

1 **Assessment of Subjective Emotional Valence and Long-Lasting Impact of Life Events:**

2 **Development and Psychometrics of the Stralsund Life Event List (SEL)**

3 Johanna König, Andrea Block, PhD, Mathias Becker, Kristin Fenske, Johannes Hertel, Sandra Van der Auwera, PhD,

4 Kathleen Zymara, Henry Völzke, MD, Harald Jürgen Freyberger, MD, Hans Jürgen Grabe, MD

5

6 **SUPPLEMENTARY MATERIAL A - Reliability Analyses**

7

8 **METHODS**

9 **Statistical Analyses**

10 For inter-rater and test-retest analyses, the SEL interviews of the reliability sample (N=19) were used. For both
11 *analyses*, an agreement of the *occurrence* of each LE was calculated. To assess the accuracy of the *occurrence time*
12 *coding*, several LEs were selected. To achieve a suitable statistical power, the selection was based on the criteria
13 that a given LE had to be present at both time points of *the* reliability measurement in at least half of the
14 interviewees (n=5). These selection criteria resulted in 23 LEs for the inter-rater reliability *assessment* and 16 LEs
15 for the test-retest reliability *assessment* (supplementary table S4). The accuracy of the *occurrence time coding* was
16 calculated by assessing the agreement of the 5-year period coded using weighted Cohen's kappa *values* [1–3]. *To*
17 *assess* the inter-rater and test-retest reliability of the SEL scores, averaged intraclass correlations (ICC) were
18 calculated *based on* answer agreements between the averages of ratings *from* several raters [1, 4, 5].

19

20 **RESULTS**

21 *Inter-Rater Reliability.* Ten participants (20.0% male; age: M=40.9, SD=12.2) *from* the reliability sample were re-
22 interviewed by a second interviewer two days (M=1.8, SD=0.4) after the initial interview. *On* average, the
23 *occurrence* agreement between the two interviews was 94.3%, with 42 LEs (60%) being reported *as* identical in
24 both interviews (table S3). Except for four items (H6, B16, L31 and T69), the weighted Cohen's kappa (κ) *for*
25 *occurrence time coding* was higher than 0.65 (Table S4), which indicates a substantial to almost perfect accuracy [1,
26 2]. As depicted in table 4, except *for* the *present impact* of past positive LEs (ICC=0.41, p=0.208), all inter-rater ICCs

27 for the SEL scores were higher than 0.69 and thus can be described as good; all scores for negative LEs were higher
28 than 0.75 and were, thus, excellent [1]. For more detailed information, see tables 4, S3 and S4.

29
30 *Test-Retest Reliability.* Nine participants (55.6% male; M=45.2, SD=13.3) from the reliability sample were re-
31 interviewed by the same interviewer 28 days (M=28.7, SD=4.7) after the initial interview. On average, the
32 occurrence agreement between the two interviews was 94.3%, with 28 LEs (40%) being reported as identical in
33 both interviews (table S3). The mean κ for the *occurrence time coding* was 0.68 (Table S4). Only the items H6
34 ($\kappa=0.36$, $p=0.097$) and W52 ($\kappa=0.24$, $p=0.232$) had a κ lower than 0.4, which was defined as moderate [1]. Except for
35 the *present impact* of past negative LEs (ICC=0.29, $p=0.308$), all ICCs of the SEL scores (0.73-0.97) indicated
36 excellent test-retest reliabilities [1]. For more detailed information, see tables 4, S3 and S4.

37

38 DISCUSSION

39 *Inter-Rater and Test-Retest Reliability.* The lowest inter-rater agreement (60%) was observed for item U78 “serious
40 accident, terrible experiences or catastrophes as listed”. Reporting traumatic experiences might be emotionally
41 demanding. Thus, the responses to item U78 might depend not only on the pure facts of the occurrence but also on
42 the relationship with and faith in the interviewer. Accordingly, the test-retest agreement of item U78 was much
43 higher (89%) than the inter-rater agreement. Low inter-rater reliabilities were also observed for item H6 “move out
44 of childhood home” and L31 “complete traineeship/pass examinations”. Both items had long latency periods
45 between the LE occurrence and the interview, which has been demonstrated to reduce recall accuracy [6]. Cohen’s
46 kappa for item W52 was critically low only in the test-retest analysis, which indicates a high test-retest-interval
47 sensitivity. Overall, inter-rater reliability, which features the comparability and equality of the ratings of different
48 interviewers, was high and thus indicated high implementation and evaluation objectivity. The test-retest reliability
49 was comparable to the inter-rater reliability for the SEL scores. Please note that the test-retest and inter-rater
50 reliability analyses were based on a small clinical sample. Hence, these preliminary results need to be interpreted
51 with caution and require future research using larger sample sizes.

52

53 *Limitations.* The preliminary test-retest and inter-rater reliability analyses were based on a small clinical sample
54 with all subjects suffering from mental disorders and receiving psychiatric or psychotherapeutic treatments at the
55 time of the interviews. As these interventions might specifically have changed the emotional evaluations of past LEs
56 and their impact on present wellbeing between the interviews, we would expect even better reliability in healthy
57 subjects. However, it is highly important to test this hypothesis in future research. Moreover, according to Brown
58 and Harris (1982) [6], a “fall-off” regarding the number of reported LEs is expected with increasing time between
59 interviews, particularly for LEs that are low in threat and salience. Nevertheless, our preliminary reliability indices
60 were good but need to be supported by further reliability studies in healthy subjects with larger sample sizes.

61

62 SUPPLEMENTARY REFERENCES

- 63 1. Hallgren KA. Computing Inter-Rater Reliability for Observational Data: An Overview and Tutorial.
64 Tutor Quant Methods Psychol. 2012;8:23–34.
- 65 2. Landis JR, Koch GG. The measurement of observer agreement for categorical data. Biometrics.
66 1977:159–74.
- 67 3. Cohen J. Weighted kappa: Nominal scale agreement provision for scaled disagreement or partial
68 credit. Psychol Bull. 1968;70:213–20.
- 69 4. Rousson V, Gasser T, Seifert B. Assessing intrarater, interrater and test-retest reliability of
70 continuous measurements. Stat Med. 2002;21:3431–46. doi:10.1002/sim.1253.
- 71 5. Stata Corp. Stata 12 Base Reference Manual. College Station, Texas: Stata Press; 2011.
- 72 6. Brown GW, Harris TO. Fall-off in the reporting of life events. Soc Psychiatry. 1982;17:23–8.
73 doi:10.1007/BF00583889.