

¹ Comparing the Temporal Stability and Convergent Validity of Risk
² Preference Measures: A Meta-Analytic Approach

³ Alexandra Bagaïni¹, Yunrui Liu¹, Madlaina Kapoor¹, Gayoung Son²,
⁴ Paul-Christian Bürkner³, Loreen Tisdall¹, and Rui Mata¹

⁵¹Center for Cognitive and Decision Sciences, University of Basel

⁶²Department of Psychology, University of Bern

⁷³Department of Statistics, TU Dortmund University

⁸ Version date: July 11, 2023

⁹ **Author Note**

¹⁰ Correspondence concerning this article should be addressed to Alexandra

¹¹ Bagaïni, Center for Cognitive and Decision Sciences, Faculty of Psychology, University

¹² of Basel, Missionsstrasse 60-62, 4055 Basel, Switzerland, E-mail:

¹³ alexandra.bagaini@unibas.ch

14

Abstract

15 Understanding whether risk preference represents a stable, coherent trait is central to
16 efforts aimed at explaining, predicting, and preventing risk-related behaviours. We help
17 characterise the nature of the construct by adopting a meta-analytic approach to
18 summarise the temporal stability and convergent validity of over 300 risk preference
19 measures (51 samples, 29 panels, >500.000 respondents). Our findings reveal significant
20 heterogeneity across and within measure categories (propensity, frequency, behaviour),
21 domains (e.g., investment, occupational, alcohol consumption), and sample
22 characteristics (e.g., age). Specifically, while self-reported propensity and frequency
23 measures of risk preference show a higher degree of stability relative to behavioural
24 measures, these patterns are moderated by domain and age. Crucially, an analysis of
25 convergent validity reveals a low agreement across measures, questioning the idea that
26 they capture the same underlying trait. Our results raise concerns about the coherence
27 and measurement of the risk preference construct.

28 **keywords:** risk preference, test-retest, age differences, life span

29 **Comparing the Temporal Stability and Convergent Validity of Risk**

30 **Preference Measures: A Meta-Analytic Approach**

31 Risk permeates all domains and stages of life. Consequently, preferences towards risk
32 may fundamentally shape individuals' health, wealth, and happiness. Risk
33 preference—an umbrella term used to reflect the individual's appetite for risk (Mata
34 et al., 2018; Schonberg et al., 2011)—not only has been related to personal decisions
35 (e.g., timing of marriage and parenthood; Schmidt, 2008), but may also be used as an
36 indicator to match individuals with products, services, and suitable careers (Breivik
37 et al., 2019; Caliendo et al., 2014; Financial Services Authority, 2011; Jin et al., 2020).
38 Because of its broad significance, risk preference is central to many theories and
39 applications in the behavioural sciences (Barseghyan et al., 2018; Steinberg, 2013).

40 Despite the importance of the construct, there is considerable discussion about
41 its central characteristics, including whether risk preference represents a stable,
42 coherent trait or rather a contextual and/or domain-specific disposition
43 (Schildberg-Hörisch, 2018; Schonberg et al., 2011; Stigler & Becker, 1977). One crucial
44 source of the confusion surrounding the nature of risk preference is the many ways it
45 has been operationalised. Specifically, the assessment of risk preference spans three
46 measurement traditions that can be classified into three broad categories of measures:
47 propensity, frequency, and behavioural measures, which, in turn, can differ in the
48 domain (e.g., health, financial) and mode of assessment (e.g., ratings, choices; cf. Table
49 1). Crucially, past work suggests that different measures do not speak with one voice
50 (e.g., Frey et al., 2017; Schonberg et al., 2011). As a consequence, resolving the debate
51 about whether risk preference shares two central characteristics of a trait, namely
52 stability and coherence, cannot be done without acknowledging the central role of
53 measurement. Standing in the way of clarity, however, is the piecemeal approach taken
54 in much past research, whereby single or few measures are adopted in any given study,
55 making it difficult to obtain an overview across measures. Our work aims to help
56 resolve this issue by taking a meta-analytic approach to investigate both the temporal
57 stability and convergent validity of extant measures of risk preference.

58 A first focus of our work is to quantify the temporal stability of risk preference
59 measures. This goal aligns with a key objective of discerning the sources of stability and
60 change in human psychology and behaviour (Fraley & Roberts, 2005), and mirrors
61 existing research into other traits (Anusic & Schimmack, 2016; Bleidorn et al., 2022;
62 Elliott et al., 2020; Enkavi et al., 2019). Although some studies in economics and
63 psychology have already probed the temporal stability of risk preference (e.g., Chuang
64 & Schechter, 2015; Mata et al., 2018; Schildberg-Hörisch, 2018), there is a lack of a
65 comprehensive comparison across measures with at least three significant gaps in
66 existing research. First, previous work found higher stability for propensity and
67 frequency measures than behavioural measures (Frey et al., 2017; Mata et al., 2018) but
68 did not fully consider the role of domain (e.g., health, financial; Mata et al., 2018),
69 leading to an oversimplified picture of the stability of measures. Second, there is little
70 consideration of how the stability of different psychological constructs varies across the
71 lifespan (Anusic & Schimmack, 2016; Bleidorn et al., 2022). Early life and young
72 adulthood, which are marked by significant biological, cognitive, and social changes,
73 usually show lower rank-order stability (Seifert et al., 2022) but past syntheses of the
74 stability of risk preference did not account for age differences (e.g., Chuang & Schechter,
75 2015; Mata et al., 2018). Third, previous research has not employed theoretically
76 grounded models to analyse temporal stability patterns across different categories of
77 measures, domains, or populations, hindering comparison with other constructs, such as
78 major personality traits, that have been studied using formal models (Anusic &
79 Schimmack, 2016).

80 A second focus of our work is to quantify the convergent validity of risk
81 preference measures. The issue of convergence is central to the goal of mapping
82 theoretical constructs to specific measures and many efforts in the behavioural sciences
83 aim to empirically estimate these links (Duckworth & Kern, 2011; Eisenberg et al.,
84 2019; Frey et al., 2017). The issue is also of practical importance because many studies
85 investigating predictors or correlates of risk preference, for example, neuroimaging and
86 genome-wide association studies (Karlsson Linnér et al., 2019; Karlsson Linnér et al.,

87 2021), are often able to use only a single or limited set of measures to capture risk
88 preference. To the extent that different measures do not speak with one voice, however,
89 these should not be used interchangeably and need to be carefully selected to match the
90 construct of interest. Previous work on risk preference reports a relatively low
91 convergence between measures, albeit propensity and frequency measures may exhibit
92 moderate convergent validity among themselves, whereas behavioural measures show
93 comparatively low convergent validity, in terms of both observable behaviour and
94 computational parameters (Frey et al., 2017; Pedroni et al., 2017). We note three key
95 gaps in extant work on the convergent validity of risk preference measures. First,
96 studies typically employ only a few different measures, thus limiting the extent to which
97 a comprehensive assessment of convergence between many measures can be performed
98 in a single study. Second, the adoption of few measures in single studies often implies
99 that the influence of measure (e.g., category, domain) or respondent characteristics
100 (e.g., age) cannot be ascertained as moderating variables that can impact the
101 convergence of measures. Third, and finally, studies have not been able to assess the
102 extent to which low convergent validity is a direct result of poor reliability of specific
103 measures (Dang et al., 2020; Strickland & Johnson, 2021).

104 The present study tackles these outstanding gaps by examining the temporal
105 stability and convergent validity of a comprehensive set of risk preference measures. For
106 this purpose, we conducted a systematic search for longitudinal data sets comprising
107 many different measures of risk preference, including propensity, frequency, and
108 behavioural measures. The curated database represents a large data trove comprising
109 29 longitudinal panels, split into 51 different samples, capturing over 300 different
110 measures of risk preference. To further enhance the comprehensiveness of this newly
111 curated data, we conducted an extensive categorisation of measures (e.g., category,
112 domain) and associated respondents (e.g., age, gender).

113 Equipped with these data, we conducted a number of analyses to gain an
114 overview of the temporal stability and convergent validity of risk preference measures.
115 First, to comprehensively examine temporal stability, we performed a variance

116 decomposition analysis that provides a picture of the amount of variance that can be
117 accounted for in temporal stability by measure, respondent, and panel-related
118 predictors. We also adopted a formal modelling approach using the meta-analytic
119 stability and change model (MASC; Anusic & Schimmack, 2016) to capture the
120 temporal stability of risk preference measures while distinguishing between domains
121 (e.g., investment, gambling, smoking, ethical). We further employed MASC to
122 re-analyse longitudinal panel data for other pertinent psychological constructs,
123 including personality and affect, thus providing a direct comparison between our results
124 and those for other major psychological constructs. Second, to comprehensively
125 examine convergent validity, we performed variance decomposition analysis to quantify
126 to what extent measure, respondent, and panel-related predictors account for the
127 heterogeneity observed between inter-correlations. Crucially, because it has been
128 suggested that the reliability of individual measures creates boundary conditions for
129 their convergence (Dang et al., 2020), we consider measure reliability as a
130 measure-related predictor in these analyses. We further report meta-analytic syntheses
131 of the empirical relation across measures both between and within category and domain
132 pairs. All in all, we hope that by clarifying the two central characteristics of measures
133 of risk preference—temporal stability and convergent validity—we can contribute to
134 improving its measurement, describing its life course patterns, and, ultimately, its
135 utility as a construct in the behavioural sciences.

136

137 Results

138 Overview of the Longitudinal Data

139 We used a systematic approach to identify a comprehensive set of longitudinal
140 samples suitable for estimating the temporal stability and convergent validity of risk
141 preference measures. Figure 1 depicts the flow of steps starting from the identification
142 of panels, screening for eligibility, and, finally, the data available for the temporal
143 stability and convergent validity analyses. Please note that we distinguish between

144 panels and samples because if panels included data from several countries, we treated
145 these as separate samples to avoid confounding within-and cross-country differences. As
146 per our inclusion criteria, all the samples had to contain at least one propensity
147 measure. This criterion was implemented to enable comparisons between propensity
148 measures, the most prevalent category in the literature on risk preference, to other
149 categories (i.e., frequency, behaviour) as well as to similar measurement approaches in
150 personality research (cf. Anusic & Schimmack, 2016). From the initial pool of 101
151 panels (157 samples) identified in our search, we were able to include 29 panels (51
152 samples) that allowed computing test-retest information for at least one measure of risk
153 preference, and 26 panels (45 samples) that allowed computing inter-correlations
154 between two or more measures of risk preference. Finally, for each risk preference
155 measure, sample, age group, and gender, we calculated test-retest correlations between
156 all measurement wave combinations for temporal stability analyses, and all possible
157 inter-correlations between measures for convergent validity analysis. This process
158 yielded over 72,000 test-retest correlation coefficients for temporal stability (Figure 2A)
159 and over 61,000 inter-correlations for convergent validity analyses (Figure 2B). As a
160 whole, the dataset covers over 300 different measures of risk preference spanning three
161 measure categories (i.e., propensity, frequency, behaviour).

162 Informed by previous work that has distinguished between different domains of
163 risk, we conducted an extensive categorisation of measures to distinguish between 14
164 different domains (e.g., general health, financial, recreational, driving), thus allowing a
165 fine-grained classification sorely lacking in the risk preference literature. Crucially, this
166 categorisation makes clear that there are important differences across, and also gaps
167 between, the domains investigated in each category. As can be seen in Figure 2C, while
168 propensity measures capture the majority, albeit not all, of the domains detected in our
169 data (9 out of 14), frequency measures capture a large but different subset of these (8
170 out of 14). In turn, behavioural measures capture only a small minority of
171 finance-related domains, such as investment and gambling (4 out of 14). This imbalance
172 is ultimately due to the different traditions spanning the psychology, economics, and

173 public health literature that have investigated risk preference using different
174 measurement categories. In what follows, we provide a fine-grained comparison of the
175 measures' temporal stability.

176 **Temporal Stability**

177 We first obtained an overview of the temporal stability data by visualising the
178 number of measures by category and retest interval as well as a breakdown of the
179 test-retest correlations by measure category (propensity, frequency, behaviour; see
180 Figure S1A). We should note that there are substantial differences in the amount of
181 data concerning different categories, with most measures being classified as propensity
182 or frequency measures and only a minority as behavioural measures. The
183 under-representation and overall shorter test-retest intervals for behavioural measures
184 observed in our sample is a product of there being overall fewer samples that have
185 included (repeatedly) such measures in their assessment batteries, likely due to the
186 additional burden of deploying behavioural measures which typically require extensive
187 instructions, multiple choices, and, potentially, incentivisation. Figure S1B provides a
188 first impression of the distributions of retest correlations across time and measure
189 categories that conveys considerable heterogeneity between measures that we explore
190 quantitatively in more detail below.

191 ***Variance Decomposition of Test-Retest Correlations***

192 Our first main question concerns the relative contribution of measure,
193 respondent, and panel characteristics in accounting for patterns of temporal stability in
194 different measures of risk preference. For this purpose, we adopted a Shapley
195 decomposition approach, a method that estimates the average marginal contribution of
196 different predictors to the variance in an outcome of interest (Grömping, 2007), in our
197 case, the test-retest correlations of risk preference measures. We were particularly
198 interested in the role of specific measure- and respondent-related predictors that have
199 been either hypothesised or shown to account for some variance in temporal stability in
200 past work on risk preference (e.g. Frey et al., 2017; Josef et al., 2016) or other

201 psychological constructs (e.g. Anusic & Schimmack, 2016). For measure-related
202 predictors, we focused on the category (i.e., propensity, frequency, behaviour), domain
203 (e.g., general health, recreational), the scale type (e.g., ordinal, open-ended), and length
204 of the test-retest interval (e.g., 6 months, 1 year, 5 years). For respondent-related
205 predictors, we considered age group, gender, and number of respondents. Finally, we
206 also included panel as a predictor to capture the role of unobserved panel characteristics
207 (e.g., quality of data collection or data entry) that can influence test-retest reliability.

208 We first conducted an omnibus analysis to assess to what extent measure,
209 respondent, and panel predictors explained differences across all test-retest correlations.
210 Altogether, a model considering all predictors captures 49.7% of the observed variance.
211 As can be seen in Figure 3A, we find that a large portion of the variance could be
212 explained by measure-related predictors, domain (13.7%), category (4.3%), retest
213 interval (6.8%), and scale type (0.5%). In turn, we find that some of the variance could
214 be explained by respondent-related predictors, in particular, age (5.2%). Finally, panel
215 captured a large portion of the variance (18.7%), suggesting that there are a number of
216 (unobserved) panel characteristics that also contribute to systematic differences in the
217 observed temporal stability of measures.

218 Given our focus on comparing measure categories, we further explored the
219 differences between the contribution of these predictors to propensity, frequency, and
220 behavioural measures separately. The models conducted separately by measurement
221 category explained 23.7%, 46.6%, and 16.6% of the total variance for propensity,
222 frequency, and behavioural measures, respectively. The results of this analysis are
223 depicted in Figure 3B. There are four main insights that can be drawn from the
224 comparison between measure categories. First, domain explained a significant
225 percentage of the variance for frequency (12.5%) relative to propensity (1.3%) and
226 behavioural (5.6%) measures. This suggests considerable heterogeneity within some
227 categories as a function of domain, in particular, in the frequency category, something
228 we will explore in more detail when analysing the temporal trajectories by domain
229 below. Second, retest interval contributed to more explanatory power for propensity

(5.2%) and frequency (6.9%) measures relative to behavioural measures (1.0%), suggesting that the temporal patterns are less pronounced for the latter. Third, concerning respondent-related predictors, we find that age explained a significant percentage of the variance in the test-retest correlations, but, in particular, for frequency (8.4%) relative to propensity (2.3%) and behavioural (0.8%) measures. These results seem to indicate some specificity regarding the effects of age by measure category. Fourth, as in the omnibus analysis, a number of (unobserved) panel characteristics seem to contribute to systematic differences between panels, albeit this effect is most pronounced for frequency measures. In what follows, we explore these results in more detail by adopting a formal modelling approach that distinguished between the different measure categories and domains.

241 Meta-Analyses of Temporal Stability

We used the Meta-Analytic Stability and Change model (MASC; Anusic & Schimmack, 2016) to capture the trajectory of test-retest correlations across measures of risk preference and compare these to other psychological constructs. MASC uses three parameters to represent different properties of temporal trajectories: reliability (proportion of between-person variance excluding random error), change (proportion of variance that is subject to changing factors), and stability of change (the rate at which change occurs over time). In our work, we adopted a sampling-based Bayesian estimation procedure to obtain full posterior distributions for each model parameter for specific measure categories (propensity, frequency, behaviour) and domains (e.g., recreational, general health, smoking, investment).

Figure 4 shows model predictions for the trajectory of test-retest correlations separately for the three measure categories and distinguishing further between domains (e.g., recreational, general health, smoking, investment) and respondent groups (age groups, gender). Figures 4A-C show the distributions of the predictions for each of the model parameters, while Figures 4D-I show the corresponding trajectories in test-retest correlations as a function of retest interval for different age groups (panels D, F, H), as well as the (equivalent) age trajectories as a function of different retest intervals (panels

259 E, G, I). While the trajectories in test-retest correlations as a function of retest interval
260 are particularly helpful to compare to similar trajectories found for other psychological
261 constructs (Anusic & Schimmack, 2016), the trajectories by age for different retest
262 intervals help visualise a potential inverted U-shape function across the life span in
263 patterns of reliability found in past work using propensity measures of risk preference
264 (Josef et al., 2016) and major personality traits (Bleidorn et al., 2022).

265 We find a ranking in the overlapping reliability estimates for the three measure
266 categories, with the highest reliability found for propensity measures ($M: 0.51$, 95%
267 HDI: [0.42, 0.61]), followed by frequency measures ($M: 0.47$, 95% HDI: [0.33, 0.63]), and
268 behavioural measures ($M: 0.30$, 95% HDI: [0.20, 0.40]). Crucially, relative to propensity
269 and behavioural measures, the reliability of frequency measures varies widely by
270 domain, with a wide range evident between the highest reliability for smoking ($M: 0.84$,
271 95% HDI: [0.78, 0.90]) and the lowest for the ethical domain ($M: 0.11$, 95% HDI: [0.04,
272 0.18]). In comparison, the ranges found for propensity measures, spanning from ethical
273 ($M: 0.64$, 95% HDI: [0.36, 0.91]) to occupational ($M: 0.41$, 95% HDI: [0.32, 0.49]), and
274 behavioural measures, spanning from investment ($M: 0.36$, 95% HDI: [0.24, 0.49]) to
275 insurance ($M: 0.26$, 95% HDI: [0.17, 0.36]), are considerably smaller. Concerning the
276 patterns of change and associated stability, the different measure categories and
277 domains appear comparable and seem to mimic those found in the temporal stability
278 literature characterised by steep changes yet some long-term stability (Anusic &
279 Schimmack, 2016; Fraley & Roberts, 2005).

280 Concerning age-related patterns, we note clear trends for propensity and
281 frequency measures but not behavioural ones. Specifically, as can be seen in Figure 4C,
282 when considering longer retest intervals (>2 years) for propensity measures, and
283 consistent with previous work (Josef et al., 2016), we note an inverse U-shape
284 association between retest-correlations and age, indicating that temporal stability peaks
285 in middle-age. Also, this pattern is observed for most domains covered by propensity
286 measures (Figures S7-S9). For frequency, the overall pattern observed in Figure 4G is
287 more mixed but we should note that this appears due to heterogeneity between domains

288 within the frequency category, as we observe an inverse-U shape with age for both
289 alcohol consumption and smoking domains. In turn, the driving, ethical, and sexual
290 intercourse domains do not show the same pattern (Figures S10-S11). For behavioural
291 measures, as seen in Figure 4I, we do not observe noticeable association between
292 temporal stability and age, and this is reflected across the individual domains (Figure
293 S12). Concerning gender, we did not identify any substantial differences, suggesting
294 males and females show comparable stability trajectories across the board.

295 Finally, we assessed where risk preference stands within the consistency
296 hierarchy of psychological constructs (Conley, 1984), by comparing the temporal
297 stability of risk preference to that of personality, life satisfaction, self-esteem, and affect
298 using data of Anusic and Schimmack (2016). Our results obtained using a Bayesian
299 framework largely replicate those of Anusic and Schimmack (cf. Figure S13) but allow
300 us to compare directly the estimates for different constructs using the same modelling
301 approach. Our reanalysis show highest reliability for personality traits ($M: 0.73$, 95%
302 HDI: [0.68, 0.77]), followed by self-esteem ($M: 0.62$, 95% HDI: [0.54, 0.71]), life
303 satisfaction ($M: 0.60$, 95% HDI: [0.55, 0.64]), and affect ($M: 0.56$, 95% HDI: [0.50,
304 0.61]). In line with the results for risk preference given above, this suggests that the
305 average stability of risk preference as captured by propensity and frequency measures,
306 is, on average, lower than that of major psychological constructs albeit it overlaps with
307 that for affect. In turn, the reliability of behavioural measures is lower than any of the
308 four constructs, suggesting a qualitative difference between this category and the
309 constructs considered. Of course, as suggested above, for frequency measures, some
310 domains show considerably higher/lower levels of stability; consequently, while
311 frequency measures in the smoking and alcohol domains rival the temporal stability of
312 major personality traits, others, like ethical and driving, show some of the lowest
313 reliability estimates observed, suggesting these do not have the same stable quality.

314 All in all, the results on temporal stability support the notion that different risk
315 preference measures show markedly different temporal stability signatures. In what
316 follows, we explore further differences between measures by evaluating their

317 inter-correlations.

318 **Convergent Validity**

319 ***Variance Decomposition of Correlations Between Measures***

320 We first obtained an overview of the convergent validity data by visualising the
321 distributions of inter-correlations of measures separately for different measure pairs

322 (Figure S14). The resulting pattern speaks to the large heterogeneity in correlations
323 between measures as well as possible differences between and within measure categories.

324 We used variance decomposition to provide a quantitative summary of correlations as a
325 function of several measure and respondent-related characteristics, as well as panel.

326 Specifically, concerning measure characteristics we included dummy-coded predictors to
327 code for the matching (e.g., propensity-propensity) or mismatching category (e.g.,
328 propensity-frequency), domain, and scale type. Further, using the results from the
329 temporal stability analyses above, we computed the average reliability of each pair of
330 measures and included this in our predictors to assess the extent to which measures'
331 reliability contribute to their convergence.

332 The variance decomposition analysis suggests that a model considering all
333 predictors captures 26.6% of the variance in inter-correlations. More substantively, as
334 shown in Figure 5, the variance decomposition analysis suggests that category and
335 domain play a considerable role: More than half of the explained variance was
336 accounted for by whether or not the pair of measures matched in terms of category
337 (7.5%) and domain (10.0%). In turn, we find that measure reliability accounted for less
338 than 1% of the variance, thus indicating little support for the idea that poor reliability
339 of risk preference measures is the main driver of their (lack of) convergence. Finally,
340 respondent-related effects offer little to no contribution, while panel characteristics seem
341 to account for some amount of variance, suggesting that unobserved panel
342 characteristics capture relevant, systematic variance in the correlation between
343 measures. In sum, the variance decomposition analysis suggests that measure
344 characteristics, specifically, category and domain, capture important aspects of measure

345 convergence. In what follows, we provide a more detailed overview of the role of these
346 factors by providing a meta-analytic correlation matrix across pairs of measures that
347 distinguishes between category and domain.

348 ***Meta-Analyses of Convergent Validity***

349 We adopted a meta-analytic approach to map out the convergent validity of risk
350 preference measures across categories and domains. For that purpose, we conducted
351 separate meta-analyses at different levels of aggregation. A meta-analysis across all
352 available inter-correlations, suggests an average meta-analytic inter-correlation of .16,
353 95% HDI: [0.13, 0.18]. However, this value hides considerable heterogeneity. As can be
354 seen in Figure 6A, across pairs of categories and domains, we observe a large range of
355 inter-correlations, from around -.2 to circa .8. The meta-analytic correlation matrix also
356 shows evidence of overall higher average correlations along the diagonal, signalling that
357 matching both category and domain leads to typically higher inter-correlations relative
358 to matching only across domains or categories. Importantly, as can be seen in Figure
359 6B, when considering aggregation at the category level, there is a clear ranking of the
360 average inter-correlations within category, with this being highest for propensity ($M =$
361 0.41, 95% HDI: [0.40, 0.43]), followed by frequency ($M = 0.20$, 95% HDI: [0.18, 0.22]),
362 and behavioural measures ($M = 0.19$, 95% HDI: [0.15, 0.23]). Finally, and more
363 importantly, there is evidence of little convergence between categories, with
364 cross-category meta-analytic correlations being around or smaller than 0.1.

365 All in all, considering jointly the results on both temporal stability and
366 convergent validity, one is left with the impression that different risk preference
367 measures can show very different psychometric signatures, including patterns of
368 temporal stability and convergent validity, supporting the notion that measurement
369 issues plague clarity concerning the nature of the construct.

371

Discussion

372 Our aim was to contribute to the ongoing debate about whether risk preference
373 represents a stable and coherent trait by adopting a meta-analytic approach to assess
374 the temporal stability and convergent validity of a large set of risk preference measures.
375 We curated an extensive collection of previously un(der)utilised longitudinal samples,
376 providing data for over 300 unique measures rooted in measurement traditions that are
377 aligned with the adoption of three broad categories of measures—propensity, frequency,
378 and behavioural measures—and covering various life domains (e.g., driving, alcohol,
379 smoking, social, ethical, recreational, occupational, gambling). Our work provides the
380 first encompassing meta-analytic syntheses of the trajectories of temporal stability and
381 convergent validity across these major measure categories while accounting for central
382 measure (e.g., domain) and respondent (e.g., age) characteristics. Crucially, we do so by
383 adopting a formal model of temporal stability that allows comparing the temporal
384 stability trajectories results both between measures of risk preference as well as with
385 other major psychological constructs.

386 Our analyses of the temporal stability of risk preference measures suggest some
387 average differences in the reliability of measures from the three measurement traditions.

388 Overall, propensity measures exhibited the highest average reliability, followed by
389 frequency measures, while behavioural measures showed the lowest average reliability.

390 Crucially, we observe considerable overlap between categories and substantial
391 heterogeneity within categories as a function of domain. Particularly, and most
392 profoundly for frequency measures, reliability varies widely between domains, with
393 smoking showing the highest reliability, while others, such as the ethical domain (i.e.,
394 violent or delinquent behaviour) showing the lowest. Concerning respondent

395 characteristics, we find that age affects the temporal stability patterns found for
396 propensity and frequency measures, but not behavioural ones. Specifically, test-retest
397 correlations were lower in younger and older age groups compared with middle-aged
398 groups for propensity measures, a common pattern found for other personality
399 constructs (Anusic & Schimmack, 2016; Bleidorn et al., 2022; Seifert et al., 2022). For

frequency measures, age patterns are more heterogeneous across domains. For example, while the smoking and alcohol consumption domains show increasing stability across adulthood followed by some decline in old age, the ethical and sexual domains shows patterns consistent with decreased stability in adolescence and young adulthood. We note that this heterogeneity maps onto distinct pathways for age-specific versus lifelong trajectories of these different behaviours (Ahun et al., 2020; Moffitt, 2018). The heterogeneity in the temporal patterns of risk preference measures poses a problem for its comparison with that of other psychological constructs. Nevertheless, one conclusion that emerges is that propensity measures show somewhat lower test-retest stability but similar age-related (inverted-U) trends compared to that of major personality traits. Frequency measures are more heterogeneous and, therefore, not easily compared as a whole, some domains, like smoking and alcohol consumption, approach the stability of personality constructs and show similar age patterns. In contrast, other domains captured by frequency measures, such as driving and ethical domains, show very low stability and most change occurring in adolescence and young adulthood. Behavioural measures show considerably lower stability compared to the other categories (i.e., propensity, frequency) or psychological constructs (e.g., self-esteem) and do not seem to capture any life span trends. As a whole, these results suggest that different measurement traditions are characterised by distinct temporal and age-related trajectories, emphasising the important role of measurement, domain, and age in moderating the patterns of temporal stability concerning risk preference measures.

Our analyses of convergent validity showed that, overall, convergence between risk preference measures was low, albeit highlighting substantial heterogeneity between measure categories. Convergence was highest for propensity measures, while frequency and behavioural ones showed lower convergence, somewhat matching previous results from individual studies (Eisenberg et al., 2019; Frey et al., 2017). One should note that frequency measures covered a considerably larger set of domains spanning health, occupational, and gambling domains compared to behavioural measures that shared a focus on financial domains (i.e., investment, gambling, insurance), which may present a

429 confound when estimating differences between these two categories. Unfortunately,
430 frequency measures did not cover these financial domains well, making a direct
431 comparison between frequency and behavioural measures impossible. Crucially, we
432 found that relatively little variance in convergence between measures was explained by
433 their average reliability, suggesting that there may be something more fundamental
434 about measure characteristics that contributes to their lack of convergence. To sum up,
435 somewhat mirroring the temporal stability analyses, the results on convergence suggest
436 that different measurement traditions do not speak with one voice but, rather, show
437 unique patterns by category and, particularly for frequency measures, are largely
438 moderated by domain. In contrast with the temporal analyses, however, age did not
439 seem to be a strong determinant of measure convergence. These results suggest the
440 different measures cannot be used interchangeably to capture individual differences in
441 risk preference and call into question the coherence of the risk preference construct.

442 Before we address the implications of our findings for our understanding of risk
443 preference and its measurement, several limitations of our study should be noted
444 concerning our 1) search and inclusion criteria, our 2) coding of predictor variables, and
445 other 3) analytical choices. First, despite conducting an extensive search for panels,
446 there may be additional ones that were missed by our independent research effort.
447 Exploring yet more panels could lead to the discovery of additional measures that could
448 further improve the scope of our findings. Further, our focus on comparison between
449 measurement traditions as well as other psychological constructs led us to consider
450 samples only if they included at least one propensity measure, which likely contributed
451 to over-representation of this category relative to others (e.g., behavioural), as well as
452 the domains represented across categories. A promising solution involves pursuing even
453 more comprehensive efforts, for example by leveraging crowd sourcing or coordinated
454 analyses across multiple research teams. By tapping into the collective expertise and
455 resources of a broader community, one could make the efforts of mapping risk preference
456 measures yet more exhaustive. Second, to assess the role of a set of theoretically
457 relevant predictors for temporal stability and convergence, we meticulously coded

458 relevant information about the measures (e.g., category, domain, test-retest interval,
459 scale type, pair type for convergence) as well as the respondents (e.g., age, gender).
460 While we recognise the value of additional information (e.g., measure incentivisation,
461 respondents' socioeconomic status), it proved challenging to obtain sufficient data to
462 allow including more fine-grained comparisons in our analysis or ensure comparability
463 across samples. Another coding issue concerns our use of panel as a predictor, which
464 could have been broken down further (e.g., main data collection mode, language) but
465 proved unfeasible to model in our framework. In light of these constraints, our coding
466 scheme and analyses were geared towards including maximally informative predictors
467 while ensuring computational feasibility. Perhaps future efforts including additional
468 data can help resolve the role of additional moderating factors. Third, our workflow
469 required making a number of analytical choices, including the binning of age groups, or
470 the selection of statistical metrics and model priors in our Bayesian framework.
471 Whenever possible, we made principled decisions informed by past work. To deal with
472 this issue, we conducted multiverse analyses to assess the robustness of our results
473 whenever possible. Finally, given the complexity of the data curation process we did not
474 pre-register our analysis but we make our data and scripts publicly available which we
475 hope will allow the research community to collaborate on future efforts to examine the
476 psychometric characteristics of risk preference measures.

477 Our findings provide a new empirical overview on the status of many extant risk
478 preference measures. We would like to point out four main implications of these
479 findings for current theorising and empirical research on risk preference. First, our
480 results indicate we need to invest new energy into developing theoretical frameworks
481 that help us make sense of the observed convergence as well as divergence across
482 measures. One factor leading to the gap between measures we have documented may
483 arise from fundamentally different concepts of risk taking being captured by different
484 measures (e.g., Bran & Vaidis, 2020). Specifically, propensity measures aim to capture
485 individuals' attitudes towards risk, while frequency measures aim to capture actual
486 risky behaviour, which will often be a product of both individuals' appetite for risk as

well as other considerations, including the opportunity to engage in these risks. In this sense, the gap observed between propensity and frequency measures could be interpreted as a special case of the classic intention-behaviour gap. This explanation, however, leaves the lower reliability of behavioural measures and their low convergence with propensity and frequency measures largely unresolved. Some researchers have pointed out current limitations of behavioural measures that can contribute to this state of affairs. For example, behavioural measures may require many trials to obtain reliable estimates of the underlying latent trait, something that is more easily and naturally accomplished by integrating behavioural episodes from memory (e.g., Haines et al., 2020). One other more general factor contributing to the gap between measures concerns the levels of granularity adopted. For example, while propensity measures are typically general, covering a broad domain (e.g., health) and time span ("in general"), frequency measures are more specific (e.g., "number of cigarettes") and constrained in time (e.g., "in the last 30 days"), and behavioural measures could perhaps be thought as yet more specific (e.g., about specific types of monetary choices). The lack of a direct match in levels of granularity can contribute to lower reliability because individuals may think of different aspects when answering general questions or even provide different answers depending on the cues that happen to come to mind on any given occasion (Arslan et al., 2020; Steiner et al., 2021). We would like to note that the effort to understand how these factors contribute to gaps between measures should not be seen as a simply methodological one. Clarifying the conceptual and empirical relations between constructs and how these are operationalised is central to achieving conceptual clarity in the behavioural sciences (e.g., Bringmann et al., 2022). Consequently, it should also be seen as part of a larger effort to integrate risk preference in the larger context of psychological constructs and associated ontologies (Eisenberg et al., 2019; Norris et al., 2019).

Second, in line with the focus on theory development, our results emphasise the need to understand the temporal stability of risk preference as a function of life span changes in a heterogeneous set of contexts or domains. Many extant theories make

516 valuable contributions to explaining the complex nature of stability and change in
517 personality traits (Möttus et al., 2019) and behaviours, such as antisocial (Moffitt,
518 2018) or health behaviours (Ahun et al., 2020). Transactional models appear
519 particularly promising in that they emphasise the interplay between individual
520 characteristics and environmental factors in determining phenotypic change across the
521 lifespan (Möttus et al., 2019). Our results suggest that such transactional models could
522 be helpful in reconciling the idea of stable individual risk preferences with differential
523 patterns across domains that are shaped by changing affordances and goals (Ravert
524 et al., 2019) as well as individuals' life experiences (Beck & Jackson, 2022).

525 Third, from a more methodological perspective, our findings suggest it is
526 important to streamline and replenish our methodological resources by focusing on
527 principled measure validation and development. Regarding validation, we should strive
528 for more comprehensive comparisons of existing measures. This can be achieved
529 through meta-analytic research, similar to our current approach, as well as primary
530 studies that explore previously overlooked measure categories, domains, and their
531 combinations (Richmond-Rakerd et al., 2020). We also need to engage more actively
532 with particular behaviours, and conduct targeted explorations of domains using
533 multiple measures across different categories (cf. risky driving, Das & Ahmed, 2022).
534 Regarding measure development, recent technological development suggests that there
535 are new forms of measurement on the horizon that could help anchor measures of risk
536 preference in more real-world experience, for example, through the use of virtual reality
537 (Roberts et al., 2021), or text-based analysis facilitated by large language models (Wulff
538 & Mata, 2022), as well as biology, through the use of advanced imaging methods that
539 track structural aspects of neural processing of reward (Tisdall et al., 2022).

540 Fourth, and finally, we need to combine the improvements awaiting us in the
541 development and validation of both theories and measures to focus on prediction. Three
542 centuries ago, the topic of risk preference emerged from Daniel Bernoulli's interest in
543 solving practical problems, aiming to use mathematical formalisation to help
544 understand how individuals make consequential decisions regarding gambling, financial

investment, and insurance (Bernoulli, 1954). Principled prediction requires a good understanding of the anticipated mechanisms as well as an informed selection of measures (on the side of both outcomes and predictors). Future work will need to integrate objective measures in the domains of health (e.g., inflammation markers, visits to the emergency department), investment (e.g., stock portfolios), and ethics (e.g., arrest records, number of speeding tickets) to assess the predictive value of different risk preference measures. We hope this focus on prediction will ultimately fuel a better understanding of what risk preference means for whom and at what stage in their life thus buttressing the utility of the construct for predicting important life outcomes and ultimately improving individuals' health, wealth, and happiness.

555

Methods

556 **Identification of Samples**

557 We used a systematic method to find a comprehensive set of longitudinal data
558 that include measures of risk preference (Figure 1). We started by identifying
559 longitudinal panels by 1) performing searches on general-purpose search engines, survey
560 listings, and data repositories (i.e., Google Database, Gateway to Global Aging Data,
561 Gesis, IZA, ICPSR, CNEF, UK Data service) using relevant terms (e.g., "risk
562 preference", "risk aversion", "risk attitude", "take risks", "survey", "panel", "longitudinal";
563 cf. Table S1 for a list of our search terms), 2) consulting past literature for references to
564 longitudinal panels or studies that have estimated the temporal stability of
565 psychological constructs (i.e., Anusic and Schimmack, 2016; Chuang and Schechter,
566 2015; Graham et al., 2020; Mata et al., 2018; Orth, 2018), and 3) informal requests to
567 colleagues for suggestions concerning panels or specific studies. This search led to
568 identifying 101 longitudinal panels (157 samples; Table S2). It is important to note that
569 we differentiate between panels and samples, such that samples have their origin in a
570 panel. For example, if a panel (e.g., SHARE) included data from multiple countries
571 (e.g., SHARE-Switzerland, SHARE-Germany, SHARE-Belgium), we treated the latter
572 as distinct samples to prevent confusion between differences within and across countries.
573 To determine the relevance of each of the 157 samples for our analyses, we adopted a
574 set of screening criteria (Table S3). In brief, we included a sample in our analyses if it
575 1) was publicly available, 2) included data on at least one consistently formatted
576 propensity measure of risk preference with responses from the same respondents across
577 at least two time points, and 3) included data on the gender and age of the respondents.
578 This procedure led to the creation of a comprehensive data trove comprising 51 samples
579 from 29 longitudinal panels (Table S4). For each sample, we included data that was
580 available as of May 2023.

581 Categorisation of Measures

582 To further add to the comprehensiveness of the newly curated data set, we
583 conducted a categorisation of each risk preference measure. The following measure
584 characteristics are particularly relevant to our analysis: measure category (e.g.,
585 propensity, frequency, behaviour), domain (e.g., investment, general health, social,
586 recreational), and scale type (e.g., open or closed questions). Table S5 presents
587 descriptions of risk preference measures that are representative of the variety of
588 measures included in the samples used for our analyses. With regards to the domains
589 captured by different risk preference measures, we included measures covering as many
590 domains as possible, that is, we did not exclude measures in pre-specified domains.
591 Further, we adopted a bottom-up, data-driven approach mostly to distinguish between
592 domains. We felt this approach was best suited for our purpose, as this allowed us to 1)
593 scope extant work and systematically identify the domains most commonly assessed in
594 the risk preference literature, and 2) provide the most comprehensive assessment to
595 date of temporal stability and convergent validity while systematically investigating the
596 role of domain at a high level of granularity. Overall, we identified 14 domains: *alcohol,*
597 *driving, drugs, ethical, gambling, general health, general risk, insurance, investment,*
598 *occupational, recreational, sexual intercourse, smoking, and social.* Our labelling scheme
599 has considerable overlap with terminology commonly used to group contexts or
600 situations within which risk taking can occur, albeit it makes fine-grained distinctions
601 within domains, such as distinguishing between smoking or alcohol consumption from a
602 more general health domain. We provide additional detail concerning an assessment of
603 measure characteristics in the Supplementary Information.

604 Temporal Stability

605 In what follows, we give an overview of steps involved in computing test-retest
606 correlations, conducting variance decomposition of test-retest correlations, and the
607 modelling of temporal stability using the Meta-Analytic Stability and Change model
608 (MASC; Anusic and Schimmack (2016)). We provide additional information concerning

609 each step in the Supplementary Information.

610 ***Computing Correlations***

611 To compute test-retest correlations, we followed a similar approach as Anusic
612 and Schimmack (2016) and Enkavi et al. (2019). For each panel we included the data
613 from all the respondents, regardless of whether or not they provided responses on all
614 measurement waves. Within each sample and for each risk preference measure, we
615 calculated test-retest correlation coefficients for each possible wave combination. For
616 example, for a sample with Waves 1, 2 and 3, we calculated three sets of test-retest
617 correlations: between Wave 1 and 2, between Wave 2 and 3, and between Wave 1 and 3.
618 More importantly, we computed test-retest correlations separately for females and males
619 as well as for respondents of different age groups (defined by binning age at the time of
620 the first data collection point into 10-year bins). Robustness checks (cf. Enkavi et al.,
621 2019) suggested high correlations between test-retest correlations computed using
622 different metrics and using (non)transformed data (Figures S2 and S3). Consequently,
623 we report results using Pearson's r correlation coefficients for non-transformed data. To
624 obtain reasonable estimates, test-retest correlations calculated from less than 30
625 responses were excluded from the main analyses. Further, we restricted the data set to
626 correlations with a retest interval of up to 20 years. This resulted in a set of 72,963
627 test-retest correlations.

628 ***Variance Decomposition***

629 To estimate the proportion of variance in the 72,963 test-retest correlations that
630 could be explained by measure-related, respondent-related, and panel predictor
631 variables, we used Shapley Decomposition (Grömping, 2007). First, we obtained the
632 adjusted R^2 value from each of the 2^8 subsets of linear regression models (2^7 regression
633 models for the category-specific variance decomposition). Second, we estimated the
634 variance explained by each predictor by calculating the weighted average change in
635 adjusted R^2 resulting from its inclusion in the model. Third, using 100 re-sampled data
636 sets we generated 100 bootstrapped estimates for each prediction and from which we

637 computed bootstrapped confidence intervals (e.g., Sharapov et al., 2021).

638 ***Meta-Analytic Stability and Change Model***

Model Description. The Meta-Analytic Stability and Change model (MASC) is a non-linear model introduced by Anusic and Schimmack (2016) to capture the trajectory of test-retest correlations over time. In this model, the test-retest correlation r_{t2-t1} at a specific *time* interval is a function of the proportion of reliable between-person variance, *rel*, the proportion of this reliable variance explained by changing factors, *change*, and the stability of these changing factors over time (per year), *stabch*. This is formalised as

$$r_{t2-t1} = \text{rel} \times (\text{change} \times (\text{stabch}^{\text{time}} - 1) + 1)$$

639 Figure S4A describes the model, and Figure S4B illustrates how different model
640 parameterisations alter the shape of the curve.

641 **Aggregation of Test-Retest Correlations.** To minimise potential
642 convergence issues that arise from meta-analysing 72,963 test-retest correlations using
643 MASC, we aggregated the test-retest correlations. We obtained these aggregates by first
644 grouping the test-retest correlations by sample, measure category, domain, and retest
645 interval, as well as respondent gender and age group. We then calculated the average
646 test-retest correlation for each of these groupings, using inverse-variance weighting and
647 accounting for the dependency between these correlations. This resulted in 7,996
648 aggregated correlations.

649 **Bayesian Model Specification.** We set up the MASC model such that for
650 each parameter (i.e., *rel*, *change* and *stabch*) we accounted for the effects of domain,
651 linear age, quadratic age and gender, as well as the interaction between linear and
652 quadratic age with domain. In addition, we included *sample* as a random factor for the
653 *rel* parameter. Importantly, to obtain meta-analytic estimates we additionally specified
654 the (aggregate) standard errors of each correlation. Lastly, to best capture
655 domain-specific effects within each category, we fitted the model separately for each

656 measure category using their respective aggregated retest correlations and aggregated
657 standard errors.

658 To estimate the parameters of this non-linear hierarchical model we used a
659 Bayesian approach to account for the large differences between sample sizes and retest
660 intervals encountered in such a large set of data sources. We specified weakly
661 informative priors on the model parameters and hierarchical standard deviations so as
662 to include values reported previously in the literature (e.g., Anusic and Schimmack,
663 2016; Frey et al., 2017; Mata et al., 2018).

664 Analyses were conducted in the R statistical environment (R Core Team, 2021),
665 using the *brms* package (Bürkner, 2017, 2018, 2021) which provides a high-level
666 interface to fit hierarchical models in Stan (Carpenter et al., 2017).

667 **Construct Comparison.** To compare the temporal stability and reliability of
668 risk preference to that of other psychological constructs (e.g., personality), we
669 re-analysed the set of correlations included in Anusic and Schimmack (2016) using a
670 Bayesian estimation procedure and set of MASC model specifications to maximise
671 comparability to the analyses conducted for risk preference.

672 **Convergent Validity**

673 In what follows, we give an overview of the main steps involved in computing
674 inter-correlations between measures, variance decomposition of inter-correlations, and
675 the meta-analyses of convergent validity. We provide additional information concerning
676 each step in the Supplementary Information.

677 ***Computing Correlations***

678 For the assessment of the convergence of risk preference measures, we started
679 with the set of samples used to assess the temporal stability of risk preference, but
680 selected only those samples that included two or more measures of risk preference
681 within at least one wave, and for which the same set of respondents had provided
682 answers. As a result, we conducted our convergent validity analyses for 45 samples from
683 26 panels (Figure 1), retaining the same three measure categories and 14 domains used

684 in the temporal stability analyses. First, for each sample, we computed the correlations
685 between every possible pair of measures within the same data collection point. We
686 computed these correlations separately for females and males as well as respondents of
687 different ages. We excluded inter-correlations computed from the responses of less than
688 30 respondents. This resulted in a data set of 61,644 inter-correlations. Robustness
689 checks (cf. Enkavi et al., 2019) suggested high correlations between inter-correlations
690 computed using different metrics and using (non)transformed data (Figures S5 and S6).
691 Here, we report results using Spearman's rho correlation coefficients for
692 non-transformed data and which were based on a minimum of 30 responses.

693 To avoid model convergence issues when running the meta-analysis, we grouped
694 the inter-correlations (e.g., by type of pair, age, gender, panel), and then aggregated the
695 inter-correlations within these groupings, resulting in 5,038 aggregated
696 inter-correlations.

697 ***Variance Decomposition***

698 To estimate the proportion of variance in inter-correlations between risk
699 preference measures that could be explained by measure-related, respondent-related,
700 and panel predictor variables, we used Shapley Decomposition (Grömping, 2007). We
701 followed the same approach used for the test-retest correlations obtaining the adjusted
702 R^2 value from each of the (2^8) models, estimating the variance explained by each
703 predictor by calculating the weighted average change in adjusted R^2 resulting from its
704 inclusion in the model, and using a bootstrapping procedure to compute confidence
705 intervals.

706 ***Meta-Analysis***

707 To obtain the overall meta-analytic estimate of the convergence of risk preference
708 measures, we first fitted a Bayesian hierarchical intercept-only model. Second, to obtain
709 meta-analytic estimates for the convergence between specific pairs of measure categories
710 and domains, we fitted Bayesian hierarchical (robust) regression models that included a
711 predictor coding for the different types of measure pairs.

712 Multiverse Analyses

713 We conducted a series of multiverse analyses with alternative data sets resulting
714 from different data pre-processing and various alternative analytic choices. We find
715 overall qualitatively similar patterns of results across the multiverse of choices
716 considered. We provide additional details concerning these analyses and results in the
717 Supplementary Information.

718 Data and Code Availability

719 All the data are made publicly available through the original data repositories
720 and need to be accessed by following the providers' data access policies. We provide
721 more detailed overview of data, analysis, and code in a companion website
722 (<https://cdsbasel.github.io/temprisk/>) and make the estimated test-retest correlations
723 and inter-correlations from the primary data sources as well as all analysis scripts
724 publicly available in an online repository (<https://osf.io/5kzgd/>).

725

Acknowledgements

726 This work was supported by grants from the Swiss National Science Foundation
727 to R.M. (<https://data.snf.ch/grants/grant/204700>,
728 <https://data.snf.ch/grants/grant/177277>). The authors thank Laura Wiles for editing
729 the manuscript.

730

Author contributions

731 A.B.: Conceptualization, Data curation, Formal analysis, Investigation,
732 Methodology, Project administration, Visualization, Writing - original draft, and
733 Writing - review & editing.

734 Y.L.: Data curation.

735 M.K.: Data curation.

736 G.S.: Data curation.

737 P.-C.B.: Formal analysis, Methodology, and Writing - review & editing.

738 L.T.: Conceptualization, Visualization, Writing - original draft, and Writing -
739 review & editing.

740 R.M.: Conceptualization, Formal analysis, Funding acquisition, Investigation,
741 Methodology, Project administration, Supervision, Visualization, Writing - original
742 draft, and Writing - review & editing.

743

Competing interests

744 The authors declare no competing interests.

References

- 745 Ahun, M. N., Lauzon, B., Sylvestre, M.-P., Bergeron-Caron, C., Eltonsy, S., &
746 O'Loughlin, J. (2020). A systematic review of cigarette smoking trajectories in
747 adolescents. *International Journal of Drug Policy*, 83, 102838.
748
749 <https://doi.org/10.1016/j.drugpo.2020.102838>
- 750 Anusic, I., & Schimmack, U. (2016). Stability and change of personality traits,
751 self-esteem, and well-being: Introducing the meta-analytic stability and change
752 model of retest correlations. *Journal of Personality and Social Psychology*,
753 110(5), 766–781. <https://doi.org/10.1037/pspp0000066>
- 754 Arslan, R. C., Brümmer, M., Dohmen, T., Drewelies, J., Hertwig, R., & Wagner, G. G.
755 (2020). How people know their risk preference. *Scientific Reports*, 10(1), 15365.
756
757 <https://doi.org/10.1038/s41598-020-72077-5>
- 758 Barseghyan, L., Molinari, F., O'Donoghue, T., & Teitelbaum, J. C. (2018). Estimating
759 risk preferences in the field. *Journal of Economic Literature*, 56(2), 501–564.
<https://doi.org/10.1257/jel.20161148>
- 760 Beck, E. D., & Jackson, J. J. (2022). A mega-analysis of personality prediction:
761 Robustness and boundary conditions. *Journal of Personality and Social
762 Psychology*, 122(3), 523–553. <https://doi.org/10.1037/pspp0000386>
- 763 Bernoulli, D. (1954). Exposition of a New Theory on the Measurement of Risk.
764
765 *Econometrica*, 22(1), 23–36. <https://doi.org/10.2307/1909829>
- 766 Bleidorn, W., Schwaba, T., Zheng, A., Hopwood, C. J., Sosa, S. S., Roberts, B. W., &
767 Briley, D. A. (2022). Personality stability and change: A meta-analysis of
768 longitudinal studies. *Psychological Bulletin*, 148, 588–619.
<https://doi.org/10.32614/RJ-2018-017>
- 769 Bran, A., & Vaidis, D. C. (2020). Assessing risk-taking: What to measure and how to
770 measure it. *Journal of Risk Research*, 23(4), 490–503.
771
<https://doi.org/10.1080/13669877.2019.1591489>

- 772 Breivik, G., Sand, T. S., & Sookermany, A. M. (2019). Risk-taking and sensation
773 seeking in military contexts: A literature review. *SAGE Open*, 9(1),
774 2158244018824498. <https://doi.org/10.1177/2158244018824498>
- 775 Bringmann, L. F., Elmer, T., & Eronen, M. I. (2022). Back to Basics: The Importance
776 of Conceptual Clarification in Psychological Science. *Current Directions in*
777 *Psychological Science*, 31(4), 340–346.
778 <https://doi.org/10.1177/09637214221096485>
- 779 Brodbeck, J., Duerrenberger, S., & Znoj, H. (2009). Prevalence rates of at risk,
780 problematic and pathological gambling in Switzerland. *The European Journal of*
781 *Psychiatry*, 23(2), 67–75.
- 782 Bürkner, P.-C. (2017). Brms: An R package for Bayesian multilevel models using Stan.
783 *Journal of Statistical Software*, 80, 1–28. <https://doi.org/10.18637/jss.v080.i01>
- 784 Bürkner, P.-C. (2018). Advanced Bayesian multilevel modeling with the R package
785 brms. *The R Journal*, 10(1), 395–411.
- 786 Bürkner, P.-C. (2021). Bayesian item response modeling in R with brms and Stan.
787 *Journal of Statistical Software*, 100, 1–54. <https://doi.org/10.18637/jss.v100.i05>
- 788 Caliendo, M., Fossen, F., & Kritikos, A. S. (2014). Personality characteristics and the
789 decisions to become and stay self-employed. *Small Business Economics*, 42(4),
790 787–814. <https://doi.org/10.1007/s11187-013-9514-8>
- 791 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M.,
792 Brubaker, M., Guo, J., Li, P., & Riddell, A. (2017). Stan: A probabilistic
793 programming language. *Journal of Statistical Software*, 76, 1–32.
794 <https://doi.org/10.18637/jss.v076.i01>
- 795 Chapman, J., Dean, M., Ortoleva, P., Snowberg, E., & Camerer, C. (2018).
796 *Econographics* (tech. rep. w24931). National Bureau of Economic Research.
797 Cambridge, MA. <https://doi.org/10.3386/w24931>
- 798 Chuang, Y., & Schechter, L. (2015). Stability of experimental and survey measures of
799 risk, time, and social preferences: A review and some new results. *Journal of*

- 800 *Development Economics*, 117, 151–170.
- 801 <https://doi.org/10.1016/j.jdeveco.2015.07.008>
- 802 Conley, J. J. (1984). The hierarchy of consistency: A review and model of longitudinal
803 findings on adult individual differences in intelligence, personality and
804 self-opinion. *Personality and Individual Differences*, 5(1), 11–25.
805 [https://doi.org/10.1016/0191-8869\(84\)90133-8](https://doi.org/10.1016/0191-8869(84)90133-8)
- 806 Dang, J., King, K. M., & Inzlicht, M. (2020). Why are self-report and behavioral
807 measures weakly correlated? *Trends in Cognitive Sciences*, 24(4), 267–269.
808 <https://doi.org/10.1016/j.tics.2020.01.007>
- 809 Das, A., & Ahmed, M. M. (2022). Structural equation modeling approach for
810 investigating drivers' risky behavior in clear and adverse weather using SHRP2
811 naturalistic driving data. *Journal of Transportation Safety & Security*.
812 <https://doi.org/10.1080/19439962.2022.2155744>
- 813 Dohmen, T., Falk, A., Huffman, D., Sunde, U., Schupp, J., & Wagner, G. G. (2011).
814 Individual risk attitudes: Measurement, determinants, and behavioral
815 consequences. *Journal of the European Economic Association*, 9(3), 522–550.
816 <https://doi.org/10.1111/j.1542-4774.2011.01015.x>
- 817 Duckworth, A. L., & Kern, M. L. (2011). A meta-analysis of the convergent validity of
818 self-control measures. *Journal of Research in Personality*, 45(3), 259–268.
819 <https://doi.org/10.1016/j.jrp.2011.02.004>
- 820 Eisenberg, I. W., Bissett, P. G., Enkavi, A. Z., Li, J., MacKinnon, D. P., Marsch, L. A.,
821 & Poldrack, R. A. (2019). Uncovering the structure of self-regulation through
822 data-driven ontology discovery. *Nature Communications*, 10(1), 2319.
823 <https://doi.org/10.1038/s41467-019-10301-1>
- 824 Elliott, M. L., Knott, A. R., Ireland, D., Morris, M. L., Poulton, R., Ramrakha, S.,
825 Sison, M. L., Moffitt, T. E., Caspi, A., & Hariri, A. R. (2020). What is the
826 test-retest reliability of common task-functional MRI measures? New empirical
827 evidence and a meta-analysis. *Psychological Science*, 31(7), 792–806.
828 <https://doi.org/10.1177/0956797620916786>

- 829 Enkavi, A. Z., Eisenberg, I. W., Bissett, P. G., Mazza, G. L., MacKinnon, D. P.,
830 Marsch, L. A., & Poldrack, R. A. (2019). Large-scale analysis of test-retest
831 reliabilities of self-regulation measures. *Proceedings of the National Academy of
832 Sciences*, 116(12), 5472–5477. <https://doi.org/10.1073/pnas.1818430116>
- 833 Falk, A., Becker, A., Dohmen, T., Enke, B., Huffman, D., & Sunde, U. (2018). Global
834 evidence on economic preferences. *The Quarterly Journal of Economics*, 133(4),
835 1645–1692. <https://doi.org/10.1093/qje/qjy013>
- 836 Financial Services Authority. (2011). *Assessing suitability: Establishing the risk a
837 customer is willing and able to take and making a suitable investment selection*
838 (tech. rep.). Financial Services Authority.
- 839 Fraley, R. C., & Roberts, B. W. (2005). Patterns of continuity: A dynamic model for
840 conceptualizing the stability of individual differences in psychological constructs
841 across the life course. *Psychological Review*, 112(1), 60–74.
842 <https://doi.org/10.1037/0033-295X.112.1.60>
- 843 Frey, R., Pedroni, A., Mata, R., Rieskamp, J., & Hertwig, R. (2017). Risk preference
844 shares the psychometric structure of major psychological traits. *Science
845 Advances*, 3(10), e1701381. <https://doi.org/10/gb2xrw>
- 846 Graham, E. K., Weston, S. J., Gerstorf, D., Yoneda, T. B., Booth, T., Beam, C. R.,
847 Petkus, A. J., Drewelies, J., Hall, A. N., Bastarache, E. D., Estabrook, R.,
848 Katz, M. J., Turiano, N. A., Lindenberger, U., Smith, J., Wagner, G. G.,
849 Pedersen, N. L., Allemand, M., Spiro, A., ... Mrocze, D. K. (2020).
850 Trajectories of Big Five personality traits: A coordinated analysis of 16
851 longitudinal samples. *European Journal of Personality*, 34(3), 301–321.
852 <https://doi.org/10.1002/per.2259>
- 853 Grömping, U. (2007). Estimators of relative importance in linear regression based on
854 variance decomposition. *The American Statistician*, 61(2), 139–147.
855 <https://doi.org/10.1198/000313007X188252>
- 856 Haines, N., Kvam, P. D., Irving, L. H., Smith, C., Beauchaine, T. P., Pitt, M. A.,
857 Ahn, W.-Y., & Turner, B. (2020). Learning from the reliability paradox: How

- 858 theoretically informed generative models can advance the social, behavioral, and
859 brain sciences. <https://doi.org/https://doi.org.10.31234/osf.io/xr7y3>
- 860 Hall, B. D., Liu, Y., Jansen, Y., Dragicevic, P., Chevalier, F., & Kay, M. (2022). A
861 survey of tasks and visualizations in multiverse analysis reports. *Computer*
862 *Graphics Forum*, 41(1), 402–426. <https://doi.org/10.1111/cgf.14443>
- 863 Harrison, G. W. (2014). Real choices and hypothetical choices. In S. Hess & A. Daly
864 (Eds.), *Handbook of choice modelling* (pp. 236–254). Edward Elgar Publishing.
- 865 Hedges, L. V., & Olkin, I. (1985). *Statistical methods for meta-analysis* (1st Ed.).
866 Academic Press.
- 867 Holt, C. A., & Laury, S. K. (2002). Risk aversion and incentive effects. *The American*
868 *Economic Review*, 92, 1644–1655. <https://doi.org/10.1257/000282802762024700>
- 869 Jin, H., Cui, M., & Liu, J. (2020). Factors affecting people's attitude toward
870 participation in medical research: A systematic review. *Current Medical Research*
871 *and Opinion*, 36(7), 1137–1143. <https://doi.org/10.1080/03007995.2020.1760807>
- 872 Josef, A. K., Richter, D., Samanez-Larkin, G. R., Wagner, G. G., Hertwig, R., &
873 Mata, R. (2016). Stability and change in risk-taking propensity across the adult
874 life span. *Journal of Personality and Social Psychology*, 111(3), 430–450.
875 <https://doi.org/10.1037/pspp0000090>
- 876 Karlsson Linnér, R., Biroli, P., Kong, E., Meddends, S. F. W., Wedow, R.,
877 Fontana, M. A., Lebreton, M., Tino, S. P., Abdellaoui, A., Hammerschlag, A. R.,
878 Nivard, M. G., Okbay, A., Rietveld, C. A., Timshel, P. N., Trzaskowski, M.,
879 de Vlaming, R., Zünd, C. L., Bao, Y., Buzdugan, L., ... Beauchamp, J. P.
880 (2019). Genome-wide association analyses of risk tolerance and risky behaviors in
881 over 1 million individuals identify hundreds of loci and shared genetic influences.
882 *Nature Genetics*, 51(2), 245–257. <https://doi.org/10.1038/s41588-018-0309-3>
- 883 Karlsson Linnér, R., Mallard, T. T., Barr, P. B., Sanchez-Roige, S., Madole, J. W.,
884 Driver, M. N., Poore, H. E., de Vlaming, R., Grotzinger, A. D., Tielbeek, J. J.,
885 Johnson, E. C., Liu, M., Rosenthal, S. B., Ideker, T., Zhou, H., Kember, R. L.,
886 Pasman, J. A., Verweij, K. J. H., Liu, D. J., ... Dick, D. M. (2021). Multivariate

- analysis of 1.5 million people identifies genetic associations with traits related to self-regulation and addiction. *Nature Neuroscience*, 24(10), 1367–1376.
<https://doi.org/10.1038/s41593-021-00908-3>
- Lejuez, C. W., Read, J. P., Kahler, C. W., Richards, J. B., Ramsey, S. E., Stuart, G. L., Strong, D. R., & Brown, R. A. (2002). Evaluation of a behavioral measure of risk taking: The Balloon Analogue Risk Task (BART). *Journal of Experimental Psychology. Applied*, 8(2), 75–84. <https://doi.org/10.1037/1076-898x.8.2.75>
- Levy, I., Snell, J., Nelson, A. J., Rustichini, A., & Glimcher, P. W. (2010). Neural Representation of Subjective Value Under Risk and Ambiguity. *Journal of Neurophysiology*, 103(2), 1036–1047. <https://doi.org/10.1152/jn.00853.2009>
- Mata, R., Frey, R., Richter, D., Schupp, J., & Hertwig, R. (2018). Risk preference: A view from psychology. *Journal of Economic Perspectives*, 32(2), 155–172.
<https://doi.org/10.1257/jep.32.2.155>
- Moffitt, T. E. (2018). Male antisocial behavior in adolescence and beyond. *Nature Human Behaviour*, 2, 177–186. <https://doi.org/10.1038/s41562-018-0309-4>
- Möttus, R., Briley, D. A., Zheng, A., Mann, F. D., Engelhardt, L. E., Tackett, J. L., Harden, K. P., & Tucker-Drob, E. M. (2019). Kids becoming less alike: A behavioral genetic analysis of developmental increases in personality variance from childhood to adolescence. *Journal of Personality and Social Psychology*, 117(3), 635–658. <https://doi.org/10.1037/pspp0000194>
- Norris, E., Finnerty, A. N., Hastings, J., Stokes, G., & Michie, S. (2019). A scoping review of ontologies related to human behaviour change. *Nature Human Behaviour*, 3(2), 164–172. <https://doi.org/10.1038/s41562-018-0511-4>
- Orth, U. (2018). Development of self-esteem from age 4 to 94 years: A meta-analysis of longitudinal studies. *Psychological Bulletin*, 144(10), 1045.
<https://doi.org/10.1037/bul0000161>
- Pedroni, A., Frey, R., Bruhin, A., Dutilh, G., Hertwig, R., & Rieskamp, J. (2017). The risk elicitation puzzle. *Nature Human Behaviour*, 1(11), 803–809.
<https://doi.org/10.1038/s41562-017-0219-x>

- 916 R Core Team. (2021). R: A language and environment for statistical computing.
- 917 Ravert, R. D., Murphy, L. M., & Donnellan, M. B. (2019). Valuing risk: Endorsed risk
918 activities and motives across adulthood. *Journal of Adult Development*, 26(1),
919 11–21. <https://doi.org/10/gmndcz>
- 920 Richmond-Rakerd, L. S., D’Souza, S., Andersen, S. H., Hogan, S., Houts, R. M.,
921 Poulton, R., Ramrakha, S., Caspi, A., Milne, B. J., & Moffitt. (2020). Clustering
922 of health, crime and social-welfare inequality in 4 million citizens from two
923 nations. *Nature Human Behaviour*, 4(3), 255–264.
924 <https://doi.org/doi:10.1038/s41562-019-0810-4>
- 925 Roberts, D. K., Alderson, R. M., Betancourt, J. L., & Bullard, C. C. (2021).
926 Attention-deficit/hyperactivity disorder and risk-taking: A three-level
927 meta-analytic review of behavioral, self-report, and virtual reality metrics.
928 *Clinical Psychology Review*, 97, 102039.
929 <https://doi.org/10.1016/j.cpr.2021.102039>
- 930 Schildberg-Hörisch, H. (2018). Are risk preferences stable? *Journal of Economic
931 Perspectives*, 32(2), 135–154. <https://doi.org/10.1257/jep.32.2.135>
- 932 Schmidt, L. (2008). Risk preferences and the timing of marriage and childbearing.
933 *Demography*, 45(2), 439–460. <https://doi.org/10.1353/dem.0.0005>
- 934 Schonberg, T., Fox, C. R., & Poldrack, R. A. (2011). Mind the gap: Bridging economic
935 and naturalistic risk-taking with cognitive neuroscience. *Trends in Cognitive
936 Sciences*, 15(1), 11–19. <https://doi.org/10.1016/j.tics.2010.10.002>
- 937 Schönbrodt, F. D., & Perugini, M. (2013). At what sample size do correlations
938 stabilize? *Journal of Research in Personality*, 47(5), 609–612.
939 <https://doi.org/10.1016/j.jrp.2013.05.009>
- 940 Seifert, I. S., Rohrer, J. M., Egloff, B., & Schmukle, S. C. (2022). The development of
941 the rank-order stability of the Big Five across the life span. *Journal of
942 Personality and Social Psychology*, 122, 920–941.
943 <https://doi.org/10.1037/pspp0000398>

- 944 Sharapov, D., Kattuman, P., Rodriguez, D., & Velazquez, F. J. (2021). Using the
945 SHAPLEY value approach to variance decomposition in strategy research:
946 Diversification, internationalization, and corporate group effects on affiliate
947 profitability. *Strategic Management Journal*, 42(3), 608–623.
948 <https://doi.org/10.1002/smj.3236>
- 949 Stan Development Team. (2022). Stan user's guide. Version 2.29.
- 950 Steegen, S., Tuerlinckx, F., Gelman, A., & Vanpaemel, W. (2016). Increasing
951 transparency through a multiverse analysis. *Perspectives on Psychological
952 Science*, 11(5), 702–712. <https://doi.org/10.1177/1745691616658637>
- 953 Steinberg, L. (2013). The influence of neuroscience on US Supreme Court decisions
954 about adolescents' criminal culpability. *Nature Reviews Neuroscience*, 14(7),
955 513–518. <https://doi.org/10/gdcf6b>
- 956 Steiner, M. D., Seitz, F. I., & Frey, R. (2021). Through the window of my mind:
957 Mapping information integration and the cognitive representations underlying
958 self-reported risk preference. *Decision*, 8, 97–122.
959 <https://doi.org/10.1037/dec0000127>
- 960 Stigler, G. J., & Becker, G. S. (1977). De gustibus non est disputandum. *The American
961 Economic Review*, 67(2), 76–90.
962 <https://doi.org/http://www.jstor.org/stable/1807222>
- 963 Strickland, J. C., & Johnson, M. W. (2021). Rejecting impulsivity as a psychological
964 construct: A theoretical, empirical, and sociocultural argument. *Psychological
965 review*, 128(2), 336–361. <https://doi.org/10.1037/rev0000263>
- 966 Tisdall, L., MacNiven, K. H., Padula, C. B., Leong, J. K., & Knutson, B. (2022). Brain
967 tract structure predicts relapse to stimulant drug use. *Proceedings of the
968 National Academy of Sciences of the United States of America*, 119(26),
969 e2116703119. <https://doi.org/10.1073/pnas.2116703119>
- 970 Williams, D. R., Rast, P., & Bürkner, P.-C. (2018). Bayesian meta-analysis with weakly
971 informative prior distributions. <https://doi.org/10.31234/osf.io/7tbrm>

- 972 Wulff, D. U., & Mata, R. (2022). On the semantic representation of risk. *Science*
973 *Advances*, 8(27), eabm1883. <https://doi.org/10.1126/sciadv.abm1883>

Table 1

Descriptions and examples of different categories of risk preference measures

Category	Description	Example
Propensity	self-report measures; individuals indicate on a (ordinal) scale to what extent they identify as someone who likes or is willing to take risks in general or in specific domains.	<i>Are you generally a person who is willing to take risks or do you try to avoid taking risks?</i> (Dohmen et al., 2011)
Frequency	self-report measures; individuals indicate on a scale or in an open field to what extent or how often they partake in activities in specific life domains.	<i>How many times in the last seven days have you had an alcoholic drink?</i> (Frey et al., 2017)
Behavioural	behavioural measures; individuals are asked to decide between two or more options typically offering different (hypothetical or real) monetary gains and/or losses with varying probability; an index of risk preference is typically derived based on a combination of choices or actions.	Mean number of pumps in a simulated balloon-pumping task (Lejuez et al., 2002); percentage of risky choices in a lottery task (Holt & Laury, 2002)

Figure 1

Flowchart of systematic search.

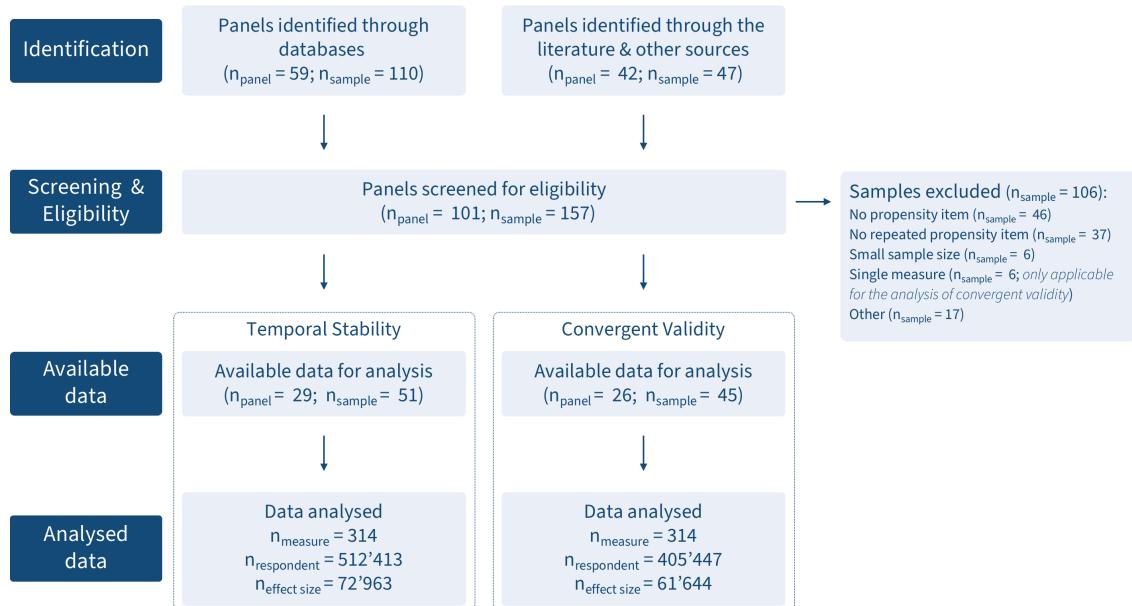


Figure 2

Overview of data. A) Two Dimensional density plot of test-retest correlations as a function of retest interval ($k = 72,963$). B) Distribution of all inter-correlations ($k = 61,644$). C) Number of unique measures split by category (propensity, frequency, behaviour), and domain.

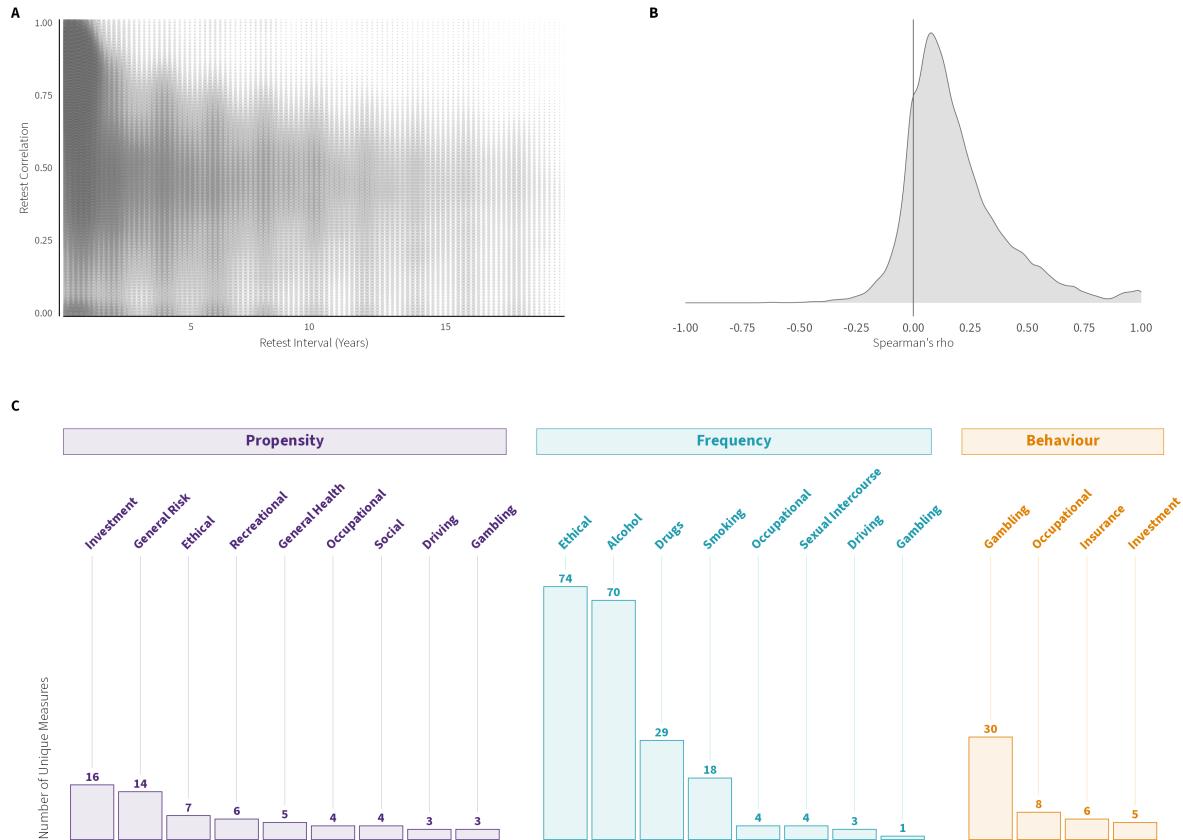


Figure 3

Variance decomposition of temporal stability. A) Relative contribution of measure, respondent, and panel-related predictors to the adjusted R^2 in regression models predicting test-retest correlations of all risk preference measures ($k = 72,963$). B) Relative contribution of measure, respondent, and panel predictors to the adjusted R^2 in regression models predicting test-retest correlations of propensity ($k = 23,936$), frequency ($k = 47,490$), and behavioural ($k = 1,537$) measures. Estimate (dot) and bootstrapped (coloured area) 95%, 80%, and 50% confidence intervals.

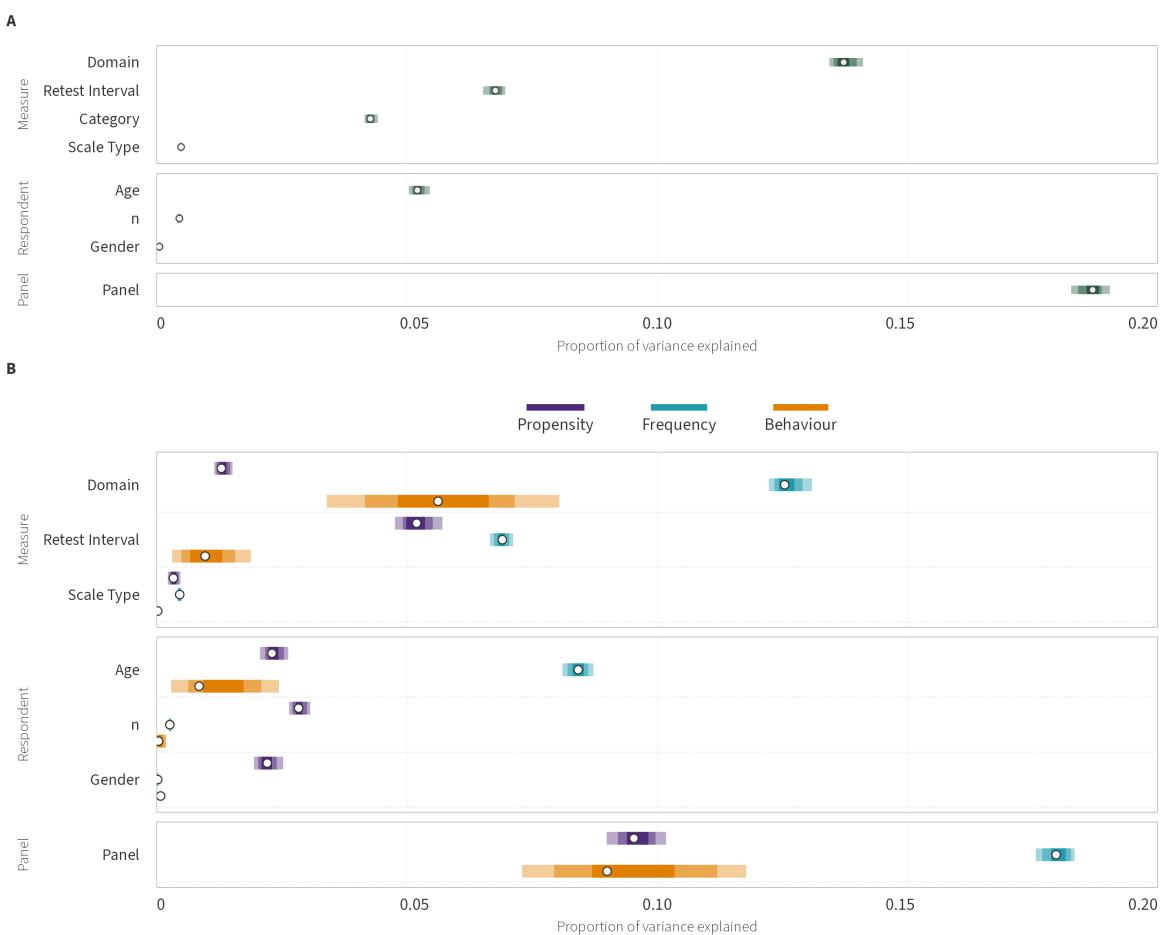


Figure 4

Meta-Analytic Stability and Change model (MASC) results. The figure shows parameter estimates for A) propensity ($k = 3,706$), B) frequency ($k = 3,678$), and C) behavioural measures ($k = 612$) of risk preference. In A-C, circles represent the mean estimate, shaded uncertainty bands represent the 50%, 80%, and 95% HDI. D-I show predictions of retest trajectories given MASC parameters as a function of retest interval (D,F,H) or age (E,G,I) across all domains (shaded uncertainty bands, 50%, 80%, and 95% HDI) as well as a selection of two domains per category (individual, annotated lines)

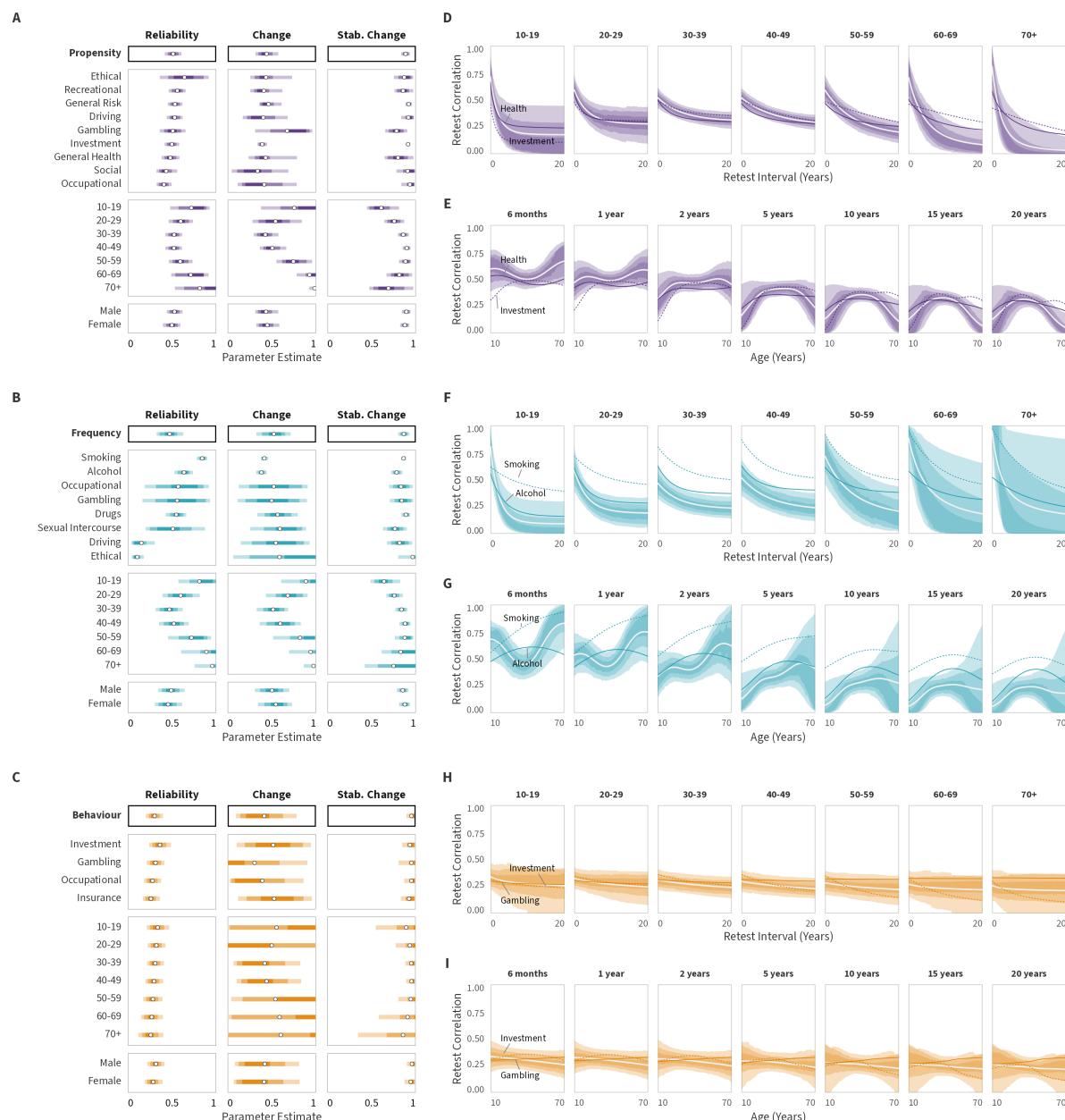


Figure 5

Variance decomposition of convergence between measures. Relative contribution of measure, respondent, and panel-related predictors to the adjusted R^2 in regression models predicting inter-correlations between measures of risk preference ($k = 61,644$). Estimate (dot) and bootstrapped (coloured area) 95%, 80%, and 50% confidence intervals.

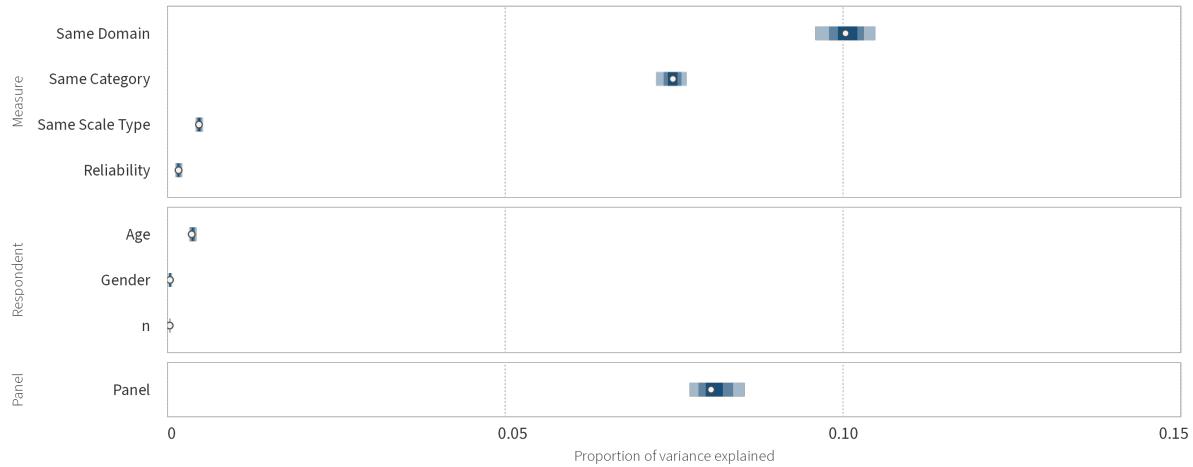
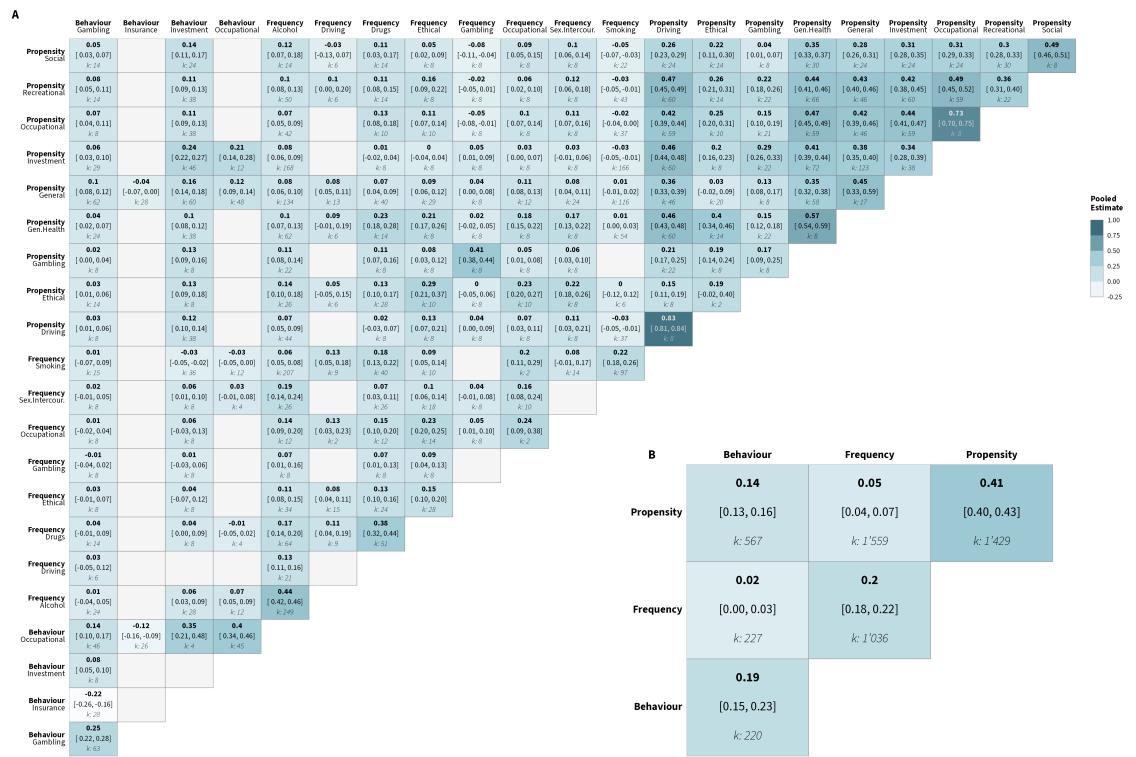


Figure 6

Meta-analytic correlation matrices. The matrices depicts the results of the meta-analyses of inter-correlations between measures of risk preference ($k = 5,038$), with each cell representing the meta-analytic result for the specific measurement pair of A) measure domains or B) measure categories. Empty grey cells are due to lack of data availability to estimate the respective correlation.



Supplementary Information

*

974	Contents	
975	Identification of Samples	S2
976	Categorisation of Measures	S2
977	Data Pre-Processing	S3
978	Data Set Information	S3
979	Risk Preference Measures	S3
980	Variable recoding based on question dependency.	S3
981	Reverse coding.	S4
982	Composite measures.	S4
983	Harmonising variable names	S4
984	Sample Demographics	S4
985	Data Processing	S5
986	Temporal Stability	S5
987	Computing test-retest correlations	S5
988	Aggregating test-retest correlations	S7
989	Convergent Validity	S7
990	Computing inter-correlations	S7
991	Aggregating inter-correlations	S8
992	Analysis	S9
993	Temporal Stability	S9
994	Variance decomposition	S9
995	Meta-Analytic Stability and Change model (MASC)	S10
996	Re-analysis of the Anusic and Schimmack (2016) data set	S13
997	Convergent Validity	S14
998	Variance decomposition	S14
999	Meta-analyses	S14
1000	Multiverse Analyses	S16

1001 Identification of Samples

1002 To find as many longitudinal panels and associated samples with risk preference
1003 measures, we devised a list of search terms related to risk (e.g., risk*, gambl*, smok*,
1004 gambl*; Table S1). This list reflects the definition of risk from the economics and
1005 psychology literatures and covers many different areas of life. It was developed by
1006 consulting the questionnaires of multi-measure studies (Arslan et al., 2020; Chapman
1007 et al., 2018; Eisenberg et al., 2019; Enkavi et al., 2019; Falk et al., 2018; Frey et al.,
1008 2017) as well as previously identified longitudinal samples (e.g., SOEP, USOC). As
1009 presented in the main paper, this search led us to identify a large number of panels
1010 (101) and associated samples (157) (Table S2), which we checked for possible inclusion
1011 in our study. We excluded sample or measures from our study using a clear set of
1012 inclusion/exclusion criteria (Table S3). Each sample was documented differently, thus,
1013 whenever available, we used the computerised (online) variable search engine to search
1014 for the risk-related terms, otherwise, we manually searched the codebooks and/or
1015 questionnaires available. Our systematic approach to search and screening resulted in
1016 the inclusion of 51 unique samples from 29 panels (Table S4).

1017 Categorisation of Measures

1018 We conducted extensive coding and categorisation of each risk preference
1019 measure that met our inclusion criteria. Specifically, we coded the following information:
1020 the name of the panel it originated from, the measure category (i.e., propensity,
1021 frequency, or behaviour), the domain (e.g., recreational, smoking), the type of scale
1022 used (i.e., ordinal, discrete, composite or open ended) and, if ordinal or discrete (with a
1023 clear range of possible response values), the number of options or points in the scale. In
1024 addition, we included information that was specific to each category of risk preference
1025 measure. Specifically, for frequency measures, we specified the number of days over
1026 which a certain behaviour had to be reported (e.g., Over the last week/month/year how
1027 many times were you intoxicated?). For behavioural measures, we recorded whether the
1028 decision was incentivised or hypothetical (cf., Harrison, 2014). Please note that we do

1029 not include these category-specific characteristics in our analyses because they are not
1030 instrumental to the comparison between categories. Nevertheless, we provide this
1031 categorisation for completeness and future possible uses of these data that control for or
1032 examine the role of such characteristics. Overall, we identified 314 unique measures
1033 stemming from 51 longitudinal samples. We provide a detailed definition, coding and
1034 description of each type of measure in Table S5, as well as a complete list of the risk
1035 preference measures in the main code book available in the online repository.

1036 Data Pre-Processing

1037 Prior to computing test-retest correlations, we pre-processed the data from each
1038 sample to create homogeneous data sets with regards to the data set information, risk
1039 preference measures, and sample demographics. We provide details concerning each
1040 step below.

1041 Data Set Information

1042 From each data set, we extracted the wave identifiers and data collection dates
1043 (i.e., day-month-year). If these dates were missing, we determined for each wave a
1044 standard date by referring to the sample's data collection timeline and choosing the
1045 half-way point (e.g., if data collection took place between January and June of 2020, the
1046 15th of March 2020 was selected as the date). In the case that only the year could be
1047 retrieved, we set June 15th as the default day. If the data collection date was missing
1048 for certain respondents within the wave of a panel, this date was filled by the mean of
1049 the available dates.

1050 Risk Preference Measures

1051 **Variable recoding based on question dependency.** Depending on the
1052 design of the questionnaires/interviews, for some samples, respondents were not asked
1053 certain questions because of their response to previous (filter) questions. This was
1054 particularly the case for frequency measures. For instance, if an individual answered the
1055 question “*Are you currently a smoker?*” with “*No*”, the follow-up question “*How many*

1056 *cigarettes a day do you smoke?"* would not be asked and would automatically receive a
1057 "missing" or "not applicable" code. By ignoring dependencies between questions,
1058 valuable information on the consistency of an individual's behaviour is missed, as
1059 instances of when behaviours might be interrupted and taken up again (e.g.,
1060 quitting/taking up smoking) are unaccounted for. To deal with this, for each sample,
1061 we took into account responses to filter-type questions and replaced invalid/missing
1062 codes in subsequent related questions by an appropriate response. In the case of the
1063 above example, for all the participants who answered "*No*", we replaced the invalid or
1064 missing code for the number of cigarettes smoked in a day with a "*0*" or "*None*". To
1065 make such replacements possible, we only included measures in our analyses that had
1066 scales that offered the possibility of a 0 value or Never/None answer (Table S3).

1067 **Reverse coding.** Whenever appropriate, we reversed the scales of measures
1068 such that higher values corresponded to greater risk-taking.

1069 **Composite measures.** We define a composite measure as a measure which
1070 represents an index of risk taking that is calculated by combining two or more
1071 individual risk preference measures. This was particularly the case for behavioural
1072 measures. If a composite measure was not available in the raw data set of the sample,
1073 we aggregated the set of available single responses using similar methods as that of
1074 studies with comparable tasks (e.g., proportion of risky choices). We provide a
1075 description of how these have been calculated for specific measures in the risk
1076 preference measure code book.

1077 **Harmonising variable names.** We standardised the names of the measures
1078 such that the same risk preference measure (or highly similarly worded measure with
1079 the same response format and scale) included in different samples shared the same
1080 variable name.

1081 ***Sample Demographics***

1082 We recorded the age and gender of each respondent. Age was calculated at the
1083 time of each data collection point. If the respondent's birth year was available in the
1084 data set, we used that to calculate their age, if not, we used the value of the

1085 pre-computed age in the data set. Further, if only age group or age range information
1086 was available (e.g., 20-30), we defined age as the midpoint value (e.g., 25). Only data
1087 from respondents between the ages of 10 and 90 years were included in the analyses.
1088 We coded gender as a binary variable (0 = male and 1 = female). For data quality
1089 purposes, we did not include in our analyses the responses of respondents whose year of
1090 birth, age (i.e., if the age difference and time difference between first and last wave of
1091 participation differed by more than 2 years) or gender was inconsistently reported
1092 across waves. Additionally, if either the year of birth, age or gender was missing and
1093 could not be retrieved or estimated based on previous waves, the respondent was
1094 excluded from the analyses.

1095 **Data Processing**

1096 ***Temporal Stability***

1097 **Computing test-retest correlations.** To address our main research
1098 objectives, for each panel and risk preference measure, we calculated for all possible
1099 wave combinations test-retest correlations (Figure S1). Correlations were calculated
1100 separately for females and males of different age groups. We computed separate sets of
1101 test-retest correlations for different age group configurations: 5, 10 and 20-year age
1102 bins. Akin to Enkavi et al. (2019), we estimated test-retest correlations using three
1103 different metrics: Pearson's r, Spearman's rho and intra-class correlations (ICC(2,1)).
1104 The correlation between these different metrics ranged between 0.58 and 0.99 (Figure
1105 S2). Further, the response distributions of some measures were highly skewed, thus we
1106 additionally computed test-retest correlations using log-transformed data. As shown in
1107 Figure S3, these were highly correlated with the test-retest correlations computed using
1108 the non-transformed data ($r = 0.91 - 0.98$). As a consequence, we report our main
1109 results using the Pearson's r correlation coefficient for the non-transformed data.
1110 Furthermore, when computing the test-retest correlations we obtained negative
1111 estimates (3.95% of the data set used for analysis); for ease of interpretation, we
1112 replaced these values with zeroes prior to any analysis or aggregation procedures (cf.,

1113 Enkavi et al., 2019).

1114 **Additional metrics.** In addition to these correlation metrics, for each
1115 test-retest correlation coefficient we recorded the following (variables with an asterisk
1116 were included in our main analyses, the rest were included for data quality assessment
1117 and data exploration):

- 1118 • Respondent information: sample size, maximum age, minimum age, mean age*,
1119 median age, standard deviation of age, proportion of female respondents*,
1120 proportion of sample lost between the first and second data collection point (i.e.,
1121 attrition rate)
- 1122 • Retest interval: minimum, maximum, mean*, median and standard deviation of
1123 the number of years between the first and second data collection point
- 1124 • Response properties: the coefficient of variation and skewness of the responses at
1125 both time points

1126 When calculating the time interval between the first and second data collection
1127 point, we noted that for panels that collected data for different surveys simultaneously
1128 (e.g. American Life Panel), not all respondents completed the surveys in the same
1129 order; some respondents would complete a more recent survey prior to an older survey
1130 (based on the mean data collection date), resulting in a negative retest interval.

1131 Therefore, for a very small number of correlations (0.17%) the minimum retest interval
1132 was negative. However, in our analyses we use the mean time difference between waves
1133 (or surveys), which minimises this issue. One exception to this concerns the German
1134 Socioeconomic Panel, which in 2020 launched a COVID-specific survey in which data
1135 collection overlapped with the 2020 core survey. We could not adequately order this
1136 pair of waves (i.e., 2020-core and 2020-covid) as we systematically had correlations that
1137 either had a negative mean or median retest interval. Therefore, we excluded
1138 correlations that from this specific pair of waves.

1139 **Sample size.** Simulation studies have shown that large sample sizes may be
1140 needed to compute stable correlation coefficients (Schönbrodt & Perugini, 2013). On the

1141 companion website we show how the number of correlation coefficients in the data set
1142 varies for different age groups based on different minimum sample size thresholds. For
1143 some age groups a substantial number of coefficients are lost as the threshold increases.
1144 To avoid losing valuable information for certain age groups, we retained the set of
1145 test-retest correlations that had a sample size of at least 30 with age groups organised
1146 in 10-year bins. In line with the multiverse approach (Steegen et al., 2016), the
1147 companion website provides an overview of the outcome of our analysis obtained using
1148 the different minimum sample size thresholds, age bins, and other processing steps.

1149 **Aggregating test-retest correlations.** Given the high number of test-retest
1150 correlations in our data set ($N = 72,963$ correlations), it was too complex and
1151 computationally intensive to adequately estimate the Meta-Analytic Stability and
1152 Change model (MASC; Anusic and Schimmack (2016)) and capture the trajectories of
1153 the correlations over time without encountering severe model convergence issues.
1154 Therefore, we aggregated the correlations prior to fitting the MASC model. Specifically,
1155 first, we transformed each Pearson's r correlation coefficient into Fisher's z, and
1156 calculated the corresponding sampling variance. Second, we grouped the test-retest
1157 correlations by panel, measure category, measure domain, 3-month retest interval,
1158 gender, and age group. For each grouping we computed a synthesised estimate by
1159 aggregating test-retest correlation coefficients whilst accounting for the dependency
1160 between them as these were computed from the same set or subset of respondents
1161 (Hedges & Olkin, 1985). For this purpose, we used inverse-variance weighting and set
1162 the correlation of the sampling errors within subsets to .5. Lastly, these aggregated
1163 correlations and their standard errors were back transformed to Pearson's r. Given that
1164 MASC model predictions are bounded between 0 and 1, we set any negative aggregated
1165 retest correlation to zero. This process resulted in 7,996 aggregated test-retest
1166 correlation being calculated.

1167 ***Convergent Validity***

1168 **Computing inter-correlations.** Samples which contained only one measure
1169 of risk preference were excluded from these analyses ($n = 6$). For each of the remaining

samples and waves, we calculated correlations between the responses of every possible pair of measures, for every wave same time point. Similar to the test-retest correlations, inter-correlations were calculated separately for females and males of different age groups. Specifically, we computed separate sets of correlations for different age group configurations: 5, 10 and 20-year age bins. We estimated inter-correlations using three different metrics, Pearson's r, Spearman's rho and intraclass correlations (ICC(2,1)), and examined inter-correlations being computed using non-transformed or log-transformed data. As shown in Figure S5, inter-correlations computed using different metrics were highly correlated ($r = 0.84 - 0.92$), as were the inter-correlations for (non)transformed data ($r = 0.95 - 0.99$) (Figure S6).

Additional metrics. For each inter-correlation coefficient we additionally recorded the following (variables with an asterisk were included in our main analyses, the rest were included for data quality assessment and data exploration):

- Response information: sample size, maximum age, minimum age, mean age*, median age, standard deviation of age, proportion of female respondents*
- Response properties: the coefficient of variation and skewness of the responses of both measures

Aggregating inter-correlations. In an effort to reduce computational costs and the potential occurrence of divergent transitions when conducting the Bayesian meta-analysis, we aggregated the inter-correlations. We followed a similar approach as for the retest correlations, we first converted each correlation coefficient into Fisher's z, and calculated the corresponding sampling variance. We then split the set of inter-correlations by sample, gender, age group, and category-domain measure pairs. For each subset we computed a synthesised estimate by aggregating the Fisher's z values using inverse-variance weighting and accounting for the dependency between them as these were computed from the same set or subset of respondents (Hedges & Olkin, 1985). To average these correlations we used inverse-variance weighting and set the correlation of the sampling errors within subsets to .5. We conducted additional

1198 analyses in which we tested the effects of this correlation on our results by setting the
1199 correlation to 0.1 and 0.9.

1200 **Analysis**

1201 ***Temporal Stability***

1202 **Variance decomposition.** To gain a better understanding of the
1203 heterogeneity observed between test-retest correlations, we conducted a variance
1204 decomposition analysis by computing the Shapley values for the following predictors:

1205 Measure characteristics

- 1206 • Category: type of measure (i.e., propensity, frequency, behaviour)
- 1207 • Domain: life domain the measure focuses on (e.g., smoking, driving, social,
1208 ethical)
- 1209 • Scale type: type of response scale (i.e., open-ended/composite index,
1210 ordinal/discrete scales)
- 1211 • Retest Interval: number of years between T1 and T2 data collection

1212 Respondent characteristics

- 1213 • Age: age group the respondents belong in (10 year bins, e.g., 20-29, 30-39)
- 1214 • Gender: gender of the respondents (i.e., female, male)
- 1215 • Number of responses: sample size for each correlation

1216 Shapley values were computed by first estimating a linear regression for each
1217 possible combination of predictors (i.e., 2^8 models for the omnibus analysis, and 2^7
1218 models for the category-specific analyses) and extracting the adjusted R^2 value. Then,
1219 for each predictor, we computed the weighted average of the change in adjusted R^2
1220 resulting from the inclusion of that predictor in the models.

1221 To obtain bootstrapped confidence intervals, we sampled the data set of
1222 correlations 100 times, and estimated for each predictor a set of 100 Shapley values. To

1223 visualise these results, we ranked these values to determine the 50%, 80% and 95%
 1224 confidence intervals.

1225 **Meta-Analytic Stability and Change model (MASC).**

1226 ***Model specification.*** To assess the trajectory of test-retest correlations of
 1227 risk preference over time we used the MASC model developed by Anusic and
 1228 Schimmack (2016) (Figure S4). Specifically, we were interested in quantifying the effects
 1229 of gender, linear age, quadratic age, and domain, as well as the interactions between
 1230 linear and quadratic age with domain on each of the MASC model parameters (i.e., *rel*,
 1231 *change* and *stabch*).

1232 In the model, domain was a sum contrast coded factor, gender was the
 1233 proportion of female respondents (*FemaleProp*) centred at 0.5 (i.e., -0.5 = males and 0.5
 1234 = females), and age (*Age*) corresponded to the mean age of the respondents centred at
 1235 40 years and transformed into decades. Quadratic age (*Age2*) was the square value of
 1236 the *Age* predictor. Lastly, retest interval was coded as the number of decades between
 1237 waves.

1238 The samples differed from each other on multiple dimensions (e.g, country, mode
 1239 of data collection), hence, to account for such differences when estimating the MASC
 1240 model parameters, we included *sample* as a random factor. We limited the (correlated)
 1241 random effects structure to the *rel* parameter by adding a varying intercept and varying
 1242 slopes for the effects of linear age, quadratic age and gender¹. We did not include a
 1243 random effects structure for the estimation of the *change* and *stabch* parameters,
 1244 because to appropriately estimate these parameters samples should have data for a long
 1245 enough period such that the test-retest correlations assymptote (Anusic & Schimmack,
 1246 2016). In the current data set, the number of test-retest correlations per sample varied
 1247 substantially, and less than the majority of the samples (~ 40%) contained retest
 1248 correlations beyond an interval of 10 years.

1249 The values of *rel*, *change* and *stabch* are bounded between 0 and 1. The *rel* and

¹ We did not include a varying slope for the effect of domain as not every sample had data on each level of domain.

1250 *change* parameters both represent proportions (i.e., the proportion of reliable
 1251 between-person variance and the proportion of reliable variance attributable to
 1252 changing factors, respectively). For the *stabch* parameter (i.e., the rate of change) we
 1253 need to take into account that over the years changes in individuals' lives accumulate
 1254 and gradually affect their behaviour to different extents, resulting in decreasing (i.e., 0
 1255 < rate of change ≤ 1) rather than increasing (i.e., rate of change > 1) correlations
 1256 across the years (Anusic & Schimmack, 2016). Therefore, to ensure that these
 1257 parameters remained within their valid intervals, we modelled them on the logit scale
 1258 (i.e., *logitrel*, *logitchange* and *logitstabch*), and subsequently back-transformed them via
 1259 the inverse logit function (Bürkner, 2021). Such as to obtain meta-analytic estimates of
 1260 each parameter, we additionally specified in the model the corresponding standard
 1261 errors of the (aggregated) retest correlations.

1262 We used Bayesian inference to estimate the meta-analytic model and specified
 1263 weakly informative priors for the model parameters and hierarchical standard deviations
 1264 so as to include estimates reported in previous literature (e.g., Anusic and Schimmack,
 1265 2016; Frey et al., 2017; Mata et al., 2018). The Bayesian hierarchical non-linear model
 1266 described below was estimated using the probabilistic programming language Stan
 1267 (Carpenter et al., 2017; Stan Development Team, 2022) via the R package *brms*
 1268 (Bürkner, 2017, 2018, 2021). The companion website reports the summary output of
 1269 the model, sample-specific model predictions, and MCMC diagnostic plots.

$$y_i \sim StudentT(\nu, \theta_i, \sqrt{se_i^2 + \sigma^2})$$

$$\theta_i = rel_i \times (change \times (stabch^{\text{time}} - 1) + 1)$$

$$rel_i = \text{logit}^{-1}(logitrel_i)$$

$$change = \text{logit}^{-1}(logitchange)$$

$$stabch = \text{logit}^{-1}(logitstabch)$$

$$\sigma \sim Cauchy_+(0, 1)$$

$$\nu \sim Gamma(2, 0.1)$$

1270 **logitrel_i parameter**

$$\begin{aligned}
logitrel_i = & (\beta_{logitrel_0} + \beta_{logitrel_{0,sample[i]}}) + (\beta_{logitrel_1} + \beta_{logitrel_{1,sample[i]}}))Age + \\
& (\beta_{logitrel_2} + \beta_{logitrel_{2,sample[i]}})Age2 + (\beta_{logitrel_3} + \beta_{logitrel_{3,sample[i]}})FemaleProp + \\
& \beta_{logitrel_4}Domain + \beta_{logitrel_5}(Age \times Domain) + \beta_{logitrel_6}(Age2 \times Domain)
\end{aligned}$$

$$\beta_{logitrel_0}, \beta_{logitrel_1}, \beta_{logitrel_2}, \beta_{logitrel_3}, \beta_{logitrel_4}, \beta_{logitrel_5}, \beta_{logitrel_6} \sim Normal(0, 1)$$

$$\begin{bmatrix} \beta_{logitrel_{0,sample}} \\ \beta_{logitrel_{1,sample}} \\ \beta_{logitrel_{2,sample}} \\ \beta_{logitrel_{3,sample}} \end{bmatrix} \sim MVNormal \left(\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, Cov \right)$$

1271

$$Cov = \begin{pmatrix} \sigma_{\beta_{logitrel_{0,sample}}}^2 & \sigma_{\beta_{logitrel_{0,sample}}} \sigma_{\beta_{logitrel_{1,sample}}} \rho_{0,1} & \sigma_{\beta_{logitrel_{0,sample}}} \sigma_{\beta_{logitrel_{2,sample}}} \rho_{0,2} & \dots \\ \sigma_{\beta_{logitrel_{1,sample}}} \sigma_{\beta_{logitrel_{0,sample}}} \rho_{0,1} & \sigma_{\beta_{logitrel_{1,sample}}}^2 & \sigma_{\beta_{logitrel_{1,sample}}} \sigma_{\beta_{logitrel_{2,sample}}} \rho_{1,2} & \dots \\ \sigma_{\beta_{logitrel_{2,sample}}} \sigma_{\beta_{logitrel_{0,sample}}} \rho_{0,2} & \sigma_{\beta_{logitrel_{2,sample}}} \sigma_{\beta_{logitrel_{1,sample}}} \rho_{1,2} & \sigma_{\beta_{logitrel_{2,sample}}}^2 & \dots \\ \sigma_{\beta_{logitrel_{3,sample}}} \sigma_{\beta_{logitrel_{0,sample}}} \rho_{0,3} & \sigma_{\beta_{logitrel_{3,sample}}} \sigma_{\beta_{logitrel_{1,sample}}} \rho_{1,3} & \sigma_{\beta_{logitrel_{3,sample}}} \sigma_{\beta_{logitrel_{2,sample}}} \rho_{2,3} & \dots \\ \sigma_{\beta_{logitrel_{0,sample}}} \sigma_{\beta_{logitrel_{2,sample}}} \rho_{0,3} \\ \sigma_{\beta_{logitrel_{1,sample}}} \sigma_{\beta_{logitrel_{3,sample}}} \rho_{1,3} \\ \dots \sigma_{\beta_{logitrel_{2,sample}}} \sigma_{\beta_{logitrel_{3,sample}}} \rho_{2,3} \\ \sigma_{\beta_{logitrel_{3,sample}}}^2 \end{pmatrix}$$

1272

$$\sigma_{\beta_{logitrel_{0,sample}}}, \sigma_{\beta_{logitrel_{1,sample}}}, \sigma_{\beta_{logitrel_{2,sample}}}, \sigma_{\beta_{logitrel_{3,sample}}} \sim Cauchy_+(0, 1)$$

$$\rho \sim LKJCorr(1)$$

1273 logitchange parameter

$$\begin{aligned}
logitchange = & \beta_{logitchange_0} + \beta_{logitchange_1}Age + \beta_{logitchange_2}Age2 + \beta_{logitchange_3}FemaleProp + \\
& \beta_{logitchange_4}Domain + \beta_{logitchange_5}(Age \times Domain) + \beta_{logitchange_6}(Age2 \times Domain) \\
& \beta_{logitchange_0}, \beta_{logitchange_1}, \beta_{logitchange_2}, \beta_{logitchange_3}, \\
& \beta_{logitchange_4}, \beta_{logitchange_5}, \beta_{logitchange_6} \sim Normal(0, 1)
\end{aligned}$$

1274

1275 logitstabch parameter

$$\begin{aligned}
logitstabch = & \beta_{logitstabch_0} + \beta_{logitstabch_1}Age + \beta_{logitstabch_2}Age2 + \beta_{logitstabch_3}FemaleProp + \\
& \beta_{logitstabch_4}Domain + \beta_{logitstabch_5}(Age \times Domain) + \beta_{logitstabch_6}(Age2 \times Domain)
\end{aligned}$$

$$\beta_{logitstabch_0}, \beta_{logitstabch_1}, \beta_{logitstabch_2}, \beta_{logitstabch_3},$$

$$\beta_{logitstabch_4}, \beta_{logitstabch_5}, \beta_{logitstabch_6} \sim Normal(0, 1)$$

1276 **Re-analysis of the Anusic and Schimmack (2016) data set.** We
 1277 re-analysed the data that the authors made available in the study's supplementary
 1278 material. The authors collated and analysed test-retest correlations spanning 15 years
 1279 for assessments of personality traits, self-esteem, life satisfaction, and affect. Prior to
 1280 any data processing or analysis we excluded from the data set retest correlations that
 1281 were computed from samples that had missing sample size information ($n = 4$), and
 1282 where respondents were on average below 10 years of age or above 90 years of age ($n =$
 1283 31) leaving a total of 949 test-retest correlations (personality = 226, self-esteem = 196,
 1284 affect = 101, life satisfaction = 426) for analysis. To remain consistent with how we
 1285 analysed the other set of retest correlations, prior to estimating the model parameters,
 1286 we first:

1287 a) calculated the sampling variance of each correlation using the following
 1288 formula,

$$\frac{\sqrt{(1 - retest^2)^2}}{n - 1} \quad (1)$$

- 1289 b) centered the age variable at 40 years and transformed it into decades,
 1290 c) centered the proportion of females variable at 0.5, and
 1291 d) rounded the retest interval variable to .25 (i.e., 3 months bins).

1292 Given that in the data set close to 80% of the studies/samples had 4 or less
 1293 observations, to avoid poor estimation of varying intercepts and slopes as well as model
 1294 convergence issues, we did not specify a random effects structure for the *rel* parameter.

1295 By following these data processing and analysis steps we deviated from the
 1296 original study's analysis in four ways. First, we used a smaller data set. Second, we
 1297 carried out the analysis using a Bayesian instead of a Frequentist approach. Third,
 1298 when conducting the meta-analysis we accounted for the correlations' standard error.
 1299 Lastly, we changed the moderators that were included in the model by adding an
 1300 interaction between age linear and construct, between age quadratic and construct, and

1301 removing the effect of scale length on the *rel* parameter.

1302 Details of the model specification in *brms*, model fit and convergence statistics
 1303 are provided in the companion website.

1304 ***Convergent Validity***

1305 **Variance decomposition.** To gain a better understanding of the
 1306 heterogeneity in the correlation between different measures, we conducted a variance
 1307 decomposition analysis. We computed the Shapley values of the following predictors:

1308 Measure characteristics

- 1309 • Measure category match: whether or not both measures belong to the same
 1310 category (i.e., propensity, frequency, behaviour)
- 1311 • Domain match: whether or not both measures focus on the same life domain (e.g.,
 1312 smoking, driving, social, ethical)
- 1313 • Scale type match: whether or not both measures have the same type of scale (i.e.,
 1314 open-ended/composite index, ordinal/discrete scales)
- 1315 • Reliability: the average reliability of the measures (using MASC model parameter
 1316 estimates to make measure and age-specific predictions)

1317 Respondent characteristics

- 1318 • Age: age group the respondents belong in (10 year bins)
- 1319 • Gender: gender of the respondents (i.e., female, male)
- 1320 • Number of responses: sample size for each correlation

1321 **Meta-analyses.** Using the aggregated Fisher's z-transformed correlations, we
 1322 conducted a Bayesian random-effects meta-analysis to quantify the convergence across
 1323 all measures, and followed a distributional modelling approach by allowing the samples
 1324 to vary in their residual standard deviation (σ).

$$y_i \sim StudentT(\nu, \theta_i, \sqrt{se_i^2 + \sigma_i^2})$$

$$\theta_i \sim Normal(\mu_\theta, \tau_\theta)$$

$$\mu_\theta \sim Normal(0, 1)$$

$$\tau_\theta \sim Cauchy_+(0, 0.3)$$

$$\log\sigma_i \sim Normal(\mu_\sigma, \tau_\sigma)$$

$$\mu_\sigma \sim Normal(0, 2)$$

$$\tau_\sigma \sim Cauchy_+(0, 0.3)$$

$$\nu \sim Gamma(2, 0.1)$$

1325 Second, we conducted two meta-regressions with categorical covariates to
 1326 estimate the convergence between a) different pairs of measure categories (e.g.,
 1327 frequency and propensity),

$$y_i \sim StudentT(\nu, \theta_i, \sqrt{se_i^2 + \sigma_i^2})$$

$$\theta_i = \beta_{\theta_0, sample[i]} + \beta_{\theta_1} CategoryPair$$

$$\beta_{\theta_1} \sim Normal(0, 1)$$

$$\beta_{\theta_0, sample} \sim Cauchy_+(0, 0.3)$$

$$\log\sigma_i = \beta_{\sigma_0, sample[i]} + \beta_{\sigma_1} CategoryPair$$

$$\beta_{\sigma_1} \sim Normal(0, 2)$$

$$\beta_{\sigma_0, sample} \sim Cauchy_+(0, 0.3)$$

$$\nu \sim Gamma(2, 0.1)$$

1328 and, b) different domains (e.g., propensity-general and frequency-smoking).

$$y_i \sim StudentT(\nu, \theta_i, \sqrt{se_i^2 + \sigma_i^2})$$

$$\theta_i = \beta_{\theta_{0,sample[i]}} + \beta_{\theta_1} DomainPair$$

$$\beta_{\theta_1} \sim Normal(0, 1)$$

$$\beta_{\theta_{0,sample}} \sim Cauchy_+(0, 0.3)$$

$$\log\sigma_i = \beta_{\sigma_{0,sample[i]}} + \beta_{\sigma_1} DomainPair$$

$$\beta_{\sigma_1} \sim Normal(0, 2)$$

$$\beta_{\sigma_{0,sample}} \sim Cauchy_+(0, 0.3)$$

$$\nu \sim Gamma(2, 0.1)$$

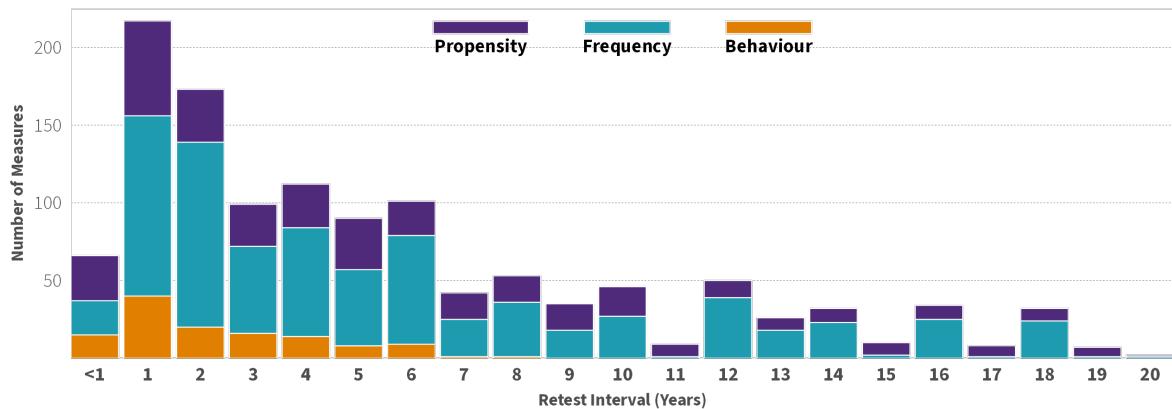
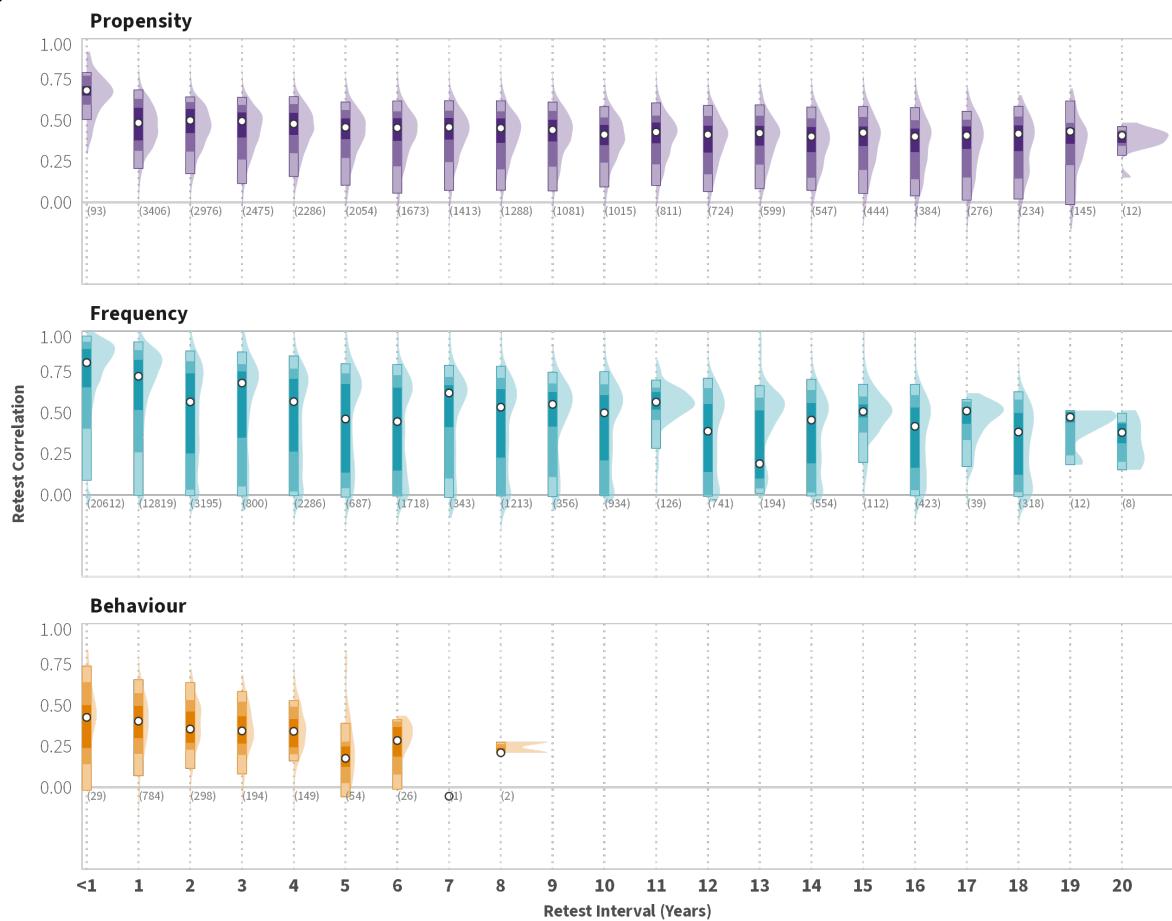
1329 In both meta-regressions we specified predictors for the residual standard
 1330 deviations, and allowed it to vary across the different levels of the categorical variables.
 1331 Based on recommendations, in all models, we used weakly informative priors (Williams
 1332 et al., 2018). Lastly, we back-transformed the results to Spearman's rho for the
 1333 reporting.

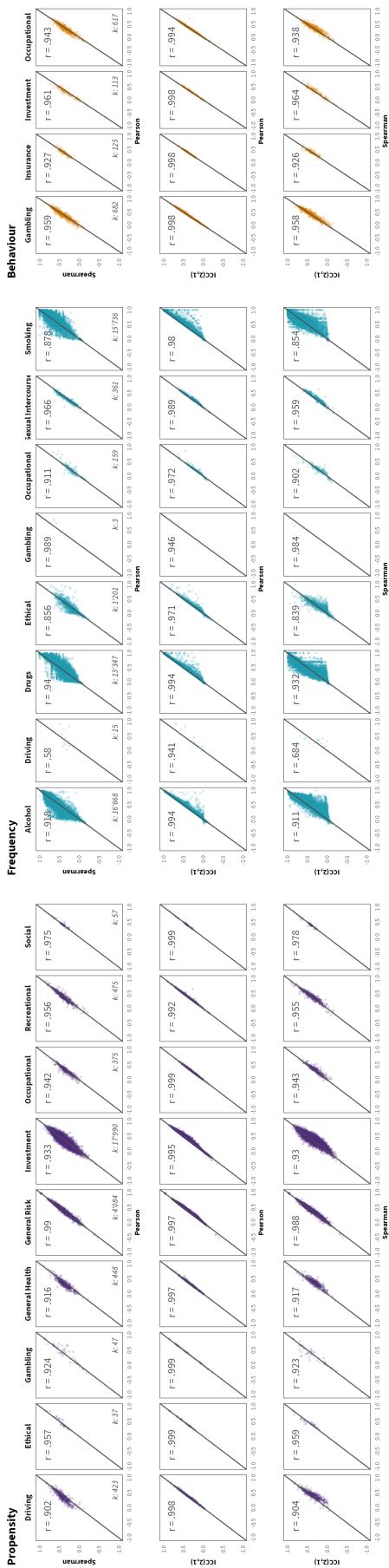
1334 Multiverse Analyses

1335 For brevity and ease of communication, we limited the reporting to a single data
 1336 set that was the result of a specific set of data pre-processing and processing choices. To
 1337 communicate transparently about our results and evaluate their robustness (i.e., how
 1338 sensitive results were to different data processing choices), we repeated our main
 1339 analyses using different data sets and model specifications (Steegen et al., 2016). On
 1340 the companion website we describe the different steps and choices that were available
 1341 when constructing and analysing the data, and include a visual summary of the
 1342 alternative results (Hall et al., 2022).

Figure S1

Overview of temporal stability measures and correlations. A) The number of measures by category (propensity, frequency, behaviour) and retest interval. B) Distributions of retest correlations as a function of retest interval for the different measure categories (propensity, frequency, behaviour).

A**B**



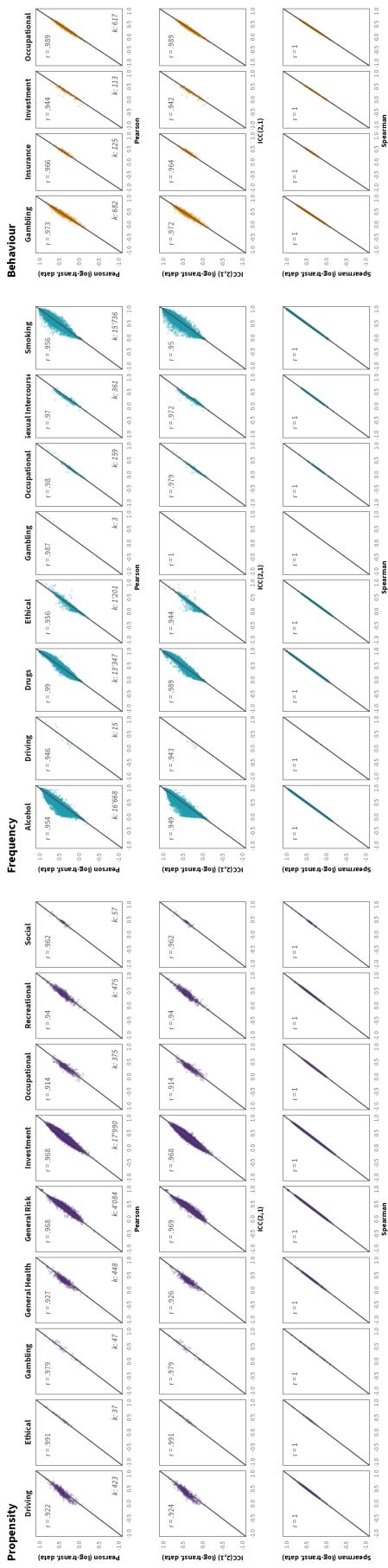


Figure S3

Scatter plots of different test-retest metrics calculated using either log-transformed or non-transformed data

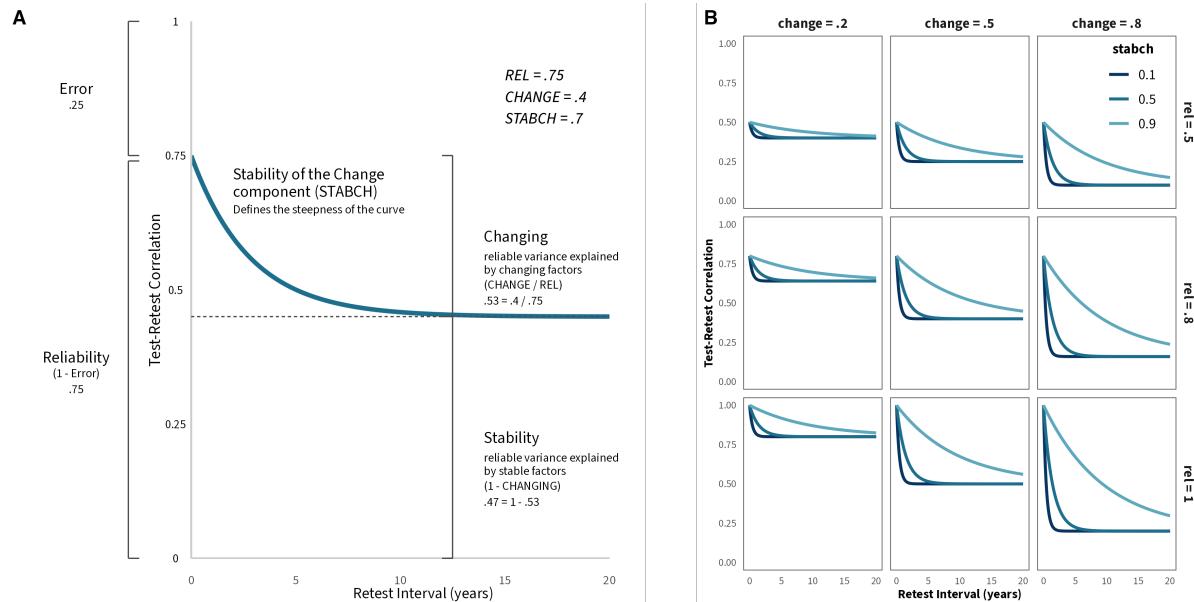
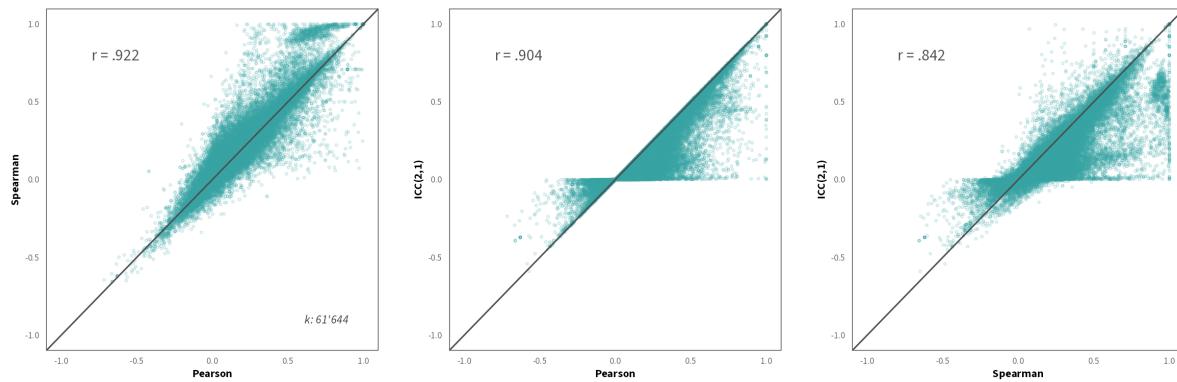
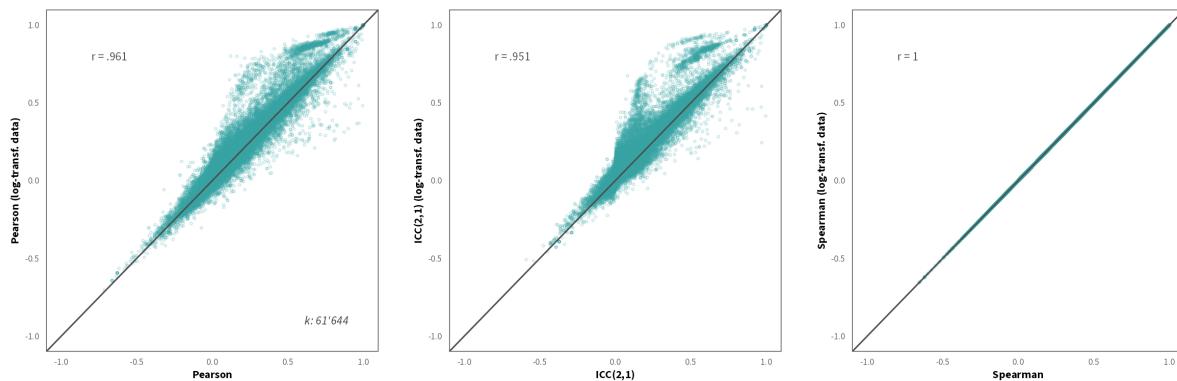


Figure S4

Depiction of the Meta-Analytic Stability and Change model (MASC). A) Visual depiction of temporal stability curve for major personality traits as estimated by Anusic and Schimmack (2016). B) Examples of different parameterisations of MASC.

**Figure S5**

Scatter plots of inter-correlations computed using Pearson's r , Spearman's ρ , or $ICC(2,1)$.

**Figure S6**

Scatter plots of different inter-correlation metrics calculated using either log-transformed or non-transformed data

Figure S7

Expected values of the posterior predictive distribution (mean, 50%, 80%, and 95% HDI) of Meta-Analytic Stability and Change model (MASC) parameters and test-retest correlations for propensity measures of risk preference in the general ($k = 1,732$), investment ($k = 1,080$), and driving ($k = 196$) domains. Left: Predicted values of the Reliability, Change, and Stability of Change parameters, split by domain, age group and gender. Right: Predicted test-retest correlations as a function of time for different age groups (upper panels) and as a function of age for different retest intervals (lower panels).

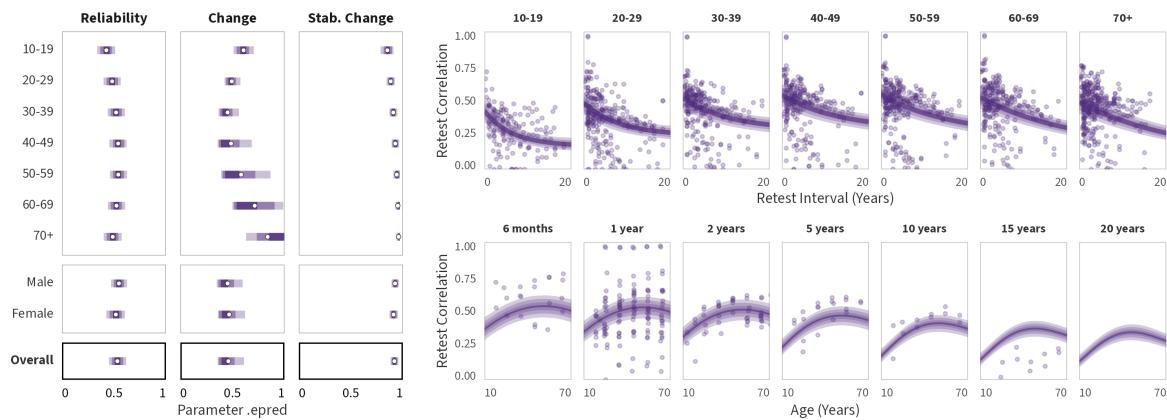
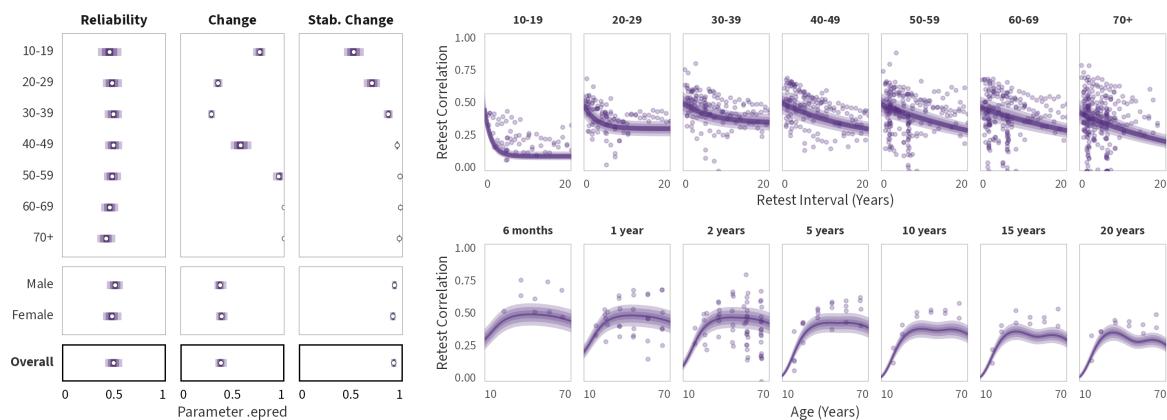
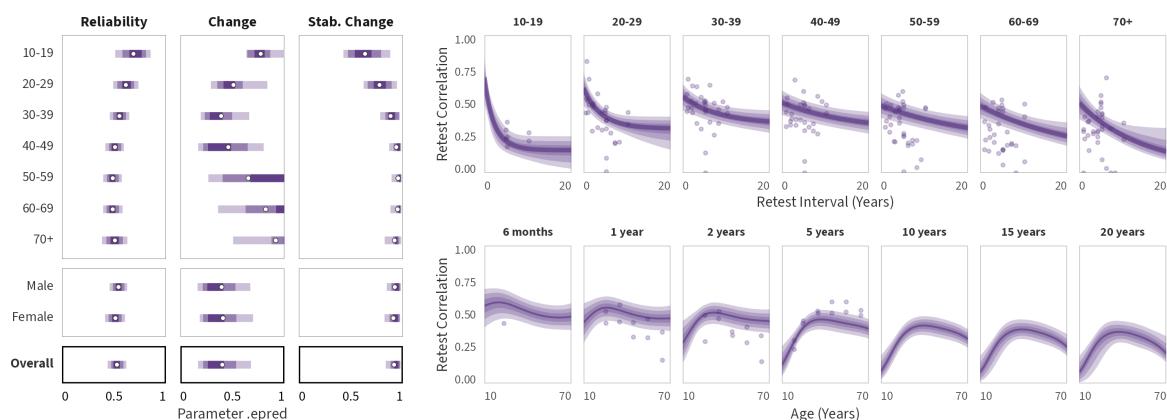
General Risk**Investment****Driving**

Figure S8

Expected values of the posterior predictive distribution (mean, 50%, 80%, and 95% HDI) of Meta-Analytic Stability and Change model (MASC) parameters and test-retest correlations for propensity measures of risk preference in the ethical ($k = 21$), gambling ($k = 38$), and general health ($k = 209$) domains. Left: Predicted values of the Reliability, Change, and Stability of Change parameters, split by domain, age group and gender. Right: Predicted test-retest correlations as a function of time for different age groups (upper panels) and as a function of age for different retest intervals (lower panels).

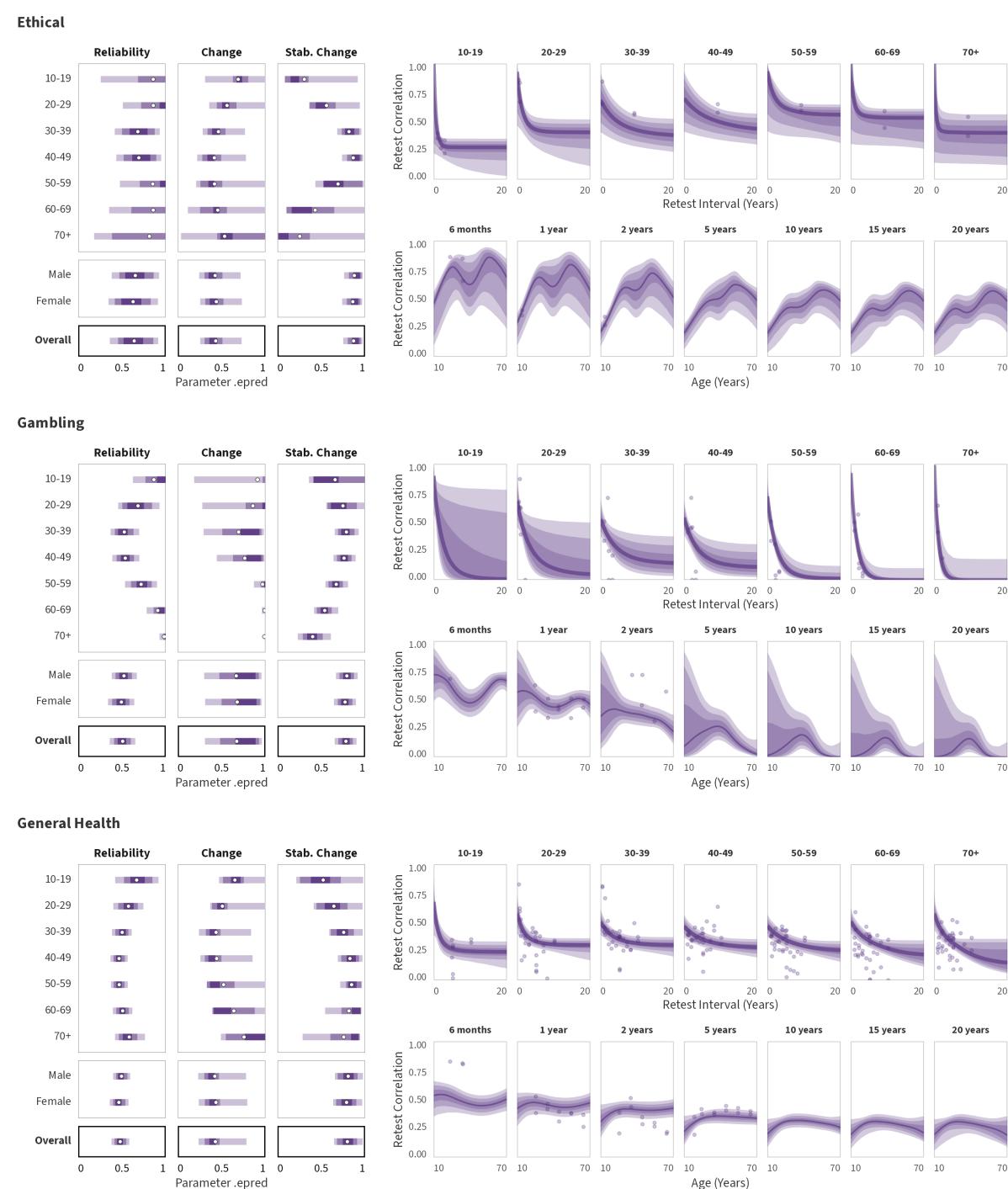


Figure S9

Expected values of the posterior predictive distribution (mean, 50%, 80%, and 95% HDI) of Meta-Analytic Stability and Change model (MASC) parameters and test-retest correlations for propensity measures of risk preference in the occupational ($k = 181$), recreational ($k = 198$), and social ($k = 51$) domains. Left: Predicted values of the Reliability, Change, and Stability of Change parameters, split by domain, age group and gender. Right: Predicted test-retest correlations as a function of time for different age groups (upper panels) and as a function of age for different retest intervals (lower panels).

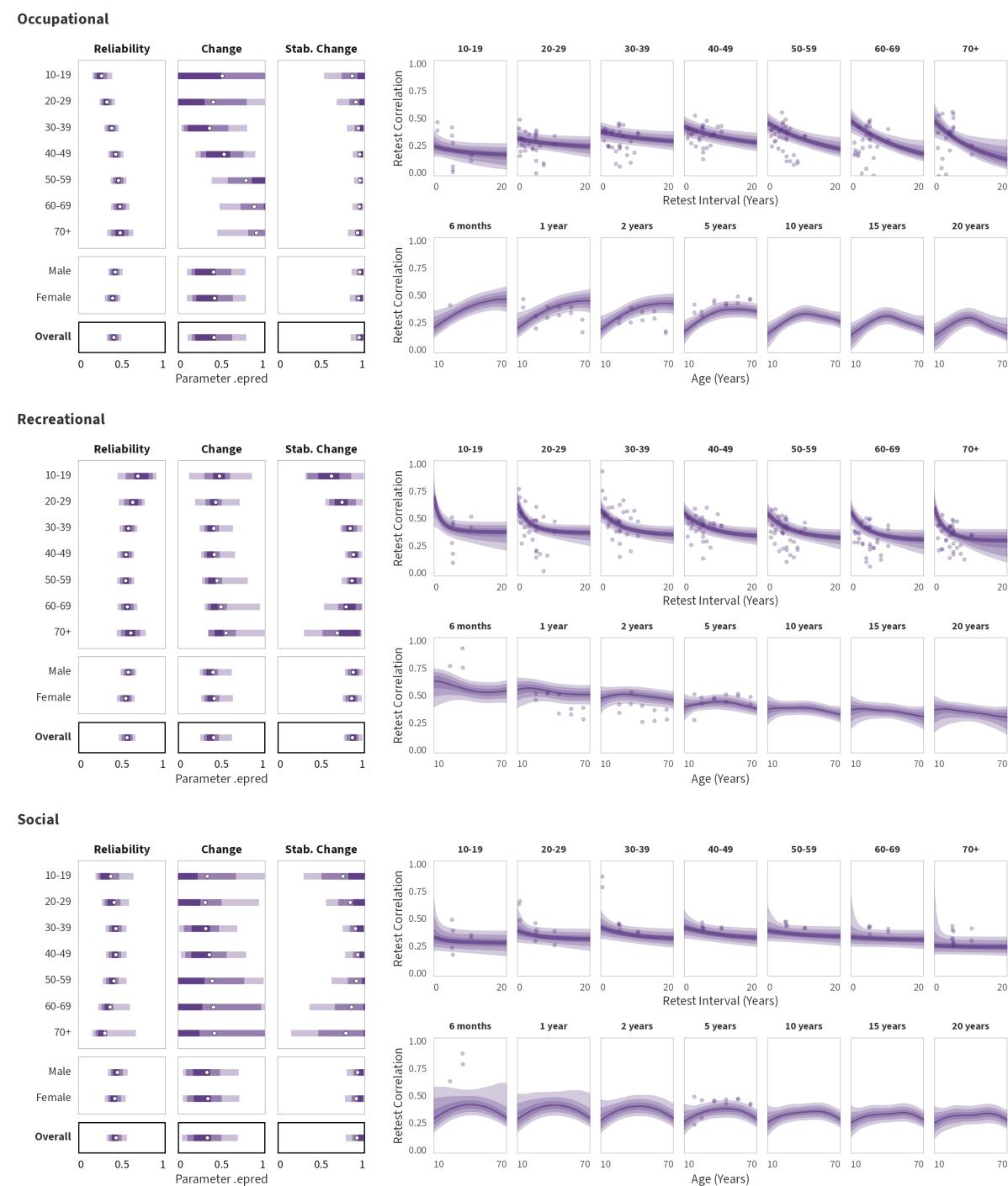


Figure S10

Expected values of the posterior predictive distribution (mean, 50%, 80%, and 95% HDI) of Meta-Analytic Stability and Change model (MASC) parameters and test-retest correlations for frequency measures of risk preference in the alcohol ($k = 1,609$), driving ($k = 15$), drugs ($k = 223$), and ethical ($k = 92$) domains. Left: Predicted values of the Reliability, Change, and Stability of Change parameters, split by domain, age group and gender. Right: Predicted test-retest correlations as a function of time for different age groups (upper panels) and as a function of age for different retest intervals (lower panels).

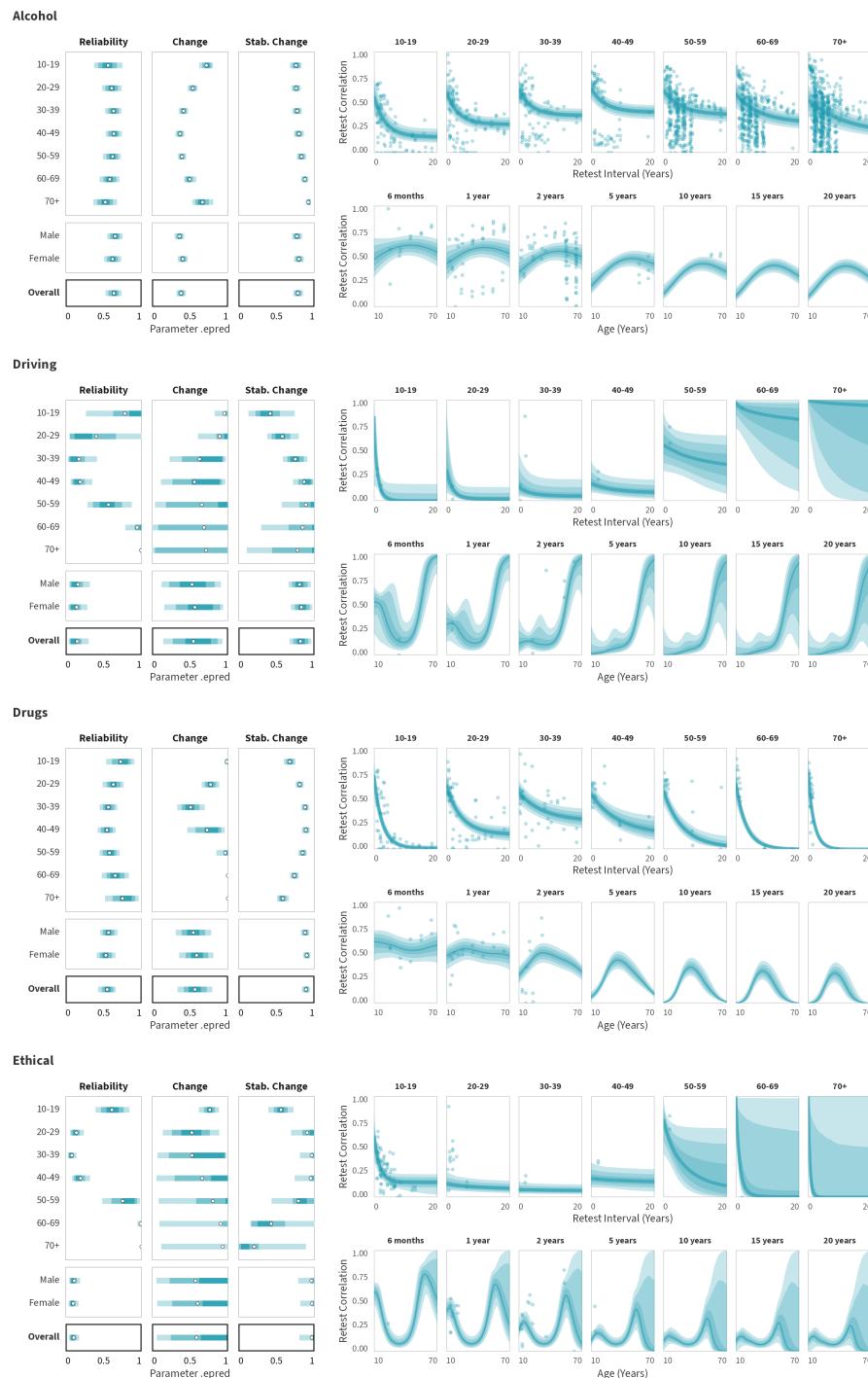


Figure S11

Expected values of the posterior predictive distribution (mean, 50%, 80%, and 95% HDI) of Meta-Analytic Stability and Change model (MASC) parameters and test-retest correlations for frequency measures of risk preference in the smoking ($k = 1,637$), sexual intercourse ($k = 82$), gambling ($k = 3$), and occupational ($k = 17$) domains. Left: Predicted values of the Reliability, Change, and Stability of Change parameters, split by domain, age group and gender. Right: Predicted test-retest correlations as a function of time for different age groups (upper panels) and as a function of age for different retest intervals (lower panels).

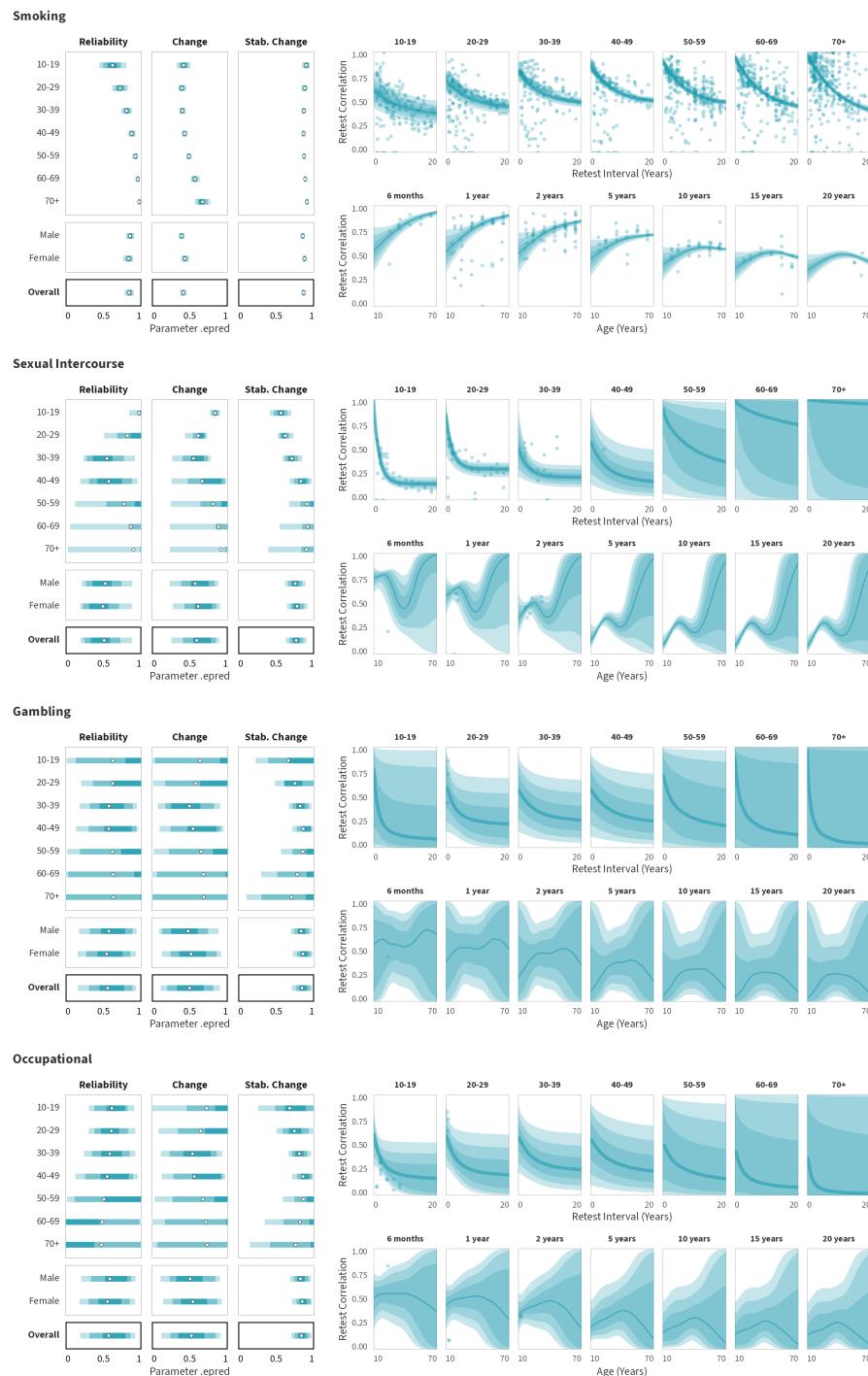


Figure S12

Expected values of the posterior predictive distribution (mean, 50%, 80%, and 95% HDI) of Meta-Analytic Stability and Change model (MASC) parameters and test-retest correlations for propensity measures of risk preference in the investment ($k = 108$), occupational ($k = 227$), gambling ($k = 197$), and insurance ($k = 80$) domains. Left: Predicted values of the Reliability, Change, and Stability of Change parameters, split by domain, age group and gender. Right: Predicted test-retest correlations as a function of time for different age groups (upper panels) and as a function of age for different retest intervals (lower panels).



Figure S13

Expected values of the posterior predictive distribution (mean, 50%, 80%, and 95% HDI) of Meta-Analytic Stability and Change model (MASC) parameters and test-retest correlations for personality ($k = 226$), affect ($k = 101$), life satisfaction ($k = 426$), and self-esteem ($k = 196$). Left: Predicted values of the Reliability, Change, and Stability of Change parameters, split by domain, age group and gender. Right: Predicted test-retest correlations as a function of time for different age groups (upper panels) and as a function of age for different retest intervals (lower panels).

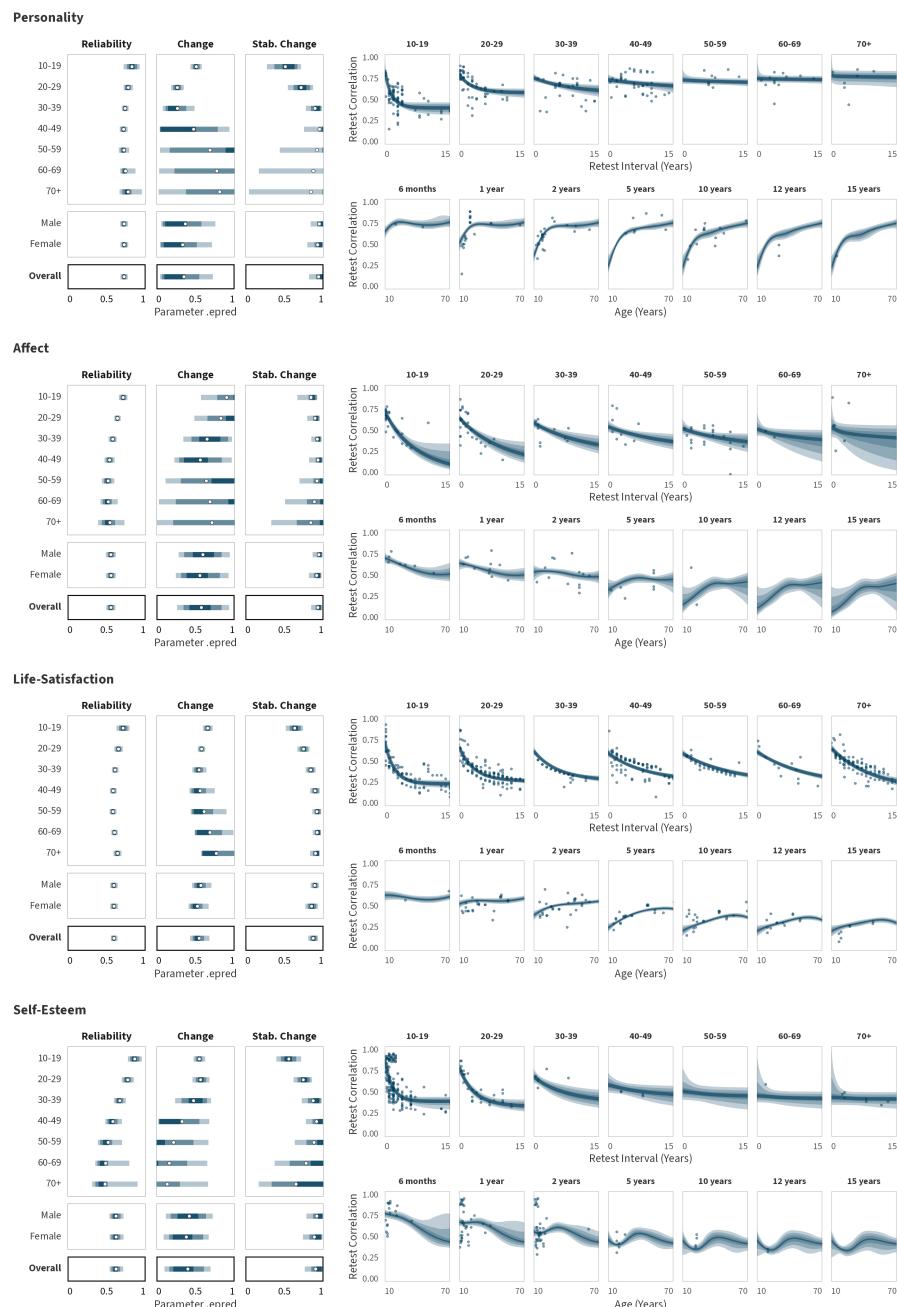


Figure S14

Convergence of risk preference measures. Distributions of inter-correlations between different risk preference measures at the same measurement occasion ($k = 61'644$), split by category-domain pairs (A), and category pairs (B).

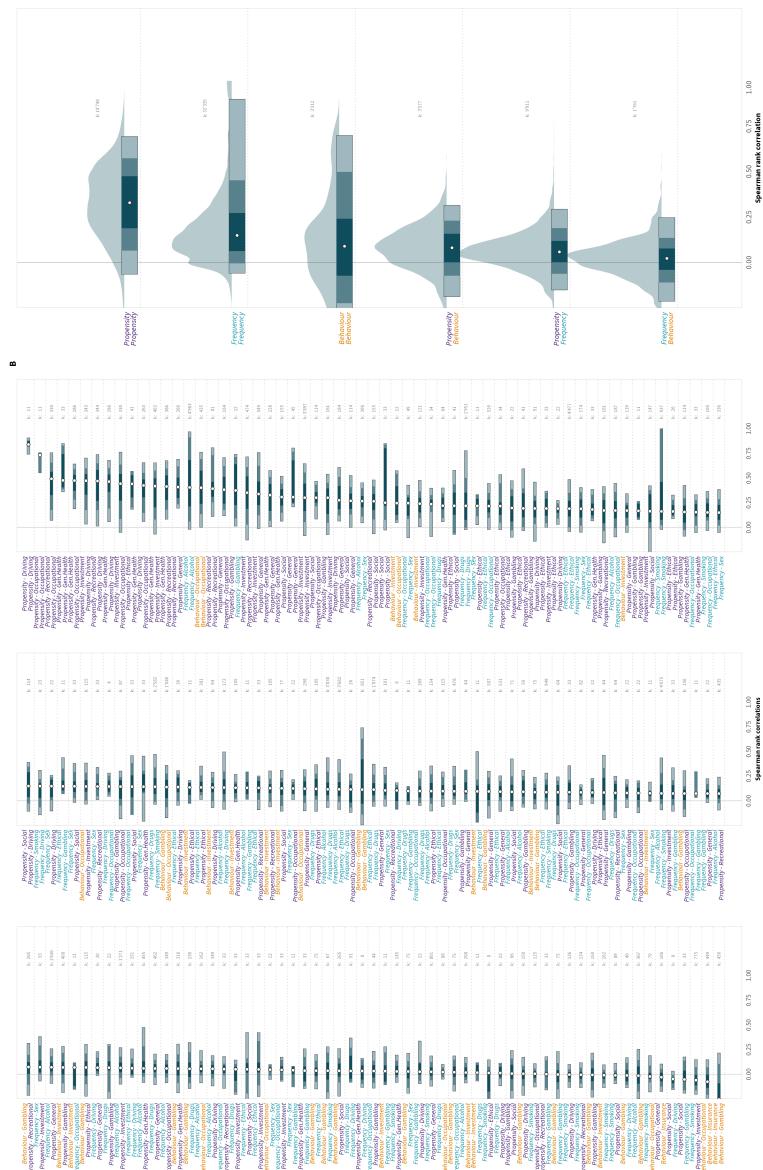


Table S1*Search terms used to identify risk preference measures*

Search terms
risk ; attitude ; loss/losing/lose ; excit* ; danger* ; avers* ; chance* ; certain ; safe* ; fear* ; adventure/venture ; impuls* ; prefer* ; careful driv* ; car ; fast ; speed ; motor* ; traffic vandal* ; damage ; cheat ; police ; convict* ; arrest* ; gun ; weapon ; shoot/shot ; troubl* ; stole/steal ; lie ; crim* ; delinquen* ; aggressive ; fight* ; assault ; violen* ; injur* ; bully ; affair ; *faith* finan* ; gambl* ; lottery ; coin ; invest* ; stocks ; bet* ; casino ; fund ; poker ; trad* ; shares ; bonds health ; drug ; alcohol ; smok* ; drink* ; cigarette ; drunk ; intox* ; marijuana/cannabis ; heroin ; meth* ; cocaine ; stimulant ; ecstasy ; hallucinogen ; tobacco ; wine ; liquor ; spirit ; beer/pint; unprotected sex/intercourse ; contracept* occupation ; career ; job ; self-employ* ; employ* ; business ; work* ; entrepreneur extreme ; sport ; bar/pub ; night* ; mountaint* ; skydiving ; bunjee ; ski ; climb* ; race stranger ; trust

Table S2
Overview of panels screened

Panel/Sample	Status	Reason for exclusion, if applicable
Adema, Nikolla, Poutvaara Sunde (2022); Economics Letters (ANPS) - Czech Republic sample	Incl.	
Adema, Nikolla, Poutvaara Sunde (2022); Economics Letters (ANPS) - India sample	Excl.	Small sample size
Adema, Nikolla, Poutvaara Sunde (2022); Economics Letters (ANPS) - Mexico sample	Excl.	Small sample size
Adema, Nikolla, Poutvaara Sunde (2022); Economics Letters (ANPS) - Spain sample	Incl.	
American Life Panel (ALP)	Incl.	
American National Election Studies (ANEES)	Excl.	Does not include propensity item (that meets criteria)
Americans' Changing Lives study (ACL)	Excl.	Does not include propensity item (that meets criteria)
Basel-Berlin Risk Study (BRRS) - Basel (From Frey et al., 2017 Science Advances)	Incl.	
Basel-Berlin Risk Study (BRRS) - Berlin (From Frey et al., 2017 Science Advances)	Incl.	
Berlin Aging Study (BASE)	Excl.	Restricted data access
Berlin Aging Study-II (BASE-II)	Excl.	Restricted data access
British Election Study 2005-2009 (BES05)	Incl.	
British Election Study 2014-2023 (BES14)	Incl.	
Bundesbank - Panel of Household Finances (PHF)	Incl.	
Bundesbank - Survey on Consumer Expectations	Excl.	Propensity item not asked repeatedly
Bundesbank Online Panel – Haushalte (BOP-HH)	Excl.	Propensity item not asked repeatedly
California Families Project (CFP)	Excl.	Restricted data access
Canadian Longitudinal Study on Aging (CLSA)	Excl.	Does not include propensity item (that meets criteria)
Cape Area Panel Study (CAPS)	Excl.	Propensity item not asked repeatedly
China Health and Retirement Longitudinal Survey (CHARLS)	Excl.	Does not include propensity item (that meets criteria)
Cognition and Aging in the USA	Excl.	Propensity item not asked repeatedly
Collaborative Studies on the Genetics of Alcoholism (COGA)	Excl.	Does not include propensity item (that meets criteria)
Costa Rican Longevity and Healthy Aging Study (CRELES)	Excl.	Does not include propensity item (that meets criteria)
Crime in the Modern City. A Longitudinal Study of Juvenile Delinquency in Münster (CMC)	Incl.	
DNB Household Survey (DHS)	Incl.	
Drichoutis Vassilopoulos (2021); Journal of Economics Management Strategy	Incl.	
Einstein Aging Study (EAS)	Excl.	Restricted data access
English Longitudinal Study of Ageing (ELSA)	Excl.	Propensity item not asked repeatedly
Enkavi et al., 2019 PNAS	Incl.	
Financial Crisis: A Longitudinal Study of Public Response (FICR)	Excl.	Does not include propensity item (that meets criteria)
Fragile Families and Child Wellbeing Study (FFCWS)	Excl.	Does not include propensity item (that meets criteria)
General Social Survey Panel (GSSP)	Excl.	Does not include propensity item (that meets criteria)
German Internet Panel (GIP)	Incl.	
German Longitudinal Election Study (GLES) - Panel 2016-2021	Incl.	

Table S2 cont.

Panel/Sample	Status	Reason for exclusion, if applicable
German Longitudinal Election Study (GLES) - Long-term Online Tracking, Cumulation	Excl.	Does not include propensity item (that meets criteria)
German Longitudinal Election Study (GLES) - Long-term Panel 2002-2005-2009	Excl.	Does not include propensity item (that meets criteria)
German Longitudinal Election Study (GLES) - Long-term Panel 2005-2009-2013	Excl.	Does not include propensity item (that meets criteria)
German Longitudinal Election Study (GLES) - Long-term Panel 2009-2013-2017	Excl.	Does not include propensity item (that meets criteria)
German Longitudinal Election Study (GLES) - Short-term Campaign Panel 2009	Excl.	Propensity item not asked repeatedly
German Longitudinal Election Study (GLES) - Short-term Campaign Panel 2013	Excl.	Does not include propensity item (that meets criteria)
German Longitudinal Election Study (GLES) - Short-term Campaign Panel 2013-2017 (repeatedly questioned respondents)	Incl.	
German Longitudinal Election Study (GLES) - Short-term Campaign Panel 2017	Excl.	Propensity item not asked repeatedly
Health and Aging in Africa: A Longitudinal Study of an INDEPTH Community in South Africa (HAALSI)	Excl.	Does not include propensity item (that meets criteria)
Health Retirement Survey (HRS)	Incl.	
Health Retirement Survey: Cognitive Economics Project (CogEcon)	Incl.	
Health, Aging, and Retirement in Thailand (HART)	Excl.	Does not include propensity item (that meets criteria)
Healthy Ageing in Scotland (HAGIS)	Excl.	Only W1/Pilot data available
High School and Beyond (HSB)	Excl.	Does not include propensity item (that meets criteria)
Household Finance and Consumption Survey (HFCS)	Excl.	Does not include propensity item (that meets criteria)
Household, Income and Labour Dynamics in Australia (HILDA)	Incl.	
Indonesia Family Life Survey (IFLS)	Excl.	Does not include propensity item (that meets criteria)
Interdisciplinary Longitudinal Study of Adult Development (ILSE and ILSE.Y)	Excl.	Restricted data access
Japan Household Panel Survey (JHPS)	Excl.	Does not include propensity item (that meets criteria)
Japanese Study of Aging and Retirement (JSTAR)	Excl.	Does not include propensity item (that meets criteria)
Korean Labour Income Panel Survey	Excl.	Does not include propensity item (that meets criteria)
Korean Longitudinal Study of Aging (KLoSA)	Excl.	Does not include propensity item (that meets criteria)
Life in Kyrgyzstan Study (LJKS)	Incl.	
Longitudinal Aging Study in India (LASI)	Excl.	Only W1/Pilot data available
Longitudinal Aging Study of Amsterdam (LASA)	Excl.	Does not include propensity item (that meets criteria)
Longitudinal Internet studies for the Social Sciences (LISS)	Excl.	Propensity item not asked repeatedly
Longitudinal Study of American Youth (LSAY)	Excl.	Does not include propensity item (that meets criteria)
Longitudinal Study of Australian Children (LSAC)	Excl.	Propensity item not asked repeatedly
Longitudinal Study of Violence Against Women - Men Sample (LSVAW-M)	Incl.	
Longitudinal Study of Violence Against Women - Women Sample (LSVAW-W)	Excl.	Does not include propensity item (that meets criteria)
Longitudinal Surveys of Australian Youth (LSAY)	Excl.	Restricted data access
Lothian Birth Cohort 1936	Excl.	Only W1/Pilot data available
Malaysia Ageing and Retirement Survey (MARS)	Incl.	
Medical Expenditure Panel Survey (MEPS)	Excl.	Propensity item not asked repeatedly
Mexican Family Life Survey (MexFLS)	Excl.	Does not include propensity item (that meets criteria)
Mexican Health and Aging Study (MHAS)	Excl.	

Table S2 cont.

Panel/Sample	Status	Reason for exclusion, if applicable
Midlife in Japan (MIDJA)	Incl.	
Midlife in the United States (MIDUS) - Milwaukee Dample	Excl.	Restricted data access
Midlife in the United States (MIDUS) - Project 1 Sample	Incl.	
Millennium Cohort Study (MCS)	Excl.	Propensity item not asked repeatedly
Monitoring the Future: Restricted-Use Panel Data	Excl.	Restricted data access
National Health and Nutrition Examination Survey (NHANES)	Excl.	Does not include propensity item (that meets criteria)
National Income Dynamics Study (NIDS)	Excl.	Does not include propensity item (that meets criteria)
National Longitudinal Study of Adolescent to Adult Health (Add Health)	Incl.	
National Longitudinal Survey of Youth 1979 (NLSY79)	Excl.	Propensity item not asked repeatedly
National Longitudinal Survey of Youth 1979 Child and Young Adult (NLSY79-CYA)	Incl.	
National Longitudinal Survey of Youth 1997 (NLSY97)	Excl.	Propensity item not asked repeatedly
National Social Life, Health, and Aging Project (NSHAP)	Incl.	
National Survey of Families and Households (NSFHS)	Excl.	Does not include propensity item (that meets criteria)
New Zealand Health, Work and Retirement Study	Excl.	Does not include propensity item (that meets criteria)
Nießen et al. (2020) . GESIS Instrument	Excl.	Small sample size
Northern Ireland Cohort for the Longitudinal Study of Ageing (NICOLA)	Excl.	Propensity item not asked repeatedly
Origin of Variance in the Oldest-Old: Octogenarian Twins (Octo-Twin)	Excl.	Does not include propensity item (that meets criteria)
Panel Study of Income Dynamics (PSID)	Excl.	Propensity item not asked repeatedly
Panel Survey of Consumer Finances 1983-1989	Excl.	Cannot match respondents
Panel Survey of Consumer Finances 2007-2009	Excl.	Does not include propensity item (that meets criteria)
Parenting Across Cultures	Excl.	Does not include propensity item (that meets criteria)
Preference Parameters Study - India (rural area) (GCOE - IN Rural)	Incl.	
Preference Parameters Study - India (urban area) (GCOE - IN)	Excl.	Does not include propensity item (that meets criteria)
Preference Parameters Study - China (urban area) (GCOE - CN)	Incl.	
Preference Parameters Study - Japan (GCOE - JP)	Incl.	
Preference Parameters Study - United States of America (GCOE - USA)	Excl.	Propensity item not asked repeatedly
Public Opinion and the Syrian Crisis in Three Democracies	Excl.	Does not include propensity item (that meets criteria)
Risky decision and happiness task: The Great Brain Experiment smartphone app	Excl.	Does not include propensity item (that meets criteria)
Rochester Adult Longitudinal Study (RALS)	Excl.	Missing documentation
Rural-Urban Migration in China and Indonesia: CHINA	Excl.	Propensity item not asked repeatedly
Rural-Urban Migration in China and Indonesia: INDONESIA	Excl.	Restricted data access
Russian Longitudinal Monitoring Survey (RLMS-HSE)	Excl.	Restricted data access
Screening Across the Lifespan Twin Study: the Younger (SALTY)	Excl.	
Seattle Longitudinal Study (SLS)	Excl.	
Socio-Economic Panel Study - Core (SOEP-Core SOEP-CoV)	Incl.	
Socio-Economic Panel Study Retest (SOEP-Retest)	Excl.	Small sample size
Sparen und Altersvorsorge in Deutschland (SAVE)	Incl.	

Table S2 cont.

Panel/Sample	Status	Reason for exclusion, if applicable
Steiner et al., (2020); Decision Studies Incl. in Entkavi et al. (2019 PNAS) meta-analysis	Excl.	Small sample size
Studies Incl. in Mata et al. (2018 JEP) meta-analysis	Excl.	No open access data
Study to Assess Risk and Resilience in Servicemembers — Longitudinal Study (STARRS)	Excl.	No open access data
Survey of Consumer Expectations (SCE)	Excl.	Propensity item not asked repeatedly
Survey of Health, Ageing and Retirement in Europe (SHARE) (Excluding the following countries: Bulgaria, Croatia, Cyprus, Finland, Greece, Hungary, Latvia, Lithuania, Luxembourg, Malta, Poland, Portugal, Romania, Slovak Republic)	Excl.	Propensity item not asked repeatedly
Survey of Health, Ageing and Retirement in Europe (SHARE) (Including the following countries: Austria, Belgium, Czech_Rep, Denmark, Estonia, France, Germany, Israel, Italy, Netherlands, Slovenia, Spain, Sweden, Switzerland)	Incl.	
Swedish Adoption/Twin Study of Aging (SATSA)	Excl.	Small sample size
Swiss Household Panel (SHP)	Excl.	Propensity item not asked repeatedly
Thailand Vietnam Socio Economic Panel (TVSEP)	Excl.	Does not include propensity item (that meets criteria)
The Brazilian Longitudinal Study of Aging (ELSI-Brazil)	Excl.	Does not include propensity item (that meets criteria)
The Irish Longitudinal Study on Ageing (TILDA)	Excl.	Does not include propensity item (that meets criteria)
Tracking Adolescents' Individual Lives Survey (TRAILS)	Excl.	Does not include propensity item (that meets criteria)
TwinLife	Incl.	
Twins of Early Development Study (TEDS)	Excl.	Propensity item not asked repeatedly
UK Biobank	Excl.	Does not include propensity item (that meets criteria)
UK Household Longitudinal Survey + British Household Panel Survey (USOC)	Excl.	Propensity item not asked repeatedly
Ukrainian Longitudinal Monitoring Survey (ULMS)	Incl.	
Understanding America Study (UAS)	Incl.	
VA Normative Aging Study (VA_NAS)	Excl.	Does not include propensity item (that meets criteria)
WHO Study on global AGEing and adult health (SAGE)	Excl.	Does not include propensity item (that meets criteria)
Wisconsin Longitudinal Study (WLSG/WLSS)	Excl.	Does not include propensity item (that meets criteria)
Work and Family Life Study	Excl.	Propensity item not asked repeatedly

End of Table

Table S3*Overview of exclusion and inclusion criteria of measures for the analyses, split by measure category*

Category	Inclusion	Exclusion	Rationale
All	1. Measures that have been asked to the same respondents across at least two time points.	1. Measures that have been asked only in one wave or only once to the respondents	1. We need responses from at least two time points two compute a test-retest correlation coefficient.
All	2. Measures where the wording and response format remained consistent across at least two time points.	2. Measures that are not consistent across at least two time points	2. Measures need to be the same across waves to accurately measure test-retest correlations
All	3. Measures that include at least 4 response options/values, or is composed of multiple (binary) measures that can be aggregated to calculate an index.	3. Measures that include less than four response options/values (e.g., yes/no, never/sometimes/always).	3. With more response options it is possible to capture more meaningful changes over time .
All	4. Measures that use an ordinal scale, discrete scale (with a clear response range) or are open-ended	4. Measures that use a nominal scale or scales with options that cannot be objectively ranked	4. Can result in subjective interpretations of what a category is and thus reduces response comparability between participants. Further if response options cannot be ranked, this can reduce the accuracy of how the test-retest correlations are computed.
Propensity	1. Measures that ask respondents about recent behaviour.	1. Measures that ask respondents about behaviour that is too far back in time or no longer relevant (e.g., asking adult respondents about their risk propensity as a child).	1. Relies on the recollection of certain events, which can result in inaccuracies. We are not capturing temporal stability based on the responses of actions that are no longer relevant .
Propensity	2. Measures that refer directly to the respondent.	2. Measures that refer to an individual other than the respondent (e.g., partner/spouse, household)	2. Another person's or group's risk preference is not necessarily reflective of the respondent's. Thus, individual changes would not be reflected in the response.
Propensity	3. Measures that can be answered by both women and men	3. Gender-specific measures (e.g., specific behaviour during pregnancy)	3. We want to collect approximately the same amount of responses from both males and females respondents to best explore gender differences.
Propensity	4. Measures that explicitly ask about risk-taking.	4. Measures that ask about ambiguity.	4. Ambiguity preference is shown to differ from risk preference (Levy et al., 2010)
Propensity	5. Measures that can be classified into a general or single life domain (e.g., general, driving, recreational)	5. Measures for which the behaviour cannot be classified into more than one pre-specified domain	5. More accurate comparison across domains
Frequency	1. Measures that ask respondents about recent or ongoing behaviour.	1. Measures that ask respondents about behaviour that is too far in time or no longer relevant (e.g., number of cigarettes smoked before quitting).	1. Relies on the recollection of certain events, which can result in inaccuracies. Asking about behaviours that are no longer taking place in the present can result in inflated correlation coefficients.
Frequency	2. Measures with a clearly specified time frame (e.g., in the last month/week how often...).	2. Measures with no clearly specified time frame or that refer to the course of the respondent's life time or that are dependent on a specific event (e.g., since you were 14 years old).	2. Such questions do not allow a proper comparison between participants as these can result in the subjective interpretation of a time frame or they are dependent on other factors (e.g., current age).

Table S3 cont.

Category	Inclusion	Exclusion	Rationale
Frequency	3. Measures that refer directly to the respondent.	3. Measures that refer to an individual other than the respondent (e.g., partner/spouse, household)	3. Another person's or group's risk preference is not necessarily reflective of the respondent's. Thus, individual changes would not be reflected in the response.
Frequency	4. Measures that use an ordinal scale, discrete scale (with a clear response range) or are open-ended	4. Measures that use a nominal scale or scales than cannot be objectively ranked	4. Can result in subjective interpretations of what a category is and thus reduces response comparability between participants. Further if response options cannot be ranked, this can reduce the accuracy of how the test-retest correlations are computed.
Frequency	5. Measures that include 0 or Never response options	6. Measures that do not include 0 or Never response options	5. It is possible to enter a response for those respondent whom this question does not apply (e.g., non-smokers smoking 0 cigarettes). Additionally, such measures help better capture changes across time (e.g., a frequent smoker at T1 but quits smoking at T2)
Frequency	5. Measures that can be answered by both women and men	6. Gender-specific measures (e.g., specific behaviour during pregnancy)	6. We want to collect the same amount of responses from both males and females respondents to best explore gender differences.
Frequency	6. Measures that can be classified into a single life domain (e.g., smoking, alcohol, driving)	6. Measures for which the behaviour can be classified into more than one life domain	6. More accurate comparison across domains
Behaviour	1. Measures with choices that vary on in terms of probabilities, or that have a clear risk component.	1. Measures with choices that not solely vary in terms of probabilities (e.g. choices dependent on the response of another individual, choices involving a dimension of time).	1. Including measures that vary on other dimensions of the choice options would result in risk preference being confounded by other preferences (e.g., social preference, time preference)
Behaviour	2. Measures with choices that involve a form of money outcome or reward.	2. Measures with choices in non-financial contexts with other forms of outcomes	2. Such measures allow a direct comparison to tasks commonly using the economics literature

End of Table

Table S4*Overview of panels included in the analyses*

Sample	Country	Collect	Oper.	Domains	N.meas.	N.waves	N.corr	N
ADDHEALTH	U.S.A.	Int.	F, P	Alc., Dri., Dru., Eth., Gen., Sex., Smo.	49	5	379	6,138
ALP	U.S.A.	Onl.	P, B	Gen., Inv., Gam., Occ.	11	18	215	3,180
ANPS-Czech-Republic	Czech Republic	Onl.	P, B	Gen., Inv.	2	2	4	230
ANPS-Spain	Spain	Onl.	P, B	Gen., Inv.	2	2	5	177
BBRS-CH	Switzerland	Lab.	F, P, B	Alc., Dri., Dru., Eth., Gam., Gen., Hea-gen., Inv., Occ., Rec., Sex., Soc.	35	2	35	34
BBRS-DE	Germany	Lab.	F, B, P	Alc., Eth., Sex., Occ., Gam., Dru., Dri., Gen., Inv., Hea-gen., Rec., Soc.	35	2	70	99
BES05	U.K.	Onl.	P	Gen.	1	2	12	3,291
BES14	U.K.	Onl.	P, B	Gen., Gam.	2	4	64	32,982
CMC	Germany	Int.	F, P	Eth., Dru., Occ.	25	4	223	2,017
COGECON	U.S.A.	Int.	P, B	Inv., Gen.	3	4	54	871
DHS	Netherlands	Int.	B, P	Gam., Gen., Inv.	7	30	14,161	10,581
DRICHOUTIS	Greece	Self-adm.	P, B	Gen., Inv.	2	3	10	113
ENKAVI	U.S.A.	Onl.	F, P, B	Alc., Dri., Dru., Eth., Gam., Hea-gen., Rec., Smo., Soc.	19	2	32	68
GCOE-CN	China	Int.	P	Gen.	1	2	10	958
GCOE-IN	India	Int.	P, B	Gen., Gam., Occ.	5	5	49	1,280
GCOE-JP	Japan	Self-adm.	P, B	Gen., Occ., Gam., Ins.	15	12	949	8,040
GCOE-USA	U.S.A.	Self-adm.	P, B	Gen., Occ., Gam., Ins.	15	9	684	7,523
GIP	Germany	Onl.	P	Gen.	1	3	32	2,129
GLES-LT	Germany	Int.	P	Gen.	1	6	130	17,320
GLES-ST	Germany	Onl.	P	Gen.	1	2	12	2,045
HILDA	Australia	Int.	P, F	Inv., Gen., Smo.	4	21	5,976	25,154
HRS-Core	U.S.A.	Int.	F, P, B	Alc., Dri., Gen., Hea-gen., Inv., Occ., Rec., Smo.	15	15	2,376	34,027
LIKS	Kyrgyzstan	Int.	F, P	Alc., Gen., Smo.	8	6	758	10,082
LSVAW-M	U.S.A.	Int.	F, P	Alc., Dru., Eth., Gen., Sex.	26	5	306	650
LSVAW-W	U.S.A.	Int.	F, P	Alc., Dru., Eth., Gen., Sex.	23	5	166	1,394
MEPS	U.S.A.	Int.	P	Gen.	1	34	272	157,599
MIDJA	Japan	Int.	P, F	Gen., Alc.	6	2	58	655
MIDUS-Project1	U.S.A.	Int.	F, P	Alc., Dru., Gen., Eth.	9	3	181	4,357
NLSY79-CYA	U.S.A.	Int.	F, P, B	Alc., Dru., Eth., Gen., Occ., Sex., Smo.	31	17	4,222	8,613
NSHAP	U.S.A.	Int.	F, P	Alc., Gen., Smo.	5	3	86	2,943
PHF	Germany	Int.	P	Inv., Gen.	2	3	56	3,566
SAVE	Germany	Self-adm.	F, P	Alc., Dri., Gam., Hea-gen., Inv., Occ., Rec.	9	10	1,895	3,758
SHARE-Austria	Austria	Int.	F, P	Alc., Inv., Smo.	7	7	148	4,863
SHARE-Belgium	Belgium	Int.	F, P	Alc., Inv., Smo.	7	7	191	6,544
SHARE-Czech-Rep	Czech-Rep	Int.	F, P	Alc., Inv., Smo.	6	6	159	5,673
SHARE-Denmark	Denmark	Int.	F, P	Alc., Inv., Smo.	8	7	183	4,249
SHARE-Estonia	Estonia	Int.	F, P	Alc., Inv., Smo.	6	4	80	6,214
SHARE-France	France	Int.	F, P	Alc., Inv., Smo.	7	7	183	5,593
SHARE-Germany	Germany	Int.	F, P	Alc., Inv., Smo.	7	7	160	5,463
SHARE-Israel	Israel	Int.	F, P	Alc., Inv., Smo.	7	5	68	2,665
SHARE-Italy	Italy	Int.	F, P	Alc., Inv., Smo.	7	7	185	5,251
SHARE-Netherlands	Netherlands	Int.	F, P	Alc., Inv., Smo.	7	5	97	3,796
SHARE-Slovenia	Slovenia	Int.	F, P	Alc., Inv., Smo.	6	4	82	3,729
SHARE-Spain	Spain	Int.	F, P	Alc., Inv., Smo.	7	7	174	6,310
SHARE-Sweden	Sweden	Int.	F, P	Alc., Inv., Smo.	7	7	167	4,869
SHARE-Switzerland	Switzerland	Int.	F, P	Alc., Inv., Smo.	7	7	170	3,442

Table S4 cont.									
Sample	Country	Collect	Oper.	Domains	N.meas.	N.waves	N. corr	N	
SOEP-Core	Germany	Int.	P, B, F	Dri., Gen., Hea-gen., Inv., Occ., Rec., Smo., Soc.	11	19	3,822	61,611	
TWINLIFE	Germany	Int.	F, P	Alc., Dri., Eth., Gen., Occ.	18	3	132	9,035	
UAS	U.S.A.	Onl.	F, P, B	Alc., Dru., Gen., Inv., Smo.	13	42	32,710	9,371	
ULMS	Ukraine	Int.	F, P, B	Alc., Dri., Gen., Hea-gen., Inv., Occ., Rec., Smo.	21	4	277	8,154	
USOC-IP	U.K.	Int.	F, B, P	Alc., Dru., Eth., Gam., Gen., Hea-gen., Inv., Smo.	12	13	493	3,707	

End of Table

1344 Notes. Mode of data collection: Onl(ine), Self-Adm(inistered), Lab(oratory), Int(erview). Measures: P(ropensity), F(requency),

1345 and B(ehaviour). Domains: Alc(ohol), Dri(ving), Dru(gs), Eth(ical), Gam(bling), Gen(erall), Hea(lth)-Gen(erall), Ins(urance), Inv(estment),

1346 Occ(upational), Rec(reative), Smok(ing), Soc(ial),

Table S5

Overview and description of the different risk preference measures included in the study, split by measure category, and domain

Category	Domain	Description	Example
Propensity	Driving	Respondents indicate on a (ordinal) scale to what extent they are likely to take risks while driving.	<i>For each of the following statements, please indicate the likelihood that you would engage in the described activity or behaviour if you were to find yourself in that situation: Not wearing a seat belt when being a passenger in the front seat. Very unlikely/Unlikely/Not sure/Likely/Very likely</i>
Propensity	Ethical	Respondents indicate on a (ordinal) scale to what extent they are likely to break rules/laws or cause harm to others or the extent to which they identify/ perceive themselves as being someone who breaks rules/laws or causes harm to others.	<i>For each of the following statements, please indicate the likelihood that you would engage in the described activity or behaviour if you were to find yourself in that situation: Taking some questionable deductions on your income tax return. Very unlikely/Unlikely/Not sure/Likely/Very likely</i>
Propensity	Gambling	Respondents indicate on a (ordinal) scale to what extent they are likely to take risks with gambling-related activities.	<i>What is the probability that you would do one of the following activities? Please rate on a scale from 0 to 10. Wagering a daily earnings on a bet. 0) very unlikely...10) very likely</i>
Propensity	Health	Respondents indicate on a (ordinal) scale to what extent they are likely to take risks with regards to their health or take part in activities or make decisions that can have detrimental consequences on their health.	<i>Please indicate the likelihood that you would engage in the described activity or behaviour if you were to find yourself in that situation: Drinking heavily at a social function. Extremely Unlikely (1) - Extremely Likely (7)</i>
Propensity	General	Respondents indicate on a (ordinal) scale to what extent they generally identify as someone who likes to take risks or is willing to take risks.	<i>Are you generally a person who is willing to take risks or do you try to avoid taking risks? Please answer on a scale from 0 to 10, where 0 means "not at all willing to take risks" and 10 means "very willing to take risks".</i>
Propensity	Investment	Respondents indicate on a (ordinal) scale to what extent they are likely to take risks with investments.	<i>Which of the statements comes closest to the amount of financial risk that you are willing to take when you save or make investments? Take substantial financial risks expecting to earn substantial returns / Take above average financial risks expecting to earn above average returns / Take average financial risks expecting to earn average returns / Not willing to take any financial risks</i>
Propensity	Occupation	Respondents indicate on a (ordinal) scale to what extent they are likely to take risks with regards to their job.	<i>Rate using a scale from 0 to 10. I don't mind taking risks in ... my professional career</i>
Propensity	Recreation	Respondents indicate on a (ordinal) scale to what extent they are likely to take risks with regards recreational actives or their likelihood of engaging in activities that involve height and/or speed and high risk of serious injury or death.	<i>For each of the following statements, please indicate your likelihood of engaging in each activity or behaviour: Going down a ski run that is beyond your ability or closed. Very unlikely/Unlikely/Not sure/Likely/Very likely</i>
Propensity	Social	Respondents indicate on a (ordinal) scale to what extent they are likely to take risks in social situations, or when trusting strangers.	<i>For each of the following statements, please indicate your likelihood of engaging in each activity or behaviour: Admitting that your tastes are different from those of your friends. Very unlikely/Unlikely/Not sure/Likely/Very likely</i>
Frequency	Alcohol	Respondents quantify the extent to which they consumed alcohol or experienced the consequences of alcohol consumption within a certain time frame.	<i>How many times in the last four weeks have you had an alcoholic drink? Most days / Once or twice a week / 2 or 3 times / Once only / Never</i>

Measure	Domain	Description	Example
Frequency	Driving	Respondents quantify the extent to which they have not been prudent whilst driving a vehicle within a certain time frame.	<i>During the past 30 days, how often did you drive a car or other vehicle when you had been drinking alcohol?</i>
Frequency	Drug	Respondents quantify the extent to which they consumed drugs or experienced the consequences of drug consumption within a certain time frame.	<i>During the last 30 days, how often, if ever, did you use these other drugs? Heroin, steroids, or MDMA (Ecstasy). 0) Never, 1) Less than once a week, 2) 1 or 2 days per week, 3) 3 or 4 days per week, 4) 5 or 6 days per week, 5) Every day</i>
Frequency	Ethical	Respondents quantify the extent to which they broke rules/laws or had issues with the law or cause harm to others within a certain time period.	<i>This is about fare dodging. How often did you do that in the last 12 months? Indicate number of times.</i>
Frequency	Gambling	Respondents quantify the extent to which they partook in gambling-related activities within a certain time frame.	<i>Pathological gambling (Brodtbeck et al., 2009)</i>
Frequency	Occupation	Respondents indicate the extent to which they have been reckless at their job/school or behaved in a way that could lead to them losing their job/get in trouble at school.	<i>This is about skipping school. How often did you do that in the last 12 months?</i>
Frequency	Sexual Intercourse	Respondents indicate the number of sexual partners or how often they had sexual intercourse without using a form of contraception within a certain time frame.	<i>With how many persons are you currently having a romantic or sexual relationship?</i>
Frequency	Smoking	Respondents quantify the extent to which they smoke cigarettes or other tobacco products within a certain time frame.	<i>About how many cigarettes or packs do you usually smoke in a day now?</i>
Frequency	Behaviour	These tasks mention a gambling-related activity/scenario or a form of game. Respondents are asked to decide between two or more options that offer different potential monetary gains and/or losses with varying probability. Also includes Willingness to Pay and Willingness to Accept tasks. Depending on the respondents' responses in such tasks, composite measures can be derived which summarise their tolerance towards risk (e.g., proportion of safe choices, risk aversion category). Such tasks can involve decision from experience or description, with hypothetical or incentivised choices.	<i>Now, imagine you have a choice between the following two options: Option A: A lottery with a 50% chance of winning 80\$ and a 50% chance of losing 50\$/ Option B: Zero dollars. Which option would you choose?</i>
Behaviour	Insurance	Tasks require respondents make choices about insurances, with hypothetical or incentivised choices. Also includes Willingness to Pay and Willingness to Accept tasks. Depending on the respondents' responses in such tasks, composite measures can be derived which summarise their tolerance towards risk (e.g., proportion of safe choices, risk aversion category).	<i>Assume that you know there is a 50% chance of losing \$1000 on a given day. You can take out insurance to cover this amount in case of loss. If an insurance policy is sold as listed below, would you purchase it?</i>

Table S5 cont.

Measure	Domain	Description	Example
Behaviour	Investment	These tasks explicitly mention an investment-related activity/scenario. Respondents can be asked how much of an endowment they wish to allocate to different options. These tasks can be hypothetical or incentivised.	<i>Imagine that you had won 100,000 euros in the lottery. Immediately after receiving your winnings you receive the following offer: You have the chance to double your money. But it is equally possible that you will lose half of the amount invested. You can participate by staking all or part of your 100,000 euros on the lottery, or choose not to participate at all. What portion of your lottery winnings would you be prepared to stake on this financially risky yet potentially lucrative lottery investment?</i>
Behaviour	Occupational	Tasks require respondents to make choices about jobs offering different salaries with different probabilities. Depending on the respondents' responses in such tasks, composite measures can be derived which summarise their tolerance towards risk (e.g., proportion of safe choices, risk aversion category).	<i>Which ONE do you prefer? Option A: A 50% chance of the salary increasing by 30%, but also a 50% chance of the salary increasing by 11%; Option B: Guaranteed salary increase of 20%.</i>

End of Table