

Similarity judgments predict N400 amplitude differences between taxonomic category members and thematic associates

Garrett Honke^{*}, Kenneth J. Kurtz, Sarah Laszlo

Binghamton University (SUNY), USA

ARTICLE INFO

Keywords:

Similarity
Semantics
Conceptual structure
Taxonomic categories
Thematic association
N400

ABSTRACT

Human similarity judgments do not reliably conform to the predictions of leading theories of psychological similarity. Evidence from the triad similarity judgment task shows that people often identify thematic associates like *DOG* and *BONE* as more similar than taxonomic category members like *DOG* and *CAT*, even though thematic associates lack the type of featural or relational similarity that is foundational to theories of psychological similarity. This specific failure to predict human behavior has been addressed as a consequence of education and other individual differences, an artifact of the triad similarity judgment paradigm, or a shortcoming in psychological accounts of similarity. We investigated the judged similarity of semantically-related concepts (taxonomic category members and thematic associates) as it relates to other task-independent measures of semantic knowledge and access. Participants were assessed on reading and language ability, then event-related potentials (ERPs) were collected during a passive, sequential word reading task that presented pseudowords and taxonomically-related, thematically-related, and unrelated word sequences, and, finally, similarity judgments were collected with the classic two-alternative forced-choice triad task. The results uncovered a correspondence between ERP amplitude and triad-based similarity judgments—similarity judgment behavior reliably predicts ERP amplitude during passive word reading, absent of any instruction to consider similarity. It was also found that individual differences in reading and language ability independently predicted ERP amplitude. This evidence suggests that similarity judgments are driven by reliable patterns of thought that are not solely rooted in the interpretation of task goals or reading and language ability.

Determining when human similarity judgments will match the predictions of psychological theories of similarity remains an unsolved problem. Similarity judgments are characteristically unstable and manipulable. The consequences of this lack of understanding of human behavior are compounded by a pressing need for better algorithmic approaches for determining conceptual similarity and semantic relatedness (Kacmajar and Kelleher, 2019). Empirical inquiries into task design and stimulus-based determinants of human similarity judgments show that individual judgment preferences can persevere in the most biasing of circumstances (Honke, 2017; Honke and Kurtz, 2019; Lin and Murphy, 2001): judgment tasks with unambiguous instructions increase the frequency of theoretically-consistent similarity judgment behavior; providing a standard for comparison increases similarity-based matching in the presence of distractors (but has the opposite effect when they are absent); the characteristics of the stimulus set (as measured by human association and similarity ratings) also have predictive value; changing the premise of the question can also affect outcomes, where

people are less likely to follow theoretical predictions under some circumstances (Lin and Murphy, 2001). Yet, these factors alone cannot consistently predict similarity judgment behavior. “Holdouts” can be found in every sample. There are always people who produce the opposite responding pattern in situations that bias the majority of the sample to produce theory-consistent or inconsistent similarity judgments.

Responding preferences are most frequently investigated with the two-alternative, forced choice triad task (see Fig. 1), where similarity judgments are solicited by providing respondents with a base concept (or standard) and two target concepts, a taxonomic category match and a thematically-associated match (Gentner and Brem, 1999; Greenfield and Scott, 1986; Honke and Kurtz, 2019; Lin and Murphy, 2001; Mirman & Graziano, 2012; Skwarchuk and Clark, 1996; Simmons and Estes, 2008; Smiley and Brown, 1979). Taxonomic category members have extensive featural overlap and similarity in relational structure (e.g., *BUTTER* and *JELLY*). Thematic associates share membership in a common

^{*} Corresponding author.

E-mail address: gthonke1@binghamton.edu (G. Honke).

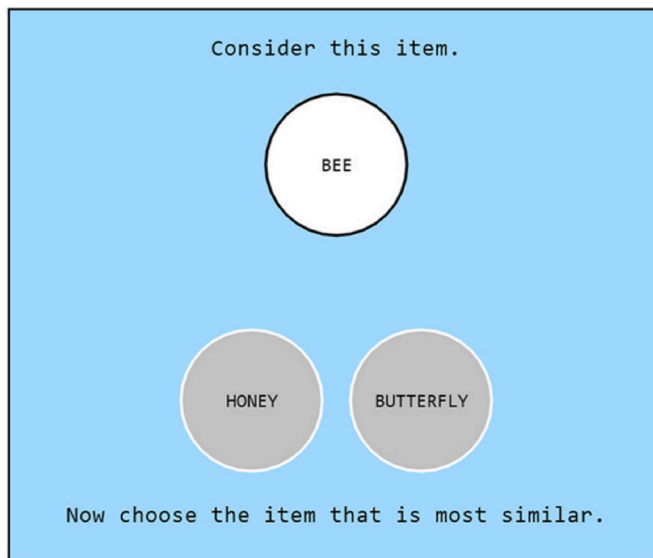


Fig. 1. Visual depiction of the triad similarity judgment task and instructions used to elicit similarity judgments.

theme (e.g., BUTTER and KNIFE).

The consistent observation of thematic intrusion during similarity judgment tasks despite strong manipulation suggests that more work is needed to understand this phenomenon (see Estes et al., 2011 for review; cf. Honke and Kurtz, 2019). It might be more effective, then, to try to predict when thematic intrusion will occur and who will be most susceptible to its effects. A critical component that remains understudied in this research area is individual variation in preference or ability to identify and distinguish between taxonomic category members and thematic associates for the purposes of judging similarity—though see Mirman and Graziano (2012), Murphy (2001), and Simmons and Estes (2008). The goal of this work is to further clarify the role of this variation in similarity judgment behavior by looking at online processing of these semantic relations under *completely unbiased* conditions (i.e., lexical decision task) and connecting this processing to behavioral response patterns from the classic forced-choice, taxonomic–thematic conflict triad task. This design directly addresses two contrasting theoretical viewpoints: Is thematic matching in the triad task the result of confusion about distinct sources of semantic relatedness (e.g., Gentner and Brem, 1999)? Or, is this behavior a result of a process that fuses thematic and taxonomic information to produce similarity judgments (e.g., Bassok and Medin, 1997; Chen et al., 2013, 2014; Simmons and Estes, 2008)? We hypothesize that an examination of the processing of these semantic relations under unbiased conditions can help to tease apart these competing hypotheses and clarify when and why deviations from psychological theories of similarity occur.

In this study, we collected electroencephalographic (EEG) data elicited by the observation of semantically-related and unrelated word (and wordform) sequences and analyzed them in relation to overt similarity judgments of the same concepts in the classic 2AFC triad task. The idea was to examine the processing of taxonomic category members and thematic associates outside of the influence of the judgment task instructions and context, and then investigate how performance in the triad task—a task shown to produce both taxonomically and thematically-biased responding—is related to an unbiased measure of semantic processing (i.e., ERP amplitude). No previous work has attempted to link ERP waveforms and similarity judgments while maintaining an EEG recording procedure free from directed, semantic-task bias. No previous work has looked at the relationship between ERPs and overt similarity judgments for the purpose of characterizing divergent, individualized activation and decision patterns. These theoretical and methodological advances increase the likelihood that

heretofore undetected ERP differences between taxonomic category members and thematic associates will be uncovered and clarify the strength of evidence for existing theoretical accounts of thematic intrusion on human similarity judgment.

1. Characterizing ERPs elicited by taxonomic similarity and thematic association

Electrophysiological research in this domain has generally fallen short of the goal of discovering *semantic* processing differences between EEG elicited by taxonomic category members and thematic associates (bounding semantic processing to the time period that the N400 component is believed to occur, roughly 300–400 ms post stimulus presentation). Thus, research has failed to increase understanding of a critical issue: What causes people to make more or less theoretically-consistent responses in similarity judgment tasks? Determining if response patterns are indicative of general patterns or biases in thinking is a key first step for understanding the role and impact of these biases on higher-order cognition. Inference from purely behavioral research is problematic because responding preferences could be an artifact of the “match-to-sample” task most frequently used to collect similarity judgments. It is not yet known if there are observable neural activation patterns that correspond to observed response biases. If this correspondence exists, however, it can provide insight into why—outside of the influence of concepts, task instructions, and design—people produce different responding patterns in similarity judgment tasks. Essentially the question is whether or not responding patterns are a consequence of semantic network organization or the 2AFC triad task.

While existing work has not adequately addressed this question, success has certainly been found in clarifying the general ERPology (i.e., the character and form of ERPs elicited by certain stimuli, see Luck, 2014, pg. 5) of the processing of these semantic relations, particularly in relation to semantically unrelated concepts. In one such study, Chen et al. (2013) recorded EEG while people performed a similarity or difference judgment task for a sequence of taxonomic and thematic category pairs. No amplitude differences between taxonomic and thematic pairs were found in the N400 time window. The analysis uncovered a reliable difference in the amplitude of the P600 component elicited by taxonomic and thematic category members—a larger (more positive) P600 for taxonomic pairs. The authors argue that this P600 difference is evidence of “less syntactic flow” in the processing of taxonomic relations (Chen et al., 2013). Another study from Chen and colleagues (Chen et al., 2014) collected EEG in a sequential concept priming experiment administered in conjunction with a lexical decision task. Study participants viewed taxonomic and thematic category pairs while indicating if the stimuli were words or non-words with a button press. No amplitude differences were found for the semantic classes of interest in the N400 time window. The study did uncover a reduced frontal negativity effect for productive thematic associations (e.g., BEE and HONEY) as compared to hierarchical relations (i.e., taxonomic category members) and other subcategories of semantic relations not relevant for this work. In other words, some qualified evidence was found for a facilitative priming effect (a reduced negative frontal activation in the 400–550 ms time window) for thematic associates as compared to taxonomic category members.

Work by Wamain et al. (2015) had more success in uncovering differences between these semantic relations. The authors found EEG amplitude differences between pictorial depictions of thematic associates and two specific sub-types of taxonomic category members (taxonomic category members that share a specific function or a general function, e.g., SAW–AXE and SAW–KNIFE, respectively). The task was to observe visual depictions of semantically related concepts and vocally name the pairs after EEG collection was finished for the trial. The difference between thematic and the subordinate taxonomic classes was reliable but this effect was only found at short inter-stimulus intervals (66 ms).

Maguire and colleagues (Maguire et al., 2010) also contribute to this effort with an ERP and ERSP (event-related spectral perturbation) based design where EEG was collected during a passive listening task. The authors found a distinction in the distribution of the power of certain frequencies across the scalp: more alpha power was found over the parietal areas of the brain for taxonomic category members and more theta power was found over the right frontal areas of the brain for thematic category members. The authors suggest that this increase in parietal alpha power is due to the fact that additional attentional resources are required to process taxonomic category members—a conclusion that dovetails with the idea that (1) processing taxonomic similarity requires an effortful comparison process (Geller et al., 2019; Kurtz et al., 2001), (2) processing taxonomic similarity is more difficult than processing thematic association (Sachs et al., 2008) and (3) less-educated people (Denney, 1974; Sharp et al., 1979; cf. Mirman and Graziano, 2012) and people who score low on the *Need for Cognition* assessment (Simmons and Estes, 2008) experience more thematic intrusion on similarity judgments.

In the closest such study to the methodology of the present experiment, Lewis et al. (2015) discovered new support for the growing body of evidence that suggests that the anterior temporal lobe (ATL) and the temporoparietal junction (TPJ) play distinct roles for the processing of taxonomic and thematic category members. In this magnetoencephalography (MEG) study, the authors show that the MEG activation pattern of these semantic classes is distinct in the anatomical regions mentioned above. Participants in the task were asked to make a physical response if the presented stimuli were related or unrelated. This is an important design choice, as the question of “relatedness” has been shown to bias responding in the triad task toward association-based responding (Skwarchuk and Clark, 1996).

Theoretical and Methodological Advances in the Present Work. A central goal of research in this area (including the studies reviewed above) has been to test for differences in *facilitative priming* between taxonomic, thematic, and unrelated word pairs as evidenced by diverging waveform amplitude roughly 300–400 ms post stimulus exposure (Kutas and Federmeier, 2011). However, there are only a few (qualified) successes and several failures in this effort to find distinctive N400 patterns between taxonomic and thematic category members in non-clinical adults samples (Chen et al., 2013, 2014; Khateb et al., 2003; Maguire et al., 2010)—notable exceptions being the work of Wamain et al. (2015), Lewis et al. (2015) and Hagoort et al. (1996). There is strong evidence that the processing of taxonomic and thematic category members occurs in different systems or networks (Lewis et al., 2015; Schwartz et al., 2011), so why do temporally-constrained and spatially-unconstrained EEG-based approaches fail to detect differences? Or stated differently, given the apparent difficulty in finding unqualified differences between taxonomic and thematic category processing, why continue to use electrophysiology to study the role of these semantic relations in similarity judgment research?

Methodological and theoretical adjustment could address several of the issues raised here. It is common in past investigations to see ERPs elicited from taxonomic and thematic category members analyzed in the aggregate (factorial analyses, e.g., ANOVA). Could it be that averaging over the sample obscures important differences in the processing of these semantic relations? Further, it has been shown that behavioral data analyzed with a factorial approach at the group level is anti-conservative (Honke, 2017; Honke and Kurtz, 2019). Whether the results are obscured by aggregation or the outcomes are anti-conservative, a major motivation of this work is to explore the use of individualized experimental design and analysis to study this (apparent) individual differences-driven phenomenon.

2. Individualized measurements for individual-differences in similarity responding patterns

Our hypothesis is that analyses that average across participants

obscure important differences—i.e., people who exhibit strong taxonomic or thematic response biases work against the calculation of a mean amplitude outcome variable. Consider that the most likely manifestation of behavioral biases (if they are detectable via electrophysiology) would be more facilitative priming (i.e., increased N400 positivity) for a specific type of semantic relationship. In this scenario, averaging across a sample of people who have reliable but opposite biases would obscure differences—thematic responders would show increased facilitative priming for thematic category members, taxonomic responders would show increased facilitative priming for taxonomic category members, and these differences would not be preserved in a measure of mean amplitude. Similarly, consider the hypothesis that people who are more susceptible to thematic intrusion produce less distinct ERPs between these semantic classes—these people are included in aggregation-based approaches as well.

There is also concern that the distinct classes of stimuli *themselves* should produce different activation patterns. Stimuli that have been well-normed would be expected to elicit different N400 activation patterns in an adequately-powered experiment simply by virtue of being different classes of semantic relations. For these reasons, the present work focuses more closely on individual differences by classifying participants based on their similarity judgment behavior and then using this classification to look at N400 amplitude differences across groups.

3. Effects of intervening tasks on ERPs and other methodological concerns

There are several methodological adjustments that can increase the likelihood that differences between taxonomic and thematic pairs can be detected. First, previous studies have often included intervening tasks directly or indirectly related to the question(s) at study during EEG recording (e.g., similarity judgments, difference judgments, button pressing). Related tasks affect the EEG signal (Luck, 2014), particularly those that require a physical response. The signal elicited by these responses cannot be distinguished from the underlying processes at study and the result is EEG data confounded by the related task. Similar to Maguire et al. (2010), the present design features passive EEG collection with no explicit task instructions or behavioral task related to the processing of the semantic relations at study. Instead, participants are asked to perform a lexical decision task where they identify pseudowords as they appear in the stimulus stream. Thus, the measuring of semantic processing does not include response potentials (trials with erroneous responses to real words are removed from analysis); the task is simply to respond if the letter string is not recognized as a word. This effectively eliminates the risk of signal contamination from the evoked response potential while ensuring that focus is maintained on the stimulus stream and not biasing participants to perform a particular semantic task.

Additionally, concepts are presented with long enough ISIs (3–3.5 s) that time-locked EEG data can be reliably attributed to the most recently presented stimulus and its semantic relationship with the preceding concept (i.e., distanced from the processing of the preceding concept itself). Results will be presented and analyzed without averaging across electrode sites, as this type of averaging carries the risk of obscuring real effects and producing anomalous patterns (Thigpen et al., 2017). Lastly, confirmatory data analysis will be restricted to the a priori hypotheses presented below—hypotheses that only relate to amplitude differences in the established time window for semantic effects (Kutas and Federmeier, 2011; Kutas and Hillyard, 1980).

4. Breadth of taxonomic and thematic category members

The types of thematic and taxonomic relations used in previous investigations have been too restrictive to make class-wide conclusions about similarity processing. This is not a problem for the particular studies we have outlined here, i.e., it is reasonable to investigate specific types of taxonomic categories (e.g., function specific taxonomic

categories, Wamain et al., 2015) or thematic relations (e.g., productive relations, Chen et al., 2014) if the research interest is in those specific sub-types. However, in this work we adopt an expansive definition where thematic category members only require temporal contiguity in an established situation and taxonomic category members are entities of the same *kind*, i.e., entities that share membership in a category of natural kinds or artifacts that is well-described by a common set of shared features and relational structure (Kurtz and Gentner, 2001; Lin and Murphy, 2001; Mirman et al., 2017).

5. The current study

The central goal of this research is to test for evidence that links the unbiased processing of taxonomic and thematic category members with similarity judgments from those same stimuli. The broad methodological hypothesis is that facilitative priming differences between these semantic relationships have been difficult to detect for the reasons outlined above. Conceptually, the question is: What if thematic intrusions are due to difficulty distinguishing between types of semantic relatedness, i.e., less distinctive EEG activation patterns between types of semantic relations? Looking for answers for these questions by averaging across an entire sample would fail if elicited waveforms have a direct correspondence with similarity judgments—which often average to a *slight* taxonomic match preference when the task goal is left ambiguous (Honke, 2017; Honke and Kurtz, 2019). The present study uses a novel experimental design to match concept similarity judgments with EEG elicited during the passive processing of those same taxonomic and thematic category members. This approach has the potential to uncover presently unknown properties of taxonomic and thematic processing and how these properties relate to similarity judgments.

5.1. Competing hypotheses for similarity judgment and their predictions

The field's understanding of the relationship between taxonomic and thematic semantic processing, the electrophysiological patterns they elicit, and their role in the formation of similarity judgments is limited. On the behavioral side, two hypotheses have been proposed to explain the effect of thematic intrusion (when concept association affects similarity judgments): the confusability account and the dual-process integration account. The confusability account posits that individuals differ in their susceptibility to thematic intrusion on their similarity judgments. The dual-process integration account suggests that associative processing cannot be excluded from similarity judgment, it is an integrated component process of the similarity judgment system.

The focal question of past EEG research has been: Are there general, sample-level differences in the N400 elicited by taxonomic and thematic category pairs? The problem is that this general approach will fail to detect N400 differences if the thematic intrusion phenomenon is better explained by the confusability account. For this reason, we focus on how differences in similarity judgment behavior might correspond to differences in semantic electrophysiology at the individual level.

What do the confusability and dual-process integration accounts predict about EEG elicited by taxonomic and thematic pairs and their corresponding similarity judgments? Chen et al. (2013) suggest that the integration account is supported by evidence that N400s elicited by taxonomic and thematic pairs are not reliably different. The authors suggest that a similarity judgment process that integrates taxonomic and thematic information should produce similar ERP amplitude in the time window of the semantically-sensitive N400 component. The lack of a reliable N400 amplitude differences between taxonomic and thematic category members is presented as support for the integration account (Chen et al., 2013).

In contrast, the confusability account suggests that some people are better than others at separating the results of distinct, semantic-relatedness processes and these people are less subject to thematic intrusion. To our knowledge, no electrophysiological correlate to this

hypothesis has been proposed.

Following from the confusability account, we hypothesize that this susceptibility to intrusion is directly attributable to differences in facilitative priming between semantic classes. We predict that reliable differences in waveforms elicited by these semantic classes at the individual level correspond to similarity-based responding in the triad task; responding that is resistant to thematic intrusion. This possibility directly relates to Gentner and Brem's argument that the similarity process is derailed when people have difficulty distinguishing between the mental output of similarity and association-based processing (Gentner and Brem, 1999). In the present study, we take the presence (or absence) of facilitative priming differences (as operationalized as N400 amplitude differences between classes) as an electrophysiological marker of the distinctiveness of the output of these processes. To measure this marker, careful consideration of confounding variables was made to minimize the possibility that the link between distinct facilitative priming and similarity judgments could be attributed to a known source of individual differences in responding behavior.

5.2. Toward characterizing individual differences in taxonomic and thematic thinking

The general approach of linking similarity judgments to measures of individual differences such as education (Denney, 1974; Sharp et al., 1979), the Need for Cognition (NFC) scale (Cacioppo and Petty, 1982; Simmons and Estes, 2008), and online processing (Mirman and Graziano, 2012) has had success in uncovering differences between people with different profiles of similarity judgment behavior. Mirman and Graziano (2012) used the visual world paradigm (VWP) eye-tracking task to investigate processing time-course and competition between taxonomic and thematic category members. They found that more competition in the VWP between taxonomic and thematic category members predicted taxonomic responding in the triad task. Assessment measures for language and reading ability were included in the current experiment to address the effect of these consequential—but non-focal—variables. Not only are these measures (exposure to print, verbal fluency, and vocabulary) effective controls for general education and language exposure variance, but they are also important for similarity judgment behavior itself.

Role of Reading Experience and Language Exposure. Three measures were collected to be used as covariates in the study: exposure to print, verbal fluency and vocabulary. The recognition of authors and magazines has been shown to predict orthographic knowledge and experience even when controlling for other measures of general aptitude (e.g., SAT scores) and domain knowledge (West and Stanovich, 1991). Vocabulary knowledge has a direct relationship with semantic priming. In children, words that are less well-known elicit stronger thematic priming than taxonomic priming. The opposite pattern is found for words that children can define and use correctly in a sentence (Ince and Christman, 2002). The relationship between verbal fluency and semantic relation processing is less clear. On one hand, the categories in our verbal fluency assessment (particularly fruits and animals) are superordinate taxonomic categories, so ease of recall of category members could be a measure of taxonomic processing ability. On the other hand, many people are successful in the task by using a free association clustering strategy (Jenkins and Russell, 1952)—like using a biome-based organization, for example, when naming living things (e.g., using the savanna biome to produce lion, elephant, antelope, rhino, zebra, etc.) or a color scheme organization to list colors (e.g., ruby, sapphire, topaz). However verbal fluency relates to the processing of taxonomic and thematic relations, the measure is predicted to help account for variance in the design that would otherwise be attributed to random error or taxonomic responding in the triad task.

Individual Differences and Similarity Judgments. Sharp et al. (1979) showed that educational attainment is related to taxonomic responding. Simmons and Estes (2008) found that triad task responding patterns

related to NFC scores, where lower scorers produced more thematic matches. [Mirman and Graziano \(2012\)](#) did not find demographic differences (i.e., education, age) to be predictive of triad responding behavior. At the least, we hypothesize that including these specific reading and language exposure assessments will allow us to disentangle the contribution of these factors and similarity judgment behavior in the analysis of EEG elicited from taxonomic and thematic category members. The outcome of these assessments was analyzed in relation to similarity judgment behavior in addition to being included in the analysis of the EEG data.

Choosing an Appropriate Task for Collecting Similarity Judgments. In an experiment on the effect of task instructions on similarity judgment behavior, we found that similarity-based instructions produced the most ambiguous responding behavior in the triad task ([Honke, 2017](#); [Honke and Kurtz, 2019](#)), i.e., “Choose the option that is most similar” (see [Fig. 1](#)). These task instructions were deliberately chosen for the similarity judgment phase of the present study. It is convenient for comparison to past work that these instructions coupled with the classic triad task are also *the most frequently used way to assess similarity judgment behavior*. They are desirable for this work because they produce a varied spread of the possible response biases. The motivation was to use a task that has the least biasing conditions in order to maximize the diversity of observed response patterns and sample roughly equal groups of participants for the EEG comparison.

6. Method

6.1. Participants

Undergraduate students ($N = 61$) from Binghamton University were recruited from the Psychology Department pool ($n = 53$) or the university community ($n = 8$) and participated for credit toward the completion of a course requirement or \$30.00 cash compensation, respectively (36 female; $\text{Age}\bar{X} = 19.0$, $\text{AgeRange} = 17\text{--}23$). Three participants were dropped due to experimenter error during the EEG collection phase. Three participants were missing data from part of the procedure; the demographics survey, the demographics survey and verbal fluency assessment, and exposure to print assessment, respectively. Where needed, these missing values were imputed with the *mi* package ([Su et al., 2011](#)) in R ([R Core Team, 2017](#)). In the analysis below, this resulted in a total of 58 participants: 56 participants with complete data and two participants with imputed values for the assessments mentioned above. The study was approved by the Internal Review Board of Binghamton University. Participants identified themselves as right-handed, monolingual English speakers with little-to-no early life exposure to any other language, normal or corrected-to-normal vision and no history of psychiatric or neurological disorders. Participants who reported recent alcohol, prescription, or recreational drug use that could affect their performance were asked to reschedule the experiment.

7. Materials

7.1. Reading and language exposure assessment

Three measures of reading and language exposure were collected prior to the EEG recording phase of the experiment. *Exposure to print* was assessed with a 160 item questionnaire consisting of real and fake authors and magazine titles following from the work of [Stanovich and West \(1989\)](#). The task was to indicate which items in the questionnaire were real while minimizing false positives. d' values were calculated for each participant as a measure of individual differences in recognition ability. *Verbal fluency* was assessed with a category member naming task where the goal was to name as many examples of a given category (fruit, colors, animals) in 60 s. The third assessment was a *vocabulary* test. It consisted

of 30 items drawn from the Verbal Reasoning section of the Graduate Records Examination (GRE) test. The concepts used in the experiment were well below the reading level of a college-aged sample, but nevertheless it is hypothesized that this measure will help to account for the differences among participants in vocabulary ability.

7.2. Concept set generation and presentation order

Concept sets ($N = 100$) were created that consisted of a standard, a taxonomic match, a thematic match, and two unrelated concepts. Concept sets were normed as follows. Similarity and association ratings, mean concreteness ratings ([Brysbaert et al., 2014](#)), and age of acquisition data ([Kuperman et al., 2012](#)) were visualized and examined for outliers. The 20 worst outliers in terms of concreteness, age of acquisition, and difference in similarity and association ratings (i.e., relatedness strength) were removed. This exclusion process resulted in 80 concept sets (see [Table 2](#) for aggregated concept set properties, comprehensive data provided in [Appendix C](#)).

Pseudowords generated from the orthographic and lexical characteristics of the experimental stimuli (i.e., frequency, length, orthographic neighborhood size, and constrained bigram frequency) were paired with concept sets in an iterative procedure that minimized the cost (difference) between the properties of the possible pseudoword matches (string length, orthographic neighborhood size, and bigram frequency) and the mean of those same properties in the real-word concept sets across 10,000 iterations of possible pseudoword–concept set combinations (pseudowords and lexical and orthographic statistics were generated from MCWord, [Medler and Binder, 2005](#)). The purpose of this process was to make sure that the pseudowords were as word-like and similar to their paired concept set as possible. Closely matching pseudowords were expected to increase the difficulty of the pseudoword identification task and thus increase attention to the word stream in the EEG recording phase ([Laszlo et al., 2012](#)).

During the EEG recording phase of the experiment, four categories of word pairs were presented with Psychtoolbox ([Brainard and Vision, 1997](#)) in a continuous stream of wordforms. Each letter string could be preceded by a member of the same taxonomic category, a member of the same thematic category, an unrelated concept, or a pseudoword (see [Fig. 2](#)). Four counter-balanced presentation orders were produced that followed three considerations: randomization of concept/letter string presentation within each set, randomization of concept set presentation across the EEG phase, and randomization of presentation of the taxonomic category member or thematic category member within each set. The latter consideration was required because the standard could not be presented multiple times in the course of EEG recording due to the possible confound of N400 repetition effects for words and non-words ([Laszlo and Federmeier, 2011](#); [Rugg and Nagy, 1989](#)).

Two randomized presentation orders were produced to satisfy the first and second considerations, where concept set order, concept order within set and taxonomic or thematic pair selection was randomly determined. To satisfy the third consideration, two additional orders were produced by replacing the randomly selected taxonomic or thematic matches with their alternatives from the same set; this process produced two sets of two randomly ordered presentation orders and four orders in total. Randomly placing the concept sets into a single stream of words and pseudowords carried the risk that unintended relationships might be produced between adjacent words. This issue was resolved within concept sets by soliciting similarity and association ratings from a separate sample of participants (results below). Between-set correspondences were handled by a team of research assistants that independently examined each counter-balanced presentation order to confirm that concepts at the boundaries between concept sets did not have incidental taxonomic or thematic relationships. When relationships were identified (independent of how weak they were perceived to be) the presentation order was altered to break up these incidental pairings.

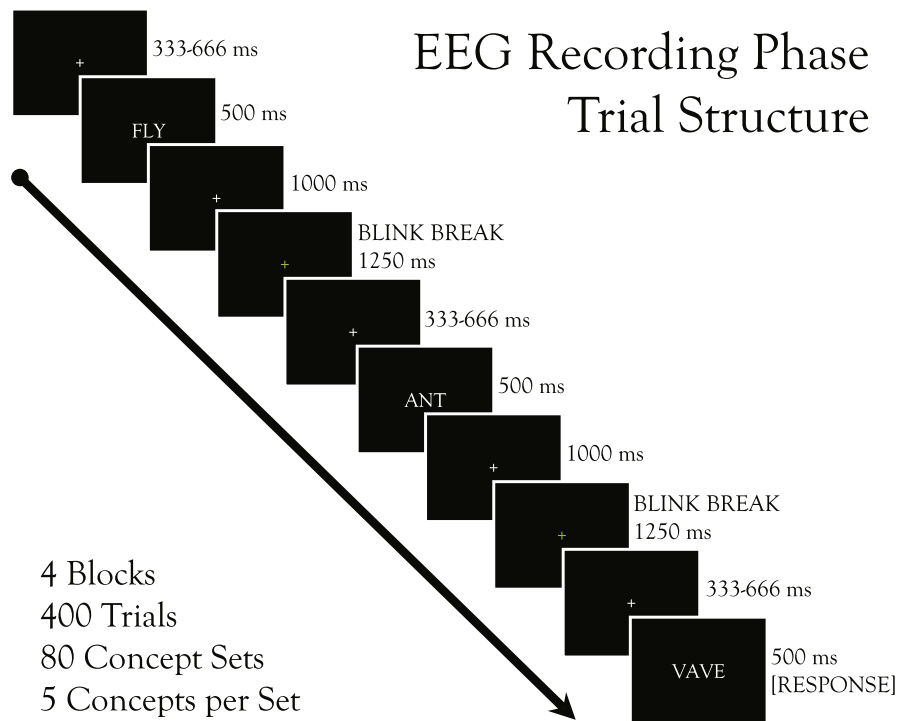


Fig. 2. Visual depiction of the trial structure for the EEG recording phase. The task goal was to observe a continuous stream of concepts and respond by pressing a button when a pseudoword appeared in the stream.

7.3. EEG recording and processing

An elastic EasyCap with 26 geodesically arranged,¹ passive amplification, ring-sintered Ag/AgCl electrodes (inter-electrode impedance maintained below 2 k Ω , see Laszlo et al., 2014) was used to record the EEG signal. Two electrodes on the outer canthi of the left and right eyes and one electrode on the suborbital ridge of the left eye were used to record the electrooculogram (EOG) and monitor blinks. The EEG and EOG were referenced to the left mastoid on-line; offline the EEG and EOG were re-referenced to the average of the left and right mastoids, the horizontal EOGs were re-referenced as a singular bipolar channel. The signal was recorded with a Brain Vision Brain Amp DC amplifier (low pass filtered at 250 Hz, high pass filtered with a 10 s time constant, sampled at 500 Hz with an A/D resolution of 16 bits).

A two-stage, offline artifact rejection procedure was applied to each participant's data (code available in the supplementary materials hosted on the Open Science Framework²). First, EEG data for each participant was filtered with a high-pass filter (0.05 Hz), ICA components were computed and components corresponding to blinks were visually identified and removed. Second, the EEG record was visually inspected with a participant-individualized amplitude threshold to identify and remove artifacts less well-identified by ICA (e.g., blocking, drift, horizontal eye movements, etc.). Exclusion criteria were as follows. Participants were candidates for exclusion from the analysis if less than 60% of all trials or less than 60% of a particular concept pair type were retained after the artifact rejection procedure (no participants met these criteria). An average of 89% of trials were retained per concept-pair type (minimum number of trials retained across concept pair types for a single participant: 70%). The EEG record was binned into concept-pair specific ERPs

¹ Geodesic placement refers to the equidistant positioning of electrodes on an approximately spherical surface—this arrangement differs from the 10–20 system that does not feature equidistant placement.

² <https://osf.io/ctzhk/>.

time-locked to stimulus onset with a 100 ms pre-stimulus baseline and a 998 ms post stimulus recording period. A band-pass filter of 0.1–20 Hz was applied to the ERPs for final analysis and presentation (e.g., Figs. 8 and 9).

7.4. Similarity judgment triad task

In the final phase of the experiment, the semantically-related concepts from the EEG phase (the standard, taxonomic match and thematic match from each set) were presented as forced-choice triads with Psychopy (Peirce, 2007). The task was identical to the classic similarity-based, 2AFC triad task (Gentner and Brem, 1999; Greenfield and Scott, 1986; Lin and Murphy, 2001; Mirman and Graziano, 2012; Skwarchuk and Clark, 1996; Simmons and Estes, 2008; Smiley and Brown, 1979). On each trial, a standard was presented first in a prioritized position followed by a taxonomic category member and a thematic category member (randomly placed at the left and right apexes of the triad below the standard). On-screen instructions directed participants to: *Consider this item [the standard] Now choose the item that is most similar*. A depiction of the task is provided in Fig. 1. Final responses, response time and all other behavior was recorded.

7.5. Procedure

Participants entered the lab and were provided with a verbal description of the complete experimental procedure. After attaining informed consent, the demographic survey and reading and language exposure assessments were administered and participants were fitted with the EEG cap. EEG recording occurred in a sound attenuated booth.³ Stimuli were presented at a distance of 75 cm on 24 inch computer monitors displaying at a resolution of 1920 x 1080. Demonstrations of

³ A subset of the sample ($n = 17$) completed the experiment in private testing rooms (not sound attenuated booths) due to lab construction.

the EEG record and the task were provided before the start of EEG collection to (1) illustrate the importance of reducing eye and body movement during EEG collection and (1) orient participants to the pseudoword identification task. Participants were instructed to maintain control of their eye and body movements and press a button as fast as possible when the image presented on the screen contained a string of letters that was not a word. This lexical decision task (LDT) was used to confirm that participants attended to the presented stimuli. The task was designed to be unrelated to the semantic relationships of interest to avoid the introduction of evoked response potentials into the EEG data of the critical trials (semantically-related and unrelated real word pairs). Concepts were presented in a continuous stream broken into four blocks that followed one of the four randomly generated and assigned counter-balanced presentation orders. Breaks were provided in between blocks (after approximately 100 trials); the task resumed when participants indicated that they were ready to start the next block.

Each trial started with a 333–666 ms fixation cross presentation that was randomly jittered in time to avoid anticipatory processing. Stimuli (images of letter strings) were presented for 500 ms followed by a 1000 ms post-stimulus fixation cross and a 1250 ms blink break. The next trial began immediately after the blink break terminated.

After the EEG recording was complete, the EEG cap was removed and participants were allowed as much time as needed before the triad similarity judgment task was started. The triad task was administered on computer and self-paced.

7.6. Statistical methods

The analyses were conducted with linear mixed-effects regression (LMER: Bates et al., 2014; Kuznetsova et al., 2015) models built in R (R Core Team, 2017) to predict amplitude with semantic pair type, word properties, concept association and similarity ratings, participant reading and language experience, similarity judgment behavior, and random effects for participant, time window and concept. Critically, the use of LMER does not require the aggregation of data across participants like factorial analysis approaches; this makes it particularly valuable for the analysis of individual differences. Mean amplitude was examined with 10 ms averaged time points constrained a priori to the time window where the N400 component is most likely to be found (300–400 ms) (Kutas and Federmeier, 2011; Kutas and Hillyard, 1980). Consistent with prior research, unaveraged EEG data collected at central, parietal and occipital electrode sites (MiCe, MiPa, LDPa, RDPa, LMOc, RMOc, LLOc, RLOc, MiOc) were used to capture the broadly distributed N400 effect. A minimal (“parsimonious”) random effects structure was used due to the overall size and complexity of the models. Further, this analysis is not subject to the maxim (and general critique) to *keep it maximal* (Barr et al., 2013), as specifying the maximal random effects structure was not expected to significantly affect parameter estimation in this situation (see Stites and Laszlo, 2015).

The central goal of the analysis was to identify amplitude differences in ERPs that can be linked to (e.g., predicted by) differences in similarity judgment behavior, but the set of additional measures that were collected also have an important relationship to these behavioral patterns. Therefore, in addition to including word-based statistics (word length, frequency, orthographic neighborhood size, and constrained bigram frequency), individual differences in reading and language ability (exposure to print, verbal fluency and GRE vocabulary assessments) and concept similarity and association ratings in the modeling of the ERP waveforms, it was also important to characterize how these variables affect behavioral response patterns in the similarity judgment task. Thus, the similarity judgment data will also be analyzed in relation to these variables.

8. Results

Recall that participants completed a series of reading and language

assessments and then viewed a stream of images of letter strings where temporally adjacent strings could be taxonomic category members, thematic category members, unrelated concepts, or concept–pseudoword pairs. The session ended with a similarity judgment triad task. The results will be presented in four sections: (1) concept rating data, (2) reading and language exposure assessments, (3) behavioral task outcomes, and (4) general ERP results with specific attention to behavioral–electrophysiological correspondences.

8.1. Concept set norming

Concept set ratings were collected with a two condition, between-subjects task with a separate set of participants ($N = 259$, association question condition: $n = 132$) recruited from the Binghamton University Psychology Department Pool. The task was to provide ratings on a ratio-scale rating line (from 0 to 100) where the anchors were NOT AT ALL TO VERY SIMILAR for the taxonomic rating condition and NOT AT ALL TO VERY WELL for the thematic rating condition (where the question targeted how well the items *go together*). A depiction of the task is provided in Appendix D (see Fig. D1). Concept ratings were analyzed to confirm that taxonomic pairs were rated highest on the similarity question, thematic pairs were rated highest on the association question, and that the standardized strength of perceived similarity for a given concept set was not reliably different from the standardized strength of the thematic relationship (Fig. 3). For the latter, the goal of this approach was to confirm that the “quality” of taxonomic pairs (in terms of perceived similarity) did not systematically differ from the quality (or association strength) of the thematic pairs within each concept set. Descriptive statistics are provided in Table 1.

8.2. Similarity and association strength

Mixed-effects LMER models were built to analyze the unadjusted association and similarity ratings. The similarity rating model (pair type as a fixed effect categorical predictor and participant as a random categorical predictor) uncovered reliably higher similarity ratings for the taxonomic pairs as compared to the thematic pairs ($\hat{\beta} = 5.488$, $SE = 0.55$, $t = 9.991$, $p < .001$) and the unrelated pairs ($\hat{\beta} = 55.837$, $SE = 0.48$, $t = 1117.06$, $p < .001$). Similarly, the model built to predict

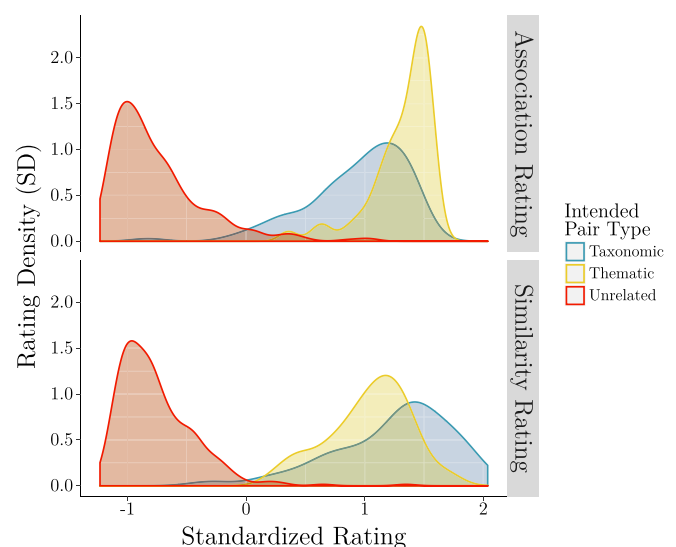


Fig. 3. Density plot of standardized ratings for the association (top) and similarity (bottom) rating tasks. Taxonomic pairs were rated as more similar, thematic pairs were rated as more associated, and unrelated pairs were rated lowest on similarity and association. Taxonomic and thematic pairs in the same concept set were not reliably different in the magnitude of their standardized similarity and association ratings.

Table 1
Concept ratings.

Pair Type	Similarity Rating	Association Rating	Similarity Rating	Association Rating
	Mean (SD)	Mean (SD)	Mean Response Time	Mean Response Time
Taxonomic	70.52 (1.24)	75.26 (0.94)	4.34 s	4.08 s
Thematic	65.05 (1.03)	88.01 (1.31)	4.32 s	3.72 s
Unrelated	14.68 (-0.75)	18.59 (-0.75)	4.38 s	4.47 s

Lexical and Orthographic Properties.

Table 2
Aggregate concept set properties.

Pair Type	Word Length	Word Frequency	Orthographic	Bigram	Similarity/Association
			Neighborhood	Frequency	Difference Score
Taxonomic	5.66	52.33	90.16	1160.26	0.35
Thematic	5.73	40.18	28.55	1133.55	0.31
Set Mean	5.79	43.94	56.72	1148.00	

association ratings showed that the thematic pairs were rated as more associated than the taxonomic pairs ($\hat{\beta} = 12.741, SE = 0.54, t = 23.43, p < .001$) and the unrelated pairs ($\hat{\beta} = 69.37, SE = 0.47, t = 147.41, p < .001$). Further, similarity and association scores within concept sets were analyzed with a paired *t*-test. Standardized similarity and association scores were calculated for the taxonomic and thematic pairs of each set—for taxonomic pairs, similarity rating *z*-scores (i.e., between-subjects normalization) were calculated and subtracted by thematic rating *z*-scores; the same process was used for the thematic pairs except that standardized association ratings were subtracted by standardized similarity ratings. This process effectively creates a measure of how much more similar or associated the semantic pairs are within a set. These similarity and association magnitude values did not produce a reliable difference ($M_{Difference} = 0.03 SD$) between the similarity scores of the taxonomic pairs and the association scores of the thematic pairs, $t(79) = 0.96, p = .34$. Thus, we cannot conclude that the concepts sets had more associated or more similar pairs (see Fig. 4). The complete similarity and association rating data is provided in Appendix B. Please see Honke and Kurtz (2019) for further description of the similarity rating procedure.

The lexical and orthographic properties of the taxonomic and thematic targets in each concept set were also analyzed to determine if there were any systematic differences between the semantic relations. Paired *t*-tests confirm no differences in word length ($M_{Difference} = -0.06, t(79) = -0.22, p = .82$), word frequency ($M_{Difference} = 11.75, t(79) = 0.41, p = .68$), average frequency (per million) of orthographic neighbors ($M_{Difference} = 61.6, t(79) = 1.12, p = .27$) and average frequency of the constrained bigrams for the wordforms ($M_{Difference} = 26.72, t(79) = 0.11, p = .91$). Lexical and orthographic statistics are provided in Appendix C. Orthographic statistics were drawn from the MCWord database (Medler and Binder, 2005) and the word frequency data came from the Shaoul and Westbury (2006) USENET corpus.

8.3. Reading and language exposure assessment

The reading and language exposure assessment data are presented in Fig. 5. Recall that exposure to print was measured with *d'*, where higher values indicate more success in identifying real magazines and authors while rejecting fake magazines and authors. The verbal fluency task was to name as many members of a category as possible in 60 s. This produced a verbal fluency score calculated by averaging the number of

distinct fruits, animals, and colors that were named in the time allotted. The GRE vocabulary assessment was a 30 item fill-in-the-blank task that was scored as a proportion correct. As mentioned above, data for one participant's verbal fluency task and one participant's exposure to print task were missing. These values were imputed in R with the mi package (Su et al., 2011).⁴ The median values from 8000 hypothetical value estimations (80 trials \times 100 hypothetical datasets) replaced the missing data points. The results of the reading and language exposure assessments are presented in Table 3. All of the measures were normally distributed according to Shapiro–Wilk tests.

8.4. Triad similarity judgment task

Similarity Judgments in the Triad Task. The taxonomic pair was selected 56.7% of the time (mean range by participant: 12.5%–98.75%)—a lower frequency of taxonomic responses than what is needed to conclude that there was a reliable taxonomic bias at the participant level (see Fig. 6). Binomial tests were conducted to classify each participant as taxonomic, thematic or ambiguous in their responding. The process resulted in 22 taxonomic biased responders, 22 thematic biased responders, and 14 ambiguous responders. When these frequency statistics are analyzed in a binomial exact test, the result is that people produce a taxonomic (or thematic) bias less frequently than would be expected by chance ($p = .087$), though this test was only marginally significant.

Response Time in the Triad Task. Overall, taxonomic matches were completed faster than thematic matches ($\hat{\beta} = 0.256, SD = 0.10, t = 2.465, p = .018$) but this effect is not found when outliers are removed ($\pm 2.5 SD$; $p = .11$). Consistent with the observation that faster responding is found for the semantic relationship that is preferred or sought out (Honke and Kurtz, 2019), people with a taxonomic bias were faster on trials where the taxonomic pair was chosen, $\hat{\beta} = -0.92, SE = 0.20, t = -4.651, p < .001$, and people with a thematic responding bias or ambiguous response preference were faster on thematic trials, $\hat{\beta} = -0.14, SE = 0.05, t = -3.032, p = .006$ and $\hat{\beta} = -0.16, SE = 0.07, t = -2.426, p = .031$, respectively. This response bias timing effect was resilient to outlier exclusion ($\pm 2.5 SD$ s).

Similarity Judgments and Reading and Language Exposure. *General Relationship between Similarity Judgments and Reading and Language Exposure.* While they had clear importance for the ERP measurement goals of the study, it was less clear how these measures might relate to similarity judgment behavior. A series of regression models were built to examine this relationship. A simple GLM built to predict taxonomic responding at the trial level that included trial and all three reading and language exposure measures uncovered reliable effects of all predictors ($ps < .001$).

A different pattern emerges when the data are analyzed with mixed effects (taking participant and concept set into account). A GLMER model built to predict trial-level taxonomic responding with fixed effect predictors for each of the reading and language exposure assessments and trial, and random effects (random intercepts for participant and concept set, and random slopes for trial) produced a reliable effect of trial, $\hat{\beta} = -0.16, SE = 0.07, t = -2.426, p = .031$ (see Fig. 7); no other reliable effects were found and allowing the terms to interact did not change this overall pattern.⁵ This result provides a replication of a newly

⁴ Parameters for the missing values were estimated at the trial level with data from the triad task and the reading and language exposure assessments (i.e., participant, trial number, concept set, trial response, response time, mean verbal fluency, exposure to print *d'* and GRE vocabulary accuracy). The ERP data were excluded from the imputation procedure due to extreme processing requirements.

⁵ Model specification: $match \sim response.bias \times exposure.to.print \times verbal.fluency \times vocab.accuracy + trial + (1 + trial|pid) + (1|concept.set)$.

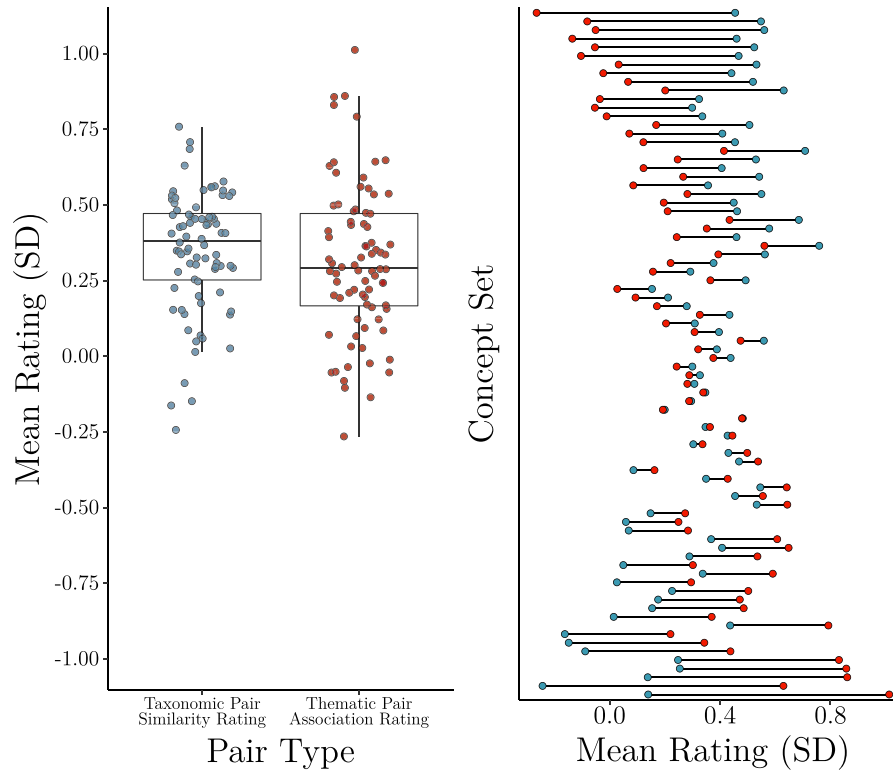


Figure 4. Visualization of the concept ratings overall (left) and paired with the match from the same concept set (right). The left panel depicts the mean similarity and association ratings for the taxonomic and thematic pairs, respectively. The right panel depicts the similarity (blue) and association (red) ratings paired for each concept set.

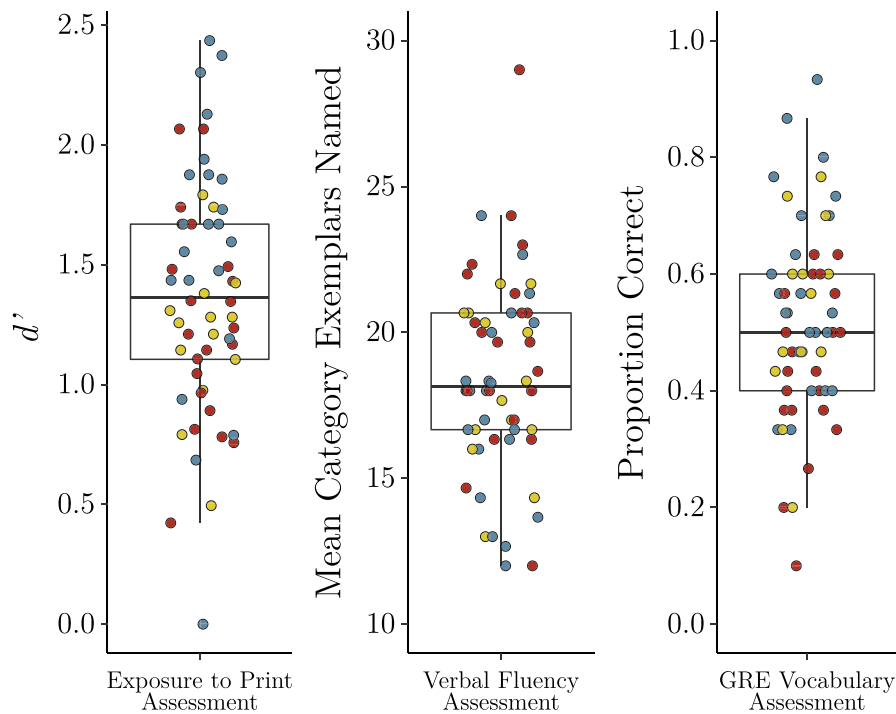


Fig. 5. Boxplots and individual data for each of the reading and language exposure tasks. Blue, red, and yellow points present people with taxonomic, thematic or ambiguous responding preferences, respectively. The data were normally distributed with no obvious outliers. Exposure to Print and GRE Vocabulary were positively related to taxonomic responding and Verbal Fluency was negatively related to taxonomic responding. (For interpretation of the references to color in this figure legend, the reader is referred to the Web version of this article.)

discovered pattern of behavior where people increase their taxonomic responding across the time-course of the experimental session (Honke and Kurtz, 2019).

No reliable differences were uncovered for exposure to print, verbal

Table 3
Behavioral descriptives.

Responding Bias	Taxonomic Responding Mean (Med.)	Exposure to Print <i>d'</i> Mean (Med.)	Verbal Fluency Mean (Med.)	GRE Vocabulary Mean Accuracy (Med.)	Pseudoword Identification Accuracy (Med.)
Taxonomic	.88 (.89)	1.57 (1.67)	17.56 (18)	.58 (.55)	.93 (76)
Ambiguous	.48 (.47)	1.23 (1.27)	18.19 (18)	.53 (.55)	.88 (73)
Thematic	.31 (.33)	1.27 (1.22)	19.53 (19.67)	.44 (.45)	.89 (73)
Mean Total	.56 (.56)	1.36 (1.39)	18.43 (18.56)	.52 (.52)	.90 (74)

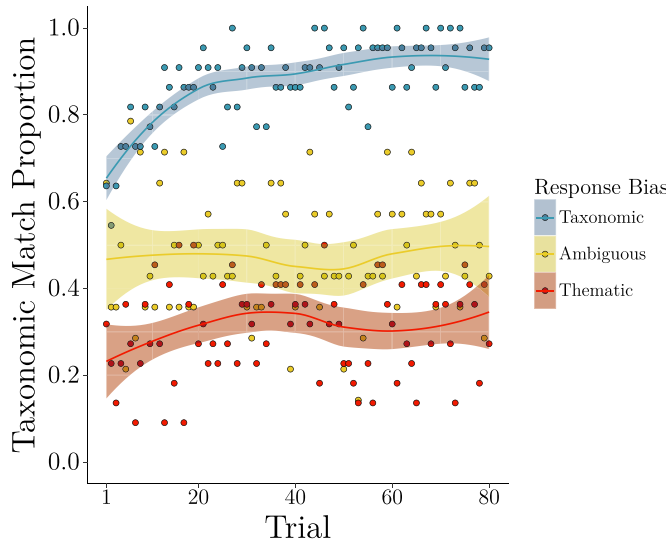


Fig. 6. Taxonomic responding frequency across trials. Points represent mean taxonomic responding by trial for response bias type.

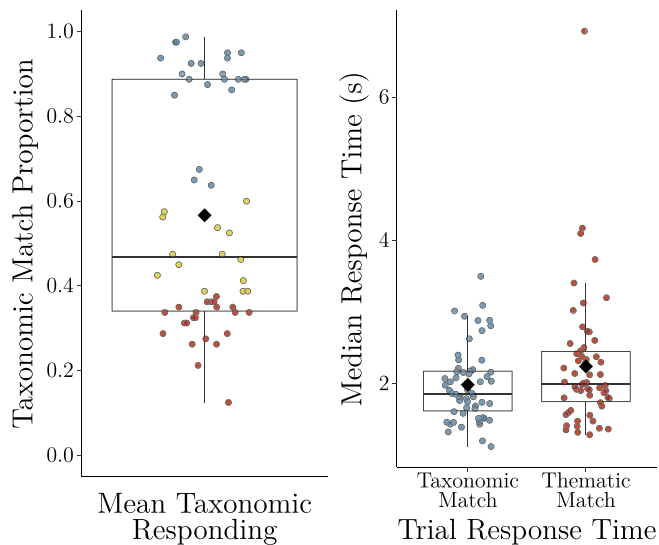


Fig. 7. Boxplots present mean taxonomic responding (left panel) and median response time for taxonomic and thematic matches (right panel) from the triad task. Individual points present participant means and medians. Diamonds present overall means. More taxonomic responding was found overall but there was no participant-level response bias majority. Trials with a taxonomic match were generally completed faster than thematic trials but the reliability of this effect turns on 2 near-outliers.

fluency or GRE vocabulary when they were analyzed in isolation (ETPWald $Z = 1.378$, $ETPp = 0.168$; VFwald $Z = -1.235$, $VFp = .22$). GRE vocabulary accuracy did approach significance as a predictor of ERP amplitude ($\hat{\beta} = 2.531$, $SE = 1.30$, Wald $Z = 1.945$, $p = .052$).⁶ See Table 3 for descriptive statistics.

The simplest explanation for the conflicting results between the simple and mixed-effects models is that people differ in similar ways in responding preferences and reading and language exposure. When the random intercept term for participant is included, this similarity is accounted for and adding predictors for the specific measures does not address significantly more variance. It is not safe to conclude that the simple GLM produced a spurious relationship between these variables, but the current results are not strong enough to make conclusions about how the predictors relate to similarity judgments overall. The patchy or bimodal distribution of mean taxonomic responding (Fig. 7) could also be playing a role in the failure to find reliable effects with the mixed-effects approach.

Individual-based Relationship between Similarity Judgments and Reading and Language Exposure. Since the overall relationship between the survey measures and taxonomic responding frequency is not clear, it might be more informative to look at this relationship with the inclusion of the response bias of each participant. In line with a central hypothesis of this paper—ERP differences are detectable between participants but not in the aggregate—it is possible that differences in the survey measures are also obscured when response bias is not accounted for. This is what was found. Note that the comprehensive model including response bias, trial and the survey measures often failed to converge. A fairly safe conclusion, then, is that response bias and the survey measures interact. When the model did converge (i.e., after many additional iterations and the use of the Nelder–Mead optimizer), this interaction was consistently reliable for the difference between taxonomic and ambiguous responding groups (e.g., $\hat{\beta} = 6.30$, $SE = 2.49$, Wald $Z = 2.534$, $p = .011$). The parameter estimates for the interaction between the taxonomic and thematic responding groups, however, were quite volatile across model initializations, $p \approx .002-.4$.

We also conducted the taxonomic responding analysis within each response bias group (as opposed to using response bias as a predictor). Again, these analyses were plagued with convergence failures. Nevertheless, an interesting pattern emerged that is worthy of mentioning even under this caveat. It was found that the survey measures and their interaction predicted taxonomic responding for the taxonomic ($ps < .001$) and ambiguous ($ps = .005-.028$) bias groups. No survey measure, however, was found to be reliable for the thematic bias group. Any conclusions taken from the results of these models should be made with extreme caution. We take this model behavior as evidence that the regression models suffer from overdispersion in the outcome variable, i. e., variability in trial-level responding that is not being sufficiently addressed by the predictors of these models.

Pseudoword Identification. The sole purpose of the pseudoword task was to confirm that participants were paying close attention to the word

⁶ The measure-isolated models only differed from the comprehensive model in that a single predictor was included from the reading and language exposure assessments (as opposed to all three measures).

stream during EEG collection, but it is possible that the ability to detect pseudowords is related to taxonomic—thematic processing (as was the case with the reading and language exposure assessments). Overall, participants did quite well in identifying pseudowords ($M = 72.2$; 90%). The correct identification of pseudowords was a reliable predictor of taxonomic responding, $\hat{\beta} = 0.08, SE = 0.03, \text{Wald } Z = 2.522, p = .012$. The analysis featured pseudoword identification and trial number as fixed effects, participant as a random intercept, trial as its random participant-level slope, and concept set as a random intercept. The effect was not reliable when the pseudoword accuracy predictor was included as a fixed-effect predictor in the mixed-effects model that included the reading and language exposure surveys (see Footnote 3.3.3 for the model specification save the fixed-effect pseudoword accuracy predictor).

8.5. Electrophysiological responses to taxonomic and thematic category members

Ideally, a comprehensive model of the EEG data (i.e., amplitude across time bins for the target channels) would be constructed that included all behavioral data and stimulus characteristics that have been collected and presented in this report, i.e., similarity judgments, reading and language exposure outcomes, and lexical and orthographic properties of the materials. Building and presenting a model with this level of complexity is prohibitive due to technical demands, difficulty of interpretation and increased false positive rate (Luck and Gaspelin, 2017). Consistent with the presentation of results thus far, the ERP analysis is divided to present specific aspects of the problem with models that address subsets of the possible predictors. First, a general analysis of the ERPs is presented that includes no similarity judgment data. The idea here is to start with a model similar to what has been used in past research to attempt to detect differences in ERP amplitude between taxonomic and thematic pairs across an entire sample. Next, a model of similarity judgment behavior, reading and language skill and orthographic and lexical variables is presented to determine if these factors predict unique variance in N400 amplitude. Finally, the simplest possible models of the relationship between similarity judgment behavior and N400 amplitude are presented.

General Properties of ERPs Elicited by Taxonomic and Thematic Category Members. We start with a comprehensive model of the ERPs without the effects of similarity behavior. An LMER model was built to predict average ERP amplitude at central–posterior electrode sites from lexical and orthographic characteristics, similarity and association rating difference scores and semantic pair type.⁷ The model uncovered reliable effects for semantic pair type but similarity ratings, word frequency, word length, orthographic neighborhood and bigram frequency were not reliable predictors.

In the aggregate, thematic category members elicited waveforms with more positive N400s than taxonomic category members ($\hat{\beta} = .337, SE = 0.048, t = 7.02$), and unrelated concepts ($\hat{\beta} = 0.205, SE = 0.055, t = 3.72$) when accounting for other sources of stimulus-based variance (see Fig. 8). Taxonomic category members elicited more negative N400s than unrelated category members ($\hat{\beta} = -0.132, SE = 0.055, t = -2.38$).

Similarity Judgments, Reading and Language Exposure and ERP Waveforms. As mentioned above, it is difficult to specify a single model that can comprehensively assess the contributions of the predictors in this design; a comprehensive analysis would include a series of models

⁷ The model predicted ERP amplitude from un-averaged, trial-level data at MiCe, MiPa, LDPa, RDPa, LMOc, RMOc, LLOc, RLOc, and MiOc with the following model specification: N400 amplitude \sim similarity.rating + frequency + length + orthographic.neighborhood + bigram.frequency + pair.type + (1 + time.window|participant) + (1|word.stimulus).

referenced to different combinations of the categorical predictors, including a large number of predictor terms in each. Therefore, we started by constructing a model that included all of the predictor terms necessary to address a question not answerable with fewer terms: Do the key variables of interest—similarity judgment behavior, reading and language exposure measures, and semantic pair type—interact to predict mean N400 amplitude while accounting for the variance of task engagement (pseudoword identification accuracy) and concept properties (similarity ratings, length, frequency, bigram frequency, orthographic neighborhood size); the random effects structure and target electrode sites were identical to the previous model. If so, further investigation of these effects would be warranted. In other words, a reliable interaction between these variables would help to validate the use of less sophisticated models without the concern that, for example, reading and language exposure can explain the effect. Reliable interactions in this general model⁸ would provide evidence against the interpretation that N400 amplitude differences are not (at least partially) related to similarity judgment behavior in the triad task.

The baseline reference levels for the analysis were taxonomic pairs for the semantic pair type variable and taxonomic responding bias for the response bias variable. A reliable interaction (exposure to print $d' \times$ verbal fluency mean \times vocabulary assessment accuracy \times response bias \times semantic pair type) was found for all pair type by response bias combinations. The variables interacted to reliably predict amplitude differences between the taxonomic and thematic bias group for taxonomic pairs vs. thematic pairs ($\hat{\beta} = -0.724, SE = 0.30, t = -2.41$) and unrelated pairs ($\hat{\beta} = -0.785, SE = 0.24, t = -3.31$) and between the taxonomic and ambiguous bias group for taxonomic pairs vs. thematic pairs ($\hat{\beta} = 2.82, SE = 0.39, t = 7.28$) and unrelated pairs ($\hat{\beta} = 4.32, SE = 0.31, t = 14.10$). The categorical reference level for semantic pair type was set to unrelated pairs to examine the effect of the interaction for unrelated and thematic pairs between the taxonomic and ambiguous bias groups. The interaction was found to be a reliable predictor of N400 amplitude, $\hat{\beta} = -1.50, SE = 0.31, t = -4.88$.

To address the remaining comparisons, the categorical reference levels for the model were set to thematic pairs and thematic response bias and the model was recalculated. The interaction was reliable between the thematic and ambiguous bias groups for thematic pairs vs. taxonomic pairs ($\hat{\beta} = -3.54, SE = 0.44, t = -8.10$) and unrelated pairs ($\hat{\beta} = 1.56, SE = 0.34, t = 4.52$). To analyze the final interaction effect for unrelated and taxonomic pairs between the ambiguous and thematic bias groups, the categorical semantic pair type reference level was set to unrelated pairs and the analysis was repeated. This interaction was also reliable, $\hat{\beta} = -5.10, SE = 0.34, t = -14.703$.

To reiterate, the interaction of similarity judgment behavior, reading and language ability assessments and semantic pair type was found to reliably predict N400 amplitude differences for every response bias–semantic pair type comparison. Similarity ratings, word length, word frequency, orthographic neighborhood, bigram frequency and pseudoword identification were not reliable predictors in the model.

Closer Examination of ERPs and Similarity Judgment Behavior. The models above suggest that similarity judgments and reading and language ability interact to predict differences in N400 amplitude across semantically related and unrelated concept pairs. A critical question that remains unresolved is how *exactly* these variables affect ERP amplitude. Models were built that held the categorical semantic pair type and response bias variables constant to determine (1) how semantic pairs differed in ERP amplitude within response bias groups and (2) how

⁸ The model structure was specified as: N400 amplitude \sim similarity.rating + length + frequency + orthographic.neighborhood + bigram.frequency + pseudoword.accuracy + pair.type \times response.bias \times exposure.to.print \times verbal.fluency \times vocabulary.accuracy + (1 + time.window|participant) + (1|word.stimulus).

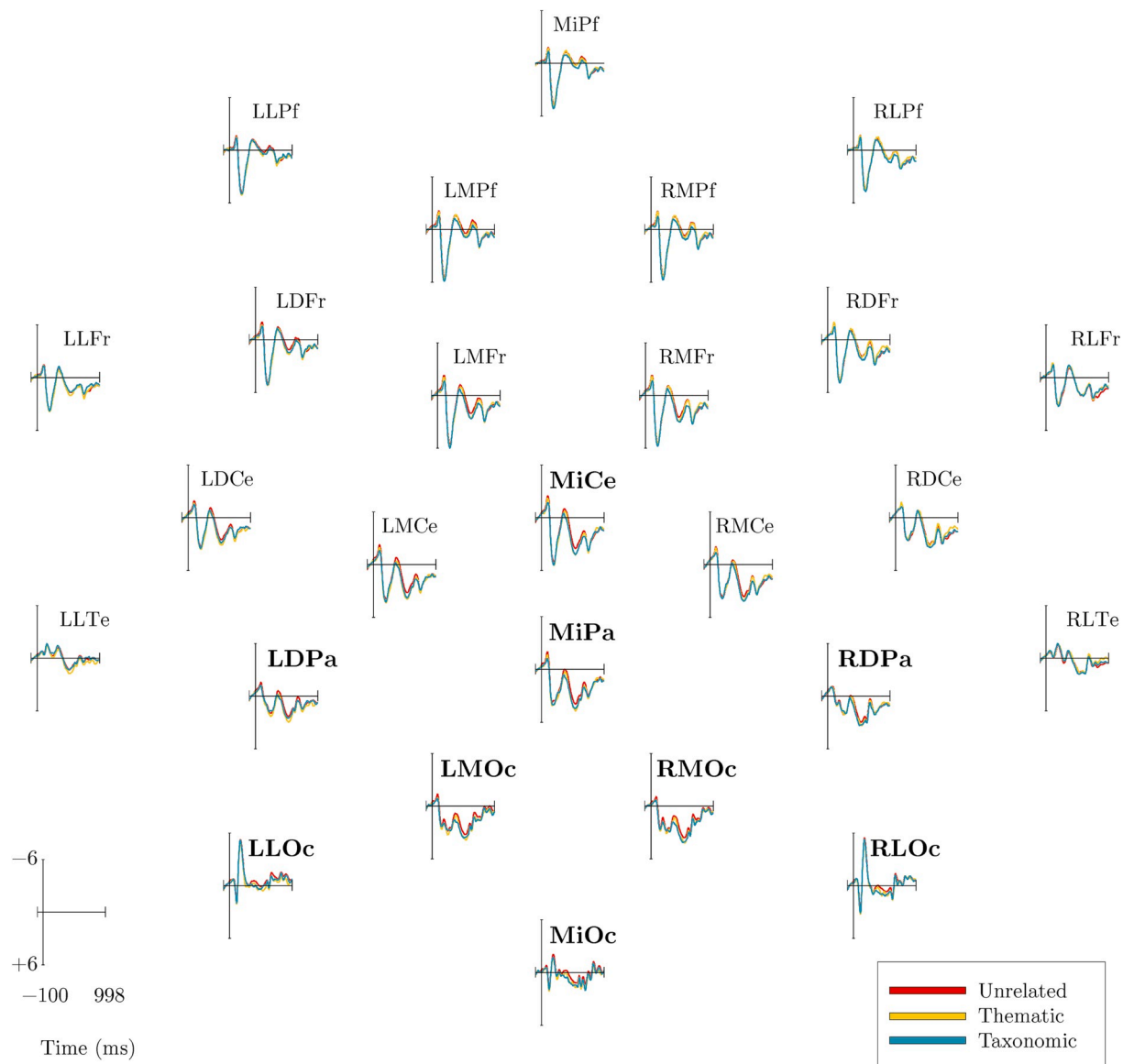


Fig. 8. Grand averaged ERP waveforms elicited in response to taxonomic, thematic and unrelated word pairs (pseudoword trials excluded). Unrelated, thematic and taxonomic pairs are presented in red, yellow and blue, respectively. The data are presented baselined and filtered with bandpass filtering at 0.1–20 Hz. Target electrode sites in bold. (For interpretation of the references to color in this figure legend, the reader is referred to the Web version of this article.)

response bias groups differed in ERP amplitude for each semantic pair (the random effects structure and target electrode sites were identical to the previous models). First, models built for each response bias group are presented to examine the differences between semantic pair types. Second, models built to examine differences across the response bias groups for each semantic pair type are presented. A depiction of these effects is presented in Fig. 9.

Semantic Pair Differences within Response Bias Groups. The mean amplitude of ERPs elicited by semantically related and unrelated pairs in the 300–400 ms time window was analyzed within each response bias group (taxonomic, thematic and ambiguous) with LMER.⁹ The goal of this analysis was to determine how the elicited waveforms of semantic pair types differed for people who produced ambiguous responding, majority taxonomic responding and majority thematic responding. The results showed that people who made more taxonomic matches in the

triad task also produced N400s that were reliably different for taxonomic and thematic pairs ($\hat{\beta} = -0.967, SE = 0.08, t = -12.13$), taxonomic and unrelated pairs ($\hat{\beta} = -0.24, SE = 0.10, t = -2.415$) and thematic and unrelated pairs ($\hat{\beta} = 1.21, SE = 0.10, t = 12.16$). People who produced more thematic matches in the triad task produced different N400s for thematic and unrelated pairs ($\hat{\beta} = 0.17, SE = 0.09, t = 2.01$), marginally different N400s ($p \approx .077$) between taxonomic and unrelated pairs ($\hat{\beta} = 0.151, SE = 0.09, t = 1.77$), and no difference between taxonomic and thematic pairs ($t = 0.29$). People who did not produce a reliable match preference (ambiguous responders) followed this same general pattern, no difference between taxonomic and thematic pairs ($t = -0.199$), but differences between unrelated pairs and thematic ($\hat{\beta} = -0.275, SE = 0.13, t = -2.245$) and taxonomic ($\hat{\beta} = -0.294, SE = 0.13, t = -2.099$) category members.

To sum, taxonomic responders were the only group to produce ERP waveforms that were reliably different for taxonomic and thematic pairs. Thematic and ambiguous responders only showed evidence of differentiation between semantically related and unrelated words (and

⁹ Simple semantic pair model (for each response bias group): $\text{amplitude} \sim \text{pair.type} + (1 + \text{time.window}|\text{participant}) + (1|\text{word.stimulus})$.

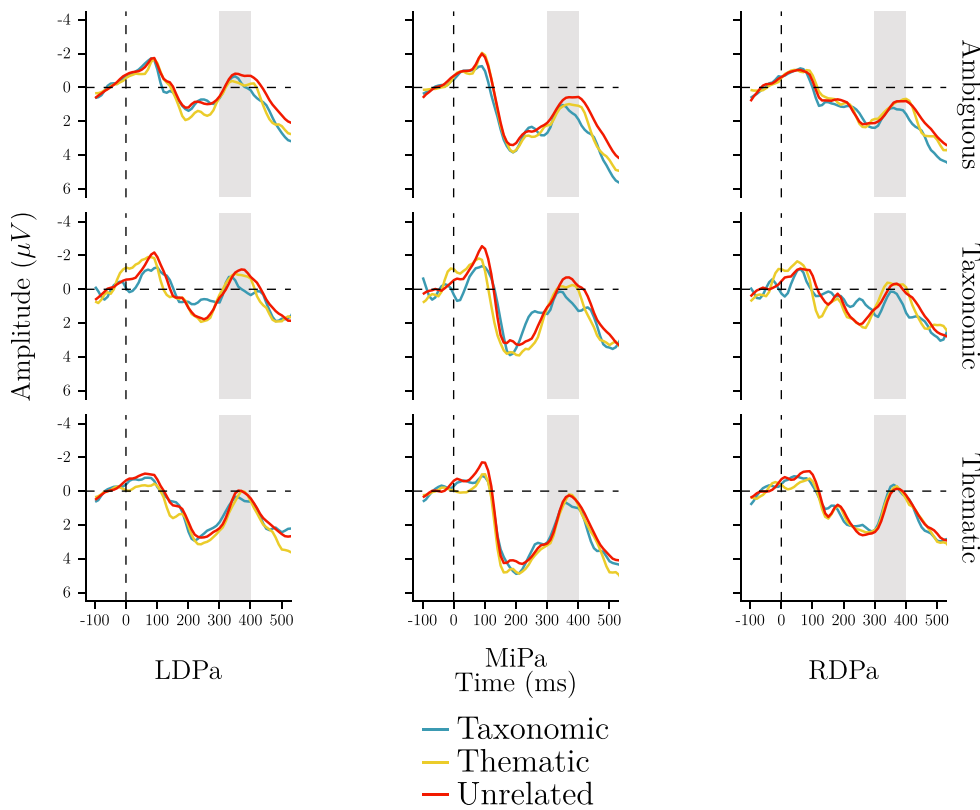


Fig. 9. ERPs elicited from taxonomic, thematic, and unrelated word pairs. Horizontally-aligned panels present response bias groups. Vertically-aligned panels present data from LDPa, MiPa and RDPa. N400s elicited by taxonomic and thematic pairs were reliably different for the taxonomic bias group only, i.e., the group that produced reliably more taxonomic responding in the similarity judgment task was the only group to produce reliably different N400s for taxonomic and thematic pairs.

only marginally so in the case of the taxonomic–unrelated comparison for the thematic bias group).

Response Bias Differences within Semantic Pairs. Similar to the previous analysis, LMER models were built that held one component constant (semantic pair type) to examine possible differences across the other factor, i.e., comparing amplitude across response bias groups for each semantic pair type.¹⁰ No reliable differences were found across response bias groups; no response bias group produced N400s of different amplitude for any semantic pair type comparison. This shows that it's not a difference in raw amplitude between response groups that drives the relationships presented thus far. Instead, it is the within-subject difference that is driving the effect.

N400 Amplitude Predicted by Semantic Pairs and Taxonomic Responding Frequency. One possible issue with the analysis above is that the cutoff for being classified as having a particular bias (α) is an arbitrary criterion—it turns on the difference between 49 ($p = .056$) and 50 ($p = .033$) consistent responses in an 80 trial experiment. Motivated by these concerns, a final set of models was constructed where mean amplitude for the N400 component was predicted by the interaction of semantic pair type and the proportion of taxonomic responses produced in the triad task (with random effects structures and electrode sites identical to the models above). The models uncovered a reliable interaction between taxonomic match proportion and semantic pair type where taxonomic responding produced reliable differences for the comparison of taxonomic pairs to thematic pairs ($\hat{\beta} = -1.25, SE = 0.14, t = -9.09$) and thematic pairs to unrelated pairs ($\hat{\beta} = -1.16, SE = 0.11, t = -10.61$), but not taxonomic and unrelated pairs ($t = 0.85$). Thus, taxonomic matching reliably interacted with pair type to predict amplitude differences between taxonomic and thematic pairs and thematic and

unrelated pairs.

9. General discussion

9.1. Summary of Results

The results of the EEG analyses show that taxonomic and thematic category members produce N400s with reliably different amplitudes at the group level. Similarity judgment behavior predicts N400 amplitude differences at the individual level. Reading and language ability (as measured by the exposure to print, verbal fluency and vocabulary assessments) predicts similarity judgment behavior. All of these variables predict unique N400 amplitude variance—similarity judgment behavior remains a reliable predictor and interacts with reading and language ability to predict N400 amplitude differences for taxonomic, thematic, and unrelated word pairs.

Looking closer at the specific relationship between similarity judgments and mean N400 amplitude, we found that people who produced particular response biases differed in systematic ways. The taxonomic bias group produced N400s that were reliably different for taxonomic and thematic pairs. This difference was not found in the thematic and ambiguous responding bias groups. The effect was also found when the response bias group variable was replaced with proportion of taxonomic responses. The results of this analysis suggest that people who show differences in their processing of taxonomic and thematic pairs are less likely to be subject to thematic intrusion and, per the confusability hypothesis, more likely to produce matches based on taxonomic similarity in the triad task.

Lastly, patterns were uncovered in the behavioral analysis of judgments and covariates that support several long-studied correlations with thematic intrusion on similarity processing. As in past work, reading and language ability could reliably predict similarity responding frequency at the group level. When the model included terms for individual-level variance *and* response bias, reliable effects were only found for the

¹⁰ Simple response bias group model (for each semantic pair type): amplitude \sim response.bias + (1 + time.window|participant) + (1|word.stimulus).

ambiguous and taxonomic response groups. Why is it that thematically-biased responders do not show the effect? This pattern of results could be taken as convergent evidence that the thematic bias is a lower-level consequence of semantic network organization and priority.

9.2. Conclusion

We set out to test two hypothesis: that (1) the failure to detect differences between ERPs elicited by taxonomic and thematic category members was caused by aggregating across individuals' response biases and (2) ERP waveforms elicited from an unbiased reading task could help support existing hypotheses for why human similarity judgments do not reliably conform to accepted theoretical accounts of psychological similarity.

The results provide support for the first hypothesis in that different patterns of N400 amplitude were found to relate to response bias groupings and taxonomic response frequency. Contrary to this hypothesis (and prior research), however, *general* differences in N400s elicited by taxonomic and thematic pairs were also found. This suggests that part of the problem in past studies could have been statistical power. In the present design, more participants were recruited in an attempt to sample adequately-sized groups with different similarity judgment biases; the size of each response bias group was comparable to the size of entire samples in studies in this research area. Sample size is likely more important for EEG studies on taxonomic and thematic categories because stimulus creation cannot be automated—the semantic pairs are hand-curated and cannot be procedurally generated online like other research domains—and the result is a smaller stimulus set than EEG investigations in other areas. Regardless of the general pattern, a novel conclusion of this work is that ERPs elicited by unbiased, passive reading of taxonomic and thematic category members reliably correspond with similarity judgments of those same concepts in the classic forced-choice triad task.

People who produced more taxonomic matches in the triad task produced N400s with reliably different amplitudes between taxonomic and thematic category members (Fig. 9). People who produced mostly thematic matches or responded ambiguously did not show this pattern; N400s elicited by thematic and taxonomic pairs in these groups only differed from unrelated pairs.

To our second hypothesis on the cause of thematic intrusion on human similarity judgments, the evidence supports many early conclusions in this research area. Supporting the claims about education and individual differences, reading and language exposure predicted similarity judgments and N400 amplitude. This work also generated evidence against the claim that apparent behavioral deviations from theoretical definitions of similarity can be attributed to the triad task. Finally, this work suggests that the dual-process integration hypothesis (as presented in Chen et al., 2013) is not an adequate explanation of the thematic intrusion effect. Increased taxonomic responding reliably predicts larger amplitude differences between semantic types. Taxonomic and thematic category members evoked reliably different N400 amplitude in the aggregate—results that contrast with the outcome and argument made in Chen et al. (2013) where a failure to find N400 differences was presented as evidence for the dual-process integration hypothesis. Note that our results suggest that samples that do not exhibit a taxonomic response bias will not produce differentiable waveforms. The similarity judgments for taxonomic and thematic pairs did not differ in Chen et al. (2013) and there is an interesting yet conflicting body of

evidence that triad responding is affected by cultural factors (Ji et al., 2004; Saalbach and Imai, 2007).

On the other hand, the confusability account remains viable as an explanation for individuals; susceptibility to thematic intrusion. Again, taxonomic responding predicted amplitude differences between taxonomic and thematic pairs even when accounting for reading and language exposure, engagement, and orthographic and lexical properties of the stimuli. We suggest that when the assessment of possible matches in the similarity judgment task resolved, some participants were better able to discriminate between sources of semantic relatedness, and those individuals were more likely to produce similarity-based responding. Ambiguous and thematic responders did not show a difference in facilitative priming between taxonomic and thematic pairs and more frequently produced thematic matches. This suggests that the cause of thematic responding in the triad task is not a preference for thematic thinking but rather less effective discrimination.

In sum, these results suggest that electrophysiological patterns elicited by the passive processing of semantically related and unrelated concept pairs are a reliable predictor of similarity judgment behavior. More reading and language skill (higher exposure to print *d'* and vocabulary assessment accuracy) predicts taxonomic matching and N400 amplitude differences. Even when accounting for individual differences, similarity judgment behavior remains reliable as a predictor of variance. We conclude that the tendency to produce fewer similarity-based matches in the triad task is directly tied to a lack of difference in facilitative priming between taxonomic and thematic pairs. ERPs that don't show differentiation between taxonomic and thematic category members are evidence of more difficulty in perceiving differences between taxonomic and thematic matches when making similarity judgments. Future work will focus on the extent to which these patterns of thought and behavior exhibit stability across testing sessions and the lifespan.

Author note

Correspondence concerning this article should be addressed to Garrett Honke, Department of Psychology, Binghamton University (SUNY), PO Box 6000, Binghamton University, Binghamton, NY 13902–6000. E-mail: ghonke1@binghamton.edu. This work was conducted at the Brain and Machine (BAM) and Learning and Representation in Cognition (LaRC) Laboratories at Binghamton University. It was supported by awards to SL from NSF CAREER-1252975, NSF TWC SBE-1422417, and NSF TWC SBE-1564046. SL and GH are now at X, The Moonshot Factory, Mountain View, California. The supplemental materials are hosted at <https://osf.io/ctzhk/> and include all experimental materials, code, and data to reproduce this document. We thank Dan Mirman, Vladimir Miskovic, Aira Domingo, Elizabeth Sacchi, the members of the LaRC and BAM Laboratories, and two anonymous reviewers for their comments on this article.

Author contribution statement

Dr. Honke devised the study with advice and support from Dr. Laszlo and Dr. Kurtz. Dr. Honke created the stimuli, programmed the studies, analyzed the data, and wrote the initial drafts of the manuscript. The data were collected in Dr. Laszlo's Brain and Machine Laboratory by graduate and undergraduate researchers in the BAM Lab. Dr. Laszlo and Dr. Kurtz provided invaluable comments and revisions during the entire process.

Appendix A. Concept Sets

Index	Standard	Taxonomic	Thematic	Unrelated	Unrelated	Pseudoword
1	CIGARETTES	ALCOHOL	LUNGS	CARPET	OUTLET	LURDUGE
2	WAITRESS	STEWARDESS	RESTAURANT	CALCIUM	SWAN	CHATAGHT
3	BEE	BUTTERFLY	HONEY	PLIERS	RECORD	INVOMBLY
4	TOOTHBRUSH	COMB	FLOSS	APE	GLASSES	RELEFUT
5	CUP	BOWL	TEA	BARBER	PHONE	SURNGE
6	SKI	SNOWBOARD	BOAT	FLOOR	STOMACH	WHICE
7	DOG	CAT	BONE	HOOD	POND	YOMECHED
8	RECEPTIONIST	HOSTESS	TELEPHONE	HAND	PARK	PAIT
9	RABBI	PASTOR	TEMPLE	DRIVEWAY	UNDERWEAR	SETIVITE
10	CABLE	CORD	TELEVISION	POT	ROCK	COSTEDED
11	GOAT	BUFFALO	FARM	CHALK	SKY	PINIER
12	FIELD	COURT	FLOWER	SCHOOL	TOAD	BEANERED
13	MINT	LOLLIPOP	BREATH	FELONY	STALLION	INYWERED
14	COOKIE	PIE	CHOCOLATE	FUR	WAVE	COLOUST
15	HORNET	WASP	STINGER	PADLOCK	RICE	BURTH
16	LAWNMOWER	SCISSORS	YARD	AUNT	BOMB	LEPELF
17	VINEYARD	ORCHARD	WINE	BEAD	DRIVER	ABOUE
18	PANDA	RACCOON	BAMBOO	LAW	WHIP	NUEENG
19	BEER	JUICE	PARTY	CARRIAGE	SHOP	LOYWED
20	SPOON	LADLE	SOUP	LION	STEREO	REIEMBLY
21	HORSE	PIG	GRASS	HOTEL	MUTANT	SUEPANED
22	CAMEL	ANTELOPE	DESERT	COFFIN	ENGINE	EATENDLY
23	BLANKET	COMFORTER	PILLOW	CUCUMBER	TAR	MOUNCTE
24	TURKEY	CHICKEN	STUFFING	LETTER	SQUARE	TOMSTED
25	SHOTGUN	PISTOL	SHELL	ARK	BELT	RERANING
26	PACKAGE	CRATE	DELIVERY	CHILD	TROUT	INTH
27	SHAMPOO	BLEACH	SHOWER	CIRCLE	PIGEON	REATOWER
28	TOE	FINGER	SANDAL	MARBLE	SPIKE	HARN
29	TRUCK	BUS	TRAILER	CACTUS	CLUB	AMILES
30	BICYCLE	CAR	HELMET	BASEMENT	SKIN	NOSTE
31	BOOTS	HEELS	SHOELACE	BALCONY	BRAIN	REARAROD
32	SAXOPHONE	HARP	JAZZ	HAIR	SODA	FOMPERED
33	OYSTER	SCALLOP	PEARL	BACTERIA	LEATHER	COSSENG
34	CRIB	BED	BABY	FERRY	PATIO	LEIGS
Index	Standard	Taxonomic	Thematic	Unrelated	Unrelated	Pseudoword
35	POLICE	FIREMAN	HANDCUFFS	CRAB	LAUNDRY	INYOPT
36	RABBIT	SQUIRREL	CARROT	BARBELL	MOTEL	TREARDE
37	MILK	LEMONADE	COW	GUITAR	WINDOW	REEROT
38	BOTTLE	CAN	INFANT	BERRY	CLOCK	YEVER
39	BIRD	BAT	NEST	CRIMINAL	PLAYGROUND	SHUR
40	ROCKET	MISSILE	ASTRONAUT	CHEESE	SINK	GERMAL
41	SHIP	CANOE	SAILOR	GLAND	UMBRELLA	STUTABLY
42	PLATE	TRAY	NAPKIN	ANKLE	CHAUFFEUR	COOWENUL
43	CROWN	HAT	KING	NOSE	SHOVEL	LERSE
44	HURRICANE	BLIZZARD	FLOOD	BADGE	FOSSIL	GAE Aid
45	LOCKER	CLOSET	JERSEY	PAINT	SPY	WAGHT
46	HEARSE	LIMOUSINE	GRAVEYARD	EYE	KITCHEN	SOLVY
47	NEEDLE	PIN	THREAD	HYDRANT	WRIST	LELICT
48	CELEBRITY	PLUMBER	FILM	FORTRESS	NECTAR	WARAENE
49	MONKEY	BEAR	BANANA	HAMMER	TOOTH	PRILY
50	OVEN	MICROWAVE	PAN	CONVICT	SCREEN	WOOUT
51	SKYSCRAPER	TOWER	ELEVATOR	HEART	HITCHHIKER	RUTISES
52	SURGEON	BUTCHER	KIDNEY	DYNAMITE	GALAXY	ISKERT
53	CHISEL	KNIFE	SCULPTURE	MIRROR	MIRROR	MEDERAN
54	SHOE	GLOVE	FOOT	TIGER	WALL	SUNICED
55	FOOTBALL	BASEBALL	QUARTERBACK	NECKLACE	PLANT	SWILUARY
56	ENVELOPE	PARCEL	STAMP	MUSCLE	YOGURT	FREANDE
57	JELLY	MARMALADE	JAR	BOOK	NAIL	ACHITIED
58	SALT	PEPPER	SEA	KNUCKLE	SAW	BERFFER
59	CASKET	BOX	GRAVE	JEWEL	STREET	HARY
60	FLY	ANT	WINGS	CEREAL	CONCRETE	VAVE
61	DOOR	GATE	KNOB	FLAG	LIQUID	VINS
62	PENGUIN	GOOSE	ICE	BRICK	HEAD	COMORVED
63	CAKE	DONUT	CANDLE	ACTRESS	BROCHURE	COREWAL
64	OWL	HAWK	MOON	CIRCUIT	DIARY	CHOURN
65	HOSE	TUBE	WATER	MOTHER	RODEO	FOVIND
66	SWEATER	HOODIE	MITTENS	BATHROOM	CHALKBOARD	MARMIGLY
67	SEDAN	BIKE	SEATBELT	COTTON	SHRIMP	FEEPPER
68	PENCIL	PEN	ERASER	FLUTE	SHEEP	HALY
Index	Standard	Taxonomic	Thematic	Unrelated	Unrelated	Pseudoword
69	BACKPACK	SUITCASE	NOTEBOOK	BUTTER	PAINTING	BROURD
70	SEAGULL	DUCK	PIER	BEDROOM	POWDER	SHERT
71	VENOM	POISON	SNAKE	GRAFFITI	RASPBERRY	TURICAFI
72	TORTILLA	BREAD	BEANS	COLD	WIRE	BREATED
73	COMPUTER	TABLET	MOUSE	ATHLETE	COUCH	CEEY

(continued on next page)

(continued)

Index	Standard	Taxonomic	Thematic	Unrelated	Unrelated	Pseudoword
74	CHAIR	SOFA	LEGS	ANCHOVY	BALL	AGATENG
75	BISCUITS	TOAST	GRAVY	DANCE	SNAIL	RENCTRY
76	FLOUR	CORNMEAL	DOUGH	BUTTON	SMOG	BEVERSS
77	SHIRT	BLOUSE	COLLAR	BRIDGE	POOL	QUMES
78	PATHWAY	SIDEWALK	GRAVEL	BABYSITTER	TYPEWRITER	SOOBRRARE
79	SNOW	RAIN	SLED	CEMETERY	NOVEL	KITSSSES
80	CITY	VILLAGE	AIRPORT	NECK	WHALE	SQUGED

Note: Unrelated words were only presented in the EEG recording phase.

Appendix B. Concept Set Ratings

Index	Standard	Taxonomic	Thematic	Unrelated	Unrelated	Thematic	Taxonomic	Tax.-Unr.	The.-Unr.	Tax.-The.
						Rating	Rating	Rating	Rating	Rating Difference
1	CIGARETTES	ALCOHOL	LUNGS	CARPET	OUTLET	-0.09	0.44	-0.38	-0.86	-0.53
2	WAITRESS	STEWARDESS	RESTAURANT	CALCIUM	SWAN	0.58	0.35	-1.04	-1.05	0.23
3	BEE	BUTTERFLY	HONEY	PLIERS	RECORD	0.35	0.43	-0.95	-0.92	-0.08
4	TOOTHBRUSH	COMB	FLOSS	APE	GLASSES	0.3	-0.05	-1.13	-1.16	0.35
5	CUP	BOWL	TEA	BARBER	PHONE	0.29	0.54	-0.76	-0.61	-0.25
6	SKI	SNOWBOARD	BOAT	FLOOR	STOMACH	0.21	0.09	-0.88	-1.01	0.12
7	DOG	CAT	BONE	HOOD	POND	0.14	1.01	-0.94	-0.9	-0.87
8	RECEPTIONIST	HOSTESS	TELEPHONE	HAND	PARK	0.69	0.43	-0.56	-0.73	0.25
9	RABBI	PASTOR	TEMPLE	DRIVEWAY	UNDERWEAR	0.46	-0.14	-1.18	-1.12	0.6
10	CABLE	CORD	TELEVISION	POT	ROCK	0.56	-0.05	-0.79	-0.86	0.61
11	GOAT	BUFFALO	FARM	CHALK	SKY	0.14	0.86	-0.99	-1.02	-0.72
12	FIELD	COURT	FLOWER	SCHOOL	TOAD	0.15	0.03	-0.92	-1	0.13
13	MINT	LOLLIPOP	BREATH	FELONY	STALLION	0.71	0.41	-0.78	-0.92	0.29
14	COOKIE	PIE	CHOCOLATE	FUR	WAVE	0.09	0.16	-1	-1.12	-0.08
15	HORNET	WASP	STINGER	PADLOCK	RICE	0.52	0.07	-0.94	-1.04	0.45
16	LAWN MOWER	SCISSORS	YARD	AUNT	BOMB	0.43	0.44	-1.02	-1.23	-0.02
17	VINEYARD	ORCHARD	WINE	BEAD	DRIVER	0.29	0.16	-0.92	-1.05	0.14
18	PANDA	RACON	BAMBOO	LAW	WHIP	0.55	0.28	-0.46	-0.33	0.27
19	BEER	JUICE	PARTY	CARRIAGE	SHOP	0.48	0.48	-0.32	-0.39	0
20	SPOON	LADLE	SOUP	LION	STEREO	0.49	0.37	-1.13	-0.9	0.13
21	HORSE	PIG	GRASS	HOTEL	MUTANT	0.25	0.86	-1.16	-1.13	-0.6
22	CAMEL	ANTELOPE	DESERT	COFFIN	ENGINE	-0.24	0.63	-1.03	-1.08	-0.87
23	BLANKET	COMFORTER	PILLOW	CUCUMBER	TAR	0.56	0.39	-1.02	-1.16	0.17
24	TURKEY	CHICKEN	STUFFING	LETTER	SQUARE	0.07	0.28	-0.37	-0.4	-0.21
25	SHOTGUN	PISTOL	SHELL	ARK	BELT	0.45	0.2	-0.97	-1.12	0.25
26	PACKAGE	CRATE	DELIVERY	CHILD	TROUT	0.03	0.3	-0.77	-0.77	-0.27
27	SHAMPOO	BLEACH	SHOWER	CIRCLE	PIGEON	0.06	0.25	-0.96	-0.9	-0.19
28	TOE	FINGER	SANDAL	MARBLE	SPIKE	0.35	0.34	-0.99	-0.92	0.01
29	TRUCK	BUS	TRAILER	CACTUS	CLUB	0.44	-0.02	-0.77	-0.94	0.46
30	BICYCLE	CAR	HELMET	BASEMENT	SKIN	0.37	0.61	-1.02	-1.13	-0.24
31	BOOTS	HEELS	SHOELACE	BALCONY	BRAIN	0.3	0.24	-0.98	-0.9	0.06
32	SAXOPHONE	HARP	JAZZ	HAIR	SODA	0.28	0.17	-1.16	-1.22	0.11
33	OYSTER	SCALLOP	PEARL	BACTERIA	LEATHER	0.15	0.27	-0.8	-0.85	-0.13
34	CRIB	BED	BABY	FERRY	PATIO	0.56	0.47	-0.76	-0.7	0.08
35	POLICE	FIREMAN	HANDCUFFS	CRAB	LAUNDRY	0.31	0.28	-0.96	-1.12	0.03
36	RABBIT	SQUIRREL	CARROT	BARBELL	MOTEL	0.55	0.64	-0.63	-0.68	-0.1
37	MILK	LEMONADE	COW	GUITAR	WINDOW	0.4	0.31	-0.75	-0.87	0.09
38	BOTTLE	CAN	INFANT	BERRY	CLOCK	0.53	0.64	-1.04	-1.15	-0.11
39	BIRD	BAT	NEST	CRIMINAL	PLAYGROUND	0.44	0.38	-0.78	-0.83	0.06
40	ROCKET	MISSILE	ASTRONAUT	CHEESE	SINK	0.35	0.36	-0.86	-0.6	-0.02
Index	Standard	Taxonomic	Thematic	Unrelated	Unrelated	Taxonomic	Thematic	Tax.-Unr.	The.-Unr.	Tax.-The.
						Rating	Rating	Rating	Rating	Rating Difference
41	SHIP	CANOE	SAILOR	GLAND	UMBRELLA	0.34	0.59	-1.08	-1.17	-0.25
42	PLATE	TRAY	NAPKIN	ANKLE	CHAUFFEUR	0.39	0.32	-0.99	-0.96	0.07
43	CROWN	HAT	KING	NOSE	SHOVEL	0.63	0.2	-1.14	-0.9	0.43
44	HURRICANE	BLIZZARD	FLOOD	BADGE	FOSSIL	0.47	-0.1	-0.84	-0.84	0.57
45	LOCKER	CLOSET	JERSEY	PAINT	SPY	0.76	0.56	-0.86	-0.69	0.2
46	HEARSE	LIMOUSINE	GRAVEYARD	EYE	KITCHEN	0.34	-0.01	-0.8	-0.68	0.35
47	NEEDLE	PIN	THREAD	HYDRANT	WRIST	0.45	0.12	-1.12	-0.98	0.33
48	CELEBRITY	PLUMBER	FILM	FORTRESS	NECTAR	0.54	0.27	-0.73	-0.84	0.28
49	MONKEY	BEAR	BANANA	HAMMER	TOOTH	0.15	0.49	-0.79	-1	-0.33
50	OVEN	MICROWAVE	PAN	CONVICT	SCREEN	0.31	0.2	-0.18	-0.83	0.1
51	SKYSCRAPER	TOWER	ELEVATOR	HEART	HITCHHIKER	0.51	0.17	-0.6	-0.48	0.34
52	SURGEON	BUTCHER	KIDNEY	DYNAMITE	GALAXY	0.3	0.29	-0.29	-0.6	0.01
53	CHISEL	KNIFE	SCULPTURE	HATCH	MIRROR	0.05	0.3	-1.06	-0.88	-0.25
54	SHOE	GLOVE	FOOT	TIGER	WALL	0.18	0.47	-1.08	-0.93	-0.3
55	FOOTBALL	BASEBALL	QUARTERBACK	NECKLACE	PLANT	0.36	0.09	-0.95	-1.02	0.27
56	ENVELOPE	PARCEL	STAMP	MUSCLE	YOGURT	-0.16	0.22	-0.54	-0.47	-0.38
57	JELLY	MARMALADE	JAR	BOOK	NAIL	0.45	0.55	-1.01	-0.68	-0.1

(continued on next page)

(continued)

Index	Standard	Taxonomic	Thematic	Unrelated	Unrelated	Taxonomic	Thematic	Tax.-Unr.	The.-Unr.	Tax.-The.
						Rating	Rating	Rating	Rating	Rating Difference
58	SALT	PEPPER	SEA	KNUCKLE	SAW	-0.15	0.34	-0.83	-0.7	-0.49
59	CASKET	BOX	GRAVE	JEWEL	STREET	0.45	-0.26	-0.78	-0.92	0.72
60	FLY	ANT	WINGS	CEREAL	CONCRETE	0.53	0.25	-1.11	-1.12	0.28
61	DOOR	GATE	KNOB	FLAG	LIQUID	0.41	0.07	-1.12	-1.16	0.34
62	PENGUIN	GOOSE	ICE	BRICK	HEAD	0.44	0.79	-0.67	-0.9	-0.36
63	CAKE	DONUT	CANDLE	ACTRESS	BROCHURE	0.41	0.65	-0.77	-0.61	-0.24
64	OWL	HAWK	MOON	CIRCUIT	DIARY	0.33	0.29	-0.95	-0.77	0.04
65	HOSE	TUBE	WATER	MOTHER	RODEO	0.41	0.12	-1	-0.99	0.28
66	SWEATER	HOODIE	MITTENS	BATHROOM	CHALKBOARD	0.53	0.03	-0.95	-1	0.5
67	SEDAN	BIKE	SEATBELT	COTTON	SHRIMP	0.2	0.19	-0.9	-1	0.01
68	PENCIL	PEN	ERASER	FLUTE	SHEEP	0.38	0.22	-1.07	-1.04	0.16
69	BACKPACK	SUITCASE	NOTEBOOK	BUTTER	PAINTING	0.01	0.37	-1.09	-1.09	-0.35
70	SEAGULL	DUCK	PIER	BEDROOM	POWDER	0.25	0.83	-0.64	-0.44	-0.58
71	VENOM	POISON	SNAKE	GRAFFITI	RASPBERRY	0.43	0.33	-1.04	-1.1	0.11
72	TORTILLA	BREAD	BEANS	COLD	WIRE	0.43	0.5	-0.75	-0.96	-0.07
73	COMPUTER	TABLET	MOUSE	ATHLETE	COUCH	0.23	0.5	-0.58	-0.66	-0.28
74	CHAIR	SOFA	LEGS	ANCHOVY	BALL	0.47	0.54	-1.16	-1.18	-0.07
75	BISCUITS	TOAST	GRAVY	DANCE	SNAIL	0.3	0.34	-1.09	-1.17	-0.03
76	FLOUR	CORNMEAL	DOUGH	BUTTON	SMOG	0.32	-0.04	-1.1	-1.13	0.36
77	SHIRT	BLOUSE	COLLAR	BRIDGE	POOL	0.46	0.21	-0.52	-0.62	0.25
78	PATHWAY	SIDEWALK	GRAVEL	BABYSITTER	TYPEWRITER	0.52	-0.05	-0.98	-1.1	0.58
79	SNOW	RAIN	SLED	CEMETERY	NOVEL	0.46	0.24	-0.64	-0.85	0.22
80	CITY	VILLAGE	AIRPORT	NECK	WHALE	0.55	-0.08	-0.87	-0.82	0.63

Appendix C. Concept Set Properties

Index	Standard	Length		Frequency		Neighborhood		Bigram	
		Tax.	Them.	Tax.	Them.	Tax.	Them.	Tax.	Them.
1	CIGARETTES	7	5	18.7	15.3	0	1.1	229.3	219.2
2	WAITRESS	10	10	3.8	33.1	0	0	99.6	449.6
3	BEE	9	5	5.2	20.8	0	64.6	499.7	820.7
4	TOOTHBRUSH	4	5	5.7	1.2	150.2	2.4	1677.7	1114.3
5	CUP	4	3	30.7	89.5	3.7	45.6	1110.3	294.8
6	SKI	9	4	NA	55.6	0	15.3	129.5	5651.2
7	DOG	3	4	43.3	28.2	132.3	54.8	1462.1	1749.3
8	RECEPTIONIST	7	9	9.6	102.9	2.6	0.1	927.1	350.3
9	RABBI	6	6	3.6	24.5	0	0	714	1499.3
10	CABLE	4	10	8.2	104	58.2	0	2397.3	819.5
11	GOAT	7	4	7.3	69.4	0	68.2	137.8	1391.9
12	FIELD	5	6	128.1	28	80.8	5.6	3111.1	1716.4
13	MINT	8	6	0.4	57.9	0	5.9	253.5	572.3
14	COOKIE	3	9	12.9	13.4	21.1	0	138.4	253.6
15	HORNET	4	7	2.5	0.4	7.5	0.7	1275	1215.6
16	LAWN MOWER	8	4	4.5	37.6	0	49.1	262.3	1014.5
17	VINEYARD	7	4	5.5	75.6	0	50.1	315.5	4733.2
18	PANDA	6	6	NA	6.2	0	0	884.3	348.4
19	BEER	5	5	21.5	373.5	3.8	15.3	1854.3	1250.5
20	SPOON	5	4	1.2	20.6	0	16.5	800.5	1713.4
21	HORSE	3	5	18.7	87	26.1	25.1	211.7	1201.9
22	CAMEL	8	6	4	40.5	0	0	362.1	828.3
23	BLANKET	9	6	1.7	14.5	5.7	2.1	877.3	593.2
24	TURKEY	7	8	31.1	4.2	1.2	0.9	613.8	1769.5
25	SHOTGUN	6	5	15.1	29.7	1.6	66.5	586.2	3226

Index	Standard	Length		Frequency		Neighborhood		Bigram	
		Tax.	Them.	Tax.	Them.	Tax.	Them.	Tax.	Them.
26	PACKAGE	5	8	2.8	15.2	1.4	2.1	1072.2	602.4
27	SHAMPOO	6	6	1.9	18.1	3.4	58.9	396.9	2287.2
28	TOE	6	6	51.8	1.1	2.9	0.4	2021.8	898.6
29	TRUCK	3	7	65.1	3.2	597.9	2.8	2755.8	1242.9
30	BICYCLE	3	6	274.9	9.5	168	0.1	1786.5	666.7
31	BOOTS	5	8	19	0.4	8.1	0	765.8	337.7
32	SAXOPHONE	4	4	2.5	6.7	40.4	0	2758.5	80.3
33	OYSTER	7	5	1	5.4	0	3.9	241.3	1699.9
34	CRIB	3	4	254.4	191.2	42.7	1.2	484.3	811.3
35	POLICE	7	9	0.7	2.3	4	0	424.2	111.8
36	RABBIT	8	6	3.7	2.6	0	2.6	417.5	779.4
37	MILK	8	3	3	23.3	0	128.6	234.2	1566.7
38	BOTTLE	3	6	1954.3	21.4	95.6	0	2766	500.3
39	BIRD	3	4	10.5	13.6	280.3	94.9	502.9	2975.5
40	ROCKET	7	9	27.3	1	0.3	0	791.3	120.1

(continued on next page)

(continued)

Index	Standard	Length		Frequency		Neighborhood		Bigram	
		Tax.	Them.	Tax.	Them.	Tax.	Them.	Tax.	Them.
41	SHIP	5	6	3.9	5.9	4.3	1.4	432.4	464.7
42	PLATE	4	6	21	4.9	6	0	658.3	258.4
43	CROWN	3	4	54.5	91.7	409.1	59.7	4629.3	1483.4
44	HURRICANE	8	5	2.6	15.6	0	158.3	158.7	806.4
45	LOCKER	6	6	10.5	13	55.9	0	796.4	607.8
46	HEARSE	9	9	2.7	4	0	0	714.8	192.8
47	NEEDLE	3	6	13.6	11.2	13.6	64.3	111.8	866.3
48	CELEBRITY	7	4	2.1	76.5	1.3	47.8	788.3	1526.5
49	MONKEY	4	6	63.8	4.3	73.6	0	2940.5	481.2
50	OVEN	9	3	2.1	26.7	0	156.4	170.2	1730.2
51	SKYSCRAPER	5	8	49	8.9	41.1	0	2327.6	251.4
52	SURGEON	7	6	5.6	4.9	0.1	0	1415.6	433
53	CHISEL	5	9	38.8	22	0	0	266.2	208.8

Index	Standard	Length		Frequency		Neighborhood		Bigram	
		Tax.	Them.	Tax.	Them.	Tax.	Them.	Tax.	Them.
54	SHOE	5	4	4.9	101.1	6.8	30.4	793.6	1986.8
55	FOOTBALL	8	11	6.5	NA	0	0	174.6	30.3
56	ENVELOPE	6	5	8.4	13.8	0	2.7	756	1207.1
57	JELLY	9	3	2.6	11.8	0	85.3	188.9	651.5
58	SALT	6	3	7	166	0.7	152.1	1990.7	1001.7
59	CASKET	3	5	78.8	31.2	23.7	8.6	232.8	1091.1
60	FLY	3	5	4	29.6	4303.6	6.5	14878.9	464.6
61	DOOR	4	4	50.9	3.7	61.9	381.9	768	957.7
62	PENGUIN	5	3	6.2	54.4	11.9	4.1	1854.2	61.6
63	CAKE	5	6	NA	8	0	21.7	943.2	1682.9
64	OWL	4	4	4.2	54.8	1	31.7	2092.1	3092.8
65	HOSE	4	5	15.2	447.9	5.2	55.2	125.1	3313.5
66	SWEATER	6	7	NA	0.8	0.3	3.3	286.5	1086.7
67	SEDAN	4	8	8.3	NA	177.4	0	1975.9	350
68	PENCIL	3	6	19.8	0.3	64	0.8	702.9	1848.7
69	BACKPACK	8	8	13	7.7	0	0	363	213.4
70	SEAGULL	4	4	9.9	5.8	9.9	3.1	1469.1	915.8
71	VENOM	6	5	12.6	15.1	66.6	6.8	838.4	147.8
72	TORTILLA	5	5	77	18.3	30.2	28.2	1289.3	1902.1
73	COMPUTER	6	5	2.9	8.4	16.2	71.3	871.7	3653
74	CHAIR	4	4	21.4	117.7	32	32.9	989	610.6
75	BISCUITS	5	5	15.4	3.9	20.7	31.2	1086.3	863.3
76	FLOUR	8	5	NA	10.9	0	15.2	678.5	2330.5
77	SHIRT	6	6	8.9	19.1	0	12.5	1079.5	900
78	PATHWAY	8	6	6.2	11	0	14.4	101.8	587.3
79	SNOW	4	4	74.2	0.8	35.2	11	1825	554.8
80	CITY	7	7	140	53.8	0.4	0	706.4	389.5

Appendix D. Rating Task

Consider how similar these items are.

WORD_1

WORD_2

Rate how similar the items are below.

I don't know
this word:
>[WORD_1]

I don't know
this word:
>[WORD_2]

Fig. D1. Figure presents a depiction of the similarity rating task. Participants were allowed to choose any point on the rating line to provide their rating. Association rating task not pictured.

Appendix E. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.neuropsychologia.2020.107388>.

References

- Barr, D.J., Levy, R., Scheepers, C., Tily, H.J., 2013. Random effects structure for confirmatory hypothesis testing: keep it maximal. *J. Mem. Lang.* 68 (3), 255–278.
- Bassok, M., Medin, D.L., 1997. Birds of a feather flock together: similarity judgments with semantically rich stimuli. *J. Mem. Lang.* 36 (3), 311–336.
- Bates, D., Maechler, M., Bolker, B., Walker, S., 2014. lme4: linear mixed-effects models using eigen and s4. R package version 1 (7), 1–23.
- Brainard, D.H., Vision, S., 1997. The psychophysics toolbox. *Spatial Vis.* 10, 433–436.
- Brysbaert, M., Warriner, A.B., Kuperman, V., 2014. Concreteness ratings for 40 thousand generally known English word lemmas. *Behav. Res. Methods* 46 (3), 904–911.
- Cacioppo, J.T., Petty, R.E., 1982. The need for cognition. *J. Pers. Soc. Psychol.* 42, 116–131.
- Chen, Q., Li, P., Xi, L., Li, F., Lei, Y., Li, H., 2013. How do taxonomic versus thematic relations impact similarity and difference judgments? An ERP study. *Int. J. Psychophysiol.* 90 (2), 135–142.
- Chen, Q., Ye, C., Liang, X., Cao, B., Lei, Y., Li, H., 2014. Automatic processing of taxonomic and thematic relations in semantic priming—differentiation by early N400 and late frontal negativity. *Neuropsychologia* 64, 54–62.
- Denney, N.W., 1974. Evidence for developmental changes in categorization criteria for children and adults. *Hum. Dev.* 17 (1), 41–53.
- Estes, Z., Golonka, S., Jones, L.L., 2011. 8 thematic thinking: the apprehension and consequences of thematic relations. *Psychology of Learning and Motivation—Advances in Research and Theory* 54, 249–294.
- Geller, J., Landrigan, J.-F., Mirman, D., 2019. A pupillometric examination of cognitive control in taxonomic and thematic semantic memory. *J. cognition* 2 (1).
- Gentner, D., Brem, S.K., 1999. Is snow really like a shovel? distinguishing similarity from thematic relatedness. In: Hahn, M., Stones, S.C. (Eds.), *Proceedings of the 21st Annual Meeting of the Cognitive Science Society*. Erlbaum, Mahwah, NJ, pp. 179–184.
- Greenfield, D.B., Scott, M.S., 1986. Young children's preference for complementary pairs: evidence against a shift to a taxonomic preference. *Dev. Psychol.* 22 (1), 19–21.
- Hagoort, P., Brown, C.M., Swaab, T.Y., 1996. Lexical–semantic event-related potential effects in patients with left hemisphere lesions and aphasia, and patients with right hemisphere lesions without aphasia. *Brain* 119 (2), 627–649.
- Honke, G., 2017. Causes and Predictors of Thematic Intrusion on Human Similarity Judgments (Doctoral Dissertation).
- Honke, G., Kurtz, K.J., 2019. Similarity is as similarity does? a critical inquiry into the effect of thematic association on similarity. *Cognition* 186, 115–138.
- Ince, E., Christman, S.D., 2002. Semantic representations of word meanings by the cerebral hemispheres. *Brain Lang.* 80 (3), 393–420.
- Jenkins, J.J., Russell, W.A., 1952. Associative clustering during recall. *J. Abnorm. Soc. Psychol.* 47 (4), 818–821.
- Ji, L.-J., Zhang, Z., Nisbett, R.E., 2004. Is it culture or is it language? examination of language effects in cross-cultural research on categorization. *J. Pers. Soc. Psychol.* 87 (1), 57–65.
- Kacmador, M., Kelleher, J.D., 2019. Capturing and measuring thematic relatedness. *Comput. Humanit.* 1–38.
- Khateb, A., Michel, C.M., Pegna, A.J., O'Dochartaigh, S.D., Landis, T., Annoni, J.-M., 2003. Processing of semantic categorical and associative relations: an ERP mapping study. *Int. J. Psychophysiol.* 49 (1), 41–55.
- Kuperman, V., Stadthagen-Gonzalez, H., Brysbaert, M., 2012. Age-of-acquisition ratings for 30,000 English words. *Behav. Res. Methods* 44 (4), 978–990.
- Kurtz, K.J., Gentner, D., 2001. Kinds of kinds: sources of category coherence. In: Moore, J., Keith, S. (Eds.), *Proceedings of the 23rd Annual Conference of the Cognitive Science Society*. Erlbaum, Mahwah, NJ, pp. 522–527.
- Kurtz, K.J., Miao, C.-H., Gentner, D., 2001. Learning by analogical bootstrapping. *J. Learn. Sci.* 10 (4), 417–446.
- Kutas, M., Federmeier, K.D., 2011. Thirty years and counting: finding meaning in the N400 component of the event-related brain potential (ERP). *Annu. Rev. Psychol.* 62, 621–647.
- Kutas, M., Hillyard, S.A., 1980. Event-related brain potentials to semantically inappropriate and surprisingly large words. *Biol. Psychol.* 11 (2), 99–116.
- Kuznetsova, A., Brockhoff, P.B., Christensen, R.H.B., 2015. Package lmerTest, 2. R package version, 0.
- Laszlo, S., Federmeier, K.D., 2011. The N400 as a snapshot of interactive processing: evidence from regression analyses of orthographic neighbor and lexical associate effects. *Psychophysiology* 48 (2), 176–186.
- Laszlo, S., Ruiz-Blondet, M., Khalifian, N., Chu, F., Jin, Z., 2014. A direct comparison of active and passive amplification electrodes in the same amplifier system. *J. Neurosci. Methods* 235, 298–307.
- Laszlo, S., Stites, M., Federmeier, K.D., 2012. Won't get fooled again: an event-related potential study of task and repetition effects on the semantic processing of items without semantics. *Lang. Cognit. Process.* 27 (2), 257–274.
- Lewis, G.A., Poeppel, D., Murphy, G.L., 2015. The neural bases of taxonomic and thematic conceptual relations: a meg study. *Neuropsychologia* 68, 176–189.
- Lin, E.L., Murphy, G.L., 2001. Thematic relations in adults' concepts. *J. Exp. Psychol. Gen.* 130 (1), 3–28.
- Luck, S.J., 2014. *An Introduction to the Event-Related Potential Technique*. MIT press.
- Luck, S.J., Gaspelin, N., 2017. How to get statistically significant effects in any erp experiment (and why you shouldn't). *Psychophysiology* 54 (1), 146–157.
- Maguire, M.J., Brier, M.R., Ferree, T.C., 2010. Eeg theta and alpha responses reveal qualitative differences in processing taxonomic versus thematic semantic relationships. *Brain Lang.* 114 (1), 16–25.
- Medler, D., Binder, J., 2005. Mcword: an On-Line Orthographic Database of the English Language. <http://www.neuro.mcw.edu/mcword/>.
- Mirman, D., Graziano, K.M., 2012. Individual differences in the strength of taxonomic versus thematic relations. *J. Exp. Psychol. Gen.* 141 (4), 601–609.
- Mirman, D., Landrigan, J.-F., Britt, A.E., 2017. Taxonomic and thematic semantic systems. *Psychol. Bull.* 143 (5), 499–520.
- Murphy, G.L., 2001. Causes of taxonomic sorting by adults: a test of the thematic-to-taxonomic shift. *Psychon. Bull. Rev.* 8 (4), 834–839.
- Peirce, J.W., 2007. Psychopy—psychophysics software in python. *J. Neurosci. Methods* 162 (1), 8–13.
- R Core Team, 2017. *R: A Language and Environment for Statistical Computing [Computer Software Manual]*. Vienna, Austria. Retrieved from. <https://www.R-project.org/>.
- Rugg, M.D., Nagy, M.E., 1989. Event-related potentials and recognition memory for words. *Electroencephalogr. Clin. Neurophysiol.* 72 (5), 395–406.
- Saalbach, H., Imai, M., 2007. Scope of linguistic influence: does a classifier system alter object concepts? *J. Exp. Psychol. Gen.* 136 (3), 485–501.
- Sachs, O., Weis, S., Krings, T., Huber, W., Kircher, T., 2008. Categorical and thematic knowledge representation in the brain: neural correlates of taxonomic and thematic conceptual relations. *Neuropsychologia* 46 (2), 409–418.
- Schwartz, M.F., Kimberg, D.Y., Walker, G.M., Brecher, A., Faseyitan, O.K., Dell, G.S., Coslett, H.B., 2011. Neuroanatomical dissociation for taxonomic and thematic knowledge in the human brain. *Proc. Natl. Acad. Sci. Unit. States Am.* 108 (20), 8520–8524.
- Shaoul, C., Westbury, C., 2006. *Usenet Orthographic Frequencies for 111. University of Alberta, Edmonton, AB, p. 627 english words (2005–2006)*. <http://www.psych.ualberta.ca/%7Ewestburylab/downloads/wlfreq.download.html>.
- Sharp, D., Cole, M., Lave, C., Ginsburg, H.P., Brown, A.L., French, L.A., 1979. Education and cognitive development: the evidence from experimental research. *Monographs of the society for research in child development*, pp. 1–112.
- Simmons, S., Estes, Z., 2008. Individual differences in the perception of similarity and difference. *Cognition* 108 (3), 781–795.
- Skwarchuk, S., Clark, J.M., 1996. Choosing category or complementary relations: prior tendencies modulate instructional effects. *Canadian J. Experimental Psychology/ Revue canadienne de psychologie expérimentale* 50 (4), 356–370.
- Smiley, S.S., Brown, A.L., 1979. Conceptual preference for thematic or taxonomic relations: a nonmonotonic age trend from preschool to old age. *J. Exp. Child Psychol.* 28 (2), 249–257.
- Stanovich, K.E., West, R.F., 1989. Exposure to print and orthographic processing. *Read. Res. Q.* 402–433.
- Stites, M.C., Laszlo, S., 2015. How do random effects structures impact LMER outcomes in an ERP study? *Psychophysiology* 52, S116–S116.
- Su, Y.-S., Gelman, A., Hill, J., Yajima, M., 2011. Multiple imputation with diagnostics (mi) in R: opening windows into the black box. *J. Stat. Software* 45 (2), 1–31.
- Thigpen, N.N., Kappenman, E.S., Keil, A., 2017. Assessing the internal consistency of the event-related potential: an example analysis. *Psychophysiology* 54 (1), 123–138.
- Wamain, Y., Pluciennicka, E., Kaléline, S., 2015. A saw is first identified as an object used on wood: ERP evidence for temporal differences between thematic and functional similarity relations. *Neuropsychologia* 71, 28–37.
- West, R.F., Stanovich, K.E., 1991. The incidental acquisition of information from reading. *Psychol. Sci.* 2 (5), 325–330.