

# Inattentive responding can induce spurious associations between task behavior and symptom measures

Samuel Zorowitz<sup>1,\*</sup>, Johanne Solis<sup>2</sup>, Yael Niv<sup>1,3</sup>, Daniel Bennett<sup>4</sup>

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

<sup>2</sup>Rutgers-Princeton Center for Computational Cognitive Neuropsychiatry, Rutgers University, NJ, USA

<sup>3</sup>Department of Psychology, Princeton University, NJ, USA

<sup>4</sup>School of Psychological Sciences, Monash University, Victoria, Australia

\*Corresponding author (zorowitz@princeton.edu)

## Abstract

Although online samples have many advantages for psychiatric research, some potential pitfalls of this approach are not widely understood. Here, we detail circumstances in which spurious correlations may arise between task behavior and symptom scores. The problem arises because many psychiatric symptom surveys have asymmetric score distributions in the general population, meaning that careless responders on these surveys will show apparently elevated symptom levels. If these participants are similarly careless in their task performance, this may result in a spurious association between symptom scores and task behavior. We demonstrate this pattern of results in two samples of participants recruited online (total  $N = 779$ ) who performed one of two common cognitive tasks. False-positive rates for these spurious correlations increase with sample size, contrary to common assumptions. Excluding participants flagged for careless responding on surveys abolished the spurious correlations, but exclusion based on task performance alone was less effective.

# Introduction

In recent years, online labor markets (e.g., Amazon Mechanical Turk, Prolific, CloudResearch) have become increasingly popular as a source of research participants in the behavioral sciences [1], in no small part due to the ease with which these services allow for recruitment of large, diverse samples. The advantages of online data collection have also begun to be recognized in psychiatric research [2], where this method offers several distinct advantages over traditional approaches to participant recruitment. The ability to assess psychiatric symptom severity in large general-population samples makes possible large-scale transdiagnostic analysis [3, 4], and facilitates recruitment from difficult-to-reach participant populations [5]. Online labor markets also facilitate re-recruitment, making them an attractive option for validating the psychometric properties of assessment tools [6] or studying clinical processes longitudinally [7].

With the advantages of online data collection also come specific drawbacks. Since participants recruited from online labor markets are typically completing experiments in their homes, they may be more likely to be distracted or multi-tasking during an experiment. They may also be more likely to use heuristic response strategies with the intention to minimize expenditure of time and cognitive effort (e.g., responding randomly on self-report surveys or behavioral tasks). Here, we will refer to such inattentive or low-effort behaviors as careless/insufficient effort (C/IE) responding [8, 9]. Among researchers using online labor markets, a common view is that poor-quality data resulting from C/IE responding can simply be treated as a source of unsystematic measurement error that can be overcome with increased sample sizes [3, 10]. Common practice in online behavioral research is to mitigate poor-quality data using the same screening methods that are typically used in in-person data collection (e.g., excluding participants who perform at- or below-chance on behavioral tasks). However, these methods may be specifically inappropriate for online psychiatry studies, as we detail below.

Here we wish to draw special attention to an underappreciated feature of psychiatric research using self-report symptom surveys. In such surveys, participants rate their endorsement of various psychiatric symptoms and, since most individuals in the general population tend to endorse no or few symptoms in many symptom domains, the resulting ground-truth symptom score distributions tend to be heavily positively skewed [11, 12]. In this situation, the assumption that C/IE responding merely increases unsystematic measurement noise becomes untenable. Because of the positive skew in the ground-truth symptom distribution, participants who respond carelessly to the symptom survey are more likely to report higher levels of symptom endorsement relative to participants who complete the survey attentively [10, 13, 14]. Consequently, unless C/IE survey responses are carefully identified and removed, a considerable proportion of putatively symptomatic individuals in an online sample may, in fact, be participants who have not engaged with the experiment with sufficient attention or effort.

When participants complete both symptom surveys and behavioral tasks—a common study design in computational psychiatry—this artifact has the potential to induce spu-

rious correlations between symptom self-report scores and task behavior. That is, while C/IE behavior is traditionally thought of as a source of noise that can result in type II (false negative) errors, here we suggest that in large-scale online psychiatric studies it can instead result in type I (false positive) errors. Concretely, if the same participants who engage in C/IE responding on surveys (and who therefore inaccurately report high levels of psychiatric symptoms) also respond with insufficient effort on behavioral tasks, this can cause experimenters to observe an entirely spurious correlation between greater symptom severity and worse task performance (see Figure 1). A similar effect has been well documented in personality psychology, where the presence of C/IE responding can induce correlations between questionnaires, and bias factor estimation in factor analysis [8, 10, 15–17].

Here, we demonstrate the real risk that C/IE responding can lead to spurious symptom-task correlations in computational psychiatry research. First, we asked to what extent recent studies in computational psychiatry screen participants based on self-report symptom data. We found that the majority of these studies did not screen participants’ survey data at all, and that very few followed best-practice recommendations for survey data screening. We then asked whether behavioral screening alone was sufficient to identify participants engaging in C/IE responding on psychiatric symptom surveys. In two new datasets from two separate online labor markets, we found that screening based on task behavior fails to completely identify participants engaging in C/IE responding on surveys. Lastly, we investigated whether, under these circumstances, C/IE responding led to spurious correlations between symptom severity and task performance for positively-skewed symptom measures. Consistent with the logic set out above, we confirmed that failure to appropriately screen out C/IE survey responding in the proof-of-concept datasets that we collected would have produced a number of spurious correlations between task behavior and self-reported symptoms that are abolished when data are screened more thoroughly.

## Results

### Narrative review of task and self-report screening practices

First, we sought to what extent recent online studies screen participants in a way that would reduce the risk of spurious correlations due to C/IE participants. We performed a narrative literature review of 49 online human behavioral studies, and evaluated whether and how each study performed task and self-report data screening (see Methods for details of the literature search).

Among studies that we reviewed, approximately 80% (39/49) used at least one method to identify C/IE responding in task behavior (Table 1). Of these, just over half relied on a single screening method, with considerable heterogeneity in behavior screening methods across studies. Most common (46% of all studies) was identifying participants whose performance was statistically indistinguishable from chance-level on some measure of

accuracy. Almost as common (38%) was screening based on low response variability (i.e., excluding participants who predominantly responded in the same fashion across trials, such as using only a single response key).

In contrast, only a minority (19/49, or 39%) of studies screened for C/IE responding in self-report symptom measures. The most common survey screening method was the use of attention checks, which are prompts for which most responses are unlikely given attentive responding. Participants who do not give the correct response to these prompts are therefore likely to be engaged in C/IE responding. Attention checks can be subdivided into instructed items (in which participants are explicitly told which response to select; e.g., ‘Please select “Strongly Agree”’), and infrequency items (in which some responses are logically invalid or exceedingly improbable; e.g., endorsing ‘Agree’ for the question ‘I competed in the 1917 Summer Olympic Games’). Of those studies that specified what type of attention check was used, instructed items were the most common method. As we discuss further below, this is notable because best-practice recommendations for data collection in personality psychology explicitly counsel *against* using instructed-item attention checks [18–20]. Only a handful of studies employed statistical or so-called unobtrusive screening methods such as outlier detection or personal consistency.

In sum, whereas screening for C/IE responding in task behavior was relatively common for online behavioral studies, screening of self-report survey data was far less prevalent. Although this pattern may seem troubling, low rates of survey data screening are not necessarily an issue if screening on task behavior alone is sufficient to remove participants engaging in C/IE responding. That is, screening on survey data may be redundant if there is a high degree of correspondence between task- and survey-based screening methods.

In the next section, we explicitly test this hypothesis in a large sample of online participants completing a battery of self-report surveys and a behavioral task. Specifically, we measure the empirical correspondence between common task- and survey-based screening methods—as identified in our literature review—so that results are informative with respect to typical study designs in online psychiatry research.

## **C/IE participants appear psychiatric when symptoms are rare**

To measure the correspondence of screening measures estimated from task and self-report behavior, we conducted an online behavioral experiment involving a simple decision-making task and a battery of commonly used self-report psychiatric symptom measures (see Methods). A final sample of 386 participants from the Amazon Mechanical Turk (N=186) and Prolific (N=200) online labor markets completed a probabilistic reversal-learning task and 5 self-report symptom measures. The reversal-learning task required participants to learn through trial-and-error which of three options yielded reward most often, and was modeled after similar tasks used to probe reinforcement-learning deficits in psychiatric disorders [21, 22]. The five self-report measures were the 7-up (which

measures symptoms of hypomania), the 7-down (which measures symptoms of depression), the GAD-7, (which measures generalized anxiety symptoms), the BIS/BAS (which measures reward and punishment motivations), the SHAPS (which measures anhedonia symptoms), and the PSWQ (which measures worry symptoms), and were chosen based on previous literature to have a variety of expected response distributions (symmetric and asymmetric). In line with current best-practice recommendations [23], each self-report instrument included one ‘infrequency’ item that could be used to identify C/IE responses in survey data (see Methods for a list of infrequency items). The entire experiment (surveys and task) was designed to require 10 minutes on average to complete (observed mean = 10.28 minutes). To minimize any influence of fatigue on survey responding, participants completed the surveys prior to beginning the task.

To assess the overall quality of the data, we examined the number of participants flagged by the choice accuracy and infrequency item screening measures. Only 26 participants (7%) were flagged as exhibiting choice behavior at or below statistically chance levels in the reversal-learning task. In contrast, 85 participants (22%) endorsed a logically invalid or improbable response on one or more of the infrequency items when completing the self-report symptom measures. This discrepancy in the proportion of participants flagged by each method is consistent with previous research, which found varying levels of sensitivity to C/IE responding across screening methods [24]. The proportion of participants flagged for C/IE responding was marginally but significantly greater on Mechanical Turk compared to Prolific for both task (MTurk:  $N=18/186$ ; Prolific:  $N=8/200$ ; two-tailed, two-sample proportions test:  $z(384) = 2.224$ ,  $p = 0.026$ ,  $h = 0.230$ , 95% CI = [0.006, 0.107]) and survey data (MTurk:  $50/186$ ; Prolific:  $35/200$ ; two-tailed, two-sample proportions test:  $z(384) = 2.223$ ,  $p = 0.026$ ,  $h = 0.227$ , 95% CI = [0.011, 0.176]).

We hypothesise that spurious behavior-symptom correlations may emerge due to a mean-shift in the average level of symptom endorsement in participants engaging in C/IE responding relative to attentive participants. In turn, a mean-shift is expected to occur when the overall rate of symptom endorsement is low; that is, comparably higher scores are more likely for C/IE participants responding at random on a questionnaire with a right-skewed score distribution. In line with our predictions, the average level of symptom endorsement was noticeably exaggerated in C/IE-responding participants for the symptom measures where symptom scores were most positively-skewed (7-up, 7-down, GAD-7; Figure 2). In contrast, where there was higher rates of symptom endorsement overall, the distributions of symptom scores between the two groups of participants were less noticeably distinct. Permutation testing confirmed that observed mean-shifts in symptom scores for C/IE participants were statistically significant for the majority of symptom measures (Table 2).

Hereafter, we use the infrequency-item method as a primary means of identifying C/IE responding in our data. To verify this approach, we conducted three validation analyses. The first analysis compared estimated internal consistency of self-report measures between the C/IE and attentive groups. The logic is that, if C/IE responding manifests as a tendency to respond randomly, we should expect to see a decrease in the consistency of a

measure in the C/IE responding group [24–26]. In line with this reasoning, we observed a reduction in Cronbach’s  $\alpha$  in the C/IE group for the majority of survey instruments (Table 2). A permutation test confirmed that the average decrease in internal consistency across measures was greater than would be expected by chance given the difference in participant numbers between groups (two-tailed, paired-samples  $t$ -test:  $t(6) = -3.689$ ,  $p = 0.021$ ,  $d = 1.506$ , 95% CI = [-0.048, -0.141]).

Second, we quantified the degree to which participants responded to self-report symptom surveys in a stereotyped fashion; that is, we determined if participants exhibited patterns in their responses that were independent of the contents of the survey items. We fit a random-intercept item factor-analysis model [27] to self-report data (see Methods), and for each participant we estimated an intercept parameter that quantified their bias towards using responses on the left or right side of the response scale, regardless of what that response signifies for a particular self-report measure (e.g., low on one symptom scale versus high on another). We observed a credible difference between the average value of this intercept for the two groups ( $\Delta\text{intercept} = -0.67$ , 95% HDI = [-0.78, -0.55]), such that C/IE participants were biased towards using the right-half of survey response options. This translates to a tendency to endorse *more severe* symptoms on the 7-up/7-down and GAD-7 scales (where the rightmost options indicate greater frequency of symptoms) but *less extreme* symptoms or personality traits on the SHAPS and BIS (where the rightmost options indicate lower frequency of symptoms or personality traits) despite these inventories measuring strongly correlated constructs (i.e., depression and anhedonia, anxiety and behavioral inhibition).

Finally, we compared the proportion of participants meeting the cutoff for clinical levels of psychopathology before and after excluding participants based on their responses to the infrequency items. Previous studies have found that applying such measures reduced the prevalence of clinical symptomology in online samples towards ground truth rates from epidemiological studies [13]. On the most positively-skewed measures, the fraction of participants reaching clinical levels of symptom endorsement prior to screening was greater than what would be expected (Table 2). For example, 13.0% of participants scored at or above clinical thresholds for (hypo)mania on the 7-up scale in our sample prior to screening, compared with a 12-month prevalence of 5% in the general population [28, 29], but this rate was reduced to 4.0% (in line with the population prevalence estimates) after exclusion of C/IE respondents. We observed a similar pattern for both major depressive disorder (MDD) and anxiety (population prevalence estimates of 7% and 5% respectively; [11, 30, 31]). Interestingly, the proportion of participants meeting threshold on the GAD-7 was elevated compared to previous literature. We suspect this may reflect elevated rates of state anxiety during the COVID-19 pandemic [32], when these data were collected. In line with previous research, we interpret these inflated rates of clinical symptomology in our sample prior to screening as suggestive of C/IE responding [13].

## Low agreement between task and self-report screening measures

Next, we evaluated the degree of correspondence between behavioral and self-report screening measures in order to determine whether screening on behavior alone was sufficient to identify and remove careless participants. In line with the literature review, we computed multiple measures of C/IE responding from each participant’s task behavior and survey responses (see Methods for description of measures). To measure the degree of correspondence between these behavioral and self-report screening measures, we performed two complementary analyses. First, we computed pairwise correlations on the unthresholded (continuous) measures using Spearman’s rank correlation. The resulting pairwise similarity matrices are presented in Figure 3 (left panel). After correcting for multiple comparisons, there were few significant correlations between the behavioral and self-report screening measures. Only choice accuracy showed significant associations with any self-report measure (specifically, the infrequency and Mahalanobis distance measures). Crucially, the sizes of these observed correlations were roughly half those observed for the correlations between the self-report measures. This is worrisome as it suggests that, although there is some relationship between C/IE responding on tasks and self-report inventories, the relationship is not strong enough to ensure reliable detection of careless participants using task data alone.

Second, we used the Dice similarity coefficient to quantify agreement between different screening methods in the set of participants flagged for exclusion (Figure 3, right panel). This approach quantifies the degree of overlap between the set of would-be excluded participants based on different screening measures under a common exclusion rate. Though some measures have relatively clear threshold cutoffs (e.g., chance level performance for task accuracy), the majority of the measures evaluated here do not. As such, we evaluated the measures with respect to the top 10% of “suspect” participants flagged by each measure, corresponding roughly to the fraction of participants having performed at chance levels on the reversal-learning task. (Results of the same analysis repeated for the top 25% of “suspicious” participants — corresponding roughly to the fraction of participants flagged by the infrequency-item measure — produced similar results; see Table S5.) Results were largely consistent with the correlation analysis: few pairs of task and self-report screening measures achieved levels of agreement greater than what would be expected by chance. The only significant cross-modality pair identified — between the infrequency item and choice accuracy measures — has a Dice similarity coefficient less than 0.4. In other words, when these two measures are used to identify the top 10% of participants most strongly suspected of C/IE responding, they agree on only two out of every five participants. Screening on choice accuracy alone (the most common method identified in our literature review) would fail to identify the majority of participants most likely engaging in C/IE responding as determined by the infrequency items.

Taken together, these results suggest that measures of C/IE responding in task and self-report data do not identify the same set of participants. This means that solely excluding participants on the basis of poor behavioral performance—the most common approach



in online studies—is unlikely to identify participants who engage in C/IE responding on self-report surveys.

## C/IE responding yields spurious symptom-behavior correlations

Here we examine the potential consequences of screening only on task behavior in our data. To do this, we estimated the pairwise correlations between the symptom scores of each of the self-report measures and several measures of performance on the reversal learning task. This analysis emulated a typical computational psychiatry analysis, in which the results of primary interest are the correlations between task behavior and self-reported psychiatric symptom severity.

For each participant, we computed both descriptive and computational-model-based measures of behavior on the reversal learning task (see Methods). To understand the effects of applying different forms of screening, we estimated the correlations between each unique pairing of a self-report symptom measure and measure of behavior under four different conditions: no screening, screening only on task behavior (i.e., only participants whose choice accuracy was above chance), screening only on self-report responses (i.e., only participants who responded correctly on all infrequency items), or both. The resulting pairwise behavior-symptom correlations following each screening procedure are presented in Figure 4. We note that we did not correct these correlation analyses for multiple comparisons, since our purpose was to demonstrate the extent of this issue across multiple behavioral measures and self-report symptoms. Any one of these correlations considered individually can be thought of as emulating a conventional analysis where fewer statistical tests would be performed.

When no rejections based on C/IE responding was applied (i.e., all participants were included in the analysis; Figure 4A), many significant correlations emerged between measures of task behavior and symptom scores, in particular for 4 of the self-report instruments (7-up, which measures symptoms of hypomania; 7-down, which measures symptoms of depression; GAD-7, which measures generalized anxiety symptoms; and BIS, which measures tendencies related to behavioral inhibition). Consistent with our predictions, the majority of these correlations involved symptom measures with asymmetric score distributions. Attending to only the most skewed measures (i.e., 7-up, 7-down, GAD-7), symptom endorsement was correlated with almost every behavioral measure. That is, significant correlations were not restricted only to general behavioral measures often used as proxies for participant effort (e.g., accuracy, inverse temperature  $\beta$ ) but also to measures of specific theoretical interest, such as asymmetry of learning from positive and negative reward prediction errors ( $\kappa$ ). Conversely, we found few significant correlations among symptom measures with more symmetric distributions. This is despite the fact these scales measure similar symptoms and syndromes (e.g., anxiety as measured by the GAD-7 and worry as measured by the PSWQ; depression as measured by the 7-down and anhedonia as measured by the SHAPS).



Next, we excluded participants from analysis based on task-behavior screening (i.e., choice accuracy, removing the 7% of participants exhibiting behavior indistinguishable from chance; Figure 4B). The pattern of correlations was largely unchanged: we again found many significant correlations between measures of behavior and asymmetric symptom measures, but almost no significant correlations involving symmetric symptom measures. This suggests that rejection of participants based on the most common form of behavioral screening (i.e., performance accuracy) had little effect on behavior-symptom correlations as compared to no screening.

In stark contrast, when we rejected participants based on self-report screening (removing 22% of participants who endorsed one or more invalid or improbable responses on the infrequency items; Figure 4C), the number of significant correlations was markedly reduced, particularly for several of the most skewed symptom measures (7-down, GAD-7) and proxy measures of task attentiveness (e.g., accuracy, inverse temperature). This pattern of correlations was largely similar when rejections were applied based on both task and self-report screening measures (Figure 4D). We also note that with stricter screening, the remaining significant correlations were, for the mostly but not always, weaker (Tables S6–S9).

These findings suggest that many of the significant behavior-symptom correlations observed without strict participant screening may indeed be spurious correlations driven by C/IE responding. Importantly, screening based on task behavior alone did not adequately protect against spurious symptom-behavior correlations in the presence of skewed distributions of symptom endorsement. For instance, consider the 7-down scale, a measure of trait depression: had we not screened participants based on infrequency items, we would have erroneously concluded that there were many significant associations between reversal-learning task performance and self-reported depression. Screening on self-report data allowed us to identify that each of these depression-behavior correlations was likely to be spurious.

One possible objection to this interpretation is that the reduction in significant correlations following self-report screening was a result of the reduced sample size after removal of C/IE respondents (which comprised over 20% of the sample). To test this alternative hypothesis, we performed the same correlation analysis after removing random subsets of participants, fixing the sample size to that obtained after excluding C/IE respondents. In this case, the pattern of significant correlations was more similar to that before screening than after screening using the infrequency measure (two-tailed, paired-samples  $t$ -test:  $t(4999) = 262.490$ ,  $p < 0.001$ ,  $d = 3.713$ , 95% CI = [0.136, 0.138]; Figure S2, compare to Figure 4A). Thus, the reduction in significant correlations following screening was unlikely to be driven solely by a reduction in statistical power.

Next, we investigated how spurious correlations depended on sample size. To do so, we performed a bootstrapping analysis where we held fixed the proportion of participants engaging in C/IE responding (i.e., 5%, 10%, 15%, 20%) and increased the total number of participants. Across all analyses, we measured the correlation and between the 7-down depression scale and learning-rate asymmetry ( $\kappa$ ), which we previously identified as likely

exhibiting a spurious association. (The following results are not specific to learning-rate asymmetry and generalize to other pairs of variables; Figure S3).

The outputs of the bootstrapping analysis are presented in Figure 5. We found that, although estimated correlation magnitudes were independent of sample size (x-axis, left panel), the absolute magnitude of the behavior-symptom correlation increased with the proportion of C/IE participants (different coloured circles, left panel). Crucially, we found false-positive rates for spurious correlations *increased* with increases in sample size in our data for all but the smallest rates of C/IE responding (right panel). This runs counter to a common assumption that larger sample sizes are protective against spurious correlations because they serve to mitigate measurement error. Although this assumption is correct for unsystematic measurement error, it no longer holds in the regime of systematic measurement error (where larger sample sizes reduce the variance of estimates, but do not alter their bias). Instead, our results suggest that, except for low rates of C/IE responding, the false-positive rate for behavioral-symptom correlations will become increasingly inflated as the sample size increases.

## Findings replicate in second study with alternative measures

One possible concern with the results presented so far is that they are specific to one instantiation of our experimental design. With more stringent quality assurance protocols during participant recruitment, or perhaps a different task or set of self-report measures, one might wonder if spurious correlations would remain such a threat.

To evaluate the generalizability of our findings, we therefore conducted a conceptual replication experiment in which an independent sample of N=393 participants (N=193 from MTurk using CloudResearch, N=200 from Prolific) completed a more difficult cognitive task, the well-known “two-step task” [33], and an alternate set of self-report measures (see Supplementary Materials B for details). Importantly, participants were recruited *after* CloudResearch and Prolific implemented new protocols to improve data quality on their respective platforms. As a final control measure, participants completed self-report symptom measures as before, but also personality measures with no hypothesized relationship to model-based planning behavior on the two-step task.

For the sake of brevity, we report here only the main pattern of findings (all results are reported in Supplementary Materials B). In the replication sample, 55 out of 393 participants (14%) endorsed a logically invalid or improbable response on one or more of the infrequency items when completing the self-report measures. This is roughly two-thirds of the fraction of participants who were flagged for C/IE responding in the original study, suggesting that the newer quality assurance protocols used by the online platforms are at least partially effective.

In the self-report symptom measures, we replicated the finding that total scores were noticeably exaggerated in participants suspected of C/IE responding, but only for symptom measures where overall rates of symptom endorsement were the lowest (Figure S7;

Table S11). Similarly, we again found that task-based screening and self-report screening measures showed low correspondence (Figure S8; Tables S12–S13); that is, excluding participants on the basis of poor behavioral performance would not have identified and removed participants who engaged in C/IE responding on self-report surveys.

Finally, when we did not apply any exclusions, we observed spurious correlations between performance on the two-step task and total scores for both symptom and personality self-report measures with a mean-shift in scores between attentive participants and participants suspected of C/IE responding (Figure S9). In contrast with our original findings, however, we found that excluding participants based on self-report *or* task screening measures was sufficient to abolish these spurious correlations.

In sum, we replicated most of the main findings from the original study in an independent sample of participants completing a different task and other self-report measures. Although we found that screening on task behavior was sufficient to protect against spurious correlations in the replication sample, it is difficult to generalize and predict when or why this might be the case for other datasets. As such, we still believe that screening for C/IE responding in both task and self-report measures is the best approach to protect oneself against the possibility of spurious correlations.

## **Patients with depression do not fail attention checks more often**

One major concern with performing rigorous screening and exclusion of participants based on C/IE detection methods is that we might inadvertently introduce an overcontrol bias [34]. That is, to this point we have treated the tendency towards C/IE responding as independent from psychopathology. However, to the extent that C/IE responding reflects lack of motivation [35], avoidance of effort [36, 37], or more frequent lapses of attention [38, 39], one might hypothesise a true underlying association between psychopathology and careless responding in online studies. It is thus plausible that rigorous screening of C/IE responding might lead to the differential exclusion of truly symptomatic participants.

To explore this possibility, we embedded attention checks into the self-report measures of two studies of patients with major depressive disorder (see Supplementary Materials C for details). Specifically, N=35 psychiatric patients (confirmed to meet criteria for a diagnosis of major depressive disorder through a structured clinical interview) across 45 unique testing sessions and N=17 healthy controls across 20 unique testing sessions, all recruited through the Rutgers-Princeton Center for Computational Cognitive Neuropsychiatry (i.e., not via online labor platforms), completed a series of self-report symptom measures, online, on their computers from the comfort of their homes. In total, 16 of 65 (24.6%) participants failed one or more attention checks. Subdivided by group, 6 of 20 (30%) healthy participants and 10 of 45 (22%) MDD patients were flagged for C/IE responding.

Using these data, we computed pairwise Bayes factors comparing three candidate mod-

els: attention check failure rates are equal between healthy and MDD patients (M1); failure rates are greater in MDD patients (M2); and failure rates are greater in healthy participants (M3). The model assuming equal rates of failure between healthy and MDD participants was 2.88 times more likely than the model assuming greater rates for MDD patients. In turn, the model assuming lower rates of failure for MDD patients was 1.27 times more likely than the model assuming equal rates. Finally, the model assuming lower rates of failure for MDD patients was 3.65 times more likely than the model assuming higher rates for MDD patients. Only the final comparison exceeds the cutoff value of 3, which is conventionally treated as the minimal amount of evidence required to treat a difference in model fit as meaningful. Although the size of the sample precludes any definitive conclusion, it is noteworthy that the model least consistent with the data was the one where MDD patients are more likely to fail infrequency item attention checks. These data suggest, therefore, that it is unlikely that individuals with high depression symptom severity were disproportionately flagged for C/IE responding in the main analyses. Accordingly, we tentatively conclude that the screening measures we are suggesting are not likely to result in overcontrol bias and false-negative correlations between tasks and symptom measures, at least in the case of individuals with depression. It remains possible that other psychiatric symptoms might be associated with a different pattern of results.

## Discussion

In this study, we highlighted a particular set of circumstances, common in computational psychiatry research done on large online samples, in which spurious correlations may arise between task behavior and self-reported symptomology. When the ground-truth prevalence of a symptom is low in the general population, participants who respond carelessly on measures assessing this symptom may erroneously appear as symptomatic. Careless responding on tasks used to measure cognitive constructs can then masquerade as a correlation between individual differences in these constructs and symptom dimensions. We found repeated evidence for this pernicious pattern in two samples of participants recruited from two popular online labor platforms. False-positive rates for these spurious correlations *increased* with sample size, because the correlations are due to measurement bias, not measurement noise. Importantly, we found that screening on task behavior alone was often insufficient to identify participants engaging in C/IE responding and prevent the false-positive correlations. Unfortunately, a literature review identified this type of screening as the most common practice in online computational psychiatry studies. We recommend instead to screen and exclude participants based on responding on surveys, a practice that abolished many spurious behavior-symptom correlations in our data.

One way of conceptualizing our results is through the lens of rational allocation of mental effort [40]. In any experiment, attentive responding is more effortful than careless responding. As such, participants completing an online task must perform a cost-benefit analysis—implicitly or otherwise—to decide how much effort to exert in responding. The

variables that factor into such calculations are presumably manifold and likely include features of the experiment (e.g., task difficulty, monetary incentives), facets of the participant (e.g., subjective effort costs, intrinsic motivation, conscientiousness), and features of the online labor market itself (e.g., opportunity costs, repercussions for careless responding).

Viewed from the perspective of effort expenditure, our results suggest that participants appraised the cost/benefit trade-off differently for behavioral tasks and self-report surveys. Specifically, we found that only 7% of participants in the first study were at chance-level performance in the task, compared to more than 22% of participants who failed one or more attention-check items in the self-report surveys (a finding that qualitatively replicated in a second study involving a different task). Moreover, different measures of C/IE responding were weakly or not at all correlated between task behavior and self-report responses. This suggests the motivation for effortful responding was greater in the behavioral tasks, though precisely why is unclear. One possibility is that we gave participants a monetary incentive for attentive responding only during the tasks (a common practice, according to our literature review). A second possibility is that participants expected fewer consequences for C/IE responding during the self-report surveys, a reasonable assumption in light of how infrequently previous experiments have screened self-report data. Alternatively, participants may have found the gamified behavioral tasks more engaging or the self-report inventory more tedious. Regardless of the reason, this discrepancy reinforces our observations concerning the inadequacy of behavioral-task screening as a stand-alone method for identifying C/IE responding. Since, in general, participants may appraise costs and benefits of effortful responding differently for behavioral tasks and self-report surveys, screening for C/IE responding on one data modality may in general be unsuitable for identifying it in the other. We therefore recommend screening on each component of an experiment.

One complicating factor for our argument is that C/IE responding may manifest in other ways than simply random responding for both behavioral tasks and self-report surveys. Indeed, there are more ways to respond carelessly than to respond attentively to a task or self-report inventory (e.g., random response selection, straight-lining, zig-zagging, acquiescence bias) [9]. The specific response strategy a participant adopts is likely to reflect the idiosyncratic integration of multiple perceived benefits (e.g., time saved, effort avoided) and costs (e.g., loss of performance bonuses, risk of detection and forfeited pay). As has been previously documented [24], the presence of multiple response strategies makes it clear why certain screening measures are more or less likely to correlate. For example, the inter-item standard deviation and personal reliability measures are both sensitive to statistically random responding, but less sensitive to straight-lining. Most importantly, a diversity of heuristic response strategies highlights the need for many screening measures of C/IE responding, each sensitive to different heuristic strategies.

Here we have focused on the potential for C/IE responding to result in spurious symptom-behavior correlations when rates of symptom endorsement are low, a case common to online computational psychiatry research. Beyond this, we should emphasize that a

diversity of heuristic response strategies entails that there is more than one mechanism by which spurious correlations can emerge. To the extent that the only prerequisite is a mean-shift between attentive and careless participants, ours is not the only situation where one might expect spurious correlations to emerge [16]. For example, random responding on items with *high* base-rate endorsement could yield spurious correlations with precisely the opposite pattern observed here. Conversely, straight-lining may actually suppress correlations when symptom endorsement is low. In sum, without more understanding about the various types of heuristic responding and when each is likely to occur in a sample, it is difficult to predict *a priori* the patterns of systematic bias that may arise for a given study. This is further impetus for experimenters to be wary of C/IE responding and to use a variety of screening measures to detect it.

One objection to the rigorous screening and exclusion of participants based on C/IE detection methods is that we might inadvertently introduce an overcontrol bias. That is, to the extent that C/IE responding might reflect symptoms common to psychopathology (e.g., low motivation, effort avoidance, inattentiveness), rigorous screening of C/IE responding might lead to the differential exclusion of truly symptomatic participants. To explore this possibility, we embedded attention checks into the self-report measures of two studies of patients with major depressive disorder. Though our final sample was small, we did not find evidence that depressed patients were more likely to fail attention checks than healthy controls (if anything, healthy participants were more likely to be flagged by C/IE screening). These results provide preliminary evidence that rigorous C/IE screening is unlikely to result in overcontrol bias. However, further research with larger samples is necessary to validate attention checks in depressed and other patient populations.

Given that the results of our patient study are preliminary and warrant further investigation, researchers might still be wary of possible overcontrol bias. However, when using self-report questionnaires for screening, for overcontrol to seriously impact results it would have to be the case that symptomatic participants frequently endorse improbable or impossible responses to infrequency-item checks (e.g., responding ‘Agree’ to “I competed in the 1917 Olympic Games”). In this case, and even if such participants truly are experiencing severe symptoms of motivation or attention, there is likely to be limited utility in measuring these symptoms using a self-report measure that they are unable to complete veridically. A similar rationale underlies the widespread use of semi-structured interviews and other clinician-report measures rather than self-report measures for in-clinic psychiatric research. We would therefore argue that, if the psychiatric phenomenon being studied is such that this issue warrants concern, the research question may be better suited to an in-person study design involving participants in the clinic who meet full diagnostic criteria than a correlational design involving an online convenience sample.

Notwithstanding the above, one response to this legitimate concern is to take a graded approach to screening and excluding participants [41]. That is, participants could be screened with respect to a multitude of measures and only the consistently flagged participants be removed, thereby reducing the risk of inducing bias. Another possibility is to use sensitivity analysis as an alternative to exclusion, testing whether full-sample ob-



served correlations are robust to the exclusion of participants flagged by measures of C/IE responding. We note that the strict screening approach used in the present study did not preclude us from identifying symptomatic participants or behavior-symptom correlations. Indeed, we found in our sample roughly 10% of participants endorsing symptoms consistent with clinical levels of depression, and approximately 20% consistent with clinical levels of acute anxiety. These estimates are within the realm of epidemiological norms [11, 30, 32]. (We should note, however, that some studies have found elevated rates of psychiatric symptomology in online participants even after controlling for C/IE responding [13].) We also observed some positive correlations between anxiety and choice behavior that were consistent with effects found in previous literature [42–44]. For example, we found higher lose-shift rates and higher learning rates following negative prediction errors correlated with self-reported anxiety. This suggests that the screening methods we employed were not so aggressive as to attenuate behavior-symptom correlations that would be expected from the literature.

There are several notable limitations to this proof-of-concept study. We used a small set of screening measures, and did not employ other recommended procedures (e.g., logging each key/mouse interaction during survey administration to detect form-filling software or other forms of speeded responding [45]). Thus, we cannot be confident that all of the flagged participants were indeed engaging in C/IE responding; similarly, we cannot be certain that we correctly excluded all participants engaged in C/IE responding. We studied behavior-symptom correlations for only two tasks and two sets of self-report instruments. It remains to be seen how generalizable our findings are, although our study design was inspired by experiments prevalent in the online computational psychiatry literature. As suggested above, future studies may find greater correspondence between task and self-report screening measures for more difficult behavioral experiments. Finally, we should note that, unlike previous studies in which some participants were explicitly instructed to respond carelessly [45], we do not have access to “ground truth” regarding which participants were engaging in C/IE responding. Future work testing the efficacy of different screening metrics for identifying instructed C/IE responding may help to identify some of the issues that we have identified here.

This study highlights the need for more research on the prevalence of C/IE responding in online samples and its interactions with task-symptom correlations. Many open questions remain, including under what conditions task- and symptom-screening measures might better correspond, what screening measures are most effective and when, and under what conditions spurious correlations are more likely to arise. For example, we found that screening on task behavior alone was insufficient to prevent putatively spurious correlations for one task (reversal learning) but was sufficient for another task (the two-step task). This discrepancy may reflect differences in the tasks (e.g., the two-step task may be more challenging and thus more sensitive to C/IE responding) or differences in the screening measures (e.g., choice accuracy across 90 trials may be a noisier measure than win-stay lose-shift choice behavior across 200 trials).

One especially pressing question is how sample size affects the likelihood of obtaining



spurious correlations. The results of a bootstrapping analysis in our data suggest that false positive rates are likely to increase with sample size. As computational psychiatry studies move towards larger samples to characterize heterogeneity in symptoms (and to increase statistical power), it will be important to understand how sample size may exaggerate the effects of systematic error. It will also be important to understand how this is moderated by overall C/IE responding rates, which we observed to vary across platforms and time, and which will presumably continue to evolve with changing labor platform and researcher screening practices.

We conclude with a list of concrete recommendations for future online studies involving correlations between task behavior and self-report instruments. We note that these recommendations are not limited to computational psychiatry studies, but are applicable to any online individual-differences cognitive science research involving similar methods (e.g., behavioral economics, psycholinguistics).

Moving forward, we strongly recommend that experimenters employ some form of self-report screening method, preferably one recommended by the best-practices literature (e.g., [9, 13, 16, 19, 24]). Our literature review found that, to date, the majority of online studies assessing behavior-symptom correlations have not used self-report screening, and our results demonstrate that stand-alone task-behavior screening is not necessarily sufficient to prevent spurious symptom-behavior correlations induced by C/IE responding. We therefore encourage experimenters to use a variety of data-quality checks for online studies and to be transparent in their reporting of how screening was conducted, how many participants were flagged under each measure, and what thresholds were used for rejection.

When collecting self-report questionnaire data, we encourage experimenters to use screening methods sensitive to multiple distinct patterns of C/IE responding (e.g., random responding, straight-lining, side bias) and, if possible, to log all page interactions (e.g., mouse clicks, keyboard presses). We specifically recommend experimenters use infrequency-item attention checks rather than instructed-item checks, as multiple studies have now shown that online participants are habituated to and circumvent the latter (e.g., [18–20]; Supplementary Materials B). Participants flagged by suspicious responses on attention-check items should either be excluded from further analysis, or assessed using sensitivity analyses to ensure that observed full-sample correlations are robust to their exclusion.

We found that spurious correlations predominantly affected self-report instruments for which the expected distributions of symptom scores were asymmetric (either positively or negatively skewed). As such, all else equal, symmetrically-distributed measures of a given construct should be preferred to asymmetrically-distributed measures (though this will often be infeasible given that the prevalence of many psychiatric symptoms in the general population is typically small). Scales with reverse-coded items can be used to quantify the consistency of participants’ responses between reverse-coded and non-reverse-coded measures of the same latent construct. With some care, this may be used to identify C/IE responding even for measures that do not include attention-check items [46]. Similarly, it may be beneficial to include multiple self-report surveys of the same

construct to measure consistency across scales.

In our experience, we have found it instructive to review discussions on public forums for participants of online labor markets (e.g., at the time of writing, Reddit, TurkNation). Doing so helps an experimenter identify what screening methods would-be participants are already aware of and prepared to answer correctly. (Several examples of workers discussing common attention checks can be found at the Github repository for this project.)

More broadly, we encourage experimenters in computational psychiatry to be mindful of the myriad reasons why participants may perform worse on a behavioral task. Whenever possible, researchers are encouraged to design experiments where the signature of some psychiatric syndrome could not also be explained by C/IE responding (e.g., [47, 48]). Experimenters should also carefully consider whether an online study is truly appropriate for the research question. In particular, if the project aims to study syndromes associated with considerable difficulty in task or survey engagement (e.g., severe ADHD, acute mania), symptomatic participants are likely to produce responses that cannot be distinguished from C/IE responding. In such a case, correlational research in online samples is likely not the best approach for the research question. Finally, we conclude by noting that it is preferable to prevent C/IE responding than to account for it after the fact [49]. As such, we recommend researchers take pains to ensure their experiments promote engagement, minimize fatigue and confusion, and compensate participants fairly and ethically.

## Methods

### Experiment

#### Sample

409 total participants were recruited to participate in an online behavioral experiment in late June - early July, 2020. Specifically, 208 participants were recruited from Amazon Mechanical Turk (MTurk) and 201 participants were recruited from Prolific. This study was approved by the Institutional Review Board of Princeton University, and all participants provided informed consent. Total study duration was approximately 10 minutes per participant. Participants received monetary compensation for their time (rate USD \$12/hr), plus an incentive-compatible bonus up to \$0.25 based on task performance.

Participants were eligible if they resided in the United States or Canada; participants from MTurk were recruited with the aid of CloudResearch services [50]. (Note: This study was conducted prior to the introduction of CloudResearch’s newest data quality filters [51]). Following recent recommendations [52], MTurk workers were not excluded based on work approval rate or number of previous jobs approved. No other exclusion criteria were applied during recruitment. It is important to note that both CloudResearch and

Prolific use a number of tools (e.g., IP-address screening) to filter out the lowest quality participants. In addition, our custom experiment delivery software (NivTurk; see below) has bot-checking functionality built into it, and rejects from the start participants who are likely to not be human. We are therefore confident that our study is not strongly affected by participants using software to automatically complete the experiment.

Data from several participants were excluded prior to analysis. Three participants (all MTurk) were excluded due to missing data. In addition, we excluded 20 participants who disclosed that they had also completed the experiment on the other platform. This left a final sample of  $N=386$  participants (MTurk:  $N=186$ , Prolific:  $N=200$ ) for analysis. The demographics of the sample split by labor market is provided in Table S1. Notably, the participants recruited from MTurk were older (mean difference = 7.7 yrs), two-tailed, two-sample  $t$ -test:  $t(384) = 6.567$ ,  $p < 0.001$ ,  $d = 0.669$ , 95% CI = [5.4, 10.0]) and comprised of fewer women (two-tailed, two-sample proportions test:  $z(384) = 2.529$ ,  $p = 0.011$ ,  $h = 0.258$ , 95% CI = [0.030, 0.228]).

## Experimental Task

Participants performed a probabilistic reversal learning task, explicitly designed to be similar to previous computational psychiatry studies [21, 22]. On every trial of the task, participants were presented with three choice options and were required to choose one. After their choice, participants were presented with probabilistic feedback: a reward (1 point) or a non-reward (0 points). On any trial one choice option dominated the others. When chosen, the dominant option yielded reward with 80% probability; the subordinate options yielded reward with only 20% probability. The dominant option changed randomly to one of the two previously subordinate options every 15 trials. Participants completed 90 trials of the task (1 learning block, 5 reversal blocks).

As a cover story, the probabilistic reversal learning task was introduced to participants as a fishing game in which each choice option was a beach scene made distinguishable by a colored surfboard with unique symbol. Participants were told they were choosing which beach to fish at. Feedback was presented as either a fish (1 point) or trash (0 points). Participants were instructed to earn the most points possible by learning (through trial-and-error) and choosing the best choice option. Participants were also instructed that the best option could change during the task, but were not informed about how often or when this would occur (see Supplementary Materials A for the complete instructions). Prior to beginning the experiment, participants had to correctly answer four comprehension questions about the instructions. Failing to correctly answer all items forced the participant to start the instructions over.

The task was programmed in jsPsych [53] and distributed using custom web-application software. All experiment code is publicly available (see Code Availability statement). A playable demo of the task is available at <https://nivlab.github.io/jspsych-demos/tasks/3arm/experiment.html>.

## Symptom Measures

Prior to completing the reversal learning task, participants completed five self-report symptom and personality-trait measures. The symptom measures were selected for inclusion based on their frequency in clinical research, and for having an expected mixture of symmetric and asymmetric score distributions.

**Seven-Up/Seven-Down.** The Seven-Up/Seven-Down (7u/7d; [54]) scale is a 14-item measure of lifetime propensity towards depressive and hypomanic symptoms. It is an abbreviation of the General Behavior Inventory [55], wherein only items that maximally discriminated between depression and mania were included. Items are scored on a 4-point scale from 0 (“Never or hardly ever”) to 3 (“Very often or almost constantly”). Total symptom scores on both subscales range from 0 to 21, and are usually strongly right-skewed, with few participants exhibiting moderate to high levels of symptom endorsement.

**Generalized Anxiety Disorder-7.** the Generalized Anxiety Disorder-7 (GAD-7; [56]) is a 7-item measure of general anxiety. The GAD-7 assesses how much a respondent has been bothered by each of seven core anxiety symptoms over the last 2 weeks. Items are scored on a 4-point scale from 0 (“not at all”) to 3 (“nearly every day”). Total scores on the GAD-7 range from 0 to 21, and are usually right-skewed, with few participants exhibiting moderate to high levels of symptom endorsement.

**Behavioral Inhibition/Behavioral Activation Scales.** the Behavioral Inhibition and Behavioral Activation Scales (BIS/BAS; [57]) are a measure of reward and punishment sensitivity. The original 42-item measure was recently abbreviated to a 14-item measure [58], which we use here. Items are scored on a 4-point scale from 1 (“very true for me”) to 4 (“very false for me”). Total scores on the BAS subscale range from 8 to 32, whereas total scores on the BIS subscale range from 4 to 16. Previous reports have found total scores to be symmetrically distributed [59]. Importantly, in order to maintain presentation consistency with the other symptom measures, the order of the BIS/BAS response options was reversed during administration such that “very false for me” and “very true for me” were the left- and rightmost anchors, respectively.

**Snaith-Hamilton Pleasure Scale.** the Snaith-Hamilton Pleasure Scale is a 14-item measure of anhedonia [60]. Items are scored on a 4-point scale from 0 (“strongly agree”) to 3 (“strongly disagree”), where higher scores indicate greater pathology. Total scores on the SHAPS range from 0 to 42, and have previously been found to be somewhat right-skewed [61, 62], with only the minority of participants exhibiting moderate to high levels of symptom endorsement. Importantly, as with the BIS/BAS, the order of the SHAPS response options was reversed during administration such that “strongly disagree” and “strongly agree” were the left- and rightmost anchors, respectively.

**Penn State Worry Questionnaire.** the Penn State Worry Questionnaire is a measure of worry symptoms [63]. The original 16-item was recently abbreviated to a 3-item measure [64], which we use here. Items are scored on a 5-point scale from 0 (“not at all typical of me”) to 4 (“very typical of me”), where higher scores indicate greater pathology.

Total symptom scores range from 0 to 12 and are usually uniformly distributed.

## Analysis

All statistical models fit as part of the analyses (described in detail below) were estimated within a Bayesian framework using Hamiltonian Monte Carlo as implemented in Stan (v2.26) [65]. For all models, four separate chains with randomised start values each took 2000 samples from the posterior. The first 1500 samples from each chain were discarded. As a result, 2000 post-warmup samples from the joint posterior were retained. Unless otherwise noted, the  $\hat{R}$  values for all parameters was less than 1.1, indicating acceptable convergence between chains, and there were no divergent transitions in any chain.

### Validation analyses

To validate the infrequency items as a sensitive measure of C/IE responding, we performed three complimentary analyses. We describe each in turn below.

**Cronbach’s  $\alpha$ .** We compared the average Cronbach’s  $\alpha$ , a measure of internal consistency, between attentive and C/IE participants. To control for the unbalanced numbers of participants in these groups, we performed a permutation test. First, we estimated Cronbach’s  $\alpha$  was estimated for each subscale and group. Next, we computed the average difference in Cronbach’s  $\alpha$  between the two groups. Then we created a null distribution for this statistic by repeating the same analysis but permuting group membership (i.e., randomly assigning participants to either group), holding fixed the sizes of both groups. This procedure was performed 5000 times. To compute a p-value, we tallied the number of null statistics equal to or (absolutely) greater than the observed test statistic.

**Random intercept item factor analysis.** We employed random intercept item factor analysis [27] to detect heuristic patterns of responding. In the model, the probability of observing response level  $k$  (of  $K$  total levels) from participant  $i$  on item  $j$  is defined as:

$$p(y_{ij} = k) = \begin{cases} 1 - \text{logit}^{-1}(\mu_i + x_j \cdot \theta_i - c_{j,1}) & \text{if } y = 1 \\ \text{logit}^{-1}(\mu_i + x_j \cdot \theta_i - c_{j,y-1}) - \text{logit}^{-1}(\mu_i + x_j \cdot \theta_i - c_{j,y}) & \text{if } 1 < y < K \\ \text{logit}^{-1}(\mu_i + x_j \cdot \theta_i - c_{j,K-1}) - 0 & \text{if } y = K \end{cases}$$

where  $\mu_i$  is an intercept for participant  $i$ ;  $\theta_i$  is a vector of latent factor scores for participant  $i$ ;  $x_j$  is a vector of factor loadings for item  $j$ ;  $c_j$  is a vector of ordinal cutpoints for item  $j$ ; and  $y_{ij}$  is the observed response for participant  $i$  on item  $j$ .

In this analysis, we did not estimate the factor loadings but instead treated them as observed. Specifically, we defined the factor loading for each item as a one-hot vector where the only nonzero entry denoted that item’s corresponding subscale. That is, all of the items from a given subscale were assigned to their own unique factor (which was fixed to one). As such, the model estimated one factor score per participant and subscale (akin to the 1-parameter ordinal logistic model).

Crucially, each participant’s responses were also predicted by a random intercept term,  $\mu_i$ , which was not factor specific but instead was fit across all items. This intercept then reflects a participant’s overall bias towards a response level. In our analysis, we coded the response levels such that the smallest value indicated endorsing the leftmost anchor (irrespective of semantic content) and the largest value indicated endorsing the rightmost anchor (irrespective of semantic content). Because the leftmost response option corresponds to symptomology on some scales (SHAPS), and a lack of symptomology for others (GAD-7, 7-up/7-down), we would not expect a consistent nonzero bias in this random intercept term for an attentive participant.

**Clinical cutoffs.** We compared the proportion of participants in our sample reaching the threshold for clinical symptomology before and after applying exclusions. For the GAD-7, previous research has suggested a clinical cutoff score of 10 or higher [11, 31]. Though the 7-up/7-down scales do not have firmly established clinical cutoffs recent work has suggested a cutoff score of 12 or higher [66], which we use here. Finally, the original authors of the SHAPS scale recommended as a cutoff a score of 3 or more when the items are binarized (1, ‘Strongly disagree’ or ‘Disagree’; 0, ‘Strongly agree’ or ‘Agree’). We use this scoring approach in Table 2.

## Correspondence of screening measures

To measure the correspondence of task- and self-report-based screening measures, we estimated a number of standard measures of data quality from each participant’s task behavior (four in total) and self-report responses (five in total). Beginning first with the self-report data, we describe each below.

**Self-report screening measure: Infrequency items.** Infrequency items are questions for which all or virtually all attentive participants should provide the same response. We embedded four infrequency items across the self-report measures. Specifically, we used the following questions:

1. Over the last two weeks, how much time did you spend worrying about the 1977 Olympics? (Expected response: *Not at all*)
2. Have there been times of a couple days or more when you were able to stop breathing entirely (without the aid of medical equipment)? (Expected response: *Never or hardly ever*)

3. I would feel bad if a loved one unexpectedly died. (Expected response: *Somewhat true for me* or *Very true for me*)
4. I would be able to lift a 1 lb (0.5 kg) weight. Expected response: *Agree* or *Strongly agree*)

Prior to conducting the study, the infrequency items were piloted on an independent sample of participants to ensure that they elicited one dominant response. In the main study, we measured the number of suspicious responses made by each participant to these questions. For thresholded analyses, participants were flagged if they responded incorrectly to one or more of these items.

**Self-report screening measure: Inter-item standard deviation.** The inter-item standard deviation (ISD) is an estimate of a participant’s response consistency on a self-report measure [67], defined as:

$$ISD = \sqrt{\frac{\sum_{i=1}^k (y_i - \bar{y})^2}{k - 1}}$$

where  $y_i$  is a participant’s response to item  $i$ ,  $\bar{y}$  is a participant’s average score across all items, and  $k$  is the total number of items for a self-report measure. A composite ISD measure was estimated per participant by summing across each of the seven self-report scales. Larger ISD values indicate lower response consistency.

**Self-report screening measure: Personal reliability.** The personal reliability coefficient is an estimate of a participant’s response consistency on a self-report measure, estimated by correlating the average scores from split-halves of their responses. To avoid any item-order bias, a participant’s personal reliability coefficient for a particular self-report measure was computed from the average correlation from 1000 random split-halves. A composite reliability measure was generated per participant by averaging across each of the seven self-report scales. Smaller reliability coefficients indicate lower response consistency.

**Self-report screening measure: Mahalanobis D.** The Mahalanobis distance is a multivariate outlier detection measure, which estimates how dissimilar a participant is relative to all others. For a participant  $i$ , the Mahalanobis D is defined as:

$$D = \sqrt{(X_i - \bar{X})^T \cdot \Sigma_{XX}^{-1} \cdot (X_i - \bar{X})}$$

where  $(X_i - \bar{X})$  represents the vector of mean-centered item responses for participant  $i$  and  $\Sigma_{XX}^{-1}$  represents the inverted covariance matrix of all items. Greater Mahalanobis D values indicate larger deviations from the average pattern of responding.



**Self-report screening measure: Reading time.** The reading time is the total number of seconds spent filling out a particular self-report measure, adjusted for that measure’s total number of items [13]. A total reading time estimate was estimated for each participant by summing across the adjusted time for each of the seven self-report measures. Shorter scores are indicative of less time having been spent on each item.

**Task-based screening variable: Choice variability.** Choice variability was defined as the fraction of trials of the most used response option per participant. Choice variability could range from 0.33 (all response options used equally) to 1.00 (only one response option used). Values closer to 1.00 are indicative of more careless responding during the task.

**Task-based screening variable: Choice accuracy.** Choice accuracy was defined as the fraction of choices of the reward-maximizing response option. For a task with 90 trials and three response options, a one-tailed binomial test at  $\alpha = 0.05$  reveals chance-level performance to be 37 or fewer correct choices (41%). Lower accuracy values are indicative of more inattentive responding during the task.

**Task-based screening variable: Win-Stay Lose-Shift.** Win-stay lose-shift (WSLS) measures a participant’s tendency to stay with a choice option following a reward versus shifting to a new choice option following a non-reward. WSLS thus measures a participant’s sensitivity to reward feedback on the screen. WSLS was estimated per participant via regression, where the current choice (stay, switch) predicted by the previous trial’s outcome (reward, non-reward) and a stationary intercept. Here we used the first (slope) term to represent a participant’s WSLS tendency. Lower values of this term indicate less sensitivity to reward feedback and are thus indicative of more careless responding during the task.

**Task-based screening variable: Response times.** “Suspicious response time” was defined as the proportion of trials with an outlier response time, here measured as responses faster than 200ms. Greater proportions of outlier response times are indicative of more careless responding during the task.

**Correspondence Analysis.** We measured the correspondence of the above screening measures via two complimentary approaches. First, we computed pairwise correlations on the unthresholded (continuous) measures using Spearman’s rank correlation. Second, we estimated the pairwise rate of agreement on the binarized measures using the Dice similarity coefficient (looking at the top 10% and 25% most suspicious respondents for each measure). The former approach estimates two measures’ monotonic association, whereas the latter approach estimates their agreement as to which participants were most likely engaging in C/IE responding. For significance testing, we used permutation testing wherein a null distribution of similarity scores (i.e., Spearman’s correlation, Dice

coefficient) was generated for each pair of screening measures by iteratively permuting participants' identities within measures and re-estimating the similarity. P-values were computed by comparing the observed score to its respective null distribution. We corrected for multiple comparisons using family-wise error rates [68].

## Correlations between behavior and symptom measures

To quantify the effects of both task and self-report data screening on behavior-symptom correlations, we estimated the pairwise correlations between the symptom scores of each of the self-report measures and several measures of performance on the reversal learning task. For each participant, we computed both descriptive and model-based measures of behavior on the reversal learning task. We describe each in turn below.

**Descriptive measures.** Descriptive task measures included the following: accuracy (the fraction of choices of the reward-maximizing response option), points (the total number of points accumulated over the game), win-stay rates (the fraction of trials on which a participant repeated the previous trial's choice following a reward outcome), lose-shift rates (the fraction of trials on which a participant deviated from the previous trial's choice following a non-reward outcome), and perseveration (the number of trials on which a participant continued to choose the previously dominant response option following a reversal in task contingencies).

**Model-based measures.** The model-based measures were derived from a common reinforcement learning model of choice behavior, the risk-sensitive temporal difference learning model [69]. In this model, the expected value of a choice option,  $Q(s)$ , is learned through cycle of choice and reward feedback. Specifically, following a decision and reward feedback, the value of the chosen option is updated according to:

$$Q_{t+1}(s) = Q_t(s) + \eta \cdot \delta_t$$

where  $\eta$  is the learning rate bounded in the range  $[0, 1]$  (controlling the extent to which value reflects the most recent outcomes) and  $\delta$  is the reward prediction error, defined as:

$$\delta_t = r_t - Q_t(s)$$

where  $r_t$  is the observed reward on trial  $t$ . In the risk-sensitive temporal difference learning model, there are separate learning rates for positive and negative prediction errors, such that positive and negative prediction errors have asymmetric effects on learning. For example, the effect of negative prediction errors on learned values is larger than that of positive errors if  $\eta_p < \eta_n$ , and vice versa if  $\eta_p > \eta_n$ .

Finally, decision-making according to the model is dictated by a softmax choice rule:

$$p(y_t = s) = \frac{\exp(\beta \cdot Q(s))}{\sum_i^S \exp(\beta \cdot Q(s))}$$

where  $\beta$  is the inverse temperature, controlling a participant’s sensitivity to the expected value of the choice options. In sum then, the model-based approach describes a participant’s choice behavior as a function of three parameters  $(\beta, \eta_p, \eta_n)$ .

We fit the reinforcement learning model to each participants’ choice behavior using Stan (details above). Notably, 11 participants (3% of sample) had parameter estimates with poor convergence, i.e.,  $\hat{R} > 1.1$ ; their parameters were removed from the correlation analysis. Participants’ parameters were fit individually (i.e., not hierarchically) so as to prevent bias during parameter estimation from partial-pooling between attentive and C/IE participants. Parameters were sampled using non-centred parameterisations (i.e., all parameters were sampled separately from a unit normal before being transformed to the appropriate range). Of note, the learning rates were estimated via an offset method such that  $\eta_p = \eta + \kappa$  and  $\eta_n = \eta - \kappa$ , where  $\kappa$  is an offset parameter controlling the extent of an asymmetry between the two learning rates. This parameter was also entered into the behavior-symptom correlation analyses.

We confirmed the model adequately fit participants’ choice behavior through a series of posterior checks (Figure S5). In particular, we confirmed the model recapitulated the group-average learning curves for each block of the experiment. Moreover, we confirmed that the model was able to recover reasonably well the choice accuracy for each participant.

The model-based measures included for analysis were: choice sensitivity ( $\beta$ , inverse temperature), positive learning rate ( $\eta_p$ ), negative learning rate ( $\eta_n$ ), and learning rate asymmetry ( $\kappa = \frac{\eta_p - \eta_n}{\eta_p + \eta_n}$ , normalized difference between  $\eta_p$  and  $\eta_n$ ). We chose these measures as they have been previously used to assess performance in clinical samples [22, 42, 70, 71].

**Correlation analysis.** Behavior-symptom correlations (after various forms of screening and exclusion) were estimated using Spearman’s rank correlation. Significance testing was performed using the percentile bootstrap method [72] so as to avoid making any parametric assumptions. These correlation analyses were not corrected for multiple comparisons, since our overarching purpose was to demonstrate the extent of this issue across multiple behavioral measures and self-report symptoms. Any one of these correlations considered individually can be thought of as emulating a conventional analysis where fewer statistical tests would be performed.

## Literature Review

To characterize common data screening practices in online computational psychiatry studies, we performed a narrative literature review [73]. We identified studies for inclusion

through searches on Google Scholar using permutations of query terms related to online labor platforms (e.g., “mechanical turk”, “prolific”, “online”), experimental paradigms (e.g., “experiment”, “cognitive control”, “reinforcement learning”), and symptom measures (e.g., “psychiatry”, “mental illness”, “depression”). We note that it was not feasible to conduct a systematic review, which requires the use of a publication database with reproducible search, because we required Google Scholar’s full-text search in order to identify papers by recruitment method (e.g., Mechanical Turk). We included in the review studies that (a) recruited participants online through a labor platform, (b) measured behavior on at least one experimental task, and (c) measured responses on at least one self-report symptom measure. Through this approach, we identified for inclusion 49 studies spanning 2015 through 2020. The complete list of studies, and search terms used to find them, are included in the Github repository for this study.

Two of the authors (S.Z., D.B.) then evaluated whether and how each of these studies performed data quality screening for both the collected task and self-report data. Specifically, we confirmed whether a study had performed a particular type of data screening, with screening categories determined based on previous taxonomies of screening methods (e.g., [9]). In addition, we assessed the total number of screening measures each study used and if monetary bonuses were paid to participants. This review was not meant to be systematic, but instead to provide a representative overview of common practices in online behavioral studies.

## Data availability

The data that support the findings of this study are openly available on Github at <https://github.com/nivlab/sciops>.

## Code availability

All code for data cleaning and analysis associated with this study is available at <https://github.com/nivlab/sciops>. The experiment code is available at the same link. The custom web-software for serving online experiments is available at <https://github.com/nivlab/nivturk>.

## Acknowledgements

The authors are grateful to Agnes Norbury, Alexandra Pike, and Oliver Robinson for helpful discussion. Research reported in this manuscript was supported in part by the National Institute of Mental Health (R01MH119511; YN), and by the National Center for Advancing Translational Sciences (UL1TR003017; YN). The content is solely the respon-

sibility of the authors and does not necessarily represent the official views of the National Institutes of Health. SZ was supported by an NSF Graduate Research Fellowship. DB was supported by an Early Career Fellowship from the Australian National Health and Medical Research Council (#1165010). The funders had no role in study design, data collection and analysis, decision to publish or preparation of the manuscript.

## Author Contributions Statement

**SZ:** Conceptualization (equal); Software development (lead); Data collection - online (lead); Formal analysis (lead); Writing - Original Draft (lead); Writing - Review & Editing (supporting); Visualization (lead). **JS:** Software development (supporting); Data collection - clinical (lead); Writing - Review & Editing (supporting). **YN:** Writing - Review & Editing (equal); Funding acquisition. **DB:** Conceptualization (equal); Software development (supporting); Data collection - online (supporting); Formal analysis (supporting); Writing - Review & Editing (equal); Visualization (supporting).

## Competing Interests Statement

The authors declare no competing interests.

## Tables

Frequency	Task Screening		Self-Report Screening	
	N=39 (80%)		N=19 (39%)	
Measure	Accuracy	18 (37%)	Attention Check	17 (35%)
	Variability	15 (31%)	Instructed	10 (20%)
	Response Time	7 (14%)	Infrequency	2 (4%)
	Comprehension Check	5 (10%)	Unspecified	5 (10%)
	Other	16 (33%)	Unobtrusive	4 (8%)

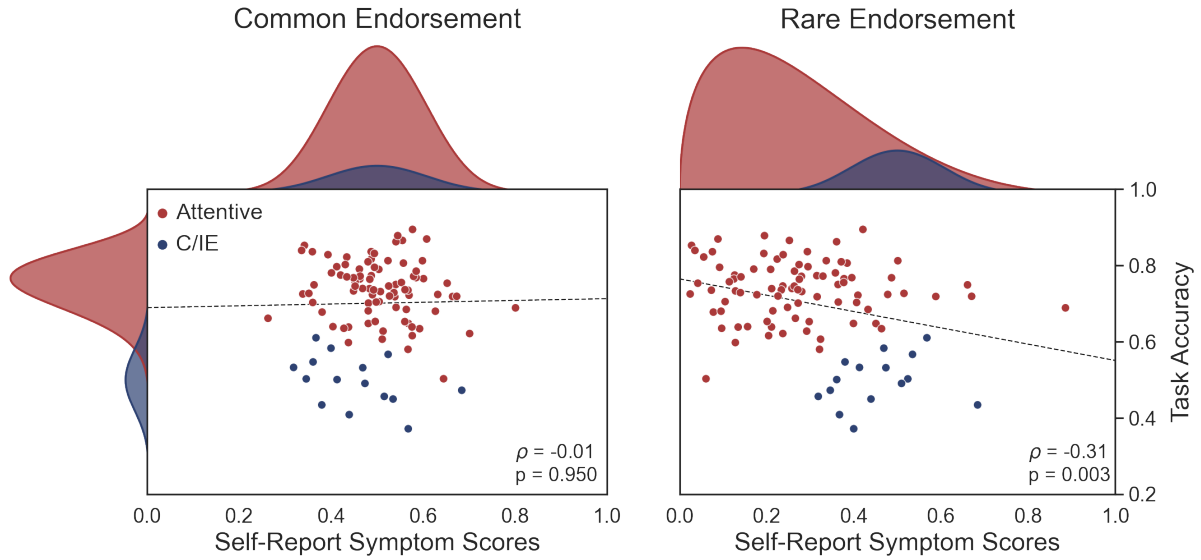
*Table 1:* The prevalence and types of task and self-report data screening practices in a sample (N=49) of recent online behavioral studies.

Subscale	Skew	Total Score				Cronbach's $\alpha$		% Clinical Cutoff	
		Attentive	C/IE	$t$ -value	$p$ -value	Attentive	C/IE	Before	After
7-up	0.806	3.9	10.2	-13.312	<0.001	0.84	0.84	13.0%	4.0%
7-down	0.759	4.8	10.7	-9.987	<0.001	0.94	0.88	17.4%	9.3%
GAD-7	0.753	4.9	9.7	-7.881	<0.001	0.92	0.87	25.9%	17.3%
BIS	0.780	7.7	7.9	-0.542	0.612	0.83	0.62	-	-
BAS	0.171	15.7	16.2	-0.912	0.357	0.84	0.71	-	-
SHAPS	0.256	8.0	10.8	-4.043	<0.001	0.90	0.81	17.9%	14.6%
PSWQ	0.193	4.8	6.7	-4.784	<0.001	0.93	0.81	7.3%	7.0%

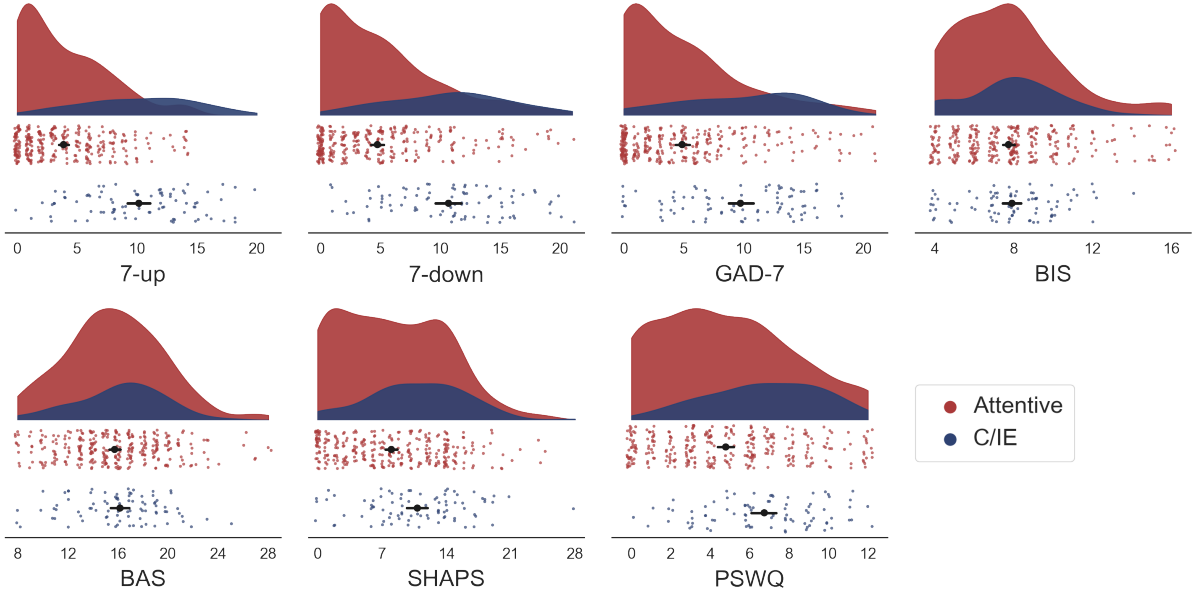
*Table 2:* Descriptive statistics of the self-report symptom measures between attentive and C/IE participants. Skew: the empirical skewness of the distribution of total symptom scores. Total score: the average symptom score across attentive and C/IE participants. Scores were compared using a two-sample  $t$ -test ( $df = 384$ ,  $\alpha = 0.05$ , two-tailed, not corrected for multiple comparisons). Cronbach's  $\alpha$ : a measure of response consistency, where values closer to 1 indicate greater consistency in responses. % Clinical Cutoff: the percentage of participants reaching threshold for clinical symptomology before and after screening based on the infrequency measure. The BIS/BAS scales do not have clinical thresholds.



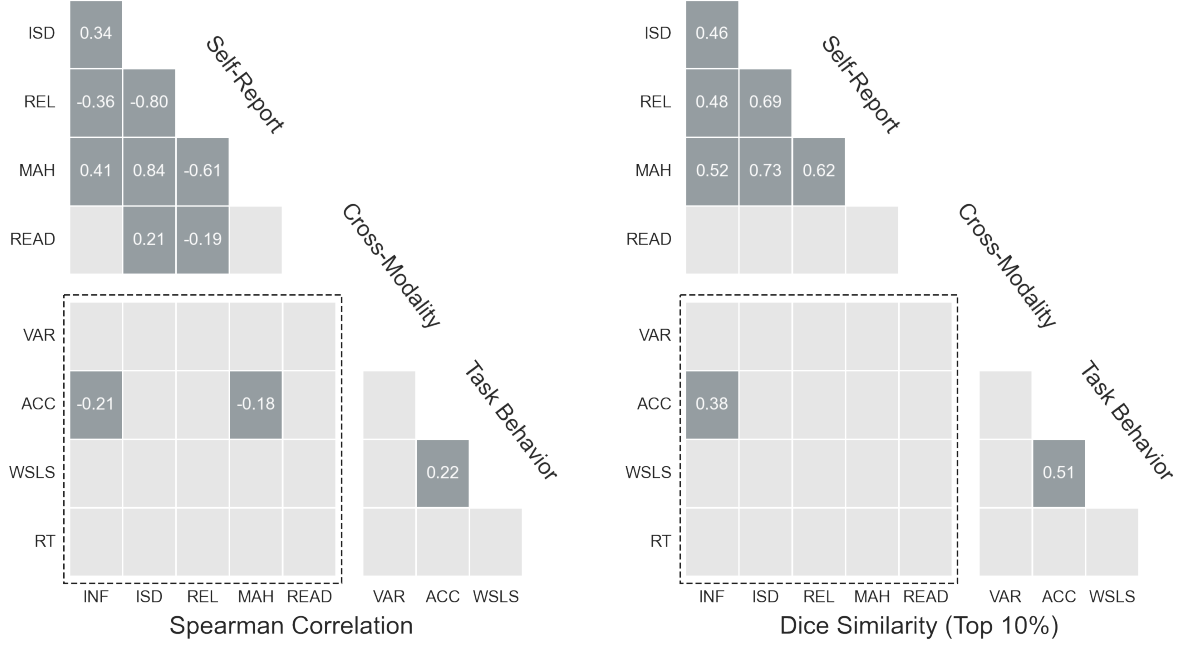
# Figures



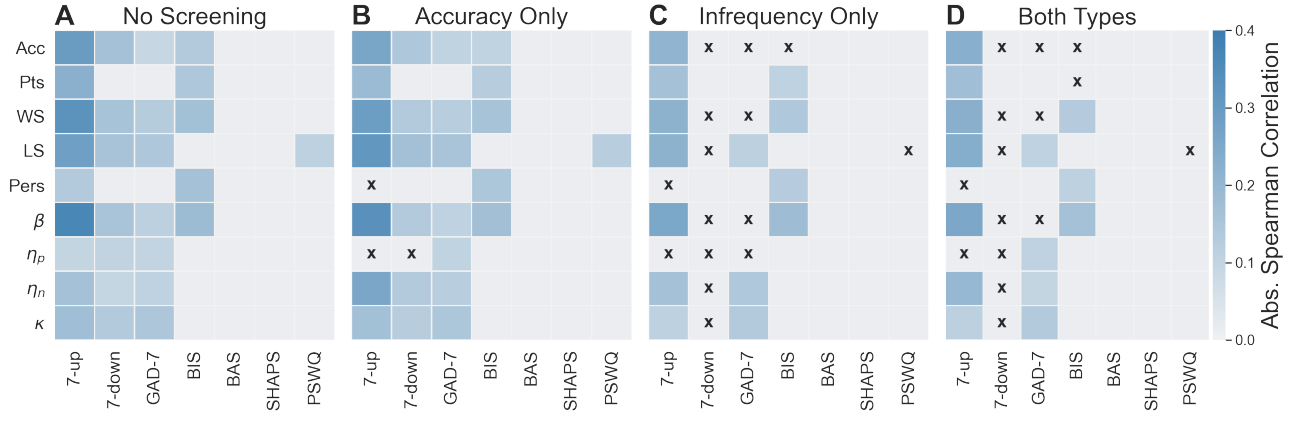
*Figure 1:* Simulated example of how spurious behavior-symptom correlations can arise when symptom endorsement is rare. *Left:* When symptoms are moderately common in the general population, C/IE respondents (blue) are indistinguishable from attentive participants (red) in self-report measures (x-axis, marginal distribution shown on top). Despite the worse task performance of C/IE respondents (y-axis), no correlation arises between symptom scores and task performance (dots are participants drawn from the shown distributions, with 15% C/IE participants; dashed line shows the (lack of) Spearman rank correlation.) *Right:* When symptoms are rare in the general population, careless respondents appear symptomatic in self-report measures. As a result, self-report symptom scores show a significant Spearman rank correlation (two-sided) with task performance.



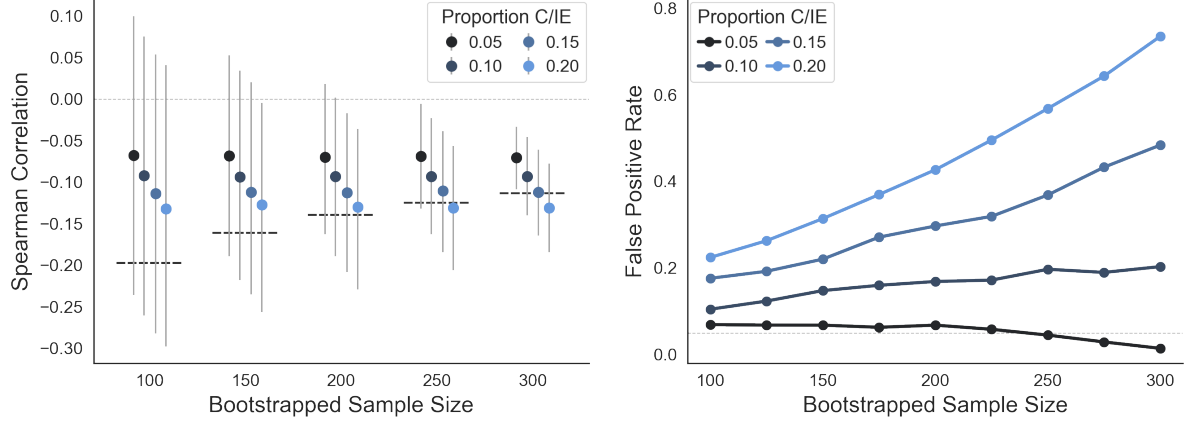
*Figure 2:* Raincloud plots of total symptom scores in attentive ( $N = 301$ ; red) and C/IE ( $N = 85$ ; blue) participants. Each colored dot represents the symptom score for one participant. Black circles: average score within each group (error bars denote 95% bootstrap confidence interval). Shaded plots: distribution of scores for each group of participants. The scales are ordered approximately according to their estimated skew (see Table 2) from top-left (7-up) to bottom-right (PSWQ). The average level of symptom endorsement is most markedly different between groups in symptom measures with the lowest overall rates of endorsement.



*Figure 3:* Similarity of task and self-report data screening measures. Each tile corresponds to the Spearman rank correlation (left) and Dice similarity coefficient (right) between two screening measures across participants ( $N = 386$ ). Similarity indices are thresholded such that only the magnitude of statistically-significant associations (permutation test,  $p < 0.05$ , two-sided, corrected for multiple comparisons) are shown. (Unthresholded values are presented in Tables S3–S5.) Cross-modality correlations between task (y-axis) and self-report screening measures (x-axis) are in the dashed rectangle. Acronyms: INF = infrequency item; ISD = inter-item standard deviation; REL = personal reliability; MAH = Mahalanobis distance; READ = reading time; VAR = choice variability; ACC = choice accuracy; WSLS = win-stay lose-shift rate; RT = suspicious response times.



*Figure 4:* Absolute Spearman rank correlations between task behavior (y-axis) and symptom measures (x-axis) under different regimes of data screening and participant exclusions. (A) No Screening = no exclusions ( $N = 386$ ). (B) Accuracy Only = exclusions based on chance-level performance in the reversal-learning task ( $N = 352$ ). (C) Infrequency Only = exclusions based on invalid or improbable responses to infrequency items ( $N = 301$ ). (D) Both Types = exclusions based on the previous two measures ( $N = 283$ ). Only statistically significant correlations are shown ( $p < 0.05$ , two-sided, not corrected for multiple comparisons; signed correlations are shown in Figure S1 and Tables S6–S9). Black Xs indicate significant correlations abolished under screening. Acronyms: Acc = choice accuracy; Pts = total points earned; WS = win-stay rate; LS = lose-shift rate; Pers = perseveration errors;  $\beta$  = inverse temperature;  $\eta_p$  = positive learning rate;  $\eta_n$  = negative learning rate;  $\kappa$  = learning rate asymmetry.



*Figure 5:* False positive rates for spurious correlations *increase* with sample size. *Left:* Spearman rank correlations and 95% bootstrap confidence intervals between learning rate asymmetry ( $\kappa$ ) and depression scores (7-down) as a function of sample size and proportion of C/IE participants. The thick dashed lines indicate the threshold for statistical significance for the Spearman correlation at the corresponding sample size. Markers are jittered along the x-axis for legibility. *Right:* False positive rates for learning rate asymmetry ( $\kappa$ ) and depression scores (7-down) as a function of sample size and proportion of C/IE participants. False positive rate was calculated as the proportion of bootstrap samples in which the Spearman rank correlation between  $\kappa$  and 7-down was statistically significant ( $p < 0.05$ , two-sided). The horizontal dotted line denotes the expected false positive rate at  $\alpha = 0.05$ .

## References

1. Stewart, N., Chandler, J. & Paolacci, G. Crowdsourcing samples in cognitive science. *Trends in cognitive sciences* **21**, 736–748 (2017).
2. Chandler, J. & Shapiro, D. Conducting clinical research using crowdsourced convenience samples. *Annual review of clinical psychology* **12** (2016).
3. Gillan, C. M. & Daw, N. D. Taking psychiatry research online. *Neuron* **91**, 19–23 (2016).
4. Rutledge, R. B., Chekroud, A. M. & Huys, Q. J. Machine learning and big data in psychiatry: toward clinical applications. *Current opinion in neurobiology* **55**, 152–159 (2019).
5. Strickland, J. C. & Stoops, W. W. The use of crowdsourcing in addiction science research: Amazon Mechanical Turk. *Experimental and Clinical Psychopharmacology* **27**, 1 (2019).
6. 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).
7. Kothe, E. & Ling, M. Retention of participants recruited to a one-year longitudinal study via Prolific. *PsyArXiv* (2019).
8. Huang, J. L., Curran, P. G., Keeney, J., Poposki, E. M. & DeShon, R. P. Detecting and deterring insufficient effort responding to surveys. *Journal of Business and Psychology* **27**, 99–114 (2012).
9. Curran, P. G. Methods for the detection of carelessly invalid responses in survey data. *Journal of Experimental Social Psychology* **66**, 4–19 (2016).
10. Chandler, J., Sisso, I. & Shapiro, D. Participant carelessness and fraud: Consequences for clinical research and potential solutions. *Journal of Abnormal Psychology* **129**, 49 (2020).
11. Lowe, B. *et al.* Validation and standardization of the Generalized Anxiety Disorder Screener (GAD-7) in the general population. *Medical care*, 266–274 (2008).
12. Tomitaka, S. *et al.* Distributional patterns of item responses and total scores on the PHQ-9 in the general population: data from the National Health and Nutrition Examination Survey. *BMC psychiatry* **18**, 1–9 (2018).
13. Ophir, Y., Sisso, I., Asterhan, C. S., Tikochinski, R. & Reichart, R. The turker blues: Hidden factors behind increased depression rates among Amazon’s Mechanical Turkers. *Clinical Psychological Science* **8**, 65–83 (2020).
14. King, K. M., Kim, D. S. & McCabe, C. J. Random responses inflate statistical estimates in heavily skewed addictions data. *Drug and alcohol dependence* **183**, 102–110 (2018).
15. Robinson-Cimpian, J. P. Inaccurate estimation of disparities due to mischievous responders: Several suggestions to assess conclusions. *Educational Researcher* **43**, 171–185 (2014).

16. Huang, J. L., Liu, M. & Bowling, N. A. Insufficient effort responding: Examining an insidious confound in survey data. *Journal of Applied Psychology* **100**, 828 (2015).
17. Arias, V. B., Garrido, L., Jenaro, C., Martinez-Molina, A. & Arias, B. A little garbage in, lots of garbage out: Assessing the impact of careless responding in personality survey data. *Behavior Research Methods*, 1–17 (2020).
18. Barends, A. J. & de Vries, R. E. Noncompliant responding: Comparing exclusion criteria in MTurk personality research to improve data quality. *Personality and individual differences* **143**, 84–89 (2019).
19. Thomas, K. A. & Clifford, S. Validity and Mechanical Turk: An assessment of exclusion methods and interactive experiments. *Computers in Human Behavior* **77**, 184–197 (2017).
20. Hauser, D. J. & Schwarz, N. Attentive Turkers: MTurk participants perform better on online attention checks than do subject pool participants. *Behavior research methods* **48**, 400–407 (2016).
21. Waltz, J. A. & Gold, J. M. Probabilistic reversal learning impairments in schizophrenia: further evidence of orbitofrontal dysfunction. *Schizophrenia Research* **93**, 296–303 (2007).
22. Mukherjee, D., Filipwicz, A. L. S., Vo, K., Satterthwaite, T. D. & Kable, J. W. Reward and punishment reversal-learning in major depressive disorder. *Journal of Abnormal Psychology* **129**, 810–823 (2020).
23. Huang, J. L., Bowling, N. A., Liu, M. & Li, Y. Detecting insufficient effort responding with an infrequency scale: Evaluating validity and participant reactions. *Journal of Business and Psychology* **30**, 299–311 (2015).
24. DeSimone, J. A. & Harms, P. Dirty data: The effects of screening respondents who provide low-quality data in survey research. *Journal of Business and Psychology* **33**, 559–577 (2018).
25. Maniaci, M. R. & Rogge, R. D. Caring about carelessness: Participant inattention and its effects on research. *Journal of Research in Personality* **48**, 61–83 (2014).
26. DeSimone, J. A., DeSimone, A. J., Harms, P. & Wood, D. The differential impacts of two forms of insufficient effort responding. *Applied Psychology* **67**, 309–338 (2018).
27. Maydeu-Olivares, A. & Coffman, D. L. Random intercept item factor analysis. *Psychological methods* **11**, 344 (2006).
28. Merikangas, K. R. *et al.* Lifetime and 12-month prevalence of bipolar spectrum disorder in the National Comorbidity Survey replication. *Archives of general psychiatry* **64**, 543–552 (2007).
29. Merikangas, K. R. & Lamers, F. The ‘true’ prevalence of bipolar II disorder. *Current opinion in psychiatry* **25**, 19–23 (2012).
30. Kessler, R. C., Petukhova, M., Sampson, N. A., Zaslavsky, A. M. & Wittchen, H.-U. Twelve-month and lifetime prevalence and lifetime morbid risk of anxiety and mood disorders in the United States. *International journal of methods in psychiatric research* **21**, 169–184 (2012).



31. Hinz, A. *et al.* Psychometric evaluation of the Generalized Anxiety Disorder Screener GAD-7, based on a large German general population sample. *Journal of affective disorders* **210**, 338–344 (2017).
32. Yarrington, J. S. *et al.* Impact of the COVID-19 Pandemic on Mental Health among 157,213 Americans. *Journal of Affective Disorders* (2021).
33. Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P. & Dolan, R. J. Model-based influences on humans’ choices and striatal prediction errors. *Neuron* **69**, 1204–1215 (2011).
34. Elwert, F. & Winship, C. Endogenous selection bias: The problem of conditioning on a collider variable. *Annual review of sociology* **40**, 31–53 (2014).
35. Barch, D. M., Pagliaccio, D. & Luking, K. Mechanisms underlying motivational deficits in psychopathology: similarities and differences in depression and schizophrenia. *Behavioral neuroscience of motivation*, 411–449 (2015).
36. Cohen, R., Lohr, I., Paul, R. & Boland, R. Impairments of attention and effort among patients with major affective disorders. *The Journal of neuropsychiatry and clinical neurosciences* **13**, 385–395 (2001).
37. Culbreth, A., Westbrook, A. & Barch, D. Negative symptoms are associated with an increased subjective cost of cognitive effort. *Journal of abnormal psychology* **125**, 528 (2016).
38. Kane, M. J. *et al.* Individual differences in the executive control of attention, memory, and thought, and their associations with schizotypy. *Journal of Experimental Psychology: General* **145**, 1017 (2016).
39. Robison, M. K., Gath, K. I. & Unsworth, N. The neurotic wandering mind: An individual differences investigation of neuroticism, mind-wandering, and executive control. *The Quarterly Journal of Experimental Psychology* **70**, 649–663 (2017).
40. Kool, W. & Botvinick, M. Mental labour. *Nature human behaviour* **2**, 899–908 (2018).
41. Kim, D. S., McCabe, C. J., Yamasaki, B. L., Louie, K. A. & King, K. M. Detecting random responders with infrequency scales using an error-balancing threshold. *en. Behav. Res. Methods* **50**, 1960–1970 (Oct. 2018).
42. Huang, H., Thompson, W. & Paulus, M. P. Computational dysfunctions in anxiety: Failure to differentiate signal from noise. *Biological psychiatry* **82**, 440–446 (2017).
43. Harlé, K. M., Guo, D., Zhang, S., Paulus, M. P. & Yu, A. J. Anhedonia and anxiety underlying depressive symptomatology have distinct effects on reward-based decision-making. *PloS one* **12**, e0186473 (2017).
44. Garrett, N., González-Garzón, A. M., Foulkes, L., Levita, L. & Sharot, T. Updating beliefs under perceived threat. *Journal of Neuroscience* **38**, 7901–7911 (2018).
45. Buchanan, E. M. & Scofield, J. E. Methods to detect low quality data and its implication for psychological research. *Behavior research methods* **50**, 2586–2596 (2018).

46. Emons, W. H. Detection and diagnosis of person misfit from patterns of summed polytomous item scores. *Applied Psychological Measurement* **33**, 599–619 (2009).
47. Eldar, E. & Niv, Y. Interaction between emotional state and learning underlies mood instability. *Nature communications* **6**, 1–10 (2015).
48. Hunter, L. E., Meer, E. A., Gillan, C. M., Hsu, M. & Daw, N. D. Increased and biased deliberation in social anxiety. *Nature Human Behaviour* **6**, 146–154 (2022).
49. Ward, M. & Meade, A. W. Applying social psychology to prevent careless responding during online surveys. *Applied Psychology* **67**, 231–263 (2018).
50. Litman, L., Robinson, J. & Abberbock, T. TurkPrime. com: A versatile crowdsourcing data acquisition platform for the behavioral sciences. *Behavior research methods* **49**, 433–442 (2017).
51. Litman, L. *New Solutions Dramatically Improve Research Data Quality on MTurk* <https://www.cloudresearch.com/resources/blog/new-tools-improve-research-data-quality-mturk/>. (Accessed: 2021-02-23).
52. 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).
53. De Leeuw, J. R. jsPsych: A JavaScript library for creating behavioral experiments in a web browser. *Behavior Research Methods* **47**, 1–12 (2015).
54. Youngstrom, E. A., Murray, G., Johnson, S. L. & Findling, R. L. The 7 Up 7 Down Inventory: A 14-item measure of manic and depressive tendencies carved from the General Behavior Inventory. *Psychological Assessment* **25**, 1377–1383 (2013).
55. Depue, R. A. *et al.* A behavioral paradigm for identifying persons at risk for bipolar depressive disorder: a conceptual framework and five validation studies. *Journal of Abnormal Psychology* **90**, 381–437 (1981).
56. Spitzer, R. L., Kroenke, K., Williams, J. B. & Lowe, B. A brief measure for assessing generalized anxiety disorder: the GAD-7. *Archives of Internal Medicine* **166**, 1092–1097 (2006).
57. Carver, C. S. & White, T. L. Behavioral inhibition, behavioral activation, and affective responses to impending reward and punishment: the BIS/BAS scales. *Journal of Personality and Social Psychology* **67**, 319–333 (1994).
58. Pagliaccio, D. *et al.* Revising the BIS/BAS Scale to study development: measurement invariance and normative effects of age and sex from childhood through adulthood. *Psychological assessment* **28**, 429–442 (2016).
59. Cooper, A., Gomez, R. & Aucote, H. The behavioural inhibition system and behavioural approach system (BIS/BAS) scales: Measurement and structural invariance across adults and adolescents. *Personality and individual differences* **43**, 295–305 (2007).
60. Snaith, R. *et al.* A scale for the assessment of hedonic tone: the Snaith–Hamilton Pleasure Scale. *The British Journal of Psychiatry* **167**, 99–103 (1995).

61. Franken, I. H., Rassin, E. & Muris, P. The assessment of anhedonia in clinical and non-clinical populations: further validation of the Snaith–Hamilton Pleasure Scale (SHAPS). *Journal of affective disorders* **99**, 83–89 (2007).
62. Leventhal, A. M. *et al.* Measuring anhedonia in adolescents: a psychometric analysis. *Journal of personality assessment* **97**, 506–514 (2015).
63. Meyer, T. J., Miller, M. L., Metzger, R. L. & Borkovec, T. D. Development and validation of the Penn State Worry Questionnaire. *Behaviour Research and Therapy* **28**, 487–495 (1990).
64. Kertz, S. J., Lee, J. & Bjorgvinsson, T. Psychometric properties of abbreviated and ultra-brief versions of the Penn State Worry Questionnaire. *Psychological Assessment* **26**, 1146–1154 (2014).
65. Stan Development Team. *Stan Modeling Language Users Guide and Reference Manual* <https://mc-stan.org>.
66. Youngstrom, E. A., Perez Algorta, G., Youngstrom, J. K., Frazier, T. W. & Findling, R. L. Evaluating and Validating GBI Mania and Depression Short Forms for Self-Report of Mood Symptoms. *Journal of Clinical Child & Adolescent Psychology*, 1–17 (2020).
67. Marjanovic, Z., Holden, R., Struthers, W., Cribbie, R. & Greenglass, E. The inter-item standard deviation (ISD): An index that discriminates between conscientious and random responders. *Personality and Individual Differences* **84**, 79–83 (2015).
68. Winkler, A. M., Ridgway, G. R., Webster, M. A., Smith, S. M. & Nichols, T. E. Permutation inference for the general linear model. *Neuroimage* **92**, 381–397 (2014).
69. Niv, Y., Edlund, J. A., Dayan, P. & O’Doherty, J. P. Neural prediction errors reveal a risk-sensitive reinforcement-learning process in the human brain. *Journal of Neuroscience* **32**, 551–562 (2012).
70. Brolsma, S. C. *et al.* Challenging the negative learning bias hypothesis of depression: reversal learning in a naturalistic psychiatric sample. *Psychological Medicine*, 1–11 (2020).
71. Ritschel, F. *et al.* Neural correlates of altered feedback learning in women recovered from anorexia nervosa. *Scientific reports* **7**, 1–10 (2017).
72. Wilcox, R. R. & Rousselet, G. A. A guide to robust statistical methods in neuroscience. *Current protocols in neuroscience* **82**, 8–42 (2018).
73. Grant, M. J. & Booth, A. A typology of reviews: an analysis of 14 review types and associated methodologies. *Health information & libraries journal* **26**, 91–108 (2009).