

One and only SNARC? How Flexible are The Flexibility of Spatial-Numerical Associations?

A Registered Report on the SNARC's Range Dependency

Lilly Roth¹, ~~Gáspár Lukács²~~, John Caffier², Ulf-Dietrich Reips², Hans-Christoph Nuerk^{1,4,5*},
Krzysztof Cipora³

¹Department of Psychology, University of Tübingen, Germany

²Department of Psychology, University of Konstanz, Germany

³Centre for Mathematical Cognition, Loughborough University, United Kingdom

⁴LEAD Graduate School & Research Network, University of Tübingen, Germany

⁵[German Center for Mental Health \(DZPG\)](#)

*Corresponding author:

hc.nuerk@uni-tuebingen.de

Abstract

Numbers are associated with space, but it is unclear how flexible these associations are. In this study, we will investigate whether the SNARC effect (Spatial-Numerical Association of Response Codes; Dehaene et al., 1993), which describes faster responses to small/large number magnitude with the left/right hand, respectively, is fully flexible (and depends only on relative ~~number~~ magnitude within a stimulus set), or not (and depends on absolute ~~number~~ magnitude as well). Evidence for relative-magnitude dependency comes from studies observing that numbers 4 and 5 were associated with the right when presented in a 0 – 5 range but with the left in a 4 – 9 range (Dehaene et al., 1993; Fias et al., 1996). However, this important conclusion was drawn solely from the absence of evidence for absolute-magnitude dependency in frequentist analysis in underpowered studies. A closer inspection of ~~that~~ those descriptive data suggests ~~that descriptively~~ absolute magnitude might also matter. Hence, we will conduct a close replication of Dehaene et al.'s (1993) Experiment 3 and a conceptual replication considering recent advances in SNARC research, investigating absolute- and relative-magnitude dependency by comparing intercepts and SNARC slopes across ranges with Bayesian statistics. ~~Besides traditional null hypothesis testing, we will also report Bayes Factors.~~ To achieve a power of .90 for detecting moderate evidence (Bayes Factor above 3 or below 1/3) for a smallest effect size of interest of Cohen's $d = 0.1520$ ~~at a significance level of $\alpha = .01$~~ , we will conduct each experiment online with maximum 376800 participants, but run sequential analyses with optional stopping at moderate evidence. We hypothesize that both absolute and relative magnitude influence spatial-numerical associations, suggesting ~~that~~ the SNARC effect operates on flexible and absolute number representations simultaneously.

Keywords: spatial-numerical associations, SNARC effect, mental number line, replication, online experiment, high statistical power

One and only SNARC? How Flexible are The Flexibility of Spatial Numerical Associations.?

A Registered Report on the SNARC's Range Dependency

Numbers are highly relevant in everyday life. Therefore, much research has been devoted to understanding how we process and represent them in our minds. Interestingly, various aspects of numerical information such as cardinality and ordinality are systematically associated with different aspects of space such as extensions or directions (Cipora et al., 2020; Cipora, [Schroeder](#) et al., 2018; Patro et al., 2014). This broad range of phenomena is referred to under the umbrella term Spatial-Numerical Associations, SNAs (Fischer & Shaki, 2014; Toomarian & Hubbard, 2018). Investigating these associations is fundamental for models of number representation and – considering the bigger picture – for models of human cognition.

The hallmark directional SNA is the Spatial-Numerical Association of Response Codes (SNARC) effect, which denotes that in left-to-right reading cultures, participants respond faster to small/large magnitude numbers on the left/right side, respectively (Dehaene et al., 1993). Interestingly, the SNARC effect can be observed in a parity judgment task, in which the magnitude of the numbers is not task-relevant. This effect has been replicated using different modalities, setups and tasks (see Cipora et al., 2019, for an online replication; Fias et al., 1996; Toomarian & Hubbard, 2018, for a recent review; Wood et al., 2008, for a meta-analysis). The SNARC effect is typically quantified using the repeated-measures regression originally proposed by Lorch and Myers (1990) and applied to the SNARC effect by Fias et al. (1996). In the first step mean differences in reaction times (RTs) between the right and left hand (dRTs) are regressed on numerical magnitude for each participant separately. A negative slope indicates an increasing right-hand advantage with increasing number magnitude (the more negative the so-called SNARC slope, the stronger the SNARC effect). Subsequently, to check for the SNARC effect at the group level, individual SNARC slopes are tested against zero with a one-sample *t*-test.

Interestingly, several studies have documented that the SNARC effect is not fixed but might be prone to several types of manipulation (Cipora, [Patro](#), & [Nuerk-et-al.](#), 2018, for a taxonomy), for instance, changing the number range of the used stimuli, which has been classified as representational, intra-experimental manipulation. The spatial mental number representation seems to be adapted to fit the task at hand. In this work we focus on the extent to which the SNARC effect flexibly adjusts to the specific range of the numbers being used in the task set.

Relative-magnitude dependency of the SNARC effect

The seminal paper by Dehaene et al. (1993) has already demonstrated in Experiment 3 that the SNARC effect depends on the relative rather than the absolute magnitude of numbers. They found the SNARC effect in two different numerical intervals ranging from 0 to 5 and from 4 to 9. In the lower interval, responses to numbers 4 and 5 were faster with the right hand (typical response pattern for large numbers), whereas within the higher interval, responses to these numbers were faster with the left hand (typical response pattern for small numbers). This finding was replicated by Fias et al. (1996, Experiment 1). It suggests that the SNARC effect dynamically adapts to the current task set (i.e., numbers being used) and is determined by the relative magnitude of the number within the set rather than its absolute magnitude. We refer to this claim about the SNARC effect as relative-magnitude dependency (RMdependency).

The RMdependency is considered as one of the crucial features of the SNARC effect and is taken for granted since these early findings. The results of Dehaene et al.'s (1993) and Fias et al.'s (1996) experiments are widely cited as an argument for the SNARC being dependent on the given number range (e.g., by Antoine & Gevers, 2016; Deng et al., 2016; Ginsburg et al., 2014; Ginsburg & Gevers, 2015; Schwarz & Keus, 2004; Pinhas et al., 2013). The RMdependency of the SNARC effect has been demonstrated by several other studies even going beyond a basic setup comprising judgments on single digit numbers. For instance, Tlauka (2002) found a SNARC effect both when using the two numbers 1 and 100 and when using the

two numbers 100 and 900. The number 100 was associated to the right/left when it was the larger/smaller of the two numbers, respectively. Ben Nathan et al. (2009) go even further, showing that the SNARC effect is not only RMdependent on the task level but built up on a trial-to-trial basis. They found the right- and left-key response speed advantages in magnitude judgment tasks to depend on the relative magnitude in comparison to the ever-changing reference number. What is more, evidence for RMdependency of the SNARC-like effects goes beyond numerical stimuli. Wühr and Richter (2022) found a SNARC-like effect (association of physically smaller/larger stimuli with the left/right, respectively) to depend on relative rather than absolute stimulus size.

Importantly, RMdependency has also been used as a methodological tool to show that a spatial-numerical phenomenon is in fact the SNARC effect. For instance, Rugani et al. (2015), Di Giorgio et al. (2019), and Giurfa et al. (2022) demonstrated the RMdependency to claim that a certain effect they observed in newly hatched chickens, in newborn children, and in honeybees is of the same nature as the SNARC effect. To sum up, there is evidence for the RMdependency of the SNARC effect in various tasks and setups, and it has even been used to validate SNAs.

RMdependency in the light of number-representation models

RMdependency fits well with most theoretical accounts of number representation. The seminal work of Restle (1970) outlining the Mental Number Line (MNL) account, which has been proposed as the first explanation for the SNARC effect (Dehaene et al., 1993), postulates that the MNL is flexible and dynamically adapts to the task demands. In line with this, Pinhas et al. (2013) claim that the resolution of the MNL can be adjusted to the numerical context. The accounts of verbal-spatial coding (Gevers et al., 2010) and polarity correspondence (Proctor & Cho, 2006) are on the one hand in line with RMdependency, but on the other hand they do not make clear statements about relative magnitude being the *only* decisive factor determining the SNARC effect. Crucially, both accounts assume that long-term number representations underlie the SNARC effect, which hardly justifies the SNARC effect's flexibility (Ginsburg & Gevers,

2015; van Dijck et al., 2015). The working memory account (Fias & van Dijck, 2016; van Dijck & Fias, 2011) originally claimed that the SNARC effect does not rely on long-term number representations, but is instead constructed during task execution, which speaks in favor of pure RMdependency. However, Ginsburg et al. (2014) argue that short-term number representations do not always fully overrule long-term number representations. This idea has been incorporated in the hybrid account proposed by van Dijck et al. (2015) as well, and it allows the coexistence of RMdependency and dependency of the SNARC effect on absolute number magnitude (henceforth AMdependency). Furthermore, concurrent RMdependency and AMdependency would also be in line with the idea that multiple number representations and multiple spatial reference frames can be activated and operated simultaneously (Weis et al., 2018). To conclude, the assumption that absolute magnitude plays no role can hardly be derived from theoretical accounts of the SNARC effect.

Hints towards AMdependency of the SNARC effect

In addition to the prominent claims on the RMdependency of the SNARC effect, the literature also provides hints towards an AMdependency of the SNARC effect. It is important to note that AMdependency can, on the one hand, influence the strength of the SNARC effect (reflected by the SNARC slope), and on the other, the location of numbers on the MNL in absolute terms (reflected by the intercept of the regression line and by dRTs of critical numbers that are part of both number ranges). Crucially, the SNARC effect seemed to be stronger in the lower than in the higher number range in both initial studies demonstrating the RMdependency (-20.1 ms vs. -10.9 ms in Dehaene et al., 1993; and -10.18 ms vs. -7.19 ms in Fias et al., 1996), suggesting AMdependency as well. In Fias et al.'s (1996) results, the observed slope difference had approximately an effect size of Cohen's $d = 0.16$ (i.e., the slope difference of 2.99 divided by the pooled standard deviation of 18.34 ms, which has been calculated with $SD = 15.1$ ms and $SD = 11.2$ ms for the lower and higher number ranges, assuming a rather conservative correlation between them of $r = 0.05$, which corresponds to the correlation we have observed

in our previous color judgment tasks). Moreover, the results pointed towards an overall shift of small/large numbers to the left/right on the MNL, respectively, since the smallest-number intercept (i.e., the predicted dRT for the smallest number magnitude of the range, which was 0/4 in the lower/higher range, respectively) was larger in the lower than in the higher range (37.52 ms vs. 14.03 ms in Dehaene et al., 1993; and 15.43 ms vs. 8.82 ms in Fias et al., 1996). However, the mean-number intercepts (i.e., the predicted dRT for the mean number magnitude of the range, which was 2.5/6.5 in the lower/higher range, respectively) did not differ much in Fias et al.'s results (-10.02 ms vs. -9.16 ms), ~~leading to a result pattern as illustrated in Scenario 4 in Figure 1~~. In Dehaene et al.'s results, this intercept seemed to be smaller in the higher number range, but it cannot be calculated exactly based on the data reported in the paper. ~~In that study, the observed result pattern looked like Scenario 5 in Figure 1~~.

Methodological limitations of the two initial studies demonstrating RMdependency

Even if we use the two original studies as a guidance for further investigations, their findings are not very reliable because of several important limitations regarding the design and the interpretation of the results. Both Dehaene et al. (1993) and Fias et al. (1996) found a significant two-way interaction of response side (left vs. right) and magnitude (small vs. medium vs. large), ~~indicating the SNARC effect¹~~. Apart from the repeated-measures regression approach, the SNARC effect can also be quantified as a two-way interaction of response side and magnitude (for methodological considerations, see Fias et al., 1996) or as linear contrast in an ANOVA (Tzelgov et al., 2013). However, the three-way interaction of response side and magnitude with interval (0 to 5 vs. 4 to 9) remained non-significant in both studies. In Fias et al.'s (1996) additional repeated-measures regression the resulting SNARC slopes differed significantly from zero in both intervals in a one-sample *t*-test, and the difference in SNARC

¹ Apart from the repeated measure regression approach below, one can also quantify the SNARC as a two-way interaction of number magnitude × response side (see Fias et al., 1996 for methodological considerations) or as linear contrast in an ANOVA (Tzelgov et al., 2013).

slopes between both intervals remained non-significant in a t -test for two dependent samples. Crucially, the strong conclusion of pure RMdependency that has been derived from these null results is dangerously close to mistaking absence of evidence for evidence of absence. Importantly, no Bayesian analysis was conducted to test whether the null results supported the null hypothesis (and it is not possible to run a post-hoc Bayesian analysis due to the lacking report of the exact t -statistic). What is more, neither Dehaene et al. (1993) nor Fias et al. (1996) tested whether the dRT pattern for the same number differed significantly between number ranges – even if the right-hand advantage (reflected by negative dRTs) for numbers 4 and 5 in the range from 0 to 5 and the left-hand advantage (reflected by positive dRTs) for these numbers in the range from 4 to 9 are often cited. Also, the smallest-number intercepts and the mean-number intercepts were not compared between ranges.

Moreover, the design was most likely underpowered for the relevant statistical comparisons in both studies (see below for calculations). On the one hand, this was due to the relatively low sample sizes ($n = 12$ in Dehaene et al., 1993; and $n = 24$ in Fias et al., 1996). On the other, only 15 repetitions per experimental cell (i.e., per number magnitude and response-key assignment) were used. Later methodological studies proposed to use at least 20 repetitions and 20 participants to detect the SNARC effect, and even more repetitions and participants to detect differences in the size of the SNARC effect (Cipora & Wood, 2017). Following the *effect-size sensitivity approach* (Giner-Sorolla et al., 2020), we have run power calculations to determine SNARC slope differences between the two number ranges that are detectable in a t -test for two dependent samples at different adequate power levels (adapting Monte-Carlo simulations by Wickelmaier, 2022 using the R package *pwr* by Champely et al., 2018). For the sample size used by Fias et al. (1996) and with the standard deviations they observed, our calculations revealed that at power levels of .80, .90, and .95, only SNARC slope differences between the two number ranges of minimum 11.0 ms ($d = 0.60$), 12.7 ms ($d = 0.69$) and 14.1 ms ($d = 0.77$) could have been detected, respectively. Note that we ran these calculations within

the frequentist framework, which corresponds to the data analysis by Fias et al. (for power calculations in both the frequentist and the Bayesian framework, see <https://doi.org/10.17605/OSF.IO/Z43PM>~~Roth, Lukács, et al., 2022~~, created using the R packages *rmarkdown* by Allaire et al., 2022;~~;~~ and *knitr* by Xie, 2022; and *BayesFactor* by Morey et al., 2015). However, such differences in SNARC slopes are very unlikely, even in case of AMdependency, because they would be larger than the typically observed SNARC slopes themselves. Because of the lack of related information in Dehaene et al.'s (1993) paper, we were not able to run such power calculations for their results; but because their sample was even smaller, they could have detected only even larger differences.

Moreover, the stimuli used in both studies (0, 1, 2, 3, 4, 5 and 4, 5, 6, 7, 8, 9) lead to two problems. First, the average number magnitude in both number ranges is larger for odd than for even numbers (3 vs. 2 in the lower and 7 vs. 6 in the higher number range). This can lead to a confound with the MARC (Linguistic Markedness of Response Codes) effect that denotes a left/right-hand advantage when responding to odd/even numbers, respectively (Nuerk et al., 2004). Such a confound may decrease the SNARC effect (Tzelgov et al., 2013; Zohar-Shai et al., 2017). The association of small/large numbers to the left/right side, respectively, should be weaker if small/large numbers are more often even/odd, respectively. More recent studies have addressed this issue by using stimuli sets in which number magnitude and contrast-coded parity are orthogonal (e.g., Cipora et al., 2019). Typically, it is done by using the number set 1, 2, 3, 4, 6, 7, 8, 9, which importantly also excludes zero (see below).

Second, using the number zero is problematic due to its special status shown in several studies: Reading time for zero is significantly longer than for any other single digit number and is not predicted by factors determining reading time of other single digit numbers (Brysbaert, 1995). Nuerk et al. (2004) and Nieder (2016) provide further empirical evidence that zero may not be represented on the MNL along with other numbers (but see Pinhas & Tzelgov, 2012, for another conclusion). Additionally, quite often participants have problems understanding the parity

status of zero (Levenson et al., 2007). Using zero also turned out problematic in SNARC studies: The RTs and dRTs for the number zero do not strongly correlate with the RTs and dRTs of other even numbers (Nuerk et al., 2004). Later studies on the SNARC effect have excluded zero from the stimuli set (e.g., Cipora et al., 2019; Cleland & Bull, 2018; Deng et al., 2016; Gevers et al., 2010, Gökyaydin et al., 2018). Ultimately, both the parity status and the presence of zero might have confounded the results of the previous studies (see Table 2). Therefore, in addition to the replication that we will conduct as close as possible to the original studies by Dehaene et al. (1993) and Fias et al. (1996), we will also conduct a conceptual replication using a suitable stimulus set to disentangle these potential confounds and tackle all the above-mentioned limitations.

Can the SNARC effect operate on two reference frames at once?

As we laid out so far, there is a general tendency to interpret the SNARC effect as entirely flexible based on the findings of RMdependency and on the inference-statistical null effects concerning AMdependency (in underpowered studies). However, the SNARC effect could be operating concurrently in both relative and absolute terms. Indeed, one of us has proposed that the SNARC effect operates on multiple number lines in previous work (Weis et al., 2018). However, that paper is not about whether the SNARC effect operates on multiple number lines in terms of RMdependency and AMdependency, but instead it used two-digit numbers as stimuli to see whether separate number lines are activated for decade and unit numbers. The operations on different number ranges are for decade and unit digits of one two-digit number (i.e., the same number, but different digits of its decomposition). Thus, the paper by Weis et al. provides the principal account that the SNARC effect could operate on multiple reference frames at once. The current study goes beyond their findings because it seeks to demonstrate that both RMdependent and AMdependent spatial mappings are concurrently present in the same digit.

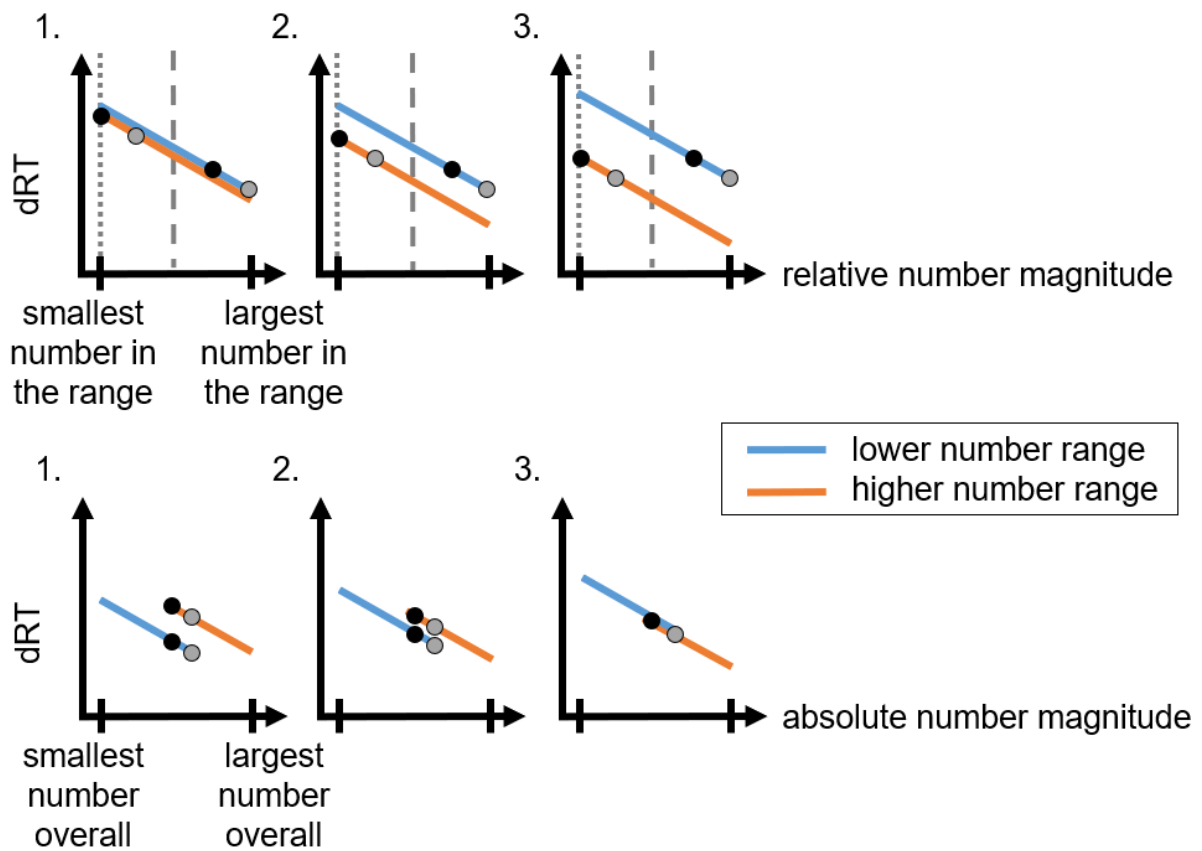
How could absolute magnitude affect the SNARC effect?

Apart from the regression slope that quantifies the strength of the SNARC effect, the smallest-number intercept (when relative magnitude of the numbers in both ranges is matched, i.e., the predicted dRT for 0 and 4 in Experiment 1 and for 1 and 4 in Experiment 2) and the mean-number intercept (i.e., the predicted dRT for 2.5 and 6.5 in Experiment 1 and for 3 and 6 in Experiment 2) can be determined in order to investigate the number mapping on the MNL. When discussing RMdependency and AMdependency of the SNARC effect, the following scenarios are possible (see Figures 1 and 2 and Table 1, for detailed elaboration of these ~~complex~~ scenarios):

1. RMdependency of the number mapping on the MNL, but no difference in the strength of the SNARC effect between number ranges (i.e., different dRTs of critical numbers that are part of both number ranges, namely 4 and 5, ~~but same SNARC slopes, same smallest-number intercepts and same mean-number intercepts~~)
2. Both RMdependency and AMdependency of the number mapping on the MNL, but no difference in the strength of the SNARC effect between number ranges (i.e., different dRTs of critical numbers, different smallest-number intercepts and different mean-number intercepts, ~~but same SNARC slopes~~)
3. AMdependency of the number mapping on the MNL, but no difference in the strength of the SNARC effect between number ranges (i.e., different smallest-number intercepts and different mean-number intercepts, ~~but same dRTs of critical numbers and same SNARC slopes~~) – note that concluding RMdependency of the number mapping on the MNL from finding a significant SNARC effect in both number ranges without testing dRTs of critical numbers is incorrect, ~~because the SNARC effect and the MNL do not flexibly adjust to the used number range in this scenario~~

Figure 1

Possible Scenarios of RMdependency and AMdependency of the number mapping on the MNL

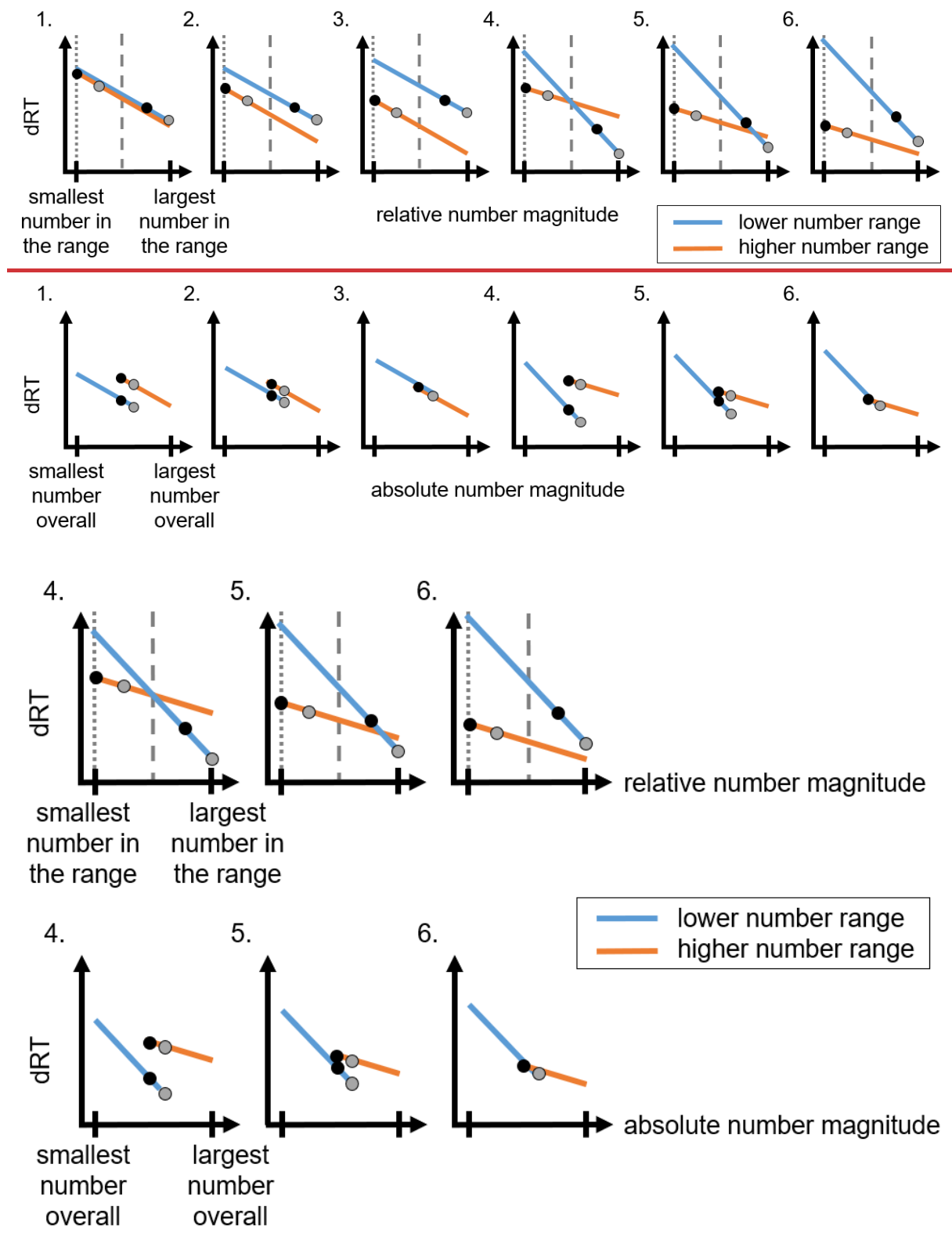


Note. This figure (retrieved from <https://doi.org/10.17605/OSF.IO/Z43PM>) illustrates Scenarios 1, 2, and 3, with the regression lines for the lower and higher number ranges being represented in blue and orange, respectively. In the upper part of the figure, relative number magnitudes are used for the x-axis, so that the regression lines for both number ranges start at their smallest and end at their largest number magnitude. For example, in Experiment 1, the dRTs for 0 (smallest number in the lower number range) and 4 (smallest number in the higher number range) are on the very left, and the dRTs for 5 (largest number in the lower number range) and 9 (largest number in the higher number range) are on the very right. In the lower part of the figure, the same scenarios are illustrated, but absolute number magnitudes are used for the x-axis. In our study, the absolute number magnitudes will be 0 to 5 and 4 to 9 in Experiment 1, and 1 to 5 (excluding 3) and 4 to 8 (excluding 6) in Experiment 2. For example, the dRTs for numbers 4 and 5 are on the very same spot of the x-axis for both the lower and the higher range, because they have the same absolute magnitude. The dotted line in the upper part of the figure depicts the intercept for the smallest number magnitude, and the dashed line depicts the intercept for the mean number magnitude in the respective number range. The black and the gray dots indicate the critical numbers being part of both the lower and the higher number range (i.e., 4 and 5).

1. AMdependency of the strength of the SNARC effect, and RMdependency of the number mapping on the MNL (i.e., different SNARC slopes, different dRTs of critical numbers, different smallest-number intercepts, ~~but same mean-number intercepts~~), as in Fias et al. (1996)
2. AMdependency of the strength of the SNARC effect, and both RMdependency and AMdependency of the number mapping on the MNL (i.e., different SNARC slopes, different dRTs of critical numbers, different smallest-number intercepts, and mean-number intercepts), as in Dehaene et al. (1993)
3. AMdependency of the strength of the SNARC effect and of the number mapping on the MNL (i.e., different SNARC slopes, different smallest-number intercepts, and different mean-number intercepts, ~~but same dRTs of critical numbers~~)

Figure 21

Possible Scenarios of RMdependency and AMdependency of the *strength of the* SNARC Effect



Note. This figure (retrieved from <https://doi.org/10.17605/OSF.IO/Z43PM>) illustrates the six possible described scenarios 4, 5, and 6. For an explanation of magnitudes, with the regression lines for the lower and higher number

ranges being represented in blue and orange, respectively. In the upper part of the figure, relative number magnitudes are used for the on the x-axis, so that the regression lines for both number ranges start at their smallest and end at their largest as well as concrete examples for data points, see *Note* of Figure 1. number magnitude. For example, in Experiment 1, the dRTs for 0 (smallest number in the lower number range) and 4 (smallest number in the higher number range) are on the very left, and the dRTs for 5 (largest number in the lower number range) and 9 (largest number in the higher number range) are on the very right. In the lower part of the figure, the same scenarios are illustrated, but absolute number magnitudes are used for the x axis, namely 0 to 5 and 4 to 9 in Experiment 1, and 1 to 5 (excluding 3) and 4 to 8 (excluding 6) in Experiment 2. For example, the dRTs for numbers 4 and 5 are on the very same spot of the x axis for both the lower and the higher range, because they have the same absolute magnitude. The dotted line in the upper part of the figure depicts the intercept for the smallest number magnitude, and the dashed line depicts the intercept for the mean number magnitude in the respective number range. The black and the gray dots indicate the critical numbers being part of both the lower and the higher number range (i.e., 4 and 5). From Roth, Lukács, et al. (2022).

Table 1

Possible Scenarios of RMdependency and AMdependency of the SNARC Effect

Scenario	1	2	3	4	5	6
Significant SNARC effect in both ranges	yes	yes	yes	yes	yes	yes
Different dRTs for critical numbers (4 and 5)	yes	yes	no	yes	yes	no
Different smallest-number intercept	no	yes	yes	yes	yes	yes
Different mean-number intercept	no	yes	yes	no	yes	yes
Different SNARC slopes	no	no	no	yes	yes	yes

Note. This table summarizes the characteristics of the six possible scenarios of RMdependency and AMdependency of the SNARC effect, which are described above and illustrated in Figures 1 and 2. The crucial distinction consists in whether dRTs, intercepts and slopes differ between the two ranges in both experiments.

The current study

In this study, we aim to answer the question whether the SNARC effect depends only on relative magnitude or whether absolute magnitude plays a role as well. First, we will replicate Experiment 3 by Dehaene et al. (1993), which has also been replicated in Experiment 1 by Fias et al. (1996), where we will also use the number ranges from 0 to 5 and from 4 to 9. Second, we will conduct a conceptual replication, which is meant to address confounds due to the unequal distribution of odd and even numbers and due to the presence of zero in both stimuli sets, where we will use the number ranges 1 to 5 (excluding 3) and 4 to 8 (excluding 6). The middle number of the range is also excluded in most SNARC studies using the typical set from 1 to 9. Moreover, the critical numbers that appear in both ranges are then the same in both experiments, namely 4 and 5. Table 2 gives an overview of the number ranges we will use and of confounds between number parity and number magnitude in Experiment 1 that will be avoided in Experiment 2.

In both of our replication experiments, a high statistical power will be obtained by testing much larger samples than Dehaene et al. (1993) and Fias et al. (1996) and by increasing the number of repetitions per experimental cell from 15 to 25. ~~To be able to Furthermore, analyses will be complemented by Bayes Factors so that we can~~ quantify evidence both for differences between number ranges and for lack of such differences, we will use the Bayesian instead of frequentist approach (for the interpretation of different values for the Bayes Factors, we will follow the recommendations by Dienes, 2021). Online experiments offer the possibility to collect data from large samples and therefore reach high statistical power (Reips, 2000, 2002). The SNARC effect has been successfully replicated in online settings (e.g., Cipora et al., 2019; Gökyaydin et al., 2018; Koch et al., 2021). The measurement in the online setup showed a similar reliability and magnitude compared to the SNARC effect that is typically observed in lab studies. Further, it seems to be valid as regards the correlations of the SNARC

effect with mean RTs and standard deviations of RTs, which are similar compared to lab studies.

In this study, we expect to replicate the findings by Dehaene et al. (1993) and by Fias et al. (1996) as concerns RMdependency. However, we also expect to find evidence towards AMdependency of the number mapping on the MNL and of the strength of the SNARC effect. Previous studies have indicated tendencies that cannot be explained by RMdependency alone. More precisely, we expect to observe Scenario 4 or 5 (see Figure 24 and Table 1) and hypothesize:

1. A SNARC effect in all used number ranges, replicating the results by Dehaene et al. (1993) and Fias et al. (1996).
2. Both RMdependency and AMdependency of the number mapping on the MNL, such that small/large numbers in relative absolute terms are shifted towards the left/right, respectively.
3. AMdependency of the strength of the SNARC effect, such that it is stronger in the lower than in the higher ranges.

Method

This study has been approved by the ethics committee of the University of Tübingen's Department of Psychology.

Statistical power considerations and sample size determination

In this study we decided to power for Cohen's $d = 0.150.20$ as the minimal effect size of interest in Hypothesis 3, because the most crucial aim of the present study is to find out whether AMdependency of the strength of the SNARC effect exists or not. By choosing this minimal effect size of interest, we will be able to find evidence for or against the SNARC slope differences between number ranges that were descriptively reported in the original studies that we wish to replicate here. Due to the lacking report of standard deviations, it is not possible to

calculate Cohen's d for the slope difference of 9.2 ms found by Dehaene et al. (1993), but the slope difference of 2.99 ms with a pooled standard deviation of 18.34 ms found by Fias et al. (1996) corresponds to an effect size of $d = 0.16$. Note that in the two original studies, the symmetric confidence intervals for these estimates must also include at least the double slope difference and effect size due to their non-significance. Hence, in case of AM dependency of the strength of the SNARC effect, the true effect size might in fact be larger than $d = 0.15$.

~~, because it corresponds to a small effect size (Cohen, 1988) and to 1% of explained variance (calculated according to Ruscio, 2008, using the conversion formula assuming equal-sized groups, see Table 2 there). Smaller effects that explain less variance are not practically meaningful. We decided not to choose the minimal effect size of interest based on effect sizes found in previous studies, because the true effect is often overestimated this way (because of type M errors resulting from too small samples, see Gelman and Carlin, 2014, and because of the publication bias).~~

To ensure a statistical power of .90 for the detection of finding moderate evidence (i.e., BF_{10} greater than 3, according to Dienes, 2021) for an effect size of Cohen's $d = 0.15$ ~~at a significance level of $\alpha = .01$~~ for t -tests, the sample needs to consist of at least 376800 participants (for power calculations, see <https://doi.org/10.17605/OSF.IO/Z43PM> ~~Roth, Lukács, et al., 2022~~). For this calculation, we used ~~the standard deviation~~ $SD = 15.1$ ms and $SD = 11.2$ ms for the lower and higher number ranges, as reported by Fias et al. (1996), although the standard deviation in our previous color judgment experiments were only $SD = 4.2$ ms and $SD = 3.9$ ms. ~~will probably be smaller because of a higher number of repetitions per experimental cell.~~ Hence, our calculations are rather conservative, and the statistical power thus is most probably even higher.

Importantly, as we run Bayesian instead of frequentist analyses, we will make use of the "Sequential Bayes Factor with maximal n " (SBF+maxN) approach as described by Schönbrodt & Wagenmakers (2018) and define an optional stopping threshold to make our data collection

more efficient. Namely, we use moderate evidence in favor of Hypothesis 3 ($BF_{10} > 3$) or against it ($BF_{10} < 1/3$) as thresholds. More precisely, for each experiment, we will first recruit 200 participants and compute the BF_{10} for the SNARC slope difference between ranges. As long as the BF_{10} does not reach any of the two thresholds yet, we will collect another 20 datasets and recalculate the BF_{10} until we reach moderate evidence. If no threshold is reached with our maximal sample size of 800 participants, we will stop the sequential recruiting of participants in any case. Nevertheless, we aim to collect 376 valid datasets.

Participants

For each experiment, adults ~~(i.e., minimum age of 18 years)~~ will be recruited via the recruiting platform Prolific. To comply with our ethics proposal, they must be at least 18 years old, and because of possible age differences in RTs, we set the maximum age to 40 years. As the experiments will be conducted in English, participation is only possible for native English speakers (as per Prolific's screening based on self-reports). Participation will take approximately 20 minutes and will be compensated with ~~£~~€5 (partial payment for partial participation).

Design and experimental task

In the parity judgment task with binary response-key setup, participants will have to indicate as fast and as accurately as possible whether the number presented on the screen is odd or even. The parity judgment task is widely used in numerical cognition and the standard task to investigate the SNARC effect (see Toomarian & Hubbard, 2018, for a review, and Wood et al., 2008, for a meta-analysis). We will assign participants randomly to one of our two experiments. In Experiment 1 (close replication of Dehaene et al., 1993, and Fias et al., 1996), the numbers from 0 to 5 will be used in the lower number range and the numbers from 4 to 9 will be used in the higher number range. In Experiment 2 (conceptual replication), the numbers from 1 to 5 (excluding 3) will be used in the lower and the numbers from 4 to 8 (excluding 6)

in the higher number range, eliminating confounds between number parity and number magnitude (see Table 2) and special influences of zero.

Table 2*Stimulus sets and their characteristics*

Experiment 1 (close replication: number ranges used by Dehaene et al., 1993, and Fias et al., 1996)				Experiment 2 (conceptual replication)			
Lower range		Higher range		Lower range		Higher range	
<u>M</u> Absolut <u>e</u> <u>m</u> agnitude <u>p</u> redictor	<u>C</u> ontrast- <u>c</u> oded <u>p</u> Parity <u>p</u> redictor	<u>M</u> Absolute <u>m</u> agnitude <u>p</u> redictor	<u>C</u> ontrast- <u>c</u> oded <u>p</u> Parity <u>p</u> redictor	<u>M</u> Absolute <u>m</u> agnitude <u>p</u> redictor	<u>C</u> ontrast- <u>c</u> oded <u>p</u> Parity <u>p</u> redictor	<u>M</u> Absolute <u>m</u> agnitude <u>p</u> redictor	<u>C</u> ontrast- <u>c</u> oded <u>p</u> Parity <u>p</u> redictor
0	+0.5	4	+0.5	1	-0.5	4	+0.5
1	-0.5	5	-0.5	2	+0.5	5	-0.5
2	+0.5	6	+0.5	4	+0.5	7	-0.5
3	-0.5	7	-0.5	5	-0.5	8	+0.5
4	+0.5	8	+0.5				
5	-0.5	9	-0.5				
Mean number magnitude depending on number parity:							
$M_{even} = 2$		$M_{even} = 6$		$M_{even} = 3$		$M_{even} = 6$	
$M_{odd} = 3$		$M_{odd} = 7$		$M_{odd} = 3$		$M_{odd} = 6$	
Correlation between number magnitude and number parity:							
$r = -.293$				$r = 0$			

Note. This table gives an overview of the stimulus set we will use in the two experiments. It shows the confound between number parity and number magnitude in both number ranges of Experiment 1 and illustrates how we will avoid it in both number ranges of Experiment 2, such that number parity and number magnitude are uncorrelated (i.e., they are orthogonal to each other as predictors in regression models). Number parity is typically contrast-coded with -0.5 for odd and +0.5 for even numbers when measuring the MARC effect. The number 0 is included in Experiment 1, but we will not use it in the conceptual replication in Experiment 2 because of its special features and irregular mental representation (as outlined in the Introduction). The numbers 4 and 5, which are written in bold in the table, are present in each of the number ranges.

In both experiments, we will use 25 repetitions per number magnitude in each number range (lower vs. higher) and each response-key assignment (MARC congruent, i.e., left-hand responses to odd and right-hand responses to even numbers, vs. MARC incongruent, i.e., right-hand responses to odd and left-hand responses to even numbers). This leads to a total of 600 trials for Experiment 1 and 400 trials for Experiment 2. In each experiment, the trials will be equally divided into four blocks (one per combination of number range and response-key assignment), and a break of minimum 30 seconds must be taken between them. Participants will be randomly assigned to one of four ~~between-subjects conditions that differ in~~ block orders (see Figure ~~32~~). The order of stimulus presentation within blocks will be fully randomized. Each trial will start with a square (extended ASCII 254 with the font size 72px) as eye fixation point (300 ms). Then the number (Open Sans font, size 72px) will be presented until a response is given. A blank screen (500 ms) will conclude the trial. Stimuli as well as fixation squares will be presented in black (0, 0, 0 in RGB notation), while the background remains gray (150, 150, 150 in RGB notation). The time course of an exemplary trial is illustrated in Figure 4. Each block will be preceded by a short practice session in which each number will be presented twice (i.e., 12 practice trials before each block in Experiment 1 and eight practice trials before each block in Experiment 2, respectively). Accuracy feedback will appear during practice sessions only.

Figure ~~32~~

~~Between-subjects conditions~~ Counterbalancing block orders in Experiments 1 and 2

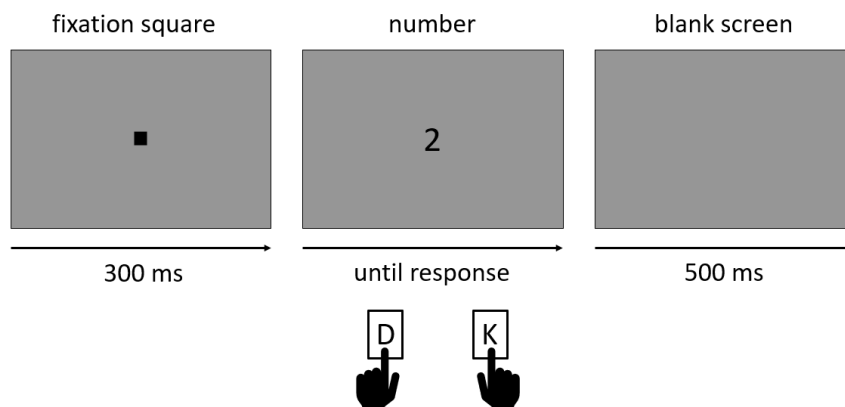
	Condition 1	Condition 2	Condition 3	Condition 4
Block 1	lower range MARC incongruent	lower range MARC congruent	higher range MARC incongruent	higher range MARC congruent
Block 2	lower range MARC congruent	lower range MARC incongruent	higher range MARC congruent	higher range MARC incongruent
Block 3	higher range MARC incongruent	higher range MARC congruent	lower range MARC incongruent	lower range MARC congruent
Block 4	higher range MARC congruent	higher range MARC incongruent	lower range MARC congruent	lower range MARC incongruent

	Block order 1	Block order 2	Block order 3	Block order 4
Block 1	lower range MARC incongruent	lower range MARC congruent	higher range MARC incongruent	higher range MARC congruent
Block 2	lower range MARC congruent	lower range MARC incongruent	higher range MARC congruent	higher range MARC incongruent
Block 3	higher range MARC incongruent	higher range MARC congruent	lower range MARC incongruent	lower range MARC congruent
Block 4	higher range MARC congruent	higher range MARC incongruent	lower range MARC congruent	lower range MARC incongruent

Note. This figure shows the four ~~between-subjects conditions depending on~~ block orders resulting from the combination of range (lower range vs. higher range) and ~~on~~ response-key assignment (MARC congruent, i.e., odd-left and even-right, vs. MARC incongruent, i.e., even-left and odd-right). Each block will be preceded by two repetitions per number as practice trials (12 trials for Experiment 1 and eight trials in Experiment 2), consist of 25 repetitions per number as experimental trials (150 trials for Experiment 1 and 100 trials in Experiment 2) and be followed by a break.

Figure 4

Time course of an exemplary trial



Procedure

The experiments have been set up with WEXTOR (<https://wextor.eu>; Reips & Neuhaus, 2002) in its HTML and JavaScript framework and adapted (see demo version for Experiment 1 at https://luk.uni-konstanz.de/numcog_3/?demo&e1 and for Experiment 2 at https://luk.uni-konstanz.de/numcog_3/?demo&e2). Our previous experiments (for preregistrations, see <https://doi.org/10.17605/OSF.IO/F2GB8>, and <https://doi.org/10.17605/OSF.IO/VBA7N>) have demonstrated that this software is suitable for detecting the SNARC effect in an online setup. At the very beginning of the experiment, a seriousness check (e.g., Reips, 2009⁷) will be applied and participants will be asked whether they want to participate seriously. Participants will be asked to take part only if they wish to give their informed consent, if they use a desktop computer or laptop, and if they are at least between 18 and 40 years old. Then, participants will be asked to provide basic demographic data such as age, gender (man, woman, other), first native language (English and potentially others), handedness (right-handed, left-handed, ambidextrous), and finger-counting habits (starting hand: left hand, right hand, or does not know or no preference; and stability: always, usually, does not know or no preference; in order to replicate findings by Hohol et al., 2022). ~~For each question, participants will have the opportunity to click on~~ For each of the above-mentioned questions, we also provided the option “I prefer not to answer:” to respect some participants’ unwillingness to share information with us and to not force them to choose any option that might not reflect the truth (Jenadeleh et al., 2023; Stieger et al., 2007). Note that in earlier studies, only very few participants chose this option in any of the above-mentioned questions. Next, if not already the case for the default response keys D and K, participants may choose response keys for the experimental task which are to be located in the same row and about one hand width apart from each other on their keyboard. Then, instructions will be displayed, and the first block of the experimental task will start with its practice trials.

After completion of the whole experimental task, participants will be asked to self-rate their math skills compared to people of their age on a visual analogue scale from *very bad* to *very good*. Next, data quality will be assessed by asking participants how they would describe their environment during participation (*silent, very quiet, fairly quiet, fairly noisy, very noisy, or extremely noisy*), whether there were any major distractions during participation (*none, one, or multiple*), and whether there were any difficulties during participation (*yes or no*, text field for comments). Moreover, we will ask participants whether they have used their left index finger for the left response key and their right index finger for the right response key throughout the experiment (yes, partly, or no). Participants will be provided a completion code for Prolific and contact information of our research team. To prevent search engine bots (e.g., Googlebot) from submitting data on our experiment, we will equip the experiment materials with a standardized "noindex, nofollow" meta tag, which prompts search engine bots not to index the experiment pages and also not to visit subsequent pages (see Reips, 2007, p. 379). Further, we will restrict participation to devices over 600 pixel screen width. In addition, to exclude multiple submissions we will perform checks based on User-Agents and IP addresses during data evaluation.

Data preprocessing

We will use the same analysis pipeline as in another of our studies, except for not applying any color vision check (for preregistrations, see <https://doi.org/10.17605/OSF.IO/F2GB8>Roth, Caffier, et al., 2022a, and <https://doi.org/10.17605/OSF.IO/VBA7N>Roth, Caffier, et al., 2022b). This pipeline is similar to that used by Cipora, van Dijck, et al. (2019) in an extensive re-analysis of 18 datasets and permits to reliably detect the SNARC effect. Specifically, only datasets of participants who indicate to be at least between 18 and 40 years old and to seriously participate will be analyzed. Datasets will be excluded if participants describe their environment as very/extremely noisy, ~~or~~ if they report multiple major distractions, or if they report that they were not using their left/right

index finger for the left/right response key, respectively. Practice trials and incorrectly answered trials will not be analyzed in the main analysis. Only trials with RTs between of minimum 200 and 1500 ms will be included in the analysis, because parity judgments faster than 200 ms are very unlikely and faster responses can therefore be treated as anticipations. Moreover, only trials with RTs of maximum 1500 ms will be included, because healthy educated adults should be capable to judge the parity status of single-digit numbers in less than 1500 ms, so that slower responses are unlikely to reflect only the mental process underlying parity judgment but instead might be caused by distractions. Further outliers will be removed in an iterative trimming procedure for each participant separately, such that only RTs that are maximum 3 SDs above or below the individual mean RT of all remaining trials will be considered. This procedure permits to exclude RTs that are unlikely for each given participant and accounts for the right-skewed distribution of RTs, where the means would otherwise be largely overestimated. Finally, only datasets of participants with at least 75% valid remaining trials will be included in the analysis. ~~and~~ Finally, only datasets of participants without any empty experimental cell (number magnitude per response side) in both number ranges will be considered, because an empty cell causes a missing dRT, which in turn makes the calculation of the SNARC slope problematic.

~~Main~~ **Data analysis**

All data analyses will be performed in the statistical computing software R (R Core Team, 2022). An overview of all hypotheses, corresponding tests, and interpretations of possible outcomes is given in the Study Design Table. Instead of frequentist analysis, we decided to take the Bayesian approach. For this, we will determine the BF_{10} associated with the corresponding Bayesian t -test to obtain evidence for both null and alternative hypotheses (using the R package *BayesFactor* by Morey et al., 2015, with a default r -scale of 0.707 as uninformed

~~prior using Cauchy distribution). More specifically, we will calculate Bayesian t -tests and extract the respective BF_{10} . Importantly, considering a BF_{10} larger than 3 as evidence against the null hypothesis is more conservative than rejecting a null hypothesis with a conventional significance level of $\alpha = .05$ in the frequentist approach (Wetzels et al., 2011). As explained above, we will apply the SBF+maxN approach for sequential data analysis with optional stopping in case of at least moderate evidence for or against Hypothesis 3. First, as dropout is more frequent in online than offline research, we will perform a dropout analysis. Differences in dropout rates between experimental conditions would indicate motivational confounding or technical differences (Reips, 2002). Dropout will be analyzed and visualized with dropR (<http://dropr.eu>, Reips & Bannert, 2015).~~

~~For each test described below, a significance level of $\alpha = .01$ will be used. The reason for using a rather conservative significance level is that we will conduct multiple tests per hypothesis, as we investigate AMdependency and RMdependency in two different experiments. Importantly, the significance level does not need to be corrected for the total number of conducted tests in this study, because the tests belong to different test families and because different theoretical inferences can be drawn from their results (Lakens, 2016). Moreover, we will look at each result individually and not generalize from one single significant result within a test family to the presence of an effect in both experiments and in all possible number ranges, so that our interpretations will not inflate the familywise error rate.~~

~~Additionally to all t tests described below, we will determine the Bayes Factor associated with the corresponding Bayesian t test to obtain evidence for both null and alternative hypotheses (using the R package *BayesFactor* by Morey et al., 2015, with a default r -scale of 0.707 as uninformed prior using Cauchy distribution). More specifically, we will calculate the probability of the data under the alternative hypothesis compared to their probability under the null hypothesis and refer to this ratio with BF_{10} in the following. A resulting BF_{10} greater than 3 or 10 will be treated as moderate or strong evidence for the~~

~~alternative hypothesis compared to the null hypothesis, respectively, while a resulting BF_{10} smaller than 1/3 or 1/10 will be treated as moderate or strong evidence for the null hypothesis compared to the alternative hypothesis, respectively (Dienes, 2021).~~

The key dependent variable will be the mean difference between RTs of the right hand minus left hand (dRT), which will be calculated for each number separately per participant and per number range. RTs will be measured as the time from the onset of the number presentation on the screen until the participant's response. A potential SNARC effect can be determined by regressing dRTs on the number magnitude (Fias et al., 1996). One regression will be calculated for each participant and for each number range. Our first dependent measure will be SNARC slopes resulting from the regression of dRTs on number magnitude, which represent the change in relative advantage of right-hand compared to left-hand responses in ms per increase by one in the number magnitude (the more negative the slope, the stronger the SNARC effect). Moreover, we will calculate smallest-number intercepts and mean-number intercepts (when relative magnitude of the numbers in both ranges is matched) as well as dRTs for critical numbers that are part of both number ranges (i.e., 4 and 5). An overview of how the following hypothesis tests can help us distinguish the scenarios is given in Table 1.

First, to test the presence of the SNARC effect on group level (Hypothesis 1), SNARC slopes will be tested against zero with two-sided Bayesian one-sample t -tests in each number range in each experiment. This procedure corresponds to the repeated-measures regressions described by Lorch and Myers (1990) and applied to the SNARC effect by Fias et al. (1996) and accounts for the within-subject design. Although we do not expect reversed, but instead regular SNARC effects reflected by negative slopes (as in each of the six scenarios described above and shown in Figures 1 and 2 and Table 1), we will use two-sided tests here to stay consistent ~~and we will use the same conservative significance level on the relevant side for all hypothesis tests~~ within this study. The lack of conclusive evidence as regards Not finding the SNARC effect or even finding evidence against it in one of the four ranges ~~despite our large~~

~~sample~~ would speak against its robustness, but we consider this to be highly unlikely because the SNARC effect in parity judgment has been shown in numerous studies using different number ranges within the interval from 0 to 9.

Second, to investigate RMdependency of the number mapping on the MNL, we will test whether dRTs for critical numbers (i.e., 4 and 5) differ between the lower and the higher number range (Hypothesis 2) with one two-sided paired Bayesian *t*-test per number in each experiment. ~~A significant~~ Evidence for a difference would imply that the SNARC effect and the MNL are (at least partly) flexible and adapt to the number range used in a task (as in Scenarios 1, 2, 4, and 5). This would be in line with the literature claiming that numbers 4 and 5 are associated with the right side in the number range from 0 to 5 and with the left side in the number range from 4 to 9. However, this finding would not fully rule out AMdependency. ~~The lack of a significant~~ Evidence against a difference ~~(supported by conclusive Bayesian evidence for the null hypothesis)~~ would indicate that the SNARC effect and the MNL are at least not fully flexible (as in Scenarios 3 and 6).

Third, to test AMdependency of the number mapping on the MNL, we will test whether the smallest-number intercepts differ between the lower and the higher number range (Hypothesis 2) with one two-sided paired Bayesian *t*-test in each experiment. ~~A significant result~~ Evidence for a difference would lead to the conclusion that small/large numbers are overall shifted to the left/right on the MNL, respectively (as in Scenarios 2, 3, 5, and 6). In other words, this would imply that the SNARC effect and the MNL are not fully flexible. ~~The lack of a significant difference (supported by conclusive Bayesian evidence for the null hypothesis)~~ Evidence against a difference would indicate that the SNARC effect and the MNL are at least partly flexible (as in Scenarios 1 and 4).

Fourth, to investigate AMdependency of the strength of the SNARC effect, we will compare SNARC slopes between the number ranges (Hypothesis 3) with one two-sided paired Bayesian *t*-test in each experiment. ~~Significantly~~ Evidence for steeper SNARC slopes in the

lower than in the higher number range can be interpreted as stronger SNARC effect within (in absolute terms) smaller than larger numbers (as in Scenarios 4, 5, and 6). This result would lead to the conclusion that the spatial mental representation seems to be more pronounced for small than for large numbers. ~~Evidence against a such difference The lack of a significant difference (supported by conclusive Bayesian evidence for the null hypothesis)~~ would indicate that the strength of the SNARC effect and of the spatial mental representation does not differ between number ranges (as in Scenarios 1, 2, and 3).

Positive controls and manipulation checks

~~To control the data quality in our study, we have implemented a seriousness check (Aust et al., Reips, 2009, review in Reips, 2021) as well as a self-assessment of noise, distractions, and other difficulties. To make sure that we will only analyze trials that reflect mental processes in correctly executed parity judgment, we will exclude incorrectly answered trials and trim RTs (as described in the data preprocessing pipeline). Also, we will exclude full datasets of participants with less than 75% valid trials to only build our results on participants who have understood and followed the task instructions. Moreover, we assess whether participant comply with the instructions to use their left and right index fingers for the left and right response keys, respectively, and only include their datasets into our analysis if they comply with the instructions.~~

~~Last, we will check for the presence of the Odd Effect (Hines, 1990; i.e., overall faster reactions to even than to odd numbers, irrespective of the response side). The Odd Effect is quite robust in the parity judgment task, but independent from the SNARC effect (as it is independent from number magnitude and from its mapping onto space and only considers parity). Therefore, we can consider its investigation as a manipulation check, and in case of its presence we will have a positive control for our experiment. For this, we will subtract the average RT for even numbers from the average RT for odd numbers per participant and test the~~

differences (one per participant) against zero in two-sided Bayesian one-sample t -tests (one per number range, with positive estimates indicating the Odd Effect).

Follow-up analysis

~~If a SNARC effect is observed in all number ranges and dRTs for critical numbers, smallest number intercepts, and SNARC slopes differ significantly between the lower and the higher number range, evidence for all three hypotheses is provided. However, the exact pattern of AMdependency of the number mapping on the MNL (Hypothesis 2) would remain unclear in this case. In other words, the true dRT pattern for a parity judgment task with single digit numbers depending on used the number range could then either look like in Scenario 4 (as observed by Fias et al., 1996) or like in Scenario 5 (as observed by Dehaene et al., 1993). Hence, we will compare mean number intercepts between number ranges in a two-sided paired t test in case all previously mentioned t tests yield a significant difference (i.e., significant SNARC effect in all ranges, different dRTs for critical numbers, different smallest number intercepts, and different slopes, which speaks in favor of Scenarios 4 and 5). A significant result would provide evidence for AMdependency of the number mapping on the MNL (Scenario 5), and a non-significant result (supported by conclusive Bayesian evidence for the null hypothesis) would indicate that the MNL is at least partly flexible (Scenario 4).~~

Exploratory data analysis

~~Moreover, we will investigate the MARC effect in Experiment 2 by regressing dRTs on the contrast-coded number parity (-0.5 for odd and $+0.5$ for even numbers; see Cipora et al., 2019). Number parity can be included as a second predictor in the regression without changing the parameter estimate for number magnitude because number parity and number magnitude are orthogonal to each other in Experiment 2 (as opposed to Experiment 1, see Table 2). The resulting regression slopes for number parity (MARC slopes) in Experiment 2 will represent the change in relative advantage of right-hand responses compared to left-hand responses in ms for even compared to odd numbers, respectively (the more negative the slope, the stronger the~~

~~MARC effect). To investigate the presence of a MARC effect in each number range, MARC slopes will be tested against zero in two-sided t tests for the lower and for the higher number range in Experiment 2 (just as SNARC slopes in the main data analysis for Hypothesis 1). To find out whether MARC slopes are steeper in one than in the other number range, they will be compared between the two ranges used in Experiment 2 with two-sided paired t tests (just as SNARC slopes in the main data analysis for Hypothesis 3). We hypothesize to find a MARC effect in both number ranges but do not expect it to differ between number ranges.~~

~~Finally, we will compare SNARC slopes between different subsamples: first between left and right handers, second between participants starting to count with the fingers of their left and right hand, and third between participants with stable and unstable finger counting habits among both left and right starters. For this purpose, we will run two-sided paired t test for independent samples separately for the lower and higher range in each experiment. Last, we will run two-sided t tests for Pearson's product moment correlation coefficient between self-estimated math skills and SNARC slopes for the lower and higher range separately in each experiment.~~

Possible limitations and unexpected outcomes

~~Not finding a Finding evidence against the SNARC effect in one of the four ranges (Experiment 1: 0 to 5 and 4 to 9; Experiment 2: 1 to 5 [excluding 3] and 4 to 8 [excluding 6]) would be an unexpected outcome which we would not have any explanation for. However, because the SNARC effect in the parity judgment task has been shown in plenty of studies (including online setups) using different number ranges within the interval from 0 to 9 and because our large sample and a high number of repetitions ensure high statistical power to detect even small effects, it seems highly unlikely not to observe a SNARC effect in every of the four ranges. In any case, all further hypothesis tests will be meaningful even if the SNARC effect is not found in all ranges.~~

Even though our Experiment 1 aims to be a direct replication of Dehaene et al.'s (1993) and Fias et al.'s (1996) study, we decided to use 25 instead of 15 repetitions per experimental cell. First, we thereby increase statistical power and measurement precision (Luck, 2019); second, we follow methodological recommendations (Cipora & Wood, 2017); and third, we ensure the comparability with our conceptual replication in Experiment 2. However, because of this methodological improvement, our experiment is therefore strictly speaking not a direct replication.

Just as the original two experiments, our Experiment 1 would have the limitation of the MARC effect being confounded with the SNARC effect because number parity and number magnitude are not orthogonal predictors in the regression model. Therefore, we can only calculate the MARC effect for the data resulting from our Experiment 2. Moreover, because of the special features and an irregular mental representation of the number zero, including it in the stimulus set could drive responses in our Experiment 1. However, we tackle these limitations in our Experiment 2 by using another stimulus set.

Further procedure

Data collection is estimated to last less than one month. Data analysis is expected to be finished within two months after data collection. We plan to write up the full article within three further months for the stage 2 submission.

Data and code availability

Anonymized data and analysis scripts will be available via the Open Science Framework (<https://doi.org/10.17605/OSF.IO/Z43PM> Roth, Lukács, et al., 2022).

Competing interests

The authors declare no financial or non-financial conflicts of interest with the content of this article.

Author contributions

All the authors have full access to all the data and take responsibility for the integrity of the data and the accuracy of the data analysis. *Conceptualization*: K. Cipora, H.-C. Nuerk, U.-D. Reips; *Data Curation*: K. Cipora, ~~G. Lukács~~, H.-C. Nuerk, U.-D. Reips, L. Roth; *Formal Analysis*: K. Cipora, H.-C. Nuerk, U.-D. Reips, L. Roth; *Funding Acquisition*: K. Cipora, H.-C. Nuerk, U.-D. Reips.; *Investigation*: K. Cipora, ~~G. Lukács~~, H.-C. Nuerk, U.-D. Reips, L. Roth; *Methodology*: K. Cipora, H.-C. Nuerk, U.-D. Reips, L. Roth; *Project Administration*: H.-C. Nuerk, U.-D. Reips, L. Roth; *Resources*: H.-C. Nuerk, U.-D. Reips; *Software*: J. Caffier, ~~G. Lukács~~, U.-D. Reips; *Supervision*: K. Cipora, H.-C. Nuerk, U.-D. Reips; *Validation*: K. Cipora, H.-C. Nuerk, U.-D. Reips, L. Roth; *Visualization*: L. Roth; *Writing – original draft*: L. Roth; *Writing – review and editing*: J. Caffier, K. Cipora, ~~G. Lukács~~, H.-C. Nuerk, U.-D. Reips.

Acknowledgements

This research was supported by the DFG project “Replicability of Fundamental Results on Spatial-Numerical Associations in Highly Powered Online Experiments (e-SNARC)” (NU 265/8-1) granted to Hans-Christoph Nuerk and Ulf-Dietrich Reips supporting Lilly Roth as well as John Caffier ~~and Gáspár Lukács~~, with the assistance of Krzysztof Cipora as a cooperation partner. The authors would like to thank Sebastian Sandbrink for English proofreading of this Registered Report.

References

- Allaire, J., Xie, Y., McPherson, J., Luraschi, J., Ushey, K., Atkins, A., Wickham, H., Cheng, J., Chang, W., & Iannone, R. (2022). Rmarkdown: Dynamic Documents for R (version 2.18). R package. <https://github.com/rstudio/rmarkdown>
- Antoine, S., & Gevers, W. (2016). Beyond left and right: Automaticity and flexibility of number-space associations. *Psychonomic Bulletin & Review*, 23(1), 148-155. <https://doi.org/10.3758/s13423-015-0856-x>
- [Aust, F., Diedenhofen, B., Ullrich, S., & Musch, J. \(2013\). Seriousness checks are useful to improve data validity in online research. *Behavior Research Methods*, 45\(2\), 527–535. <https://doi.org/10.3758/s13428-012-0265-2>](https://doi.org/10.3758/s13428-012-0265-2)
- Ben Nathan, M., Shaki, S., Salti, M., & Algom, D. (2009). Numbers and space: Associations and dissociations. *Psychonomic Bulletin & Review*, 16(3), 578-582. <https://doi.org/10.3758/PBR.16.3.578>
- Brysbaert, M. (1995). Arabic number reading: On the nature of the numerical scale and the origin of phonological recoding. *Journal of Experimental Psychology: General*, 124(4), 434-452. <https://doi.org/10.1037/0096-3445.124.4.434>
- Champely, S., Ekstrom, C., Dalgaard, P., Gill, J., Weibelzahl, S., Anandkumar, A., Ford, C., Volcic, R., & De Rosario, H. (2018). Pwr: Basic functions for power analysis (version 1.3-0). R package. <https://CRAN.R-project.org/package=pwr>
- Cipora, K., He, Y., & Nuerk, H.-C. (2020). The spatial–numerical association of response codes effect and math skills: why related? *Annals of the New York Academy of Sciences*, 1477(1), 5-19. <https://doi.org/10.1111/nyas.14355>
- Cipora, K., Patro, K., & Nuerk, H.-C. (2018). Situated influences on spatial–numerical associations. In T. Hubbard (Ed.), *Spatial Biases in Perception and Cognition* (pp. 41–59). Cambridge: Cambridge University Press. <https://doi.org/10.1017/9781316651247>

Cipora, K., Schroeder, P. A., Soltanlou, M., & Nuerk, H.-C. (2018). More space, better mathematics: Is space a powerful tool or a cornerstone for understanding arithmetic? In K. S. Mix & M. T. Battista (Eds.), *Visualizing Mathematics. Research in Mathematics Education* (pp. 77-116). Springer, Cham. https://doi.org/10.1007/978-3-319-98767-5_4

Cipora, K., Soltanlou, M., Reips, U.-D., & Nuerk, H.-C. (2019). The SNARC and MARC effects measured online: Large-scale assessment methods in flexible cognitive effects. *Behavior Research Methods*, *51*(4), 1676-1692. <https://doi.org/10.3758/s13428-019-01213-5>

Cipora, K., van Dijck, J.-P., Georges, C., Masson, N., Goebel, S. M., Willmes, K., Pesenti, M., Schiltz, C., & Nuerk, H.-C. (2019). A Minority pulls the sample mean: On the individual prevalence of robust group-level cognitive phenomena – the instance of the SNARC effect. PsyArXiv. <https://doi.org/10.31234/osf.io/bwyr3>

Cipora, K., & Wood, G. (2017). Finding the SNARC instead of hunting it: A 20*20 Monte Carlo investigation. *Frontiers in Psychology*, *8*, 1194. <https://doi.org/10.3389/fpsyg.2017.01194>

Cleland, A. A., & Bull, R. (2019). Automaticity of access to numerical magnitude and its spatial associations: The role of task and number representation. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *45*(2), 333-348. <https://doi.org/10.1037/xlm0000590>

~~Cohen, J. (1988). The effect size index: d. *Statistical power analysis for the behavioral sciences*. Mahwah, New Jersey: Lawrence Erlbaum Associates, pp. 284-288.~~

Dehaene, S., Bossini, S., & Giraux, P. (1993). The mental representation of parity and number magnitude. *Journal of Experimental Psychology: General*, *122*(3), 371-396. <https://doi.org/10.1037/0096-3445.122.3.371>

- Deng, Z., Chen, Y., Zhu, X., & Li, Y. (2017). The effect of working memory load on the SNARC effect: Maybe tasks have a word to say. *Memory & Cognition*, *45*(3), 428-441. <https://doi.org/10.3758/s13421-016-0676-x>
- Di Giorgio, E., Lunghi, M., Rugani, R., Regolin, L., Dalla Barba, B., Vallortigara, G., & Simion, F. (2019). A mental number line in human newborns. *Developmental Science*, *22*(6), e12801. <https://doi.org/10.1101/159335>
- Dienes, Z. (2021). How to use and report Bayesian hypothesis tests. *Psychology of Consciousness: Theory, Research, and Practice*, *8*(1), 9–26. <https://doi.org/10.1037/cns0000258>
- Fias, W., Brysbaert, M., Geypens, F., & d'Ydewalle, G. (1996). The importance of magnitude information in numerical processing: Evidence from the SNARC effect. *Mathematical Cognition*, *2*(1), 95-110. <https://doi.org/10.1080/135467996387552>
- Fias, W., & van Dijck, J.-P. (2016). The temporary nature of number—space interactions. *Canadian Journal of Experimental Psychology / Revue Canadienne de Psychologie Expérimentale*, *70*(1), 33–40. <https://doi.org/10.1037/cep0000071>
- Fischer, M. H., & Shaki, S. (2014). Spatial associations in numerical cognition – From single digits to arithmetic. *Quarterly Journal of Experimental Psychology*, *67*(8), 1461-1483. <https://doi.org/10.1080/17470218.2014.927515>
- ~~Gelman, A., & Carlin, J. (2014). Beyond Power Calculations: Assessing Type S (Sign) and Type M (Magnitude) Errors. *Perspectives on Psychological Science*, *9*(6), 641–651.~~
- Gevers, W., Santens, S., Dhooge, E., Chen, Q., Van den Bossche, L., Fias, W., & Verguts, T. (2010). Verbal-spatial and visuospatial coding of number–space interactions. *Journal of Experimental Psychology: General*, *139*(1), 180-190. <https://doi.org/10.1037/a0017688>
- Giner-Sorolla, R., Aberson, C. L., Bostyn, D. H., Carpenter, T., Conrique, B. G., Lewis, N. A., & Soderberg, C. (2019). Power to detect what? Considerations for planning and evaluating sample size. Open Science Framework. <https://osf.io/jnmya/>

- Ginsburg, V., & Gevers, W. (2015). Spatial coding of ordinal information in short-and long-term memory. *Frontiers in Human Neuroscience*, 9, 1-10. <https://doi.org/10.3389/fnhum.2015.00008>
- Ginsburg, V., van Dijck, J.-P., Previtali, P., Fias, W., & Gevers, W. (2014). The impact of verbal working memory on number–space associations. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 40(4), 976-986. <https://doi.org/10.1037/a0036378>
- Giurfa, M., Marcout, C., Hilpert, P., Thevenot, C., & Rugani, R. (2022). An insect brain organizes numbers on a left-to-right mental number line. *Proceedings of the National Academy of Sciences*, 119(44), e2203584119. <https://doi.org/10.1073/pnas.2203584119>
- Gökaydin, D., Brugger, P., & Loetscher, T. (2018). Sequential effects in SNARC. *Scientific Reports*, 8(1), 1-13. <https://doi.org/10.1038/s41598-018-29337-2>
- Hines, T. M. (1990). An odd effect: Lengthened reaction times for judgments about odd digits. *Memory & Cognition*, 18, 40-46. <https://doi.org/10.3758/BF03202644>
- Hohol, M., Wołoszyn, K., & Cipora, K. (2022). No fingers, no SNARC? Neither the finger counting starting hand, nor its stability robustly affect the SNARC effect. *Acta Psychologica*, 230, 103765. <https://doi.org/10.1016/j.actpsy.2022.103765>
- Jenadeleh, M., Zagermann, J., Reiterer, H., Reips, U.-D., Hamzaoui, R., & Saupe, D. (2023). *Relaxed forced choice improves performance of visual quality assessment methods. Proceedings of the 15th International Conference on Quality of Multimedia Experience (QoMEX), Ghent, June 2023. ArXiv. <https://doi.org/10.48550/arXiv.2305.00220>*
- Koch, N., Huber, J., Lohmann, J., Cipora, K., Butz, M. V., & Nuerk, H.-C. (2021). Mental number representations are spatially mapped both by their magnitudes and ordinal positions. PsyArXiv. <https://doi.org/10.31234/osf.io/p89h3>
- ~~Lakens, D. (2016, February 14). Why you don't need to adjust your alpha level for all tests you'll do in your lifetime. The 20% Statistician. Retrieved from:~~

<https://daniellakens.blogspot.com/2016/02/why-you-dont-need-to-adjust-you-alpha.html>

Levenson, E., Tsamir, P., & Tirosh, D. (2007). Neither even nor odd: Sixth grade students' dilemmas regarding the parity of zero. *The Journal of Mathematical Behavior*, 26(2), 83-95. <https://doi.org/10.1016/j.jmathb.2007.05.004>

Lorch, R. F., & Myers, J. L. (1990). Regression analyses of repeated measures data in cognitive research. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 16(1), 149-157. <https://doi.org/10.1037/0278-7393.16.1.149>

[Luck, S. \(2019, February 19\). Why experimentalists should ignore reliability and focus on precision. Luck Lab. Retrieved from https://lucklab.ucdavis.edu/blog/2019/2/19/reliability-and-precision](https://lucklab.ucdavis.edu/blog/2019/2/19/reliability-and-precision)

Mitchell, T., Bull, R., & Cleland, A. A. (2012). Implicit response-irrelevant number information triggers the SNARC effect: Evidence using a neural overlap paradigm. *Quarterly Journal of Experimental Psychology*, 65(10), 1945-1961. <https://doi.org/10.1080/17470218.2012.673631>

Morey, R. D., Rouder, J. N., Jamil, T., Urbanek, S., Forner, K., & Ly, A. (2015). BayesFactor: Computation of Bayes factors for common designs (version 0.9.12-4.4). R package. <https://CRAN.R-project.org/package=BayesFactor>

Nieder, A. (2016). Representing something out of nothing: The dawning of zero. *Trends in Cognitive Sciences*, 20(11), 830-842. <https://doi.org/10.1016/j.tics.2016.08.008>

Nuerk, H.-C., Iversen, W., & Willmes, K. (2004). Notational modulation of the SNARC and the MARC (linguistic markedness of response codes) effect. *The Quarterly Journal of Experimental Psychology Section A*, 57(5), 835-863. <https://doi.org/10.1080/02724980343000512>

- Patro, K., Nuerk, H.-C., Cress, U., & Haman, M. (2014). How number-space relationships are assessed before formal schooling: A taxonomy proposal. *Frontiers in Psychology, 5*, 419. <https://doi.org/10.3389/fpsyg.2014.00419>
- Pinhas, M., & Tzelgov, J. (2012). Expanding on the mental number line: Zero is perceived as the “smallest”. *Journal of Experimental Psychology: Learning, Memory, and Cognition, 38*(5), 1187-1205. <https://doi.org/10.1037/a0027390>
- Pinhas, M., Pothos, E. M., & Tzelgov, J. (2013). Zooming in and out from the mental number line: Evidence for a number range effect. *Journal of Experimental Psychology: Learning, Memory, and Cognition, 39*(3), 972-976. <https://doi.org/10.1037/a0029527>
- Proctor, R. W., & Cho, Y. S. (2006). Polarity correspondence: A general principle for performance of speeded binary classification tasks. *Psychological Bulletin, 132*(3), 416-442. <https://doi.org/10.1037/0033-2909.132.3.416>
- R Core Team (2022). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>
- Reips, U.-D. (2000). The Web experiment method: Advantages, disadvantages, and solutions. In M. O. Birnbaum (Ed.), *Psychological experiments on the internet* (pp. 89-117). Academic Press. <https://doi.org/10.1016/B978-012099980-4/50005-8>
- Reips, U.-D. (2002). Standards for Internet-based experimenting. *Experimental Psychology, 49*(4), 243-256. <https://doi.org/10.1026/1618-3169.49.4.243>
- ~~Reips, U.-D. (2009). Internet experiments: Methods, guidelines, metadata. *Human Vision and Electronic Imaging XIV, Proceedings of SPIE, 7240*, Article 724008. <https://doi.org/10.1117/12.823416>~~
- ~~Reips, U. D. (2007). The methodology of Internet-based experiments. In A. Joinson, K. McKenna, T. Postmes, & U. D. Reips (Eds.), *The Oxford handbook of internet psychology* (pp. 373-390). Oxford University Press.~~
- ~~Reips, U. D., & Bannert, M. (2015). dropR: Analyze dropout of an experiment or survey (version 0.9). R package. <http://dropr.eu>~~

- Reips, U.-D., & Neuhaus, C. (2002). WEXTOR: A Web-based tool for generating and visualizing experimental designs and procedures. *Behavior Research Methods, Instruments, & Computers*, 34(2), 234-240. <https://doi.org/10.3758/BF03195449>
- Restle, F. (1970). Speed of adding and comparing numbers. *Journal of Experimental Psychology*, 83(2), 274-278. <https://doi.org/10.1037/h0028573>
- ~~Roth, L., Caffier, J., Cipora, K., Reips, U. D., & Nuerk, H. C. (2022a). SNARC Automaticity: Categorical Color Judgment (blue vs. yellow). Open Science Framework. <https://doi.org/10.17605/OSF.IO/F2GB8>~~
- ~~Roth, L., Caffier, J., Cipora, K., Reips, U. D., & Nuerk, H. C. (2022b). SNARC Automaticity: Color Intensity Judgment (light cyan vs. dark cyan). Open Science Framework. <https://doi.org/10.17605/OSF.IO/VBA7N>~~
- ~~Roth, L., Lukács, G., Reips, U. D., Nuerk, H. C., & Cipora, K. (2022). SNARC Flexibility: Range (In)Dependency. Open Science Framework. <https://doi.org/10.17605/OSF.IO/Z43PM>~~
- Rugani, R., Vallortigara, G., Priftis, K., & Regolin, L. (2015). Number-space mapping in the newborn chick resembles humans' mental number line. *Science*, 347(6221), 534-536. <https://doi.org/10.1126/science.aaa1379>
- ~~Ruscio, J. (2008). A probability-based measure of effect size: Robustness to base rates and other factors. *Psychological Methods*, 13(1), 19-30.~~
- ~~Schönbrodt, F. D., & Wagenmakers, E.-J. (2018). Bayes factor design analysis: Planning for compelling evidence. *Psychonomic Bulletin & Review*, 25(1), 128-142. <https://doi.org/10.3758/s13423-017-1230-y>~~
- Schwarz, W., & Keus, I. M. (2004). Moving the eyes along the mental number line: Comparing SNARC effects with saccadic and manual responses. *Perception & Psychophysics*, 66(4), 651-664. <https://doi.org/10.3758/BF03194909>

- Stieger, S., Reips, U.-D., & Voracek, M. (2007). Forced-response in online surveys: Bias from reactance and an increase in sex-specific dropout. *Journal of the American Society for Information Science and Technology*, 58, 1653-1660. <http://doi.org/10.1002/asi.20651>
- Tlauka, M. (2002). The processing of numbers in choice-reaction tasks. *Australian Journal of Psychology*, 54(2), 94-98. <https://doi.org/10.1080/00049530210001706553>
- Toomarian, E. Y., & Hubbard, E. M. (2018). On the genesis of spatial-numerical associations: Evolutionary and cultural factors co-construct the mental number line. *Neuroscience & Biobehavioral Reviews*, 90, 184-199. <https://doi.org/10.1016/j.neubiorev.2018.04.010>
- Tzelgov, J., Zohar-Shai, B., & Nuerk, H.-C. (2013). On defining quantifying and measuring the SNARC effect. *Frontiers in Psychology*, 4, 302. <https://doi.org/10.3389/fpsyg.2013.00302>
- van Dijck, J.-P., & Fias, W. (2011). A working memory account for spatial–numerical associations. *Cognition*, 119(1), 114-119. <https://doi.org/10.1016/j.cognition.2010.12.013>
- van Dijck, J.-P., Ginsburg, V., Girelli, L., & Gevers, W. (2015). Linking numbers to space: From the mental number line towards a hybrid account. In R. C. Kadosh & A. Dowker (Eds.), *The Oxford handbook of numerical cognition* (pp. 89–105). Oxford University Press.
- Weis, T., Nuerk, H.-C., & Lachmann, T. (2018). Attention allows the SNARC effect to operate on multiple number lines. *Scientific Reports*, 8(1), 1-13. <https://doi.org/10.1038/s41598-018-32174-y>
- Wetzels, R., Matzke, D., Lee, M. D., Rouder, J. N., Iverson, G. J., & Wagenmakers, E.-J. (2011). Statistical Evidence in Experimental Psychology: An Empirical Comparison Using 855 t Tests. *Perspectives on Psychological Science: A Journal of the Association for Psychological Science*, 6(3), 291–298. <https://doi.org/10.1177/1745691611406923>

Wickelmaier, F. (2022). *Simulating the Power of Statistical Tests: A Collection of R Examples*.

ArXiv. <https://doi.org/10.48550/arXiv.2110.09836>

Wühr, P., & Richter, M. (2022). Relative, not absolute, stimulus size is responsible for a correspondence effect between physical stimulus size and left/right responses. *Attention, Perception, & Psychophysics*, 84(4), 1342-1358. <https://doi.org/10.3758/s13414-022-02490-7>

Xie Y. (2022). knitr: A General-Purpose Package for Dynamic Report Generation in R (version 1.41). R package. <https://yihui.org/knitr/>

Yu, S., Li, B., Zhang, S., Yang, T., Jiang, T., Chen, C., & Dong, Q. (2018). Does the spatial-numerical association of response codes effect depend on digits' relative or absolute magnitude? Evidence from a perceptual orientation judgment task. *The Journal of General Psychology*, 145(4), 415-430. <https://doi.org/10.1080/00221309.2018.1532391>

Zohar-Shai, B., Tzelgov, J., Karni, A., & Rubinsten, O. (2017). It does exist! A left-to-right spatial–numerical association of response codes (SNARC) effect among native Hebrew speakers. *Journal of Experimental Psychology: Human Perception and Performance*, 43(4), 719–728. <https://doi.org/10.1037/xhp0000336>

PCI Study Design Table: How Flexible are Spatial-Numerical Associations?

A Registered Replication Report by L. Roth, ~~G. Lukács~~, J. Caffier, U.-D. Reips, H.-C. Nuerk, and K. Cipora

Question	Hypothesis	Sampling plan	Analysis Plan	Rationale for deciding the sensitivity of the test for confirming or disconfirming the hypothesis	Interpretation given different outcomes	Theory that could be shown wrong by the outcomes
Can a SNARC effect be observed in all number ranges?	<p><i>Hypothesis 1:</i> A robust SNARC effect is expected in all used number ranges, i.e., <u>we expect to find at least moderate evidence for SNARC slopes</u> (one per participant and per number range, calculated by regressing dRTs on number magnitude) <u>are expected to be significantly smaller than zero in each number range.</u></p>	<p>To reach the desired statistical power of .90 for the detection of <u>finding moderate evidence (i.e., $BF_{10}^* > 3$) for an effect size of Cohen's $d = 0.1520$ at $\alpha = .01$ in two-sided Bayesian one-sample t-tests or in two-sided Bayesian paired t-tests, at least 376 700 participants need to be tested (for power calculations, see https://doi.org/10.17605/OSF.IO/Z43PM Roth et al., 2022).</u></p>	<p>Four two-sided Bayesian one-sample t-tests of SNARC slopes against zero in each number range separately (Experiment 1: 0 – 5 and 4 – 9; Experiment 2: 1 – 5 excluding 3 and 4 – 8 excluding 6);</p> <p>Bayes Factor associated with two-sided one-sample t-test</p>	<p><u>The most crucial aim of the present study is to find out whether AMdependency of the strength of the SNARC effect exists (Hypothesis 3).</u> The <u>minimally relevant effect size value</u> of $d = 0.1520$ was chosen because <u>it corresponds to the SNARC slope difference of 2.99 ms between number ranges (with a pooled standard deviation of 18.34 ms) that was descriptively found but remained non-</u></p>	<p>Finding <u>moderate or even strong evidence for a SNARC slope that is smaller than 0 in a Bayesian t-test and a high BF_{10}^* in each number range would provide evidence for our hypothesis and be in line with results from previous studies.</u> The significance level of $\alpha = .01$ does not need to be adjusted for multiple tests, because results will be interpreted for each number range separately.</p>	<p>The SNARC effect in the parity judgment task has been shown in numerous studies using different number ranges within the interval from 0 to 9 (as in all scenarios, see Figure 1 and Table 1 in the manuscript). <u>We therefore expect to find at least moderate evidence for it in all four number ranges. Finding at least moderate evidence against the SNARC in any of the four number ranges would be highly surprising</u></p>

		<p><u>However, we will employ the SBF+maxN approach as described by Schönbrodt & Wagenmakers (2018). More precisely, we will first recruit 200 participants and then calculate the BF₁₀ after each added 20 participants. In case the BF₁₀ reaches the threshold of 1/3 or of 3 (i.e., moderate evidence for or against the null hypothesis) before getting to the sample size of 700 participants, we will stop recruiting earlier.</u></p>		<p><u>significant in the original study by Fias et al. (1996) that we wish to replicate here. Note that due to the lacking report of standard deviations, it is not possible to calculate</u></p>		<p><u>given that it is a robust effect in the parity judgment task. Not finding it in one of the four ranges despite our large sample would speak against the robustness of the SNARC effect.</u></p>
<p>Does the number mapping on the MNL³ depend on whether it is the lowest vs. highest number in the current number range?</p>	<p><i>Hypothesis 2:</i> For the same number, a left-/right-hand advantage is expected when it is the lowest/highest number in the current number range, respectively. We hypothesize RMdependency¹ (and possibly AMdependency² as well, see below) of the number mapping on the MNL³.</p>		<p>Four two-sided paired Bayesian <i>t</i>-tests of dRTs for the same number in lower vs. higher number range (i.e., for 4 and 5 in each experiment);</p> <p>Bayes Factor associated with two-sided paired <i>t</i>-test</p>	<p><u>Cohen's <i>d</i> for the slope difference of 9.2 ms found by Dehaene et al. (1993). it corresponds to 1% of the explained variance ($r^2 = .01$). According to the common convention introduced by Cohen (1988), $d = 0.20$ corresponds to a small effect size. Any effect sizes smaller than Cohen's small effect size are usually not practically meaningful.</u></p>	<p><u>If the dRT-Finding moderate or even strong evidence for a different pattern for numbers that appear in both number ranges differs significantly between the lower and the higher number range in a <i>t</i>-test and the BF₁₀ is high*, this would provide evidence for RMdependency¹ of the SNARC effect.</u></p> <p><u>Finding moderate or even strong evidence against a</u></p>	<p>Evidence for RMdependency¹ would indicate flexibility of the MNL³, such that its resolution adapts to the context and that relative magnitude plays a role for spatial-numerical associations. However, this does not rule out the possibility that absolute magnitude plays a role as well (see below).</p> <p>Evidence for AMdependency² would indicate that</p>

					<p>different <u>If the dRT pattern does not differ significantly and the BF₁₀ is low*, this-would indicate AMdependency² of the number mapping on the MNL³.</u></p>	<p>the MNL³ is at least not fully flexible.</p> <p>Full RMdependency is illustrated in Scenarios 1 and 4, full AMdependency is shown in Scenarios 3 and 6, and a combination of both corresponds to Scenarios 2 and 5 in Figure 1.</p>
<p>Does the mapping of numbers on the MNL³ depend on whether they are small vs. high numbers in absolute terms?</p>	<p><i>Hypothesis 2:</i> A left-/right-hand advantage could be observed for small/large numbers in absolute terms, respectively (on top of RMdependency², see above). However, we cannot derive any clear hypothesis from the literature about whether dRTs are lower for</p>		<p>Two two-sided paired <u>Bayesian</u> <i>t</i>-tests of smallest-number intercept in lower vs. higher number range (one test per experiment);</p> <p>Bayes Factor associated with two-sided paired <i>t</i>-test</p>		<p><u>Finding moderate or even strong evidence for different</u> If the smallest-number intercepts differ significantly between in the lower and the higher number range in a <u>Bayesian</u> <i>t</i>-test and the BF₁₀ is high*, this would <u>provide evidence for indicate</u> AMdependency¹</p>	<p>Evidence for AMdependency¹ would indicate that the MNL³ and the SNARC effect are not fully flexible and that absolute magnitude plays a role for spatial-numerical associations. However, this does not rule out the possibility that relative magnitude plays a role as well (see above).</p>

	the smallest number in a higher than in a lower range (as observed by Dehaene et al., 1993, but not by Fias et al., 1996).				of the number mapping on the MNL ³ . <u>Finding moderate or even strong evidence against different smallest-number intercepts</u> If no significant difference is detected and the BF₁₀ is low*, this would indicate RMdependency ² of the number mapping on the MNL ³ .	Evidence for RMdependency ² would indicate that the MNL ³ is at least partly flexible.
Does the strength of the SNARC effect depend on absolute number magnitudes in the used range?	<i>Hypothesis 3:</i> The SNARC effect is expected to be stronger in the lower than in the higher number ranges.		Two two-sided paired <u>Bayesian</u> <i>t</i> -tests of SNARC slopes in lower vs. higher number range (one test per experiment); Bayes Factor associated with two-sided paired <i>t</i>-test		<u>Finding moderate or even strong evidence for a more negative SNARC slope</u> If the SNARC slope is significantly more negative in one of the two number ranges with a high BF₁₀* , <u>would indicate that</u> the SNARC effect seems to be stronger in this	If-Finding the SNARC effect is to be stronger in the lower than in the higher number range, <u>would indicate that</u> the spatial mental representation of small numbers seems to be is more pronounced than for large numbers (as in Scenarios 4, 5, 6 in Figure 1).

number range than in the other.

If the SNARC effect does not differ between number ranges, no evidence can be provided for the strength of the SNARC effect to depend on absolute number magnitudes (as in Scenarios 1, 2, 3).

Positive control or manipulation check:
Can we observe the Odd Effect (Hines, 1990), irrespective of the response side?

We expect to observe the Odd Effect, that is at least moderate evidence for the differences in RTs between odd and even numbers (i.e., average RT for odd numbers minus average RT for even numbers per participant) to be positive.

Four two-sided Bayesian one-sample *t*-tests of differences in RTs between odd and even numbers against zero for each number range separately (Experiment 1: 0 – 5 and 4 – 9; Experiment 2: 1 – 5 excluding 3 and 4 – 8 excluding 6)

Finding moderate or even strong evidence for the Odd Effect would provide evidence for our hypothesis and be in line with results from previous studies. What is even more, we consider this as a *positive control or manipulation check*, such that evidence for the Odd Effect would indicate that response patterns we observe are typical for the

The Odd Effect is quite robust in the parity judgment task, and we therefore expect to find at least moderate evidence for it in all four number ranges. Finding at least moderate evidence against the Odd Effect in any of the four number ranges would be highly surprising.

					<u>parity judgment task.</u>	
<p><i>Follow-up analysis, in case previous tests yield significant results and support all above hypotheses: Does the mapping of numbers on the MNL³, on top of the strength of the SNARC effect, depend on whether they are small vs. high numbers in absolute terms?</i></p>	<p>A left-/right-hand advantage could be observed for small/large numbers in absolute terms, respectively (on top of RMdependency²; see above). However, we cannot derive any clear hypothesis from the literature about whether the mean dRT for the lower range is larger than for the higher range (as in Dehaene et al., 1993) or whether the mean dRTs do not differ between ranges (as in Fias et al., 1996).</p>		<p>Two two-sided paired <i>t</i>-tests of mean-number intercepts (with the middle of each interval being coded as 0, i.e., centered to 2.5 and 6.5 in Experiment 1, and centered to 3 and 6 in Experiment 2) in lower vs. higher number range;</p> <p>Bayes Factor associated with two-sided paired <i>t</i>-test</p>		<p>If the mean-number intercepts differ significantly between number ranges in a <i>t</i> test, such that it is larger for the lower than for the higher number range, and if the BF₁₀ is high*, this would provide evidence for AMdependency¹ of the number mapping on the MNL³.</p> <p>If no significant difference is detected and the BF₁₀ is low*, this would indicate RMdependency¹ of the number mapping on the MNL³.</p>	<p>AMdependency¹ of the number mapping on the MNL would support the hypothesis that absolute (in addition to relative) magnitude plays a role for the spatial mental representation of numbers. In this case, the overall dRT pattern would look similar to Fias et al.'s (1996; see Scenario 4 in Figure 1).</p> <p>Lack of evidence for AMdependency¹ of the number mapping on the MNL would not support the hypothesis that absolute (in addition to</p>

					relative) magnitude plays a role for the spatial mental representation of numbers. In this case, the overall dRT pattern would look similar to Dehaene et al.'s (1993; see Scenario 5).
<i>Exploratory, only for Experiment 2:</i> Can a MARC effect be observed in both number ranges?	A robust MARC effect is expected in both number ranges, i.e., MARC slopes (one per participant and per number range, calculated by regressing dRTs on contrast-coded number parity, where odd numbers are coded with -0.5 and even numbers are coded with +0.5) are expected to significantly differ from zero in each number range.		Two two-sided one-sample t tests of MARC slopes against zero in both number ranges separately (1—5 excluding 3 and 4—8 excluding 6); Bayes Factor associated with two-sided one-sample t test (Note that we will only calculate the MARC effect in Experiment 2 because number parity and number magnitude		Finding a MARC slope that differs from 0 in a t test and a high BF_{10}^* in both number ranges would indicate that the MARC effect can even be found when using only four instead of the typically used eight numbers single digit numbers. The significance level of $\alpha = .01$ does not need to be adjusted for multiple tests, because results will be interpreted

			correlate with each other in Experiment 1.)		for both number ranges separately.	
<i>Exploratory, only for Experiment 2:</i> Does the MARC effect differ between number ranges?	The MARC effect is not expected to differ between number ranges.		One two-sided paired <i>t</i> -test of MARC slopes in lower (1–5 excluding 3) vs. higher (4–8 excluding 6) number range; Bayes Factor associated with two-sided paired <i>t</i> -test		If the MARC slope is significantly more negative in one of the two number ranges with a high BF_{10}^* , the MARC effect seems to be larger in this number range than in the other.	The goal of this exploratory analysis is not to show a theory wrong.
<i>Exploratory:</i> Does the SNARC effect differ between subsamples?	The SNARC effect is not expected to depend on (1) handedness, (2) finger-counting habits, or (3) the stability of finger-counting habits.		Four two-sided paired <i>t</i> -tests of SNARC slopes comparing (1) left vs. right-handers, (2) participants starting to count with their left vs. right hand, and (3) participants with stable vs. unstable finger-counting habits within each starting hand, for each number range and each experiment;		If the SNARC slopes are significantly more negative in one subsample than in the other with a high BF_{10}^* , results can be interpreted as a hint towards differences in the strength of the SNARC effect depending on handedness, finger-counting habits and their stability. However, depending on the	The goal of this exploratory analysis is not to show a theory wrong, but to replicate earlier findings.

			Bayes Factor associated with two-sided paired t -test		size of the subsamples, the power might not be sufficient to detect these differences.	
<i>Exploratory: Does the SNARC effect correlate with self-estimated math skills?</i>	The SNARC effect is not expected to correlate with self-estimated math skills.		Four two-sided t -tests of Pearson's product-moment correlation coefficient between SNARC slopes and self-estimated math skills (one test per number range per experiment); Bayes Factor associated with two-sided correlation t -test		If the correlation is significantly negative with a high BF_{10}^* , the SNARC effect seems to be stronger the higher the self-estimated math skills are. If the correlation is significantly positive with a low BF_{10}^* , the SNARC effect seems to be weaker the higher the self-estimated math skills are.	The goal of this exploratory analysis is not to show a theory wrong.

¹**RMdependency [Relative-Magnitude dependency]:** The SNARC effect dynamically adapts to the stimulus set used in the task and is determined by the relative magnitude of the numbers within the set.

²**AMdependency [Absolute-Magnitude dependency]:** The SNARC effect depends on the absolute magnitude of the numbers.

³**MNL [Mental Number Line]:** The MNL has been proposed as the first explanation for the SNARC effect.

* The BF_{10} is the Bayes Factor defined as the probability of the obtained data under the alternative hypothesis compared to their probability under the null hypothesis. A resulting BF_{10} greater than 3 or 10 will be treated as moderate or strong evidence for the alternative hypothesis compared to the null

hypothesis, respectively, while a resulting BF_{10} smaller than 1/3 or 1/10 will be treated as moderate or strong evidence for the null hypothesis compared to the alternative hypothesis, respectively (Dienes, 2021).

Reference for power calculations:

~~Roth, L., Lukács, G., Reips, U. D., Nuerk, H. C., & Cipora, K. (2022). SNARC Flexibility: Range (In)Dependency. Open Science Framework. <https://doi.org/10.17605/OSF.IO/Z43PM>~~

One and only SNARC? The Flexibility of are Spatial-Numerical Associations. A Registered Report on the SNARC's Range Dependency - Power simulations

Lilly Roth

Version 2: 17th July 2023

This script provides all power calculations that we have run for our Registered Report on the flexibility of spatial-numerical associations and the SNARC's range dependency. It includes Monte-Carlo power simulations from Wickelmaier (2022, <https://doi.org/10.48550/arXiv.2110.09836>) adapted from the frequentist to the Bayesian approach. We calculated Bayes Factors with the R package *BayesFactor* by Morey et al. (2015, <https://CRAN.R-project.org/package=BayesFactor>) for all relevant *t*-tests. At the end of this script, we provide an illustration of the power depending on the used sample size and the true effect size within a plot.

This script was created with the R packages *rmarkdown* by Allaire et al. (2023, <https://cran.r-project.org/web/packages/rmarkdown/index.html>) and *knitr* by Xie et al. (2023, <https://cran.r-project.org/web/packages/knitr/index.html>). This script (Version 2: 17th July 2023) as well as the previous one (Version 1: 28th November 2022) can be downloaded from <https://doi.org/10.17605/OSF.IO/Z43PM>.

```
rm(list = ls())
library("rmarkdown")
library("knitr")
library("BayesFactor")
set.seed(123)
```

Power simulations for the original studies

Most Parameters were taken from from Fias et al. (1996): sample size, standard deviations for low (0,1,2,3,4,5) and high (4,5,6,7,8,9) number range, and standard deviation of the difference between both ranges

For the missing parameter of Pearson product-moment correlation within participants between two blocks, we chose $r = .05$. We have observed this value in our two SNARC automaticity experiments in color judgment tasks (for preregistrations, see <https://doi.org/10.17605/OSF.IO/F2GB8> and <https://doi.org/10.17605/OSF.IO/VBA7N>), which was surprisingly low and might be higher in the parity judgment task. We prefer to take a rather conservative value in order not to overestimate the power. Please note that if the correlation turns out to be higher, the standard deviation is lower (see formula for **SD.Fias.diff** below), so that the corresponding power will be higher than estimates provided here.

```
r <- .05
n.Fias <- 24
SD.Fias.low <- 15.1
SD.Fias.high <- 11.2
SD.Fias.diff <- sqrt(SD.Fias.low^2 + SD.Fias.high^2 - 2*r*SD.Fias.low*SD.Fias.high)
```

```
obs.diff.Fias <- 7.19-10.18
obs.diff.Dehaene <- 10.9-20.1
rep_n <- 5000 # number of repetitions for our simulations
```

Power-determination analysis (Giner-Sorolla et al., 2019)

applied to Fias et al. (1996)

Given the sample size used by Fias and colleagues (i.e., 24 participants), what is the power to detect a given population effect size (e.g., a difference of 10, 5, or 1 in the SNARC slopes)?

Note that we run the following calculations both within the *Bayesian framework*, to ensure comparability between the simulations for the original studies and for our current study, and within the *frequentist framework*, because the original studies were run within the frequentist framework.

```
BF <- replicate(rep_n, {
  d <- rnorm(n = n.Fias, mean = 10, sd = SD.Fias.diff)
  # d = random sample of 24 differences d from normal distribution around 10
  extractBF(ttestBF(d, mu = 0, alternative = "two.sided"))$bf
})
power.10.Bayes <- round(mean(BF > 3), 3)

pval <- replicate(rep_n, {
  d <- rnorm(n.Fias, mean=10, sd=SD.Fias.diff)
  # d = random sample of 24 differences d from normal distribution around 10
  t.test(d, mu=0, alternative = "two.sided", conf.level = .95)$p.value
})
power.10.freq <- round(mean(pval < .05), 3)
```

0.553 power to find moderate evidence ($BF_{10} > 3$, Bayesian framework) and 0.731 power to detect a significant effect ($p < .05$, frequentist framework) for a difference of 10 in the SNARC slopes (i.e., increase of right- hand advantage in ms per magnitude unit) between ranges in a *t*-test with $n = 24$ and $sd = 15.1$ ms for the lower and $sd = 11.2$ ms for the higher range

```
BF <- replicate(rep_n, {
  d <- rnorm(n = n.Fias, mean = 5, sd = SD.Fias.diff)
  # d = random sample of 24 differences d from normal distribution around 5
  extractBF(ttestBF(d, mu = 0, alternative = "two.sided"))$bf
})
power.5.Bayes <- round(mean(BF > 3), 3)

pval <- replicate(rep_n, {
  d <- rnorm(n.Fias, mean=5, sd=SD.Fias.diff)
  # d = random sample of 24 differences d from normal distribution around 5
  t.test(d, mu=0, alternative = "two.sided", conf.level = .95)$p.value
})
power.5.freq <- round(mean(pval < .05), 3)
```

0.132 power to find moderate evidence ($BF_{10} > 3$, Bayesian framework) and 0.24 power to detect a significant effect ($p < .05$, frequentist framework) for a difference of 5 in the SNARC slopes (i.e., increase of right- hand advantage in ms per magnitude unit) between ranges in a *t*-test with $n = 24$ and $sd = 15.1$ ms for the lower and $sd = 11.2$ ms for the higher range

```

BF <- replicate(rep_n, {
  d <- rnorm(n = n.Fias, mean = 1, sd = SD.Fias.diff)
  # d = random sample of 24 differences d from normal distribution around 1
  extractBF(ttestBF(d, mu = 0, alternative = "two.sided"))$bf
})
power.1.Bayes <- round(mean(BF > 3), 3)

pval <- replicate(rep_n, {
  d <- rnorm(n.Fias, mean=1, sd=SD.Fias.diff)
  # d = random sample of 24 differences d from normal distribution around 1
  t.test(d, mu=0, alternative = "two.sided", conf.level = .95)$p.value
})
power.1.freq <- round(mean(pval < .05), 3)

```

0.022 power to find moderate evidence ($BF_{10} > 3$, Bayesian framework) and 0.062 power to detect a significant effect ($p < .05$, frequentist framework) for a difference of 1 in the SNARC slopes (i.e., increase of right-hand advantage in ms per magnitude unit) between ranges in a *t*-test with $n = 24$ and $sd = 15.1$ ms for the lower and $sd = 11.2$ ms for the higher range

To sum up, with the standard deviations observed by Fias et al. (1996), their sample was not large enough to find evidence for SNARC slope differences of 10 (power of 0.553 in a Bayesian and 0.731 in a frequentist analysis), 5 (0.132 Bayesian and 0.24 frequentist), or (0.022 Bayesian and 0.062 frequentist) between the number ranges.

Effect-size sensitivity approach: (Giner-Sorolla et al., 2019)

applied to Fias et al. (1996)

Given the sample size used by Fias and colleagues (i.e., 24 participants) and a desired power level (e.g., 0.8, 0.9, or 0.95), what is the minimum population effect size that can be detected?

Note that we run the following calculations both within the *Bayesian framework*, to ensure comparability between the simulations for the original studies and for our current study, and within the *frequentist framework*, because the original studies were run within the frequentist framework.

```

BF <- replicate(rep_n, {
  d <- rnorm(n = n.Fias, mean = 12.8, sd = SD.Fias.diff)
  extractBF(ttestBF(d, mu = 0, alternative = "two.sided"))$bf
})
power.Fias.80.Bayes <- round(mean(BF > 3), 3)

pval <- replicate(rep_n, {
  d <- rnorm(n.Fias, mean=11.0, sd=SD.Fias.diff)
  t.test(d, mu=0, alternative = "two.sided", conf.level = .95)$p.value
})
power.Fias.80.freq <- round(mean(pval < .05), 3)

```

0.80 power to find moderate evidence ($BF_{10} > 3$, Bayesian framework) for a difference of 12.8 ms (i.e., $d = 12.8 / SD.Fias.diff = 0.70$) or to find a significant effect ($p < .05$, frequentist framework) for a difference of 11.0 ms (i.e., $d = 11.0 / SD.Fias.diff = 0.60$) in the SNARC slopes between ranges with $n = 24$ and $sd = 15.1$ ms for the lower and $sd = 11.2$ ms for the higher range

```

BF <- replicate(rep_n, {
  d <- rnorm(n = n.Fias, mean = 14.6, sd = SD.Fias.diff)
  extractBF(ttestBF(d, mu = 0, alternative = "two.sided"))$bf
})
power.Fias.90.Bayes <- round(mean(BF > 3), 3)

pval <- replicate(rep_n, {
  d <- rnorm(n.Fias, mean=12.7, sd=SD.Fias.diff)
  t.test(d, mu=0, alternative = "two.sided", conf.level = .95)$p.value
})
power.Fias.90.freq <- round(mean(pval < .05), 3)

```

0.90 power to find moderate evidence ($BF_{10} > 3$, Bayesian framework) for a difference of 14.6 ms (i.e., $d = 14.6 / SD.Fias.diff = 0.80$) or to find a significant effect ($p < .05$, frequentist framework) for a difference of 12.7 ms (i.e., $d = 12.7 / SD.Fias.diff = 0.69$) in the SNARC slopes between ranges with $n = 24$ and $sd = 15.1$ ms for the lower and $sd = 11.2$ ms for the higher range

```

BF <- replicate(rep_n, {
  d <- rnorm(n = n.Fias, mean = 16.0, sd = SD.Fias.diff)
  extractBF(ttestBF(d, mu = 0, alternative = "two.sided"))$bf
})
power.Fias.95.Bayes <- round(mean(BF > 3), 3)

pval <- replicate(rep_n, {
  d <- rnorm(n.Fias, mean=14.1, sd=SD.Fias.diff)
  t.test(d, mu=0, alternative = "two.sided", conf.level = .95)$p.value
})
power.Fias.95.freq <- round(mean(pval < .05), 3)

```

0.95 power to find moderate evidence ($BF_{10} > 3$, Bayesian framework) for a difference of 16.0 ms (i.e., $d = 16.0 / SD.Fias.diff = 0.87$) or to find a significant effect ($p < .05$, frequentist framework) for a difference of 14.1 ms (i.e., $d = 14.1 / SD.Fias.diff = 0.77$) in the SNARC slopes between ranges with $n = 24$ and $sd = 15.1$ ms for the lower and $sd = 11.2$ ms for the higher range

To sum up, with the standard deviations observed by Fias et al. (1996) and with the sample size they used, only unreasonably large SNARC slope differences (i.e., larger than typical SNARC slopes themselves) could have been detected at adequate power levels.

Power simulations for the current study

Parameters for Monte-Carlo power simulations

In the following, we will simulate the statistical power for different sample sizes and different SNARC slope differences between number ranges that can be considered to be minimally relevant effects.

r = Pearson product-moment correlation of unstandardized SNARC slopes between two blocks of around .05, as in our two SNARC automaticity experiments


```
r <- .05
```

s = standard deviation for slopes between participants in each range

In our two SNARC automaticity experiments, the standard deviations for slopes were 4.21 and 3.93, and in Fias et al. (1996), the pooled standard deviation for slopes in the lower and higher ranges was 13.29. Although we do not think that there generally is a higher variability of slopes in the parity judgment task as compared to color judgment tasks, and although our planned online study will have high measurement precision, so that we expect rather small standard deviations in the current study, we use the pooled standard deviation from Fias et al. (1996) here, which is more conservative.

```
s <- sqrt( (SD.Fias.low^2 + SD.Fias.high^2) / 2 )
```

sxy = covariance of unstandardized SNARC slopes between two blocks

```
sxy <- r*s*s
```

necessary_n = sample size

to be determined for the current study and to be varied for illustrating different power scenarios in a plot

mean_d = minimal effect size of interest:

```
mean_d <- 0.15
```

effect and effect size detected by Fias et al. (1996)

```
E.Fias <- 7.19-10.18
```

```
ES.Fias <- E.Fias/SD.Fias.diff
```

Sample size calculation for the current study

```
necessary_n <- 800
BF <- replicate(rep_n, {
  d <- rnorm(n = necessary_n, mean = mean_d, sd = 1)
  extractBF(ttestBF(d, mu = 0, alternative = "two.sided"))$bf
})
power.samplesize <- round(mean(BF > 3), 3)
```

In order to achieve a power of 0.900 to find moderate evidence ($BF_{10} > 3$) for the minimally relevant effect size of $d = 0.15$, 800 datasets need to be collected.

Monte-Carlo power simulations with different sample sizes for $BF_{10} > 3$

In the following, we simulate the power for various differences in SNARC slope (i.e., between -15 and 15 in steps of 0.25 ms per number magnitude) and with different sample sizes (i.e., 20, 40, 80, 160, 320, 640).

```

difference <- seq(from = -15, to = 15, by = 0.25)
data.frame <- data.frame(matrix(ncol = 3, nrow = length(difference)))
colnames(data.frame) <- c("samplesize", "difference", "simulatedpower")

power.simulation <- function(samplesize){
  for (i in seq_along(difference)){
    data.frame$samplesize[i] <- samplesize
    data.frame$difference[i] <- difference[i]

    BF <- replicate(rep_n, {
      d <- rnorm(n = samplesize, mean = difference[i],
                sd = sqrt(s^2 + s^2 - 2*sxy))
      extractBF(ttestBF(d, mu = 0, alternative = "two.sided"))$bf
    })
    data.frame$simulatedpower[i] <- mean(BF > 3)
  }
  return(data.frame)
}

power020 <- power.simulation(samplesize = 20)
power040 <- power.simulation(samplesize = 40)
power080 <- power.simulation(samplesize = 80)
power160 <- power.simulation(samplesize = 160)
power320 <- power.simulation(samplesize = 320)
power640 <- power.simulation(samplesize = 640)

power.all <- rbind(power020, power040, power080, power160, power320, power640)

# setw("...")
write.table(power.all, file = "RegisteredReport_Study3_SNARC-Flexibility_Roth_Power_v2.txt",
            sep = "\t", dec = ".", quote = FALSE, row.names = FALSE)

```

Power plot

In the following, we create a plot illustrating the results we obtained in the above power simulations.

```

rm(list = ls())
power.all <- read.table("RegisteredReport_Study3_SNARC-Flexibility_Roth_Power_v2.txt",
                       header = TRUE)

power020 <- power.all[power.all$samplesize == 20,]
power040 <- power.all[power.all$samplesize == 40,]
power080 <- power.all[power.all$samplesize == 80,]
power160 <- power.all[power.all$samplesize == 160,]
power320 <- power.all[power.all$samplesize == 320,]
power640 <- power.all[power.all$samplesize == 640,]

pdf("RegisteredReport_Study3_SNARC-Flexibility_Roth_Power.pdf",
    height = 6, width = 6, pointsize = 13)

```

```

par(mgp = c(2, .7, 0), mai = c(.8, .8, .1, .1))

plot(simulatedpower ~ difference, power020,
     type = "l", lty = 1,
     xlim = c(-15, 15), ylim = c(0,1),
     xlab = "SNARC slope difference (in ms) per magnitude unit",
     ylab = "Simulated power")

points(simulatedpower ~ difference, data = power040, type = "l", lty = 1)
points(simulatedpower ~ difference, data = power080, type = "l", lty = 1)
points(simulatedpower ~ difference, data = power160, type = "l", lty = 1)
points(simulatedpower ~ difference, data = power320, type = "l", lty = 1)
points(simulatedpower ~ difference, data = power640, type = "l", lty = 1)

points(simulatedpower[seq(from = 1, to = nrow(power020), by = 4)]
       ~ difference[seq(from = 1, to = nrow(power020), by = 4)],
       data = power020, type = "p", pch = 0)
# only draw points for every fourth simulated point to make the graph not too crowded

points(simulatedpower[seq(from = 1, to = nrow(power040), by = 4)]
       ~ difference[seq(from = 1, to = nrow(power040), by = 4)],
       data = power040, type = "p", pch = 1)
points(simulatedpower[seq(from = 1, to = nrow(power080), by = 4)]
       ~ difference[seq(from = 1, to = nrow(power080), by = 4)],
       data = power080, type = "p", pch = 4)
points(simulatedpower[seq(from = 1, to = nrow(power160), by = 4)]
       ~ difference[seq(from = 1, to = nrow(power160), by = 4)],
       data = power160, type = "p", pch = 15)
points(simulatedpower[seq(from = 1, to = nrow(power320), by = 4)]
       ~ difference[seq(from = 1, to = nrow(power320), by = 4)],
       data = power320, type = "p", pch = 16)
points(simulatedpower[seq(from = 1, to = nrow(power640), by = 4)]
       ~ difference[seq(from = 1, to = nrow(power640), by = 4)],
       data = power640, type = "p", pch = 16)

# desired power level of 0.900
abline(h = 0.9, lty = 2, lwd = 2)

# difference in SNARC slopes (high range - low range) descriptively observed by
# Fias et al. (1996)
obs.diff.Fias <- 7.19-10.18
abline(v = obs.diff.Fias, lty = 2, lwd = 2, col = "deepskyblue")
# difference in SNARC slopes (high range - low range) descriptively observed by
# Dehaene et al. (1993)
obs.diff.Dehaene <- 10.9-20.1
abline(v = obs.diff.Dehaene, lty = 2, lwd = 2, col="green")

legend(x = -3.5, y = 1.08, expression("Fias et al. \n(1996)"), title = "",
      bty = "n", text.col = "deepskyblue", cex = 0.9, pt.cex = 1)
legend(x = -16, y = 0.24, expression("Dehaene \net al. \n(1993)"), title = "",
      bty = "n", text.col = "green", cex = 0.9, pt.cex = 1)

# probability of inconclusive/misleading evidence (BF10 < 3) despite true effect

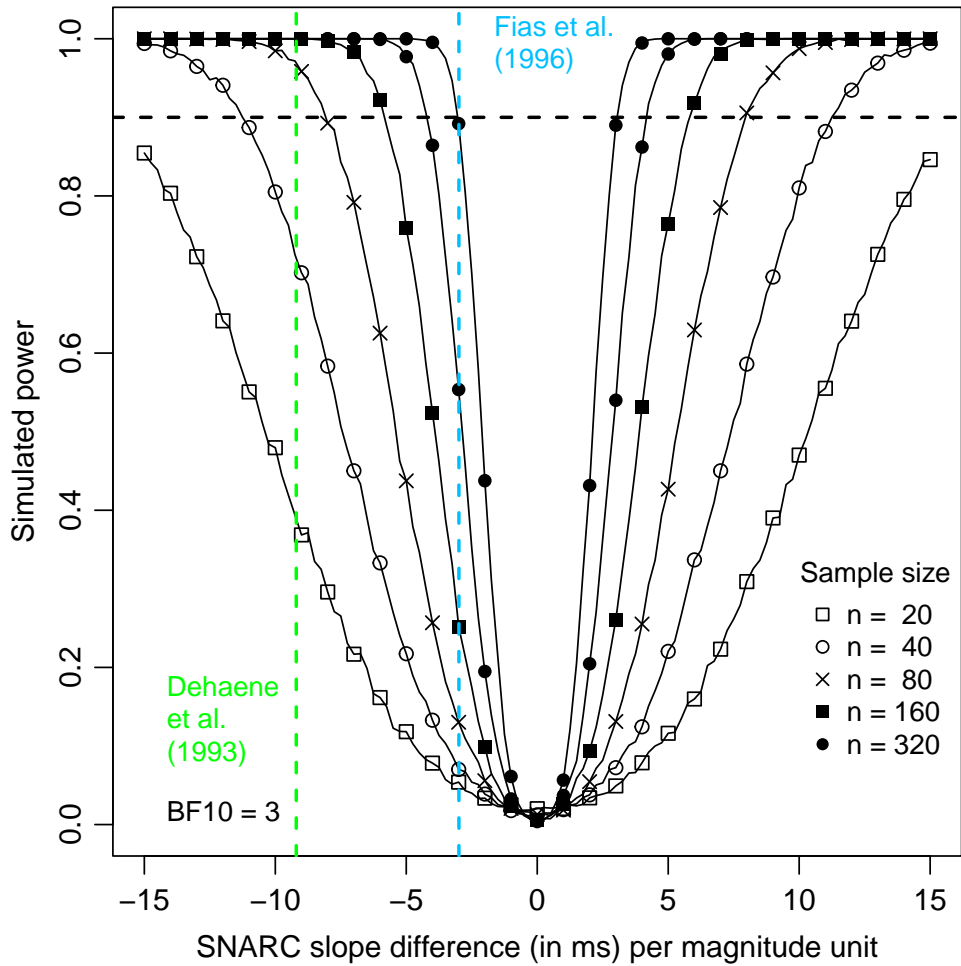
```

```

abline(h = 3, lty = 2)
legend(x = -16, y = 0.10, expression("BF10 = 3"), title = "", bty="n",
      cex = 0.9, pt.cex = 1)

legend(x = 9.8, y = 0.35, expression("n = 20", "n = 40", "n = 80", "n = 160", "n = 320"),
      title = "Sample size", pch = c(0, 1, 4, 15, 16), bty = "n", cex = 0.9, pt.cex = 0.9)
dev.off()

```



This plot shows the simulated power (y-axis) to find moderate evidence ($BF_{10} > 3$) for different SNARC slope differences (x-axis) depending on sample size (20, 40, 80, 160, 320, 640), while assuming a pooled standard deviation of $s = 13.29$ like in Fias et al. (1996) and a Pearson product-moment correlation of unstandardized SNARC slopes between two blocks of around $r = .05$ as in our two SNARC automaticity experiments in color judgment tasks (for preregistrations, see <https://doi.org/10.17605/OSF.IO/F2GB8> and <https://doi.org/10.17605/OSF.IO/VBA7N>).

The observed effect sizes found in the two original studies are shown as green (Dehaene et al., 1993) and blue (Fias et al., 1996) dashed vertical lines. The desired power level of 0.9 is shown as dashed horizontal line.

To sum up, the sample sizes used by Dehaene et al. (1993) and by Fias et al. (1996) were not large enough to detect plausible SNARC slope differences. To achieve a power of 0.90

to detect a minimally relevant effect size of $d = 0.15$, we will therefore collect data from 800 participants.