

## University of Groningen

### Test-Retest Reliability of the STRAQ-1

Dujols, Olivier; Klein, Richard A.; Lindenberg, Siegwart; Van Lissa, Caspar J.; Ijzerman, Hans

*Published in:*  
Collabra: Psychology

*DOI:*  
[10.1525/collabra.122155](https://doi.org/10.1525/collabra.122155)

**IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.**

*Document Version*  
Publisher's PDF, also known as Version of record

*Publication date:*  
2024

[Link to publication in University of Groningen/UMCG research database](#)

*Citation for published version (APA):*  
Dujols, O., Klein, R. A., Lindenberg, S., Van Lissa, C. J., & Ijzerman, H. (2024). Test-Retest Reliability of the STRAQ-1: A Registered Report. *Collabra: Psychology*, 10(1), Article 122155.  
<https://doi.org/10.1525/collabra.122155>

#### **Copyright**

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.


#### **Take-down policy**

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

*Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.*

## Social Psychology

# Test-Retest Reliability of the STRAQ-1: A Registered Report

Olivier Dujols<sup>1</sup><sup>a</sup>, Richard A. Klein<sup>2</sup>, Siegwart Lindenberg<sup>2,3</sup>, Caspar J. Van Lissa<sup>2</sup>, Hans IJzerman<sup>1,4</sup>

<sup>1</sup> Univ. Grenoble Alpes, Univ. Savoie Mont Blanc, LIP/PC2S, F-38000, Grenoble, France, <sup>2</sup> Tilburg University, The Netherlands, <sup>3</sup> Rijksuniversiteit Groningen, The Netherlands, <sup>4</sup> Institut Universitaire de France (IUF), France

Keywords: Test-Retest, Longitudinal Measurement Invariance, Attachment Theory, Social Thermoregulation, Registered Report

<https://doi.org/10.1525/collabra.122155>

---

## Collabra: Psychology

Vol. 10, Issue 1, 2024

---

This Registered Report provides the first test of measurement invariance across time points and estimates of test-retest reliability for the Social Thermoregulation, Risk Avoidance Questionnaire (STRAQ-1, Vergara et al., 2019). The scale was developed and validated to understand the physiological drives underlying interpersonal bonding, measured by four constructs: the desire to socially regulate one's temperature, the desire to solitary regulate one's temperature, the sensitivity to higher temperatures, and the desire to avoid risk. Previous studies with large samples across 12 countries showed that the STRAQ-1 has a stable factorial structure, satisfying internal consistencies for the temperature subscales, and expected correlations in its nomological network. However, to date, this instrument has no estimates of test-retest reliability. Throughout four academic years (from 2018 to 2022),  $N = 183$  French student participants took the STRAQ-1 at least two times. Out of the four STRAQ-1 subscales, only two were longitudinally invariant across two-time points. Thus, half of the STRAQ-1 constructs and latent scores were dissimilar and incomparable across time. We then conducted test-retest reliability using Intra Class Correlation coefficient (ICC) for the *Social Thermoregulation*, *Solitary Thermoregulation*, *High-Temperature Sensitivity*, and *Risk Avoidance* subscales. ICCs estimates were respectively for agreement and consistency: .70, .70 overall moderate to good, .62, .62 overall moderate, .67, .67 overall moderate, and .53, .53 overall poor to moderate, respectively. Our study suggests that test-retest reliability was insufficient for psychological diagnosis, and that future studies should address the problem of low generalizability of the constructs.

In the psychological literature, how people engage in interpersonal relationships is often understood through the prism of attachment theory, which proposes that individuals seek relational closeness to feel secure (Bowlby, 1969–1982). But, while the importance of the physical safety of human infants is recognized in infant care in hospitals (e.g., temperature regulation), much less attention has been devoted to its self-report measurement in adults. Indeed, adult attachment measures focus primarily on the self-reported feelings of emotional safety and leave aside the issue of physical safety (e.g., Brennan et al., 1998; Fraley et al., 2000).

A notable exception to this is the Social Thermoregulation and Risk Avoidance Questionnaire (STRAQ-1) developed and validated by Vergara et al. (2019). The STRAQ-1 measures physical safety and the physiological drives underlying interpersonal bonding through four constructs: the desire to socially regulate one's temperature, the desire to solitary regulate one's temperature, the sensitivity to higher temperatures, and the desire to avoid risk. Previous

studies in a large sample across 12 countries showed that the STRAQ-1 has a stable factorial structure, acceptable internal consistencies for the temperature subscales, and expected correlations in its nomological network. However, to date, no assessment of the test-retest reliability, crucial for the scale psychometrics and future use (e.g., evaluation of the impact of an intervention), has been conducted. In this article, we first assess longitudinal measurement invariance of the STRAQ-1 across time points, followed by an analysis of the test-retest reliability.

### Attachment and its Measurement

Bowlby (1969–1982) proposed that social relationships are essential and adaptive to a child's survival since they are not able to survive by themselves. He postulated that a motivational system - the behavioral attachment system - drives the child to seek protection and support from the adult through crying and clinging behaviors. This behavioral system binds the child to the caregiver(s) so that they

---

a Correspondence: [dujols.ol@gmail.com](mailto:dujols.ol@gmail.com)

become attachment figure(s). Based on the availability and reliability of the care, the child will construct a mental representation (a working model) of the ability of their attachment figure to provide security, that in return, will impact their behaviors and feelings of security (Bretherton & Munholland, 2008; De Wolff & van IJzendoorn, 1997). These attachment patterns - the extent to which the child is secure or insecure in its relationship with the attachment figure - have been found to vary between individuals, and to be relatively stable from infancy to adulthood - when the main attachment figure becomes the romantic partner (Dugan & Fraley, 2022; Fraley, 2019; Fraley et al., 2021).

To measure attachment and identify how children differ based on it, an initial three-style classification was derived from observations of children: avoidant attachment, secure attachment, and anxious attachment (Ainsworth, 1979). This classification was later expanded to include disorganized attachment (Main & Solomon, 1986, 1990). In adults, the most widely used and currently psychometrically most sound instrument for measuring adult attachment styles is the Experiences in Close Relationships Inventory (ECRI; Brennan et al., 1998), which has since been revised (Experiences in Close Relationships Revised, ECR-R; Fraley et al., 2000, measured, for instance through the level of agreement with the statement “*I am very comfortable being close to romantic partners.*”). However, neither the ECRI, the ECR-R, nor the measures that preceded these adult attachment scales considered physical safety, such as protection against the cold, which is one essential aspect of survival proposed by Bowlby (1969–1982)<sup>1</sup>.

### Social-Thermoregulation-Based Attachment

The theory of social-thermoregulation-based attachment was based on observations of non-human animals that found that when the temperature decreases, both infants and adults tend to move closer to their conspecifics to save energy and increase survival fitness (for example, through huddling, see Gilbert et al., 2010). The importance of physical proximity has also been studied in humans, demonstrating a determining role of thermoregulation in newborns. For instance, Bystrova et al. (2007) found that the mother’s temperature was related to that of their infant and increased after the birth of the infant (and even more so with skin-to-skin contact and early breastfeeding). In adulthood, attachment moderate people’s responses to temperature: securely attached people think of their loved ones when they are cold (versus warm), whereas this effect flips for those who are insecurely attached (IJzerman et al., 2018; see also Rocha IJzerman, 2021).

But existing attachment measures often do not map onto the concept of social-thermoregulation-based attachment. To better measure inter-individual differences in the regulation of temperature and risk through social relationships, Vergara et al. (2019) developed the Social Thermoregulation and Risk Avoidance Questionnaire (STRAQ-1). Across 12 countries and 1,510 participants, they found that the STRAQ-1 had a four-factor structure: (1) *Social Thermoregulation* (5 items;  $\omega_t = .83$ ; reflecting the desire to warm up physically with close others), (2) *Solitary Thermoregulation* (8 items;  $\omega_t = .77$ ; reflecting a desire to regulate temperature alone), and (3) *High-Temperature Sensitivity* (7 items;  $\omega_t = .83$ ; reflecting a preference for colder temperatures and a distaste for hotter temperatures, and (4) *Risk Avoidance* (3 items;  $\omega_t = .57$ ; reflecting the tendency to avoid - social - exploration). The items of each subscale are presented in the Supplementary Materials (<https://osf.io/6px2u>).

In several French samples, the internal consistencies of the subscales were similar to those of the original validation study (Sarda et al., 2021; Vidal et al., 2022; Wittman et al., 2022<sup>2</sup>). Vergara et al. (2019) also investigated the nomological network of the STRAQ-1. We provide the most relevant correlations in [Table 1](#). Again, the correlations with attachment have been replicated (excluding the *Risk Avoidance* subscale) in a French sample, and showed a similar pattern, with the addition of a relationship to loneliness (Wittmann et al., 2022). However, despite evidence of the STRAQ-1 factorial structure, sufficiently high internal consistencies (except for the risk subscale), and validity through the nomological network, to date, no test-retest reliability has been conducted.

Test-retest reliability is crucial for scale psychometrics and for its future use. Without acceptable test-retest reliability, it is possible to confound artifacts of the measurement with true pre- and post-intervention differences in the rating of the scale or miss the true effects of an intervention. Thus, test-retest reliability is necessary for theory development and to use the scale for interventions (cf., IJzerman et al., 2017). Therefore, the main purpose of this article is to examine the test-retest reliability of the STRAQ-1. Before doing so, we provide an assessment of longitudinal measurement invariance of the STRAQ-1 across time points as it is a prerequisite for test-retest analysis. In the Supplementary Materials (<https://osf.io/mr8n3/>) we also provide internal consistency (Alpha and Omega), per time point. We expected psychometrics in our current French samples similar to the original finding of Vergara et al. (2019). This research was conducted in line with the CO-RE Lab Lab Philosophy v6 (Goncharova et al., 2019).

1 See Vergara et al. (2019) for a more in-depth review of existing adult attachment measures and the motivation behind the development of the STRAQ-1.

2 The data in the Wittman et al. (2022) study is (partly) the same data that we are using in this project. The data have thus been previously observed by the main author Adrien Wittman and the person that did the code review (Mae Braud), but not in a way that would be related to the proposed analyses of this project. None of the analyses that we intend to conduct in the current have been run on the data.

**Table 1. Correlation in the nomological network of STRAQ-1.**

	Attachment Anxiety	Attachment Avoidance	Health	Stress	Self Control	Network Size
Social Thermoregulation	<i>n.s.</i>	-.31	.10	<i>n.s.</i>	<i>n.s.</i>	.10
Solitary Thermoregulation	.08	<i>n.s.</i>	<i>n.s.</i>	.11	<i>n.s.</i>	<i>n.s.</i>
High-Temperature Sensitivity	.10	<i>n.s.</i>	-.11	.15	-.17	-.10
Risk Avoidance	.17	<i>n.s.</i>	-.11	.24	<i>n.s.</i>	-.12

Note. In the table, the reported correlations are all significant, the interested reader can refer to the Supplementary Materials via our OSF page (<https://osf.io/86qdx>) for the complete nomological network of the scale investigated in the original development paper (Vergara et al., 2019).

## Method

The Stage 1 version of the manuscript associated with this Registered Report was granted in-principle acceptance on July 2023 the 18th. The original accepted Stage 1 manuscript, unchanged after the in-principle acceptance, and the associated open review process may be viewed at this link: <https://rr.peercommunityin.org/articles/rec?id=419>. Following in-principle acceptance after the stage 1 review, we conducted the planned analyses.

## Participants

A pool of psychology students replied to the STRAQ-1 for four academic years, in 2018-2019 (from October 15<sup>th</sup> to 27<sup>th</sup>,  $N_{2018} = 505$ ), 2019-2020 (from September 16<sup>th</sup> to October 01<sup>th</sup>,  $N_{2019} = 298$ ), 2020-2021 (from March 2<sup>th</sup> to 28<sup>th</sup>,  $N_{2020} = 236$ ), and 2021-2022 (from January 12<sup>th</sup> to February 11<sup>th</sup>,  $N_{2021} = 400$ ), as part of a larger “test week”. We merged the participant responses across the four academic years based on their pseudo-anonymized code. In total,  $N = 183$  French students took the STRAQ-1 at least two times (161 females, 19 males, 3 others,  $M_{age\_T1} = 19.70$ ,  $SD_{age\_T1} = 2.74$ ;  $M_{height\_T1} = 166.00$ ,  $SD_{height\_T1} = 7.42$ ,  $M_{weight\_T1} = 59.60$ ,  $SD_{weight\_T1} = 11.30$ ), from which  $N = 25$  participants took the STRAQ-1 at least three times (24 females, 1 males,  $M_{age\_T1} = 19.00$ ,  $SD_{age\_T1} = 1.54$ ;  $M_{height\_T1} = 166.00$ ,  $SD_{height\_T1} = 7.43$ ,  $M_{weight\_T1} = 57.70$ ,  $SD_{weight\_T1} = 8.65$ ), and  $N = 4$  participants took the STRAQ-1 four times (4 females,  $M_{age\_T1} = 19.50$ ,  $SD_{age\_T1} = 2.08$ ;  $M_{height\_T1} = 165.00$ ,  $SD_{height\_T1} = 12.00$ ,  $M_{weight\_T1} = 61.00$ ,  $SD_{weight\_T1} = 12.30$ ).

Because we did not have specific hypotheses about the impact of specific academic years, we decided to label the first STRAQ-1 score that we had for a participant T1, with the second STRAQ-1 score T2 (and so forth). T1 and T2 thus do not reflect a specific academic year. The gap between T1 and T2 could vary between one to three years. For example, if a participant took the STRAQ-1 in 2019, 2020, and 2021, then this participant has three-time points: T1 then corresponds to the score of 2019, T2 to 2020, and T3 would be dropped from the pre-registered analysis (but would be included in the exploratory analysis).

## R Packages

We used the following R packages to conduct the analysis: rio (Chan et al., 2021), janitor (Firke, 2021), tidyverse (Wickham et al., 2019), psych (Revelle, 2022), GPARotation (Coen & Robert, 2005), EFA.dimensions (O’Connor, 2022), lavaan (Rosseel, 2012), semPlot (Epskamp, 2022), semTools (Jorgensen, 2021), energy (Rizzo & Szekely, 2022), semPower (Moshagen & Erdfelder, 2016), ICC.Sample.size (Zou, 2012).

## Power Analysis

As we relied on secondary data, we did not conduct an a priori power analysis, but instead we conducted a sensitivity power analysis. Based on the number of participants ( $N = 183$ ) that answered the STRAQ-1 at least twice, we calculated projected power to detect desired effect size. There are two recommendations for sample size for longitudinal measurement invariance analyses: five (Dimitrov, 2014) versus ten (Kline, 2016). In the former, we would need 5 (participants) \* 8 (items) \* 2 (time points) = 80 participants. In the latter, we would need 10 (participants) \* 8 (items) \* 2 (time points) = 160 participants. We also computed power for a general configural longitudinal measurement invariance models (CFA) models. We set power to 80%, alpha to .05, the amount of misfit to correspond to an RMSEA of at least .05, and the degrees of freedom to 100. The result of this analysis was that 164 participants would be required. In either case, our sample size was slightly above the required sample for detecting longitudinal measurement invariance over two time points.

We relied on Intra-Class Correlation (ICC) estimates for the test-retest reliability. Given that we had 183 participants with at least two time points, we had 80% power to detect an ICC of 0.2 with a pre-specified value of alpha of 0.05 (Bujang & Baharum, 2017). This means that if a small test-retest reliability exists (ICC = 0.2) we would have an 80% chance to detect it. However, we expected our subscales to present test-retest reliability between moderate (ICC between 0.5 and 0.75) and good (ICC between 0.75 and 0.9). Because researchers have argued that detection of non-zero ICC scores may not be sufficient and meaningful (see for example, Parsons et al., 2019), we also conducted a power analysis to estimate the 95% CI width that our sam-



ple will provide as a function of different ICC values. This power analysis suggested that based on our sample size  $N = 183$ , we could estimate any ICC above .30 with a 0.2 width of the 95% CI, and any ICC above .80 with a 0.1 width of the 95% CI. Hence, we had sufficient power to detect our expected ICC. The R code associated with our power analysis is available in the Supplementary Materials via our OSF page: <https://osf.io/mr8n3/>.

## Measure

Participants rated the four subscales of the questionnaire STRAQ-1 on a Likert type scale ranging from 1 = Strongly Disagree to 5 = Strongly Agree. Our cut-off for the selection of the labels for internal consistency was: above or equal to .70 for acceptable, and under .70 for poor. This cut-off is often used in the literature and is based on Nunnally & Bernstein (1994), even if it was not intended as a gold standard for acceptable internal consistency. The *Social Thermoregulation subscale* presented acceptable internal consistency: McDonald's  $\omega_{t1} = .85$ ,  $\omega_{t2} = .85$  (e.g., "I prefer to warm up with someone rather than with something"). The *Solitary Thermoregulation subscale* presented acceptable internal consistency: McDonald's  $\omega_{t1} = .76$ ,  $\omega_{t2} = .77$  (e.g., "When it is cold, I more quickly turn up the heater than others"). The *High-Temperature Sensitivity subscale* presented acceptable internal consistency: McDonald's  $\omega_{t1} = .76$ ,  $\omega_{t2} = .73$  (e.g., "I am sensitive to heat"). The *Risk Avoidance subscale* presented poor internal consistency: McDonald's  $\omega_{t1} = .49$ ,  $\omega_{t2} = .58$  (e.g., "I try to maintain myself in familiar places"). For each subscale, we averaged their items into a mean score.

## Results

### Confirmatory Analyses

The main goal of the analysis was to examine the test-retest reliability of the STRAQ-1, but longitudinal measurement invariance across time points must be established before conducting test-retest reliability (Chen, 2008). We first assessed the longitudinal measurement invariance of each of the four STRAQ-1 subscales across two-time points. We then ran the test-retest analysis. The R scripts of the analysis are available on the project's OSF page: <https://osf.io/mr8n3/>. All reported analyses were preregistered unless specified otherwise.

### Longitudinal Measurement Invariance

The main goal of the analysis was to ensure that the nature of the construct had not changed substantially over time. In longitudinal studies, the nature or meaning of a construct may change over time, resulting in longitudinal measurement non-invariance (Chen, 2008). Confirmatory Factor Analysis (CFA) is a common method for evaluating the level of invariance across time points (Drasgow & Kanfer, 1985; Widaman et al., 2010). Our procedure to test for longitudinal measurement invariance was to compare progressively more constrained CFA models. These models test

incremental levels of measurement invariance across our two-time points (T1-T2). The levels of longitudinal measurement invariance have different implications for the construct: (a) if the configural level holds, then the structure of the measure is similar between T1 and T2; (b) if the metric level hold, then the structure of the measure and the constructs are similar between T1 and T2; (c) if the scalar level hold then the structure of the measure and the constructs are similar and the mean differences between T1 and T2 can be compared. Longitudinal scalar invariance is thus the minimal level required for our planned ICC analysis that uses the means scores of T1 and T2 (Kline, 2016; Mackinnon et al., 2022).

To investigate whether the variables in our dataset followed a multivariate normal distribution, we used the function ``mvnorm.etest`` from the Energy package. The analysis showed that our data does not follow a multivariate distribution ( $E = 3.76$ ,  $p < .001$ ). A priori, we had already decided to use the WLSMV estimator instead of ML or MLR as arguments in the `cfa` function in lavaan to compute our CFA model, irrespective of the outcome of the test for multivariate normality. Our measure is a 5-point Likert type scale, the label are (1) "Strongly disagree", (2) "Disagree" (3) "Neutral", (4) "Agree", (5) "Strongly agree". But the numbers do not necessarily represent equal intervals or differences in magnitude between the ordered labels. Consequently, data obtained from a Likert scale are generally considered as ordinal, rather than continuous (where the intervals are equal between values). The WLSMV is the preferred solution when (a) the data is ordinal, and (b) if data is potentially not normally distributed, as it makes no distribution assumptions (see Flora & Curran, 2004; Kline, 2016; Li, 2016). Then, we reported the robust weighted least squares fit for each model. We also verified the absence of Heywood cases (factor loading  $> 1$  or negative variances). For the registered analysis, when we detected residual correlations above  $r = .10$  in a model, we choose not to apply modification indices. But we did apply them in the exploratory part of the analysis. We then tested configural invariance, freely estimating the parameters and thresholds for T1 and T2, to verify whether the same latent factor structures held across time points. Our criteria for configural invariance were comparative fit index  $< .95$ , root mean square error of approximation  $< .06$  (CI 90% upper bond  $< .10$ , and non-significant  $p$ -value), and standardized root mean square residual  $< .05$  (Kline, 2016).

Following our configural invariance test, we tested metric invariance, constraining the factor loadings and thresholds to be the same between T1 and T2, to verify whether the latent constructs were similar across time points. Then we tested scalar invariance, constraining the items' intercepts and thresholds to be the same between T1 and T2, to ensure that the latent score at T1 and T2 were comparable. Finally, we tested residual invariance, further constraining the residual variances to be the same between T1 and T2, to ensure strict invariance of the latent score between T1 and T2. Residual invariance has been described to be hard to reach for most psychological measurement instruments (Kline, 2016; Van De Schoot et al., 2015). We thus consid-

ered the subscales that reached scalar invariance as longitudinally invariant.

To identify which level of longitudinal measurement invariance holds for each model, we followed the recommendation of Mackinnon et al. (2022). Mackinnon et al. (2022) provided several criteria to assess model fit for measurement invariance, one of these is the delta CFI (of .01) which is also recommended by a simulation study (Cheung & Rensvold, 2002). We decided to rely only on a  $\Delta$ CFI of  $-.01$  or more to conclude that the model with the largest CFI should be chosen. This means that if the  $\Delta$ CFI is inferior or equal to  $-.01$  we will choose the more parsimonious model and conclude for the longitudinal invariance of the specific level (metric, or scalar, or residual). Before pre-registration, we made choices about which metrics and cut-offs we would base our conclusion and interpretation of the subscale's performance. But we acknowledge a lack of clear norms in the field about which metric to choose for our planned analyses. So, in addition to our pre-registered metric and cut-offs, we report in Table 2 the results of other fit metrics even though we did not plan to use them for inferences and did not preregister any cut-of-value for them. This process will allow other researchers, who would prefer other indicators or cut-offs than ours, to be able to evaluate our models according to their criteria.

Out of the four STRAQ-1 subscales, two reached longitudinal scalar invariance across two-time points. Table 2 provides a complete description of the fits of all the models. Based on the results of the longitudinal CFA models, we considered longitudinally invariant the subscales that reached scalar invariance. The *Social Thermoregulation* (Configural-Metric  $\Delta$ CFI:  $+.014$ ; Metric-Scalar  $\Delta$ CFI  $< .001$ ) and *High-Temperature Sensitivity* (Configural-Metric  $\Delta$ CFI:  $+.012$ ; Metric-Scalar  $\Delta$ CFI  $< .001$ ) subscales met our criteria to reach scalar invariance, and thus are considered longitudinally invariant across two-time points. By contrast, the *Risk Avoidance* and *Solitary Thermoregulation* subscales were considered longitudinally non-invariant across two time points. The *Risk Avoidance* subscale failed to reach metric invariance (Configural-Metric  $\Delta$ CFI =  $-.027$ ). The configural model of the *Solitary Thermoregulation* subscale had insufficient fit to the data ( $\chi^2 = 158.05$ , CFI =  $.899$ , RMSEA =  $.061$ , 90% CI RMSEA =  $[.043, .077]$ , SRMR =  $.072$ ). Based on these analyses, the *Social Thermoregulation* and *High-Temperature Sensitivity* constructs are thus respectively similar across two-time points and their latent scores can be meaningfully compared in our dataset. The *Risk Avoidance* and *Solitary Thermoregulation* constructs are thus respectively dissimilar across two-time points and their latent scores comparison may not be meaningful in our dataset.

### Test-Retest Reliability

The main goal of the analysis was to investigate the test-retest reliability of the four STRAQ-1 subscales (*Social Thermoregulation*, *Solitary Thermoregulation*, *High-Temperature Sensitivity*, and *Risk Avoidance*), using Intraclass Correlation Coefficient (ICC) analysis. The ICC analysis compares the variation across different ratings of the same

individuals to the variation across all ratings and all individuals. An ICC close to 1 indicates that the scores from the same individual are highly similar. An ICC close to zero shows that the scores from the same individual are not similar. Koo & Li (2016) defined standards for the ICC with reliability being poor at  $ICC < 0.5$ ; moderate at  $0.5 < ICC < 0.75$ ; good at  $0.75 < ICC < 0.9$ ; and excellent at  $ICC > 0.9$ . These are the cut-off values that we used for labeling our results. If the 95% confidence interval of an ICC estimate was in between two labels, we used both (for example, if the 95% CI interval would have been  $[.83, .94]$ , the level of reliability would have been regarded as “good” to “excellent”; see Koo & Li, 2016). We recognize that the discussion around cut-offs is contentious and that cut-offs are often arbitrarily chosen, which may make our values equally arbitrary (see e.g., Watson, 2004). The resulting labels (e.g., “good”) are considered as one of many means to assess the validity of a measure (Rodebaugh et al., 2016) and a first step towards defining a normative range of reliability estimates for a scale that will be applied across samples or contexts.

We computed and report ICC(2,1), to evaluate absolute agreement between participants at two time points, and ICC(3,1), to evaluate consistency. Both of these ICCs are calculated through two-way mixed-effect models. ICC(2,1) accounts for systematic and random error by specifying the time of measurement as a random effect in the model. ICC(3,1) only accounts for random error because the time of measurement is not specified as a random effect in the model (Koo & Li, 2016). The STRAQ-1 subscales' test-retest reliability between the two time points was estimated with intraclass correlation coefficients (ICCs) using the psych package in R (Revelle, 2022). The analysis code is available on the OSF: <https://osf.io/mr8n3/>.

For the High-Temperature Sensitivity subscale, the estimated agreement was  $.70$ , 95% CI =  $[.62, .77]$ , and the estimated consistency was  $.70$ , 95% CI =  $[.62, .77]$ . For the Social Thermoregulation subscale, the estimated agreement was  $.62$ , 95% CI =  $[.52, .70]$ , and the estimated consistency was  $.62$ , 95% CI =  $[.52, .70]$ . For the Solitary Thermoregulation subscale, the estimated agreement was  $.67$ , 95% CI =  $[.60, .74]$ , and the estimated consistency was  $.67$ , 95% CI =  $[.60, .74]$ . Finally, for the Risk Avoidance subscale, the estimated agreement was  $.48$ , 95% CI =  $[.36, .59]$ , and the estimated consistency was  $.49$ , 95% CI =  $[.37, .59]$ . We found the overall test-retest reliability over two-time points of the High-Temperature Sensitivity subscale to be “moderate” to “good”, of the Social and Solitary Thermoregulation subscales to be “moderate”, and of the Risk Avoidance subscale to be “poor” to “moderate”.

### Exploratory Analysis (not pre-registered)

We computed non pre-registered extra analyses that are labeled as exploratory either because of the relative degree of flexibility they introduce in the analysis (partial invariance), or because we did not have enough power to be sure of the effects (test-retest on more than two-time points), or because we did not have a priori hypotheses but wanted to

**Table 2. CFA fits of the longitudinal invariance models.**

Model name	Configural model	Metric model	Scalar model	Measurement invariance
Social Thermoregulation	$\chi^2 = 44.20$ CFI = .972 RMSEA = .054 90% CI RMSEA = [.015, .084] SRMR = .046	$\chi^2 = 41.78$ CFI = .986 RMSEA = .036 90% CI RMSEA = [<.001, .068] SRMR = .054	$\chi^2 = 40.66$ CFI = .986 RMSEA = .036 90% CI RMSEA = [.000, .068] SRMR = .050	Configural-Metric $\Delta CFI = +.014$ Metric-Scalar $\Delta CFI < .001$ Final decision: Scalar invariance
Solitary Thermoregulation	$\chi^2 = 158.05$ CFI = .899 RMSEA = .061 90% CI RMSEA = [.043, .077] SRMR = .072	$\chi^2 = 155.06$ CFI = .916 RMSEA = .053 90% CI RMSEA = [.035, .070] SRMR = .077	$\chi^2 = 153.80$ CFI = .917 RMSEA = .053 90% CI RMSEA = [.035, .070] SRMR = .073	Configural-Metric $\Delta CFI = +.017$ Metric-Scalar $\Delta CFI = +.001$ Final Decision: Non invariance
High Temperature Sensitivity	$\chi^2 = 103.85$ CFI = .957 RMSEA = .052 90% CI RMSEA = [.029, .073] SRMR = .050	$\chi^2 = 101.09$ CFI = .969 RMSEA = .042 90% CI RMSEA = [.012, .063] SRMR = .061	$\chi^2 = 99.94$ CFI = .969 RMSEA = .042 90% CI RMSEA = [.012, .064] SRMR = .058	Configural-Metric $\Delta CFI = +.012$ Metric-Scalar $\Delta CFI < .001$ Final decision: Scalar invariance
Risk Avoidance	$\chi^2 = 4.72$ CFI = 1.000 RMSEA < .001 90% CI RMSEA = [<.001, .100] SRMR = 0.26	$\chi^2 = 11.82$ CFI = .973 RMSEA = .051 90% CI RMSEA = [<.001, .109] SRMR = .047	$\chi^2 = 10.57$ CFI = .975 RMSEA = .053 90% CI RMSEA = [<.001, .114] SRMR = .042	Configural-Metric $\Delta CFI = - .027$ Metric-Scalar $\Delta CFI = +.002$ Final decision: Configural invariance

check the robustness of our confirmatory analyses (effect of the “academic year”).

**Exploratory Partial Longitudinal Measurement Invariance**

We explored the partial longitudinal invariance and modification indices of the scales that did not reach at least scalar longitudinal invariance in the confirmatory analysis. We found the Risk avoidance subscale to reach longitudinal scalar invariance across two-time points after freeing the loadings of the items “*I don’t trust people I have not met before*” across the two time points. Additionally, we found the Solitary Thermoregulation subscale to reach longitudinal scalar invariance across two-time points after applying modification indices (correlated residuals of some items) on the configural model. We do not consider the scale to provide longitudinal invariance because of (i) the small number of items in the Risk avoidance – three – included in the Risk Avoidance subscale, the partial invariance corresponds to 1/3 of the items (ii) these modifications of the models are post-hoc and were not pre-registered.

**Exploratory Intra Class Correlation**

The next ICC analyses were exploratory because we did not have the power to test for longitudinal measurement invariance for three and four-time points. To select the label for overall excellent/good/moderate/poor we took the worst ICC between ICC(2,1) and ICC(3,1). Based on the cut off values defined by Koo & Li (2016) the labels were: poor

at  $ICC < 0.5$ ; moderate at  $0.5 < ICC < 0.75$ ; good at  $0.75 < ICC < 0.9$ ; and excellent at  $ICC > 0.9$ . We computed ICCs estimates including only the 25 participants with three time points. We computed the ICCs for *Social Thermoregulation*, *Solitary Thermoregulation*, *High-Temperature Sensitivity*, and *Risk Avoidance*. The ICCs were respectively .65, .44, .60, .45 for agreement (ICC 2,1), and .65, .44, .61, .45 for consistency (ICC 3,1). For this test we had 90% power to detect an ICC of 0.4. Thus two subscales (*Social Thermoregulation* and *High-Temperature Sensitivity*) presented “moderate”, and two subscales (*Solitary Thermoregulation* and *Risk Avoidance*) presented “poor” test-retest reliability across at least three time points. We also computed ICCs in models including only the four participants that did the STRAQ-1 four-time. We computed the ICC for *Social Thermoregulation*, *Solitary Thermoregulation*, *High-Temperature Sensitivity*, and *Risk Avoidance*. the ICCs were respectively .59, .42, .65, .95 for agreement (ICC 2,1), and .59, .42, .65, .95 for consistency (ICC 3,1). For this test we had 80% power to detect an ICC of 0.8. These exploratory results indicated that one subscale presented “excellent”, and two subscales presented “moderate” and one “poor” test-retest reliability across more than two-time points. But all these ICCs were underpowered except for the *Risk Avoidance* subscale.

**Exploratory Effect of the Academic Year (Over Four-Time Points)**

As a robustness analysis, we investigated whether the “academic year” could determine differences in STRAQ-1

scores (e.g., because of the onset of the COVID-19 pandemic or temperature changes over the years<sup>3</sup>). These analyses were exploratory because we did not have a priori hypotheses about the effect of the academic year (2018, 2019, 2020, 2021). Also, in case of an effect of the academic year, we would not have been able to say anything about the cause of the effect, and we would only have been able to speculate about why this effect occurred.

We used a linear mixed model to compute both ICCs estimates (ICC 2,1 and ICC 3,1) from four linear mixed models in which the academic year was specified as a random effect of each of the STRAQ-1 scores. We computed the ICCs estimates for *Social Thermoregulation*, *Solitary Thermoregulation*, *High-Temperature Sensitivity*, and *Risk Avoidance* over four-time points. As in the confirmatory analysis section, we excluded participants and consider them outliers only if their Cook's D or Lever presents "gaps" (value at least three times the Cook D or Lever of the previous value for the highest value) or when the Studentized residual absolute value was above four. We did not find a large effect of the academic year on the STRAQ-1 response. The standard deviations of the random effects of the academic year were respectively  $>.001$ ,  $.064$ ,  $.135$ ,  $.019$  and the ICCs were respectively  $.62$ ,  $.67$ ,  $.70$ ,  $.49$  for agreement (ICC 2,1) and  $.62$ ,  $.66$ ,  $.71$ ,  $.49$  for consistency (ICC 3,1). For this test, we had 90% power to detect an ICC of 0.4. The results indicated that the minor random effects of the academic year were minor and minor changes in the ICCs induced by the inclusion of the four times points.

## Discussion

We provided the first test-retest reliability across time points of the STRAQ-1 subscales. The assessment of test-retest reliability was necessary for the psychometrics of the scale and its future use, but also for theory development (cf. IJzerman et al., 2017). In addition, we assessed the internal consistency of the scales in our sample and the longitudinal measurement invariance of the STRAQ-1 subscales across two time points. Overall, we found in our data that the STRAQ-1 subscales had relatively low test-retest reliability, acceptable and similar reliability compared to previous studies, and evidence of longitudinal invariance across two time points for only two out of four subscales.

In our sample, the Social Thermoregulation subscale, the Solitary Thermoregulation subscale, and the Risk Avoidance subscale had similar internal consistency to what was originally reported by Vergara et al. (2019). For each subscale respectively the internal consistencies that we found in our sample for T1 and T2 compared to the ones found by Vergara et al. (2019) were:  $\omega_t = .85-.85 / .83$ ;  $\omega_t =$

$.76-.77 / .77$ ;  $\omega_t = .49-.58 / .57$ . Interestingly, we found a small discrepancy between the internal consistency in our sample and that of Vergara and colleagues for High-Temperature Sensitivity  $\omega_t = .76-.73 / .83$ . The discrepancy may be explained by the fact that Vergara et al. (2019) relied on a much more geographically diverse sample.

We concluded that two of the STRAQ-1 (Social Thermoregulation and High-Temperature Sensitivity) out of the four subscales were longitudinally invariant across two-time points in our sample. The current data, suggest that test-retest reliability was insufficient for psychological diagnosis, and that future studies should address the problem of low measurement invariance (see COTAN standards, Evers et al., 2015). The development of new scales including more culturally suitable items may resolve the low generalizability pointed out by our analyses. The Social Thermoregulation, Risk Avoidance and Eating Questionnaire – 2 (STRAEQ-2, Dujols et al., 2024) includes new scales developed at 53 sites in 32 countries. The STRAEQ-2 is currently in validation and could potentially resolve the measurement invariance issues pointed out by the current study. Future studies should test for longitudinal measurement invariance of the STRAEQ-2.

In our sample, the STRAQ-1 subscales show relatively low stability across two time point separated by at least one year. According to Vergara et al. (2019), the STRAQ-1 was supposed to measure stable – trait – constructs that are unlikely to change rapidly in adulthood. A recent meta-analysis about personality trait development across the lifespan showed (similarly to previous meta-analysis, see Roberts & DelVecchio, 2000) that – after young adulthood – traits are indeed stable: they found the average rank-order stability to be  $r = .60$ , but with a large heterogeneity across studies (Bleidorn et al., 2022). Nevertheless, life events (for instance, attachment traumas) are known to introduce changes in personality traits and can be linked differently to different traits (Bleidorn et al., 2018, Bühler et al., 2023). But because no test-retest of a scale to assess this had been conducted when we conducted the study, we did not have any strong a priori hypothesis (i) about how life events could induce changes in participant responses to the STRAQ-1, and (ii) about the timeframe in which such change in the measured personality traits could occur.

Additionally, by using in combination ICC(2,1) and ICC(3,1), we did not find systematic error between our measurement time points - independently of the result longitudinal invariances in our sample. The values of our two ICCs showed near equality for all the STRAQ-1 subscales. This qualitative indicator shows an absence of – or at least a very low – systematic bias between our measurement

<sup>3</sup> Our mean and standard deviation of temperature in degrees Celsius (in degrees Fahrenheit in between parentheses) in Grenoble for the years included in the sample were for 2018  $M_{temp} = 12.3^{\circ}\text{C}$  (54.14°F),  $min_{temp} = -11.5^{\circ}\text{C}$ ,  $max_{temp} = 35.5^{\circ}\text{C}$ ; for 2019  $M_{temp} = 12.2^{\circ}\text{C}$  (53.96°F),  $min_{temp} = -11.5^{\circ}\text{C}$ ,  $max_{temp} = 35.5^{\circ}\text{C}$ ; for 2020  $M_{temp} = 12.4^{\circ}\text{C}$  (54.32°F),  $min_{temp} = -6.1^{\circ}\text{C}$ ,  $max_{temp} = 37.3^{\circ}\text{C}$ ; for 2021  $M_{temp} = 10.9^{\circ}\text{C}$  (51.62°F),  $min_{temp} = -9.9^{\circ}\text{C}$ ,  $max_{temp} = 33.3^{\circ}\text{C}$ ; for 2022  $M_{temp} = 12.9^{\circ}\text{C}$  (55.22°F),  $min_{temp} = -7.7^{\circ}\text{C}$ ,  $max_{temp} = 37.2^{\circ}\text{C}$ . In addition, the mean temperature in degrees Celsius (again, Fahrenheit in between parentheses) of the month(s) we conducted the study were  $16^{\circ}\text{C}$  (60.8°F) for 2018,  $18^{\circ}\text{C}$  (64.4°F) for 2019,  $10^{\circ}\text{C}$  (50°F) for 2020,  $6.5^{\circ}\text{C}$  (43.7°F) for 2021.



points (Liljequist et al., 2019). This was further confirmed in our robustness exploratory analyses. But there is random error in our measurement: it can be qualitatively observed from the omega values of our subscales. Overall, the internal consistencies of the subscales are acceptable, but far from excellent, and are likely to have reduced the correlation obtained from our test-retests, since low internal consistency is known to reduce the observed correlation between constructs (Reis & Judd, 2000).

### Constraints On Generality

We conducted our study on mostly female students (89.44%) who were 19.70 years old on average. Age is known to be an important predictor of people's thermoregulatory abilities, especially in older age (Khan et al., 1992). Thus, a different sample of older participants might affect how the STRAQ-1 items would be perceived and how participants would respond to them. Thus, different samples, including older people for example, would likely result in different findings compared to the ones provided in our study. Future studies using the STRAQ-1, or closely related constructs, such as the STRAEQ-2 (Dujols et al., 2024), should investigate if our result replicates in significantly different samples to further explore the psychometrics of the measure.

We measured the STRAQ-1 subscales over long periods of time and in similar contexts (an online questionnaire in the spring-to-winter period). Future studies could further explore the stability of the STRAQ-1 using intra-individual designs, including more repeated measures: using, for example, ecological moment assessment, to further investigate potential seasonality in people's response to the STRAQ-1. People's ratings, for example, of the High-Temperature sensitivity subscale could vary according to different moments of the same day, or between summer and winter.

### Conclusion

This Registered Report provides the first test of measurement invariance across two time points (separated by approximately a year or more) and estimates of test-retest reliability over the same period for the Social Thermoregulation, Risk Avoidance Questionnaire (STRAQ-1, Vergara et al., 2019). In our sample, the Social Thermoregulation subscale, the Solitary Thermoregulation subscale, and the Risk Avoidance subscale had similar internal consistencies to those reported by Vergara et al. (2019). We concluded that only two of the STRAQ-1 subscales were longitudinally

invariant across two-time points. Additionally, we found that test-retest reliability was overall moderate to good for *Social Thermoregulation*, overall moderate for *Solitary Thermoregulation* and *High-Temperature Sensitivity*, and overall poor to moderate for *Risk Avoidance*. Our study suggests that test-retest reliability was insufficient for psychological diagnosis, and that future studies should address the problem of low generalizability.

### Competing Interests

The authors of this article declare that they have no financial conflict of interest with the content of this article. None of the authors are recommenders at PCI Registered Reports at the time of the submission of this article.

### Author Note

The approved STAGE 1 version of the manuscript is available on OSF at this link: <https://osf.io/qz5eg>.

### Data Accessibility Statement

The data is available at this link: <https://osf.io/u58vk/>; and materials at this link: <https://osf.io/4wapd/>.

### Author Contributions

Conceptualization: Olivier Dujols (Lead), Siegwart Lindenberg (Equal), Hans IJzerman (Equal). Data curation: Olivier Dujols (Lead). Formal Analysis: Olivier Dujols (Lead), Hans IJzerman (Supporting). Funding acquisition: Olivier Dujols (Supporting), Hans IJzerman (Lead). Investigation: Olivier Dujols (Supporting), Hans IJzerman (Lead). Methodology: Olivier Dujols (Lead), Hans IJzerman (Equal). Project administration: Olivier Dujols (Supporting), Hans IJzerman (Lead). Resources: Olivier Dujols (Lead), Hans IJzerman (Equal). Software: Olivier Dujols (Lead). Visualization: Olivier Dujols (Lead). Writing – original draft: Olivier Dujols (Lead), Hans IJzerman (Supporting). Writing – review & editing: Olivier Dujols (Lead), Siegwart Lindenberg (Equal), Hans IJzerman (Equal). Supervision: Siegwart Lindenberg (Equal), Hans IJzerman (Lead). Validation: Caspar J. Van Lissa (Lead).

Submitted: July 02, 2024 PDT, Accepted: July 05, 2024 PDT



This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License (CCBY-4.0). View this license's legal deed at <http://creativecommons.org/licenses/by/4.0> and legal code at <http://creativecommons.org/licenses/by/4.0/legalcode> for more information.

## References

- Ainsworth, M. S. (1979). Infant–mother attachment. *American Psychologist*, *34*(10), 932–937. <https://doi.org/10.1037/0003-066X.34.10.932>
- Bleidorn, W., Schwaba, T., Zheng, A., Hopwood, C. J., Sosa, S. S., Roberts, B. W., & Briley, D. A. (2022). Personality stability and change: A meta-analysis of longitudinal studies. *Psychological Bulletin*, *148*(7–8), 588–619. <https://doi.org/10.1037/bul0000365>
- Bowlby, J. (1969–1982). Attachment and Loss, Vol. 1: Attachment. In *Attachment and loss*. Basic Books.
- Brennan, K. A., Clark, C. L., & Shaver, P. R. (1998). Self-report measurement of adult romantic attachment: An integrative overview. In J. A. Simpson & W. S. Rholes (Eds.), *Attachment theory and close relationships* (pp. 46–76). Guilford Press.
- Bretherton, I., & Munholland, K. A. (2008). Internal working models in attachment relationships: Elaborating a central construct in attachment theory. In J. Cassidy & P. R. Shaver (Eds.), *Handbook of attachment: Theory, research, and clinical applications* (pp. 102–127). Guilford Press.
- Bühler, J. L., Orth, U., Bleidorn, W., Weber, E., Kretzschmar, A., Scheling, L., & Hopwood, C. J. (2023). Life Events and Personality Change: A Systematic Review and Meta-Analysis. *European Journal of Personality*, *1*(1). <https://doi.org/10.1177/08902070231190219>
- Bujang, M. A., & Baharum, N. (2017). A simplified guide to determination of sample size requirements for estimating the value of intraclass correlation coefficient: a review. *Archives of Orofacial Science*, *12*(1).
- Bystrova, K., Matthiesen, A. S., Vorontsov, I., Widström, A. M., Ransjö-Arvidson, A. B., & Uvnäs-Moberg, K. (2007). Maternal axillar and breast temperature after giving birth: effects of delivery ward practices and relation to infant temperature. *Birth*, *34*(4), 291–300. <https://doi.org/10.1111/j.1523-536X.2007.00187.x>
- Chan, C., Chan, G. C. H., Leeper, T. J., & Becker, J. (2021). *Rio: A swiss-army knife for data file I/O* (R package version 0.5.29) [Computer software].
- Chen, F. F. (2008). What happens if we compare chopsticks with forks? The impact of making inappropriate comparisons in cross-cultural research. *Journal of Personality and Social Psychology*, *95*(1), 1005–1018. <https://doi.org/10.1037/a0013193>
- Cheung, G. W., & Rensvold, R. B. (2002). Evaluating goodness-of-fit indexes for testing measurement invariance. *Structural Equation Modeling*, *9*(2), 233–255. [https://doi.org/10.1207/S15328007SEM0902\\_5](https://doi.org/10.1207/S15328007SEM0902_5)
- Coen, A. B., & Robert, I. J. (2005). Gradient projection algorithms and software for arbitrary rotation criteria in factor analysis. *Educational and Psychological Measurement*, *65*(1), 676–696. <https://doi.org/10.1177/0013164404272507>
- De Wolff, M. S., & van IJzendoorn, M. H. (1997). Sensitivity and Attachment: A Meta-Analysis on Parental Antecedents of Infant Attachment. *Child Development*, *68*(1), 571–591. <https://doi.org/10.1111/j.1467-8624.1997.tb04218.x>
- Dimitrov, D. M. (2014). *Statistical methods for validation of assessment scale data in counseling and related fields*. John Wiley & Sons.
- Dragow, F., & Kanfer, R. (1985). Equivalence of psychological measurement in heterogeneous populations. *Journal of Applied Psychology*, *70*(4), 662–680. <https://doi.org/10.1037/0021-9010.70.4.662>
- Dugan, K. A., & Fraley, R. C. (2022). The roles of parental and partner attachment working models in romantic relationships. *Journal of Social and Personal Relationships*, *39*(7), 2154–2180. <https://doi.org/10.1177/02654075221075254>
- Dujols, O., Klein, R. A., Lindenberg, S., STRAEQ-2 team, & IJzerman, H. (2024). *Development and validation of the Social Thermoregulation, Risk Avoidance, and Eating Questionnaire - 2 (STRAEQ-2)*. <https://osf.io/ggbzk>
- Epskamp, S. (2022). *SemPlot: Path diagrams and visual analysis of various SEM packages' output* (R package version 1.1.6) [Computer software].
- Evers, A., Lucassen, W., Meijer, R., & Sijtsma, K. (2015). COTAN review system for evaluating test quality. *NIP*. <https://psynip.nl/wp-content/uploads/2022/05/COTAN-review-system-for-evaluating-test-quality.pdf>
- Firke, S. (2021). *Janitor: Simple tools for examining and cleaning dirty data* (R package version 2.1.0) [Computer software].
- Flora, D. B., & Curran, P. J. (2004). An empirical evaluation of alternative methods of estimation for confirmatory factor analysis with ordinal data. *Psychological Methods*, *9*(4), 466–491. <https://doi.org/10.1037/1082-989X.9.4.466>
- Fraley, R. C. (2019). Attachment in adulthood: Recent developments, emerging debates, and future directions. *Annual Review of Psychology*, *70*(1), 401–422. <https://doi.org/10.1146/annurev-psych-010418-102813>
- Fraley, R. C., Dugan, K. A., Thompson, R. A., Simpson, J. A., & Berlin, L. J. (2021). The consistency of attachment security across time and relationships. In R. A. Thompson, J. A. Simpson, & L. J. Berlin (Eds.), *Attachment: The fundamental questions*. Guilford Press.
- Fraley, R. C., Waller, N. G., & Brennan, K. A. (2000). An item response theory analysis of self-report measures of adult attachment. *Journal of Personality and Social Psychology*, *78*(2), 350–365. <https://doi.org/10.1037/0022-3514.78.2.350>
- Gilbert, C., McCafferty, D., Le Maho, Y., Martrette, J. M., Giroud, S., Blanc, S., & Ancel, A. (2010). One for all and all for one: the energetic benefits of huddling in endotherms. *Biological Reviews*, *85*(3), 545–569. <https://doi.org/10.1111/j.1469-185X.2009.00115.x>

- Goncharova, M., Silan, M. A., Dujols, O., Stoianova, T., Sparacio, A., Adetula, A., & IJzerman, H. (2019). *The CO-RE Lab Lab Philosophy*. <https://psyarxiv.com>
- IJzerman, H., Heine, E. C., Nagel, S. K., & Pronk, T. M. (2017). Modernizing relationship therapy through Social Thermoregulation Theory: Evidence, hypotheses, and explorations. *Frontiers in Psychology*, 8(1), 635. <https://doi.org/10.3389/fpsyg.2017.00635>
- IJzerman, H., Neyroud, L., Courset, R., Schrama, M., Post, J., & Pronk, T. (2018). Socially thermoregulated thinking: How past experiences matter in thinking about our loved ones. *Journal of Experimental Social Psychology*, 79(1), 349–355. <https://doi.org/10.1016/j.jesp.2018.08.008>
- Khan, F., Spence, V. A., & Belch, J. J. F. (1992). Cutaneous vascular responses and thermoregulation in relation to age. *Clinical Science*, 82(1), 521–528. <https://doi.org/10.1042/cs0820521>
- Kline, R. B. (2016). *Principles and practice of structural equation modeling* (Fourth edition). Guilford Press.
- Koo, T. K., & Li, M. Y. (2016). A guideline of selecting and reporting intraclass correlation coefficients for reliability research. *Journal of Chiropractic Medicine*, 15(2), 155–163. <https://doi.org/10.1016/j.jcm.2016.02.012>
- Li, C. H. (2016). Confirmatory factor analysis with ordinal data: Comparing robust maximum likelihood and diagonally weighted least squares. *Behavior Research Methods*, 48(1), 936–949. <https://doi.org/10.3758/s13428-015-0619-7>
- Liljequist, D., Elfving, B., & Skavberg Roaldsen, K. (2019). Intraclass correlation—A discussion and demonstration of basic features. *PloS One*, 14(7), e0219854. <https://doi.org/10.1371/journal.pone.0219854>
- Mackinnon, S., Curtis, R., & O'Connor, R. (2022). A tutorial in longitudinal measurement invariance and cross-lagged panel models using lavaan. *Meta-Psychology*, 6(1). <https://doi.org/10.31234/osf.io/tkzrb>
- Main, M., & Solomon, J. (1986). Discovery of an insecure-disorganized/disoriented attachment pattern. In T. B. Brazelton & M. Yogman (Eds.), *Affective development in infancy* (pp. 95–124). Ablex Publishing.
- Main, M., & Solomon, J. (1990). Procedures for identifying infants as disorganized/disoriented during the Ainsworth Strange Situation. In *Attachment in the preschool years: Theory, research, and intervention* (pp. 121–160). University of Chicago Press.
- Moshagen, M., & Erdfelder, E. (2016). A new strategy for testing structural equation models. *Structural Equation Modeling*, 23(1), 54–60. <https://doi.org/10.1080/10705511.2014.950896>
- Nunnally, J. C., & Bernstein, I. H. (1994). *Psychometric theory*. McGraw-Hill.
- O'Connor, B. P. (2022). *EFA.dimensions: Exploratory Factor Analysis functions for assessing dimensionality* (R package version 0.1.7.4) [Computer software].
- Parsons, S., Kruijt, A. W., & Fox, E. (2019). Psychological science needs a standard practice of reporting the reliability of cognitive-behavioral measurements. *Advances in Methods and Practices in Psychological Science*, 2(4), 378–395. <https://doi.org/10.1177/2515245919879695>
- Reis, H. T., & Judd, C. M. (Eds.). (2000). *Handbook of research methods in social and personality psychology*. Cambridge University Press.
- Revelle, W. (2022). *Psych: Procedures for personality and psychological research* (R package version 2.2.9) [Computer software]. Northwestern University.
- Rizzo, M., & Szekely, G. (2022). *Energy: E-statistics, multivariate inference via the energy of data* (R package version 1.7-10) [Computer software].
- Roberts, B. W., & DelVecchio, W. F. (2000). The rank-order consistency of personality traits from childhood to old age: a quantitative review of longitudinal studies. *Psychological Bulletin*, 126(1), 3–25. <https://doi.org/10.1037/0033-2909.126.1.3>
- Rocha IJzerman, H. (2021). *Heartwarming: How our inner thermostat made us human*. WW Norton & Company.
- Rodebaugh, T. L., Scullin, R. B., Langer, J. K., Dixon, D. J., Huppert, J. D., Bernstein, A., & Lenze, E. J. (2016). Unreliability as a threat to understanding psychopathology: The cautionary tale of attentional bias. *Journal of Abnormal Psychology*, 125(6), 840. <https://doi.org/10.1037/abn0000184>
- Rosseel, Y. (2012). Lavaan: An R package for structural equation modeling. *Journal of Statistical Software*, 48(1), 1–36. <https://doi.org/10.18637/jss.v048.i02>
- Sarda, E., Klein, R. A., Dujols, O., & IJzerman, H. (2021). Validation of the ISP131001 sensor for mobile peripheral body temperature measurement. *International Review of Social Psychology*, 34(1), Article12. <https://doi.org/10.5334/irsp.409>
- Van De Schoot, R., Schmidt, P., De Beuckelaer, A., Lek, K., & Zondervan-Zwijnenburg, M. (2015). Measurement invariance. *Frontiers in Psychology*, 6(1), 1064. <https://doi.org/10.3389/fpsyg.2015.01064>
- Vergara, R. C., Hernández, C., Jaume-Guazzini, F., Lindenberg, S., Klein, R. A., & IJzerman, H. (2019). Development and Validation of the Social Thermoregulation and Risk Avoidance Questionnaire (STRAQ-1). *International Review of Social Psychology*, 32(1), 18. <https://doi.org/10.5334/irsp.222>
- Vidal, N., Costello, J., Ribotta, B., Gurgand, L., & IJzerman, H. (2022). Assessing the reliability of an infrared thermography protocol to assess cold-induced Brown Adipose Tissue activation in young adults. *Social Psychological Bulletin*. In press. <https://psyarxiv.com/rkde9/download?format=pdf>
- Watson, D. (2004). Stability versus change, dependability versus error: Issues in the assessment of personality over time. *Journal of Research in Personality*, 38(4), 319–350. <https://doi.org/10.1016/j.jrp.2004.03.001>

- Wickham, H., Averick, M., Bryan, J., Chang, W., D'Agostino McGowan, L., François, R., Golemund, G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Lin Pedersen, T., Miller, E., Milton Bache, S., Müller, K., Ooms, J., Robinson, D., Paige Seidel, D., Spinu, V., & Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), Article1686. <https://doi.org/10.21105/joss.01686>
- Widaman, K. F., Ferrer, E., & Conger, R. D. (2010). Factorial invariance within longitudinal structural equation models: Measuring the same construct across time. *Child Development Perspectives*, 4(1), 10–18. <https://doi.org/10.1111/j.1750-8606.2009.00110.x>
- Wittmann, A., Braud, M., Dujols, O., Forscher, P., & IJzerman, H. (2022). Individual differences in adapting to temperature in French students are only related to attachment avoidance and loneliness. *Royal Society Open Science*, 9(5), Article201068. <https://doi.org/10.1098/rsos.201068>
- Zou, G. Y. (2012). Sample size formulas for estimating intraclass correlation coefficients with precision and assurance. *Statistics in Medicine*, 31(29), 3972–3981. <https://doi.org/10.1002/sim.5466>

## Supplementary Materials

### Peer Review Correspondence

Download: [https://collabra.scholasticahq.com/article/122155-test-retest-reliability-of-the-straq-1-a-registered-report/attachment/240087.docx?auth\\_token=RoxnzcEXIjYExa\\_IOMBb](https://collabra.scholasticahq.com/article/122155-test-retest-reliability-of-the-straq-1-a-registered-report/attachment/240087.docx?auth_token=RoxnzcEXIjYExa_IOMBb)

---