

1

2 **This paper is now published in *Neuroscience & Biobehavioral Reviews***
3 **(<https://www.sciencedirect.com/science/article/pii/S0149763421000221>). This is the submitted**
4 **version of the paper. Please cite as**

5

6 Moriarity, D. P., & Alloy, L. B. (2021). Back to basics: The importance of measurement

7 properties in biological psychiatry. *Neuroscience & Biobehavioral Reviews*, 123, 72-82.

8

9

10

11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

Back to Basics:
The Importance of Measurement Characteristics in Biological Psychiatry

Daniel P. Moriarity*¹, Lauren B. Alloy¹

¹ Temple University

Daniel P. Moriarity was supported by National Research Service Award F31MH122116. Lauren B. Alloy was supported by National Institute of Mental Health grants R01MH077908 and R01MH101168.

*Corresponding author. E-mail: Daniel.moriarity@temple.edu.

Declarations of interest: None.

29

Abstract

30 Biological psychiatry is a major funding priority for organizations that fund mental health
31 research (e.g., National Institutes of Health). Despite this, some have argued that the field has
32 fallen short of its considerable promise to meaningfully impact the classification, diagnosis, and
33 treatment of psychopathology. This may be attributable in part to a paucity of research about key
34 measurement properties (“physiometrics”) of biological variables as they are commonly used in
35 biological psychiatry research. Specifically, study designs informed by physiometrics are more
36 likely to be replicable, avoid measurement concerns that drive down effect sizes, and maximize
37 efficiency in terms of time, money, and the number of analyses conducted. This review describes
38 five key psychometric principles (internal consistency, dimensionality, method-specific variance,
39 temporal stability, and temporal specificity), illustrates how lack of understanding about these
40 characteristics imposes meaningful limitations on research, and reviews examples of
41 psychometric studies featuring a variety of popular biological variables to illustrate how this
42 research can be done and substantive conclusions drawn about the variables of interest.

43

44 Keywords: Biological psychiatry, measurement, methods, reliability, internal consistency,
45 dimensionality

46

Introduction

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

The integration of biological and psychopathological research into the field of biological psychiatry is prioritized highly at the National Institutes of Health. Whereas there is substantial discussion and standard reporting of certain types of measurement characteristics (e.g., dimensionality, retest reliability) for self-report questionnaires, less work has been done to investigate these measurement features for many relevant biological constructs and they are less frequently reported (Hajcak and Patrick, 2015). This is not to say that there has not been important investigation and regular reporting of measurement characteristics specific to biological variables (e.g., intra-assay coefficients of variation). Rather, several metrics key to common methodological and statistical practices in psychiatry research have not received comparable attention for biological variables. This may be due to greater confidence in the measurement of that which is directly observable (e.g., concentrations of analytes in blood). However, the ease with which a construct is operationally defined and measured does not directly translate to measurement qualities suitable for common statistical approaches.

It is important to remember Cronbach and Meehl's (1955) admonition, "One does not validate a test, but only a principle for making inferences" (p. 297). Confidence that a test can measure a variable accurately is not sufficient to know that the test facilitates the inferences tested in statistical models. For that, there is need for a thorough analysis of measurement characteristics germane to the intended data collection and statistical procedures. Armed with information about key measurement characteristics (henceforth referred to as "physiometrics"; Segerstrom & Smith, 2012), researchers can design more cost-effective and well-powered studies that are better indicators of the true associations between variables of interest.

The Perils of a Paucity of Physiometric Research

69 Variables with poor or unknown psychometrics impose multiple limitations to meaningful
70 research. Thus, to ensure that biological psychiatry research reaches its maximum potential
71 utility, it is important to evaluate measurement qualities key to typical methods used in
72 biological psychiatry research to determine what study designs and analytic techniques are best
73 suited to various biomarkers. In this section, we outline some of the risks and constraints
74 imposed by research using variables with poor or unknown measurement characteristics.

75 **Internal Consistency**

76 Many theories in biological psychiatry are about multifaceted biological constructs (e.g.,
77 reward processing, inflammation, etc.); however, studies commonly test multiple individual
78 indices of these larger constructs (Segerstrom and Smith, 2012). Given concerns about the
79 reliability of single-item measures and issues with multiple statistical comparisons, increased use
80 of composite biological variables might benefit replicability in biological psychiatry. When used
81 thoughtfully, composite measures also have the benefit of accentuating variance shared between
82 components and reducing the impact of measurement error. When using composite measures, it
83 is important to report internal consistency, which indicates the level of shared variance between
84 component variables (“true score”) relative to unshared (“error”) variance (Cortina, 1993).
85 Typically, researchers have hypotheses about the relationship between two constructs (e.g.,
86 inflammation and depression); consequently, it is beneficial to maximize the “true score” of their
87 constructs of interest. Although reporting internal consistency for self-report questionnaires is
88 standard practice, it is infrequently reported for applicable biological variables. For example,
89 internal consistency is reported inconsistently for measures involving the creation of a single
90 score from several trials of a task (e.g., error related negativity (ERN)), despite providing insight
91 regarding consistent performance across the task and having implications for effect size (Hajcak

92 et al., 2017). *Thus, whenever aggregate variables are used, it is important to report a measure of*
93 *internal consistency (e.g., Cronbach's α , coefficient Ω).*

94 **Dimensionality**

95 Another important consideration when working with aggregate measures is the concept
96 of dimensionality. Dimensionality refers to the degree to which a set of variables indicates the
97 presence of one or more higher-order constructs. For example, under traditional
98 conceptualizations of psychopathology, all behaviors on a depression questionnaire are
99 associated with the construct of depression. Similarly, an assortment of biological variables (e.g.,
100 different proinflammatory proteins) could serve as markers of a higher-order construct (e.g.,
101 inflammation). It also is important to consider potential construct heterogeneity, the possibility
102 that several lower-order constructs (e.g., pro- and anti-inflammatory processes) might comprise a
103 larger construct of interest (e.g., inflammation).

104 Empirical evaluation of dimensionality is possible with dimension reduction techniques
105 such as exploratory factor analysis (EFA) and principal components analysis (PCA). Both
106 approaches investigate the structure of data with the logic that if all component variables are
107 indicators of the same process, they should be strongly associated with one another (i.e., have
108 high internal consistency, Clark & Watson, 1995, 2019; Loevinger, 1957). As such, dimension
109 reduction approaches can help identify whether sets of variables are unidimensional or
110 multidimensional in nature as well as components that might not load onto any of these
111 processes (Tabachnick and Fidell, 2013). The primary theoretical distinction between the two is
112 that the dimensions found in EFA are theorized to cause the variables, whereas the dimensions
113 found in PCA are simply aggregates of observed variables. Statistically, only shared variance is
114 analyzed in an EFA, but all variance is analyzed in a PCA.

115 Modeling decisions uninformed by dimensionality research can have negative
116 implications. Assuming unidimensionality that is not present (i.e., aggregating unrelated
117 components) reduces internal consistency and, consequently, the maximum observable effect
118 size (Hajcak et al., 2017). Relatedly, if only some dimensions/indicators are related to a criterion
119 of interest, aggregating them with unrelated variables might wash out true effects. Alternatively,
120 falsely assuming multidimensionality reduces power via failure to aggregate shared variance of
121 interest. Further, it introduces issues with multiple comparisons.

122 However, these techniques are not appropriate for all datasets. It is important to consider
123 that the maximum number of dimensions is constricted by the number of indicator variables
124 tested. In other words, there needs to be enough variables per dimension to statistically anchor
125 each dimension. Further, datasets with lower numbers of variables, higher dimensionality, and
126 weaker associations between the variables and the dimensions require higher sample sizes to
127 produce stable results (Guadagnoli and Velicer, 1988). Additionally, it is ill-advised to draw
128 conclusions about dimensionality without thoughtful consideration of biological plausibility.
129 *Consequently, it is important to consider dimensionality when multiple indicators of a broader*
130 *construct of interest are collected before proceeding with hypothesis testing involving that*
131 *construct. However, modeling decisions should be informed both by empirical investigation (if*
132 *appropriate in the context of the dataset used) and biological plausibility.*

133 **Method-specific Variance**

134 Although not a “metric” in the sense of something explicitly testable and reportable like
135 the other characteristics reviewed here, a critical measurement issue for biological psychiatry is
136 method-specific variance. In addition to the “random” variance that contributes to measurement
137 error, there is variability associated with the specific method of measurement (e.g., self-report,

138 behavioral, psychophysiological) that is unrelated to the true construct of interest (Patrick et al.,
139 2013). Consequently, two measures of the same construct using different methods will have
140 smaller associations compared to two measures using similar modalities (e.g., self-report
141 correlated with biological vs. self-report correlated with self-report). Given that biological
142 psychiatry is, by definition, a multimodal field, this is a pervasive issue that needs to be
143 considered when designing studies and interpreting results. *Thus, method-specific variance*
144 *should be considered for all studies including multiple measurement modalities. This issue*
145 *should inform power analyses, measurement error-adjusted analytic techniques, and*
146 *consideration of aggregating multimethod assessments of the same construct. For a more*
147 *detailed review of this issue and strategies to address it, see Patrick et al. (2019).*

148 **Temporal Stability**

149 Whereas a measure given to multiple people at a single time point has two sources of
150 variance (between-person differences and measurement error), a measure given multiple times
151 introduces a third source of variability: within-person variance. Measures with low within-person
152 variability (small changes over time) have high temporal stability. Temporal stability is most
153 frequently quantified using retest Pearson correlations (correlating scores on a measure at two
154 different time points) and intraclass correlation coefficients (ICCs, which quantify the proportion
155 of stable between-person differences across multiple time points). It is standard practice to report
156 (or at least cite other work about) the temporal stability of self-report measures, but it is reported
157 less consistently for biological variables (e.g., Moriarity et al., 2020b). This is concerning, given
158 that information about temporal stability is necessary to interpret the probability with which a
159 score at baseline will be similar to the score at follow-up. It is important to note that highly stable
160 measures are not always the goal; many biological constructs would be expected to have both

161 trait (relatively stable) and state (varying across time and situational factors) components. Target
162 temporal stability should be informed by the conceptual stability of the construct in question
163 (e.g., few would expect mood to be 100% stable in a community sample over the course of a
164 year). *Temporal stability should be reported for all longitudinal studies. It should be calculated*
165 *in the sample when repeated measures are available, or estimates reported from existing studies*
166 *when calculation within the sample is impossible.*

167 **Temporal Specificity**

168 Somewhat related is the concept of temporal specificity. Longitudinal data are necessary
169 to establish directionality of associations; however, time between data points is an important
170 methodological consideration. For example, the relationship between eating a hot pepper and
171 experiencing pain after a couple minutes would not be as strong days after the meal. Thus,
172 exploratory analyses are necessary to evaluate how the relationships between variables might
173 fluctuate as a function of time (including potential developmental considerations). Temporally-
174 informed study designs could improve replicability, provide information about when changes in
175 biological risk factors manifest behaviorally (and vice-versa), and inform treatment studies given
176 expected delays between interventions and symptom reduction (e.g., anti-inflammatory
177 treatments for depression). *Thus, the field would benefit from more exploratory studies*
178 *investigating the temporal specificity of associations of interest to identify optimal time lags*
179 *between measurements.*

180 **Artificial Effect Size Deflation and Power**

181 The practical implications of many biological psychiatry studies are often questioned
182 because they frequently have small effect sizes, which could be directly impacted by the use of
183 measures uninformed by their psychometrics (such as those reviewed above). To illustrate,

184 consider the formula for the maximum correlation between two variables as a function of their
185 reliability: $r_{xy}(\text{max}) = \sqrt{r_{xx}r_{yy}}$ where r_{xy} represents the maximum possible correlation
186 between variables x and y , r_{xx} represents the reliability of variable x and r_{yy} represents the
187 reliability of variable y (Davidshofer and Murphy, 2005). Only if two measures are perfectly
188 reliable (both r_{xx} and $r_{yy} = 1$) can the maximum correlation = 1. As reliability decreases, so does
189 the maximum observable correlation. Consider two research teams testing the same hypothesis
190 and using the same measure for variable x ($r_{xx} = .70$), but different measures for variable y ($r_{yy} =$
191 $.70$ for Team A but $r_{yy} = .30$ for Team B). The maximum observable correlation is $.70$ for Team
192 A, but only $.46$ for Team B. Similar results have been found concerning the relationship between
193 internal consistency and effect sizes (Hajcak et al., 2017).

194 This penalty is magnified in more complex designs. For example, many variables in
195 biological psychiatry (e.g., inflammation) are theorized to be mediators between stress and
196 psychopathology (e.g., Moriarity et al., 2018; Slavich and Irwin, 2014). Mediation analyses
197 involve calculating the product of the association between i) the focal predictor and the mediator
198 (a' pathway) and ii) the mediator and the outcome variable (b' pathway). Thus, unreliability of
199 the mediator will dampen both estimates. Consequently, the downward bias introduced by poor
200 reliability is effectively squared when calculating their product.

201 This bias also exists for group comparisons, which often occur in biological psychiatry in
202 the form of case-control studies (e.g., Ng et al., 2019). The test statistics for these analyses
203 (independent samples t -tests and between-subjects ANOVAs) are a ratio of the magnitude of the
204 group difference divided by a variance component. Poor reliability inflates variability,
205 decreasing the maximum observable effect. For example, consider a researcher using an
206 independent samples t -test to compare levels of interleukin (IL)-6 between participants with

207 Major Depressive Disorder (MDD) and non-depressed controls. The formula for an independent
 208 samples t-test is $t = \frac{M_1 - M_2}{SE}$. Suppose the difference in IL-6 for individuals with MDD vs. non-
 209 depressed controls ($M_1 - M_2$) is .30. In scenario A, the standard error of this difference (SE) is .15,
 210 and the t -score will = 2. The critical value that the t -score must be above to be significant at $p <$
 211 .05 is 1.96, so the researchers have a significant result. Now imagine scenario B, in which the
 212 group difference is the same, but the SE of this difference increases to .2 because of less reliable
 213 IL-6 measurement. Now the t -score is 1.5, which is not significant, despite having the same
 214 observed difference between the groups. The same logic applies for standardized (but not
 215 unstandardized) measures of effect size (e.g., Cohen's $d = \frac{M_1 - M_2}{SD_{pooled}}$). Given the same difference
 216 between two means, as the standard deviation increases, d decreases. However, this does not
 217 mean that measurement error always results in attenuated effect sizes. Although it is true that the
 218 median standardized effect size will be lower when estimated with vs. without error, random
 219 error variance also can result in over-estimates (Segerstrom and Boggero, 2020), leading to false
 220 positives that could inspire misguided studies and intervention efforts.. *Thus, inflated variability*
 221 *caused by unreliable measures can cause true effects to be overlooked both in terms of*
 222 *probability under null-hypothesis testing as well as their substantive implications via*
 223 *standardized effect sizes.* Given the importance of individual differences research in the Research
 224 Domain Criteria (RDoC; Cuthbert & Kozak, 2013) initiative, this is a key (and addressable)
 225 source of bias in popular analytic strategies for NIH-funded research.

226 **Examples of Psychometric Research in Biological Psychiatry**

227 Below, several examples of psychometric research investigating a variety of biological
 228 variables are reviewed to illustrate the techniques used and conclusions about the variables of
 229 interest.

230 **Internal Consistency**

231 As previously discussed, strong internal consistency is evidence that various components
232 of a measure are responded to similarly. To illustrate the importance of investigating internal
233 consistency for neural measures, Hajcak and colleagues (2017) evaluated error-related negativity
234 (ERN) averaged across multiple trials as a function of the number of trials completed by
235 participants in two groups (with and without generalized anxiety disorder). The study reported
236 two measures of internal consistency: Cronbach's α (how representative one trial was of all
237 trials) and split-half reliability (correlating the average scores from the odd and even trials). They
238 found that α increased sharply between four and eight trials, and modestly until approximately
239 fourteen trials, after which α only increased subtly. Cronbach's α reached a maximum of .75 -
240 .85, which was comparable to the Spearman-Brown corrected split-half reliability ($r_{sb} = .71-.75$).
241 The lack of reliability when fewer trials were included is an expected feature of Cronbach's α ,
242 and dovetails with concerns about the reliability of single-item/few-item indicators. Further, the
243 diminishing returns of increased trials reflects that more trials only decreases random error, not
244 systematic error (e.g., error introduced by data collection techniques). These results can help
245 researchers plan the ideal number of trials to minimize participant burden without resulting in
246 data with subpar measurement qualities and, consequently, limited utility. Additionally, they
247 highlight one way of comparing different methods of data collection. For example, comparing
248 the trajectories and plateaus of internal consistency as number of trials increases could provide
249 insight on ratios of random vs. systematic error for two different ERN measures.

250 Kaye, Bradford, and Curtin (2016) present a thorough investigation of several
251 measurement qualities (internal consistency, temporal stability, and effect size stability, the latter
252 two will be discussed later) of acoustic startle (defensive reflex in response to brief, startling

253 noise probes) and corrugator responses (reaction of the corrugator muscle associated with
254 frowning) during a no-shock/predictable shock/unpredictable shock (NPU) task, an affective
255 picture viewing task, and resting state task over two study visits (approximately one-week apart).
256 Specifically, they evaluated Spearman-Brown corrected split-half reliability between odd and
257 even trials as a measure of internal consistency. Further, the authors compared performance of
258 within-person standardized (Bradford et al., 2015) vs. unstandardized scores for startle
259 potentiation and the time domain and frequency domain for corrugator potentiation. For the sake
260 of brevity, this review will focus on startle potentiation. For the NPU task, the internal
261 consistency for raw scores was higher than standardized scores for both predictable and
262 unpredictable startle responses, with scores ranging from good to adequate ($r_{sb} = .81, .64, .57,$
263 $.52$, respectively). For the affective picture viewing task, internal consistency for startle
264 modulation was poor for all scores, but standardized scores were better for pleasant, and raw
265 scores were better for unpleasant, startle modulation (raw pleasant $r_{sb} < .00$, standardized
266 pleasant $r_{sb} = .16$, raw unpleasant $r_{sb} = .14$, standardized unpleasant $r_{sb} < .00$). Because within-
267 subject standardized scores would have no utility for the resting state task, only internal
268 consistency was reported for raw scores ($r_{sb} = .95$). In addition to their descriptive value,
269 comparison of different types of responses and the influence of within-person standardization
270 across several tasks is informative for the establishment of best-practices for these behavioral
271 tasks.

272 Given the rise in popularity and high cost of functional magnetic resonance imaging
273 (fMRI) in biological psychiatry, investigation of measurement properties of these methods is
274 crucial. Luking and colleagues (2017) evaluated the split-half internal consistency for ERPs and
275 blood oxygen level-dependent (BOLD) responses to monetary gain and loss feedback (an fMRI

276 measure) within the ventral striatum and medial and/or lateral prefrontal cortex using Spearman-
277 Brown corrected split-half reliability (comparing odd/even trials). Similar to Kaye et al. (2016),
278 they compared several scoring methods: raw scores, difference scores (gain – loss), and residual
279 scores (gain controlling for loss). Raw BOLD responses across all regions and ERPs to both gain
280 and loss feedback demonstrated high internal consistency ($.66 \geq r_{sb} \geq .86$). Raw scores had
281 consistently higher internal consistency than residual scores ($.26 \geq r_{sb} \geq .50$), which had
282 uniformly higher internal consistency than difference scores ($.02 \geq r_{sb} \geq .36$). Thus, although
283 residual scores may not have ideal internal consistency, they might be preferable over
284 subtraction-based difference scores for studying between-person differences in within-person
285 processes with these measures.

286 Instead of concluding that difference scores (common in many areas of biological
287 psychiatry) are universally unreliable, it is important to consider *why* reliability was lowest for
288 the difference scores, and under what context difference scores have utility. First, when variance
289 associated with one variable is removed from another (either via subtraction or creating a
290 residual term), the variance removed will be from the reliable variance because it is impossible
291 for two variables to share *random* error. This reduction in reliability is greater when the two raw
292 variables are highly correlated (Thomas and Zumbo, 2012). However, as emphasized in the
293 discussion of temporal stability above, reliability needs to be considered in light of the expected
294 true reliability. For reasons beyond the technical scope of this review (see Rogosa and Willett,
295 1983), when the individual differences in the difference score are not small, the reliability of the
296 difference score will be more similar to the reliability of the raw scores. There also is evidence
297 that BOLD difference scores that contrast win and loss conditions vs. neutral, instead of
298 comparing win to loss conditions, can result in more reliable estimates (Holiga et al., 2018;

299 Plichta et al., 2012), but the appropriateness of this approach depends on the research question at
300 hand. Alternatively, many have argued that polynomial regression is a preferable technique to
301 using difference scores altogether (Edwards, 2001).

302 It is important to note that residual/difference scores also hold the potential to isolate
303 theoretically relevant variance in certain designs. For example, consider a study that compared
304 P3 amplitudes (an event related potential) to aversive vs. neutral stimuli (used to index general
305 reactivity) as predictors of threat sensitivity, finding the split-half reliability excellent for both
306 conditions ($r_{sb} = .92$ and $.90$, respectively; Perkins et al., 2017). Split-half reliability for the
307 difference between the two conditions (aversive-neutral) was poor ($r_{sb} = .29$). Recalling that
308 variance removed when creating a difference score always comes from true variability, never
309 random error, this decrease in reliability is not a surprise. As would be expected considering the
310 relationship between reliability and correlations described above, the absolute value of the
311 correlation between the difference score and threat sensitivity ($r = -.12$) was smaller than the
312 correlation between general reactivity and threat sensitivity ($r = .16$). However, a larger
313 proportion of the systematic variance (true score) in the difference score was associated with
314 threat sensitivity (i.e., $(-.12^2/.29) * 100 = 5.00\%$) compared to general reactivity (i.e., $(.16^2/.92) * 100 = 2.78\%$). This approach was particularly important when considering that the association
315 between general reactivity and threat sensitivity was positive, but that the association between
316 the variance unique to the aversive condition and threat sensitivity was negative. Thus, the
317 variance from general reactivity could washout the association unique to the aversive condition if
318 it were not removed from the variable. Consequently, it is important to consider how variables
319 with modest reliability, but that include substantial amounts of criterion-related variance, can be
320 informative.

322 **Dimensionality**

323 Recall the example of inflammation as a complex construct often indexed by several
324 indicators (Segerstrom and Smith, 2012). One study of atherosclerosis (Egnot et al., 2018)
325 assessed the dimensionality of several inflammatory proteins and coagulation biomarkers
326 (specifically, CRP, IL-6, fibrinogen, Lp(a), sICAM-1, PTX-3, and D-dimer). The results of the
327 EFA found a two-factor solution: Factor 1 consisted of CRP, IL-6, and fibrinogen; Factor 2
328 consisted of D-dimer and PTX-3, whereas sICAM-1 and Lp(a) did not load on either factor.
329 Factor 1 was interpreted to represent a non-specific inflammatory process, whereas Factor 2 was
330 interpreted to indicate coagulation burden. The authors then tested the factors as predictors of
331 several outcomes, finding some associations unique to only one of the two factors. For example,
332 although both factors were positively associated with risk for low ankle brachial index, higher
333 levels of coagulation burden (Factor 2), but not inflammation (Factor 1), were associated with
334 elevated common femoral artery intima-media thickness, suggesting that coagulation burden
335 might be a better indicator of subclinical peripheral artery disease than inflammation.

336 Independent component analysis (ICA) is a technique for investigating dimensionality
337 primarily used with neuroimaging and EEG data. Kakeda et al. (2020) used ICA as a data-driven
338 approach to identify brain regions that might differ in grey matter volume between individuals
339 with depression and controls, and whether the volume in these regions correlated with serum
340 TNF α . Specifically, they used source-based morphometry (which applies an ICA to a segmented
341 image) to arrange the voxels into common morphological features of grey matter concentration
342 among participants. Results indicated fourteen independent structural components; however,
343 based on previous work (Williams, 2016), Kakeda and colleagues excluded four primarily
344 cerebellar networks. Of the ten remaining components, two (a prefrontal network and an insula-

345 temporal network) had less grey matter volume in a group of participants with depression
346 compared to controls. Of these two, serum TNF α was significantly negatively correlated with the
347 prefrontal network, but was not significantly correlated with the insula-temporal network.

348 **Method-specific Variance**

349 As described earlier, a major obstacle for biological psychiatry research is domain-
350 specific method variance, the systematic tendency for two measures of the same construct using
351 different modalities (e.g., self-report vs. biological vs. behavioral) to have smaller associations
352 than two measures using the same modality. Ostensibly, one reason for this is that measures from
353 disparate modalities each contribute unique method-specific error (variance related to the
354 measurement method and unrelated to the construct of interest; Patrick et al., 2013). This
355 suggests that the integration of indices of a construct across multiple methods of measurement
356 into single variables, described as the “cross-domain approach” (Patrick et al., 2013; Venables et
357 al., 2018), might accentuate the shared variance related to the construct of interest, improving
358 utility and construct validity.

359 To illustrate this, Nelson, Patrick, and Bernat (2011) measured three event-related
360 potential (ERP) measures (ERN and P3 response to target stimuli from a flanker task and P3
361 response to feedback stimuli from a gambling feedback task) and investigated a) whether these
362 measures represent overlapping indicators of externalizing proneness, and b) whether they index
363 a shared neural process that accounts for their individual associations with externalizing
364 proneness. Results of an EFA suggested that a single factor accounted for the covariance among
365 all three variables, and that all three variables contributed similarly to this shared factor. To
366 evaluate whether this factor represented brain processes associated with externalizing proneness,
367 Nelson and colleagues (2011) ran another EFA including the three ERP measures as well as a

368 self-report measure of externalizing proneness, again finding a single factor. Results of analyses
369 using the aggregated ERP factor found that the aggregate measure had stronger correlations with
370 the majority of physiological and psychometric externalizing proneness criterion variables tested
371 than did the individual ERP measures. In fact, the composite factor out-performed comparison
372 ERP measures (not included in the composite) in predicting externalizing proneness, likely due
373 to the composite variable accentuating the shared externalizing proneness-related variance in the
374 individual ERP variables. However, as described above (and discussed by the authors), a factor
375 analysis on three ERP components and a self-report measure is not enough to provide a
376 convincing evaluation of the true structure of these measures or provide enough options to
377 support alternative models. In other words, there were not enough components to anchor more
378 than one factor, so the factor analytic solution could, at most, feature one aggregate measure
379 and/or unrelated variables. Still, this study serves as an example of how variable aggregation can
380 result in variables with stronger predictive validity than the component parts.

381 To extend this work, Venables and colleagues (2018) first ran EFAs on several indices of
382 inhibition-disinhibition within specific measurement domains (self-report, behavioral
383 performance, brain response). Consistent with the ERP study above, indices within discrete
384 measurement domains revealed single factor solutions. All possible pairwise correlations
385 between these three domain factors were significantly positively correlated. Next, two
386 confirmatory factor analyses (CFA) were estimated: the first specifying all indices across the
387 three measurement domains loading onto a single factor, and the second specifying three lower
388 order factors corresponding with each measurement method that, in turn, load onto a higher order
389 *cross-domain* factor. The former demonstrated poor model fit, but the cross-domain factor model
390 fit the data well. Further, comparative fit indices found significant differences in model fit,

391 suggesting that inhibition-disinhibition is best represented by a cross-measurement domain,
392 hierarchical factor structure. Additionally, the cross-domain factor frequently demonstrated
393 significant correlations with the vast majority of criterion variables tested, whereas
394 measurement-domain specific scores were less likely to be correlated with criterion variables
395 from other measurement domains. Thus, these results demonstrate how thoughtful investigation
396 of dimensionality in biological psychiatry can improve the construct validity of variables by the
397 creation of cross-measurement domain composites that ameliorate concerns about a) the
398 reliability of single-item measures (which are common in biological psychiatry) and b)
399 downward-biased estimates due to measurement domain-specific variability.

400 **Temporal Stability**

401 Out of all the psychometric characteristics described above, biological psychiatry
402 probably has done the best with assessing and reporting temporal stability (the reliability of a
403 measure between different time points). However, there are many constructs of interest for which
404 there is a paucity of research on this topic, especially when considering the wide breadth of study
405 durations seen in behavioral health research. Before reviewing some examples of temporal
406 stability research in biological psychiatry, it is important to emphasize that temporal stability
407 estimates are only informative for the duration in which they are studied. Unfortunately, across
408 all disciplines of behavioral health research, it is commonplace for previous work to be cited as
409 evidence that a measure has sound temporal stability with no reference to the duration for which
410 the measure's stability originally was assessed. Further, it also is essential to reiterate that having
411 low temporal stability is not always indicative of a poor measure. The temporal stability of a
412 measure is dependent on, and constrained by, stability of the construct under question. If one
413 evaluated the 6-month temporal stability of depressed mood and height in a sample of adults, one

414 would expect height to be more stable. Other contextual concerns, such as age, also are important
415 to consider. For example, one would expect relatively lower 6-month temporal stability of height
416 in a sample of 10-year-olds than a sample of adults.

417 The most straightforward metric of temporal stability is retest reliability using Pearson's
418 r , the correlation between a measure at two different time points. In addition to internal
419 consistency metrics, Kaye et al. (2016) (described above) also investigated one-week temporal
420 stability of startle and corrugator responses to three tasks (NPU, affective picture viewing, and
421 resting state) comparing raw vs. within-person standardized scores (Bradford et al., 2015) as well
422 as differences in the effect size of task manipulations (predictable and unpredictable potentiation
423 for the NPU task and pleasant and unpleasant modulation for the affective picture viewing task)
424 between the two sessions. Similar to above, this review only will cover startle responses for the
425 sake of brevity.

426 Temporal stability was higher for raw scores for both predictable and unpredictable
427 startle potentiation during the NPU task (both $r = .71$) compared to within-person standardized
428 scores ($r = .58$ and $.49$, respectively). When comparing the effect size of NPU manipulations
429 between study visits, no significant differences were observed for raw or standardized
430 predictable startle potentiation and raw unpredictable startle potentiation (all $\eta_p^2 = .001-.033$, $p >$
431 $.05$), but the standardized startle potentiation was smaller at the second visit ($\eta_p^2 = .04$, $p = .03$),
432 suggesting that the manipulation lost potency over time. Regarding the affective picture viewing
433 task, one-week temporal stability was poor for both raw and standardized scores for pleasant
434 startle modulation ($r < .00$ and $= .08$, respectively), but was higher for the unpleasant startle
435 modulation ($r = .50$ for raw, $r = .40$ for standardized). The effect sizes for the raw pleasant and
436 unpleasant startle modulations were not significantly different after one week ($\eta_p^2 = .02$, $p = .10$;

437 $\eta_p^2 = .03; p = .09$, respectively). It is interesting to note that the effect sizes for the standardized
438 pleasant and unpleasant startle modulations differed between testing sessions ($\eta_p^2 = .05, p = .02$;
439 $\eta_p^2 = .10, p < .001$, respectively), but in opposite directions (Visit 2 was smaller for pleasant
440 startle modulation, but larger for unpleasant). As mentioned above, standardized scores for the
441 resting state task have no utility, but the raw scores had high one-week temporal stability ($r =$
442 $.89$) and scores were smaller at Visit 2 ($\eta_p^2 = .21, p < .001$, respectively). There was no
443 manipulation during (and consequently, no effect size for) the resting state task. In sum, these
444 results demonstrate how different analytic approaches (i.e., raw vs. within-person standardized
445 scores) can influence important temporal dynamics of behavioral tasks such as stability and the
446 potency of the manipulation, which have important implications for designing and interpreting
447 research using repeated measures of these tasks.

448 Temporal stability also can be influenced by how extreme values are handled, as
449 evidenced by Landau et al. (2019), a study investigating salivary CRP. Immunoassays use
450 standard concentrations of an analyte to generate a standard curve, on which sample values are
451 interpolated. Many samples have values that are flagged by the procedure as too high or low to
452 fit onto the standard curve. In “strict” standard curve datasets, these extreme values are excluded;
453 in “relaxed” standard curve datasets, they are extrapolated outside the standard curve range.
454 There are several techniques currently used to handle these values: list-wise deletion, pair-wise
455 deletion, multiple imputation (extreme values replaced with multiply imputed values), and
456 winsorization (extreme values replaced with the most extreme value on the standard curve).
457 Landau and colleagues (2019) applied each of these four techniques to a strict and a relaxed
458 dataset, resulting in eight total datasets. Additionally, they compared the reliability of samples
459 taken in the morning compared to the evening, given evidence of diurnal variation in CRP (Out

460 et al., 2012). The average two-day Pearson r was .49 for morning samples and .60 for evening
461 samples, suggesting that evening samples might be more stable. Winsorization of extreme values
462 resulted in the highest temporal stability, regardless of time of day (mean winsorized morning $r =$
463 .61, mean winsorized evening $r = .77$, mean nonwinsorized morning $r = .45$, mean
464 nonwinsorized evening $r = .54$) or whether the dataset was strict or relaxed (mean winsorized
465 strict $r = .70$, mean winsorized relaxed $r = .68$, mean nonwinsorized strict $r = .47$, mean
466 nonwinsorized relaxed $r = .52$). On average, relaxed datasets had higher stability than strict
467 datasets (mean $r = .56$ vs. $.52$). However, it is important to always consider data management
468 techniques in the context of one's specific dataset. For example, winsorization might be less
469 appropriate when there are many extreme cases in a dataset. Further, the decision to modify
470 observed values should always involve contemplation about how "extreme" values are defined,
471 the likelihood that they are valid (not the result of measurement error), and the influence
472 "extreme" values would have on planned analyses (e.g., assumptions of normality, sensitivity to
473 outliers).

474 It will come as no surprise that, in addition to statistical procedure, measurement
475 procedure can influence temporal stability as well. In addition to the actual method of data
476 collection (e.g., specific self-report measure, particular imaging scanner model), some biological
477 variables can be measured from different sources. For example, inflammatory proteins most
478 frequently are measured via assaying blood samples (Moriarity et al., 2020a; Muscatell et al.,
479 2016), but salivary measures have been increasing in popularity because they are less expensive
480 and invasive than blood-based methods. However, the utility and comparability of these methods
481 has been questioned as salivary markers of inflammation might reflect local, rather than
482 systemic, immune function (Riis et al., 2015). Out and colleagues (2012) made an important

483 contribution to this discussion by comparing the one- and two-year retest reliabilities of both
484 plasma and salivary measures of CRP in a sample of adult women. Plasma CRP had higher one-
485 year retest reliability than saliva CRP between years 2 and 3 ($r = .70$ vs. $.57$), but lower
486 reliability between years 1 and 2 ($r = .53$ vs. $.61$). Plasma CRP also had higher two-year
487 reliability ($r = .58$ vs. $.46$). Thus, results indicate comparable, but not identical, one and two-year
488 retest stabilities when using these two methods to measure CRP.

489 Another important factor to consider when assessing temporal stability is the role of
490 human development. Particularly for youth undergoing drastic growth and developmental
491 changes, it is plausible that temporal stabilities of many biological variables will differ compared
492 to adults. Riis and colleagues (2014) extended the previous study to a sample of adolescent girls
493 using a similar design (i.e., 3 yearly measurements of plasma and saliva inflammatory analytes).
494 This study assessed nine cytokines, but did not measure CRP, so results cannot be directly
495 compared. Controlling for age, the average year 1 to year 2, year 2 to year 3, and year 1 to year 3
496 reliabilities were higher for serum compared to saliva (average $r_s = .61$ vs. $.30$, $.33$ vs. $.25$, and
497 $.40$ vs. $.34$, respectively). However, when comparing the stability of individual proteins, a more
498 complex picture emerged. One-year retest reliability was uniformly higher for plasma between
499 years 1 and 2 ($r_s = .39 - .75$ vs. $.21 - .38$). However, this discrepancy was less consistent between
500 years 2 and 3 in which plasma reliability was higher for only four of the seven analytes (plasma
501 $r_s = .10 - .54$; saliva $r_s = .09 - .36$) and for two-year reliability, for which saliva reliability was
502 higher for four of the analytes (plasma $r_s = .16 - .57$; saliva $r_s = .19 - .46$). Thus, although these
503 two studies suggest that serum measures of inflammation might be more stable than salivary
504 measures, there might be important protein-level differences in ideal measurement methods.
505 Also, the mouth is home to a complex microbiome that might introduce more confounding

506 factors compared to circulating blood (Giannobile et al., 2009). Thus, future research
507 establishing best practices for salivary methods of collection might find different estimates of
508 temporal stability.

509 Another popular way to quantify temporal stability is intra-class correlation coefficients
510 (ICCs), which assess the proportion of total variance (between-person + within-person) that is
511 attributable to between-person differences. Thus, higher ICCs indicate less relative within-person
512 variability and greater temporal stability. Conventionally, ICCs less than .40 are considered poor,
513 between .40 and .59 are considered fair, between .60 and .74 are considered good, and above .75
514 are considered excellent indicators of temporal stability (Cicchetti, 1993). An important
515 distinction between ICCs and retest reliability indexed by Pearson's r is that correlations
516 primarily reflect rank-order stability (i.e., an individual will have the same relative ranking in a
517 sample at Time 1 and Time 2), whereas ICCs reflect rank-order stability *and* mean-level changes
518 between time points. Thus, ICCs are a preferable measure when evaluating how stable a given
519 score is over time.

520 Continuing the discussion of inflammation, Shields and colleagues (2019) reported ICCs
521 (in their supplemental material) for seven different salivary inflammatory proteins (CRP, IL-6,
522 IL-8, IL-18, IL-1 β , TNF α , MCP). They report stability estimates for two different durations: 120
523 minutes apart during the same testing session ("short-term reliability") and an 18-month follow-
524 up ("long-term stability"). Importantly, testing stability of salivary analytes within the same
525 testing session can help identify how many measurements of these proteins would be necessary
526 to achieve a specific level of reliability. Short-term reliability ICCs ranged from .37 (for IL-8) to
527 .80 (for CRP). To reach a goal short-term reliability of $r = .80$ using the Spearman-Brown
528 prophecy formula, between one (CRP) and four measurements (IL-8 and IL-18) were needed.

529 The number of measurements needed to reach a goal short-term reliability indexed by ICCs was
530 not reported. ICCs were low for all 7 proteins at the 18-month follow-up (all ICCs < .28),
531 suggesting lower temporal stability of salivary inflammatory proteins using ICCs compared to
532 Pearson's *r*. Conceptually, this indicates that salivary inflammatory proteins might be more
533 stable in terms of their person-level rank-order than their actual value.

534 Given the relative expense of much biological psychiatry research (e.g., neuroimaging),
535 many studies are cross-sectional and prospective studies typically have small sample sizes. Thus,
536 meta-analyses pooling the results of multiple studies together have the potential to be very useful
537 in investigating the temporal stability of various measures. Elliot and colleagues (2020)
538 evaluated temporal stability of task-related fMRI measures in regions of interest (ROIs) using a
539 meta-analysis of 90 substudies (N = 1,008 and 1,146 ICC estimates). When selecting articles, the
540 authors noticed that several of the studies reported thresholded ICCs (i.e., only reported ICCs
541 above a threshold, comparable to only reporting effect sizes for results with $p < .05$). Due to
542 concerns this might inflate estimates of reliability, meta-analyses were conducted separately for
543 studies reporting unthresholded vs. thresholded ICCs. These concerns were supported by results
544 showing that the average ICC for unthresholded results (77 substudies) was poor (mean ICC =
545 .397; 95% CI, .330 - .460), whereas the average stability for tasks in thresholded substudies (13
546 substudies) was moderate (mean ICC = .705; 95% CI, .628 - .768). Further, a moderation
547 analysis including all substudies confirmed that the decision to report thresholded ICCs was
548 associated with significantly higher ICCs. Importantly, test-retest interval (the duration between
549 the two points of measurement) was not found to be a significant moderator of temporal stability,
550 although the authors do not provide information on the average test-retest interval or variability
551 in the intervals between studies. The authors highlight several methodological limitations of their

552 meta-analysis (e.g., different, potentially outdated scanners, different pre-processing and analysis
553 pipelines).

554 These results suggest lower than ideal temporal stability for the study of individual
555 differences. Importantly, the authors highlight that these tasks were created to robustly result in
556 group-level changes, not to assess between-person differences in these changes. Therefore, the
557 problem is not necessarily in the measures, but how researchers have extended their use to
558 research questions they were not built to address. It also is important to highlight that this study
559 only investigated ROIs. Similar analyses examining whole brain patterns might be more
560 temporally stable. Additionally, some common ROIs not included in this paper (e.g., left nucleus
561 accumbens and right anterior insula activity) have better temporal stability (e.g., ICC > .5) at
562 large intervals (> 2.5 years) during the monetary incentive delay task included in Elliot et al.
563 (2020) (Wu et al., 2014). In response to Elliot and colleagues (2020), Kragel et al. (2020, note
564 this is a pre-print that has not undergone peer review) describe nine recent studies demonstrating
565 strong short-term stability (i.e., less than five weeks) for task-based fMRI measures. They
566 conclude that studies aggregating information across multiple brain regions (rather than ROIs)
567 and/or aggregation across similar tasks, with larger samples, more data per participant (i.e., more
568 time in the scanner), and shorter retest intervals paint a more promising picture of temporal
569 stability for fMRI task measures than Elliot et al. (2020). Thus, further work is needed to identify
570 best practices for individual differences research using various fMRI measures.

571 Recall that measures taken across multiple time points for multiple people have three
572 sources of variability: between-person, within-person, and measurement error. Generalizability
573 theory (Shavelson and Webb, 1991) is an extension of these principles that estimates what
574 proportion of a single assessment is generalizable to other time points by separating variance due

575 to stable individual differences, measurement occasions, and the interaction between the two.
576 Results of generalizability analyses then can be used to inform the design of later studies with the
577 goal of achieving a desired reliability. Segerstrom and colleagues (2014) applied this theory to
578 investigate how many days of sampling would be needed to reliably characterize between-person
579 differences and within-person changes in three cortisol metrics: diurnal mean, diurnal slope, and
580 area under the curve (AUC) in two separate samples. Sample 1 consisted of young adults who
581 provided five cortisol samples per day, for three consecutive days, across five separate occasions
582 (mean time after previous occasion; Time 2: 44 days, Time 3: 57 days, Time 4: 36 days, Time 5:
583 29 days). Results indicated that three days were necessary for adequate reliability to facilitate
584 individual differences research (defined as $r = .60$ in this study) for the diurnal mean, four days
585 for the AUC, and 11 days for diurnal slope. Further, reliable measurement of within-person
586 changes would require three days of data for the mean, four for AUC, and eight for slope.
587 Correlations comparing slopes calculated with 2, 3, and 4 time points per day suggested that
588 collecting two samples per day (taken during the morning and evening) were excellent at
589 reproducing slope estimates using four samples ($r = .97$), suggesting that collecting more than
590 two samples per day does not substantively improve measurement. To evaluate whether these
591 results replicate in a demographically different sample, a second study was conducted in older
592 adults that resulted in comparable estimates. These results suggest that collecting two samples
593 per day for several days will provide more reliable estimates than collecting more samples, but
594 across fewer days.

595 **Temporal Specificity**

596 In addition to temporal stability, temporal specificity of effects is integral to advance
597 longitudinal research. To illustrate this, consider the following studies of inflammation as a risk

598 factor for depression. Miller and Cole (2012) reported that CRP predicted depression symptoms
599 at a six-month follow-up, but only in female adolescents exposed to childhood adversity.
600 Gimeno et al. (2009) found that CRP and IL-6 predicted depression symptoms 12 years in the
601 future. However, neither van den Biggelaar et al. (2007; five years of annual follow-ups) nor
602 Stewart, Rand, Muldoon, and Kamarck (2009; six-year follow-up) found significant associations
603 between IL-6 and future depression symptoms, but van der Biggelaar and colleagues found that
604 CRP predicted future depression. Further, Copeland and colleagues (2012) did not find that CRP
605 predicted future depression in a sample of adolescents with up to nine assessments over a 12-
606 year period. Although there might be (and likely are) many moderators influencing this
607 heterogeneity in results, time to follow-up is a plausible candidate that could inform design of
608 future, and interpretation of past, studies.

609 Moriarity and colleagues (2019) explored this possibility in a sample of 201 adolescents
610 with a baseline blood draw and a total of 582 assessments of depression symptoms (time to
611 follow-up ranged from .07 – 30.49 months). Using hierarchical linear models, they tested main
612 effects models of five inflammatory proteins on change in depression symptoms as well as five
613 exploratory models testing interactions between the five biomarkers, sex, and time to follow-up.
614 The only protein with a significant unconditional main effect was CRP; however, three of the
615 four remaining proteins demonstrated significant three-way interactions. Specifically, both IL-6
616 and TNF α had stronger, more positive associations with change in depression symptoms as time
617 to follow-up increased, but only for females (e.g., Figure 1). Conversely, IL-8 had a stronger
618 association with change in depression symptoms for males as time to follow-up increased, but
619 the association was negative. These results highlight how associations might not replicate
620 between samples with different demographic characteristics (e.g., sex) or different intervals

621 between assessments. This line of inquiry might be particularly important during adolescence,
622 which is both a time of elevated risk for first onset of many psychopathologies (e.g., depression;
623 Cummings et al., 2014) as well as a time of rapid social, biological, and psychological
624 development.

625 The rise in popularity of intensive longitudinal designs allows for a wealth of new
626 opportunities to investigate temporal specificity on a smaller time scale. For example, Graham-
627 England and colleagues (2018) measured serum levels of seven inflammatory proteins
628 (combined into an inflammatory composite) and CRP (analyzed individually) after a 14-day
629 ecological momentary assessment (EMA) protocol. Before starting the EMA protocol,
630 participants completed questions about recalled positive and negative affect “over the past
631 month”. Then, participants completed questions about experienced positive and negative affect
632 five times per day for 14 days leading up to the blood draw. Neither the inflammatory composite
633 nor CRP were significantly predicted by positive or negative affect “over the past month” or
634 aggregated positive or negative affect over the 14-day EMA protocol. However, when the affect
635 variables were separated by week, Week 2 (closest to the blood draw), but not Week 1, negative
636 affect significantly predicted the inflammatory composite variable. Exploratory analyses found
637 that the association between negative affect and inflammation consistently increased in strength
638 as the lag between measurements shortened. Thus, these two studies illustrate how it is possible
639 to leverage longitudinal studies of different time scales to identify whether risk factors for
640 psychopathology operate on a proximal or distal time scale, providing important insight to study
641 design and intervention efforts.

642 **Artificial Effect Size Deflation and Power**

643 As reviewed in the conceptual portion of this paper, all of the psychometric examples
644 reviewed thus far have implications for model performance; however, some researchers have
645 empirically tested the relationship between psychometrics and effect size/power in biological
646 psychiatry. For example, Hajcak and colleagues' (2017) paper on how internal consistency of
647 ERN changes as a function of trials completed in two groups of participants with, and without,
648 generalized anxiety disorder (reviewed above) also tested how between-group effect sizes were
649 related to internal consistency. Cohen's d increased almost parallel to increases in internal
650 consistency as the number of trials increased ($r = .94$). Given that two primary goals of
651 biological psychiatry are understanding i) group differences between those with and without
652 mental illness, and ii) the between-person variability in within-person effects contributing to
653 psychiatric risk, resilience, and treatment, this is noteworthy.

654 Simulation studies present a powerful option to evaluate the state of current measurement
655 practices. Segerstrom and Boggero (2020) used 212 study designs included as part of a meta-
656 analysis (Boggero et al., 2017) on the relationship between various psychosocial correlates and
657 cortisol awakening response to investigate the probability of misestimates using these data.
658 100,000 data sets were simulated for each study design with sample sizes and reliability
659 estimates extracted from the original studies. Boggero and colleagues (2020) found a meta-
660 analytic effect size of less than $r = 0.10$, which was used as the "true" effect size for the purposes
661 of the simulation study. Two types of misestimates were assessed: 1) sign errors (i.e., when the
662 association was negative, instead of positive like the meta-analytic effect); and 2) magnitude
663 errors (i.e., when the estimate was more than .10 away from the meta-analytic effect). Consistent
664 with literature reviewed above, more days of sampling in cortisol studies are associated with
665 higher reliability. More days of sampling (and, by extension, reliability) was, in turn, consistently

666 negatively correlated with both sign and magnitude errors in the simulations. Given that results
667 found that around 20% of all simulations resulted in sign errors, and nearly 40% in magnitude
668 errors, this study highlights increased cortisol sampling as a way to increase reliability and
669 overall study quality.

670 **The Promise of Biological Psychiatry**

671 Biological psychiatry has the potential to enhance both physical and mental health
672 through the investigation of the reciprocal associations between the body and mind. However,
673 this potential only can be realized with carefully crafted theory and rigorous methodology. Many
674 have argued that the field has fallen short of its promise to meaningfully impact psychiatric
675 classification, diagnosis, prevention, and treatment so far (Kapur et al., 2012; Miller, 2010;
676 Venkatasubramanian and Keshavan, 2016). One important reason for this may be that a lack of
677 sufficient attention to key measurement properties of biological variables has constrained the
678 utility of these data in statistical modeling, and thus, inference generation, despite rapid
679 technological advances allowing for more precise data acquisition in many biological subfields.

680 Although the physiometric characteristics covered in this review are far from exhaustive,
681 we would like to reiterate five steps that would improve biological psychiatry research: 1)
682 thoughtful investigation of the dimensionality of complex biological constructs in datasets
683 including multiple indicators of these constructs; 2) standardized reporting of internal
684 consistency when using aggregate measures; 3) careful consideration of the implications of
685 method-specific variance; 4) standardized reporting of temporal stability, preferably calculated
686 with the sample being analyzed or at least a reference to previous research with a similar time
687 frame; and 5) increased exploration into the temporal specificity of associations between
688 biological and behavioral phenomena. Further, it is imperative to keep in mind how the results of

689 these investigations might be contingent on other analytic choices (e.g., handling of extreme
690 values; Landau et al., 2019) and sample characteristics (e.g., sex; Moriarity et al., 2019).

691 A physiometric awakening in biological psychiatry would promote a wide array of
692 benefits to the field and those whom this work is intended to benefit. Projects uninformed by
693 basic measurement principles germane to their study methods risk inflating the noise-to-signal
694 ratio in statistical models. As a result, there is an increased risk for false-negatives and false-
695 positives, hindering the actual progress of the field as well as belief in its utility relative to the
696 associated costs. Further, many standardized effect sizes between biological and psychological
697 variables likely are biased downward due to less than ideal matching of measures to procedures
698 and method specific variance, weakening the appearance of their practical implications.
699 Thoughtful application of measurement principles can reduce error-related variability in future
700 studies via improvement of both study design and statistical modeling, resulting in improved
701 replicability of findings and less biased effect sizes.

702 Moreover, physiometric studies can provide guidance about which variables have the
703 most utility, under what research designs they operate well, and how to optimally model
704 constructs of interest. To illustrate this, consider designing a study of experienced negative affect
705 as a predictor of inflammatory and coagulatory markers in adolescents. Having read Nelson and
706 colleagues (2011), you know that aggregating variables containing overlapping variance can
707 accentuate the shared variance related to other variables, increasing power. You originally
708 considered the same panel of biomarkers as Egnot et al. (2018), but you decided not to assay and
709 analyze sICAM-1 and Lp(a) because neither loaded onto either of the two factors in their study.
710 This decision saves you money, enabling recruitment of more participants, hiring additional
711 staff, or purchasing other supplies. Additionally, because Engeland and colleagues (2018) found

712 that the association between negative affect and inflammation was stronger at shorter intervals,
713 you might plan a one-week EMA protocol rather than a two-week protocol, saving money, time,
714 and participant burden. However, instead of testing separate regressions for each day of negative
715 affect, you could improve statistical rigor of this comparison by testing for moderations by time
716 interval using multilevel models like Moriarity et al. (2019).

717 In addition to improving study design, thoughtful application of various statistical
718 approaches holds the potential to ameliorate psychometric issues in biological psychiatry. One
719 example is structural equation modeling (SEM), a powerful tool for reducing the impact of poor
720 reliability on statistical models. SEM allows the estimation of latent factors from the shared
721 variance between items, removing measurement error associated with individual observed
722 variables and accentuating shared variance between biomarkers of interest. However, SEM
723 models require larger samples than traditional models. Thus, multi-study collaborations might be
724 necessary to permit model testing for more expensive measures.

725 As described in Perkins et al. (2017), many physiological variables of interest are
726 associated with many different psychological constructs. Thus, when possible, researchers
727 should carefully consider whether building statistical models that can isolate portions of variance
728 relevant to one trait vs. another would be beneficial. However, we would like to underscore that
729 the suitability of various variance isolation techniques is context dependent. As described above,
730 variance removed from a variable always comes from the “true” and reliable variance, never
731 from error variance. Thus, difference scores or predictors with variance partialled out for
732 covariates are almost always less reliable and have a lower signal-to-error ratio (Lynam et al.,
733 2006). This is amplified when the predictors are highly correlated (Thomas and Zumbo, 2012).
734 Finally, it also is critical to remember that difference scores (or predictors with variance

735 partialled out in multiple regression) are conceptually different than the raw variables. These
736 interpretive concerns are more extreme with more heterogenous (lower internal consistency)
737 measures, because it is more likely that the variance removed might only be associated with a
738 subset of the components of the original variable.

739 Additionally, most of this article has discussed psychometric work anchored in classical
740 test theory. Future work could utilize generalizability theory, an extension of classical test theory
741 described above in the review of Segerstrom et al. (2014). Alternatively, item response theory
742 (IRT) estimates reliability for varying levels of a continuum rather than the entire range of a
743 measure. Typically, IRT requires binary or polytomous indicators, but continuous response
744 models (CRM) are an extension of IRT models that allow for continuous variables (Samejima,
745 1973). Psychometric research utilizing these approaches might lead to useful insight for how to
746 best collect and model biological data.

747 Increasing the efficiency of study design and statistical modeling will improve the ability
748 to accurately detect associations and their effect sizes. These advancements have the potential to
749 smooth the transition from basic research to the improvement of interventions and policy via
750 increasing confidence in results and the ability to gauge their utility. Importantly, with lower
751 rates of false positives, there is a reduced chance that ineffective biological interventions may be
752 explored that have little to no real-world utility.

753 Fortunately, as reviewed above, some researchers are working to arm the rest of the field
754 with this crucial information. As more psychometric work is published, the value of
755 comprehensive reviews of this literature increases. Recently, Segerstrom (2020) and Gloger et al.
756 (2020) published reviews of salivary and serum biomarker psychometrics, respectively, but many

757 more topics would benefit from a focused psychometric review (e.g., neuroimaging, ERP, heart
758 rate variability).

759 However, it is critical to admonish the dangers of treating particular levels of
760 psychometric characteristics as benchmarks to hit, without careful consideration of what they
761 mean in relation to the constructs being studied. Several methodologists have warned that
762 primarily focusing on creating measures with high internal consistency can result in the removal
763 of items/components that contribute to lower internal consistency, but would help capture the
764 true breadth of the construct of interest (Clark and Watson, 2019; Cronbach and Meehl, 1955).
765 This sacrifices construct validity for higher internal consistency and faux-unidimensionality.
766 Further, internal consistency increases as a function of the number of components included in its
767 calculation, potentially resulting in larger, but not better, measures. Additionally, although there
768 are many contexts in which high temporal stability can be beneficial, it is critical to avoid
769 overvaluing components of larger constructs (e.g., brain regions for neuroimaging studies) with
770 higher reliability. Rather, there should be reciprocal interplay between methodology and theory.

771 Creating a solid psychometric foundation for biological psychiatry is not without
772 obstacles. First and foremost, biological variables often are more expensive to measure than
773 psychological variables, some of which can be measured via self-report questionnaires
774 administered online from the comfort of participants' homes. Measurement research and
775 construct validation are, by their nature, iterative processes, amplifying the associated cost of this
776 work. However, it is crucial to appreciate that good psychometric research is an investment; it
777 will result in increased statistical power and better study design in the future, saving money and
778 time. This requires investment both on the part of researchers as well as funding agencies.
779 Fortunately, there is a lot of important work that can be done with existing data sets. Any study

780 with repeated measures of a variable can estimate its temporal stability. Any study using an
781 aggregate measure can assess the internal consistency of its components. In fact, there are many
782 publicly available data sets that offer great opportunities for physiometric research (e.g., the
783 Human Connectome Project; Van Essen et al., 2013).

784 Finally, this work can, at times, be statistically intensive and conceptually abstract. One
785 of the strengths of biological psychiatry is that, by nature, it is an interdisciplinary pursuit with
786 experts along the biology—psychology spectrum. Collaboration with statisticians and
787 measurement specialists can serve as a catalyst for the efficient, high-quality research that is
788 needed for biological psychiatry to reach its full academic, clinical, and policy-informing
789 potential.

790 **Conclusion**

791 It is important to end on a clarification that the issues highlighted in this article should not
792 be received with apprehension or pessimism. Rather, it is an invitation to ask new questions of
793 the data collected to help the field of biological psychiatry realize its potential. Biological
794 psychiatry has been criticized for falling short of its considerable promise in advancing
795 knowledge about the interplay between biology and behavior in ways that will translate to
796 substantive impact on clinical outcomes (Kapur et al., 2012; Miller, 2010; Venkatasubramanian
797 and Keshavan, 2016). One addressable barrier to meaningfully advancing biological psychiatry
798 is an understanding and appreciation of measurement characteristics for biological variables. By
799 leveraging existing data sets and prioritizing funding for physiometric research, it is possible to
800 advance current methods to allow for more informative and replicable studies that will provide
801 greater clarity into what areas of research offer the greatest promise to make meaningful impacts
802 on mental health, and how best to integrate them into intervention efforts.

803 Acknowledgements: Thank you to Drs. Michelle Bryne, Thomas Olino, Lauren Ellman, and

804 David Smith for providing feedback on drafts of this article.

805



807

808 *Figure 1. Temporal specificity of Log IL-6 predicting change in depression symptoms by sex.*
 809 *This figure was first presented in Moriarity et al. (2019). Note: IL = interleukin, CDI =*
 810 *Children's Depression Inventory*

811

812

References

813 Boggero, I.A., Hostinar, C.E., Haak, E.A., Murphy, M.L.M., Segerstrom, S.C., 2017.

814 Psychosocial functioning and the cortisol awakening response: Meta-analysis, P-curve
 815 analysis, and evaluation of the evidential value in existing studies. *Biol. Psychol.* 129, 207–
 816 230. <https://doi.org/10.1016/j.biopsycho.2017.08.058>

817 Bradford, D.E., Starr, M.J., Shackman, A.J., Curtin, J.J., 2015. Empirically based comparisons of
 818 the reliability and validity of common quantification approaches for eyeblink startle
 819 potentiation in humans. *Psychophysiology* 52, 1669–1681.

- 820 <https://doi.org/10.1111/psyp.12545>
- 821 Cicchetti, D. V, 1993. Guidelines, criteria, and rules of thumb for evaluating normed and
822 standardized assessment instruments in psychology. *Psychol. Assess.* 6, 284–290.
823 <https://doi.org/10.1037/1040-3590.6.4.284>
- 824 Clark, L.A., Watson, D., 2019. Constructing validity: New developments in creating objective
825 measuring instruments. *Psychol. Assess.* 31, 1412–1427.
826 <https://doi.org/10.1037/pas0000626>
- 827 Clark, L.A., Watson, D., 1995. Constructing validity: Basic issues in objective scale
828 development. *Psychol. Assess.* 7, 309–319.
- 829 Copeland, W.E., Shanahan, L., Worthman, C., Angold, A., Costello, E.J., 2012. Cumulative
830 depression episodes predict later C-reactive protein levels: A prospective analysis. *Biol.*
831 *Psychiatry* 71, 15–21. <https://doi.org/10.1016/j.biopsych.2011.09.023>
- 832 Cortina, J.M., 1993. What is coefficient alpha? An examination of theory and applications. *J.*
833 *Appl. Psychol.* 78, 98–104. <https://doi.org/10.1037//0021-9010.78.1.98>
- 834 Cronbach, L.J., Meehl, P.E., 1955. Construct validity in psychological tests. *Psychol. Bull.* 52,
835 281–302.
- 836 Cummings, C., Caporino, N., Kendall, P.C., 2014. Comorbidity of anxiety and depression in
837 children and adolescents: 20 Years After. *Psychol. Bull.* 140, 816–845.
838 <https://doi.org/10.1037/a0034733>
- 839 Cuthbert, B.N., Kozak, M.J., 2013. Constructing constructs for psychopathology: The NIMH
840 research domain criteria. *J. Abnorm. Psychol.* 122, 928–937.
841 <https://doi.org/10.1037/a0034028>
- 842 Davidshofer, K.R., Murphy, C.O., 2005. *Psychological testing: principles and applications.*

- 843 Edwards, J.R., 2001. Ten difference score myths. *Organ. Res. Methods* 4, 265–287.
844 <https://doi.org/10.1177/109442810143005>
- 845 Egnot, N.S., Barinas-Mitchell, E., Criqui, M.H., Allison, M.A., Ix, J.H., Jenny, N.S., Wassel,
846 C.L., 2018. An exploratory factor analysis of inflammatory and coagulation markers
847 associated with femoral artery atherosclerosis in the San Diego Population Study. *Thromb.*
848 *Res.* 164, 9–14. <https://doi.org/10.1016/j.thromres.2018.02.003>
- 849 Elliott, M.L., Knodt, A.R., Ireland, D., Morris, M.L., Poulton, R., Ramrakha, S., Sison, M.L.,
850 Moffitt, T.E., Caspi, A., Hariri, A.R., 2020. What is the test-retest reliability of common
851 task-fMRI measures? New empirical evidence and a meta-analysis. *Psychol. Sci.* 681–700.
852 <https://doi.org/10.1101/681700>
- 853 Giannobile, W. V, Beikler, T., Kinney, J.S., Ramseier, C.A., Wong, D.T., 2009. Saliva as a
854 diagnostic tool for periodontal disease: current state and future directions. *Periodontol* 2000
855 52–64. <https://doi.org/10.1111/j.1600-0757.2008.00288.x.Saliva>
- 856 Gloger, E.M., Smith, G.T., Segerstrom, S.C., 2020. Stress physiology and physiometrics. *Handb.*
857 *Res. Methods Heal. Psychol.*
- 858 Graham-Engeland, J.E., Sin, N.L., Smyth, J.M., Jones, D.R., Knight, E.L., Sliwinski, M.J.,
859 Almeida, D.M., Katz, M.J., Lipton, R.B., Engeland, C.G., 2018. Negative and positive
860 affect as predictors of inflammation: Timing matters. *Brain. Behav. Immun.* 74, 222–230.
861 <https://doi.org/10.1016/j.bbi.2018.09.011>
- 862 Guadagnoli, E., Velicer, W.F., 1988. Relation of Sample Size to the Stability of Component
863 Patterns. *Psychol. Bull.* 103, 265–275. <https://doi.org/10.1037/0033-2909.103.2.265>
- 864 Hajcak, G., Meyer, A., Kotov, R., 2017. Psychometrics and the neuroscience of individual
865 differences: Internal consistency limits between-subjects effects. *J. Abnorm. Psychol.* 126,

- 866 823–834. <https://doi.org/10.1037/abn0000274>
- 867 Hajcak, G., Patrick, C.J., 2015. Situating psychophysiological science within the Research
868 Domain Criteria (RDoC) framework. *Int. J. Psychophysiol.* 98, 223–226.
869 <https://doi.org/10.1016/j.ijpsycho.2015.11.001>
- 870 Holiga, Š., Sambataro, F., Luzy, C., Greig, G., Sarkar, N., Renken, R.J., Marsman, J.B.C.,
871 Schobel, S.A., Bertolino, A., Dukart, J., 2018. Test-retest reliability of task-based and
872 resting-state blood oxygen level dependence and cerebral blood flow measures. *PLoS One*
873 13, 1–16. <https://doi.org/10.1371/journal.pone.0206583>
- 874 Kakeda, S., Watanabe, K., Nguyen, H., Katsuki, A., Sugimoto, K., Igata, N., Abe, O.,
875 Yoshimura, R., Korogi, Y., 2020. An independent component analysis reveals brain
876 structural networks related to TNF- α in drug-naïve, first-episode major depressive disorder:
877 a source-based morphometric study. *Transl. Psychiatry* 10. [https://doi.org/10.1038/s41398-](https://doi.org/10.1038/s41398-020-00873-8)
878 [020-00873-8](https://doi.org/10.1038/s41398-020-00873-8)
- 879 Kapur, S., Phillips, A.G., Insel, T.R., 2012. Why has it taken so long for biological psychiatry to
880 develop clinical tests and what to do about it. *Mol. Psychiatry* 17, 1174–1179.
881 <https://doi.org/10.1038/mp.2012.105>
- 882 Kaye, J.T., Bradford, D.E., Curtin, J.J., 2016. Psychometric properties of startle and corrugator
883 response in NPU, affective picture viewing, and resting state tasks. *Psychophysiology* 53,
884 1241–1255. <https://doi.org/10.1111/psyp.12663>
- 885 Kragel, P.A., Han, X., Kraynak, T.E., Gianaros, P.J., Wagner, T.D., 2020. fMRI can be highly
886 reliable, but it depends on what you measure. *PsyArXiv*.
- 887 Landau, E.R., Trinder, J., Simmons, J.G., Raniti, M., Blake, M., Waloszek, J.M., Blake, L.,
888 Schwartz, O., Murray, G., Allen, N.B., Byrne, M.L., 2019. Salivary C-reactive protein

- 889 among at-risk adolescents: A methods investigation of out of range immunoassay data.
890 *Psychoneuroendocrinology* 99, 104–111. <https://doi.org/10.1016/j.psyneuen.2018.08.035>
- 891 Loevinger, J., 1957. Objective tests as instruments of psychological theory. *Psychol. Rep.* 3,
892 635–694.
- 893 Luking, K.R., Nelson, B.D., Infantolino, Z.P., Sauder, C.L., Hajcak, G., 2017. Internal
894 consistency of functional magnetic resonance imaging and electroencephalography
895 measures of reward in late childhood and early adolescence. *Biol. Psychiatry Cogn.*
896 *Neurosci. Neuroimaging* 2, 289–297. <https://doi.org/10.1016/j.bpsc.2016.12.004>
- 897 Lynam, D.R., Hoyle, R.H., Newman, J.P., 2006. The Perils of Partialling Cautionary Tales From
898 Aggression and Psychopathy. *Assessment* 13, 328–341.
899 <https://doi.org/10.1177/1073191106290562>
- 900 Miller, G.A., 2010. Mistreating psychology in the decades of the brain. *Perspect. Psychol. Sci.* 5,
901 716–743. <https://doi.org/10.1038/jid.2014.371>
- 902 Miller, G.E., Cole, S.W., 2012. Clustering of depression and inflammation in adolescents
903 previously exposed to childhood adversity. *Biol. Psychiatry* 72, 34–40.
904 <https://doi.org/10.1016/j.biopsych.2012.02.034>.
- 905 Moriarity, D.P., Mac Giollabhui, N., Ellman, L.M., Klugman, J., Coe, C.L., Abramson, L.Y.,
906 Alloy, L.B., 2019. Inflammatory proteins predict change in depressive symptoms in male
907 and female adolescents. *Clin. Psychol. Sci.* 7, 754–767.
908 <https://doi.org/10.1177/2167702619826586>
- 909 Moriarity, D.P., McArthur, B.A., Ellman, L.M., Coe, C.L., Abramson, L.Y., Alloy, L.B., 2018.
910 Immunocognitive model of depression secondary to anxiety in adolescents. *J. Youth*
911 *Adolesc.* 47, 2625–2636. <https://doi.org/10.1007/s10964-018-0905-7>

- 912 Moriarity, D.P., Ng, T., Curley, E., McArthur, B.A., Ellman, L.M., Coe, C.L., Abramson, L.Y.,
913 Alloy, L.B., 2020a. Reward sensitivity, cognitive response style, and inflammatory response
914 to an acute stressor in adolescents. *J. Youth Adolesc.* 49, 2149–2159.
- 915 Moriarity, D.P., Ng, T., Titone, M.K., Chat, I.K., Nusslock, R., Miller, G.E., Alloy, L.B., 2020b.
916 Reward sensitivity and ruminative response styles for positive and negative affect interact to
917 predict inflammation and mood symptomatology. *Behav. Ther.* 51, 829–842.
918 <https://doi.org/10.1016/j.beth.2019.11.007>
- 919 Muscatell, K.A., Moieni, M., Inagaki, T.K., Dutcher, J.M., Jevtic, I., Breen, E.C., Irwin, M.R.,
920 Eisenberger, N.I., 2016. Exposure to an inflammatory challenge enhances neural sensitivity
921 to negative and positive social feedback. *Brain. Behav. Immun.* 57, 21–29.
922 <https://doi.org/10.1016/j.bbi.2016.03.022>
- 923 Nelson, L.D., Patrick, C.J., Bernat, E.M., 2011. Operationalizing proneness to externalizing
924 psychopathology as a multivariate psychophysiological phenotype. *Psychophysiology* 48,
925 64–72. <https://doi.org/10.1111/j.1469-8986.2010.01047.x>
- 926 Ng, T.H., Alloy, L.B., Smith, D. V, 2019. Meta-analysis of reward processing in Major
927 Depressive Disorder: Distinct abnormalities within the reward circuit? *Transl. Psychiatry* 9,
928 2–10.
- 929 Out, D., Hall, R.J., Granger, D.A., Page, G.G., Woods, S.J., 2012. Assessing salivary C-reactive
930 protein: Longitudinal associations with systemic inflammation and cardiovascular disease
931 risk in women exposed to intimate partner violence. *Brain. Behav. Immun.* 26, 543–551.
932 <https://doi.org/10.1016/j.bbi.2012.01.019>
- 933 Patrick, C.J., Iacono, W.G., Venables, N.C., 2019. Incorporating neurophysiological measures
934 into clinical assessments: Fundamental challenges and a strategy for addressing them.

- 935 Psychol. Assess. 31, 1512–1529. <https://doi.org/10.1037/pas0000713>
- 936 Patrick, C.J., Venables, N.C., Yancey, J.R., Hicks, B.M., Nelson, L.D., Kramer, M.D., 2013. A
937 construct-network approach to bridging diagnostic and physiological domains: Application
938 to assessment of externalizing psychopathology. *J. Abnorm. Psychol.* 122, 902–916.
939 <https://doi.org/10.1037/a0032807>
- 940 Perkins, E.R., Yancey, J.R., Drislane, L.E., Venables, N.C., Balsis, S., Patrick, C.J., 2017.
941 Methodological issues in the use of individual brain measures to index trait liabilities: The
942 example of noise-probe P3. *Int. J. Psychophysiol.* 111, 145–155.
943 <https://doi.org/10.1016/j.ijpsycho.2016.11.012>
- 944 Plichta, M.M., Schwarz, A.J., Grimm, O., Morgen, K., Mier, D., Haddad, L., Gerdes, A.B.M.,
945 Sauer, C., Tost, H., Esslinger, C., Colman, P., Wilson, F., Kirsch, P., Meyer-Lindenberg, A.,
946 2012. Test-retest reliability of evoked BOLD signals from a cognitive-emotive fMRI test
947 battery. *Neuroimage* 60, 1746–1758. <https://doi.org/10.1016/j.neuroimage.2012.01.129>
- 948 Riis, J.L., Granger, D.A., Dipietro, J.A., Bandeen-Roche, K., Johnson, S.B., 2015. Salivary
949 cytokines as a minimally-invasive measure of immune functioning in young children:
950 Correlates of individual differences and sensitivity to laboratory stress. *Dev. Psychobiol.* 57,
951 153–167. <https://doi.org/10.1002/dev.21271>
- 952 Riis, J.L., Out, D., Dorn, L.D., Beal, S.J., Denson, L.A., Pabst, S., Jaedicke, K., Granger, D.A.,
953 2014. Salivary cytokines in healthy adolescent girls: Intercorrelations, stability, and
954 associations with serum cytokines, age, and pubertal stage. *Dev. Psychobiol.* 56, 797–811.
955 <https://doi.org/10.1002/dev.21149>
- 956 Rogosa, D.R., Willett, J.B., 1983. Demonstrating the reliability of the difference score in the
957 measurement of change. *J. Educ. Meas.* 20, 335–343.

- 958 Samejima, F., 1973. Homogenous case of the continuous response model. *Psychometrika* 38,
959 203–2019.
- 960 Segerstrom, S.C., 2020. Psychometrics in Salivary Bioscience. *Int. J. Behav. Med.* 27, 262–266.
- 961 Segerstrom, S.C., Boggero, I.A., 2020. Expected Estimation Errors in Studies of the Cortisol
962 Awakening Response: A Simulation. *Psychosom. Med.* 82, 751–756.
963 <https://doi.org/10.1097/PSY.0000000000000850>
- 964 Segerstrom, S.C., Boggero, I.A., Smith, G.T., Sephton, S.E., 2014. Variability and reliability of
965 diurnal cortisol in younger and older adults: Implications for design decisions.
966 *Psychoneuroendocrinology* 49, 299–309. <https://doi.org/10.1016/j.psyneuen.2014.07.022>
- 967 Segerstrom, S.C., Smith, G.T., 2012. Methods, variance, and error in psychoneuroimmunology
968 research: The good, the bad, and the ugly, in: Segerstrom, S.C. (Ed.), *Oxford Handbook of*
969 *Psychoneuroimmunology*. Oxford U Press, New York, NY, pp. 421–432.
- 970 Shavelson, R.J., Webb, N.M., 1991. *Generalizability Theory: A Primer*. Sage,.
- 971 Shields, G.S., Slavich, G.M., Perlman, G., Klein, D.N., Kotov, R., 2019. The short-term
972 reliability and long-term stability of salivary immune markers. *Brain. Behav. Immun.* 81,
973 650–654. <https://doi.org/10.1016/j.bbi.2019.06.007>
- 974 Slavich, G.M., Irwin, M.R., 2014. From stress to inflammation and major depressive disorder: a
975 social signal transduction theory of depression. *Psychol. Bull.* 140, 774–815.
976 <https://doi.org/10.1037/a0035302>
- 977 Stewart, J.C., Rand, K.L., Muldoon, M.F., Kamarck, T.W., 2009. A prospective evaluation of the
978 directionality of the depression-inflammation relationship. *Brain. Behav. Immun.* 23, 936–
979 944. <https://doi.org/10.1016/j.bbi.2009.04.011>
- 980 Tabachnick, B.G., Fidell, L.S., 2013. *Using multivariate statistics*, Sixth. ed. Pearson, Boston,

- 981 MA.
- 982 Thomas, D.R., Zumbo, B.D., 2012. Difference scores from the point of view of reliability and
983 repeated-measures ANOVA: In defense of difference scores for data analysis. *Educ.*
984 *Psychol. Meas.* 72, 37–43. <https://doi.org/10.1177/0013164411409929>
- 985 van den Biggelaar, A.H.J., Gussekloo, J., de Craen, A.J.M., Frölich, M., Stek, M.L., van der
986 Mast, R.C., Westendorp, R.G.J., 2007. Inflammation and interleukin-1 signaling network
987 contribute to depressive symptoms but not cognitive decline in old age. *Exp. Gerontol.* 42,
988 693–701. <https://doi.org/10.1016/j.exger.2007.01.011>
- 989 Van Essen, D.C., Smith, S.M., Barch, D.M., Behrens, T.E., Yacoub, E., Ugurbil, K., Consortium,
990 W.-M.H., 2013. The WU-Minn Human Connectome Project: An overview. *Neuroimage* 80,
991 62–79. <https://doi.org/10.1038/jid.2014.371>
- 992 Venables, N.C., Foell, J., Yancey, J.R., Kane, M.J., Engle, R.W., Patrick, C.J., 2018. Quantifying
993 inhibitory control as externalizing proneness: A cross-domain model. *Clin. Psychol. Sci.* 6,
994 561–580. <https://doi.org/10.1177/2167702618757690>
- 995 Venkatasubramanian, G., Keshavan, M.S., 2016. Biomarkers in psychiatry – A critique. *Ann.*
996 *Neurosci.* 23, 3–5. <https://doi.org/10.1159/000443549>
- 997 Williams, L.M., 2016. Precision psychiatry: A neural circuit taxonomy for depression and
998 anxiety. *The Lancet Psychiatry* 3, 472–480. [https://doi.org/10.1016/S2215-0366\(15\)00579-](https://doi.org/10.1016/S2215-0366(15)00579-9)
999 9.Precision
- 1000 Wu, C.C., Samanez-Larkin, G.R., Katovich, K., Knutson, B., 2014. Affective traits link to
1001 reliable neural markers of incentive anticipation. *Neuroimage* 2014 84, 279–289.
1002 https://doi.org/10.1007/978-3-319-55511-9_5
- 1003