfidence Paradigm**

Moore and Healy (2008) collected data from 82 college undergraduates on many measures. I describe the relevant components here and refer the reader to Moore and Healy for complete details. Participants completed 18 10-item trivia quizzes: an easy, medium, and hard quiz on each of six topics. The quizzes were presented sequentially in six blocks. Each block contained an easy, a medium, and a hard quiz on different topics, in randomized order. In addition to other measures, for each quiz, participants: (a) provided a pre-quiz measure of expected performance, (b) took the quiz, and (c) provided a post-quiz measure of estimated performance.

Analysis of these data requires addressing two key complications. First, the quizzes differed in difficulty. If trivia quiz ability exists, performing well on one quiz should predict performing well on another, all else equal. But performing well on one quiz is a signal not only that ability may be high, but also that difficulty may be low. Indeed, the block-randomized design leads to negative autocorrelation in difficulty between successive quizzes. In expectation in the first 5 blocks, there is a 44% chance that a hard quiz is followed by an easy quiz but only an 11% chance that a hard quiz is followed by another hard quiz. This mechanically generates a negative correlation between performance on one quiz and

<sup>&</sup>lt;sup>16</sup> Moore and Healy do not make the inferential error reported here. Rather, the availability and richness of their data (<u>https://osf.io/6tecy/</u>) present an opportunity to examine whether the error can affect real inferences.

performance on the subsequent quiz. To address this design characteristic, I consider expectations, estimates, and performance for *blocks* (where each block is a triplet of quizzes) rather than for *quizzes*. A block always consists of one easy, one medium, and one hard quiz, substantially reducing the extent to which performance on one block is negatively correlated with performance on other blocks.

Second, the data provide rich within-subject data but with a modest sample size for betweensubject analyses by current standards (82 participants). To exploit the within-subject data, I consider performance on sets of sequential quizzes, clustering errors by subject. For example, to examine ability, I regress block performance on prior block performance such that each participant contributes five observations: block 2 performance as a function of block 1 performance, block 3 performance as a function of block 2 performance, etc. The analysis accounts for non-independence through clustered errors using the *lm\_robust()* function from the *estimatr* package (Blair et al., 2022).

## A Puzzle: Overconfidence Predicts Subsequent Performance

Using the first five quiz blocks to provide measures of performance and self-evaluations, I follow the approaches from the literature to construct three measures of overestimation: residualized selfevaluations, self-evaluation controlling for performance, and difference scores. I similarly construct the three corresponding measures of overplacement using Moore and Healy's (2008) overplacement measure.

Overestimation as assessed via residualized self-evaluations predicted performance in the next block, as did the partial coefficient on self-evaluations controlling for performance. This latter coefficient is necessarily equal to that on the residual, but estimated more precisely because performance accounts for additional variance in the dependent variable in the multiple regression analysis. Overestimation as assessed via the difference score did not significantly predict subsequent performance. Overplacement as assessed via residualized relative self-evaluations predicted relative performance in the next block, as did overplacement as assessed via the partial coefficient on relative self-evaluations controlling for relative performance. Overplacement as assessed via the difference score also predicted subsequent performance, but the coefficients were significantly negative. Results are detailed in Table 1.

## Table 1

Overestimation								
Measure Type	Predictor Lag	Coef (SE)	Test statistic	95% CI				
Residual Score	1 (Primary)	0.298 (0.141)	t(38) = 2.11, p = .042	[0.012, 0.584]				
	2	0.381 (0.164)	t(33) = 2.33, p = .026	[0.048, 0.715]				
	3	0.323 (0.208)	t(28) = 1.56, p = .130	[-0.102, 0.748]				
	4	0.434 (0.233)	t(20) = 1.86, p = .077	[-0.051, 0.918]				
	5	0.295 (0.403)	t(15) = 0.73, p = .476	[-0.562, 1.151]				
Multiple Pegression	1 (Primary)	0 208 (0 088)	t(38) = 3.40 n = 002	[0 121 0 476]				
Multiple Regression	$\frac{1}{2}$	0.298(0.088) 0.381(0.092)	t(33) = 3.40, p = .002 t(33) = 4.14, p < .001	[0.121, 0.470] [0.194, 0.569]				
	2 3	0.381(0.092) 0.323(0.164)	t(33) = 4.14, p < .001 t(28) = 1.07, p = .058	[0.194, 0.309]				
	J 1	0.323(0.104) 0.434(0.155)	t(20) = 1.97, p = .030 t(20) = 2.81, p = .011	[-0.012, 0.059]				
	+ 5	0.434(0.133) 0.205(0.365)	t(20) = 2.81, p = .011 t(15) = 0.81, p = .432	[0.111, 0.750]				
	5	0.293(0.303)	t(15) = 0.61, p = .452	[-0.462, 1.071]				
Difference Score	1 (Primary)	0.085 (0.121)	t(37) = 0.71, p = .485	[-0.160, 0.330]				
	2	0.137 (0.131)	t(31) = 1.05, p = .302	[-0.129, 0.404]				
	3	0.077 (0.177)	t(25) = 0.43, p = .668	[-0.287, 0.440]				
	4	0.052 (0.174)	t(17) = 0.30, p = .771	[-0.316, 0.419]				
	5	-0.001 (0.357)	t(14) = -0.00, p = .997	[-0.765, 0.763]				
		Overplacemen	ıt					
Measure Type	Predictor Lag	Coef (SE)	Test statistic	95% CI				
Residual Score	1 (Primary)	0.401 (0.103)	t(42) = 3.89, p < .001	[0.193, 0.609]				
	2	0.419 (0.106)	t(42) = 3.95, p < .001	[0.205, 0.632]				
	3	0.417 (0.120)	t(41) = 3.48, p = .001	[0.175, 0.658]				
	4	0.437 (0.138)	t(33) = 3.16, p = .003	[0.156, 0.717]				
	5	0.400 (0.249)	t(22) = 1.61, p = .122	[-0.116, 0.915]				
Multiple Regression	1 (Primary)	0.401 (0.056)	t(42) = 7.13, p < .001	[0.287, 0.514]				
	2	0.419 (0.062)	t(42) = 6.73, p < .001	[0.293, 0.544]				
	3	0.417 (0.085)	t(41) = 4.91, p < .001	[0.246, 0.588]				
	4	0.437 (0.093)	t(33) = 4.68, p < .001	[0.247, 0.626]				
	5	0.400 (0.176)	t(22) = 2.28, p = .033	[0.036, 0.764]				
Difference Score	1 (Primary)	-0.188 (0.062)	t(52) = -3.05, p = .004	[-0.311, 0.064]				
	2	-0.152 (0.062)	t(49) = -2.45, p = .018	[-0.277, -0.027]				
	3	-0.191 (0.077)	t(48) = -2.49, p = .016	[-0.346, -0.037]				
	4	-0.206 (0.073)	t(36) = -2.83, p = .007	[-0.354, -0.059]				
	5	-0 202 (0 123)	t(24) = -1.64 p = 114	[-0.456_0.052]				

Regression Coefficients Predicting Subsequent Performance

*Note*. All degrees of freedom estimated given cluster-robust standard errors. In all analyses, residuals were calculated using only observations included in the relevant analysis.

One might argue that overconfidence improves subsequent trivia quiz ability (e.g., via selfefficacy). Such explanations would imply the correlation would be stronger for adjacent blocks. Yet there is no evidence that coefficients from the residual or covariate analyses diminish with lags. See Table 1.

Instead, the conceptual model proposed above provides a parsimonious explanation: performance in the current and future blocks are both driven by ability, and the measure of overconfidence is confounded with ability. The fact that difference scores did not predict future performance (for overestimation) or negatively predicted future performance (for overplacement) may be attributable to imperfect self-insight.

To examine whether the conceptual model is consistent with the residual and multiple regression analyses, I examine whether four necessary components are in place: (1) Are there differences in ability? (2) Do participants have insight into their own ability? (3) Does trivia quiz performance contain error as a measure of trivia quiz ability? and (4) Does a proxy for ability predict self-evaluations beyond performance? I focus on the overestimation results.

## 1. Are There Differences in Ability? Yes

If performance is correlated across blocks, there is evidence of systematic differences in trivia quiz ability.<sup>17</sup> I regress performance in block *t* on prior performance in block *t*-1, clustering errors by subject. The coefficient on lagged performance was 0.754 (SE = 0.045, t(31) = 16.59, p < .001; 95% CI: [0.662, 0.847]), indicating high performance on one block is strongly associated with high performance on the next block. When an analogous approach was used with block *t*-2, *t*-3, etc., there was no evidence of a relationship that decays with lag (lag 2: b = 0.783, SE = 0.057; lag 3: b = 0.788, SE = 0.076; lag 4: b = 0.796, SE = 0.088; lag 5: b = 0.756, SE = 0.094). These results are consistent with the presence of differences in trivia quiz ability ( $\sigma_s^2 > 0$ ), which are noisily measured by each quiz.

<sup>&</sup>lt;sup>17</sup> Ability includes skill, knowledge, and other necessary inputs that remain stable during the study.

## 2. Do Participants Have Insight Into Their Own Ability? Yes

If participants can predict how they will perform on a quiz, it suggests they have some insight into their own trivia quiz ability. I regress expectations on subsequent performance, clustering errors by subject. At the time of the expectation, neither the quiz difficulty nor the quiz topic was known yet. The coefficient on performance was 0.467 (SE = 0.068, t(31) = 6.87, p < .001, 95% CI: [0.329, 0.606]) indicating participants have partial insight into how they will perform.<sup>18</sup> Given the limited information available to participants, this is most readily attributable to self-insight into their own ability ( $\rho > 0$ ).

## 3. Does Trivia Quiz Performance Contain Error as a Measure of Trivia Quiz Ability? Yes

In regressing performance on lagged performance, it comes as no surprise that indeed,  $R^2 = 0.538$ < 1, indicating that it is not the case that both measures are errorless indicators of the same construct. Performance as a measure of ability contains error ( $\sigma_{\nu}^2 > 0$ ).

## 4. Does a Proxy for Ability Predict Self-Evaluations Beyond Performance? Yes

The last required component is that self-evaluations of performance are regressive toward ability. Ability is not directly observable, but subsequent performance provides a noisy proxy; for this purpose, the only requirement is that it contains sufficient signal. I regress self-evaluations on current performance and subsequent performance, where subsequent performance serves as a noisy proxy for ability. The coefficient on current performance was 0.880 (*SE* = 0.035, *t*(49) = 24.94, *p* < .001; 95% CI: [0.809, 0.951]), indicating that participants have some idea of how well they did on each block (though this coefficient also partially captures the role of ability). Critically, the coefficient on subsequent performance was 0.099 (*SE* = 0.029, *t*(49) = 3.43, *p* = .001, 95% CI: [0.041, 0.156]): controlling for current performance, future performance is predictive of current self-evaluations. The magnitude of this coefficient did not attenuate with more intervening blocks (1 intervening block: *b* = 0.119, *SE* = 0.034; 2

<sup>&</sup>lt;sup>18</sup> One may be concerned that participants are aware of the difficulty of the third quiz in each block, thereby inflating this relationship. If the first quiz was hard and the second was medium, participants could determine that the third would be easy. The main result also holds using only the first quiz from each block (adjusted for difficulty), which was completely randomized (b = 0.233, SE = 0.048, t(31) = 4.80, p < .001, 95% CI: [0.134, 0.332]). The coefficient was no stronger when using only the third quiz (b = 0.177, SE = 0.041, t(45) = 4.33, p < .001).

intervening blocks: b = 0.101, SE = 0.052; 3 intervening blocks: b = 0.141, SE = 0.044; 4 intervening blocks: b = 0.099, SE = 0.107). This suggests that post-quiz self-evaluations are regressive toward idiosyncratic ability ( $\alpha > 0$ ) in addition to assessing performance as intended.

#### Subsequent Performance Illustrates the Problem Regarding Other Correlates of Overconfidence

Although the puzzle suggests that overconfidence predicts future performance, a more parsimonious, and perhaps more probable, explanation is that there are differences in ability, people have self-insight, performance is a noisy measure of ability, and self-evaluations pick up ability in addition to performance. Therefore the measure of overconfidence is confounded with ability and ability is what predicts future performance. A key problem is that many findings in the literature of an association between overconfidence and other correlates use an approach equivalent to that in the puzzle above, but do not sufficiently consider the relevant alternative explanation after incompletely accounting for ability.

## **Empirical Application II: Correlating Overplacement With Status and Social Class**

The first application finds that measures of overconfidence predict subsequent performance. This can be readily explained by the confound between measures of overconfidence and ability. While this indicates a problematic conclusion one might draw from the data, it does not require us to reinterpret prior findings. Perhaps the whole endeavor is a statistical curiosity with little connection to substantive claims. Using another pair of findings, I examine how the model provides a potential alternative explanation for how differences in overconfidence correlate with other constructs. The goal is to consider whether this model could account for the results, not whether it rules out the original interpretations.

In both cases, I report a set of parameters compatible with the reported correlations. Two important caveats are in order. First, given the degrees of freedom, the parameter values I report are a subset of those that fit the data, not the only ones that fit the data. Second, and more importantly, one should not interpret the specific values as precise point estimates. This model, like all models, is wrong. Even if the links are correct, it is unlikely they capture the correct functional form. Instead, it informs qualitative conclusions about the relevant components: good vs. poor self-insight, high vs. low construct validity, high vs. low ambiguity, strong vs. weak relationship. The qualitative pattern is what matters.

## **IIA: Overplacement, Perceived Competence, and Status**

In six studies, Anderson et al. (2012) propose that holding ability constant, greater confidence generates higher status in the eyes of others. Study 1 was correlational and well-suited to examine using the current framework.<sup>19</sup> Could ability alone be responsible for the correlation with status in this study?

In Study 1, participants took a geography quiz and reported a self-evaluation of their quiz performance percentile and a self-evaluation of their general geography knowledge percentile. They then repeated the task with a partner. Their partner then rated their perceived competence using a similar percentile measure as well as their status. Self-evaluation was computed as the average of the taskspecific and general-knowledge self-evaluation percentiles. Overplacement was measured as the residual of self-evaluation percentile regressed on actual percentile. Overplacement as assessed via the residual was correlated with both partner-rated perceived competence and partner-rated status. These correlations were nearly as strong as the correlations of actual performance with perceived competence and status.

Might the results instead be explainable through links between *ability* and perceived competence and status instead? I consider sets of  $\rho$ ,  $\lambda$ ,  $\alpha$ ,  $\beta_{PC}$  (i.e., the causal impact of ability on perceived competence) and  $\beta_S$  (i.e., the causal impact of ability on rated status). For simplicity, I constrained error variances such that each variable had a variance of 1. It is possible to recover nearly-identical correlations to those observed in the data with  $\rho = 1$ ,  $\lambda = .7$ ,  $\alpha = .8$ ,  $\beta_{PC} = .55$  and  $\beta_S = .45$ . Given the number of parameters, multiple configurations fit similarly well (e.g.,  $\rho = .75$ ,  $\lambda = .55$ ,  $\alpha = .75$ ,  $\beta_{PC} = .75$ ,  $\beta_S =$ .6).<sup>20</sup> The empirical correlations and the correlations implied by these two sets of parameters are given in Table 2.

This paper is sometimes referenced for its methodology as a recent example of using residual scores to assess overconfidence (cf., Belmi et al. 2020; Cheng et al. 2021; Lyons et al. 2021; Murphy et al. 2015). As a result, beyond its substantive finding, its methodological approach is influential. Although

<sup>&</sup>lt;sup>19</sup> The model does not attempt to account for the results of the remaining studies, and therefore should not be interpreted as providing a counter-explanation for the paper's complete set of results.

<sup>&</sup>lt;sup>20</sup> The paper's text also reports the correlation between actual performance and self-evaluation was .56. The correlation implied by this latter set of parameters recovers .56 ( $\lambda \rho \alpha + (1 - \alpha) = .55 * .75 * .75 + .25$ ).

the present model has enough free parameters to readily fit the data, it has fewer free parameters than

implied by the paper's analysis, given additional required links from beliefs to outcomes.

## Table 2

Measure	Perceived Competence		Status		
Actual Performance	.39		.33		
Residual	.36		.26		
$\rho = 1, \lambda = .7, \alpha = .8$		$\beta_{PC} = .55$		$\beta_S = .45$	
Actual Performance	.39		.32		
Residual	.35		.28		
$\rho = .75, \lambda = .55, \alpha = .75$		$\beta_{PC} = .75$		$\beta_S = .6$	
Actual Performance	.41		.33		
Residual	.35		.28		

Empirical Application IIA: Observed and Modeled Correlations

## **IIB: Overplacement and Social Class**

Whereas Anderson et al. (2012) propose overconfidence causes status, Belmi et al. (2020) propose social class causes overconfidence due to the pursuit of status. (The model predicts a set of correlations, no matter the direction of the causal arrow; in both cases the target studies assessed correlations.) Could social class result in greater test-taking ability instead (e.g., mere ability to perform on tests, perhaps due to tests biased by social factors or differential access to test-preparation resources)? This example provides a more-stringent test of the model's explanatory ability given a larger set of correlations to be fit, including both residual measures and difference score measures of overconfidence.

Across four studies including three different tasks, the paper uses four measures of social class (self-report, income, education, and parental education), and multiple ways of measuring overplacement on each task. The prioritized measure is the residual score measure: self-evaluated percentile is regressed on actual performance percentile, and the residual is used as a measure of overplacement. The paper also uses the difference score measure: the difference between self-evaluated and actual percentiles.<sup>21</sup>

<sup>&</sup>lt;sup>21</sup> The paper also uses Edwards' (1995) proposed approach of calculating Wilks' lambda to compare coefficients across two regressions: one predicting self-evaluated percentile from social class and the other predicting actual percentile from social class. While I do not model that analysis here, I do examine whether the model parameters successfully capture the correlations between social class and both self-evaluated and actual percentiles.

Across studies, three different tests were used to assess overplacement. In Study 1, participants assessed whether two successive images matched. In Study 2, the test was a test intended to assess general cognitive ability (the Wonderlic Personnel Test). In Studies 3 and 4, the tests were trivia quizzes. Performance relies on a combination of skills, including mere test-taking ability which may be biased by social factors (e.g., differential access to test-preparation resources).

The paper reports meta-analytic estimates in its Tables 14 and 15. Each measure of social class was more-correlated with perceived performance than actual performance, each measure of social class correlated with residualized self-evaluations, and self-reported social class (but not income, education, nor parental education) correlated with the difference between perceived and actual percentiles. I examine whether plausible parameter values are consistent with the reported estimates (without covariates). These correlations, along with their 95% confidence intervals, are reproduced in the top panel of Table 3.<sup>22</sup>

## Table 3

Empirical Application IIB: Observed and Modeled Correlations Between Social Class and Measures of Perceived and Actual Performance and Overconfidence

Maasura	Salf Papart	Incomo	Education	Dorontol
Wiedsule	Sell-Kepolt	Income	Education	Falental
Perceived	.24 [.20, .28]	.11 [.11, .12]	.10 [.01, .18]	.14 [.08, .19]
Actual	.00 [11, .12]	.02 [12, .21]	.07 [03, .16]	.00 [12, .12]
Residual	.23 [.20, .26]	.11 [.04, .18]	.08 [01, .16]	.13 [.06, .19]
Difference	.13 [.05, .22]	.06 [09, .20]	01 [09, .08]	.08 [03, .20]
$\rho = 1, \lambda = .2, \alpha = .7$	$\beta = .3$	$\beta = .15$	$\beta = .1$	$\beta = .15$
Perceived	.23	.11	.08	.11
Actual	.06	.03	.02	.03
Residual	.22	.11	.07	.11
Difference	.16	.08	.05	.08
$\rho = .6, \lambda = .05, \alpha = .95$	$\beta = .4$	$\beta = .2$	$\beta = .15$	$\beta = .2$
Perceived	.23	.11	.09	.11
Actual	.02	.01	.01	.01
Residual	.23	.11	.09	.11
Difference	.15	.08	.06	.08

<sup>&</sup>lt;sup>22</sup> There are substantive differences across studies, and a single set of parameters may not be able to account for every study's results. This is not unique to this model, as this applies to any attempt to fit the same parameters across studies. The present model's constraints are explicit.  $\lambda$ ,  $\alpha$ , and  $\rho$  need not be consistent across studies.

As with the previous analysis, I examine whether any parameter configurations could generate these results. I again constrain error variances such that each variable had variance of 1. The correlations resulting from two such sets of parameters are reported in the lower two panels of Table 3. The middle panel represents a case in which people hold accurate beliefs ( $\rho = 1$  and  $\tilde{S}_i = S_i$ ) and there is no latent overplacement. The bottom panel represents a case in which people hold correlated beliefs ( $0 < \rho < 1$ ), so while some people may exhibit overplacement and others underplacement, social class only relates to beliefs through their respective correlations with test-taking ability. Across the four measures of social class, the parameters relating test-taking ability, beliefs, performance, and self-evaluations ( $\rho, \lambda, \alpha$ ) are held constant, but the parameter relating social class to test-taking ability ( $\beta$ ) is allowed to vary.

Note several aspects of this analysis. First and most relevantly, the fitted parameters can adequately account for the entire pattern of meta-analytic results without any role for latent overplacement (middle and bottom panels) and even with perfectly-calibrated beliefs (middle panel, in which  $\rho = 1$ ). The model uses 6 (middle panel) or 7 (bottom panel) free parameters for 16 correlations, though it captures the full pattern at least as well as the proposed theory, which would add 1 to 4 more parameters (i.e., to relate latent overconfidence to each measure of social class). Every correlation is well-within the 95% confidence interval, and more than two-thirds are within the implied 67% confidence interval, and more than two-thirds are within the residual score (for every measure of social class) and the difference score (for the self-report measure of social class) implicating a relationship between overplacement and social class.

Second, this neither rules out the authors' hypothesis nor indicates an impact of social class on test-taking ability. It merely presents an alternative interpretation of the data in which social class may not be related to overplacement. It is important to reinforce that under this interpretation, ability could refer to mere test-taking ability (attributable to e.g., differential access to test-preparation resources).

Third, multiple parameter configurations are each nearly-equally compatible with data. The two in the table are not the only two possibilities. Generally speaking, compatible parameters have low  $\lambda$ , high

 $\alpha$ , and a tradeoff between  $\rho$  and  $\beta$ . Figure 4 presents a set of compatible values of  $\rho$ ,  $\lambda$ ,  $\alpha$ , and  $\beta$  for self-reported social class. Similar values may be constructed for income, education, and parental education.

Finally, one should probe whether these parameter values are plausible. For low values of  $\beta$ , these data require high  $\rho$ , low  $\lambda$ , and high  $\alpha$ . Whether these are themselves plausible depends on outside knowledge and cannot be resolved using these data alone. (For example, a small correlation between self-evaluations and actual performance is not sufficient to imply  $\rho$  is low, as it may be due to low  $\lambda$  instead.) Other approaches are required to rule out the model's null. The paper also proposes a rich network of associations with desire for status and perceived competence, which are not included in this exercise.

## Figure 4



Empirical Application IIB: Compatible Parameter Values for Self-Reported Social Class

*Note*. Dark shaded regions indicate regions where all of the four target correlations implied by the set of parameters falls within the 95% confidence interval of the corresponding correlation for self-reported social class. Cases where each value is closer to the center of the confidence interval are shaded darkest. Lightly-shaded regions in different colors indicate regions where only a subset of the target correlations with self-reported social class fall within the 95% confidence interval (blue: performance; red: self-evaluation; green: residual; orange: difference.)

Table A3 in the supplement provides analogous parameter configurations consistent with findings for 16 different articles. Some of the parameter configurations are plausible, but others are not.

#### **General Discussion**

Both in theory and in practice, widely-used measures of overconfidence are confounded with ability. This applies to both overestimation and overplacement. While the biases in the residuals and difference scores are driven by different properties (the first by error in measurement, the second by imperfect construct validity), both are attributable to the fallacy of equating performance with ability. Below I discuss additional nuances and caveats of the model and end with some recommendations.

## **Model Nuances and Caveats**

#### Model Specifications

First, I model the measure of performance as a function of ability ( $S_i$ , weighted by  $\lambda$ ) and luck ( $v_i$ ). I have largely treated luck as though it is random measurement error. But it is possible to decompose luck into task-specific stable luck (e.g., someone taking a financial literacy quiz may be lucky that in part it assesses institutional trust rather than just financial literacy) as well as temporally inconsistent random luck (e.g., someone was distracted in the moment and misread an answer). Once luck is decomposed in this way, it is apparent that addressing measurement error is unlikely to completely address the problem for the residual. In this representation, reducing measurement error to 0 does not reduce  $\sigma_v^2$  to 0.

Second, if there is only measurement error, lack of *reliability* in the performance measure leads to lack of *validity* in residual measure. If that were the case, then increasing reliability of the performance measure could increase validity of the residual measure. I address this in the second recommendation.

Finally, this model represents but one functional form. It is unlikely to be perfectly specified. As one example, it posits a simple linear relationship between ability and beliefs. Yet people who are lessskilled may be less-well-calibrated regarding their own skill (Kruger & Dunning 1999). Despite statistical critiques (Krueger & Mueller 2002; see supplement), there is evidence for the unskilled-and-unaware effect that addresses those critiques (Feld et al. 2017; Jansen et al. 2021). Such a relationship will affect the magnitude and form of the confound, but it will not alleviate the problem. This reinforces that the model here represents *one* possible null to be ruled out, not *the only* possible null to be ruled out.

#### **Residuals vs. Difference Scores**

As noted throughout, different researchers have relied on residuals, or difference scores, or difference scores that in practice act like residuals. There remains debate about the relative merits of the two (see e.g., Belmi et al. 2020; Lyons et al. 2021; Parker & Stone 2014 for recent discussions of the two approaches by researchers interested in the substantive questions). The current model reinforces that neither one can be considered right and the other wrong; their relative biases depend on the relative contributions of ability and luck to the measure of performance.

Other proposals, particularly in the domain of positive self-views and self-enhancement, have proposed the use of polynomial regression, response surface analysis, or condition-based regression analysis (e.g., Edwards 1994; Humberg et al. 2018a, 2019) to address the concerns regarding difference scores and residuals. By themselves, they are not sufficient to account for the concern because in their base form they do not account for measurement error or construct mismatch. Only in conjunction with a strategy to address measurement error will they address reliability for the residual score measure. Even then, construct validity remains a concern for both the residual score and the difference score measures.

## Is $\lambda < 1$ Just a Form of Overconfidence?

Throughout, I note that the measure of performance must fully and only measure the target construct to make use of the difference score measure. This matters because the prior to which people regress must align with the construct the performance measure assesses. A mismatch, as in the case of a financial literacy scale with items that measure trust instead, is equivalent to construct invalidity, or  $\lambda < 1$ , which leads to the core problem for difference scores. (Recall that the measure may be highly reliable even with low validity.) A potential critique is that this simply represents a different form of improper confidence: people confidently rely upon a prior that should not apply and so regress to the wrong belief as a result. I argue we cannot be so quick to attribute such a problem to the participant's updating strategy (and therefore a variant of overconfidence) rather than the researcher's inferential strategy.

Consider again the phrenologist introduced earlier. Both the phrenologist and the participant may earnestly believe that the phrenologist is generating a diagnostic measure of intelligence. If the participant is asked how they perform on this measure of intelligence, but they have no idea about their own head measurements ( $\alpha = 1$ ), they will report their beliefs about their own intelligence. Of course, their score on the phrenology examination will be unrelated to their intelligence ( $\lambda = 0$ ). If people have partial self-insight ( $\rho > 0$ ), then on average, people who think they received a higher score from the phrenologist than they truly did will be more intelligent. A skeptic may argue: "That is overconfidence! The participant is inappropriately regressing their self-evaluation of performance on a nondiagnostic test to their beliefs about their own intelligence." But in such a case, it would be inappropriate to fault the participant for regressing to the very construct the researcher claims to be measuring with a worthless instrument. Thankfully, most researchers are not phrenologists and are using instruments with greater validity. But this provides only cold comfort.

This raises a thorny question regarding whether the effects of using misleading labels for a performance task should be considered overconfidence. Overconfidence should not depend on what the researcher believes a task measures. If we do not accept the overconfidence label in the case described above (in which the participant earnestly believes the measure is measuring the same construct the researcher earnestly believes it measures), we perhaps ought to be cautious in accepting an overconfidence label in the presence of a mislabeled instrument (when the participant earnestly believes the measure is measures; Ehrlinger & Dunning, 2003). *Beyond Overconfidence* 

# These concerns regarding the use of residual and difference scores represent a broader concern. They apply anytime there are two noisy measures of correlated (or identical) constructs where one is of interest and the other is to be ruled out. This is directly relevant in the literature on self-enhancement (Taylor & Brown, 1988; Colvin et al., 1995). For example, in a meta-analysis of the relationship between self-enhancement and narcissism, Grijalva and Zhang (2016) find that the reported relationship is stronger when researchers use residual scores than difference scores. This finding is consistent with the larger bias

for the residual score than difference scores across a wide range of parameter values noted earlier. But it may also speak to other disparate applications with structural similarities, such as when a marketing researcher has two measures of product quality and is interested in the role of one, net of the other. While this manuscript addresses overconfidence, the problem is a more general one.

## Recommendations

What solutions are available? No easy ones. But the absence of an easy solution does not provide cover to carry on as though there is no problem. I present four recommendations. Used by themselves or in concert, they have the potential to reduce the extent of problematic inferences.

## 1. Use Reliable, Valid Measures

Most importantly, this serves as a call to use reliable, valid measures. This is nothing new: no one touts the use of an unreliable or invalid measure. But given strong theoretical reasons to believe there is a confound without such measures, this reinforces their importance. This is particularly important given that performance measures are often moderately reliable at best. Anderson et al. (2012) report scale reliability of 0.66 for a narrowly-scoped geography quiz. For typical trivia quizzes, Krueger and Mueller (2002) report split-half correlations ranging from 0.17 to 0.56 and Burson et al. (2006) report split-half correlations ranging from 0.17 to 0.56 and Burson et al. (2006) report split-half

All else equal, averages across larger random samples of items will tend to have better coverage and reliability than smaller samples. Consider the infinite population of items that properly capture the target construct. A measure constructed using a larger random sample of items from that population will tend to have both better construct coverage (as there will be fewer untapped facets) and less measurement error than a measure constructed using a smaller random sample of items from that population. But properly specifying the population of items from which to sample is difficult. While it may be a valuable theoretical principle, simply relying on a larger sample of items may not be feasible in practice.

## 2. Account for Measurement Error

To provide an unbiased test when controlling for performance (i.e., when using the residual measure or multiple regression measure), it is useful to recall the conditions under which there is no bias.

The bias is eliminated if any of three conditions hold: (a)  $\alpha = 0$ , meaning there is no ambiguity and participants have no reason to regress their self-evaluations towards their prior beliefs, (b)  $\rho = 0$ , meaning latent beliefs are unrelated to latent ability, or (c) when  $\frac{\lambda^2}{\lambda^2 + \sigma_v^2} = 1$ , meaning there is no error and performance is at least partially related to ability.

This latter concern is analogous to a classic problem in which measurement error in one independent variable (performance) biases both its own coefficient and the coefficients of measures of correlated latent variables. Possible solutions to address this include structural equation models (e.g., Kline 2005) and errors-in-variables (e.g., Culpepper and Aguinis 2011). These approaches only help to the extent that  $\sigma_v^2$  only represents measurement error and not other constructs (as in the case of stable luck described above). The supplement presents additional details regarding both approaches. In a third empirical application (Parker et al.'s 2012 study of inappropriate confidence and financial planning), I reanalyze the data, first using a structural equation model and then using an adjustment for errors-invariables. In this application, the standard analysis finds an apparent correlation between financial planning and overconfidence. Accounting for measurement error largely attenuates or eliminates it.

## 3. Bound the Parameter Space

Rather than attempting to rule out this alternative explanation, researchers may instead relax the strength of their claims by acknowledging the conditions under which the null may hold. Given the ability to characterize the magnitude of the bias, one can qualitatively specify parameter configurations that could or could not account for the results. In some cases, no parameter values may be able to account for the set of correlations without overconfidence. In others, there may be feasible but implausible parameter configurations: while they are mathematically plausible, they may be ruled out based on theory.

Consider Figure 4 in application IIA, depicting the parameter values that are compatible with the null model. As can be seen by the crosscutting lightly shaded areas, while many parameters may be compatible with any one correlation, a much smaller set is compatible with the full set of correlations. In this case, a set of parameter values remains. In some cases, the remaining values may be implausible.

In other cases, one can rule out the alternative explanation altogether. First, if there is no relationship between ability and the candidate correlate, then although the measure is confounded, the confound has no bite to it. Note that no correlation between *performance* and the candidate correlate is not sufficient: such a lack of correspondence could merely indicate that performance is a poor measure of ability even if it is a reliable measure of something else.

Second, if the relationship between ability and the candidate correlate and the relationship between residualized overconfidence and the candidate correlate have opposite signs, the core bias described here could not account for the results. For example, as shown in Table A3 in the supplement, while the results of Lyons et al. (2021) regarding overconfidence and false news could be accounted for with the reported parameters, it would require an unlikely (and possibly implausible) sign on  $\beta$ , suggesting it is unlikely to account for the results. The bias may still be consequential: indeed, it may suggest the magnitude of the relationship between overconfidence and the outcome measure of interest is underestimated. This is not guaranteed to hold for difference scores due to the potential negative confound as depicted in Figure 3.

To examine this, one might consider whether overconfidence on one set of items predicts performance on a separate set of items, as in application I. If it does, this would suggest either that (a) the instrument has the problems described here, or (b) that overconfidence varies with ability (an inverted Dunning-Kruger effect; cf. Burson et al. 2006). But several notes are in order. First, the items used to calculate performance must be separate from those used to calculate overconfidence. Second, this analysis can provide evidence of a problem, but cannot provide evidence of absence of a problem. Third, one must avoid inadvertently accepting the null hypothesis, particularly given uncertainty around the estimates. In establishing such bounds, it is important to consider the uncertainty regarding one's estimate, not merely the point estimate itself. Finally, and as described above, these bounds are with respect to this particular strictly-specified null model. Other null models (e.g., one in which an unskilled-and-unaware effect holds but ability is the only correlate of behavior) may not be so readily ruled out.

#### 4. Use Alternative Measurement Approaches

Finally, one may opt to use a different measurement approach altogether. Multiple measures have been introduced which may be less susceptible to the focal confound. Direct measures of overclaiming (i.e., claiming one recognizes people, objects, or events that do not exist; Paulhus et al. 2003) may reduce the problems described here. One interpretation of such measures in terms of the current model is that ability is known to be constant and minimal (i.e., no one has the requisite knowledge to recognize things that do not exist.) Yet concerns remain regarding the role of inferences in the face of ambiguity. As a result of their ambiguity regarding individual items, people likely rely on their priors, which may again lead high-ability people to be more likely to overclaim than low-ability people in certain circumstances.

Similarly, Binnendyk and Pennycook (2024) and Lawson, Larrick, and Soll (2024) each propose measures of individual differences in overconfidence. These are based on estimated performance on a task for which performance is at or near chance and difficult to ascertain (Binnendyk and Pennycook 2024) or expectations of performance on specific tasks for which one has little contextual basis to expect superior performance. One interpretation of these measures in terms of the current model is that ability at these tasks ought to be unrelated to other correlates of interest. As with the Paulhus et al. (2003) measure, there is reason to be optimistic regarding these tasks and to potentially prefer them over the other methods described here. But they may not completely address the problems laid out here. To the extent that some people accurately believe themselves to be generally more successful than others at a wide variety of tasks, the same problem could persist. In such a case, the tasks could operate as poor measures of the underlying construct that participants rely upon (see "Is  $\lambda < 1$  Just a Form of Overconfidence?").

#### **Summary and Conclusion**

Measures of overconfidence on particular tasks vary across people. Yet widely-used measures of overconfidence that are used to study its correlates are confounded with the very thing they aim to rule out: ability. This is because performance measures are imperfect, so accounting for performance is insufficient to account for ability. Given ambiguity regarding performance, measures of confidence ought to regress towards prior beliefs about ability even when they are intended to be self-evaluations of task

performance. Because performance itself is an imperfect measure, the variance of self-evaluation attributable to ability is not fully partialed out when accounting for performance. The result is that both residual and difference overconfidence measures are confounded with ability. In an idealized null model, this bias can be quantified.

These confounds imply that it is possible to observe surprising results. In one reanalysis, I find overconfidence predicts later performance even after several intervening tasks. They also provide an alternative interpretation of findings from the literature: one need not posit a role for (or even the presence of) differences in overconfidence when considering relations with social status and social class. Instead, the results could be driven by test-taking ability alone. If researchers are willing to make strong assumptions regarding construct validity and estimate or assume reliability of each measure, it is possible to address these concerns through structural equation modeling or error-in-variables adjustments. However, these partial solutions are not an automatic panacea: complications may arise regarding construct validity and unstable estimates. Instead, design-based solutions (e.g., experimental manipulations or other measurement approaches) or accepting alternative interpretations of the results (i.e., plausible parameter configurations) may ultimately prove necessary. This work serves as a stark reminder to further improve our collective attempts to measure differences in overconfidence (whether stable or transient) and their true associations with traits, decisions, and behaviors.

#### References

- Alba JW, Hutchinson JW (2000) Knowledge calibration: What consumers know and what they think they know. *Journal of Consumer Research*, 27(2), 123-156.
- Anderson A, Baker F, Robinson DT (2017) Precautionary savings, retirement planning and misperceptions of financial literacy. *Journal of Financial Economics*, *126*(2), 383-398.
- Anderson C, Brion S, Moore DA, Kennedy JA (2012) A status-enhancement account of overconfidence. *Journal of Personality and Social Psychology*, 103(4), 718-735.
- Ames DR, Kammrath LK (2004) Mind-reading and metacognition: Narcissism, not actual competence, predicts self-estimated ability. *Journal of Nonverbal Behavior*, 28(3), 187-209.
- Avdeenko A, Bohne A, Frölich M (2019) Linking savings behavior, confidence and individual feedback: A field experiment in Ethiopia. *Journal of Economic Behavior & Organization*, 167.
- Barranti M, Carlson EN, Côté S (2017) How to test questions about similarity in personality and social psychology research: Description and empirical demonstration of response surface analysis. Social Psychological and Personality Science, 8(4), 465-475.
- Belmi P, Neale MA, Reiff D, Ulfe R (2020) The social advantage of miscalibrated individuals: The relationship between social class and overconfidence and its implications for class-based inequality. *Journal of Personality and Social Psychology*, *118*(2), 254-284.
- Binnendyk J, Pennycook G (2024) Individual differences in overconfidence: A new measurement approach. *Judgment and Decision Making*, forthcoming.
- Birnbaum MH, Mellers BA (1979) Stimulus recognition may mediate exposure effects. *Journal of Personality and Social Psychology*, *37*(3), 391-394.
- Blair G, Cooper J, Coppock A, Humphreys M, Sonnet L (2022) *estimatr: Fast Estimators for Design-Based Inference*. https://declaredesign.org/r/estimatr/, https://github.com/DeclareDesign/estimatr.
- Burson KA, Larrick RP, Klayman J (2006) Skilled or unskilled, but still unaware of it: how perceptions of difficulty drive miscalibration in relative comparisons. *Journal of Personality and Social Psychology*, 90(1), 60-77.

- Campbell WK, Goodie AS, Foster JD (2004) Narcissism, confidence, and risk attitude. *Journal of Behavioral Decision Making*, 17(4), 297-311.
- Carlson JP, Vincent LH, Hardesty DM, Bearden WO (2009) Objective and subjective knowledge relationships: A quantitative analysis of consumer research findings. *Journal of Consumer Research*, *35*(5), 864-876.
- Cheng JT, Anderson C, Tenney ER, Brion S, Moore DA, Logg JM (2021) The social transmission of overconfidence. *Journal of Experimental Psychology: General*, *150*(1), 157.

Cochran WG (1968) Errors of measurement in statistics. *Technometrics*, 10(4), 637-666.

- Cohen J, Cohen P, West SG, Aiken LS (2003) *Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences.* Routledge.
- Colvin CR, Block J, Funder DC (1995) Overly positive self-evaluations and personality: negative implications for mental health. *Journal of Personality and Social Psychology*, 68(6).
- Cronbach LJ, Furby L (1970) How we should measure "change": Or should we?. *Psychological Bulletin*, 74(1), 68.
- Culpepper SA, Aguinis H (2011) Using analysis of covariance (ANCOVA) with fallible covariates. *Psychological Methods*, *16*(2), 166-178.
- Edwards JR (1994) Regression analysis as an alternative to difference scores. *Journal of Management*, 20(3), 683-689.
- Edwards JR (1995) Alternatives to difference scores as dependent variables in the study of congruence in organizational research. *Organizational Behavior and Human Decision Processes*, 64(3), 307-324.
- Edwards JR, Parry ME (1993) On the use of polynomial regression equations as an alternative to difference scores in organizational research. *Academy of Management journal*, *36*(6), 1577-1613.
- Ehrlinger J, Dunning D (2003) How chronic self-views influence (and potentially mislead) estimates of performance. *Journal of Personality and Social Psychology*, 84(1), 5-17.

Erev I, Wallsten TS, Budescu DV (1994) Simultaneous over-and underconfidence: The role of error in

judgment processes. Psychological Review, 101(3), 519.

- Fagerlin A, Zikmund-Fisher BJ, Ubel PA, Jankovic A, Derry HA, Smith DM (2007) Measuring numeracy without a math test: development of the Subjective Numeracy Scale. *Medical Decision Making*, 27(5), 672-680.
- Feld J, Sauermann J, De Grip A (2017) Estimating the relationship between skill and overconfidence. *Journal of Behavioral and Experimental Economics*, 68, 18-24.
- Fiedler K (2021) Suppressor effects in self-enhancement research: A critical comment on condition-based regression analysis. *Journal of Personality and Social Psychology*, *121*(4), 792-795.
- Fleming SM, Lau HC (2014) How to measure metacognition. *Frontiers in Human Neuroscience*, 8(443), 1-9.
- Fuller WA (1987) Measurement Error Models. John Wiley & Sons.
- Gigerenzer G, Hoffrage U, Kleinbölting H (1991) Probabilistic mental models: a Brunswikian theory of confidence. *Psychological Review*, *98*(4), 506.
- Gillen B, Snowberg E, Yariv L (2019) Experimenting with measurement error: Techniques with applications to the Caltech cohort study. *Journal of Political Economy*, *127*(4), 1826-1863.
- Griffin D, Murray S, Gonzalez R (1999) Difference score correlations in relationship research: A conceptual primer. *Personal Relationships*, 6(4), 505-518.
- Grijalva E, Zhang L (2016) Narcissism and self-insight: A review and meta-analysis of narcissists' selfenhancement tendencies. *Personality and Social Psychology Bulletin*, 42(1), 3-24.
- Grinblatt M, Keloharju M (2009) Sensation seeking, overconfidence, and trading activity. *The Journal of Finance*, 64(2), 549-578.

Healy PJ, Moore DA (2007) Bayesian overconfidence. Available at SSRN 1001820.

Humberg S, Dufner M, Schönbrodt FD, Geukes K, Hutteman R, Küfner AC, van Zalk MH, Denissen JJA, Nestler S, Back MD (2019) Is accurate, positive, or inflated self-perception most advantageous for psychological adjustment? A competitive test of key hypotheses. *Journal of Personality and Social Psychology*, 116(5), 835-859.

- Humberg S, Dufner M, Schönbrodt FD, Geukes K, Hutteman R, van Zalk MH, Denissen JJA, Nestler S,
   Back MD (2018a) Enhanced versus simply positive: A new condition-based regression analysis
   to disentangle effects of self-enhancement from effects of positivity of self-view. *Journal of Personality and Social Psychology*, 114(2), 303-322.
- Humberg S, Dufner M, Schönbrodt FD, Geukes K, Hutteman R, van Zalk MH, Denissen JJA, Nestler S,
  Back MD (2018b) Why Condition-Based Regression Analysis (CRA) is Indeed a Valid Test of
  Self-Enhancement Effects: A Response to Krueger et al. *Collabra: Psychology*, 4(1), 26.
- Humberg S, Dufner M, Schönbrodt FD, Geukes K, Hutteman R, van Zalk MH, ..., Back MD (2022) The true role that suppressor effects play in condition-based regression analysis: None. A reply to Fiedler (2021). *Journal of Personality and Social Psychology*, *123*(4), 884-888.
- Humberg S, Nestler S, Back MD (2019) Response surface analysis in personality and social psychology:
   Checklist and clarifications for the case of congruence hypotheses. *Social Psychological and Personality Science*, 10(3), 409-419.
- Jansen RA, Rafferty AN, Griffiths TL (2021) A rational model of the Dunning–Kruger effect supports insensitivity to evidence in low performers. *Nature Human Behaviour*, *5*(6), 756-763.
- John OP, Robins RW (1994) Accuracy and bias in self-perception: Individual differences in selfenhancement and the role of narcissism. *Journal of Personality and Social Psychology*, 66(1).
- Johns G (1981) Difference score measures of organizational behavior variables: A critique. Organizational Behavior and Human Performance, 27(3), 443-463.
- Juslin P (1994) The overconfidence phenomenon as a consequence of informal experimenter-guided selection of almanac items. *Organizational Behavior and Human Decision Processes*, 57(2), 226.
- Kahneman D (1965) Control of spurious association and the reliability of the controlled variable. *Psychological Bulletin*, 64(5), 326-329.
- Kahneman D (2011) Thinking, Fast and Slow. Macmillan.
- Ke D (2021) Who wears the pants? Gender identity norms and intrahousehold financial decision-making. *The Journal of Finance*, 76(3), 1389-1425.

Klayman J, Soll JB, Gonzalez-Vallejo C, Barlas S (1999) Overconfidence: It depends on how, what, and whom you ask. *Organizational Behavior and Human Decision Processes*, 79(3), 216.

Kline RB (2005) Principles and Practice of Structural Equation Modeling. 2nd ed. New York: Guilford.

- Kramer MM (2016) Financial literacy, confidence and financial advice seeking. *Journal of Economic Behavior & Organization*, 131, 198-217.
- Krueger JI, Heck PR, Asendorpf JB (2017) Self-enhancement: Conceptualization and assessment. *Collabra: Psychology*, *3*(1), 28.
- Krueger J, Mueller RA (2002) Unskilled, unaware, or both? The better-than-average heuristic and statistical regression predict errors in estimates of own performance. *Journal of Personality and Social Psychology*, 82(2), 180-188.
- Kruger J, Dunning D (1999) Unskilled and unaware of it: how difficulties in recognizing one's own incompetence lead to inflated self-assessments. *Journal of Personality and Social Psychology*, 77(6), 1121-1134.
- Landier A, Thesmar D (2008) Financial contracting with optimistic entrepreneurs. *The Review of Financial Studies*, 22(1), 117-150.
- Larkin I, Leider S (2012) Incentive schemes, sorting, and behavioral biases of employees: Experimental evidence. *American Economic Journal: Microeconomics*, *4*(2), 184-214.
- Lawson MA, Larrick RP, Soll JB (2023) Individual differences in overconfidence and their psychological bases. <u>http://dx.doi.org/10.2139/ssrn.4558486</u>
- Lawson MA, Larrick RP, Soll JB (2024) The Individual Overconfidence Test: A parsimonious measure of trait-level overconfidence. <u>http://dx.doi.org/10.2139/ssrn.4893481</u>
- Li S, Hale R, Moore DA (2024) Is overconfidence an individual difference? <u>https://doi.org/10.31234/osf.io/xprf7</u>
- Lord FM (1956) The measurement of growth. ETS Research Bulletin Series, 1956(1), i-22.
- Lord FM (1958) Further problems in the measurement of growth. *Educational and Psychological Measurement*, 18(3), 437-451.

- Lord FM (1960) Large-sample covariance analysis when the control variable is fallible. *Journal of the American Statistical Association*, 55(290), 307-321.
- Lusardi A, Mitchell OS (2017) How ordinary consumers make complex economic decisions: Financial literacy and retirement readiness. *Quarterly Journal of Finance*, 7(3), 1750008.
- Lyons BA, Montgomery JM, Guess AM, Nyhan B, Reifler J (2021) Overconfidence in news judgments is associated with false news susceptibility. *Proceedings of the National Academy of Sciences*, 118(23), e2019527118.
- MacIntyre PD, Noels KA, Clément R (1997) Biases in self-ratings of second language proficiency: The role of language anxiety. *Language Learning*, 47(2), 265-287.
- McNemar Q (1958) On growth measurement. Educational and Psychological Measurement, 18(1), 47.
- Moore DA, Healy PJ (2008) The trouble with overconfidence. *Psychological Review*, 115(2).
- Moore DA, Dev AS (2017) Individual differences in overconfidence. *Encyclopedia of Personality and Individual Differences. Springer. Retrieved from http://osf. io/hzk6q.*
- Moore DA, Schatz D (2017) The three faces of overconfidence. *Social and Personality Psychology Compass*, *11*(8), e12331.
- Moore D, Swift SA (2011) The three faces of overconfidence in organizations. In *Social Psychology and Organizations*, 179-216.
- Moorman C, Diehl K, Brinberg D, Kidwell B (2004). Subjective knowledge, search locations, and consumer choice. *Journal of Consumer Research*, *31*(3), 673-680.
- Murphy SC, von Hippel W, Dubbs SL, Angilletta Jr MJ, Wilson RS, Trivers R, Barlow FK (2015) The role of overconfidence in romantic desirability and competition. *Personality and Social Psychology Bulletin*, *41*(8), 1036-1052.
- Parker AM, Bruin De Bruin W, Yoong J, Willis R (2012) Inappropriate confidence and retirement planning: Four studies with a national sample. *Journal of Behavioral Decision Making*, 25(4), 382-389.

Parker AM, Stone ER (2014) Identifying the effects of unjustified confidence versus overconfidence:

Lessons learned from two analytic methods. *Journal of Behavioral Decision Making*, 27(2), 134-145.

- Paulhus DL, Harms PD, Bruce MN, Lysy DC (2003) The over-claiming technique: measuring selfenhancement independent of ability. *Journal of Personality and Social Psychology*, 84(4), 890.
- Rogosa D, Brandt D, Zimowski M (1982) A growth curve approach to the measurement of change. *Psychological bulletin*, 92(3), 726.
- Stankov L, Crawford JD (1996) Confidence judgments in studies of individual differences. *Personality* and Individual Differences, 21(6), 971-986.
- Taylor SE, Brown JD (1988) Illusion and well-being: A social psychological perspective on mental health. *Psychological Bulletin*, *103*(2), 193-210.
- Thomson GH (1924) A formula to correct for the effect of errors of measurement on the correlation of initial values with gains. *Journal of Experimental Psychology*, 7(4), 321.
- Wall TD, Payne R (1973) Are deficiency scores deficient?. Journal of Applied Psychology, 58(3).
- Westfall J, Yarkoni T (2016) Statistically controlling for confounding constructs is harder than you think. *PloS One*, *11*(3), e0152719.
- Zuckerman M, Knee CR (1996) The relation between overly positive self-evaluation and adjustment: a comment on Colvin, Block, and Funder (1995). *Journal of Personality and Social Psychology*, 70(6), 1250-1.