# UC Merced

**Title**

A task-general model of human randomization

**Permalink**

**Journal**

**Authors**

Yousif, Sami R
McDougle, Samuel David
Rutledge, Robb B

**Publication Date**

2022

Peer reviewed

# A task-general model of human randomization

**Sami R. Yousif (sami.yousif@yale.edu)**
**Samuel D. McDougle (samuel.mcdougle@yale.edu)**
**Robb B. Rutledge (robb.rutledge@yale.edu)**

Department of Psychology, Yale University, New Haven, CT, 06520 USA

## Abstract

Does the human mind contain a task-general 'randomization machine'? Stable biases of randomization have been identified that span multiple domains and modalities, in both lower-level perceptual tasks and in higher-level cognitive tasks. The stability of such biases indicates that the mind may rely on a stable set of properties to create and perceive randomness. But what computational principles support randomization? Here, we approach this question by building a computational model of human randomization that generalizes across spatial and numerical tasks. We show that simple computational heuristics capture higher-order properties of human-generated random sequences, in both numerical and spatial randomization tasks each with many possible options. Furthermore, we show that human behavior in both types of tasks can be approximated by the *same* low-dimensional model, implying that a domain-general set of computational principles may underlie randomization behavior in general.

**Keywords:** randomness; alternation bias; random number generation; spatial randomness; mental number line

## Introduction

Humans reliably misunderstand or misperceive randomness. For example, people tend to think that sequences of coin flips that best represent randomness are ones that alternate between heads and tails more often than truly random sequences (e.g., Kahneman & Tversky, 1972). When people are tasked with generating random sequences themselves, they also generate sequences with elevated numbers of alternations (for review, see Bar-Hillel & Wagenaar, 1991). Indeed, there seem to be stable biases of randomness that exist across cognitive tasks (for review, see Bar-Hillel & Wagenaar, 1991) and perceptual tasks (e.g., Reiner et al., 2021; Yu et al., 2018), that persist across modality (e.g., in vision and in audition; Yu et al., 2018), and that persist across domains (e.g., for numbers but also letters, coin-flips, arbitrary button presses, etc.; see, e.g., Bar-Hillel & Wagenaar, 1991). Here, we formalize a model of random behavior that (a) succeeds in capturing key features of human randomness, and (b) generalizes across distinct tasks (i.e., a 'random number' task and a 'random location' task).

## Why should we care about random behavior?

The conception and perception of randomness is one of the building blocks of human cognition: to perceive the structure of the world, the mind must differentiate signal from noise. Thus, detecting randomness is vital to learning, whether that be simple conditioning (e.g., Rescorla & Wagner, 1972),

statistical learning (for review, see Sherman et al., 2020), or even higher-level language acquisition (e.g., Kelly & Martin, 1994). Randomness is also an important factor in many human behaviors. For instance, fallacies of randomness influence not only gambling behavior, but also everyday activities like reading stock charts or interpreting weather forecasts.

## Human randomness

Perhaps the single most robust feature of human randomness is an 'over alternation bias'. This bias describes the tendency to generate sequences of coin flips with more alternations than would be expected in truly random sequences, and to perceive such sequences (i.e., ones with slightly more alternations than average) as the most random (for review, see Bar-Hillel & Wagenaar, 1991). This is true not just for coin flips, but also for other kinds of visual stimuli (e.g., grids of alternating colors) as well as auditory stimuli (Yu et al., 2018). Further, this bias is consistent with other known fallacies of subjective probability (e.g., the Gambler's fallacy; Kahneman & Tversky, 1972; Reuter et al., 2005; Wagenaar, 1988).

Binary sequences like coin-flips are of course not the only way we encounter random information. Even in many gambling scenarios, there are more than two possible options (e.g., at a blackjack or roulette table). With more options on the table, it is not obvious what the 'over alternation bias' entails. Does it reflect a tendency to avoid repeating identical choices back-to-back, or a more general tendency to depart from a previous "area" of the possibility space (in which case other nearby options should be similarly less likely)? Given this uncertainty, understanding randomness in scenarios with more than two possible options (see, e.g., Towse & Neil, 1998) can enrich our understanding of human randomization.

Here we aim to assess human randomness in three environments that are more complex than typical binary sequence tasks: (1) A random number generation task where participants iteratively generated single-digit random integers [1-9], (2) A random location generation task in which participants iteratively generated random locations along a line, and (3) A random location task in a two-dimensional plane. The goal of our study was to identify stable computational principles that underlie random behavior across disparate tasks in order to begin bridging the gap between laboratory studies of human randomization and the real world.

## Current Study

We propose a computational model of human randomization behavior that generalizes across two behavioral domains. In the first study, a random number generation task, participants were asked to generate sequences of 250 random integers between 1 and 9. In two additional studies, random location generation tasks, participants were asked to generate sequences of 250 random locations (either on a horizontal line, or in a two-dimensional square). The model aims to capture features of human behavior that generalize across these three different tasks.

## Study 1

In a first study, participants played a random number game. In this task, they were instructed to iteratively press numbers on the keyboard while trying to be "as random as possible."

## Method

**Participants** Participants recruited via Prolific (N = 200) completed the task in exchange for monetary compensation (at a rate of approximately $10/hour).

**Design and Procedure** The task was administered via custom JavaScript code. Participants were told to generate sequences of random numbers to the best of their ability by pressing the number buttons on their keyboard. Each time a number was pressed, it would appear in the center of the screen for 750ms. Once it disappeared, participants were free to select another number. There was no time limit on responses. This would continue until 250 numbers had been selected. Prior to the task, participants completed four practice trials that were identical to task trials.

**Exclusions** Participants were asked to be as random as possible, with minimal explanation. Given the nature of the task, it is possible for participants to respond in intentionally non-random ways (e.g., with a sequence such as: 1, 1, 1, 3, 3, 3, 5, 5, 5, etc.). Patterns like this in the human data would artificially inflate model predictions; the models could easily pick up on such patterns, but these patterns would not reflect the construct we are interested in (i.e., human randomization). Thus, to be conservative, we applied two exclusion criteria. First, we excluded any participant who selected any one of the possible numbers fewer than 10 times (i.e., less than 4% of the time). This number was chosen because ~99.99% of truly random (uniform) distributions of choices should include at least 10 of each number given our task design. Second, we excluded any participant whose "average numerical distance" (the mean numerical distance between successive numbers; see *Results*) was less than 2.4. Again, ~99.99% of truly random samples should be above this value. From the original sample size of 200 participants, this yielded a conservative sample of 171 participants. Note: We are intentionally excluding participants who were non-random because those participants are *easier* to explain and predict. In other words, we are trying to ensure that our model captures the behavior of participants who are *trying* to generate sequences of random numbers.

## Results & Discussion

A matrix depicting the overall pattern of transition probabilities can be seen in Figure 1A. As is evident from the figure, participants behavior is not truly random.

**Model comparisons** We compared three separate models to explain these data. The first model is a simple "stay/switch" model. It captures the tendency to repeat choices back-to-back. For example, if the participant selected choice $C^i$ on trial $t$ (where $C$ represents the space of possible choices, and $i$ indexes the specific value chosen from $C$), the probability of choosing $C^i$ for the next trial $t+1$ is adjusted based on a stay/switch 'bonus' parameter, $\varepsilon$:

$$P(C^i)_{t+1} = P(C^i)_{t+1} + \varepsilon$$

Where a positive parameter captures a tendency to repeat choices and a negative parameter captures a tendency to switch choices. The second model is a "side-switch" model. It captures the tendency to switch sides of the distribution — e.g., the likelihood that someone choosing a number greater than 5 (the middle/median of the distribution) would subsequently choose a number lower than 5, or vice versa. For example, if the choice $C^i$ at trial $t$ was greater than 5, the probability of choosing all numbers greater than 5 would be increased/decreased based on the side-switch parameter, $\eta$:

$$P(C^i_{<5})_{t+1} = P(C^i_{<5})_{t+1} + \eta$$

And the probability of any choice greater than 5 would be adjusted in an equal but opposite way:

$$P(C^i_{>5})_{t+1} = P(C^i_{>5})_{t+1} - \eta$$

After both adjustments, the probabilities for the nine options would be normalized to sum to 1 using the softmax function. The third model is a combined model that incorporates the parameters of both previous models.

Each model was fit to each participant's observed choices using maximum likelihood estimation. After fitting each single-parameter model, the single free parameter in the stay/switch model ($\varepsilon$) was found to be significantly negative (signrank test, $p<.001$), meaning that participants avoided repeating choices back-to-back. The single free parameter in the side-switch model ($\eta$) was found to be significantly positive (signrank test, $p<.001$), reflecting a side alternation bias. After fitting the combined model, the stay/switch parameter similarly came out significantly negative (signrank test, $p<.001$), and the side-switch parameter was again significantly positive (signrank test, $p<.001$).

The combined model had the lowest (best) total AIC score (183390), followed by the stay/switch model (184360), then the side-switch model (186207). This suggests participants used both heuristics to generate random numbers.
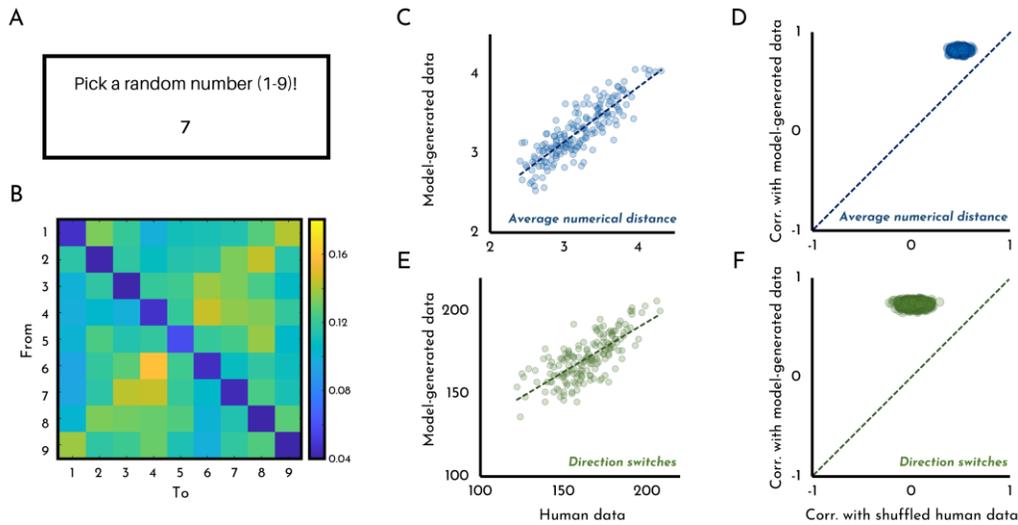
Figure 1: Results from Study 1. (A) Transition probability from each value to each other value. (B) Correlation between the *average numerical distance* and (D) the *direction switches* in the human data and model-generated data. (C) Correlations from 1000 simulations between human data and shuffled human data (x-axis) as well as correlations from 1000 simulations between human data and model-generated data (y-axis), for the *average numerical distance* metric and (E) the *direction switches* metric. Points above the identity line indicate better model performance.

**Model validation** Participants in this task were asked to generate random sequences of numbers to the best of their ability. Given the nature of this task, participants completing it optimally are impossible to predict. To validate whether our models capture meaningful aspects of the observed behavior, we can ask whether the model successfully captures higher-order properties of the human data that are not explicit in the model itself. Here, we selected two such higher-order properties: (1) *average numerical distance* (i.e., the average difference between each chosen number and the next), and (2) *direction switches* (i.e., the number of times that the 'direction' of responses changes between increasing and decreasing; e.g., the sequence [1, 3, 7, 4, 2, 9] contains two direction switches [7-4, and 2-9]; repeated values would not constitute a direction switch).

To validate the models, we asked whether data generated using the model parameters resembles human-generated sequences. We did this 1000 times for each modeled subject and asked whether the average correlation between a given property (i.e., *average numerical distance* or *direction switches*) between the model-generated and human-generated data was greater than the average correlation that would be expected in the "null" distribution generated by shuffling human data and correlating properties of that data with properties of the true human data. Thus, our critical comparisons were between two different correlations: The first between the model-generated data and the human data, and the other between the human data and that same shuffled human data. If the former is higher, that would indicate successful prediction. For example, we wanted to know whether the model is capturing the typical 'average numerical distance' of human data, and so we asked whether the correlation between the 'average numerical distance value'

for each participants' model-generated data and the true data was *higher* than the correlation between the true data and the shuffled data. The logic of this comparison is that it ensures any success of the model cannot be attributed to the base-rates with which people select certain numbers; in other words, we were testing whether the model was picking up on features that were specific not just to which numbers people selected, but the *order* in which they selected those numbers.

The average correlation for the *average numerical distance* measure between the human data and randomly shuffled human data was $r=.50$. For the combined model, the average correlation between model-generated and human data was $r=.81$ (see Figure 1B for a representative example of the model fits; see Figure 1C for the correlation comparisons across the 1000 simulations). In other words, the model is sufficient to capture most of the variation in the *average numerical distance* metric.

The average correlation for the *direction switches* measure between the model data and randomly shuffled human data was $r=.04$. For the combined model, the average correlation between model-generated and human data was $r=.73$ (see Figure 1D for a representative example of the model fits; see Figure 1E for the correlation comparisons across the 1000 simulations). Again, the model captures most of the variation in the *direction switches* metric.

Taken together, both the model fitting and model validation results suggest that both the stay/switch and side-switch heuristics are important for capturing human randomization behavior during number generation. Not only does the combined model best predict behavior, but it captures a large amount of variance in higher-order properties of human data.

## Study 2

In a second study, participants played a random *location* game. In this task, they were instructed to iteratively click locations on a line while trying to be as random as possible. Even though this task involved continuous rather than discrete responses, the task is analogous to the random number task in that participants generated random information along a single dimension. Because of this, we can conduct the same analyses across the two tasks to explore consistent patterns of behavior; in other words, we can ask whether a similar model may explain behavior in both cases.

### Method

200 additional participants completed this task. The task was identical to Study 1 except that, instead of selecting random numbers, participants clicked locations on a one-dimensional line. The line was positioned in the center of the screen. It extended 800 pixels horizontally. When participants clicked near the line, a vertical blue line appeared, centered on the line, at the horizontal position of the mouse (the vertical position of the mouse did not affect the vertical line position). The line itself was 20 pixels thick. However, the functional region in which participants could click extended 200 pixels vertically. This was done to minimize errant clicks. Similar to Study 1, a prompt indicated that participants should click randomly anytime the vertical blue line was not visible. Once a click was made, the blue line would remain visible at the clicked location for 750ms before disappearing.

Of the original sample of 200 participants, only 82 survived both of our exclusion criteria. This conservative sample is substantially smaller than that of Study 1 (see *General Discussion*). We again emphasize that excluding data points like these makes the model performance more conservative.

### Results & Discussion

To make the results of this study comparable to those of Study 1, we converted all locations into nine discrete values based on evenly spaced bins. For example, a response on the far left of the line would be coded as a 1, a response on the far right would be coded as a 9, and a response in the middle would be coded as a 5. The exclusions mentioned above (see *Method*) were calculated with respect to these discretized values, as are all subsequent analyses. A matrix depicting the overall pattern of results can be seen in Figure 2A.

**Model comparisons** We analyzed these data using the same three models we had used in the previous study. Each model was fit to each participant's observed choices using maximum likelihood estimation. After fitting each single-parameter model, the single free parameter in the stay/switch model ($\varepsilon$) was found to be significantly negative (signrank test, $p<.001$), meaning that participants avoided repeating choices back-to-back. The single free parameter in the side-switch model ($\eta$) was not significantly different from zero (signrank test, $p=.29$). After fitting the combined model, the stay/switch parameter similarly came out significantly

negative (signrank test, $p<.001$), and, in contrast to Study 1, the side-switch parameter was significantly negative (signrank test, $p<.001$).

The combined model had the lowest (best) total AIC score (89299), followed by the stay/switch model (89510) then followed by the side-switch model (89879). This suggests that, like the random number task, both parameters play a role in people's randomization.

**Model validation** The average correlation for the *average numerical distance* measure between the human data and randomly shuffled human data was $r=.34$. For the combined model, the average correlation between model-generated and human data was $r=.70$ (see Figure 2B for an example of the model fits; see Figure 2C for the correlation comparisons across the 1000 simulations). The average correlation for the *direction switches* measure between the model data and randomly shuffled human data was $r=-.01$. For the combined model, the average correlation between model-generated and human data was $r=.49$ (see Figures 2D and 2E).

## Study 3

In a final study, participants played a random, two-dimensional location game. In this task, they were instructed to iteratively click locations within a square region while trying to be as random as possible. Here, we ask whether the same model of randomness based on one-dimensional judgments can usefully predicts judgments in two dimensions.

### Method

200 additional participants completed this task. The task was identical to Study 2 except that, instead of selecting random locations in a one-dimensional space, participants clicked locations within a two-dimensional bounding square. The line was positioned in the center of the screen. It extended 600 pixels horizontally and vertically. When participants clicked anywhere in that 600 x 600 region, a dot appeared at the location of the cursor. The dot was 10 pixels in diameter. Similar to Studies 1 and 2, a prompt indicated that participants should click randomly anytime the dot was not visible. Once a click was made, the dot would remain visible at the clicked location for 750ms before disappearing.

Of the original sample of 200 participants, 105 survived exclusion in the x-dimension, 69 survived exclusion in the y-dimension, and 65 survived exclusion in both dimensions. We analyze data separately in each dimension with all participants who survived exclusion in that dimension. These conservative samples are substantially smaller than that of Study 1, but comparable to that of Study 2. We again emphasize that had we not excluded participants, the model would have performed *even better* (but without necessarily capturing anything meaningful about human randomness).

### Results & Discussion

We analyzed participants responses separately for each dimension (i.e., x vs. y). As in Study 2, we converted all
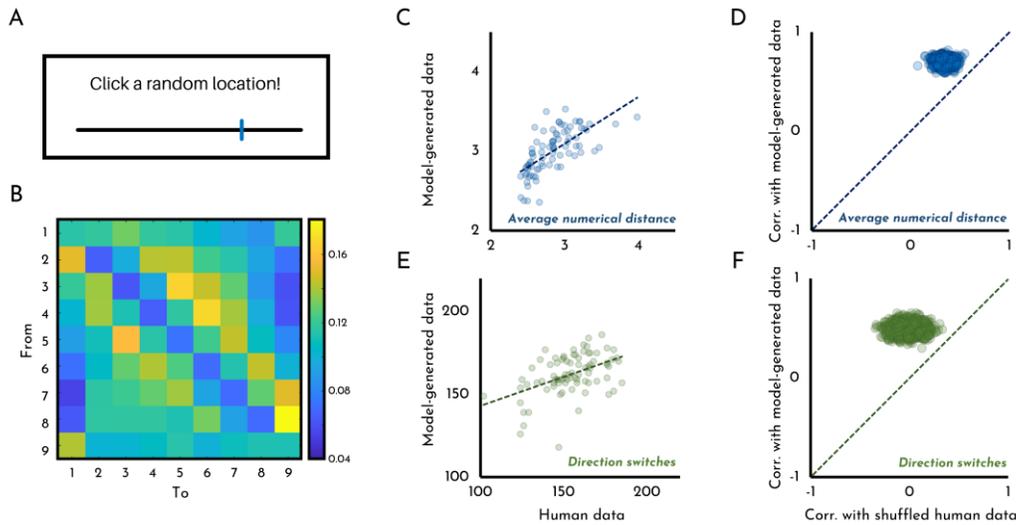
Figure 2: Results from Study 2. Here, continuous responses along a line were discretized into the values 1-9, evenly distributed across the line. (A) Transition probability from each value to each other value. (B) Correlation between the *average numerical distance* and (D) the *direction switches* in the human data and model-generated data. (C) Correlations from 1000 simulations between human data and shuffled human data (x-axis) as well as correlations from 1000 simulations between human data and model-generated data (y-axis), for the *average numerical distance* metric and (E) the *direction switches* metric. Points above the identity line indicate better model performance.

locations into nine discrete values in each dimension based on evenly spaced bins. For example, a response in the most top-left corner of the square would be coded as a 1 in both the x and y dimensions, whereas a response in the bottom-right corner would be coded as a 9 in both dimensions. The exclusions mentioned above (see *Method*) were calculated with respect to these discretized values, as are all subsequent analyses.

**Model comparisons** We separately analyzed behavior in each dimension using the same analyses as the previous two experiments. For the x-dimension, in the combined model, both the stay/switch parameter (signrank test, *p*<.001) and the side-switch parameter came out significantly positive (signrank test, *p*<.001). In other words, participants were *more* likely to stay near their previous response, but also more likely to switch sides of the distribution. The tendency to stick with the previous response is mostly driven by responses on the far-left side and far-right side of the distribution. For the y-dimension, in the combined model, the stay/switch parameter was significantly positive (signrank test, *p*=.002) and the side-switch parameter was positive, but not significant (signrank test, *p*=.15).

For the x-dimension, the combined model had the lowest (best) total AIC score (114600), followed by the stay/switch model (114880) then the side-switch model (115140). For the y-dimension, the combined model had the lowest (best) total AIC score (75380), followed by the side-switch model (75503) then the stay/switch model (75652).

**Model validation** We again compared how well model-generated data captured higher order properties (e.g., *average numerical distance* and *direction switches*) of human data.

First, we looked at x-dimension. The average correlation for the *average numerical distance* measure between the human data and randomly shuffled human data was *r*=.39. For the combined model, in comparison, the average correlation between model-generated and human data was *r*=.70. The average correlation for the *direct switches* measure between the human data and randomly shuffled human data was *r*=-.03. For the combined model, the average correlation between model-generated and human data was *r*=.66.

Then we did the same analysis for the y-dimension. The average correlation for the *average numerical distance* measure between the human data and randomly shuffled human data was *r*=.35. For the combined model, in comparison, the average correlation between model-generated and human data was *r*=.67. The average correlation for the *direct switches* measure between the human data and randomly shuffled human data was *r*=-.02. For the combined model, the average correlation between model-generated and human data was *r*=.57. For both the x- and y-dimensions, the model is a significantly better predictor of human behavior than even the same, shuffled human data.

This study demonstrates that human random behavior in a two-dimensional task can be at least partially explained by their behavior in each separate dimension. This conclusion is far from obvious: It could have been that performance in one dimension was captured by the model but not the other, or that neither dimension was well-captured by the model. That the model was able to predict behavior in both dimensions
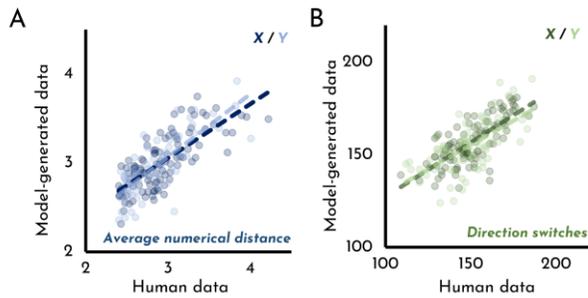
Figure 3: Results from Study 3. As in Study 2, continuous responses in each dimension of a square were discretized into the values 1-9, evenly distributed in space. (A) Correlation between the *average numerical distance* in the human data and (a representative simulation of) model-generated data. (B) Correlation between the *direction switches* in the human data and (a representative simulation of) model-generated data. The plots here depict both the x- and y-dimensions.

indicates that (a) the generative model that people use to generate random locations in two-dimensional space may act on each dimension (at least partially) independently, and (b) that generative model may strongly resemble the generative model used to generate not only random locations in a separate one-dimensional space, but also random numbers. More generally, it may be surprising that this model works at all when applied to the spatial domain given that we arbitrarily use the same number of bins (9) in both the spatial and number tasks to allow comparison to the numerical task. For some participants, using fewer bins might provide a more parsimonious description of behavior.

## General Discussion

Here we have proposed a model of human randomization that generalizes across numerical and spatial randomization tasks, as well as across one-dimensional and two-dimensional tasks. The ability of the model to capture behavior in all of these cases suggests that there are indeed stable aspects of human randomization behavior that generalize across tasks. Although some previous work has shown that some biases are stable across contexts (see, e.g., Yu et al., 2018), the present work goes one step further: it establishes a set of concrete, computational principles that can be used to generate sequences that reflect individual biases and that may generalize across a broader range of tasks — including those that involve fundamentally continuous (rather than discrete) input.

### A domain-general randomness generator?

Here, we have tested three cases of human randomization behavior: random number generation (on a limited, discrete set of values, 1-9), random one-dimensional location generation, and two-dimensional location generation. We have shown that there are computational heuristics that generalize across all three tasks. However, this is only the tip of the iceberg. There are of course many other ways that

random information can manifest in the world. For example, randomization need not be limited to options along one or two dimensions. Certain forms of abstract art, for instance, may play on regularities along numerous spatial dimensions as well as other visual dimensions (e.g., color). The spatial distribution of crowds may tell us something about where that crowd is headed, or what its goals are. The spatial distribution of leaves on the forest floor may tell us whether they had been tampered with and whether some other entity may be nearby. In each of these cases, there are multiple kinds of visual information that come together to form a cohesive percept of randomness. If the goal is to understand human randomization *in general*, then we must ask how well these computational principles generalize to more complex cases (i.e., those involving multiple dimensions of one type, or a combination of dimensions of different types).

### A mental number line account

We have thus far emphasized how the proposed model(s) generalize across two very different tasks. However, one possibility is that the similarity observed across Study 1 and Study 2 is due to a deep similarity between the two tasks. In other words: Even though one task involved discrete choices (Study 1) and the other involved continuous choices (Study 2), there is considerable evidence that humans represent numerical magnitudes along a sort of 'mental number line' (see, e.g., Aulet et al., 2021; Dehaene et al., 1993; Zorzi et al., 2002). If true, this could mean that participants were 'co-opting' spatial cognition in order to produce random behavior in the number task, even though the task is not intrinsically spatial.

### Conclusion

Does the mind contain a domain-general 'randomization machine'? Answering this question would require generalizing the present results to a variety of other tasks. However, the model proposed here promises one way forward. We have shown that a simple model can capture higher-order properties of human randomization behavior even after conservatively excluding highly non-random individuals. We believe that this approach is 'scalable' in the sense that it can be applied to a range of tasks and stimuli, thus providing a way to understand the general computational principles underlying human randomization.

## References

Aulet, L. S., Yousif, S. R., & Lourenco, S. F. (2021). Spatial–numerical associations from a novel paradigm support the mental number line account. *Quarterly Journal of Experimental Psychology*, 17470218211008733.

Dehaene, S., Bossini, S., & Giraux, P. (1993). The mental representation of parity and number magnitude. *Journal of Experimental Psychology: General*, *122*, 371-396.

Kahneman, D., & Tversky, A. (1972). Subjective probability: A judgment of representativeness. *Cognitive Psychology*, *3*, 430–454.

Kelly, M. H., & Martin, S. (1994). Domain-general abilities applied to domain-specific tasks: Sensitivity to probabilities in perception, cognition, and language. *Lingua*, *92*, 105-140.

Reiner, C., Yousif, S., Sherman, B., & Keil, F. (2021). Common structure underlying visual and non-visual judgments of randomness. *Journal of Vision*, *21*, 2254-2254.

Rescorla, R. A., & Wagner, A. R. (1972). A theory of Pavlovian conditioning: Variations in the effectiveness of reinforcement and nonreinforcement. In A. H. Black & W. F. Prokasy (Eds.), Classical conditioning II: Current theory and research (pp. 64–99). New York, NY: AppletonCentury-Crofts.

Reuter, J., Raedler, T., Rose, M., Hand, I., Galscher, J., & Buchel, C. (2005). Pathological gambling is linked to reduced activation of the mesolimbic reward system. *Nature Neuroscience*, *8*, 147–148.

Sherman, B. E., Graves, K. N., & Turk-Browne, N. B. (2020). The prevalence and importance of statistical learning in human cognition and behavior. *Current Opinion in Behavioral Sciences*, *32*, 15-20.

Towse, J. N., & Neil, D. (1998). Analyzing human random generation behavior: A review of methods used and a computer program for describing performance. *Behavior Research Methods, Instruments, & Computers*, *30*, 583-591.

Wagenaar, W. A. (1988). Paradoxes of gambling behaviour. Hillsdale, NJ: Lawrence Erlbaum.

Yu, R. Q., Gunn, J., Osherson, D., & Zhao, J. (2018). The consistency of the subjective concept of randomness. *Quarterly Journal of Experimental Psychology, 71*, 906-916.

Zorzi, M., Priftis, K., & Umiltà, C. (2002). Brain damage: Neglect disrupts the mental number line. *Nature*, *417*, 138-139.