

1 Approximating Cognitive Representations Using Space

2 John R. Starr (jrs673@cornell.edu), Ashlyn Winship,
3 & Marten van Schijndel^a

^aDepartment of Linguistics, Cornell University

4 **Abstract**

In everyday life, people intuitively use space to make meaningful distinctions between objects. In this paper, we present a novel, free-to-use on-line experimental paradigm that capitalizes on these intuitions: GRIS (Generating Representations in Space). In GRIS experiments, participants manipulate a set of objects (text, audio, images) and place them on canvases. Following an introduction to the paradigm, we present three studies which demonstrate how experiments in the GRIS paradigm can both a) replicate prior psycholinguistic results and b) reveal nuanced insights about human and computational representations.

5 *Keywords:* representations, psycholinguistics, acceptability, typicality,
6 paradigm

7 **1. Introduction**

8 In our daily lives, we use space to make and represent meaningful rela-
9 tionships between objects: we separate different kinds of clothes into different
10 compartments, read menus that spatially group items on the page according
11 to their broader classifications, and press elevator buttons that are verti-
12 cally ordered to reflect the structure of their buildings. In these ways and
13 many more, humans intuitively use space to simplify choice, perception, and
14 computation (Kirsh, 1995), allowing us to navigate and represent complex
15 structures and relations with ease.

16 One of the primary approaches to studying space in the cognitive sci-
17 ences is through its relationship with language. Previous research has shown
18 that people construct mental representations that encode spatial relation-
19 ships (Taylor and Tversky, 1992; Bryant, 1997; Kemmerer, 1999), and that
20 these relationships are marked on a schematic level of varying detail (Talmy,

1983; Landau and Jackendoff, 1993; Hayward and Tarr, 1995; Tversky and Lee, 1998). Studies in this domain often focus on how language organizes our cognitive representations of objects and their locations, both in discourse and in the real world. Other work on the relationship between space and language suggests that we transfer linguistic information onto mental spaces of the world (consisting of information of referents, their beliefs, actions, etc.), which we then blend together to understand the relevant discourse (Fauconnier, 1994; Sweetser, 1999; Fauconnier et al., 2007).

In this paper, we demonstrate that an alternative perspective on the relationship between space and language is also fruitful: space as a tool to contextualize our understanding of language, and, more broadly, human cognition. To show the utility of this alternative perspective, we present an experimental paradigm – GRIS (**G**enerating **R**epresentations **I**n **S**pace) – which a) capitalizes on the way humans intuitively use space, b) approximates representations of language and other cognitive phenomena, and c) does so in a way that is easily comparable to embedding representations from computational models, allowing us to further probe the matches and mismatches between humans and models. At a high level, participants in GRIS experiments can move objects (text, image, audio) onto a canvas and use space in a meaningful way, where information is incrementally collected about which objects were moved, when they were moved, and where they moved to. To briefly highlight our results, we demonstrate how GRIS experiments can 1) both replicate results from other psycholinguistic paradigms and provide further contextual nuance to such results, 2) develop multi-dimensional graphs that can be used for computational modeling, and 3) facilitate and simplify experimental designs that require multiple complex (pairwise) comparisons. More broadly, we argue that GRIS allows participants to use their natural intuitions about space to inform our understanding of how people represent various kinds of linguistic information.

1.1. Article Organization

In the following section, we provide further motivation for developing a paradigm that uses space meaningfully. In section 3, we outline the GRIS paradigm, introducing its key functionalities and structure. In sections 4-6, we present three GRIS experiments,¹ demonstrating how the paradigm can a)

¹All items, data, and analysis code can be found at the following anonymized link: https://osf.io/94gck/?view_only=1ed03a3757fe44ba9a036510be60b7c6.

55 replicate prior results across a number of cognitive domains, and b) capture
 56 more nuanced relationships between representations than other experimental
 57 paradigms and computational models of linguistic structure.² In section 7,
 58 we discuss the general implications of GRIS and present possible directions
 59 for future work. In section 8, we conclude.

60 2. On Space

61 2.1. Space & Psycholinguistic Paradigms

62 From a design perspective, many standard psycholinguistic paradigms³
 63 minimally engage with space: standard rating and judgment tasks often
 64 present an item in isolation (or near isolation) alongside a scale or a drop-
 65 down box, and forced-choice tasks only capture the pairwise differences be-
 66 tween one or two items. Some experimental paradigms do inherently use
 67 space as a metric for psycholinguistic effort, such as measuring how partici-
 68 pants’ eyes move to different locations on a screen when using eye-tracking in
 69 the visual world paradigm (Cooper, 1974; Tanenhaus et al., 1995), or follow-
 70 ing the trajectory of a participant’s mouse/cursor across the screen using a
 71 mouse-tracking paradigm (Freeman and Ambady, 2010; Wilcox et al., 2024)
 72 However, these paradigms do not fully capitalize on the possible utilities of
 73 space: the locations of objects in the visual world paradigm are often op-
 74 timized to be distinct and do not carry inherent meaning themselves, and
 75 mouse-tracking uses space as a proxy for processing difficulty instead of as a
 76 representational, organizational mechanism.

77 We propose that the use of space can simplify psycholinguistic designs
 78 for both researchers and participants. As an example, consider a rating task
 79 where a participant is asked to rate the difference between item pairs on
 80 a scale (according to some metric), for a total of four unique items. We
 81 visualize two possible iterations of this experiment in Figure 1.

82 In Figure 1A, participants are asked to directly measure the difference be-
 83 tween all possible combinations of the four items (either in one trial or across

²In this paper, we focus on *linguistic* representations, though GRIS can be easily ex-
 tended to approximate other kinds of cognitive structures and relationships such as in
 vision or acoustics.

³As will be described later in this article, GRIS is not designed to capture on-line
 processing. Accordingly, we do not elaborate on the use of space in psycholinguistic tasks
 such as self-paced reading (Just et al., 1982) or (Forster et al., 2009).

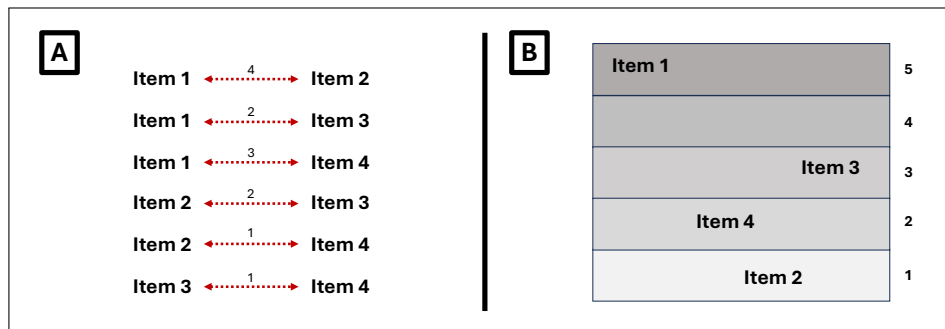


Figure 1: Sample rating tasks for four items. **(A)** Ratings in a pairwise comparison layout, where the number above the arrows reflects the scalar differences between two items. **(B)** Ratings in a space-motivated layout. The overall scalar ratings are identical between (A) and (B).

several trials). For these items, participants maintain six pairwise comparisons that are all in relation to one another. Alternatively, the experiment in Figure 1B present a version of the experiment that better capitalizes on human spatial intuitions: participants are asked to use space to distinguish between all possible combinations, where larger separation between items reflects larger differences. Note that the absolute value of the ratings are identical between both iterations of the experiment.

While both versions require the participant to make the same number of (underlying) pairwise comparisons, we offer that experiment (B) is more informative than experiment (A), for a number of reasons. First, participants are able to concretize the relative relations. Rather than needing to maintain an implicit scale of differences in experiment (A) – which may lead to possible inconsistencies as the number of comparisons increases – the relations between items are explicitly visualized in a manner that is easy to manipulate. Second, the directionality of differences is transparently coded: while Item 4 is one away from both Item 2 and Item 3, experiment (B) easily captures the direction of the effects, whereas experiment (A) does not. Third, experiment (B) contextualizes the different items and their relative relations, allowing participants to quickly set the bounds of the underlying scale(s) that they are using to distinguish between items.

2.2. Space & Computational Models

Vector spaces that are generated by modern computational models of language are often used as proxies for human linguistic structure: for example,

high-dimensional vectors for the words “cat” and “dog” are typically near one another in computational vector spaces, whether such vectors are computed using word co-occurrence statistics (e.g., Pennington et al., 2014) or more complex, contextual operations (e.g., Radford et al., 2019). Accordingly, close proximity in computational vector spaces⁴ is often interpreted as human-like similarity,⁵ at all levels of linguistic structure, including phonetic information (Parrish, 2017; Zouhar et al., 2023), phonological segments (Silfverberg et al., 2018), phrases (Passos et al., 2014), and others.

However, while this interpretation about the relationship between human and model representations holds true generally, prior work has noted some mismatches: for example, model representations have been shown to occupy a narrow region of the embedding space (a phenomenon known as *anisotropy*; Mimno and Thompson, 2017; Ethayarajh, 2019), have “rogue” dimensions that dominate similarity metrics (Timkey and van Schijndel, 2021), and fail to be robust to minor orthographic noise (Matthews et al., 2024).⁶ Moreover, no current psycholinguistic methodology – to our knowledge – approximates human representational spaces of linguistic structure in a manner that is comparable to those generated by off-the-shelf computational models, making it difficult to align human and model representations.

2.3. *Desiderata for GRIS*

Given this overview of psycholinguistic paradigms and computational representations, we present the fundamental motivations behind GRIS:

1. A flexible experimental paradigm that uses space to construct meaningful, interpretable relations between objects.
2. A tool that allows researchers to quickly build experiments.
3. An experimental interface that participants can intuitively use.
4. Output data that approximates human cognitive representations that are easily aligned to representations from computational models.

⁴We use “proximity” as a catch-all term for similarity in the vector space, given that there are a variety of metrics – Euclidean distance, cosine similarity, Spearman’s ρ , etc. – to determine representational similarity.

⁵See Apidianaki (2023) for an overview.

⁶Some research in activation & representation engineering (e.g., Turner et al., 2023; Wu et al., 2024) demonstrates how these representations can be fine-tuned to perform better on down-stream tasks; we do not discuss these approaches in detail, though we do address them in the discussion.

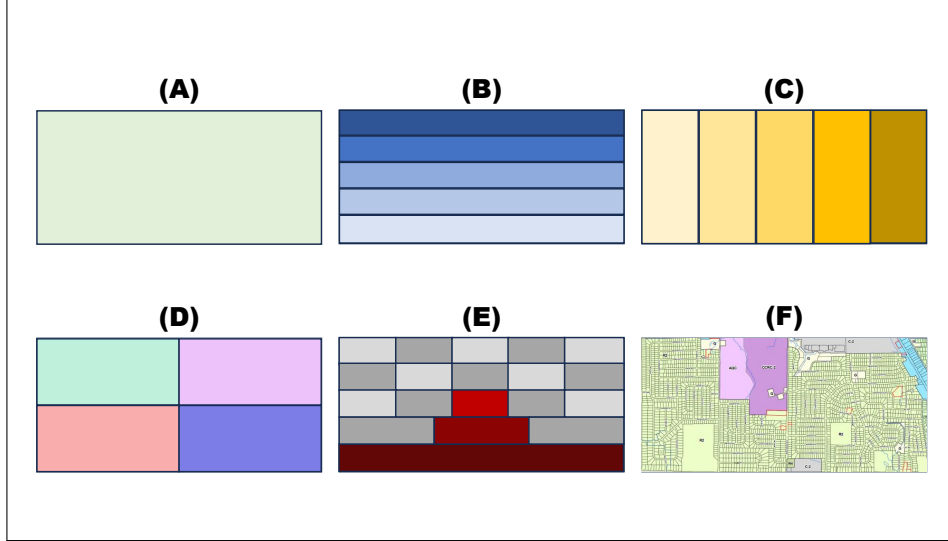


Figure 2: Sample GRIS canvases. Canvases can be blank (A), split into categories in both cartesian dimensions (B, C) simultaneously (D) and irregularly (E), or placed under an image (F).

135 3. GRIS: A Walkthrough

136 The core idea of the GRIS paradigm is to provide participants with ob-
 137 jects that they can drag and drop onto a labeled canvas. In the following two
 138 subsections, we overview the structure of a GRIS experiment and demon-
 139 strate how participants navigate through a trial.

140 3.1. Structuring a GRIS Experiment

141 GRIS experiments have two fundamental components: 1) objects that
 142 can be placed, and 2) a canvas to place objects on.

143 Objects can be text, images, or audio; these objects are distinct targets
 144 that can be individually moved. By default, objects are located in a reservoir
 145 at the bottom of the screen, though their initial positions can be changed to
 146 accommodate relevant research questions.

147 Canvases, from a participant’s perspective, can be either blank or split
 148 into different categories; some sample canvases are provided in Figure 2.
 149 From a researcher’s perspective, canvases are built of individual, labeled can-
 150 vas blocks that are either square or rectangular. By default, canvas blocks
 151 are labeled using a four-point coordinate system ($x-cat$, $y-cat$, $x-abs$, $y-abs$),
 152 where the first two dimensions are used to mark the category that the block

153 belongs to, and the last two dimensions are used to mark the absolute position
154 of the block on the overall canvas; note that canvas labels can be modified
155 to accommodate other systems. Beyond labeling, each canvas block can be
156 independently specified for height, width, and color. Finally, images can be
157 overlaid on the canvas, allowing for additional designs beyond those possible
158 by combinations of squares and rectangles.

159 For ease of use for other researchers, we have developed the GRIS toolkit,⁷
160 which provides instructions on how to build, run, and analyze GRIS experi-
161 ments.

162 3.2. *Participating in a GRIS Experiment*

163 GRIS experiments are designed to be simple and intuitive for partici-
164 pants. To explain how a participant navigates through a GRIS experiment,
165 we provide a sample, partially-completed trial in Figure 3. In this sample
166 trial, the participant has access to five objects – different shapes – which
167 begin in the reservoir (B) below the blank canvas (A). The instructions at
168 the top of the screen indicate that the participant should order the objects
169 in a line, where the leftmost shapes are the “roundest” according to their
170 intuitions. The participant first placed the star on the right boundary of
171 the screen, then placed the oval on the left boundary; the abnormal shape
172 was placed between these two shapes. Once they have placed all five ob-
173 jects on the canvas, the participant will be prompted to continue to the next
174 trial, though they can continue to re-arrange the objects at any point in time
175 throughout the trial.

176 For each drag-and-drop, GRIS collects 1) which object was moved, 2)
177 the object’s original location, 3) its new location, 4) and the timestamps for
178 both the initial drag and the final drop. Data are also collected about when
179 each trial begins and ends, as well as the final positions for all objects at the
180 conclusion of each trial.

181 3.3. *Interim Summary*

182 In summary, GRIS is a simple – yet flexible – experimental paradigm
183 that can accommodate a wide variety of research questions and designs. To

⁷The GRIS toolkit is publicly-available on GitHub at the following link: <https://anonymous.4open.science/r/gris-toolkit-demo-923F/>. Currently, GRIS experiments are run on PC Ibex (Zehr and Schwarz, 2018), a free, on-line research platform intended for experiments in psycholinguistics and cognitive science.

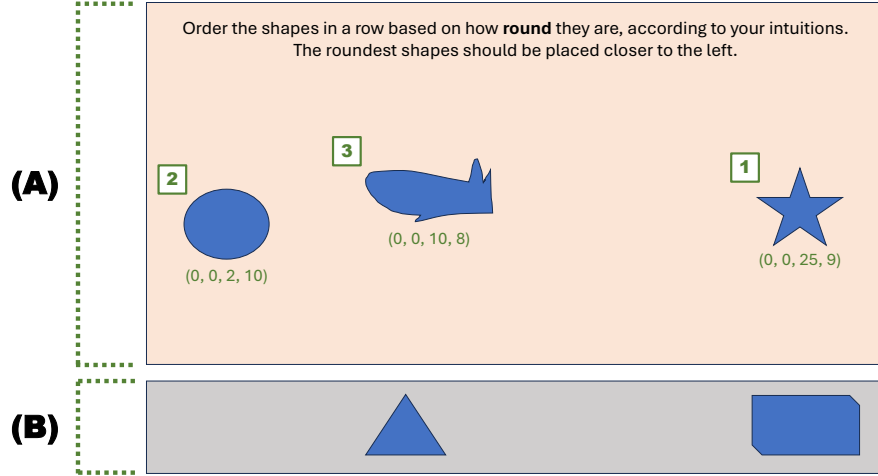


Figure 3: Sample of a partially-completed GRIS trial. Participants see a canvas (A) and a reservoir of objects (B); participants do not see anything marked in green. In this example, the objects are shapes. Participants are presented instructions – either during the trial (as in this figure) or prior to the trial – which guide how the participant should manipulate the objects. For each recorded drag-and-drop action, data is collected about the time/order of that action (boxed red numbers in the figure) and both the original and new coordinates for that object; only the dropped coordinates are presented in the figure.

184 validate GRIS’ effectiveness and demonstrate the kinds of analyses that it
 185 permits, we present a series of three experiments which capitalize on many
 186 of the features offered by the paradigm.

187 4. Experiment 1: Sentence Acceptability

188 Acceptability judgments probe what structures are (un)acceptable in a
 189 language: these structures can range from low-level judgments of phonolog-
 190 ical structure to high-level judgments of multi-sentence, multi-speaker dis-
 191 courses. In this section, we focus on sentence acceptability judgments in
 192 English.

193 For explanatory purposes, consider the sentences in (1)-(3). We adopt
 194 standard conventions for marking degrees of acceptability and grammatical-
 195 ity from linguistic research, where * indicates a sentence is ungrammaticality,
 196 and # indicates a sentence is odd or slightly marked.

197 (1) *An girls is hungry.

198 (2) Randy wanted to write a novel.

199 (3) #?Want to write, Randy did a novel.

200 While ungrammatical sentences like (1) are rated toward the boundaries
201 of the acceptability spectrum, others display more gradient judgments: for
202 example, sentence (2) is often preferred over sentence (3), even though both
203 are grammatical sentences of English. Previous research primarily collects
204 acceptability judgments using Likert scales (Gibson et al., 2011), forced-
205 choice tasks (Mahowald et al., 2016), or response times (Konieczny, 2000).
206 In isolation, such sentence acceptability judgments appear to be robust across
207 experimental paradigms, suggesting that people have consistent preferences
208 about the internal structure of their language (Sprouse, 2011; Sprouse et al.,
209 2013). However, these measures do not always capture the relative relation-
210 ship of sentence acceptability across structures. For example, people express
211 consistent preferences: generally speaking, $(2) > (3) > (1)$. But, each of these
212 pairwise preferences reflects a different underlying scale: while (1) is less ac-
213 ceptable than (3) and (3) is less acceptable than (2), the former distinction
214 is motivated by differences in grammaticality, while the latter distinction is
215 motivated by differences in frequency and syntactic complexity.

216 Moreover, isolated syntactic judgments may also conflate degrees of ac-
217 ceptability: a rating of 3 for one construction may not be comparable to
218 a rating of 3 for another construction, even though the ratings are identi-
219 cal. Capturing the contextual organization of syntactic acceptability across
220 phenomena would help us understand the broader organization of human
221 language understanding and cognition.

222 In this study, we use GRIS to replicate large-scale sentence acceptability
223 judgments from prior work, while also showing how the acceptability differ-
224 ence between sentence pairs can strongly vary depending on the context that
225 they appear in.

226 4.1. Design & Procedure

227 4.1.1. Stimuli

228 All stimuli were drawn from Sprouse et al. (2013), which randomly sam-
229 pled informal (i.e. not experimentally-tested) acceptability judgments of En-
230 glish sentence pairs from *Linguistic Inquiry*, a well-established journal in
231 theoretical linguistics. After sampling these sentence pairs, Sprouse et al.
232 (2013) collected acceptability ratings for each sentence within each pair to

233 test whether the informal judgments were valid for larger populations; we
234 will use these ratings to confirm that our findings correlate with prior work.

235 We sampled 72 pairs from the Sprouse et al. (2013) dataset. All 72 sen-
236 tence pairs were classified according to the general linguistic phenomenon
237 that their original paper tested; these classifications were drawn from the
238 abstracts of the papers themselves. By labeling the linguistic phenomenon
239 that each pair tests, we can then combine pairs of different classifications to
240 understand how different syntactic phenomena influence sentence acceptabil-
241 ity across structures, allowing us to obtain a broader understanding of the
242 organizational preferences of acceptability judgments. Some sample classifi-
243 cations of phenomena are listed below in (4):

- 244 (4) a. WORD ORDER:
245 Fred mowed the green lawn. > Fred mowed the lawn green.⁸
246 b. DEFINITES:
247 This is a table. > This is table.

248 From this set of 72 sentence pairs, we randomly selected 24 sentence pairs
249 to serve as our target pairs: all participants saw each of these 24 sentence
250 pairs. To test the impact of context on making these acceptability judgments,
251 the remaining 48 items were broken into two sets of 24 sentences, each of
252 which was paired with the 24 example items so that each target pair could
253 appear in context with different phenomena. In sum, this process led to two
254 sets of 24 items with four sentences (two pairs) each.

255 4.1.2. Procedure

256 See Figure 4 for a sample trial for Experiment 1. Participants saw four
257 sentences below a gradiently-colored canvas, where the color gradient re-
258 flected a 5-point Likert scale. Participants were instructed to move the sen-
259 tences from the bottom of the screen onto the canvas according to how “ac-
260 ceptable” the sentences were, according to their intuitions. Participants were
261 told that the “most acceptable” sentences should be placed at the top of the
262 canvas (5, on a standard Likert scale), while the “least acceptable” sentences
263 should be placed at the bottom (1, on a standard Likert scale). They were

⁸While the example provided here does introduce a resultative construction, the pri-
mary arguments of the original paper discuss the construction’s implications on word
order.

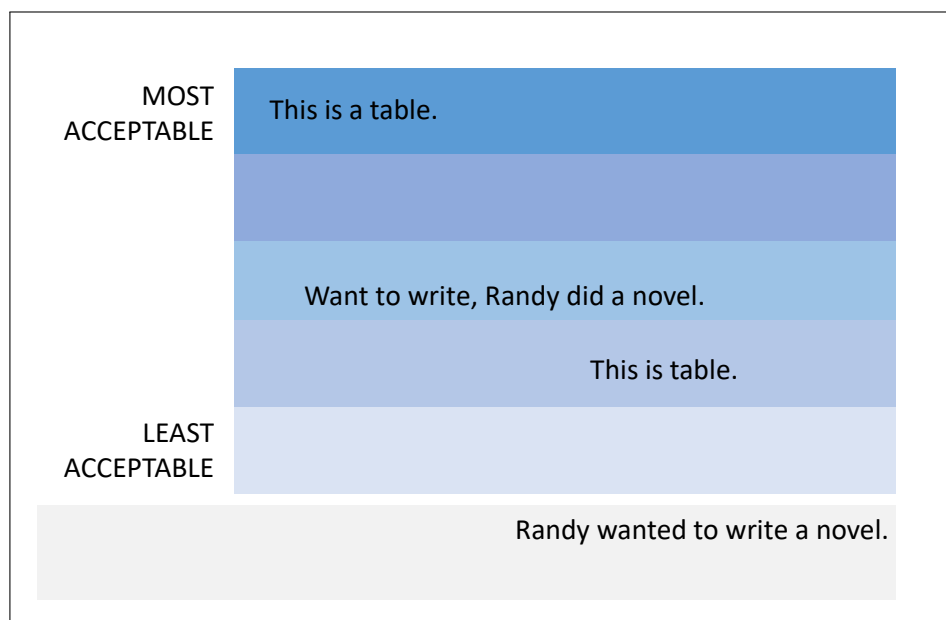


Figure 4: Sample trial for Experiment 1. Font has been enlarged for readability.

also told that multiple sentences could occupy the same level on the scale.
Sentence positions below the canvas were randomized for each item.

4.1.3. Participants

Twenty-five participants were recruited using the online research platform Prolific. Participants were all native speakers of English between the ages of 18 and 55.

4.2. Results

4.2.1. Base Acceptability

To measure sentence acceptability judgments within each trial, we collected the final position of all sentences once the trial was complete. We z -scored acceptability ratings by participant to ensure that responses were compared on similar scales.

Results for Experiment 1 are visualized in Figure 5. To test whether unacceptable sentences were rated significantly lower than acceptable ones, we fit a linear mixed-effects model to the z -scored acceptability rating, with a fixed effect of sentence TYPE (acceptable/unacceptable), and random in-

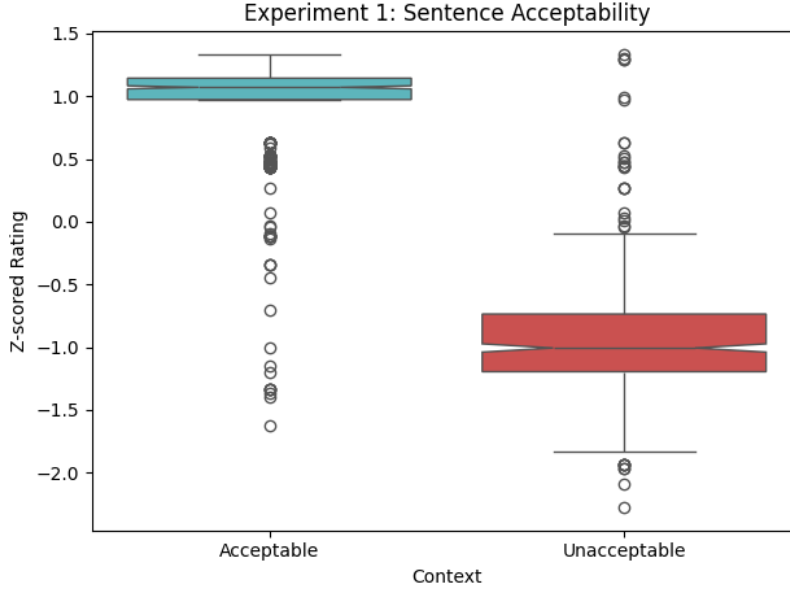


Figure 5: Base acceptability results for Experiment 1. Notches indicate 95% bootstrapped CIs.

tercepts for participants and items.⁹ Participants rated the UNACCEPTABLE sentences as significantly less acceptable than the ACCEPTABLE ones ($\hat{\beta} = -0.184$, $SE = 0.031$, $t = -58.80$, $p < 0.001$); these sentence ratings also strongly correlate ($r = 0.88$) with those found by Sprouse et al. (2013).

4.2.2. Contextual Acceptability

In addition to the basic acceptability analyses in the previous section, we measured how acceptability differences varied within each target pair according to the classification of the context pair that was present in the trial. To do so, we calculated the difference between each sentence in the target pair, then averaged the ratings within each context classification.

Results for contextual acceptability differences are shown in Figure 6. We find that some phenomena display similar levels of acceptability (< 0.4 Likert difference) regardless of context (e.g., *Agreement*, *Definites*), while others show significant variation (e.g., *Movement*, *Word Order*, *Clause*). For exam-

⁹The complete model formula was: Z-SCORED RATING \sim TYPE + (1 | item) + (1 | participant). The baseline was the “Acceptable” condition.

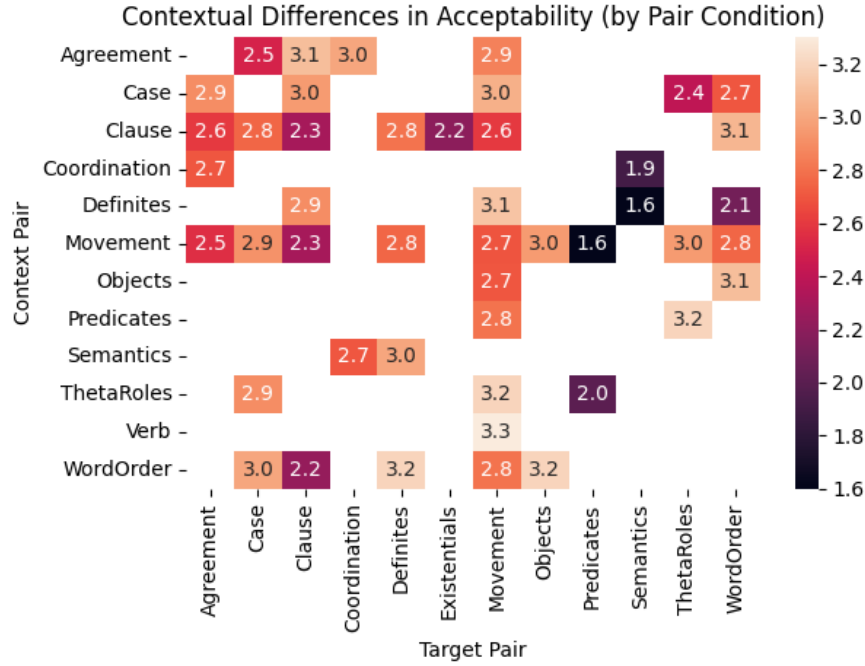


Figure 6: Contextual acceptability results for Experiment 1. X-axis represents the classification for the target pair. Y-axis represents the classification of the context pair. Cells indicate difference between acceptable and unacceptable sentences within each target pair; darker colors indicate smaller differences on a 5-point Likert scale.

294 ple, consider the *Word Order* classification for the target pair from (4-a):
 295 *Fred mowed the green lawn > Fred mowed the lawn green*. When placed in
 296 the context of a sentence pair that modulates *Definites*, the difference be-
 297 tween the *green lawn* and *lawn green* sentences was approximately 2.1 on
 298 a 5-point Likert scale; but, when placed in the context of a sentence pair
 299 that modulates *Objecthood*, the difference between the *green lawn* and *lawn*
 300 *green* sentences was approximately 3.1. These varying differences have sig-
 301 nificant consequences on how researchers interpret acceptability judgments:
 302 a difference of ~ 3 points on a 5-point Likert scale easily distinguishes an
 303 acceptable sentence (5) from an unacceptable one (2), whereas a difference
 304 of ~ 2 points could be the distinction between a totally acceptable sentence
 305 (5) and a moderately acceptable one (3).

306 4.3. Discussion

307 The results of this task show that GRIS can be used to reliably repli-
308 cate prior experimental results involving pairwise comparisons, while also
309 systematically capturing the variability of sentence acceptability in different
310 contexts. More specifically, GRIS reveals how previous sentence acceptabil-
311 ity judgments in isolation may not serve as reliable representations of overall
312 sentence acceptability in context.

313 5. Experiment 2: Category Typicality

314 Category typicality assesses how “typical” an object is within a broader
315 category (Rosch, 1975; Farmer et al., 2006). For example, “robins” and “spar-
316 rows” are found to be more typical representations of birds than “toucans”
317 and “penguins” across cognitive domains, including language (Rosch, 1975;
318 Meints et al., 1999) and vision (Maxfield et al., 2014). Traditionally, category
319 typicality has been measured using rating or decision tasks (Rosch, 1975),
320 production tasks (Rosch et al., 1976), or inductive-reasoning tasks (Osherson
321 et al., 1990), all of which ask the participant to consider a specific word in
322 relation to the broader category label. Recent computational work also sug-
323 gests that computational models of language may learn some aspects of cat-
324 egory typicality from the statistical usage distributions of everyday language
325 (Misra et al., 2021), though these analyses focus on probability estimates
326 from pre-trained language models rather than representational analyses.

327 In this experiment, we build a typicality-rating experiment using GRIS,
328 finding that manipulating words in space both 1) replicates previous category
329 typicality effects and 2) allows us to directly compare representational spaces
330 between humans and models.

331 5.1. Design & Procedure

332 5.1.1. Stimuli

333 We used eight of the original ten categories from Rosch (1975): *fruits*,
334 *vehicles*, *weapons*, *vegetables*, *tools*, *birds*, *sports*, *clothing*. All items were in
335 English. Each category has a list of approximately 50-60 words, where each
336 word has a typicality rating that was averaged across 209 subjects; we use
337 these ratings as our ground truth. To test whether the presence of different
338 words modified typicality ratings, we constructed eight items that used ten
339 words from each category; we did not use all of the words from Rosch (1975),
340 as there would be too many words for participants to move on the screen.

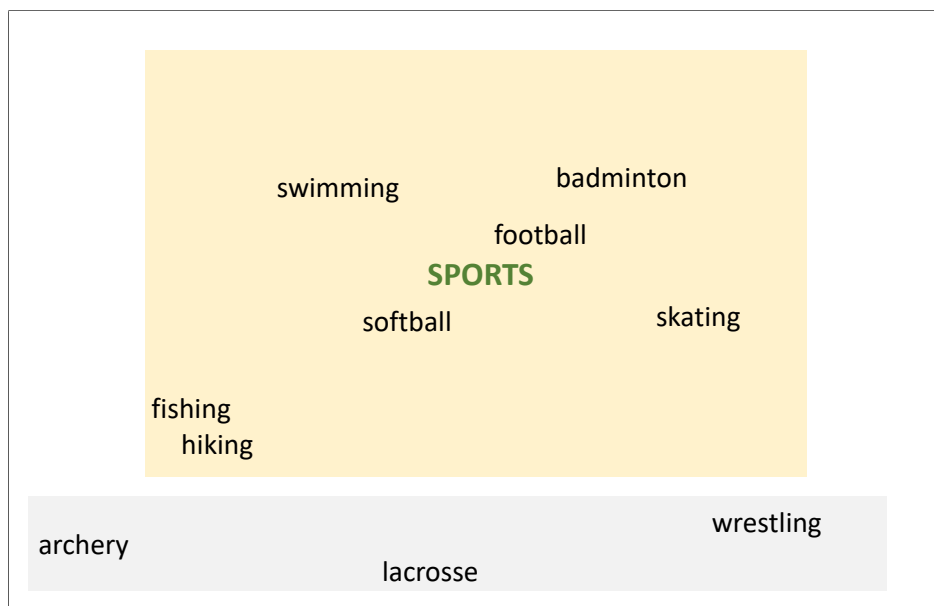


Figure 7: Sample trial for Experiment 2 (Typicality); font size enlarged to improve figure readability. Category label is marked in the center in green.

341 5.1.2. Procedure

342 A sample item for Experiment 2 is visualized in Figure 7. Participants
 343 saw a canvas with a word bank below. In the middle of the canvas was a
 344 bolded category label (i.e. **SPORTS**). Participants were told to move words
 345 from the bank onto the canvas according to how “typical” an example the
 346 word was of the category: words that were more typical examples of the
 347 category should be placed closer to the category label.

348 5.1.3. Participants

349 As in Experiment 1, twenty-five participants were recruited using the on-
 350 line research platform Prolific. Participants were all native English speakers
 351 between the ages of 18 and 55.

352 5.2. Results

353 As in Experiment 1, we collected the final positions for all words once
 354 the trial was complete. For each trial, we calculated every word’s distance
 355 from the center; we z -scored these distances by participant to ensure that all
 356 participants were comparable in how they used the space. Finally, following

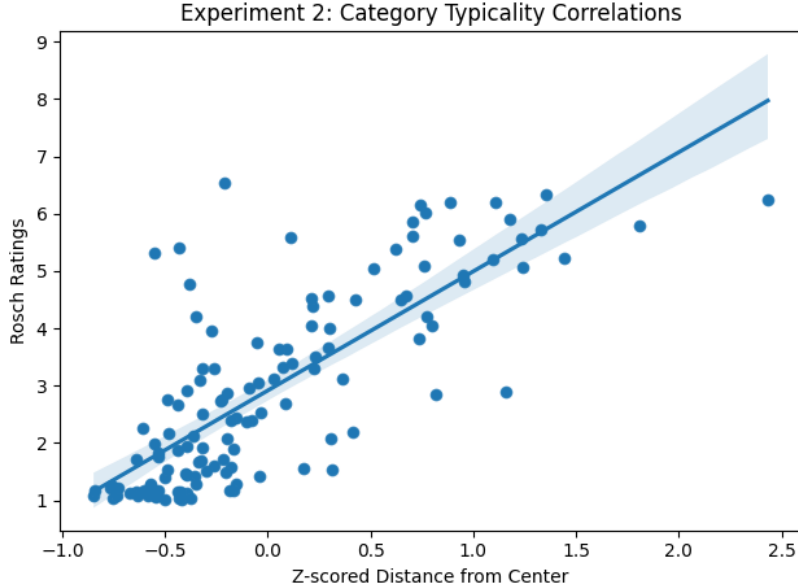


Figure 8: Correlation results for Experiment 2. X-axis indicates the Z -scored distance from center for a word. Y-axis indicates the original ratings from Rosch (1975).

the rating averaging from Rosch (1975), we meaned the distances for each word across participants.

Experimental results are visualized in Figure 8. We find a strong correlation ($r=0.78$) between the original rankings from Rosch (1975) and the distance of each word from its category label in our study, indicating that GRIS can be used to replicate prior category typicality results.

5.2.1. Computational Analyses

For our computational analyses, we extracted vector representations of words from three models: GLoVe 6B.300D (Pennington et al., 2014), BERT (Devlin, 2018), and GPT2 (Radford et al., 2019). For the non-contextual model (GLoVe), we gathered the raw vectors for both the word and the category label. Following Misra et al. (2021), for both of the contextual models (BERT & GPT2), we framed each word X with its category label Y in the following way: $A(n) X \text{ is a typical } Y.$; instead of gathering the probability of each word X in the sentence, we extracted the vector representations of both the word and the label using the `minicons` Python package (Misra, 2022). Approaching our computational analyses in this way allows us to

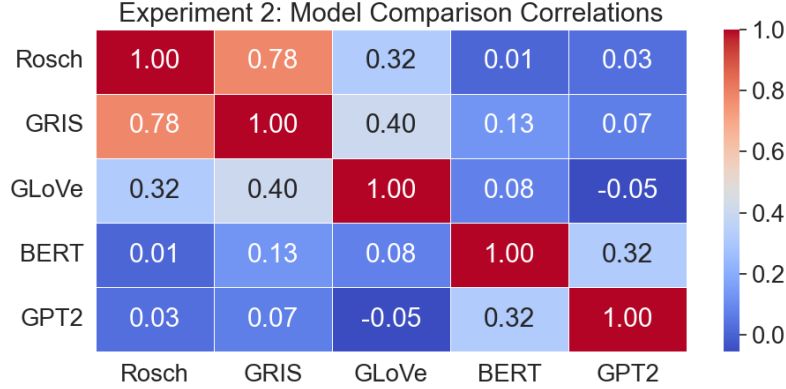


Figure 9: Correlation metrics between model representations and experimental results. Each cell corresponds to the Pearson’s correlation coefficient between the models and experimental measures on the x- and y-axes.

most directly compare the representational spaces constructed in the human experiment with those generated by computational models of language; our approach differs from that of Misra et al. (2021), in directly comparing model similarities to human similarity judgments rather than mapping model log-probabilities to human behavioral responses.

For each of the three models, we computed the Euclidean distance between the vectors for every word and its corresponding category label.¹⁰ We then calculated the Spearman’s correlation for all possible model comparisons.

Results for these multiple-correlation analyses are visualized in Figure 9. We find that GRIS is the only set of representations that connect a word to its category label in a manner that strongly correlates with the original rankings from Rosch (1975); the distances between words and their labels for GLoVe representations only weakly correlate with the original Rosch rankings, though there is a slightly stronger correlation between GloVe distances and our experimental data. We note that representational distances in BERT and GPT2 weakly correlate with one another, but fail to display any strong correlations with GLoVe or either set of experimental data. We also note that

¹⁰Analyses using standardized cosine similarity and Spearman’s rank correlation coefficient were also conducted; Euclidean distance performed best in the correlation analyses.

GRIS also has the highest average correlation coefficient across comparisons.

5.3. Discussion

In this experiment, we replicated prior typicality representations for eight categories. Experiments 1 and 2 show how GRIS can reliably replicate prior results; this experiment also demonstrates how GRIS builds constructs representational spaces more accurately than a number of well-established computational models. These findings differ from Misra et al. (2021), likely due to the fact that we are conducting *representational* analyses and not *behavioral* ones: while previous computational work has shown that behavioral measures moderately align with human behavior, our work demonstrates that studies of human representations cannot simply rely on vectors generated by these models.

6. Experiment 3: Multi-dimensional Similarity

In the previous two experiments, we demonstrated how GRIS can be used to both replicate and provide further detail about prior studies. In this experiment, we showcase how GRIS can be used to advance new questions within an established literature in cognitive science: pattern recognition.

For decades, cognitive scientists have studied how people recognize patterns across a variety of cognitive domains (Chater and Vitányi, 2003; Reed, 1972; Edelman, 1999; Edelman and Duvdevani-Bar, 1997). We contribute to this literature by examining how one form of pattern recognition – similarity assessments – arises during language processing.

Prior work suggests that the cognitive sources of similarity are a concept’s familiarity (strength in memory), association (relationships with other concepts), and inherent perceptual likeness (surface appearance); see Hiatt and Trafton (2017) for an overview. Linguistic similarity, broadly defined, has also been shown to influence pattern recognition. For example, semantic similarity is well-known to produce priming effects (McNamara, 2005; Neely et al., 1989; Shelton and Martin, 1992), and, while less studied, syntactic similarity has shown similar effects (Lester et al., 2017). Orthographic similarity improves recall accuracy in a probed serial-recall task (Lin et al., 2015), and phonological similarity has been shown to facilitate the learning of novel words (Papagno and Vallar, 1992).

While each of these features contributes to overall perception of similarity between linguistic units, how do people balance the multiple avenues of sim-

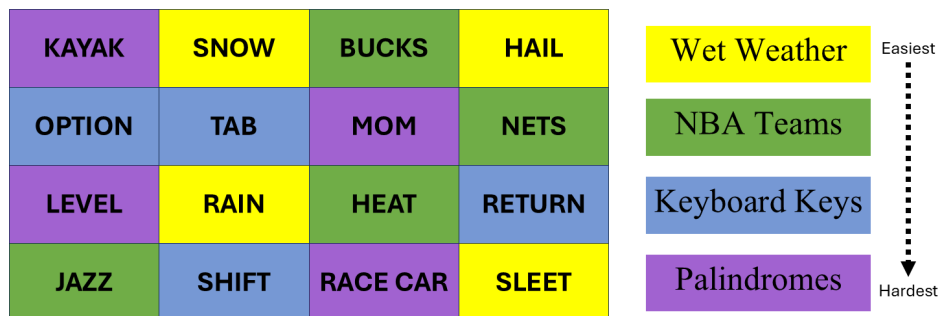


Figure 10: Sample Connections puzzle (left) with categories (right); puzzle in original format does not have colors. Colors reflect difficulty, as determined by the editors of the publication: yellow is the easiest, green is the second-easiest, blue is the second-hardest, and purple is the hardest.

427 ilarity to determine a single sense of similarity? Importantly, this research
 428 question would be difficult to test with standard paradigms, as it involves
 429 significant numbers of pair-wise comparisons that would be both costly to
 430 run and difficult to interpret. In this experiment, we demonstrate how the
 431 drag-and-drop functionality of GRIS-based experiments easily allows us to
 432 determine how different types of similarity are represented and prioritized
 433 among each other.

434 6.1. Stimuli

435 Materials for this experiment come from *Connections*, a free, publicly-
 436 available game hosted by *The New York Times*. In this game, players see a
 437 grid of 16 words and are told to separate the words into four distinct groups
 438 that are labeled; each item belongs to only one group. Importantly, each
 439 group of four words forms a labeled category, and these categories have vary-
 440 ing difficulty: yellow groups are the easiest, green groups the second-easiest,
 441 blue groups the second-hardest, and purple groups the most difficult.¹¹ A
 442 sample item and its corresponding solution are shown in Figure 10.

443 For 300 puzzles, two annotators categorized each group of words into
 444 one of three broader similarity categories: *Semantic Association* (e.g., “wet
 445 weather”: hail, rain, sleet, snow), *World Knowledge* (e.g., “NBA teams”:

¹¹These difficulties are suggested by *The New York Times*; we do not focus on whether these difficulties are accurate, instead studying the cognitive question surrounding similarity comparisons.

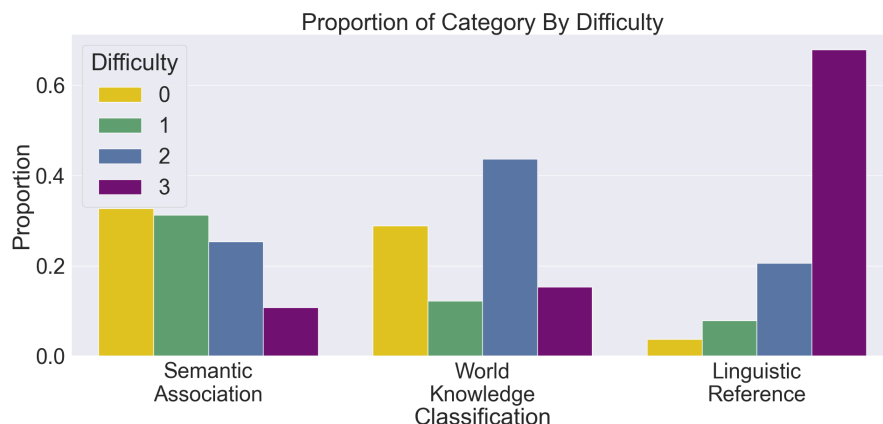


Figure 11: Distribution of similarity categories by difficulty. Difficulty levels closer to 0 are considered easier.

bucks, heat, jazz, nets), and *Linguistic Reference* (e.g., “palindromes”: kayak, level, mom, race car). As visualized in Figure 11, we see that indeed some similarities are considered more difficult than others: semantic association groups tend to occupy the easier categories, world knowledge groups tend to occupy the middle difficulties, and abstract linguistic reference groups tend to occupy the most challenging difficulties.

6.2. Design & Procedure

6.2.1. Stimuli

From our annotated data, we selected 10 puzzles that had at least two of the similarity categories. Given that we are using puzzles generated by the publication, we were unable to perfectly balance the different similarity categories across all puzzles.¹²

6.2.2. Procedure

Similar to Experiment 2, participants saw a blank canvas with a word bank of words below. Participants were instructed to move these words onto the canvas according to how similar they were; similar words should be placed

¹²Instead, categories were balanced to be approximately 40% semantic association, 30% world knowledge, and 30% linguistic reference.

462 closer together. Participants were instructed to use as much of the canvas as
463 they felt was appropriate.

464 To train them on the task but to avoid biasing their decisions, participants
465 completed two practice trials prior to the experiment where they grouped
466 both shapes and numbers.

467 6.2.3. *Participants*

468 Nineteen native speakers of English between the ages of 18 and 55 were
469 recruited on Prolific.

470 6.3. *Results*

471 For each trial, we collected the final position for all words. For every group
472 within each trial, we computed two distance comparisons. WITHIN GROUP
473 distances were computed by calculating the average distance between every
474 word within each group with other members of that same group. OUTSIDE
475 GROUP distances were computed by calculating the average distance between
476 every word within a group with every other word not in that group.

477 Results are visualized in Figure 12. To determine how people used dis-
478 tance to group similar words together, we fit a linear mixed-effects regres-
479 sion model that predicted DISTANCE, with fixed effects of COMPARISON
480 (within group/outside group), CATEGORY (semantic association/world ex-
481 perience/linguistic reference), and their full interactions, along with random
482 intercepts for participants, items, and puzzle difficulty.¹³ We find a main
483 effect of COMPARISON, such that WITHIN GROUP comparisons are signifi-
484 cantly closer together than OUTSIDE GROUP comparisons ($\hat{\beta} = -2.323$, $SE =$
485 0.772 , $t = -3.263$, $p < 0.01$). Additionally, we report a significant interaction
486 between COMPARISON and CATEGORY, such that SEMANTIC ASSOCIATION
487 groups clustered significantly closer together than LINGUISTIC REFERENCE
488 groups in the WITHIN GROUP comparison ($\hat{\beta} = -3.085$, $SE = 0.884$, $t = -3.491$,
489 $p < 0.001$).

490 6.4. *Discussion*

491 In this experiment, we showed that certain similarity patterns are easier
492 to find than others. More specifically, this experiment showed that groups

¹³The complete model formula was: $\text{DISTANCE} \sim \text{COMPARISON} * \text{CATEGORY} + (1 \mid \text{item}) + (1 \mid \text{participant}) + (1 \mid \text{difficulty})$. The baseline conditions were the OUTSIDE GROUP and LINGUISTIC REFERENCE groups, respectively.

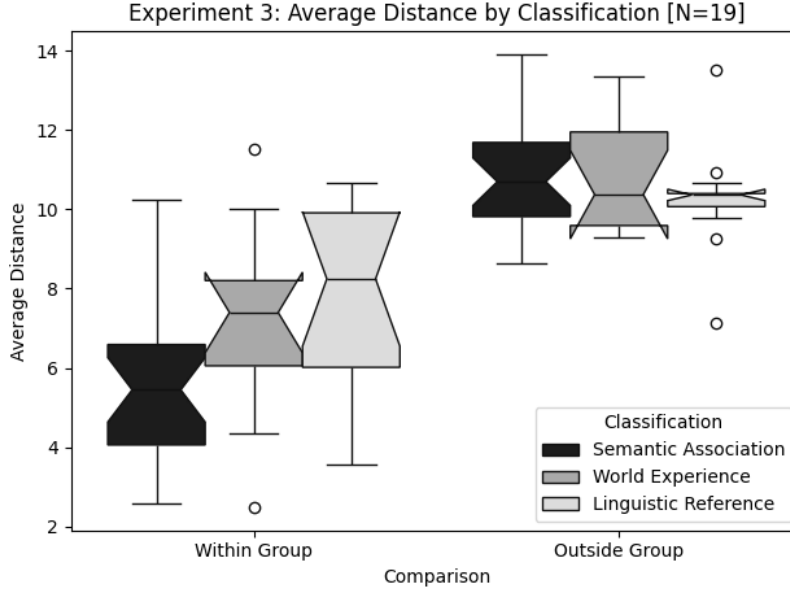


Figure 12: Average distance by category for Experiment 3. Notches indicate bootstrapped 95% CIs.

of words that pattern according to semantic association are easiest to find. These findings may derive from the fact that semantic association requires less reasoning to identify possible clusters of words, compared to other, more abstract groupings.

Beyond these results, we argue that the drag-and-drop paradigm of GRIS-based experiments works well to investigate the complex relationships between representations and reasoning: other paradigms – including rating tasks, forced-choice tasks, and priming tasks – would require significantly less transparent pairwise comparisons to accomplish the results of this study.

7. General Discussion

In this paper, we have shown how GRIS allows researchers across the cognitive sciences to use space as a way of approximating human representational spaces, allowing experimenters to model representational spaces both within class (Experiment 2) and across classes (Experiment 3), while also providing information about the relative relationships between objects on the grid across participants. These findings align with prior work which demonstrate

509 how similarity and difference is highly individualized (Simmons and Estes,
510 2008).

511 Additionally, we have shown how GRIS can the use of space can easily
512 contextualize psycholinguistic findings: we found that acceptability differ-
513 ences between sentence pairs can vary greatly according to the context that
514 they appear in (Experiment 1). We hope that future work using GRIS can
515 expand the relative comparisons between different stimuli modalities (e.g.,
516 text, image, audio).

517 *7.1. What kinds of analyses does GRIS support?*

518 In this subsection, we introduce four broad categories for analyzing future
519 GRIS data, each of which are tied to specific kinds of research questions.
520 These broad categories are:

- 521 1. Location-based
- 522 2. Graph-based
- 523 3. Timing-based
- 524 4. Trial-based

525 Location-based analyses suit questions about ordering or categorical dis-
526 tinctions between objects. For example, the sample trial in Figure 3 studies
527 the linear order of shapes, where each object’s position on the x-axis reflects
528 the object’s relative roundness, according to the participant: as a result, an
529 analysis for this sample trial would likely focus on the y-axis information for
530 each object, unless otherwise specified in the question. Canvases with cate-
531 gorical splits – like those in Figure 2(B)-(E) – also likely use location-based
532 analyses. We demonstrated location-based analyses in Experiments 1 and 2.

533 Graph-based analyses fit questions that investigate the relative relation-
534 ship between objects. Given that the tool collects information about the
535 individual position of each object over the course of the trial, each GRIS
536 trial builds a fully-connected weighted graph, where each object is a node,
537 and the distance between two objects serves as the weighted edge between
538 these objects. For example, a graph-based analysis would align with an ex-
539 periment involving unsupervised clustering of objects. We demonstrated a
540 graph-based analysis in Experiment 3.

541 Timing-based analyses address questions that involve the order of indi-
542 vidual movements and how long each movement took. For example, a timing-
543 based analysis could indicate which objects were most salient to participants

544 (i.e. which objects were moved first), or whether certain objects were more
545 difficult to place (i.e. took longer to drop) in relation to the relevant research
546 question.

547 Finally, trial-based analyses address questions about participant- and
548 item-level behaviors. For example, a trial-based analysis might study whether
549 people how similar representational spaces are between people; an analysis
550 of this kind might construct a large-scale network of object relations for each
551 participant, and then apply transformations to such networks to determine
552 if certain clusters emerge across participants.

553 8. Why Use GRIS?

554 We conclude the paper by collecting our broader arguments for how GRIS
555 can help further our understanding of the human mind.

556 First, GRIS relies on natural human intuitions around space to build
557 contextual and interpretable approximations of cognitive representations. In
558 this paper, we demonstrate three possible ways that space can be meaning-
559 fully used to advance questions in the cognitive sciences; we hope that future
560 work further develops this approach to understanding the mind.

561 Second, as has been mentioned previously, GRIS is very flexible and can
562 be used to answer a range of questions in the cognitive sciences; the paradigm
563 provides a sandbox for both researchers and participants alike to play in.
564 GRIS is supported for desktop, laptops, and tablets.

565 Third and finally, GRIS creates multi-dimensional representations that
566 are easily comparable to popular computational models of language, such
567 as Large Language Models (LLMs). These representations can be used to
568 further explore mismatches between humans and models to help understand
569 what aspects of human cognition are not determinable from data alone.

570 In summary, we note the centrality of spatial reasoning and language to
571 cognition, and how unifying them can 1) make an experiment more intuitive,
572 2) yield more holistic and contextually-relevant results, and 3) construct rep-
573 resentations that facilitate comparisons between humans and computational
574 models.

575 **References**

- 576 Apidianaki, M., 2023. From word types to tokens and back: A survey of
577 approaches to word meaning representation and interpretation. *Computa-*
578 *tional Linguistics* 49, 465–523. doi:10.1162/coli_a_00474.
- 579 Bryant, D.J., 1997. Representing space in language and perception. *Mind &*
580 *Language* 12, 239–264. doi:10.1111/j.1468-0017.1997.tb00073.x.
- 581 Chater, N., Vitányi, P., 2003. Simplicity: a unifying principle in cog-
582 nitive science? *Trends in Cognitive Sciences* 7, 19–22. doi:10.1016/
583 s1364-6613(02)00005-0.
- 584 Cooper, R.M., 1974. The control of eye fixation by the meaning of spoken
585 language: a new methodology for the real-time investigation of speech
586 perception, memory, and language processing. *Cognitive Psychology* 6,
587 84–107. doi:10.1016/0010-0285(74)90005-X.
- 588 Devlin, J., 2018. Bert: Pre-training of deep bidirectional transformers for
589 language understanding. arXiv preprint arXiv:1810.04805 .
- 590 Edelman, S., 1999. Representation and recognition in vision. MIT Press.
- 591 Edelman, S., Duvdevani-Bar, S., 1997. A model of visual recognition and
592 categorization. *Philosophical Transactions of the Royal Society of London.*
593 *Series B: Biological Sciences* 352, 1191–1202. doi:10.1098/rstb.1997.
594 0102.
- 595 Ethayarajh, K., 2019. How contextual are contextualized word represen-
596 tations? Comparing the geometry of BERT, ELMo, and GPT-2 embed-
597 dings, in: Inui, K., Jiang, J., Ng, V., Wan, X. (Eds.), *Proceedings of the*
598 *2019 Conference on Empirical Methods in Natural Language Processing*
599 *and the 9th International Joint Conference on Natural Language Process-*
600 *ing (EMNLP-IJCNLP)*, Association for Computational Linguistics, Hong
601 Kong, China. pp. 55–65. URL: <https://aclanthology.org/D19-1006/>,
602 doi:10.18653/v1/D19-1006.
- 603 Farmer, T.A., Christiansen, M.H., Monaghan, P., 2006. Phonological typical-
604 ity influences on-line sentence comprehension. *Proceedings of the National*
605 *Academy of Sciences* 103, 12203–12208. doi:10.1073/pnas.0602173103.

- 606 Fauconnier, G., 1994. Mental spaces: Aspects of meaning construction in
607 natural language. Cambridge University Press.
- 608 Fauconnier, G., et al., 2007. Mental spaces. The Oxford Handbook of Cog-
609 nitive Linguistics , 351–376.
- 610 Forster, K.I., Guerrera, C., Elliot, L., 2009. The maze task: Measuring
611 forced incremental sentence processing time. Behavior Research Methods
612 41, 163–171. doi:10.3758/BRM.41.1.163.
- 613 Freeman, J.B., Ambady, N., 2010. Mousetracker: Software for studying real-
614 time mental processing using a computer mouse-tracking method. Behavior
615 Research Methods 42, 226–241. doi:10.3758/BRM.42.1.226.
- 616 Gibson, E., Piantadosi, S., Fedorenko, K., 2011. Using mechanical turk to
617 obtain and analyze english acceptability judgments. Language and Lin-
618 guistics Compass 5, 509–524. doi:10.1111/j.1749-818X.2011.00295.x.
- 619 Hayward, W.G., Tarr, M.J., 1995. Spatial language and spatial representa-
620 tion. Cognition 55, 39–84. doi:10.1016/0010-0277(94)00643-y.
- 621 Hiatt, L.M., Trafton, J.G., 2017. Familiarity, priming, and perception in
622 similarity judgments. Cognitive Science 41, 1450–1484. doi:10.1111/cogs.
623 12418.
- 624 Just, M.A., Carpenter, P.A., Woolley, J.D., 1982. Paradigms and processes
625 in reading comprehension. Journal of Experimental Psychology: General
626 111, 228. doi:10.1037//0096-3445.111.2.228.
- 627 Kemmerer, D., 1999. “near” and “far” in language and perception. Cognition
628 73, 35–63. doi:10.1016/s0010-0277(99)00040-2.
- 629 Kirsh, D., 1995. The intelligent use of space. Artificial intelligence 73, 31–68.
630 doi:10.1016/0004-3702(94)00017-U.
- 631 Konieczny, L., 2000. Locality and parsing complexity. Journal of Psycholin-
632 guistic Research 29, 627–645. doi:10.1023/a:1026528912821.
- 633 Landau, B., Jackendoff, R., 1993. Whence and whither in spatial language
634 and spatial cognition? Behavioral and Brain Sciences 16, 255–265. doi:10.
635 1017/S0140525X00029927.

- 636 Lester, N., Feldman, L., del Prado Martín, F.M., 2017. You can take a noun
637 out of syntax...: Syntactic similarity effects in lexical priming., in: CogSci.
- 638 Lin, Y.C., Chen, H.Y., Lai, Y.C., Wu, D.H., 2015. Phonological similarity
639 and orthographic similarity affect probed serial recall of chinese characters.
640 Memory & Cognition 43, 538–554. doi:10.3758/s13421-014-0495-x.
- 641 Mahowald, K., Hartman, J., Graff, P., Gibson, E., 2016. Snap judgments:
642 A small n acceptability paradigm (snap) for linguistic acceptability judg-
643 ments. Language , 619–635doi:10.1353/lan.2016.0052.
- 644 Matthews, J., Starr, J., Schijndel, M., 2024. Semantics or spelling? probing
645 contextual word embeddings with orthographic noise, in: Ku, L.W., Mar-
646 tins, A., Srikumar, V. (Eds.), Findings of the Association for Computa-
647 tional Linguistics: ACL 2024, Association for Computational Linguistics,
648 Bangkok, Thailand. pp. 4495–4504. URL: [https://aclanthology.org/](https://aclanthology.org/2024.findings-acl.266/)
649 2024.findings-acl.266/, doi:10.18653/v1/2024.findings-acl.266.
- 650 Maxfield, J.T., Stalder, W.D., Zelinsky, G.J., 2014. Effects of target typical-
651 ity on categorical search. Journal of Vision 14, 1–1. doi:10.1167/14.12.1.
- 652 McNamara, T.P., 2005. Semantic priming: Perspectives from memory and
653 word recognition. Psychology Press.
- 654 Meints, K., Plunkett, K., Harris, P.L., 1999. When does and ostrich become
655 a bird? the role of typicality in early word comprehension. Developmental
656 Psychology 35, 1072. doi:10.1037//0012-1649.35.4.1072.
- 657 Mimno, D., Thompson, L., 2017. The strange geometry of skip-gram with
658 negative sampling, in: Palmer, M., Hwa, R., Riedel, S. (Eds.), Pro-
659 ceedings of the 2017 Conference on Empirical Methods in Natural Lan-
660 guage Processing, Association for Computational Linguistics, Copenhagen,
661 Denmark. pp. 2873–2878. URL: <https://aclanthology.org/D17-1308/>,
662 doi:10.18653/v1/D17-1308.
- 663 Misra, K., 2022. minicons: Enabling flexible behavioral and representational
664 analyses of transformer language models. arXiv preprint arXiv:2203.13112
665 .

- 666 Misra, K., Ettinger, A., Rayz, J.T., 2021. Do language models learn typicality
667 judgments from text? Proceedings of the Annual Meeting of the Cognitive
668 Science Society 43.
- 669 Neely, J.H., Keefe, D.E., Ross, K.L., 1989. Semantic priming in the lexi-
670 cal decision task: roles of prospective prime-generated expectancies and
671 retrospective semantic matching. *Journal of Experimental Psychology:*
672 *Learning, Memory, and Cognition* 15, 1003. doi:10.1037//0278-7393.
673 15.6.1003.
- 674 Osherson, D.N., Smith, E.E., Wilkie, O., Lopez, A., Shafir, E., 1990.
675 Category-based induction. *Psychological Review* 97, 185. doi:10.1037/
676 0033-295X.97.2.185.
- 677 Papagno, C., Vallar, G., 1992. Phonological short-term memory and the
678 learning of novel words: The effect of phonological similarity and item
679 length. *The Quarterly Journal of Experimental Psychology* 44, 47–67.
680 doi:10.1080/14640749208401283.
- 681 Parrish, A., 2017. Poetic sound similarity vectors using phonetic features,
682 in: Proceedings of the AAAI Conference on Artificial Intelligence and In-
683 teractive Digital Entertainment, pp. 99–106.
- 684 Passos, A., Kumar, V., McCallum, A., 2014. Lexicon infused phrase embed-
685 dings for named entity resolution. arXiv preprint arXiv:1404.5367 .
- 686 Pennington, J., Socher, R., Manning, C.D., 2014. Glove: Global vectors for
687 word representation, in: Proceedings of the 2014 Conference on Empirical
688 Methods in Natural Language Processing (EMNLP), pp. 1532–1543.
- 689 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al.,
690 2019. Language models are unsupervised multitask learners. OpenAI blog
691 1, 9.
- 692 Reed, S.K., 1972. Pattern recognition and categorization. *Cognitive Psy-*
693 *chology* 3, 382–407. doi:10.1016/0010-0285(72)90014-X.
- 694 Rosch, E., 1975. Cognitive representations of semantic categories. *Journal*
695 *of Experimental Psychology: General* 104, 192. doi:10.1037/0096-3445.
696 104.3.192.

- 697 Rosch, E., Simpson, C., Miller, R.S., 1976. Structural bases of typicality
698 effects. *Journal of Experimental Psychology: Human Perception and Per-*
699 *formance* 2, 491. doi:10.1037/0096-1523.2.4.491.
- 700 Shelton, J.R., Martin, R.C., 1992. How semantic is automatic semantic
701 priming? *Journal of Experimental Psychology: Learning, Memory, and*
702 *Cognition* 18, 1191. doi:10.1037//0278-7393.18.6.1191.
- 703 Silfverberg, M.P., Mao, L., Hulden, M., 2018. Sound analogies with phoneme
704 embeddings. *Society for Computation in Linguistics* 1.
- 705 Simmons, S., Estes, Z., 2008. Individual differences in the perception of simi-
706 larity and difference. *Cognition* 108, 781–795. doi:10.1016/j.cognition.
707 2008.07.003.
- 708 Sprouse, J., 2011. A validation of amazon mechanical turk for the collection
709 of acceptability judgments in linguistic theory. *Behavior Research Methods*
710 43, 155–167. doi:10.3758/s13428-010-0039-7.
- 711 Sprouse, J., Schütze, C.T., Almeida, D., 2013. A comparison of informal
712 and formal acceptability judgments using a random sample from linguistic
713 inquiry 2001–2010. *Lingua* 134, 219–248. doi:10.1016/j.lingua.2013.
714 07.002.
- 715 Sweetser, E., 1999. Compositionality and blending: semantic composition
716 in a cognitively realistic framework. *Cognitive Linguistics: Foundations,*
717 *Scope, and Methodology* 129162.
- 718 Talmy, L., 1983. How language structures space, in: *Spatial Orientation:*
719 *Theory, Research, and Application*. Springer, pp. 225–282.
- 720 Tanenhaus, M.K., Spivey-Knowlton, M.J., Eberhard, K.M., Sedivy, J.C.,
721 1995. Integration of visual and linguistic information in spoken language
722 comprehension. *Science* 268, 1632–1634. doi:10.1126/science.7777863.
- 723 Taylor, H.A., Tversky, B., 1992. Spatial mental models derived from survey
724 and route descriptions. *Journal of Memory and language* 31, 261–292.
725 doi:10.1016/0749-596X(92)90014-0.
- 726 Timkey, W., van Schijndel, M., 2021. All bark and no bite: Rogue dimen-
727 sions in transformer language models obscure representational quality, in:

- 728 Moens, M.F., Huang, X., Specia, L., Yih, S.W.t. (Eds.), Proceedings of
729 the 2021 Conference on Empirical Methods in Natural Language Process-
730 ing, Association for Computational Linguistics, Online and Punta Cana,
731 Dominican Republic. pp. 4527–4546. URL: [https://aclanthology.org/](https://aclanthology.org/2021.emnlp-main.372/)
732 2021.emnlp-main.372/, doi:10.18653/v1/2021.emnlp-main.372.
- 733 Turner, A.M., Thiergart, L., Leech, G., Udell, D., Vazquez, J.J., Mini, U.,
734 MacDiarmid, M., 2023. Activation addition: Steering language models
735 without optimization. arXiv e-prints , arXiv–2308.
- 736 Tversky, B., Lee, P.U., 1998. How space structures language, in: Spatial
737 cognition: An interdisciplinary approach to representing and processing
738 spatial knowledge. Springer, pp. 157–175.
- 739 Wilcox, E.G., Ding, C., Sachan, M., Jäger, L.A., 2024. Mouse tracking
740 for reading (motr): A new naturalistic incremental processing measure-
741 ment tool. Journal of Memory and Language 138, 104534. doi:10.3929/
742 ethz-b-000676855.
- 743 Wu, Z., Arora, A., Wang, Z., Geiger, A., Jurafsky, D., Manning, C.D., Potts,
744 C., 2024. Reft: Representation finetuning for language models. Advances
745 in Neural Information Processing Systems 37, 63908–63962.
- 746 Zehr, J., Schwarz, F., 2018. PennController for Internet Based Experiments
747 (IBEX).
- 748 Zouhar, V., Chang, K., Cui, C., Carlson, N., Robinson, N., Sachan, M.,
749 Mortensen, D., 2023. Pwesuite: Phonetic word embeddings and tasks
750 they facilitate. arXiv preprint arXiv:2304.02541 .