

A Multiverse Assessment of the Reliability of the Self Matching Task as a Measurement of the Self-Prioritization Effect

Zheng Liu^{1,2†}, Mengzhen Hu^{1†}, Yuanrui Zheng¹, Jie Sui³,
Hu Chuan-Peng^{1*}

^{1*}School of Psychology, Nanjing Normal University, Nanjing, China.

^{2*} School of Humanities and Social Science, The Chinese University of
Hong Kong-Shenzhen, Shenzhen, China.

^{3*}School of Psychology, University of Aberdeen, Old Aberdeen, Scotland.

*Corresponding author(s). E-mail(s): hu.chuan-peng@nnu.edu.cn;
hcp4715@hotmail.com;

[†]These authors contributed equally to this work.

Abstract

The Self Matching Task (SMT) is widely used to investigate the cognitive mechanisms underlying the Self-Prioritization Effect (SPE), wherein performance is enhanced for self-associated stimuli compared to other-associated ones. Although the SMT robustly elicits the SPE, there is a lack of quantifying the reliability of this paradigm. This ignorance is problematic, given the prevalence of the reliability paradox in cognitive tasks: many well-established cognitive tasks demonstrate relatively low reliability when used to evaluate individual differences, despite exhibiting replicable effects at the group level. To fill this gap, this preregistered study investigated the reliability of SPE derived from the SMT using a multiverse approach, combining all possible indicators and baselines reported in the literature. We first examined the robustness of 24 SPE measures across 42 datasets ($N = 2250$) using a meta-analytical approach. We then calculated the Split-Half Reliability (r) and Intraclass Correlation Coefficient (ICC2) for each SPE measure. Our findings revealed a robust group-level SPE across datasets. However, when evaluating individual differences, SPE indices derived from Reaction Time (RT) and Efficiency exhibited relatively higher, compared to other SPE indices, but still unsatisfied split-half reliability (approximately 0.5). For the reliability across multiple time points, as assessed by ICC2, RT and Efficiency demonstrated moderate levels of test-retest reliability (close to 0.5). These findings revealed the

presence of a reliability paradox in the context of SMT-based SPE assessment. We discussed the implications of how to enhance individual-level reliability using this paradigm for future study design.

Keywords: Self-Prioritization Effect (SPE), Self Matching Task (SMT), Reliability, Multiverse

1 Introduction

The Self-Prioritization Effect (SPE) reflects individuals' biased responses towards self-related information in comparison to information related to others. This phenomenon holds a central position within cognitive psychology and underscores a core facet of human cognition and self-awareness (Sui & Humphreys, 2017). SPE has been found in a broad range of cognitive tasks (e.g., Cunningham et al., 2008; Rogers et al., 1977; Sui et al., 2012). Despite SPE is often argued to be a self-specific effect, it has been challenging to be disassociated from the familiarity effect. That is, the self-related stimuli, such as own faces (Keenan et al., 2000; Kircher et al., 2000; Turk et al., 2002), own voices (S. M. Hughes & Harrison, 2013; Payne et al., 2021), or own names (Constable, Rajsic, et al., 2019) are usually more familiar to participants than those other-related stimuli. To overcome such limitation, Sui et al. (2012) introduced the Self Matching Task (SMT), where the self-relatedness (and other-relatedness) was acquired in the lab. In this task, participants first associated geometric shapes with person labels (e.g., circle = you, triangle = best friend, square = stranger) and then performed a matching task, judging whether a shape-label pair presented on the screen matched the acquired relationship. A typical pattern from this task is that shapes associated with the self exhibit a processing advantage over shapes related to others. This SPE from SMT has subsequently been replicated by many researchers (Constable, Elekes, et al., 2019; Golubickis et al., 2020; Golubickis et al., 2017; Hu et al., 2020), highlighting the robustness of the effect.

The reliability of SMT as a measurement of SPE, however, has not been examined. Here, the reliability of a cognitive task refers to its ability in producing consistent results for the same person across sessions or times (Parsons et al., 2019; Zorowitz & Niv, 2023). One common method to assess reliability is the Split-Half Reliability (r), where a test is divided into two halves, and the correlation between the data from these two halves is calculated. A high correlation suggests that the test is internally consistent and measures the same construct reliably (Pronk et al., 2022). Another widely used method is Test-retest reliability, which refers to the extent to which a measurement or assessment tool produces consistent and stable results over time when administered to the same group of individuals under identical conditions (Kline, 2015). Both methods are from classical test theory in psychometrics (Borsboom, 2005), but they are less known to experimental psychologists. In experimental research, researchers focus on the robustness of experimental effects. Robustness, in this context, pertains to the extent to which a cognitive task consistently produces the same effect at the group level across various independent participant samples. For

example, the “group effect” in the Stop-Signal Task refers to differences in Reaction time between different stop-signal delays (Hedge et al., 2018). An effect is considered robust if these differences can be consistently observed in different samples performing the Stop-Signal Task.

In recent years, driven by a growing interest in employing cognitive tasks to assess individual differences, researchers have turned their attention to evaluating the reliability of cognitive tasks (e.g., Hedge et al., 2018; Kucina et al., 2023). However, existing findings have raised concerns about the reliability of many cognitive tasks (Karvelis et al., 2023; Rouder & Haaf, 2019), with a considerable body of research highlighting the moderate to low-level reliability found in the cognitive task measurements (Clark et al., 2022; Enkavi et al., 2019; Green et al., 2016). For instance, Hedge et al. (2018) reported a range of test-retest reliabilities about frequently employed experimental task metrics (such as Stroop and Stop-Signal Task), with a notable prevalence of discrepancy between the low reliability for individual differences and the robustness of the experimental effects. This discrepancy, named the “reliability paradox” (Logie et al., 1996), has gained much attention in recent years. Like other cognitive tasks, SPMT was also employed by researchers as a measure of individual differences in SPE. For example, a recent study examined the individual differences of SPE and how these individual differences are correlated to brain network (Zhang et al., 2023). Likewise, in clinical investigation, the SMT has been incorporated to assess deviations in self-processing among specific populations, including individuals affected by autism or depression (e.g., Hobbs et al., 2023; Liu et al., 2022; Moseley et al., 2022). The findings from these studies are diverse. On one hand, research has demonstrated that behavioral data from SMT could function as a viable marker for depression screening (Liu et al., 2022). Additionally, performance in SMT has been employed to decode brain functional connectivity during resting state (Zhang et al., 2023) or understand the functions of self-associations in cognition (Scheller & Sui, 2022a, 2022b; Yankouskaya et al., 2023). These studies suggest the potential for significant individual-level variability in SMT performance. On the other hand, Hobbs et al. (2023) assessed the role of self-referencing in relation to depression using SMT but found a limited association between individuals’ performance in SMT and depression scores. Moseley et al. (2022) also found inconsistent correlations between SPE and its relationship to autistic traits, mentalizing ability and loneliness. These conflicting trends underscore the need to evaluate the reliability of SMT as a measurement of SPE.

Further, the variability in quantifying SPE using SMT calls for a comprehensive examination of the reliability of different SPE measures. As simple as the SMT, there are multiple approaches to quantify the SPE, encompassing various indicators and baselines. In a typical SMT experiment, two direct outcomes are generated: Reaction Time (RT) and choices. The RT and Accuracy (ACC) of choices are the two most widely used indicators of SPE. Several other indicators can be derived from these direct outcomes: Efficiency (η) (Humphreys & Sui, 2015; Stoeber & Eysenck, 2008), sensitivity score (d -prime, d') of Signal Detection Theory (Hu et al., 2020; Sui et al., 2012), drift rate (v) and starting point (z) estimated using the Drift-Diffusion Model (DDM) (Macrae et al., 2017; Reuther & Chakravarthi, 2017). In addition to the variability of indicators, SPE can be estimated by calculating the difference between self

condition and different baselines. Indeed, the selection of baselines varies across studies, such as “Close other” (e.g., Friend) (Navon & Makovski, 2021; Svensson et al., 2022), “Stranger” (Constable et al., 2021; Orellana-Corrales et al., 2020), “Celebrity” (e.g., “LuXun”) (Qian et al., 2020) and “Non-person” (e.g., None) (Schäfer & Frings, 2019). As a result, three pivotal questions regarding the reliability of the SMT remain unresolved: First, given the variability of indicators (RT, ACC, d' , η , v , z) and choice of baseline conditions (“Close other”, “Stranger”, “Celebrity”, and “Non-person”), which way of quantifying SPE is the most reliable one(s)? Second, is the SMT suitable for assessing individual differences in SPE? Finally, is there a reliability paradox in the assessment of SPE using SMT? Addressing these questions is crucial for SMT-based measurements, allowing for an accurate assessment of the SPE and its applications in various domains.

To address these three questions, the present study adopted a multiverse approach to investigate the reliability of SPE measures computed using different indicators under various baseline conditions in the SPMT. This was achieved by re-analysing 42 independent datasets ($N = 2250$) from 24 papers and 3 unpublished projects that employed the SMT. In order to comprehensively assess the SPE measures derived from SPMT, we created a “multiverse” of possible indicators (RT, ACC, d -prime, η , v , z) combined with various baseline conditions (“Close other”, “Stranger”, “Celebrity”, and “Non-person”). We first assessed the experimental effect across this multiverse using meta-analysis. The individual level consistency was examined using permutation-based Split-Half Reliability (r) and Intraclass Correlation Coefficient (ICC2, Two-way random effect model, absolute agreement) for assessing the consistency of task performance over time. The findings of our study provided valuable insights into the reliability of SMT and its indicators, having the potential to facilitate the future utilization of SMT in research, clinical settings, and personal performance monitoring.

2 Methods

2.1 Ethics Information

As this study is a secondary analysis of pre-existing data sourced from publicly available datasets or archived data previously collected by the author’s group, informed consent and confidentiality are not applicable.

2.2 Experimental Design

Here we provided a detailed overview of the original experimental design of SMT, as described in Experiment 1 by Sui et al. (2012). The original SMT used a 2 by 3 within-subject design. The first independent variable, labelled “Matching,” consisted of two levels: “Matching” and “Non-matching”, indicating whether the shape and label were congruent. The second independent variable, labelled “Identity”, comprised three levels: “Self”, “Friend”, and “Stranger”, representing the corresponding identity associated with the shape.

The original SMT consisted of two stages (refer to Fig. 1). In the first stage (instructional stage), participants were instructed to associate three geometric shapes (circle,

triangle and square) with three labels (self, friend, and stranger) for approximately 60 seconds. The shape-label associations were counterbalanced between participants. In the second phase (matching task), participants completed a matching task. Each trial started with a fixation cross displayed in the center of the screen for 500 ms, followed by a shape-label pairing and fixation cross for 100 ms. The screen then went blank for 800–1200 ms, or until a response was made. Participants were required to judge whether the presented shape and label matched the learned associations from the learning phase and respond as quickly and accurately as possible by pressing one of two buttons within the allowed timeframe. Prior to the formal experimental phase, participants completed a training session consisting of 24 practice trials.

After the training, participants completed six blocks of 60 trials in the matching task, with two matching types (matching/non-matching) and three shape associations, for a total of 60 trials per association. Short breaks lasting up to 60 seconds were provided after each block.

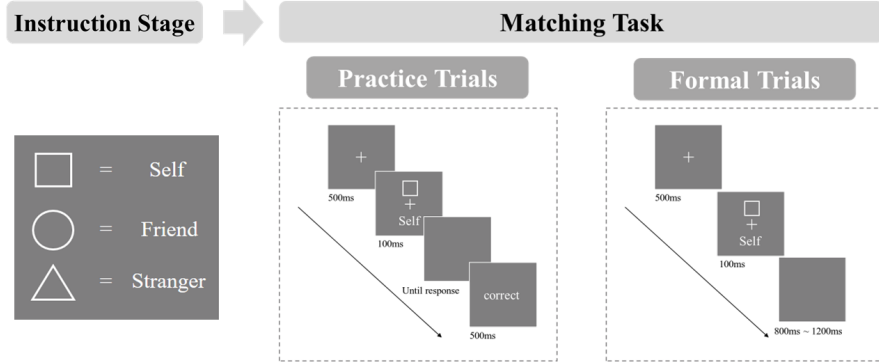


Fig. 1 Procedure of the original SPMT in Experiment 1 (Sui et al., 2012). *Note:* The relation between shape-label pairs was counter-balanced between participants.

2.3 Datasets Acquisition

Initially, two datasets that employed the SPMT were available to us: one from an unpublished project conducted in our laboratory (Hu et al., 2023), for which we provide more details in the supplementary materials (in section 1.1), and the other provided by our collaborators (Liu et al., 2023). Concurrently, we are conducting a meta-analysis on SPE using the SPMT (pre-registration available at OSF (<https://osf.io/euqmf>)). During this process, we identified an additional 24 papers with datasets potentially suitable for our present study. The detailed paper selection procedure was presented in supplementary material, Figure S2. The selection of these papers was based on specific criteria:

- 1) The paper must primarily utilize the SPMT as their method.
- 2) The experimental design should not incorporate any stimuli that could potentially trigger a familiarity effect (e.g., using self-face, self-name).
- 3) The trial-level data is either openly available or shared with us upon request, enabling us to estimate at least one reliability index.

Specifically, we identified a total of 41 papers with potentially accessible data via screening of related databases. Of these papers, 13 papers made their trial-level data publicly available. For the remaining 28 papers, we reached out to the authors and requested access to their trial-level data. Out of those 28 requests, 11 papers provided us with trial-level data. During revision, we obtained additionally two unpublished datasets (Sui 2014a, 2015). In total, our analysis comprised raw data from 24 papers and 3 unpublished projects from our laboratory and collaborators.

It is important to highlight that the research culture discourages direct replications (Makel et al., 2012). As a result, all the datasets included in our analysis underwent some degrees of modification to the original design (e.g., change shapes, modify sequence) as well as including additional independent variables (refer to Table 1 for specification). For our analysis, we focused exclusively on datasets that adhered to the design of SMT without incorporating any stimuli that could potentially trigger a familiarity effect (e.g., oneself or friends' name or face). Procedural differences from the original matching task (e.g., the timing of stimulus presentation; and the nature of stimuli used), were considered secondary to the overarching criteria of stimulus neutrality. For datasets from experiments that manipulated other independent variables (e.g., mood), we only utilized data from control conditions so that the data were close to the original design of SMT.

In the end, we were able to incorporate 42 independent datasets from the above-mentioned papers and projects. Nonetheless, not all studies incorporated retest sessions. If a publicly available dataset did not include a retest session with SMT, we excluded it from calculating the Intraclass Correlation Coefficient and only considered the split-half reliability. The details of the included studies and conditions in the datasets are described in Table 1.

Table 1 Datasets Information

Author & Publication Year	Study	Independent Variable				Sample Size	# of Trials per Condition	SPE Indices					Reliability	
		IV 1	IV 2	IV 3	IV 4			RT	ACC	d'	η	v	z	ICC
Hu et al. (2023)	1	Matching	Identity	Emotion Control , Neutral, Happy, Sad	Session 1-6	33	60	✓	✓	✓	✓	✓	✓	✓
Constable and Knoblich (2020)	1	Matching	Identity	Switch Identity Partner, Stranger	Phase 1-2	92	40	✓	✓	✓	✓	✓	✓	✓
	2	Matching	Identity Self; Stranger	—	—	56	24	✓	✓	✓	✓	✓	✓	✓
Qian et al. (2020)	1	Matching	Stranger Identity Self; Stranger; Celebrity	Cue With, Without	—	25	24	✓	✓	✓	✓	✓	✓	✓
Schäfer and Frings (2019)	2	Matching	Identity Self; Celebrity	Cue With, Without	—	32	50	✓	✓	✓	✓	✓	✓	✓
	1	Matching	Celebrity Self; Mother, Acquaintance	—	—	35	24	✓	✓	✓	✓	✓	✓	✓
Golubickis and Macrae (2021)	1	Matching	Identity	Presentation Mixed ; Blocked	—	30	30	✓	✓	✓	✓	✓	✓	✓
Navon and Makovski (2021)	1	Matching	Identity	—	—	13	60	✓	✓	✓	✓	✓	✓	✓
	3	Matching	Identity Self; Father; Stranger	—	—	28	60	✓	✓	✓	✓	✓	✓	✓
Svensson et al. (2022)	4	Matching	Stranger Identity Self; Friend	—	—	27	60	✓	✓	✓	✓	✓	✓	✓
	1	Matching	Identity Self; Friend	—	—	20	50	✓	✓	✓	✓	✓	✓	✓
Xu et al. (2021)	2	Matching	Identity Self; Friend	Frequency Self ≥ Friend	—	24	100	✓	✓	✓	✓	✓	✓	✓
	3	Matching	Identity Self; Friend	Frequency Self < Friend	—	25	100	✓	✓	✓	✓	✓	✓	✓
	1	Matching	Identity	Tasks Modified; Unmodified	—	105	60	✓	✓	✓	✓	✓	✓	✓
Woźniak et al. (2018)	1	Matching	Identity	Facial Gender Male; Female	—	18	56	✓	✓	✓	✓	✓	✓	✓
Liu et al. (2023)	2	Matching	Identity	Facial Gender Male; Female	—	18	60	✓	✓	✓	✓	✓	✓	✓
	1	Matching	Identity Self; Stranger	—	—	298	16	✓	✓	✓	✓	✓	✓	✓

Note: Study represents different studies from a single article; IV: independent variable. For IV3 and IV4, we only included the baseline conditions that are similar to the original design in Sui et al. (2012), which were highlighted in **BOLD** font. If other variables that could be counterbalanced are indicated by underscores, we will solely utilize these variables as stratification variables during the split-half process

2.4 Analysis

Analysis plans for this study were preregistered on OSF (<https://osf.io/zv628>). All analyses in this paper were performed using the statistical software R (R Core Team, 2021). The drift rate (v) and starting point (z) of the Drift-Diffusion Model (DDM) was obtained using the “RWiener” package (Wabersich & Vandekerckhove, 2014).

The road map of the current study can be found in Fig. 2 and will be further elucidated in the subsequent sections.

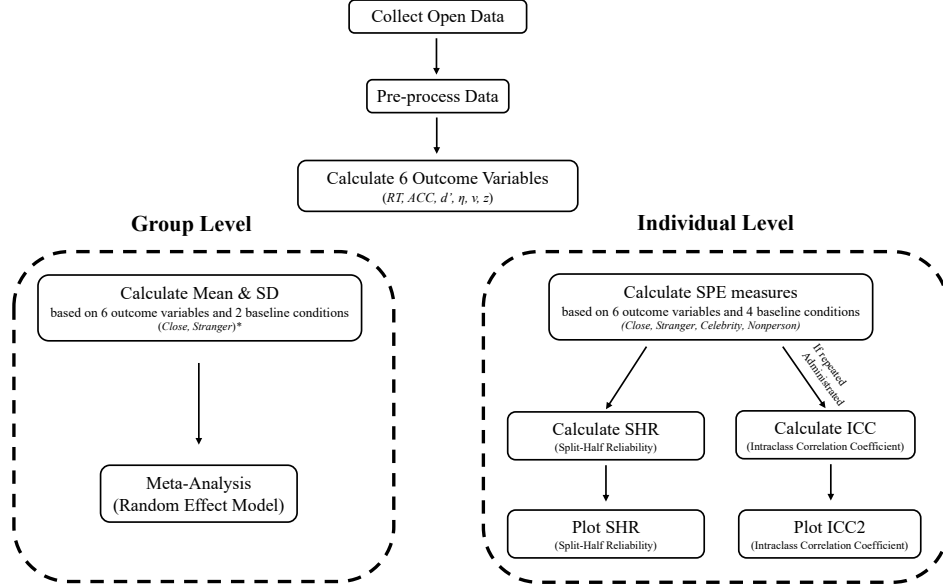


Fig. 2 Roadmap of the Current Study. *Note:* Only one paper have Celebrity and Nonpersons baseline, thus no included in the meta-analysis

2.4.1 Data Pre-processing

For all the seventeen datasets (see Table. 1), we applied the following exclusion criteria for excluding data:

1. Participant Exclusion Criteria
 - (i) Participants who had wrong trial numbers because of procedure errors is excluded from the analysis,
 - (ii) participants with an overall accuracy < 0.5 is excluded from the analysis,
 - (iii) participants with any of the conditions with zero accuracy is excluded from the analysis.
2. Trial Level Data Exclusion Criteria

- (i) Trials where the keypress occurs outside the two required keys and non-responsive trials are excluded from the analysis,
- (ii) the practice trials are excluded,
- (iii) the experimental design involved independent variables more than self-referential and matching (e.g., included valence of emotion as a third independent variable).

2.4.2 Calculating the Indicators and SPE Measures

We created a “multiverse” of SPE Measures. Specifically, for each study, we first calculated six indicators for each experimental condition: Reaction Time (RT), Accuracy (ACC), Sensitivity Score (d'), Efficiency (η), Drift Rate (v), and Starting Point (z). Reaction Time and Accuracy were obtained directly from the datasets, while sensitivity score was calculated based on choices; Efficiency was calculated based on Reaction Time and Accuracy; Drift Rate (v) and Starting Point (z) were estimated using standard DDM with Reaction Time and choice data. . As the self-prioritization effect was more often found in matching conditions (i.e., the difference between self-matching condition and other-matching condition), our calculation only focused on matching conditions, except for the d' , the calculation of which involved both matching and non-matching conditions. The SPE Measures were then computed using different indicators under available baseline conditions in the studies (see Table. 2).

Table 2 Indicators and SPE Measures Calculation

Indicators	Indicators Calculation		SPE Measures Calculation		Example Literature
Reaction Time (RT)	Total Reaction Time	/Total Responses	$RT_{\text{self-matching}} - RT_{\text{other-matching}}$		Humphreys and Sui (2015), Sui et al. (2012), and Sui and Humphreys (2017)
Accuracy (ACC)	# of Correct Responses	/Total Responses	$ACC_{\text{self-matching}} - ACC_{\text{other-matching}}$		Constable, Elekes, et al. (2019), Enock et al. (2018), and Sui et al. (2012)
d' -prime (d')	$Z(\text{Hits}) - Z(\text{False Alarms})$		$d'_{\text{self-matching}} - d'_{\text{other-matching}}$		Hu et al. (2020) and Sui et al. (2012)
Efficiency (η)	RT / ACC		$\eta_{\text{self}} - \eta_{\text{other}}$		Humphreys and Sui (2015) and Stoerber and Eysenck (2008)
Drift Rate (v)	Decomposed from RT and choice based on standard DDM		$v_{\text{self-matching}} - v_{\text{other-matching}}$		Golubickis et al. (2020) and Golubickis et al. (2017)
Starting Point (z)			$z_{\text{self-matching}} - z_{\text{other-matching}}$		Macrae et al. (2017) and Reuther and Chakravarthi (2017)

Note: $Z(\cdot)$ denotes the calculation of Z-score. In this context, "hit" refers to the ACC in matching trials, while "false alarm" refers to the error rate (1-ACC) in mismatch trials; the condition "Other" vary across contrast, we calculated the SPE for each "Other" condition. These could be the differences for "Self vs Close other", "Self vs Stranger", "Self vs Celebrity" or "Self vs Non-person".

2.4.3 Estimating the Robustness of SPE

The robustness of experimental effects (group-level effect) of SPE in SPMT was calculated using a meta-analytical approach. We employed a random effects model, given the anticipated heterogeneity among participant samples (Page et al., 2021). The effect size index used for all outcome measures was Hedges' g , a correction of Cohen's d that accounts for bias in small sample sizes (Hedges & Olkin, 1985). Hedges' g represents the magnitude of the difference between the self and baseline condition.

When calculating Hedges' g , we have reversed scored the effect size for variables with negative values (Reaction Time and Efficiency). Conversely, for all indicators, a positive effect size indicates a bias towards associating stimuli with the self as compared to baseline associations. For the estimation and interpretation of effect sizes, an effect size around 0.2 was interpreted as a small effect size, around 0.5 as a medium effect size, and around 0.8 as a large effect size (Fritz et al., 2012; Hedges & Olkin, 1985).

2.4.4 Estimating the Reliability of SPE

Split-half reliability. We assessed the split-half reliability by first splitting the trial-level data into two halves and calculating the Pearson correlation coefficients (r). To ensure methodological rigor, we used three approaches for splitting the trial-level data: first-second, odd-even and permuted (Kahveci et al., 2022; Pronk et al., 2022). The first-second approach split trials into the first half and the second half. The odd-even approach split the trials into sequences based on their odd or even numbers. The permuted approach shuffled the trial order and randomly assigned trials to two halves, iterating the process multiple times (usually thousands of times) to calculate the average and 95% confidence intervals of the split-half reliability.

In our analyses, we first stratified the trial-level data for each participant in the study based on experimental conditions. For example, in the case of a 2 by 3 within-subject design, we stratified the data based on the two independent variables: matching (matching, non-matching) and identity (self, stranger, friend). Subsequently, we applied the three splitting approaches (Pronk et al., 2022). When using the permuted approach, we randomly split the stratified data into two halves 5000 times, which resulted in 5000 pairs of two halves of the data. Next, we calculated 5000 Pearson correlation coefficients for these 5000 pairs. After that, we calculated the mean and 95% confidence intervals of the 5000 correlation coefficients. The first-second split and odd-even split only resulted in a single reliability coefficient. Finally, after computing the split-half reliability coefficients for each dataset, substantial variations were observed across the datasets.

To derive a more accurate estimation of the average split-half reliability for each SPE measure, we employed a synthesis approach for reliability coefficients using a minimum-variance unbiased aggregation method (Alexander, 1990; Olkin & Pratt, 1958). This approach corrects for the underestimation inherent in simply averaging correlations due to the specific distribution properties of correlation coefficients (Shieh,

2010). The method involves a correction and weighting of the reliability coefficients based on the number of participants. We calculated the weighted-average reliabilities using the “cormean” function within the “AATtools” Package (Kahveci, 2020). Although there is no strict criterion for defining the level of split-half reliability for psychological and educational measures, a widely accepted guideline for split-half reliability coefficient is that a value of 0.5 is “poor”, a value of 0.70 is “acceptable”, and a value greater than 0.8 means excellent reliability (Cicchetti & Sparrow, 1981).

Test-Retest Reliability (ICC). The Intraclass Correlation Coefficient (ICC) serves as a widely recognized measure for evaluating test-retest reliability (Fisher, 1992). Differing from the Pearson correlation coefficient, which primarily quantifies the linear association between two continuous variables, the ICC extends its prowess to scenarios involving multiple measurements taken on the same subjects, while also considering both the correlation and agreement between multiple measurements, making it a more comprehensive measure of test-retest reliability (Koo & Li, 2016). Since our primary aim was to evaluate the appropriateness of the SPMT in assessing individual differences and repeated administration, to achieve this objective, we assessed the test-retest reliability of the six indicators for our dataset that involved test-retest sessions using the function “ICC” in the “psych” package (Revelle, 2017). We focused on using the Two-way random effect model based on absolute agreement (ICC2) within the ICC family (Chen et al., 2018; Koo & Li, 2016; Xu et al., 2023). ICC2 gives an estimate of the proportion of total variance in measurements that is attributed to between-subjects variability (individual differences) and within-subjects variability (variability due to repeated measurements) (Xu et al., 2023). For the calculation of ICC2 estimates, the formula is:

$$ICC2 = \frac{MSBS - MSE}{MSBS + (k - 1)MSE + \left(\frac{k}{n}\right)(MSBM - MSE)}, \quad (1)$$

where $MSBS$ is the mean square between subjects, MSE is the mean square error, $MSBM$ is the mean square between measurements, k is the number of measurements, n is number of participants.

The traditional benchmarks for interpreting ICC values are as follows: ICC less than 0.50 suggests poor reliability; ICC between 0.50 and 0.75 suggests moderate reliability; ICC between 0.75 and 0.9 suggests good reliability; ICC above 0.9 suggests excellent reliability (Cicchetti & Sparrow, 1981; Kupper & Hafner, 1989).

3 Deviation from Preregistration

We adhered to our pre-registration plan as much as possible, however, there were a few differences between the current report and the pre-registration document. First, in our initial preregistration plan, we did not anticipate analyzing the group-level effect of SPE due to the perceived robustness of the effect across a diverse range of research. However, as our study progressed, we recognized the value of providing a more comprehensive assessment. Thus, we included an estimation of pooled effect sizes across the included study to represent the group-level effect. Second, we used a different algorithm for estimating the parameters of the drift-diffusion model. In the

preregistration, we planned to estimate the drift rate (v) and starting point (z) of the Drift-Diffusion Model using the “fit_ezddm” function from the “hausekeep” package (Lin et al., 2020). This function served as a wrapper for the EZ-DDM function (Wagenmakers et al., 2007). However, we observed limitations in the algorithm’s ability to accurately estimate parameter z during parameters recovery (details provided in the Supplementary Materials, section 1.2). After comparing the 5 algorithms, we found that the “RWiener” package (Wabersich & Vandekerckhove, 2014) achieved a favourable balance between accuracy, confidence interval and computational efficiency, making it the most suitable choice for our analysis. Nevertheless, for transparency, we have included the results from ezDDM in the supplementary materials (see Supplementary, Fig. S2-4). Third, we did not explicitly state in the preregistration report that we would perform a weighted average of the split-half reliabilities for all datasets. However, considering the significant impact of the number of trials on reliability (Kucina et al., 2023), during the formal analysis, we assigned different weights to each study based on the number of trials. Subsequently, we calculated a weighted average of the split-half reliabilities. Fourth, in our original preregistration, we outlined our intention to include both ICC2 and ICC2k in our data analysis. However, to obtain an overall estimate of the reliability, we weighted each study based on the number of participants. Fourth, in our original preregistration, we outlined our intention to include both ICC2 and ICC2k in our data analysis. However, as our understanding of Intraclass Correlation Coefficients (ICC) improved, we realized that ICC2 is the appropriate index for our research purpose. More specifically, ICC2k was mentioned in the preregistration as an index of robustness of group-level effect, but it turned out to be another index of reliability for individual differences. We corrected this misinterpretation of ICC2k in the final report. Fifth, we conducted exploratory analysis using the data we collected to investigate the relationship between the number of trials, permuted split-half reliability, and effect size (Hedges’ g) (refer to Supplementary Fig. S8-10). In addition, as suggested by one reviewer, we used the Spearman-Brown prediction formula based on our current data to predict the trial counts at which the SMT achieves sufficient reliability (Pronk et al., 2023). Sixth, the writing of the current manuscript was improved based on the pre-registration. For example, in our preregistration, we included different baseline conditions when calculating SPE in the method section but did not mention this in our introduction and abstract. Finally, we had incorrectly labelled the permutation method as Monte Carlo in the first version of the preprint. Thus, we corrected the misuse of the phrase in the updated version. Additionally, upon a thorough examination of the Monte-Carlo approach, we identified that its utilization could inflate reliability due to its psychometric properties (Kahveci et al., 2022). Consequently, we did not include this method in our analysis.

4 Results

Of the 42 independent datasets, 34 of them contain data for “Close other”, 34 of them contain data for “Stranger”, 1 of them has data for “Celebrity”, and 4 of them have data for “Nonperson”. Since there were only a few datasets for “Celebrity” and “Nonperson”, their results were presented in the supplementary materials.

4.1 Group Level Effect of SPE

We conducted a meta-analytical assessment to examine the robustness of SPE as measured by SPMT. We used a random effect model to synthesize the effect across different studies, with Hedges’ g as the index of effect size. We found that all measures of SPE, except the parameter z estimated from DDM, exhibited moderate to large effect sizes (see Table. 3 for numeric results for all six SPE measures, Fig. 3 for forest plots of effect sizes for RT). Our findings indicated a robust and substantial experimental effect of SPE. The I^2 value, all being greater than 75%, indicates high heterogeneity among studies, justifying the selection of the random effect model (Borenstein et al., 2021). The results for “Celebrity” and “None” as baselines were included in the supplementary materials (see Supplementary Table. S1).

Table 3 Meta-analytical Results of SPE Measures in SPMT

Baseline	Indicators	Hedges’ g [95%CI]	# of Studies	Q	p	I^2
Close other	RT	0.47 [0.30, 0.63]	14	68.67	< .001	84.94%
	ACC	0.73 [0.42, 1.03]	14	144.57	< .001	92.87%
	d'	0.44 [0.28, 0.59]	14	81.96	< .001	83.02%
	η	0.88 [0.50, 1.25]	14	128.47	< .001	94.67%
	v	0.54 [0.32, 0.76]	14	142.79	< .001	91.16%
	z	0.15[−0.03, 0.33]	14	122.30	0.11	88.95%
Stranger	RT	0.59 [0.40, 0.78]	13	55.30	< .001	83.20%
	ACC	0.78 [0.48, 1.08]	13	77.78	< .001	88.60%
	d'	0.35 [0.21, 0.50]	13	47.81	< .001	75.38%
	η	0.92 [0.56, 1.29]	13	98.79	< .001	93.30%
	v	0.44 [0.28, 0.59]	13	50.98	< .001	79.33%
	z	0.08[−0.09, 0.24]	13	70.48	0.37	84.44%

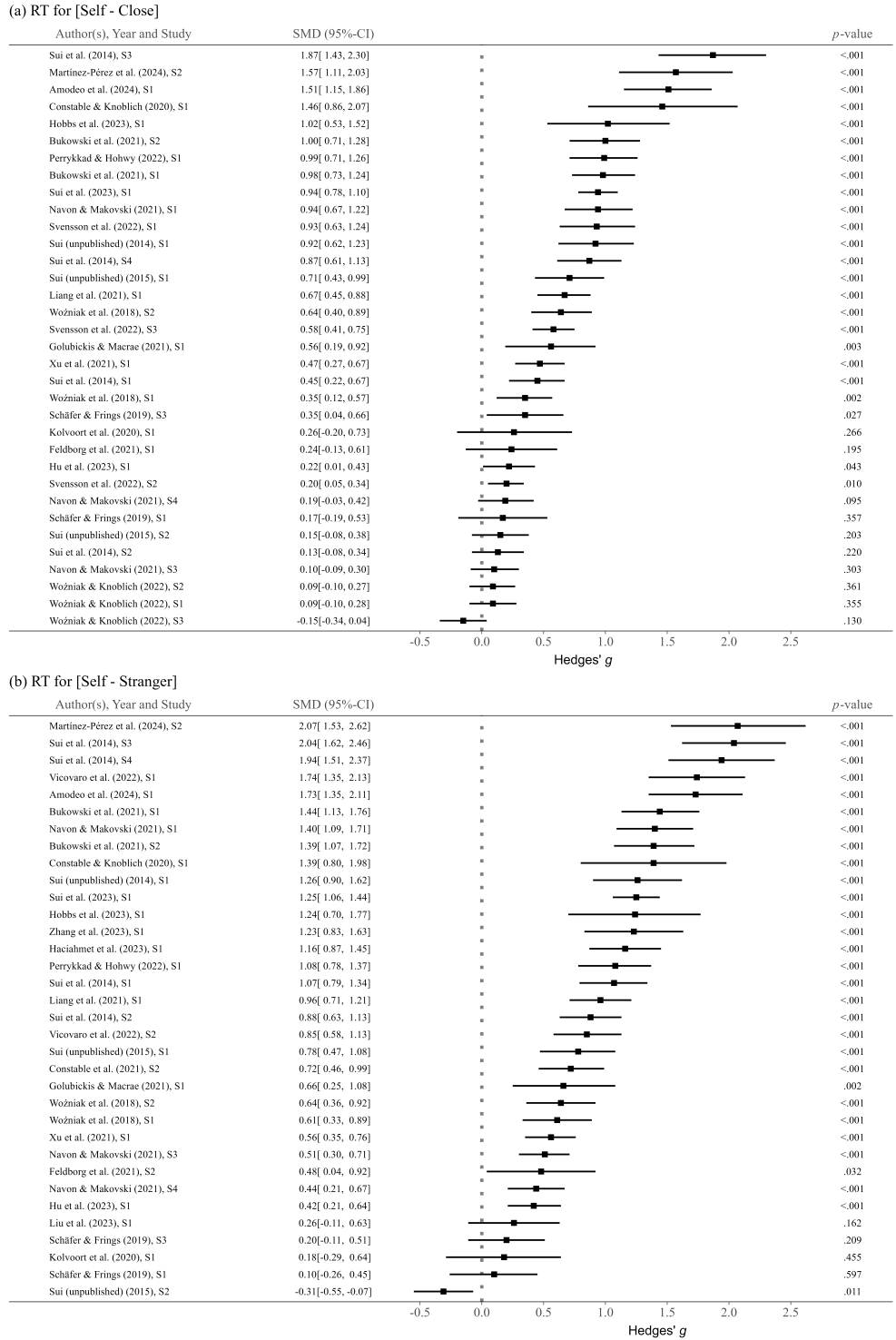


Fig. 3 Forest Plots for Group-level Self-Prioritization Effect (SPE) as Quantified by RT. (a) When “Close other” as the baseline condition for SPE, i.e., the “Self vs. Close other” contrast; (b) When “Stranger” as the baseline condition for SPE, i.e., the “Self vs. Stranger” contrast.

4.2 Split-half Reliability

We used three different approaches to split the data when calculating split-half reliability: the first-second, odd-even and permuted methods. Also, we used the weighted average split-half reliability as the overall reliability across studies. Here we only presented the results from the permuted split-half method both for clarity and for the robustness of this approach (Pronk et al., 2022) (see Fig. 4(a)). The results of the other two split-half methods can be found in the supplementary materials (see Supplementary Fig. S5).

We found that, among all SPE measures, the four with highest split-half reliabilities were as follows: Reaction Time (RT) with “Stranger” as baseline Reaction Time (RT) with “Close other” as baseline ($r = .55$, 95%CI [.38, .70]); Efficiency (η) with “Close other” as baseline ($r = .56$, 95% CI [.33, .74]); η with “Stranger” as baseline ($r = .54$, 95%CI [.33, .70]); RT with “Stranger” as baseline ($r = .49$, 95%CI [.35, .62]). These SPE measures achieved a split-half reliability of around 0.6 or higher, which is considered acceptable. For all other SPE measures, the reliability was around 0.5 or lower, indicating poor reliability. These included Accuracy (ACC), Sensitivity Score (d'), Drift Rate (v), and Starting Point (z) under four baselines. It’s worth noting that split-half reliability of z , the starting point parameter estimated from DDM, for all baselines was around 0, implying a lack of reliability. Since there is only one paper for “Celebrity” and four for “Nonperson”, their results are presented in the supplementary materials.

4.3 Test-retest Reliability

ICC could only be calculated for the dataset from our laboratory (Hu et al., 2023), which has 2 baseline conditions, the “Close other” and “Stranger”, in the experimental design. The ICC2, which measures the reliability for individual differences, aligns with the findings observed in split-half reliability estimation (see Fig. 4(b)). Specifically, when using “Close other” as baseline, the ICC2 for SPE measured by RT was .53 (95% CI [.39, .69]), and for Efficiency, it was .52 (95% CI [.38, .68]). Meanwhile, when “Stranger” was used as baseline, the ICC2 for RT was .58 (95% CI [.45, .73]), and for Efficiency, it was .35 (95% CI [.21, .52]). All other measures of SPE exhibited reliability lower than 0.5. To test the robustness of the results, we explored one additional dataset that included a re-test session but deviated strongly from the original SPMT, the result showed a similar pattern here (see Supplementary Fig. S6).

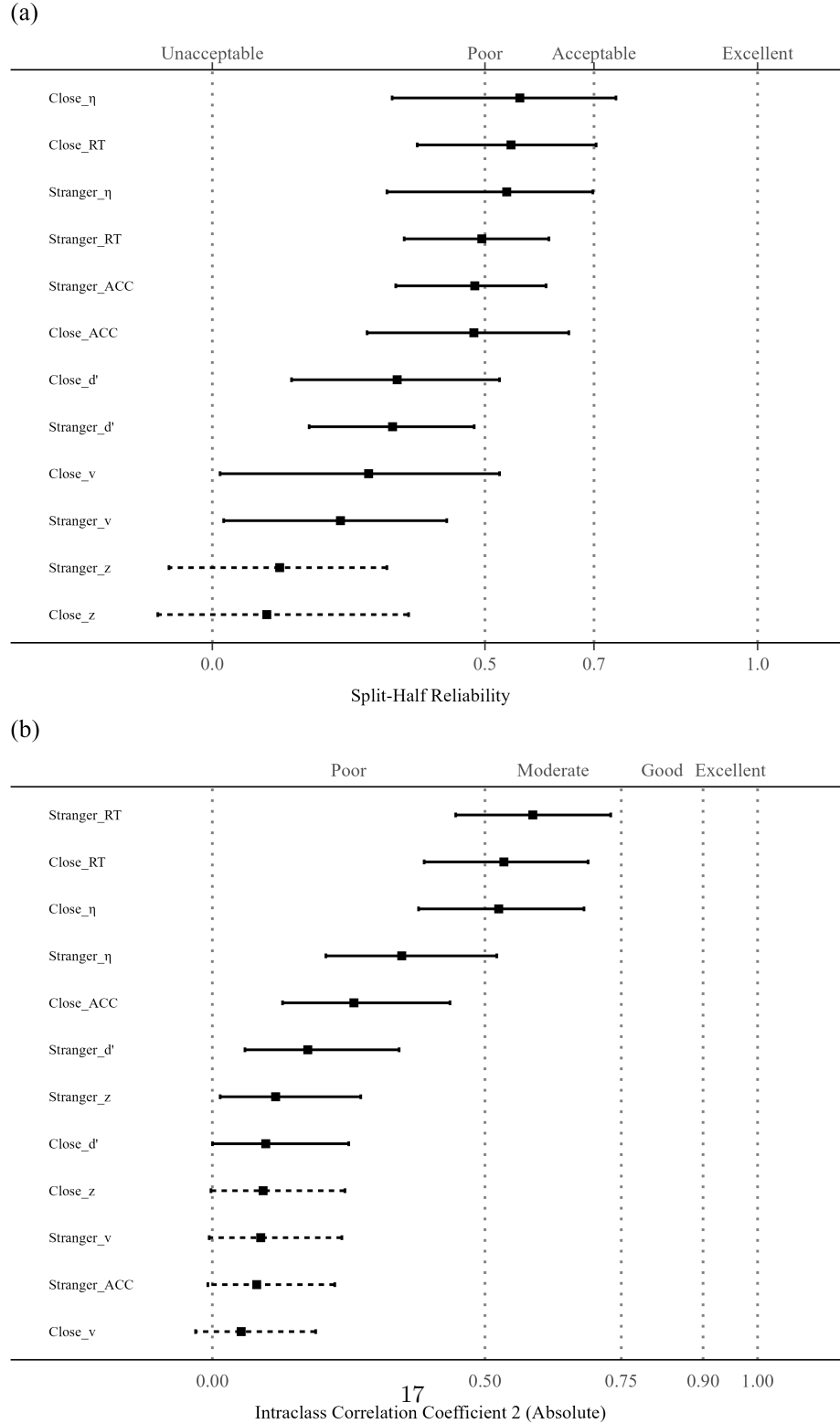


Fig. 4 Reliability for Different SPE Measures. (a) The Weighted Average Split-Half Reliability (Permutated); (b) Intraclass Correlation Coefficient (ICC2). *Note:* The vertical axis represents 12 different SPE measures, combining six indicators (RT, ACC, d' , η , v , z) and two baseline conditions (“Close other” and “Stranger”). The weighted average split-half reliability (figure a) and ICC values and their corresponding 95% confidence intervals are illustrated using points and lines. The dashed line indicates that the confidence interval for that point estimate extends across 0, implying a non-significant value. Due to the fact that there is only one paper for “Celebrity” and one for “Nonperson”, their results is presented in the supplementary materials.

5 Discussion

In this pre-registered study, we examined the reliability of various measures from the Self Matching Task (SMT) in assessing the self-prioritization effect (SPE). Our analyses revealed that except for parameters z from DDM, all the other measures exhibited robust SPE. However, when it came to reliability for individual differences, only two measures of SPE, Reaction Time and Efficiency, exhibited relatively higher but still unsatisfactory reliability, among all indicators that have been reported in the literature. Notably, this variability in reliability may not be taken into account for experimental parameters such as duration of stimulus presentation and response rules in the analyses - key parameters that influence cognitive task performance (Hedge et al., 2018). Our results revealed a “reliability paradox” for the SMT, when relying only on the matching task type and unfamiliar stimuli. These findings provided important methodological insights for using the SMT in assessing SPE at the individual level.

First, the Reaction Time (RT) and Efficiency (η) appeared to be the best measures among all the different ways to measure SPE (the other were ACC, d' , the parameter v and z from DDM). Our results revealed that the Reaction Time and Efficiency performed relatively well on both group level and individual levels. On group level, effect sizes of SPE as measured by Reaction Time and Efficiency were moderate to large effect; on individual level, SPE as measured by Reaction Time and Efficiency were higher for both split-half and test-retest reliability than other measures of SPE. These findings align with prior research (Draheim et al., 2016; M. M. Hughes et al., 2014), which also found greater within-session reliabilities for Reaction Time and accuracy composition compared to only incorporated accuracy. This is not surprising, as the difficulty of many cognitive tasks is low, making it more appropriate to focus on reaction time or a combination of reaction time and accuracy (e.g., efficiency). Similarly, the findings for the d' score are consistent with research on the reliability of other cognitive tasks (e.g., the matching task by Smithson et al. (2024); the recognition tasks by Franks and Hicks (2016)). It has been proposed that d' is heavily influenced by task difficulty, the nature of the target, and attentional factors (Vermeiren & Cleere-mans, 2012). Therefore, researchers should consider these factors when using d' to study individual differences. In addition, for different baseline conditions used for calculating SPE in the literature, “Stranger” and “Close other” (e.g., friends, or mother) are the most commonly utilized. Notably, “Stranger” produced a slightly higher effect size for most of the six indicators and demonstrated greater reliability when it came to Reaction Time. These results aligns with the preliminary results of our ongoing meta-analysis, suggesting that the selection of a baseline could be a significant moderator of the SPE. Taken together, for researchers interested in balancing between the group-level SPE and reliability, using Reaction Time and Efficiency as the indicators might be a good choice.

Second, taking the group-level robustness and individual-level results together, our findings revealed a “reliability paradox” in SPMT. We observed that the majority of the SPE measures demonstrated moderate to large effect sizes when analyzed at the group level. However, when considering individual differences, only the SPE measures derived from RT and Efficiency displayed comparatively higher values than other SPE measures but still did not meet the criteria for satisfactory split-half reliability.

Likewise, when examining the reliability across multiple time points using ICC2, RT and Efficiency still ranked the highest but only showed moderate levels of test-retest reliability. Our finding also aligned with the “reliability paradox” of cognitive tasks discovered in previous studies (Enkavi et al., 2019; Hedge et al., 2018). The precise causes behind the reliability paradox observed in SPE measurements using the SPMT warrant thorough investigation. However, one of the most plausible explanations is that the SPMT, like other cognitive tasks, tends to exhibit minimal variability among participants while maximizing the detection of SPE at the group level (Liljequist et al., 2019). Alternatively, the current finding indicates that when assessing reliability for individual differences, it is essential to consider critical experimental parameters such as stimulus presentation, response rules, stimulus onset asynchrony/inter-trial interval, and the number of practice trials. These parameters enable a fine-grained design for individual differences in SPE using the SMT. The current study sheds light on the specific types of inquiries on how to proficiently use the SMT to address both group and individual-level differences in SPE. More specifically, at the group level, the interpretations of the results remain largely consistent, even without taking into account experimental parameters such as varying response rules. However, the relatively low reliability of all the SPE measures in the current analysis without considering these design parameters calls for attention when researchers are interested in individual-level analyses, such as in clinical settings or searching for an association with data from questionnaires (e.g., Hobbs et al., 2023; Moseley et al., 2022). Nonetheless, the reliability results of reaction time (RT) measures remain generally higher, particularly in existing studies focusing on individual-level differences ((e.g., Liu et al., 2022; Zhang et al., 2023)). Future research needs to exercise greater caution and follow the standard practice to maximize reliability at the individual level in their results (Parsons et al., 2019).

Many studies have previously used the SMT to assess robust group-level SPE, more recent studies showed a burgeoning interest in using the SMT to quantify individual variability in SPE. Several approaches have recently been proposed to enhance the reliability of cognitive tasks, which may prove valuable for the SMT. These include using gamification (Friebs et al., 2020), latent model (Eisenberg et al., 2019; Enkavi et al., 2019) or generative models (Haines et al., 2020) to analyze the data. Some of these suggestions have already been validated by empirical data. For example, Kucina et al. (2023) re-designed the cognitive conflict task by incorporating more trials and gamification indeed improving the reliability compared to the traditional Stroop task alone. Our exploratory analyses of the relationship between trial numbers and reliability also suggest that increasing trial numbers may improve reliability (please refer to Supplementary section 2.4).

Finally, a surprising result is the notably low split-half and test-retest reliability observed in the parameters (v and z) derived from the drift-diffusion model. In our analyses, we applied common and easy-to-use methods to datasets, estimated parameters for each condition of each participant and then calculated the reliability. The reliability of both the drift rate (v) and the starting point (z) fell well below acceptable levels. These results contradict previous findings that drift rate (v) and starting point (z) can be used as an index of SPE. Several studies interpreted the drift rate

(v) as the index of the speed and quality of information acquisition and reported higher drift rate for self-relevant stimuli (e.g., Golubickis et al., 2020; Golubickis et al., 2017). However, the reliability of drift rate (v) is relatively low in our study. As for the starting point (z), studies also reported SPE using z and interpreted this effect as a preference for matching response when the stimuli are self-relevant (e.g., Macrae et al., 2017; Reuther & Chakravarthi, 2017). Our meta-analytical results indicated that the Hedges’ g for starting point (z) was around zero. The split-half reliability of z was also small, possibly because z fails to adequately reflect the SPE. These findings raised concerns about applying the standard drift-diffusion model to data from SPMT directly. Previous studies also found that the standard drift-diffusion model did not fit the data from matching task (Groulx et al., 2020). Additionally, the reliability of parameters derived from other cognitive models, such as reinforcement learning models (Eckstein et al., 2022), has also been found to be unsatisfactory. These findings called for a more principled approach when modelling behavioral data to more accurately capture the fundamental cognitive processes at play (e.g., Wilson & Collins, 2019), instead of applying the standard DDM blindly.

5.1 Implications of the Current Study

Our findings can offer an initial guide for researchers considering the use of SMT. Firstly, we recommend that researchers employ Reaction time and Efficiency as the indicators of SPE since they strike a balance between achieving a substantial effect size at the group level and ensuring reliability at the individual level. Second, if researchers are interested in a relatively bigger group-level effect size, using the “Self vs Stranger” contrast may prove beneficial. Third, if feasible, increase the number of trials, as it may enhance the overall reliability of the measurements. We used the Spearman-Brown prediction formula (Pronk et al., 2023) to predict the trial numbers required for different levels of reliability. The results indicated that the number of trials required for archive sufficient reliability (e.g., 0.8) varied across different SPE indices. For SPE measured by RT, approximately 180 trials are required to achieve a reliability of 0.8 (see Fig S11 for more caveats). Lastly, we caution against the careless application of the standard drift-diffusion model and instead advocate for a principled modelling approach.

5.2 Limitations

Several limitations warrant acknowledgment. Firstly, although we made efforts to enhance sample diversity by including open data when available, it is important to note that the majority of our samples still consisted of individuals from what is commonly referred to as “(W)EIRD” populations (Rad et al., 2018; Yue et al., 2023), most of the participants were recruited from universities and are healthy adults. As a result, our findings may not be fully representative of the broader population, and it is necessary to include a more diverse sample to ensure greater generalizability of the paradigm. Secondly, our results reported here assessed the robustness and reliability of SPE with the original experimental design of Sui et al. (2012), which means the

robustness and reliability of different variants of SPMT still need further investigation. For a more systematic meta-analysis of SPE measured by SPMT, please see our ongoing project (<https://osf.io/euqmf>). While this focused analysis in a small set of papers from a large pool using the SMT enabled a deeper understanding individual-level reliability of the SPE, we recognize that expanding the scope and criteria to include more papers could potentially bolster the generalizability of our findings. This implies that further investigation is necessary to assess the robustness and reliability of other variations of the SMT, as well as other tasks used to measure SPE. This is particularly crucial given findings suggesting that different cognitive measures of self-biases may exhibit considerable independence from one another (Nijhof et al., 2020). Thirdly, when assessing the intraclass correlation coefficients (ICC2), only one dataset had available data from multiple tests, which could potentially limit the representativeness of the results. This issue is mitigated by the fact that additional analysis of one dataset (see supplementary section 2.3) with different designs showed similar results as we reported in the main text.

6 Conclusion

This study provided the first empirical assessment of the reliability of the self matching task (SMT) for measuring individual differences in SPE. We found a robust self-prioritization effect for all measures of SPE, except the starting point parameter z estimated from DDM. Meanwhile, the reliability of all the SPE measures (Reaction Time, Accuracy, Efficiency, sensitivity score, drift rate and starting point) fell short of being satisfactory. The results of the current study may serve as a benchmark for the improvement of individual-level reliability using this paradigm.

Acknowledgments

The data collection from Hu et al. (2023) was supported by the National Science Foundation (China, Grant No. 31371017) to JS. The author wishes to express gratitude to Dr. Sercan Kahveci for his valuable feedback on the first version of the preprint.

Author Contributions

HCP: Conceptualization, Methodology, Investigation, Resources, Writing - Original Draft, Writing - Review & Editing, Project administration, Supervision. LZ: Methodology, Data Curation, Software, Formal analysis, Visualization, Investigation, Writing - Original Draft. HMZ: Methodology, Data Curation, Software, Formal analysis, Visualization, Investigation, Writing - Original Draft. ZYR: Software. SJ: Funding acquisition, Data Curation, Writing - Review & Editing.

Data and Material Availability

The pre-registration plan is available at OSF(<https://osf.io/zv628>). The de-identified raw data from our lab is available at Science Data Bank (<https://doi.org/10.57760/>

sciedb.08117). The simulated data is accessible on GitHub (<https://github.com/Chuan-Peng-Lab/ReliabilitySPE>).

Code Availability

Code used to simulate and analyze the data is made accessible on GitHub (<https://github.com/Chuan-Peng-Lab/ReliabilitySPE>).

Competing Interests

The authors declare no competing interests.

References

- Alexander, R. A. (1990). A note on averaging correlations. *Bulletin of the Psychonomic Society*, 28(4), 335–336.
- Borenstein, M., Hedges, L. V., Higgins, J. P., & Rothstein, H. R. (2021). *Introduction to meta-analysis*. John Wiley & Sons.
- Borsboom, D. (2005). *Measuring the mind: Conceptual issues in contemporary psychometrics*. Cambridge University Press.
- Chen, G., Taylor, P. A., Haller, S. P., Kircanski, K., Stoddard, J., Pine, D. S., Leibenluft, E., Brotman, M. A., & Cox, R. W. (2018). Intraclass correlation: Improved modeling approaches and applications for neuroimaging. *Human Brain Mapping*, 39(3), 1187–1206. <https://doi.org/10.1002/hbm.23909>
- Cicchetti, D. V., & Sparrow, S. A. (1981). Developing criteria for establishing interrater reliability of specific items: Applications to assessment of adaptive behavior. *Am J Ment Defic*, 86(2), 127–137. <https://psycnet.apa.org/record/1982-00095-001>
- Clark, K., Birch-Hurst, K., Pennington, C. R., Petrie, A. C., Lee, J. T., & Hedge, C. (2022). Test-retest reliability for common tasks in vision science. *Journal of Vision*, 22(8), 18–18. <https://doi.org/10.1167/jov.22.8.18>
- Constable, M. D., Elekes, F., Sebanz, N., & Knoblich, G. (2019). Relevant for us? we-prioritization in cognitive processing. *Journal of Experimental Psychology: Human Perception and Performance*, 45(12). <https://doi.org/10.1037/xhp0000691>
- Constable, M. D., & Knoblich, G. (2020). Sticking together? re-binding previous other-associated stimuli interferes with self-verification but not partner-verification. *Acta Psychologica*, 210, 103167. <https://doi.org/10.1016/j.actpsy.2020.103167>
- Constable, M. D., Rajsic, J., Welsh, T. N., & Pratt, J. (2019). It is not in the details: Self-related shapes are rapidly classified but their features are not better remembered. *Memory & Cognition*, 47, 1145–1157. <https://doi.org/10.3758/s13421-019-00924-6>
- Constable, M. D., Becker, M. L., Oh, Y.-I., & Knoblich, G. (2021). Affective compatibility with the self modulates the self-prioritisation effect. *Cognition and Emotion*, 35(2), 291–304. <https://doi.org/10.1080/02699931.2020.1839383>

- Cunningham, S. J., Turk, D. J., Macdonald, L. M., & Macrae, C. N. (2008). Yours or mine? ownership and memory. *Consciousness and Cognition*, 17(1), 312–318. <https://doi.org/10.1016/j.concog.2007.04.003>
- Draheim, C., Hicks, K. L., & Engle, R. W. (2016). Combining reaction time and accuracy: The relationship between working memory capacity and task switching as a case example. *Perspectives on Psychological Science*, 11(1), 133–155.
- Eckstein, M. K., Master, S. L., Xia, L., Dahl, R. E., Wilbrecht, L., & Collins, A. G. (2022). The interpretation of computational model parameters depends on the context. *Elife*, 11, e75474.
- Eisenberg, I. W., Bissett, P. G., Zeynep Enkavi, A., Li, J., MacKinnon, D. P., Marsch, L. A., & Poldrack, R. A. (2019). Uncovering the structure of self-regulation through data-driven ontology discovery. *Nature Communications*, 10(1), 2319. <https://doi.org/10.1038/s41467-019-10301-1>
- Enkavi, A. Z., Eisenberg, I. W., Bissett, P. G., Mazza, G. L., MacKinnon, D. P., Marsch, L. A., & Poldrack, R. A. (2019). Large-scale analysis of test–retest reliabilities of self-regulation measures. *Proceedings of the National Academy of Sciences*, 116(12), 5472–5477. <https://doi.org/10.1073/pnas.1818430116>
- Enock, F. E., Sui, J., Hewstone, M., & Humphreys, G. W. (2018). Self and team prioritisation effects in perceptual matching: Evidence for a shared representation. *Acta Psychologica*, 182, 107–118. <https://doi.org/10.1016/j.actpsy.2017.11.011>
- Fisher, R. A. (1992). Statistical methods for research workers. *Springer New York*. https://doi.org/10.1007/978-1-4612-4380-9_6
- Franks, B. A., & Hicks, J. L. (2016). The reliability of criterion shifting in recognition memory is task dependent. *Memory & Cognition*, 44, 1215–1227.
- Friebs, M. A., Dechant, M., Vedress, S., Frings, C., & Mandryk, R. L. (2020). Effective gamification of the stop-signal task: Two controlled laboratory experiments. *JMIR Serious Games*, 8(3), e17810. <https://doi.org/10.2196/17810>
- Fritz, C. O., Morris, P. E., & Richler, J. J. (2012). Effect size estimates: Current use, calculations, and interpretation. *Journal of Experimental Psychology: General*, 141(1), 2. <https://doi.org/10.1037/a0024338>
- Golubickis, M., Falbén, J. K., Ho, N. S., Sui, J., Cunningham, W. A., & Macrae, C. N. (2020). Parts of me: Identity-relevance moderates self-prioritization. *Consciousness and Cognition*, 77, 102848. <https://doi.org/10.1016/j.concog.2019.102848>
- Golubickis, M., Falbén, J. K., Sahraie, A., Visokomogilski, A., Cunningham, W. A., Sui, J., & Macrae, C. N. (2017). Self-prioritization and perceptual matching: The effects of temporal construal. *Memory & Cognition*, 45, 1223–1239. <https://doi.org/10.3758/s13421-017-0722-3>
- Golubickis, M., & Macrae, C. N. (2021). Judging me and you: Task design modulates self-prioritization. *Acta Psychologica*, 218, 103350. <https://doi.org/10.1016/j.actpsy.2021.103350>
- Green, S. B., Yang, Y., Alt, M., Brinkley, S., Gray, S., Hogan, T., & Cowan, N. (2016). Use of internal consistency coefficients for estimating reliability of

- experimental task scores. *Psychonomic Bulletin & Review*, 23, 750–763. <https://doi.org/10.3758/s13423-015-0968-3>
- Groulx, J. T., Harding, B., & Cousineau, D. (2020). The ez diffusion model: An overview with derivation, software, and an application to the same-different task. *The Quantitative Methods for Psychology*, 16(2), 154–174. <https://doi.org/10.20982/tqmp.16.2.p154>
- Haines, N., Kvam, P. D., Irving, L. H., Smith, C., Beauchaine, T. P., Pitt, M. A., Ahn, W.-Y., & Turner, B. M. (2020). Theoretically informed generative models can advance the psychological and brain sciences: Lessons from the reliability paradox. <https://doi.org/10.31234/osf.io/xr7y3>
- Hedge, C., Powell, G., & Sumner, P. (2018). The reliability paradox: Why robust cognitive tasks do not produce reliable individual differences. *Behavior Research Methods*, 50, 1166–1186. <https://doi.org/10.3758/s13428-017-0935-1>
- Hedges, L., & Olkin, I. (1985). *Statistical methods for meta-analysis*, academic press.
- Hobbs, C., Sui, J., Kessler, D., Munafò, M. R., & Button, K. S. (2023). Self-processing in relation to emotion and reward processing in depression. *Psychological Medicine*, 53(5), 1924–1936. <https://doi.org/10.1017/S0033291721003597>
- Hu, C.-P., Peng, K., & Sui, J. (2023). Data for training effect of self prioritization[ds/ol]. v2. *Science Data Bank*. <https://doi.org/10.57760/sciencedb.08117>
- Hu, C.-P., Lan, Y., Macrae, C. N., & Sui, J. (2020). Good me bad me: Prioritization of the good-self during perceptual decision-making. *Collabra. Psychology*, 6(1), 20. <https://doi.org/10.1525/collabra.301>
- Hughes, M. M., Linck, J. A., Bowles, A. R., Koeth, J. T., & Bunting, M. F. (2014). Alternatives to switch-cost scoring in the task-switching paradigm: Their reliability and increased validity. *Behavior research methods*, 46, 702–721.
- Hughes, S. M., & Harrison, M. A. (2013). I like my voice better: Self-enhancement bias in perceptions of voice attractiveness. *Perception*, 42(9), 941–949. <https://doi.org/10.1068/p7526>
- Humphreys, G. W., & Sui, J. (2015). The salient self: Social saliency effects based on self-bias. *Journal of Cognitive Psychology*, 27(2), 129–140. <https://doi.org/10.1080/20445911.2014.996156>
- Kahveci, S. (2020). Aattools: Reliability and scoring routines for the approach-avoidance task.
- Kahveci, S., Bathke, A., & Blechert, J. (2022). Reliability of reaction time tasks: How should it be computed? <https://doi.org/10.31234/osf.io/ta59r>
- Karvelis, P., Paulus, M. P., & Diaconescu, A. O. (2023). Individual differences in computational psychiatry: A review of current challenges. *Neuroscience & Biobehavioral Reviews*, 105137. <https://doi.org/10.1016/j.neubiorev.2023.105137>
- Keenan, J. P., Wheeler, M. A., Gallup, G. G., & Pascual-Leone, A. (2000). Self-recognition and the right prefrontal cortex. *Trends in Cognitive Sciences*, 4(9), 338–344. [https://doi.org/10.1016/S1364-6613\(00\)01521-7](https://doi.org/10.1016/S1364-6613(00)01521-7)
- Kircher, T. T., Senior, C., Phillips, M. L., Benson, P. J., Bullmore, E. T., Brammer, M., Simmons, A., Williams, S. C., Bartels, M., & David, A. S. (2000). Towards

- a functional neuroanatomy of self processing: Effects of faces and words. *Cognitive Brain Research*, 10(1-2), 133–144. [https://doi.org/10.1016/S0926-6410\(00\)00036-7](https://doi.org/10.1016/S0926-6410(00)00036-7)
- Kline, P. (2015). *A handbook of test construction (psychology revivals): Introduction to psychometric design*. Routledge.
- 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>
- Kucina, T., Wells, L., Lewis, I., de Salas, K., Kohl, A., Palmer, M. A., Sauer, J. D., Matzke, D., Aidman, E., & Heathcote, A. (2023). Calibration of cognitive tests to address the reliability paradox for decision-conflict tasks. *Nature Communications*, 14(1), 2234. <https://doi.org/10.1038/s41467-023-37777-2>
- Kupper, L. L., & Hafner, K. b. (1989). On assessing interrater agreement for multiple attribute responses. *Biometrics*, 45(3), 957–967. <https://doi.org/10.2307/2531695>
- 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>
- Lin, H., Saunders, B., Friese, M., Evans, N. J., & Inzlicht, M. (2020). Strong effort manipulations reduce response caution: A preregistered reinvention of the ego-depletion paradigm. *Psychological Science*, 31(5), 531–547. <https://doi.org/10.1177/0956797620904990>
- Liu, Song, Y., Lee, N. A., Bennett, D. M., Button, K. S., Greenshaw, A., Cao, B., & Sui, J. (2022). Depression screening using a non-verbal self-association task: A machine-learning based pilot study. *Journal of Affective Disorders*, 310, 87–95. <https://doi.org/10.1016/j.jad.2022.04.122>
- Liu, Sui, J., & Hildebrandt, A. (2023). To see or not to see: The parallel processing of self-relevance and facial expressions. *Manuscript submitted for publication*.
- Logie, R. H., Sala, S. D., Laiacona, M., Chalmers, P., & Wynn, V. (1996). Group aggregates and individual reliability: The case of verbal short-term memory. *Memory & Cognition*, 24, 305–321. <https://doi.org/10.3758/BF03213295>
- Macrae, C. N., Visokomogilski, A., Golubickis, M., Cunningham, W. A., & Sahraie, A. (2017). Self-relevance prioritizes access to visual awareness. *Journal of Experimental Psychology: Human Perception and Performance*, 43(3), 438. <https://doi.org/10.1037/xhp0000361>
- Makel, M. C., Plucker, J. A., & Hegarty, B. (2012). Replications in psychology research: How often do they really occur? *Perspectives on Psychological Science*, 7(6), 537–542. <https://doi.org/10.1177/1745691612460688>
- Moseley, R., Liu, C. H., Gregory, N., Smith, P., Baron-Cohen, S., & Sui, J. (2022). Levels of self-representation and their sociocognitive correlates in late-diagnosed autistic adults. *Journal of Autism and Developmental Disorders*, 52(7), 3246–3259.
- Navon, M., & Makovski, T. (2021). Are self-related items unique? the self-prioritization effect revisited. <https://doi.org/10.31234/osf.io/9dzm4>

- Nijhof, A. D., Shapiro, K. L., Catmur, C., & Bird, G. (2020). No evidence for a common self-bias across cognitive domains. *Cognition*, 197, 104186.
- Olkin, I., & Pratt, J. W. (1958). Unbiased estimation of certain correlation coefficients. *The annals of mathematical statistics*, 201–211.
- Orellana-Corrales, G., Matschke, C., & Wesslein, A.-K. (2020). Does self-associating a geometric shape immediately cause attentional prioritization? comparing familiar versus recently self-associated stimuli in the dot-probe task. *Experimental Psychology*, 67(6), 335. <https://doi.org/10.1027/1618-3169/a000502>
- Page, M. J., McKenzie, J. E., Bossuyt, P. M., Boutron, I., Hoffmann, T. C., Mulrow, C. D., Shamseer, L., Tetzlaff, J. M., Akl, E. A., Brennan, S. E., et al. (2021). The prisma 2020 statement: An updated guideline for reporting systematic reviews. *International Journal of Surgery*, 88, 105906. <https://doi.org/10.1136/bmj.n71>
- 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>
- Payne, B., Lavan, N., Knight, S., & McGettigan, C. (2021). Perceptual prioritization of self-associated voices. *British Journal of Psychology*, 112(3), 585–610. <https://doi.org/10.1111/bjop.12479>
- Pronk, T., Hirst, R. J., Wiers, R. W., & Murre, J. M. (2023). Can we measure individual differences in cognitive measures reliably via smartphones? a comparison of the flanker effect across device types and samples. *Behavior Research Methods*, 55(4), 1641–1652.
- Pronk, T., Molenaar, D., Wiers, R. W., & Murre, J. (2022). Methods to split cognitive task data for estimating split-half reliability: A comprehensive review and systematic assessment. *Psychonomic Bulletin & Review*, 29(1), 44–54. <https://doi.org/10.3758/s13423-021-01948-3>
- Qian, H., Wang, Z., Li, C., & Gao, X. (2020). Prioritised self-referential processing is modulated by emotional arousal. *Quarterly Journal of Experimental Psychology*, 73(5), 688–697. <https://doi.org/10.1177/1747021819892158>
- R Core Team. (2021). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. Vienna, Austria. <https://www.R-project.org/>
- Rad, M. S., Martingano, A. J., & Ginges, J. (2018). Toward a psychology of homo sapiens: Making psychological science more representative of the human population. *Proceedings of the National Academy of Sciences*, 115(45), 11401–11405. <https://doi.org/10.1073/pnas.1721165115>
- Reuther, J., & Chakravarthi, R. (2017). Does self-prioritization affect perceptual processes? *Visual Cognition*, 25(1-3), 381–398. <https://doi.org/10.1080/13506285.2017.1323813>
- Revelle, W. R. (2017). Psych: Procedures for personality and psychological research. <https://CRAN.R-project.org/package=psych>

- Rogers, T. B., Kuiper, N. A., & Kirker, W. S. (1977). Self-reference and the encoding of personal information. *J Pers Soc Psychol*, 35(9), 677–88. <https://doi.org/10.1037//0022-3514.35.9.677>
- Rouder, J. N., & Haaf, J. M. (2019). A psychometrics of individual differences in experimental tasks. *Psychonomic Bulletin & Review*, 26(2), 452–467. <https://doi.org/10.3758/s13423-018-1558-y>
- Schäfer, S., & Frings, C. (2019). Understanding self-prioritisation: The prioritisation of self-relevant stimuli and its relation to the individual self-esteem. *Journal of Cognitive Psychology*, 31(8), 813–824. <https://doi.org/10.1080/20445911.2019.1686393>
- Scheller, M., & Sui, J. (2022a). The power of the self: Anchoring information processing across contexts. *Journal of Experimental Psychology: Human Perception and Performance*, 48(9), 1001.
- Scheller, M., & Sui, J. (2022b). Social relevance modulates multisensory integration. *Journal of Experimental Psychology: Human Perception and Performance*, 48(9), 1022.
- Shieh, G. (2010). Estimation of the simple correlation coefficient. *Behavior Research Methods*, 42(4), 906–917.
- Smithson, C. J., Chow, J. K., Chang, T.-Y., & Gauthier, I. (2024). Measuring object recognition ability: Reliability, validity, and the aggregate z-score approach. *Behavior Research Methods*, 1–15.
- Stoeber, J., & Eysenck, M. W. (2008). Perfectionism and efficiency: Accuracy, response bias, and invested time in proof-reading performance. *Journal of Research in Personality*, 42(6), 1673–1678. <https://doi.org/10.1016/j.jrp.2008.08.001>
- Sui, J., He, X., & Humphreys, G. W. (2012). Perceptual effects of social salience: Evidence from self-prioritization effects on perceptual matching. *Journal of Experimental Psychology: Human Perception and Performance*, 38(5), 1105–1117. <https://doi.org/10.1037/a0029792>
- Sui, J., & Humphreys, G. W. (2017). The self survives extinction: Self-association biases attention in patients with visual extinction. *Cortex*, 95, 248–256. <https://doi.org/10.1016/j.cortex.2017.08.006>
- Svensson, S. L., Golubickis, M., Maclean, H., Falbén, J. K., Persson, L. M., Tsamadi, D., Caughey, S., Sahraie, A., & Macrae, C. N. (2022). More or less of me and you: Self-relevance augments the effects of item probability on stimulus prioritization. *Psychological Research*, 86(4), 1145–1164. <https://doi.org/10.1007/s00426-021-01562-x>
- Turk, D. J., Heatherton, T. F., Kelley, W. M., Funnell, M. G., Gazzaniga, M. S., & Macrae, C. N. (2002). Mike or me? self-recognition in a split-brain patient. *Nature Neuroscience*, 5(9), 841–842. <https://doi.org/10.1038/nn907>
- Vermeiren, A., & Cleeremans, A. (2012). The validity of d measures. *PloS one*, 7(2), e31595.
- Wabersich, D., & Vandekerckhove, J. (2014). The rwienr package: An r package providing distribution functions for the wiener diffusion model. *R Journal*, 6(1). <https://doi.org/10.32614/RJ-2014-005>

- Wagenmakers, E.-J., Van Der Maas, H. L., & Grasman, R. P. (2007). An ez-diffusion model for response time and accuracy. *Psychonomic Bulletin & Review*, 14(1), 3–22. <https://doi.org/10.3758/BF03194023>
- Wilson, R. C., & Collins, A. G. (2019). Ten simple rules for the computational modeling of behavioral data. *Elife*, 8, e49547. <https://doi.org/10.7554/eLife.49547>
- Woźniak, M., Kourtis, D., & Knoblich, G. (2018). Prioritization of arbitrary faces associated to self: An eeg study. *PloS One*, 13(1), e0190679. <https://doi.org/10.1371/journal.pone.0190679>
- Xu, Kiar, G., Cho, J. W., Bridgeford, E. W., Nikolaidis, A., Vogelstein, J. T., & Milham, M. P. (2023). Rex: An integrative tool for quantifying and optimizing measurement reliability for the study of individual differences. *Nature Methods*, 1–4. <https://doi.org/10.1038/s41592-023-01901-3>
- Xu, Yuan, Y., Xie, X., Tan, H., & Guan, L. (2021). Romantic feedbacks influence self-relevant processing: The moderating effects of sex difference and facial attractiveness. *Current Psychology*, 1–13. <https://doi.org/10.1007/s12144-021-02114-7>
- Yankouskaya, A., Lovett, G., & Sui, J. (2023). The relationship between self, value-based reward, and emotion prioritisation effects. *Quarterly Journal of Experimental Psychology*, 76(4), 942–960.
- Yue, L., Zuo, X.-N., & Chuan-Peng, H. (2023). The weird problem in a “non-weird” context: A meta-research on the representativeness of human subjects in chinese psychological research. <https://doi.org/osf.io/y9hwq>
- Zhang, Y., Wang, F., & Sui, J. (2023). Decoding individual differences in self-prioritization from the resting-state functional connectome. *NeuroImage*, 120205. <https://doi.org/10.1016/j.neuroimage.2023.120205>
- Zorowitz, S., & Niv, Y. (2023). Improving the reliability of cognitive task measures: A narrative review. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*. <https://doi.org/10.1016/j.bpsc.2023.02.004>