

# Improving the reliability of cognitive task measures: A narrative review

Samuel Zorowitz<sup>1,\*</sup>, Yael Niv<sup>1,2</sup>

<sup>1</sup>Princeton Neuroscience Institute, Princeton University, USA

<sup>2</sup>Department of Psychology, Princeton University, USA

\*Corresponding author (zorowitz@princeton.edu)

## Abstract

Cognitive tasks are capable of providing researchers with crucial insights into the relationship between cognitive processing and psychiatric phenomena. However, many recent studies have found that task measures exhibit poor reliability, which hampers their usefulness for individual-differences research. Here we provide a narrative review of approaches to improve the reliability of cognitive task measures. Specifically, we introduce a taxonomy of experiment design and analysis strategies for improving task reliability. Where appropriate, we highlight studies that are exemplary for improving the reliability of specific task measures. We hope that this article can serve as a helpful guide for experimenters who wish to design a new task, or improve an existing one, to achieve sufficient reliability for use in individual-differences research.

## 1 Introduction

Cognitive tasks hold great promise for biological psychiatry. When properly designed, such tasks are capable of isolating and measuring specific cognitive processes. Individual differences in performance on cognitive tasks can therefore provide researchers with crucial insights into the cognitive processes underlying psychiatric phenomena. Elsewhere in psychology, cognitive tasks have been useful in predicting important outcomes such as academic achievement [1] and cognitive decline [2]. Cognitive tasks, then, have the potential to be invaluable tools for refining our understanding of psychiatric symptoms and syndromes. For a cognitive task to be useful in this regard, however, it must possess sufficient measurement properties.

We define a cognitive task as any experimental paradigm that measures behavior in order to make inferences about one or more cognitive processes (e.g., Stroop task, delay discounting task, reversal-learning task). Cognitive-task measures of behavioral performance can be descriptive (e.g., proportion correct responses, average response time)

or model-based (e.g., drift rate in evidence accumulation models). The psychometric quality of a measure can be summarized by three key properties: *discriminatory power*, *validity* and *reliability* [3]. The discriminatory power of a task measure describes its ability to measure variability in participants’ performance. This is a necessary property of tasks used to study individual differences; where there is no variation in performance, there are no individual differences to study. The validity of a task measure concerns whether it actually measures what it intends to measure. Finally, the reliability of a task measure characterizes the degree to which it consistently measures some feature of participants. That is, a task measure is reliable if, assuming participants have not changed, it produces the same scores, or the same ordering of scores, for participants within a single testing session or across multiple testing sessions. This review focuses on task-measure reliability.

## The formal definition of reliability

In classical test theory [4], the variance in observed scores on a task measure  $x$  is the sum of the true score variance  $\sigma_T^2$ , reflecting real individual differences in the latent construct of interest, and measurement error  $\sigma_E^2$ , i.e.,  $\sigma_x^2 = \sigma_T^2 + \sigma_E^2$ . The reliability of a measure is defined as the proportion of variance attributable to the true score variance relative to total variance:  $\rho_{xx'} = \frac{\sigma_T^2}{\sigma_T^2 + \sigma_E^2}$ . Thus, reliability quantifies the magnitude of individual differences relative to the noisiness of a task measure; the larger the reliability of a task measure, the more it reflects true individual differences rather than noise. Reliability is therefore a prerequisite for validity: an unreliable task measure reflects measurement error and not the construct of interest. If this were not reason enough to care about reliability, the observed correlation between two measures (e.g., task performance and self-reported symptom score) is bounded by their individual reliabilities [5]:

$$\rho_{xy} = \rho_{tt} \sqrt{\rho_{xx'} \cdot \rho_{yy'}} \quad (1)$$

where  $\rho_{xx'}$  and  $\rho_{yy'}$  are the reliabilities of two measures,  $x$  and  $y$ ;  $\rho_{tt}$  is their true latent correlation; and  $\rho_{xy}$  is their observed correlation. As all reliabilities are  $< 1$ , the reliability of a measure places an upper bound on the maximum observable correlation between itself and a second measure (Figure 1). As an important corollary, as measure reliability decreases, the number of participants required to reliably detect a correlation between two measures increases [6]. Thus, poor reliability hampers our ability to investigate associations between cognitive processes, as measured by task performance, and other variables of interest.

To further complicate matters, the reliability of a task measure is not absolute – it reflects interactions between the task design, the participants, and the context in which the task is administered. Indeed, task reliability can vary as a function of experiment parameters (specific stimulus set, number of trials, time limits [8, 9]); sample populations (healthy adults, children, psychiatric patients [9, 10]); testing locations (in clinic, on-line); response modality (desktop, smartphone, virtual reality [11, 12]); scoring method (component scores, difference scores); and estimation method [13–15]. For example, a cognitive task originally designed for use with an adult population may prove too difficult for children. Their task performance may drop to chance level, thereby minimizing between-participant variance and, as a consequence, task reliability. As a second exam-

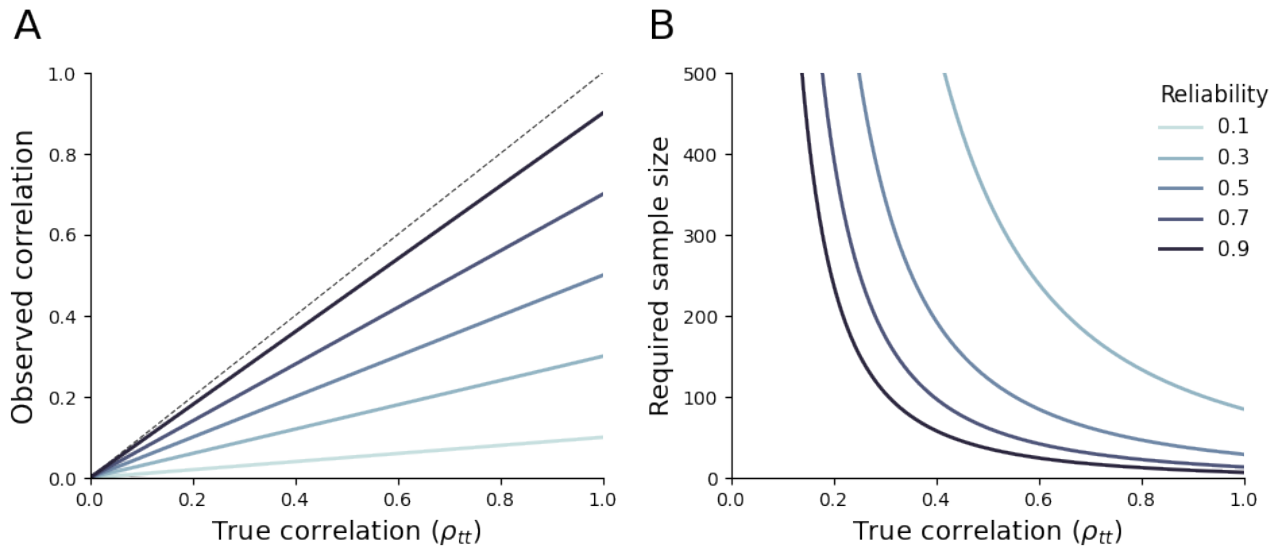


Figure 1: The relationship between measure reliability, observed correlations, and statistical power. (A) The maximum expected observed correlation between two measures as a function of their true (latent) correlation and reliability. As the reliability of two measures decreases, so too does their observed correlation. (B) Required sample size for 80% power to detect true correlations between two measures given their reliability. As the reliability of two measures decreases, the number of participants required to detect an association increases. Even with a large sample size of  $N = 500$ , two measures with moderate reliability ( $\rho_{xx'} = \rho_{yy'} = 0.5$ ) will only reliably detect true correlations above 0.3, which are likely high for individual-differences cognitive research [7].

ple, participants completing an experiment online from their homes may experience more distraction than if they participated in the lab. This may increase measurement error, leading to a concomitant decrease in reliability. Experimenters therefore cannot assume the reliability of a task measure is constant. At the very least, researchers should evaluate reliability after having made changes to a task or scoring procedure, or when administering the task to new sample populations or in new testing contexts. Ideally, researchers would investigate and report the reliability of task measures as part of any individual differences research.

Although verifying the reliability of cognitive task measures is paramount to individual differences research, the reliability of task measures is seldom reported [6, 16]. When they are reported, task measures frequently exhibit lower reliability than what is conventionally considered the minimum acceptable level for individual-differences research ( $\rho_{xx'}$  on the order of 0.7 – 0.8) and the reliability regularly achieved by self-report measures. Indeed, many studies have now found that task measures exhibit moderate-to-low reliability [17–22].

One possible explanation for this finding is the so-called “reliability paradox” of cognitive tasks [17], which states that the often lackluster reliability of tasks is a result of a mismatch in goals between experimental and individual-differences psychological research. In experimental psychology, the goal is often to demonstrate the existence of a behavioral effect. One means of increasing the power to detect an effect is to minimize between-participants variance. This is the exact opposite of what is desirable for individual differences research, where between-participants variance is essential to achieving

reliable task measures. For example, the Stroop effect is one of the most robust effects in experimental psychology; virtually everyone shows a Stroop effect [23]. However, in part due to this fact, between-participants variance on the Stroop effect is often limited [24]. Thus, the tendency in biological psychiatry to adopt the most prominent tasks in experimental psychology—the ones that most reliably demonstrate a behavioral effect—may actually hamstring efforts to study individual differences.

Regardless, we do not believe that task measures are inherently less reliable than self-report measures, or that pessimism about task-based individual-differences research is warranted. It is possible to (re)design tasks to achieve good reliability, even to the high levels dictated by conventional standards [25–28]. The purpose of the current article is to provide a narrative review of approaches to improve task-measure reliability. Specifically, we introduce a taxonomy of strategies for improving the reliability of cognitive-task measures through experiment design and analysis. Where appropriate, we highlight studies that are exemplary for improving the reliability of specific task measures. For the interested reader, we review methods for calculating the reliability of task measures in the supplementary materials, as these topics have been discussed at length elsewhere [6, 14, 16]. We hope that this article can serve as a helpful guide for experimenters designing a new task, improving an existing task, or refining their scoring methods to achieve sufficient reliability for use in individual-differences research.

## 2 Improving task reliability

As defined above, the reliability of a task measure is the proportion of variance attributable to between-participant differences relative to measurement error. Thus, the two major strategies for improving the reliability of a measure are to increase between-participant variability or decrease measurement error. In what follows, we discuss approaches for accomplishing each objective in turn. Where appropriate, we highlight studies that are exemplary for improving the reliability of a task measure by implementing a particular strategy.

### 2.1 Increasing between-participant variance

#### 2.1.1 Ceiling & floor effects

By definition, the reliability of a task measure is zero when there is no variability across participants. Thus, range restriction of task measures via ceiling or floor effects is a serious obstacle to reliability. Siegelman and colleagues [29] noted the consequences of floor effects on reliability in the context of statistical learning tasks. In such tasks, participants must learn to identify subtle patterns in the transition probabilities underlying a continuous sequence of stimuli. In reanalyzing archival datasets, Siegelman and colleagues found a majority of participants were at chance-level performance in discriminating between legitimate and foil sequence patterns; consequently, the reliability of conventional proportion correct measures suffered. In response, the authors designed a new statistical learning task involving stimulus sequences that ranged more widely in their difficulty to learn. Only a minority of participants showed chance-level performance on this new task

and, as such, the reliability of proportion correct scores improved (from  $\rho = 0.75$  to  $0.88$ ). Similarly, in developing an abbreviated working-memory task, Oswald and colleagues [30] found that they could remove the easiest trials—those with ceiling level performance—with virtually no change to task reliability. This is because those trials are incapable of differentiating ability across participants and therefore cannot contribute significantly to the reliability of the task.

Researchers administering a task to a new population should be especially wary of range-restriction effects. Cognitive tasks calibrated for one group of participants may not be adequately sensitive for others due to being too easy or difficult for a different group. For example, Arnon and colleagues [10] found that statistical learning tasks developed for adults were too difficult for young children and therefore yielded unreliable discrimination scores in that population. Similarly, Kyllonen and colleagues [31] developed a battery of fluid-reasoning measures for highly educated adults after observing ceiling effects in performance when using preexisting fluid-reasoning tasks in this population.

### 2.1.2 Repeatability & practice effects

A related issue for task reliability is practice effects, where participants' performance on a task improves with repeated administrations. Practice effects are relatively common for cognitive tasks [32, 33]. They might occur due to the attenuation of task-irrelevant nuisance factors (e.g., performance anxiety) and/or the learning of task-specific knowledge or strategies. Practice effects are not inherently an issue for reliability — especially if an experimenter is only interested in the consistency, but not the absolute agreement, of participants' performance over time — but they can become a pernicious issue if they are exhibited differentially across participants or if they are severe enough to induce ceiling effects. For example, Paredes and colleagues [34] observed large practice effects on the Pavlovian go/no-go task for short retest intervals (3 days, 14 days), which resulted in poor estimates of test-retest reliability. In developmental and lifespan studies, practice effects are potentially complicated by their interaction with age [35, 36]; that is, practice-induced ceiling effects may present in some age groups but not others.

One strategy for minimizing practice effects is simply to increase the time interval between task administrations. The more time that elapses between sessions, the greater the probability that participants will have forgotten task-specific knowledge or strategies [32, 33]. Of course, this solution may not always be possible or desirable, especially if a researcher is only able to or specifically interested in studying a behavior over a short time period. Moreover, some forms of learning do not easily dissipate with time [37].

A second strategy is to use a combination of clear instructions and practice trials to help participants reach stable performance from the start of an experiment (see the “Improving experiment designs” section below). Another strategy is to design tasks so as to prevent or discourage the formation of task-specific strategies. For example, McLean and colleagues [25] investigated the repeatability of the beads task. In the beads task, participants are presented with two jars containing beads of two colors in equal but opposite ratios. In each trial of the task, a predetermined sequence of beads is drawn from one jar. Participants must decide which jar beads were being drawn from or request to see more beads. In a typical version of the task, the same sequence of bead draws is

used across all trials. McLean and colleagues found that participants were aware that the sequence repeated across trials and, as a consequence, became more erratic in their decision to witness more bead draws with additional trials. In response, the authors developed a new version of the task that included distractor sequences of bead draws. This new design was effective in preventing participants from becoming aware of the target sequence, which in turn resulted in more consistent responding, which improved the reliability of participants’ information seeking scores (from  $\rho = 0.62$  to  $0.84$ ).

### 2.1.3 Enhancing experimental manipulations

The preceding sections described potential threats to between-participants variability, but not approaches to improve it. A primary strategy for increasing between-participants variability is to enhance the experimental manipulation. Amplifying the strength of an experimental manipulation (e.g., making a task more challenging, increasing the potency of affect induction) typically increases the range of participants’ responses to it. For example, Kucina and colleagues [26] investigated the reliability of cognitive conflict effects (as measured by response time) in new versions of several standard cognitive-control tasks (e.g., Stroop, Flanker, Simon) that amplified cognitive interference via two task design features. First, they combined multiple sources of cognitive interference in the same task, for example by combining the Stroop and Simon effects to create a “Stroopon” task. Second, for a subset of trials, they required participants to make multiple responses based on both relevant and irrelevant stimuli attributes. Compared to previous versions of these tasks, these manipulations had the effect of increasing task demands, which resulted in greater between-participants variance and, consequently, required hundreds fewer trials to achieve a reliability of  $\rho = 0.8$ . (See also Snijder and colleagues [27] for a similar redesign of classic cognitive-control tasks that improved reliability in part by increasing proactive control demands.)

A related strategy is to calibrate the difficulty of the task to the average ability of the population of interest. For example, consider a task trial with two response options. Given the Bernoulli distribution, the variance in responses on this trial will be maximal when the probability of choosing either response is equal. Thus, aggregating across many trials, between-participants variance is maximized — and task reliability is improved — when the difficulty of all items is matched to the average ability of the sample [38] (or slightly higher if participants can guess the correct response [39, 40]). Of course, this design principle is only helpful to the degree that a researcher knows the average ability level of their participants. If this is unknown or poorly characterized, then it is instead desirable to design a task with trials spanning a range of difficulty levels.

### 2.1.4 Sample population

A final strategy for increasing between-participant variance is to simply recruit more diverse samples. While convenient, undergraduate students from a single university are likely to be relatively homogeneous in their cognitive profiles. It may be worthwhile instead to recruit participants from the community or from an online labor market (e.g., Amazon Mechanical Turk, Prolific Academic, CloudResearch Panels). With regard to the latter, because online participants typically complete experiments from their homes or

other poorly-controlled environments, they are more likely to be distracted or to multi-task during an experiment [41]. Thus, when recruiting online samples, experimenters should take special care to ensure that an increase in between-participants variance is not offset by a concomitant increase in measurement noise. Separately, online participants may be more familiar with particular experimental paradigms due to previous exposure [42], which may attenuate between-participants performance variability for the reasons previously mentioned (e.g., practice effects). Thus, researchers running experiments online may want to alter task paradigms so that they appear less similar to preexisting versions and/or limit the recruitment of highly-experienced participants [43].

## 2.2 Decreasing measurement noise

### 2.2.1 Increasing trial numbers

Perhaps the most straightforward approach to decreasing measurement error, and thereby increasing reliability, is to increase the number of task trials. The relationship between reliability and the number of trials defined as:

$$\rho = \frac{\sigma_T^2}{\sigma_T^2 + \frac{\sigma_E^2}{n}} \quad (2)$$

where  $\sigma_T^2$  is the true between-participants variance,  $\sigma_E^2$  is measurement error (i.e. trial-level variance), and  $n$  is the number of trials. In practice, this relationship often holds [8, 9] though with some exceptions [44, 45]. Notably, increasing the number of task trials only benefits reliability if measurement error is random. If increasing task length results in participant fatigue or boredom, measurement noise may systematically increase and reliability will suffer. Increasing the number of trials may be impractical for other reasons, and because of diminishing marginal improvements for reliability, achieving a desired level of reliability through this means alone may require prohibitively long experiments.

### 2.2.2 Improving experiment designs

Measurement error can be reduced through improving the design of experiments, which can be accomplished in many ways. The reliability of a task measure can be improved by including in an experiment only the most discriminating stimuli. For example, in the context of an emotion recognition task, stimuli of good discriminability would be those where participants with good emotion-recognition ability consistently correctly identify the displayed emotion while participants with poor ability consistently incorrectly identify the displayed emotion [46]. In contrast, stimuli with poor discriminability — those for which performance between high- and low-ability participants is indistinguishable — will lead to more measurement noise and decreased reliability [47]. In experiments where stimuli are intended to be unique and distinguishable, improving both the linguistic and visual distinctness of stimuli may prevent participant confusion and therefore aid reliability [48].

Other design features of an experiment that are specific to the task mechanics may affect reliability. Consider, for example, the dot-probe task, in which participants must disengage attention from a distracting image on one part of the screen in order to identify and respond to the orientation of a pair of dots elsewhere on the screen. Dot-bottom trials,



where a participant must disengage from a distracting stimulus located at the top of the screen and saccade to the bottom of the screen, are more reliable than dot-top trials (e.g., dot-bottom:  $\rho = 0.33$ ; dot-top:  $\rho = 0.07$ ) [44, 49]. This has been explained by suggesting that because participants’ gazes are biased towards the top of the screen, saccading away from the top requires a stronger level of disengagement. Dot-bottom trials may therefore be better measures of attentional bias.

It is worth stressing that clear instructions are essential for task reliability. When participants are unsure of what they are intended to do in an experiment, their behavior is likely to be more variable across time (as their understanding of the task evolves) and across participants (due to different interpretation of instructions), thereby diminishing reliability. Clear instructions thus help to ensure that participants show stable behavior from the start of the experiment. (Clear instructions may also work to ensure the validity of an experiment by discouraging participants from using strategies not of interest to the experimenter.) In their “10 simple rules” paper for designing cognitive experiments [50], Barbosa and colleagues provide practical suggestions for writing task instructions. When adapting a task for a new population, experimenters should ensure that the instructions are still appropriate. Task instructions that are comprehensible to healthy adult participants may not be suitable for other populations like children [51].

Another strategy to reduce measurement error is to make use of a practice phase. Practice trials can help to minimize the effects of nuisance factors such as performance anxiety or unfamiliarity with the response modality, and give participants an opportunity to make sure they understand the task instructions. Thus, practice trials can help participants reach a “steady state” of responding, thereby reducing the noisiness of their responses across task trials and increasing reliability [52]. Practice can take the form of a standalone practice block or by designating as such and discarding (or modeling separately) the first few trials of an experiment (e.g., [25]), though the latter may not fully allow participants to explore response options during the practice.

Yet another strategy to diminish measurement error is to “gamify” an experiment. Incorporating (video) game design elements into cognitive tasks can increase participant engagement and motivation [53], countering the would-be effects of boredom and fatigue on task reliability. For example, Kucina and colleagues [26] cite task gamification as an important factor that contributed to the reliability of their cognitive-control tasks. Similarly, Verdejo and colleagues [22] partially attribute the adequate reliability of their impulsivity task battery to gamified task design (range:  $\rho = 0.52 - 0.71$ ).

### 2.2.3 Reducing parameter estimation noise

When parameters from cognitive models are used as indices of participants’ task performance, another means to improve measure reliability is to decrease estimation noise. The estimation noise of a parameter given an experiment and model can be quantified through simulation studies [54, 55]. Here, an experimenter generates artificial data for the experiment using representative model parameters and then attempts to recover the model parameters by fitting the model to the simulated data. Estimation noise is the inverse of the (relative or absolute) agreement between the true and recovered parameters. Alterations to experimental design can improve parameter recovery and estimation noise,



and multiple frameworks have been proposed for testing and improving experimental designs to aid parameter recovery [56, 57]. Parameter recovery can also be affected by the model estimation method [58, 59]. In particular, the partial pooling properties of hierarchical Bayesian models can be especially beneficial for improving parameter recovery and decreasing estimation noise [60].

A related approach is to use adaptive experimental designs [61], where the trials of an experiment are designed in real time so as to present each participant with stimuli or trial types that are matched to their particular response patterns or ability levels. Though undoubtedly a more complex experimental design, adaptive experiments have the advantage of selecting the most informative trials for resolving the ability or preference level of a participant (e.g., as measured by a cognitive model). Adaptive designs have been successfully used in cognitive research, for example to study working memory [62] and delay discounting [63]. For a detailed discussion of adaptive-design experiments, see [64].

Parameter estimation may be further improved by leveraging additional information. For example, latent variables may be more accurately measured through the inclusion of demographic variables or other covariates that, if associated to model parameters, can aid in resolving parameter estimates [65]. An extension of this idea is to utilize joint modeling of dependent variables; that is, to design models where multiple observed trial-level variables are predicted simultaneously. For example, the joint modeling of choice and response time has been found to improve the precision and reliability of estimated parameters in cognitive ability testing [66] and reinforcement learning [67, 68]. It is also possible to incorporate physiological and/or neural correlates of behavior, such as skin conductance response, fMRI BOLD signal, and EEG [69].

## 2.3 Difference scores

Difference scores deserve special treatment in the context of reliability. Difference scores subtract a measure of a participant’s performance in one condition from their performance in another. An example is the Stroop interference effect, calculated as the average reaction-time difference between congruent and incongruent trials. Difference scores are commonly used because they allow experimenters to isolate particular cognitive processes (e.g., processing cognitive conflict) while controlling for other sources of variance (e.g., perceptual processing, motor ability) through the subtraction of conditions that share that variance. The challenge is that the reliability of a difference score is a function of the reliability of each of its components *and* the correlation between the components:

$$\rho_{dd'} = \frac{\sigma_x^2 \rho_{xx'} + \sigma_y^2 \rho_{yy'} - 2\rho_{xy} \sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2 - 2\rho_{xy} \sigma_x \sigma_y} \quad (3)$$

where  $\sigma_x^2$  and  $\sigma_y^2$  are the variances of task measures  $x$  and  $y$ ,  $\rho_{xx'}$  and  $\rho_{yy'}$  are the reliabilities of task measures  $x$  and  $y$ , and  $\rho_{xy}$  is the correlation between task measures  $x$  and  $y$  [70]. When the variances of the two measures are equal, this reduces to:

$$\rho_{dd'} = \frac{\rho_{xx'} + \rho_{yy'} - 2\rho_{xy}}{2 - 2\rho_{xy}} \quad (4)$$

From the above equation, one can see that the reliability of a difference score measure is diminished to the extent that its components are correlated. Two measures derived from

the same task will often be correlated due to shared domain-general cognitive processes. Thus, it will often (if not always) be the case that difference scores derived from task measures will be less reliable than the average reliability of their components.

### 2.3.1 Enhance & purify task measures

The equations above suggest three steps experimenters can take to improve the reliability of difference scores: (1) improve the reliability of the component measures, (2) increase the relative difference between the variances of two task measures, and (3) minimize the correlation between the task measures. The first strategy has been our focus so far. The second approach deserves further comment. As discussed elsewhere [70], the reliability of a difference score measure increases as the difference (or the ratio) between the variances of the component measures increases. Intuitively, this is because as the variance of one (but not the other) component grows, so too does the proportion of unshared reliable variance. Figure 2A shows that even when the correlation of two measures is large, it is possible to achieve acceptable reliability insofar as the ratio of the variances (i.e.  $\sigma_x/\sigma_y$ ) is sufficiently different from 1. (It is also worth noting that component variances essentially function as weights in determining the overall reliability, such that the difference-score reliability reflects more the reliability of the component measure with the larger variance; Figure 2B.) This speaks to the advantage of increasing the between-participants variability of a measure of performance in one experiment condition without increasing performance variability in a second condition. This may explain how Kuchina and colleagues improved the reliability of their difference score measures after making only the incongruent trials more difficult [26].

The third approach is to purify task measures; that is, to decorrelate the components of a difference score by reducing or eliminating their shared variance. Rey-Mermet and colleagues [28] provide an interesting example in the context of executive control. In typical executive control tasks, response times on congruent and incongruent trials are highly correlated, reflecting shared variance from conflict-irrelevant processes including baseline processing speed (e.g., perceptual processing, motor speed) and performance strategies (e.g., individual differences in speed-accuracy preferences; [71]). Rey-Mermet and colleagues designed a number of “response-deadline” executive control tasks where participants had a limited amount of time to respond during a trial. The duration of the response deadline was calibrated for each participant individually, such that participants achieved a fixed accuracy level in blocks of neutral trials, which was then used as the deadline for both congruent and incongruent trials. The calibration procedure controls for individual differences in processing speed that contribute to performance in both congruent and incongruent trials. It also controls for individual differences in strategy, as regardless of whether a participant was biased towards speed or accuracy, inefficient executive control would result in lower accuracy. With this calibration procedure, the researchers found that the reliability of an accuracy difference score (incongruent minus congruent, which ranged from  $\rho = 0.58$  to  $\rho = 0.91$ ) was as good or better than what had previously been reported for executive control tasks. Thus, controlling for shared variance across measures—that is, purifying measures—can help to improve task reliability.

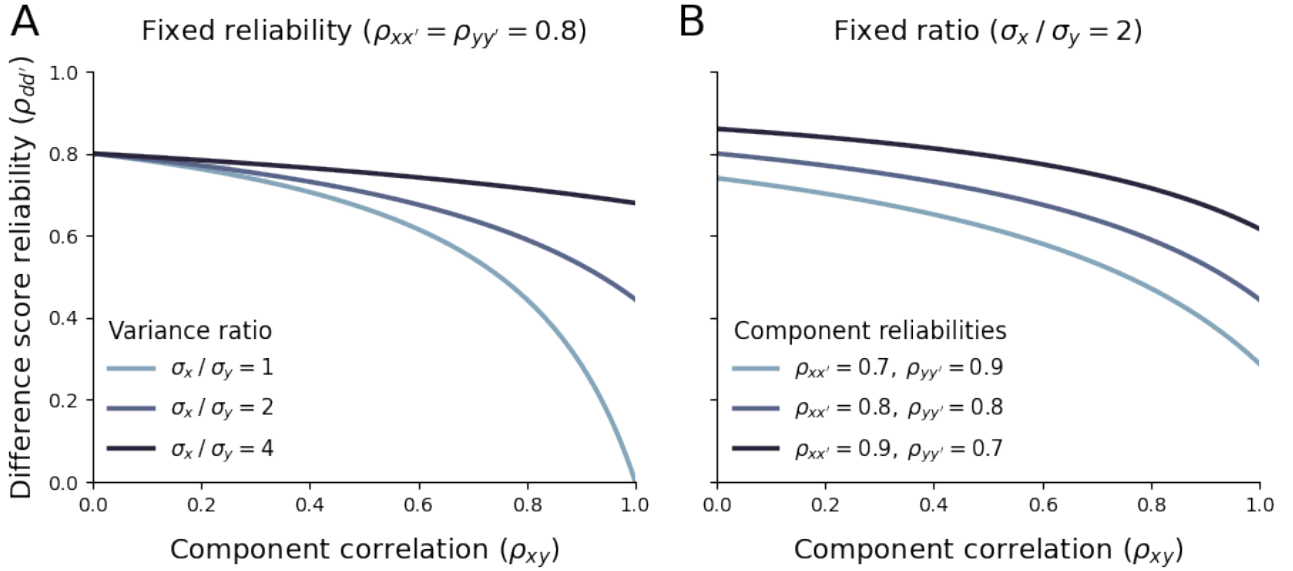


Figure 2: Difference score reliability as a function of the variances and reliabilities of its component measures. (A) Difference-score reliability as the ratio of the component measure variances ( $\sigma_x/\sigma_y$ ) increases, with component reliabilities held fixed ( $\rho_{xx'} = \rho_{yy'} = 0.8$ ). When component measures have equal variances ( $\sigma_x/\sigma_y = 1$ ), large correlations between the measures substantially diminish the reliability of a difference score measure. When the variances are unequal ( $\sigma_x/\sigma_y > 1$ ), even large correlations between the measures are less deleterious for reliability. (B) Difference score reliability as the component reliabilities change, with the ratio of the component-measure variances fixed ( $\sigma_x/\sigma_y = 2$ ). When the variances are unequal, the reliability of a difference score reflects more the component measure with the larger variance (here, component  $x$ ).

### 2.3.2 Identify alternative measures

Rather than improving the reliability of difference scores, one can simply avoid using them in the first place. This recommendation has a long history in experimental psychology. Indeed, because difference scores will virtually always be less reliable, many authors have advocated to abandon them [72–74].

What then are the alternatives to difference scores? Draheim and colleagues [71] provide an in-depth review of alternatives to difference scores in the context of response-time measures, though much of their discussion is applicable to difference scores in general. One possibility is to simply use the component measures (e.g., performance on incongruent trials in a Stroop task alone). Of course, because component scores will be contaminated other sources of variance, such as baseline performance, interpreting component scores should be done with caution. Another approach is to identify alternative measures of task performance. For example, intra-individual response time variability and cognitive efficiency have been identified as correlates of executive control that can be measured reliably [75, 76] and are altered in psychopathology [77, 78].

## 3 Conclusion

We have briefly reviewed issues and research regarding the reliability of cognitive-task measures. Specifically, we introduced a taxonomy of experiment design and analysis

strategies for improving the reliability of cognitive-task measures, highlighting exemplary studies that have successfully implemented such approaches. We hope we have made clear the importance of calculating (and reporting) the reliability of task measures intended for use in psychiatric research. We also hope we have provided a useful guide for experimenters who wish to design a new task, or to improve an existing task, in order to study individual differences in cognitive processing.

We conclude with two important points. First, although we have discussed the importance of task reliability, we have largely avoided the question of when a task measure is “reliable enough”. Though it is tempting to fall back on conventional cutoffs (e.g.,  $\rho \geq 0.7$ ), what constitutes sufficient reliability in actuality will depend on the goal(s) of the researcher. If the goal is to detect a significant individual-differences correlation, such as between a task measure and self-reported symptom measure, then a task measure with “unacceptable” reliability by conventional standards may suffice (e.g., if a researcher has the resources to collect a sample large enough to be adequately powered to detect a correlation at the attenuated magnitude). On the other hand, if a researcher intends to estimate an individual-differences correlation with high precision, or use a task measure in a high stakes setting (e.g., treatment selection for an individual patient), then high reliability may be required. We cannot overstate the value of simulation studies (e.g., [24]) for researchers trying to determine what level of reliability is required to meet their goals and risk preferences.

Second, we would like to emphasize that reliability is but one of many important considerations in the design and evaluation of cognitive task measures. Task measures may be reliable but show poor convergent validity [27, 79], raising questions about whether they actually measure the constructs they are intended to measure. Similarly, task measures may be reliable but exhibit poor ecological validity [80], thus being poor proxies for cognition in real-world settings. Task measures may also be reliable but show poor predictive validity [22], explaining little (unique) variance in other variables of interest (e.g., symptoms, treatment response). Finally, there are many other practical considerations (e.g., task duration, engagement, accessibility) to ensure cognitive tasks are able to be deployed successfully in the clinic or in naturalistic environments at scale [81].

Despite the challenges of making task measures reliable and valid, we are optimistic about their current and future use in biological psychiatry. We believe that, with further efforts towards developing, documenting, and sharing reliable task paradigms, our field can make increased strides towards understanding, predicting, and ultimately relieving psychiatric illness.

## 4 Code availability

The code used to generate the figures in this manuscript is publicly available at <https://github.com/nivlab/biopsych-reliability-review>.

## 5 Acknowledgments

We thank Rachel Bedder and Felix Jan Nitsch for helpful feedback on this manuscript. This project was made possible with support from the National Center for Advancing Translational Sciences (NCATS), a component of the National Institute of Health (NIH), under award number UL1TR003017a and a National Science Foundation Graduate Research Fellowship.

## 6 Disclosures

The authors have no conflicts of interest to declare.

## 7 Citation diversity statement

Recent work in several fields of science has identified a bias in citation practices such that papers from women and other minority scholars are under-cited relative to the number of such papers in the field [82, 83]. Here we sought to proactively consider choosing references that reflect the diversity of the field in thought, form of contribution, gender, race, ethnicity, and other factors. First, we obtained the predicted gender of the first and last author of each reference by using databases that store the probability of a first name being carried by a woman [82]. By this measure and excluding self-citations to the first and last authors of our current paper), our references contain 8.3% woman(first)/woman(last), 14.7% man/woman, 18.6% woman/man, and 58.4% man/man. This method is limited in that a) names, pronouns, and social media profiles used to construct the databases may not, in every case, be indicative of gender identity and b) it cannot account for intersex, non-binary, or transgender people. Second, we obtained predicted racial/ethnic category of the first and last author of each reference by databases that store the probability of a first and last name being carried by an author of color [84, 85]. By this measure (and excluding self-citations), our references contain 4.7% author of color (first)/author of color(last), 15.5% white author/author of color, 18.0% author of color/white author, and 61.8% white author/white author. This method is limited in that a) names and Florida Voter Data to make the predictions may not be indicative of racial/ethnic identity, and b) it cannot account for Indigenous and mixed-race authors, or those who may face differential biases due to the ambiguous racialization or ethnicization of their names. We look forward to future work that could help us to better understand how to support equitable practices in science.

## 8 References

1. Spiegel, J. A., Goodrich, J. M., Morris, B. M., Osborne, C. M. & Lonigan, C. J. Relations between executive functions and academic outcomes in elementary school children: A meta-analysis. *Psychological Bulletin* **147**, 329 (2021).

2. Hartshorne, J. K. & Germine, L. T. When does cognitive functioning peak? The asynchronous rise and fall of different cognitive abilities across the life span. *Psychological science* **26**, 433–443 (2015).
3. Kline, P. *A handbook of test construction (psychology revivals): introduction to psychometric design* (Routledge, 2015).
4. Allen, M. J. & Yen, W. M. *Introduction to measurement theory* (Waveland Press, 2001).
5. Spearman, C. “General intelligence”, objectively determined and measured. *Am. J. Psychol.* **15**, 201–293 (1904).
6. Parsons, S., Kruijt, A.-W. & Fox, E. Psychological Science Needs a Standard Practice of Reporting the Reliability of Cognitive-Behavioral Measurements. *Advances in Methods and Practices in Psychological Science* **2**, 378–395 (2019).
7. Owens, M. M. *et al.* Recalibrating expectations about effect size: A multi-method survey of effect sizes in the ABCD study. *PloS one* **16**, e0257535 (2021).
8. Paap, K. R. & Sawi, O. The role of test-retest reliability in measuring individual and group differences in executive functioning. *Journal of Neuroscience Methods* **274**, 81–93 (2016).
9. Cooper, S. R., Gonthier, C., Barch, D. M. & Braver, T. S. The role of psychometrics in individual differences research in cognition: A case study of the AX-CPT. *Frontiers in psychology* **8**, 1482 (2017).
10. Arnon, I. Do current statistical learning tasks capture stable individual differences in children? An investigation of task reliability across modality. *Behavior research methods* **52**, 68–81 (2020).
11. Pronk, T., Hirst, R. J., Wiers, R. W. & Murre, J. M. Can we measure individual differences in cognitive measures reliably via smartphones? A comparison of the flanker effect across device types and samples. *Behavior Research Methods*, 1–12 (2022).
12. Bruder, L. R., Scharer, L. & Peters, J. Reliability assessment of temporal discounting measures in virtual reality environments. *Scientific reports* **11**, 1–16 (2021).
13. Rouder, J. N. & Haaf, J. M. A psychometrics of individual differences in experimental tasks. en. *Psychon. Bull. Rev.* **26**, 452–467 (2019).
14. Haines, N. *et al.* Learning from the reliability paradox: How theoretically informed generative models can advance the social, behavioral, and brain sciences. *PsyArXiv* (2020).
15. Chen, G. *et al.* Trial and error: A hierarchical modeling approach to test-retest reliability. *NeuroImage* **245**, 118647 (2021).
16. Green, S. B. *et al.* Use of internal consistency coefficients for estimating reliability of experimental task scores. en. *Psychon. Bull. Rev.* **23**, 750–763 (2016).
17. Hedge, C., Powell, G. & Sumner, P. The reliability paradox: Why robust cognitive tasks do not produce reliable individual differences. *Behavior Research Methods* **50**, 1166–1186 (2018).



18. Frey, R., Pedroni, A., Mata, R., Rieskamp, J. & Hertwig, R. Risk preference shares the psychometric structure of major psychological traits. en. *Sci Adv* **3**, e1701381 (2017).
19. Enkavi, A. Z. *et al.* Large-scale analysis of test–retest reliabilities of self-regulation measures. *Proceedings of the National Academy of Sciences* **116**, 5472–5477 (2019).
20. Von Bastian, C. C. *et al.* *Advancing the understanding of individual differences in attentional control: Theoretical, methodological, and analytical considerations* 2020.
21. Nitsch, F. J., Lüpken, L. M., Lüscho, N. & Kalenscher, T. On the reliability of individual economic rationality measurements. en. *Proc. Natl. Acad. Sci. U. S. A.* **119**, e2202070119 (2022).
22. Verdejo-Garcia, A. *et al.* A unified online test battery for cognitive impulsivity reveals relationships with real-world impulsive behaviours. *Nature Human Behaviour* **5**, 1562–1577 (2021).
23. Haaf, J. M. & Rouder, J. N. Developing constraint in Bayesian mixed models. *Psychological methods* **22**, 779 (2017).
24. Rouder, J., Kumar, A. & Haaf, J. M. Why most studies of individual differences with inhibition tasks are bound to fail (2019).
25. McLean, B. F., Mattiske, J. K. & Balzan, R. P. Towards a reliable repeated-measures beads task for assessing the jumping to conclusions bias. *Psychiatry research* **265**, 200–207 (2018).
26. Kucina, T. *et al.* A solution to the reliability paradox for decision-conflict tasks. *PsyArXiv* (2022).
27. Snijder, J.-P., Tang, R., Bugg, J., Conway, A. R. & Braver, T. On the Psychometric Evaluation of Cognitive Control Tasks: An Investigation with the Dual Mechanisms of Cognitive Control (DMCC) Battery. *PsyArXiv* (2022).
28. Rey-Mermet, A., Gade, M., Souza, A. S., Von Bastian, C. C. & Oberauer, K. Is executive control related to working memory capacity and fluid intelligence? *Journal of Experimental Psychology: General* **148**, 1335 (2019).
29. Siegelman, N., Bogaerts, L. & Frost, R. Measuring individual differences in statistical learning: Current pitfalls and possible solutions. *Behavior research methods* **49**, 418–432 (2017).
30. Oswald, F. L., McAbee, S. T., Redick, T. S. & Hambrick, D. Z. The development of a short domain-general measure of working memory capacity. *Behavior research methods* **47**, 1343–1355 (2015).
31. Kyllonen, P. *et al.* General fluid/inductive reasoning battery for a high-ability population. *Behavior Research Methods* **51**, 507–522 (2019).
32. Hausknecht, J. P., Halpert, J. A., Di Paolo, N. T. & Moriarty Gerrard, M. O. Retesting in selection: a meta-analysis of coaching and practice effects for tests of cognitive ability. *Journal of Applied Psychology* **92**, 373 (2007).
33. Scharfen, J., Peters, J. M. & Holling, H. Retest effects in cognitive ability tests: A meta-analysis. *Intelligence* **67**, 44–66 (2018).

34. Paredes, N., Zorowitz, S. & Niv, Y. The Psychometric Properties of the Pavlovian Instrumental Transfer Task in an Online Adult Sample. *Biological Psychiatry* **89**, S132 (2021).
35. Anokhin, A. P. *et al.* Age-related changes and longitudinal stability of individual differences in ABCD Neurocognition measures. *Developmental Cognitive Neuroscience* **54**, 101078 (2022).
36. Salthouse, T. A. Influence of age on practice effects in longitudinal neurocognitive change. *Neuropsychology* **24**, 563 (2010).
37. Schiller, D. *et al.* Preventing the return of fear in humans using reconsolidation update mechanisms. *Nature* **463**, 49–53 (2010).
38. Gulliksen, H. The relation of item difficulty and inter-item correlation to test variance and reliability. *Psychometrika* **10**, 79–91 (1945).
39. Lord, F. M. The relation of the reliability of multiple-choice tests to the distribution of item difficulties. *Psychometrika* **17**, 181–194 (1952).
40. Feldt, L. S. The relationship between the distribution of item difficulties and test reliability. *Applied measurement in education* **6**, 37–48 (1993).
41. Newman, A., Bavik, Y. L., Mount, M. & Shao, B. Data collection via online platforms: Challenges and recommendations for future research. *Applied Psychology* **70**, 1380–1402 (2021).
42. Chandler, J., Mueller, P. & Paolacci, G. Nonnaivete among Amazon Mechanical Turk workers: Consequences and solutions for behavioral researchers. *Behavior research methods* **46**, 112–130 (2014).
43. Robinson, J., Rosenzweig, C., Moss, A. J. & Litman, L. Tapped out or barely tapped? Recommendations for how to harness the vast and largely unused potential of the Mechanical Turk participant pool. *PloS one* **14**, e0226394 (2019).
44. Price, R. B. *et al.* Empirical recommendations for improving the stability of the dot-probe task in clinical research. *Psychological assessment* **27**, 365 (2015).
45. Klingelhofer-Jens, M., Ehlers, M. R., Kuhn, M., Keyaniyan, V. & Lonsdorf, T. B. Robust group-but limited individual-level (longitudinal) reliability and insights into cross-phases response prediction of conditioned fear. *bioRxiv* (2022).
46. Keutmann, M. K., Moore, S. L., Savitt, A. & Gur, R. C. Generating an item pool for translational social cognition research: methodology and initial validation. *Behavior research methods* **47**, 228–234 (2015).
47. Embretson, S. E. & Reise, S. P. *Item response theory* (Psychology Press, 2013).
48. Yoo, A. H., Keglovits, H. & Collins, A. The importance of linguistic information in human reinforcement learning. *PsyArXiv* (2022).
49. Aday, J. S. & Carlson, J. M. Extended testing with the dot-probe task increases test–retest reliability and validity. *Cognitive processing* **20**, 65–72 (2019).
50. Barbosa, J. *et al.* A practical guide for studying human behavior in the lab. *Behavior Research Methods*, 1–19 (2022).

51. Hughes, C. & Graham, A. Measuring executive functions in childhood: Problems and solutions? *Child and adolescent mental health* **7**, 131–142 (2002).
52. Alexander, C., Paul, M., Michael, M., *et al.* The effects of practice on the cognitive test performance of neurologically normal individuals assessed at brief test–retest intervals. *Journal of the International Neuropsychological Society* **9**, 419–428 (2003).
53. Sailer, M., Hense, J. U., Mayr, S. K. & Mandl, H. How gamification motivates: An experimental study of the effects of specific game design elements on psychological need satisfaction. *Computers in human behavior* **69**, 371–380 (2017).
54. Wilson, R. C. & Collins, A. G. Ten simple rules for the computational modeling of behavioral data. *Elife* **8**, e49547 (2019).
55. Palminteri, S., Wyart, V. & Koechlin, E. The importance of falsification in computational cognitive modeling. *Trends in cognitive sciences* **21**, 425–433 (2017).
56. Broomell, S. B. & Bhatia, S. Parameter recovery for decision modeling using choice data. *Decision* **1**, 252 (2014).
57. Melinscak, F. & Bach, D. R. Computational optimization of associative learning experiments. *PLoS computational biology* **16**, e1007593 (2020).
58. Lerche, V. & Voss, A. Retest reliability of the parameters of the Ratcliff diffusion model. *Psychol. Res.* **81**, 629–652 (2017).
59. Waltmann, M., Schlagenhaut, F. & Deserno, L. Sufficient reliability of the behavioral and computational readouts of a probabilistic reversal learning task. *Behavior Research Methods*, 1–22 (2022).
60. Katahira, K. How hierarchical models improve point estimates of model parameters at the individual level. *Journal of Mathematical Psychology* **73**, 37–58 (2016).
61. Myung, J. I., Cavagnaro, D. R. & Pitt, M. A. A tutorial on adaptive design optimization. *Journal of mathematical psychology* **57**, 53–67 (2013).
62. Gonthier, C., Aubry, A. & Bourdin, B. Measuring working memory capacity in children using adaptive tasks: Example validation of an adaptive complex span. *Behavior Research Methods* **50**, 910–921 (2018).
63. Ahn, W.-Y. *et al.* Rapid, precise, and reliable measurement of delay discounting using a Bayesian learning algorithm. *Scientific reports* **10**, 1–10 (2020).
64. Kwon, M., Lee, S. H. & Ahn, W.-Y. Adaptive design optimization as a promising tool for reliable and efficient computational fingerprinting. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging* (2022).
65. Curran, P. J., Cole, V., Bauer, D. J., Hussong, A. M. & Gottfredson, N. Improving Factor Score Estimation Through the Use of Observed Background Characteristics. *en. Struct. Equ. Modeling* **23**, 827–844 (2016).
66. Bertling, M. & Weeks, J. P. Using response time data to reduce testing time in cognitive tests. *Psychological Assessment* **30**, 328 (2018).
67. Ballard, I. C. & McClure, S. M. Joint modeling of reaction times and choice improves parameter identifiability in reinforcement learning models. *Journal of Neuroscience Methods* **317**, 37–44 (2019).

68. Shahar, N. *et al.* Improving the reliability of model-based decision-making estimates in the two-stage decision task with reaction-times and drift-diffusion modeling. *PLoS computational biology* **15**, e1006803 (2019).
69. Palestro, J. J. *et al.* A tutorial on joint models of neural and behavioral measures of cognition. *Journal of Mathematical Psychology* **84**, 20–48 (2018).
70. Chiou, J.-s. & Spreng, R. A. The reliability of difference scores: A re-examination. *The Journal of Consumer Satisfaction, Dissatisfaction and Complaining Behavior* **9**, 158–167 (1996).
71. Draheim, C., Mashburn, C. A., Martin, J. D. & Engle, R. W. Reaction time in differential and developmental research: A review and commentary on the problems and alternatives. *Psychological Bulletin* **145**, 508 (2019).
72. Lord, F. M. The measurement of growth. *ETS Research Bulletin Series* **1956**, i–22 (1956).
73. Cronbach, L. J. & Furby, L. How we should measure “change”: Or should we? *Psychological bulletin* **74**, 68 (1970).
74. Edwards, J. R. Ten difference score myths. *Organizational research methods* **4**, 265–287 (2001).
75. Saville, C. W. *et al.* On the stability of instability: Optimising the reliability of intra-subject variability of reaction times. *Personality and individual differences* **51**, 148–153 (2011).
76. Weigard, A., Clark, D. A. & Sripada, C. Cognitive efficiency beats top-down control as a reliable individual difference dimension relevant to self-control. *Cognition* **215**, 104818 (2021).
77. Kofler, M. J. *et al.* Reaction time variability in ADHD: a meta-analytic review of 319 studies. *Clinical psychology review* **33**, 795–811 (2013).
78. Heathcote, A. *et al.* Decision processes and the slowing of simple choices in schizophrenia. *Journal of abnormal psychology* **124**, 961 (2015).
79. Eckstein, M. K. *et al.* The Interpretation of Computational Model Parameters Depends on the Context. *BioRxiv*, 2021–05 (2022).
80. Steiner, M. D. & Frey, R. Representative design in psychological assessment: A case study using the Balloon Analogue Risk Task (BART). en. *J. Exp. Psychol. Gen.* (Apr. 2021).
81. Germine, L., Strong, R. W., Singh, S. & Sliwinski, M. J. Toward dynamic phenotypes and the scalable measurement of human behavior. *Neuropsychopharmacology* **46**, 209–216 (2021).
82. Dworkin, J. D. *et al.* The extent and drivers of gender imbalance in neuroscience reference lists. *Nature neuroscience* **23**, 918–926 (2020).
83. Bertolero, M. A. *et al.* Racial and ethnic imbalance in neuroscience reference lists and intersections with gender. *BioRxiv* (2020).

84. Ambekar, A., Ward, C., Mohammed, J., Male, S. & Skiena, S. *Name-ethnicity classification from open sources* in *Proceedings of the 15th ACM SIGKDD international conference on Knowledge Discovery and Data Mining* (2009), 49–58.
85. Sood, G. & Laohaprapanon, S. Predicting race and ethnicity from the sequence of characters in a name. *arXiv preprint arXiv:1805.02109* (2018).

# Supplementary materials

## Calculating reliability

There are two measures of reliability of relevance for task measures: internal consistency and test-retest reliability. Internal consistency is the reliability of a measure in a single administration of a task. Test-retest reliability (also known as temporal stability) is the reliability of a task measure across two or more administrations. Even for cross-sectional experiments, where stability over time is of no scientific interest, test-retest reliability is a useful index. This is because estimates of internal consistency tend to be inflated due to between-participants variance from construct-irrelevant, state-dependent factors (e.g., current mood, fatigue). Given sufficient time between testing sessions, test-retest reliability should be less biased by state factors, and thus a better estimate of the true cross-sectional task-measure reliability. However, whereas high internal consistency is always desirable for individual difference studies, high test-retest reliability may not be, depending on the construct a researcher is intending to measure. For example, low one-month test-retest reliability may not be problematic for an index of a transient cognitive process (e.g., mood-dependent attentional biases), but is a problem if individual differences in a cognitive process are hypothesized to be stable across time (e.g., extraversion and social reward processing).

The test-retest reliability of a task measure can be calculated in numerous ways [1]. Perhaps the simplest approach is to compute the Pearson correlation between participants' scores from two sessions. An alternative approach is to calculate the intraclass correlation coefficient (ICC), which decomposes a task measure into true score variance and error variance. There are many formulas for calculating ICC [2], with the critical distinction being whether reliability is based on the *consistency* or *absolute agreement* of a task measure across two administrations. Consistency-based ICCs are affected only by the relative ordering of participants across time; that is, they are insensitive to systematic changes to the actual values of a task measure across time (e.g., due to practice effects on task performance). In contrast, absolute-agreement-based ICCs measure the degree to which scores are stable across time. The type of ICC to use depends on the experimenter's goals and the ultimate use of the task measure.

Calculating the internal consistency of a task measure is more complicated. The most common measure of internal consistency is Cronbach's  $\alpha$ , which is a function of the average correlation across all unique pairs of trials. However, Cronbach's  $\alpha$  is an accurate measure of reliability only under assumptions that are unrealistic for many tasks (e.g., equivalence of trials, uncorrelated measurement error; [1, 3]). As such, internal consistency for task measures is instead usually calculated via split-half reliability, where reliability is estimated after trial data have been divided into two halves. A critical challenge in calculating split-half reliability is in deciding how to partition the data, as estimates of reliability may be also biased if the data partitions violate either of the two above assumptions (for detailed discussion, see [3, 4]). For example, first-second splitting (i.e., partitioning the data into the first and second halves of an experiment) may underestimate reliability due to nonequivalence of the two partitions resulting from practice, fatigue, or other linear time effects. In contrast, odd-even splitting (i.e., partitioning the data



into odd and even trials) may cause bias when behavior across trials is non-independent (i.e., measurement error is correlated across trials), artificially inflating the similarity of data across partitions and thereby decreasing estimates of measurement noise and overestimating reliability. Therefore, where possible, a permutation-based approach to calculating split-half reliability is recommended [1, 4]. Here, reliability is averaged across many thousands of random partitions of the data into halves. (Insofar that Cronbach’s  $\alpha$  is analytically equivalent to the average of all possible split-half reliability estimates [5], permutation-based split-half reliability provides an approximation to Cronbach’s  $\alpha$  while avoiding its problematic assumptions.) For task measures derived from cognitive models, however, it may be prohibitively computationally intensive to employ such an approach due to the need to re-estimate the model for each new subset of the data. Moreover, for learning tasks that are commonly used in computational psychiatry, cognitive-model based task measures cannot be estimated from only a subset of the trials due to inherent non-independence of task behavior across trials. In these cases, one can design tasks with at least two independent blocks. The model can then be fit to each block independently and reliability calculated using the model parameters estimated from each.

As a final point, traditional sum or mean score estimates of performance (e.g., proportion correct responses, mean response time) may substantially underestimate task reliability [6–8]. This is because such summary scores are contaminated by trial-level noise that, in the absence of a sufficiently (possibly prohibitively) large number of trials, increases measurement error (and thus diminishes reliability). Instead, it may be preferable to use trial-level hierarchical models in which observations are organized hierarchically (e.g., individuals within a group, trials within an individual) with variability modeled at both levels. Hierarchical models exert a pooling or regularization effect on person-level variables, in effect correcting for measurement error and improving estimates of reliability [6–8]. The benefits of hierarchical models for estimating reliability have been multiply demonstrated [9–12], though see [13] for discussion of when these benefits may be limited. Using statistical models that more accurately characterize the latent data-generating process (e.g., using the shifted log-normal distribution to model response times) may also improve reliability estimates [7, 14]. For a detailed discussion of hierarchical and generative models in the context of task reliability, see [7, 15].

## Supplementary references

1. Parsons, S., Kruijt, A.-W. & Fox, E. Psychological Science Needs a Standard Practice of Reporting the Reliability of Cognitive-Behavioral Measurements. *Advances in Methods and Practices in Psychological Science* **2**, 378–395 (2019).
2. Koo, T. K. & Li, M. Y. A Guideline of Selecting and Reporting Intraclass Correlation Coefficients for Reliability Research. en. *J. Chiropr. Med.* **15**, 155–163 (2016).
3. Green, S. B. *et al.* Use of internal consistency coefficients for estimating reliability of experimental task scores. en. *Psychon. Bull. Rev.* **23**, 750–763 (2016).

4. Pronk, T., Molenaar, D., Wiers, R. W. & Murre, J. Methods to split cognitive task data for estimating split-half reliability: A comprehensive review and systematic assessment. en. *Psychon. Bull. Rev.* **29**, 44–54 (Feb. 2022).
5. Cronbach, L. J. Coefficient alpha and the internal structure of tests. *psychometrika* **16**, 297–334 (1951).
6. Rouder, J. N. & Haaf, J. M. A psychometrics of individual differences in experimental tasks. en. *Psychon. Bull. Rev.* **26**, 452–467 (2019).
7. Haines, N. *et al.* Learning from the reliability paradox: How theoretically informed generative models can advance the social, behavioral, and brain sciences. *PsyArXiv* (2020).
8. Chen, G. *et al.* Trial and error: A hierarchical modeling approach to test-retest reliability. *NeuroImage* **245**, 118647 (2021).
9. Snijder, J.-P., Tang, R., Bugg, J., Conway, A. R. & Braver, T. On the Psychometric Evaluation of Cognitive Control Tasks: An Investigation with the Dual Mechanisms of Cognitive Control (DMCC) Battery. *PsyArXiv* (2022).
10. Sullivan-Toole, H., Haines, N., Dale, K. & Olino, T. Enhancing the Psychometric Properties of the Iowa Gambling Task Using Full Generative Modeling. *Faculty/Researcher Works* (2022).
11. Brown, V. M., Chen, J., Gillan, C. M. & Price, R. B. Improving the reliability of computational analyses: Model-based planning and its relationship with compulsivity. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging* **5**, 601–609 (2020).
12. Waltmann, M., Schlagenhauf, F. & Deserno, L. Sufficient reliability of the behavioral and computational readouts of a probabilistic reversal learning task. *Behavior Research Methods*, 1–22 (2022).
13. Rouder, J., Kumar, A. & Haaf, J. M. Why most studies of individual differences with inhibition tasks are bound to fail (2019).
14. Price, R. B., Brown, V. & Siegle, G. J. Computational modeling applied to the dot-probe task yields improved reliability and mechanistic insights. *Biological Psychiatry* **85**, 606–612 (2019).
15. Haines, N., Sullivan-Toole, H. & Olino, T. From Classical Methods to Generative Models: Tackling the Unreliability of Neuroscientific Measures in Mental Health Research. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging* (2023).