

# AI-enhanced semantic feature norms for 786 concepts

Siddharth Suresh\*<sup>1,2,3</sup>, Kushin Mukherjee\*<sup>4</sup>, Tyler Giallanza<sup>5</sup>, Xizheng Yu<sup>6</sup>, Mia Patil<sup>2</sup>,  
Jonathan D. Cohen<sup>6,7</sup>, and Timothy T. Rogers<sup>1,3</sup>

<sup>1</sup>Dept. of Psychology, University of Wisconsin-Madison, <sup>2</sup>Dept. of Computer Sciences, University of Wisconsin-Madison,

<sup>3</sup>Wisconsin Institute for Discovery, <sup>4</sup>Dept. of Psychology, Stanford University, <sup>5</sup>Dept. of Psychology, Princeton University,

<sup>6</sup>Dept. of Computer Science, Brown University, <sup>7</sup>Princeton Neuroscience Institute

## Abstract

Semantic feature norms have been foundational in the study of human conceptual knowledge, yet traditional methods face trade-offs between concept/feature coverage and verifiability of quality due to the labor-intensive nature of norming studies. Here, we introduce a novel approach that augments a dataset of human-generated feature norms with responses from large language models (LLMs) while verifying the quality of norms against reliable human judgments. We find that our AI-enhanced feature norm dataset, *NOVA: Norms Optimized Via AI*, shows much higher feature density and overlap among concepts while outperforming a comparable human-only norm dataset and word-embedding models in predicting people’s semantic similarity judgments. Taken together, we demonstrate that human conceptual knowledge is richer than captured in previous norm datasets and show that, with proper validation, LLMs can serve as powerful tools for cognitive science research.

**Keywords:** semantic knowledge; feature listing; large language models; similarity judgments

## Introduction

The study of human conceptual knowledge has relied on semantic feature norms — representations of concepts in terms of their associated features — since their introduction by Rosch in the 1970s (Rosch, 1975). Norming studies present participants with a set of concepts and, for each, asks them to list as many characteristic properties as they can. Aggregating features across items and participants creates semantic vectors the elements of which correspond to the elicited features and the entries of which indicate whether people regularly judge the concept to possess the corresponding property. Proximity between two such feature vectors relates systematically to their perceived semantic relatedness—thus lions and tigers are viewed as similar kinds of things because they have many overlapping and few distinguishing properties. Norming datasets collected over the years (Buchanan et al., 2019; Devereux et al., n.d.; Dilkina et al., 2008; Hansen & Hebart, 2022; McRae et al., 2005; Ruts et al., 2004) have helped to answer questions about the organization of semantic memory (Ashcraft, 1978; Collins & Loftus, 1975), its degradation in semantic disorders (Cree & McRae, 2003; Farah & McClelland, 2013; Garrard et al., 2001; Rogers & McClelland, 2004), its relationship to control (Giallanza et al., 2024), and its neural bases (Clarke & Tyler, 2014; Cox et al., 2024) (see Kumar (2021) for a review).

Semantic norming requires extensive human labor both in data collection and curation/post-processing. Prior studies have met this challenge in different ways, each requiring some degree of compromise as elaborated below. Other recent work has sought alternatives to human feature norms by making use of natural language processing technologies, including word embeddings from methods such as *word2vec* and *GloVe* (Mikolov et al., 2013; Pennington et al., 2014) as well as feature norms generated artificially by large language models (LLMs) (Hansen & Hebart, 2022). However, word embeddings fail to capture the semantic structure perceived by humans as effectively as feature norms, and their dimensions lack the transparent interpretability of feature-based representations, at least for concrete objects (Suresh, Mukherjee, & Rogers, 2023; Suresh, Mukherjee, Yu, et al., 2023). LLMs can generate super-human lists of features that go far beyond what a typical person might know (and hence are non-representative of human knowledge) and frequently confabulate properties that are untrue (the well-documented ‘hallucination problem’ Huang et al., 2024).

The current work seeks a middle way between human-only and machine-only norm generation. We crowd-sourced feature lists for a modestly large and representative set of 786 concrete object concepts thus ensuring that the features included in the set are those that human participants discern. We then used LLMs to aid in the most labor-intensive parts of data curation and post-processing, resulting in a novel *AI-enhanced* set of semantic feature norms – *NOVA: Norms Optimized Via AI*. We illustrate remarkable differences between human-only and AI-enhanced norm sets, then report empirical studies designed to assess whether the AI-enhanced norms capture human-perceived semantic structure better than do human-only norms or “out-of-the-box” word embeddings.

## Study I: Building NOVA

Human feature-norming studies involve up to 4 steps, each requiring significant effort and thus subject to constraints that can limit the resulting data. Here we consider each step, limitations faced by prior studies, and the approach taken in the current work. The overall workflow for our approach is shown in Figure 1.

*Concept selection.* The structure appearing in a given dataset depends on the concepts included. Early norms

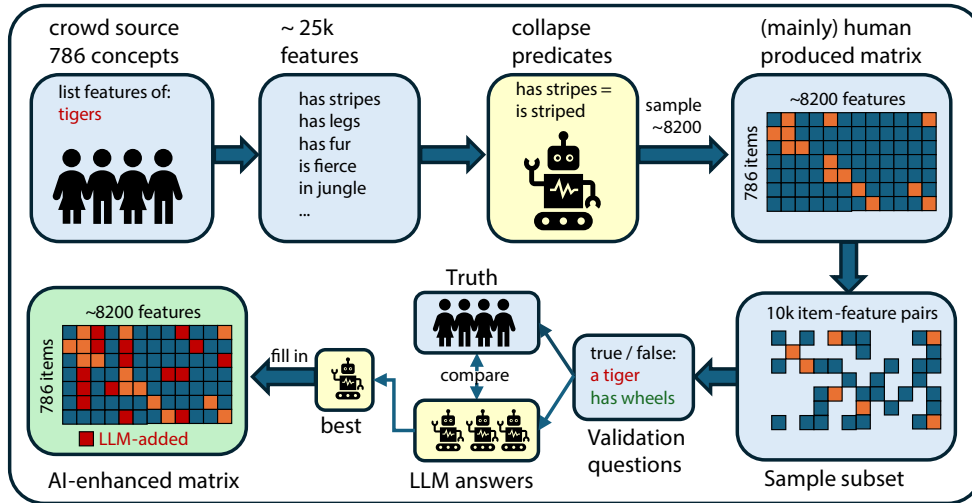


Figure 1: A schematic representation of our workflow. Features were initially crowd-sourced for 786 concepts, forming a human-generated matrix. A subset of 10,000 concept-feature pairs underwent validation via human judgments. LLM responses were compared to these human judgments to determine the best-performing strategy. Using this method, LLMs completed the matrix for all 8,200 selected features, forming the AI-augmented matrix.

used in semantic memory studies focused on hierarchically structured, easily nameable concepts (e.g., animals, plants), often excluding typical examples (e.g., robins, sparrows) in favor of atypical ones easier to name (e.g., penguins, ostriches) and omitting concepts that don't fit neatly into these hierarchies. To improve representativeness we included all 565 concepts from the Ecoset dataset Mehrer et al., 2021, which comprises frequent, unambiguous basic-level concrete object names, along with items from the McRae (McRae et al., 2005) and Leuven (De Deyne et al., 2008) norms. We also added superordinate categories (e.g., animal, vehicle) and higher-frequency subordinate names (e.g., robin, trout) to better capture domain substructure. The final set comprised **786** concrete object concepts.

*Feature elicitation.* We elicited features from participants on Amazon Mechanical Turk using procedures described below.

*Feature reduction.* Norming studies typically yield a large set of unique features, most appearing in a single concept. To manage this complexity, researchers often consolidate distinct yet semantically related properties – e.g., if *is hairy* and *is furry* are used by different participants to describe a ‘coconut’, these features may be deemed as equivalent, creating a single feature that overlaps for concepts possessing both *is hairy* (e.g., ape) and *is furry* (e.g., rabbit). While this process simplifies the feature space and enhances conceptual similarity across items, it is labor-intensive and relies on subjective human judgments. We instead performed a minimal feature collapse by using GPT-3 to extract phrase embeddings of featural descriptions (e.g., *has a furry outer layer*), then clustering these and merging only highly similar clusters. This approach collapsed phrases with variable wording but near-identical semantic content (e.g., *has a*

*furry outer layer*, *is furry*, and *feels furry*) while still distinguishing close synonyms (e.g., *is hairy* vs *is furry*). This step reduced the initial ~25k raw features to ~20k features, from which we randomly sampled ~8,200 features for subsequent analysis.

*Feature verification.* The features that participants generate in the elicitation phase typically constitute a fraction of what they actually know. For this reason, some norming studies conduct a *feature verification* step where human participants consider every concept/feature pair and judge whether the feature is true of the concept (De Deyne et al., 2008; Dilkina et al., 2008). This step greatly enriches the structure encoded in the norms. For instance, most participants list the feature *has a long neck* for giraffes and swans but for few other items. Yet when asked, most participants agree that *has a long neck* is true of items as varied as a duck, a beer bottle, and a cello. Thus, the verification phase surfaces latent knowledge that participants don't generate spontaneously. Since the number of concept/feature pairs grows exponentially, this is by far the most labor-intensive part of the process and prior studies have either employed a relatively modest set of concepts and features (Dilkina et al., 2008) or have limited verification only to specific semantic domains (De Deyne et al., 2008). We leveraged LLMs to conduct the feature-verification phase – first comparing different models and strategies in their ability to capture human judgments on a randomly-sampled set of concept-feature pairs, then using the most successful strategy to verify all ~6.5M concept/feature pairs, producing an AI-enhanced norm set.

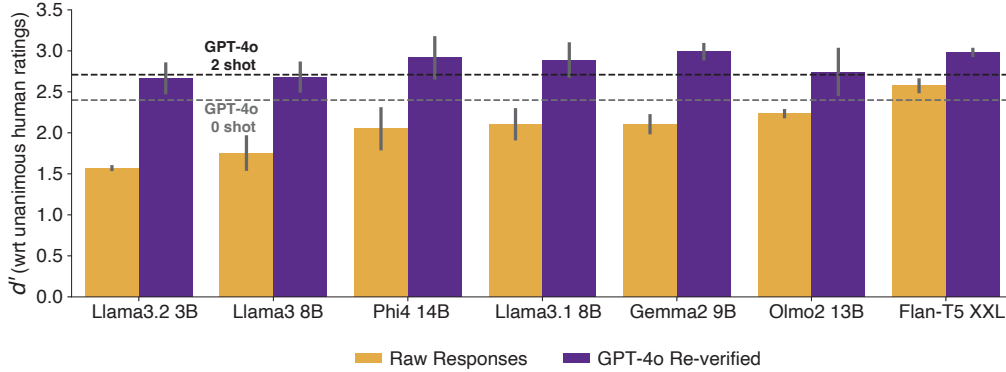


Figure 2: Models’ ability to reliably predict human feature-concept ratings measured as  $d'$  using raw responses (orange) and responses re-verified using GPT-4o. Bar heights show mean  $d'$  across the 0-shot and 2-shot experiments. Gray and black dashed lines correspond to GPT-4o’s performance in the 0-shot and 2-shot settings respectively. Errorbars correspond to bootstrapped 95% confidence intervals.

## Methods

**Human feature elicitation.** This phase provided human-elicited data for all concepts in the set, providing the raw features from which human-only and AI-enhanced norms were derived.

*Participants.* 50 participants were recruited through Amazon Mechanical Turk and were compensated \$4 for the task which would require 20 minutes to complete. The study was approved by the Princeton Internal Review Board, IRB Protocol 6079.

*Stimuli and procedure.* Stimuli were 786 concrete object nouns. Using a web-based interface, each participant viewed up to 75 different words in randomized order, and for each typed in as many different features as they could generate. The instructions emphasized generating various

types of features, including physical/perceptual features (appearance, smell), functional features (uses, contexts), and other characteristics. Participants were asked to format their responses as individual features per line using standardized phrasing (e.g., “has ears” rather than “a dog is an animal that has ears”).

**Human feature verification.** This phase had human participants verify  $\sim 10k$  concept-feature pairs, providing an empirical basis for evaluating the performance of different AI-aided approaches to feature verification.

*Participants.* 556 participants were recruited through Amazon Mechanical Turk and compensated \$1.40 for a 5-8 minute task. Participants were allowed to complete multiple sessions contingent upon maintaining satisfactory performance.

*Stimuli and procedure.* The stimuli were concept-property pairs sampled randomly from results of the feature-elicitation task. Data were collected through an online interface. Each trial paired one concept (e.g. “alligator”) with one feature randomly sampled from the full set. The sampled feature could come from any domain or item—for alligator, it could be something reasonable (e.g. “has legs”), something clearly false (e.g. “has wheels”) or something uncertain (e.g. “has ears”). For each pair participants judged whether the property is true of the item by pressing a keyboard button. The instructions emphasized that subjective properties should be evaluated based on common consensus (e.g., “cute” for “dog”), and properties that were sometimes true should be marked as true (e.g., “brown” for “dog”). Each participant made about 110 judgments, and we collected 5 or more judgments on each of 10,545 unique pairs. Participants could skip unfamiliar concepts or nonsensical properties by pressing the space bar, with skipped items replaced to maintain the required number of judgments.

### AI-enhanced feature verification.

Our ultimate goal was to use LLMs to complete the feature-verification step for all possible concept/property pairs. Since

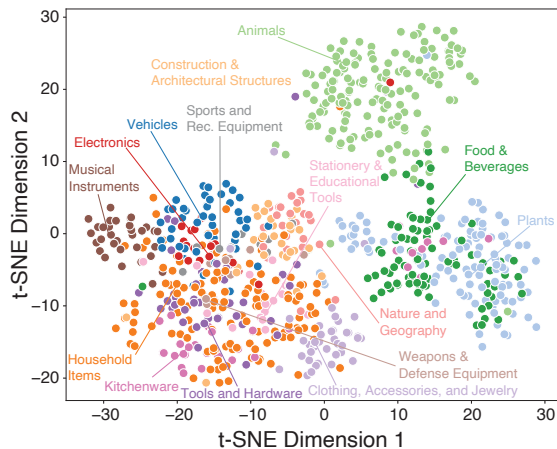


Figure 3:  $t$ -stochastic neighbor embeddings of the semantic vectors for each of 786 concepts derived from the final verified matrix. Category labels were generated by combining higher order labels from existing norm datasets and LLM-suggested categories from GPT-4o.

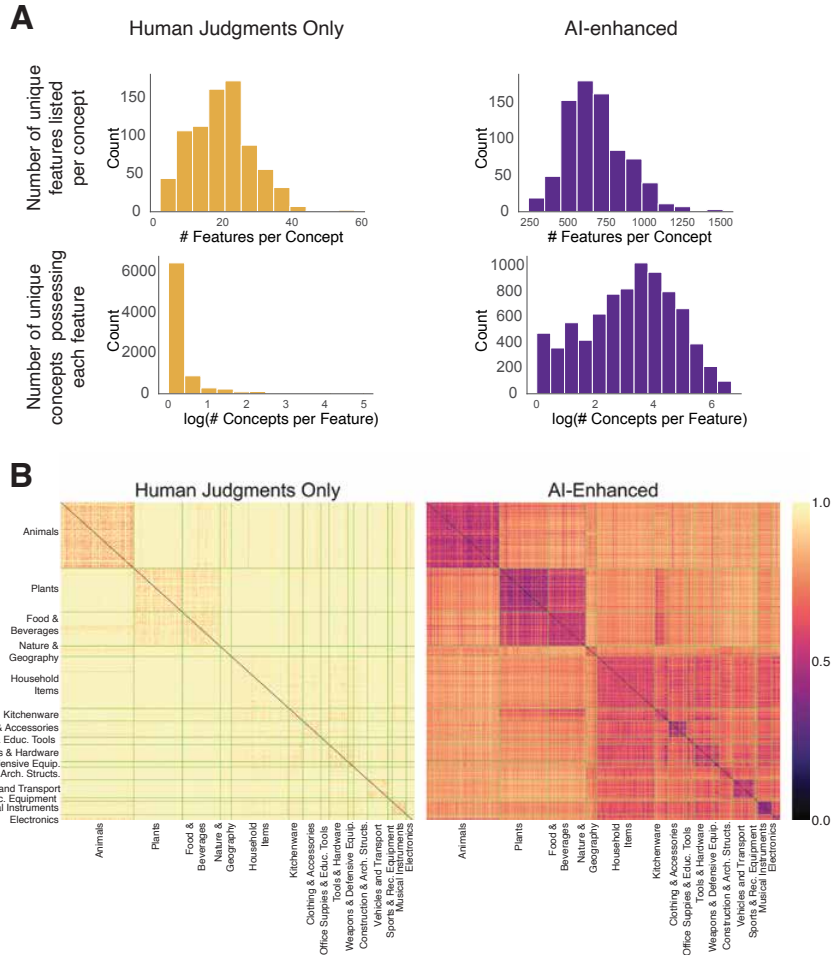


Figure 4: (A) Counts of valid features per concept and number of concepts that share common features for the reduced human-generated matrix (top row) and AI-enhanced norm matrix (bottom row). (B) Pairwise cosine dissimilarity matrices based on the reduced human-generated norm matrix (left) and AI-enhanced norm matrix (right).

there are millions of possible pairs, we first considered how well each of several different models and prompting strategies could capture real human judgments on the items collected in the human feature-verification study. In these data participants showed different opinions for about 40% of the items—thus either opinion expressed by an LLM would agree with at least one human participant for these items. We therefore selected the 6,122 concept-feature pairs for which all participants made the same decision (either all yes or all no), and used these decisions as a ground-truth for evaluating LLM performance.

*Model Suite.* We primarily focused on performant open-sourced language models because these are accessible to other researchers for replication purposes and relatively more affordable to access. We included models that have open weights, are generally high-scoring on standard LLM benchmarks (Hendrycks et al., 2020), and can be run on consumer-grade hardware. Specifically, we evaluated 3 models from Meta’s Llama family (Llama3, Llama3.1, and Llama3.2) (Dubey et al., 2024), Microsoft’s Phi-4 (Abdin

et al., 2024), Ai2’s Olmo2 (OLMo et al., 2024), and Google’s Gemma2 (Team et al., 2024) and Flan-T5 (Wei et al., 2021). We evaluated all models at full bfloat16 precision on a Nvidia H100 GPU. For comparison to a state-of-the-art closed model, we also evaluated GPT-4o via its API.

*Evaluation Protocol.* We prompted all models using the following general prompt -

In one word True or False, answer the following question question: Is the property [x] true for [y]? Answer:

...where  $x$  was a feature and  $y$  was a concept with the square brackets included in the prompt. We ran two prompting experiments: (1) a zero-shot experiment providing the models with just the question above as input, and (2) a two-shot experiment providing the models with two example feature-concept pairs, one true and one false, to potentially improve the models’ ability to perform the task via in-context learning (Brown et al., 2020). We used the same two examples for all prompts.

*Post-processing.* To extract meaningful answers from model-generated text we first restricted responses to a maximum of five tokens, then conducted a case-insensitive search of model responses for the strings ‘True’ or ‘Yes’ to indicate a positive response, and ‘False’ or ‘No’ to indicate a negative response. In rare cases where no match was found we set the model response to ‘False’.

**Results.** To measure how closely LLM responses aligned with unanimous human judgments for the 6,122 feature-concept pairs, we adopted a signal detection approach, treating human responses as the true signal and model responses as guesses. Where humans agreed the property was true of the concept, model guesses were scored as hits if they concurred and misses otherwise. Where humans agreed the property was not true of the concept, model guesses were scored as correct rejections if they concurred and false alarms otherwise. From these counts we computed hit rates and false alarm rates, then converted these to the  $d'$  measure of signal discrimination.

The average  $d'$  for both zero and two shot conditions can be seen in Figure 2 (yellow bars). Two-shot GPT-4o outperformed all open-sourced models, which varied in their match to human responses. Two-shot Flan-T5 XXL performed best amongst open-sourced models and better than the zero-shot GPT-4o. Flan-T5’s lower  $d'$  relative to GPT-4o was driven by a propensity to respond with ‘true’ to many queries, buoying its hit rate but also increasing its false-positive rate. To preserve the benefits of GPT-4o without incurring a prohibitive cost, we next considered a ‘re-verification’ approach in which the ‘true’ responses generated by a given open-source model were subsequently re-verified by GPT-4o, retaining the ‘true’ value only if both models agreed. The results are shown as purple bars in Figure 2. Re-verification improved performance for all models, surpassing GPT-4o alone. Flan-T5 XXL remained a top model, closely matched by Gemma2 9B. Given the strong baseline performance of Flan T5, we chose this model with GPT-4o re-verification to fill out the full semantic feature matrix.

**Using Flan T5 and GPT-4o to impute the AI-enhanced matrix.** In the human-only matrix, entry  $[i, j]$  has a value of 1 wherever a participant produced feature  $j$  for concept  $i$  and a 0 in all other entries. For every 0 in this matrix, we prompted Flan T5 XXL to decide whether the corresponding property is / is not true of the corresponding concept. Where the model decided ‘not true,’ the zero value was retained in the matrix. Where the model decided ‘true,’ (534,010 out of 6,436,554 possible pairs) we prompted GPT-4o with the same pair to re-verify the answer. If GPT-4o agreed the property was true, the cell value was replaced with 1, otherwise the 0 value was retained. This procedure yielded the final **AI-enhanced** norms matrix. Figure 3 shows t-SNE based embeddings of all concepts from this matrix.

The AI-enhanced matrix differed remarkably from the human-only matrix in its feature density. While the human

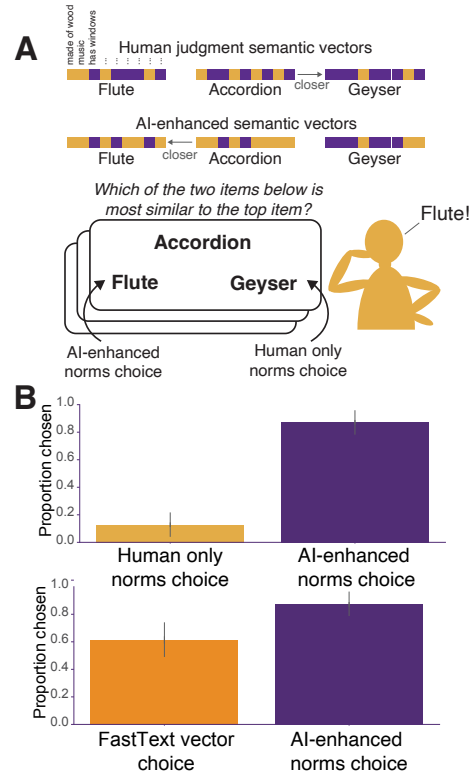


Figure 5: (A) Procedure for generating trials for the triadic judgment experiment and an example trial. (B) Proportion of human responses that aligned with the human matrix (yellow bar) vs. the AI-enhanced matrix (purple bar) and with FastText word embeddings (orange bar) vs. AI-enhanced semantic vectors (purple bar) in Experiment 2. Error bars represent standard errors of the means.

matrix has about 20 features per concept on average, the AI-enhanced matrix has about 700 (Figure 4A), and while the majority (78%) of features in the human-only matrix are true of just one concept, this is true of just 5% of features in the AI-enhanced matrix. The increased feature density produces much more richly-structured similarity relations, as shown by the heat plot of pairwise distances between concepts in Figure 4B. While some of this difference may be attributable to false-positives in the AI-enhanced dataset, the comparison to human judgments suggests that the LLM verification strategy is quite good at discriminating true positives from true negatives ( $d' > 3.0$ ). Thus the result suggests that human knowledge about features of concepts may be considerably richer than prior norming studies have suggested.

## Study 2: Using the new norms dataset to predict human semantic judgments

To assess whether the AI-enhanced norms in NOVA capture information about semantic structure beyond human-only norms or other approaches, we compared different approaches in their ability to predict human behavior in a triadic similarity judgment task (Hebart et al., 2023; Jamieson

et al., 2015; Sievert et al., 2023). In this task, participants must decide which of two option concepts is semantically more similar to a target concept. A candidate semantic embedding can “predict” human decisions by selecting whichever option word lies closer to the target word in the embedding space. We can assess the quality of the embedding by comparing how often the predicted response agrees with actual human decisions. In this study, we compared the predictions of NOVA embeddings to those based on the human-only feature norms and to those generated by a common word-embedding approach (FastText).

We selected triplets designed to maximally discriminate NOVA and human-only feature norms. Thus for each trial, one of the option items was closer to the target in the human-only space while the other was closer in the AI-enhanced space (see Figure 5). We then computed how often the majority-vote across human participants agreed with the predictions of each embedding (AI-enhanced, human-only, FastText). If the AI-enhanced norms in NOVA contain information irrelevant to human-perceived semantics, their predictions should agree with human judgments less often than do those of the human-only norms. Furthermore, if either set of norms simply recapitulates the semantic structure evident in word embeddings, then predictions from the norms should do about as well as predictions from the FastText embeddings.

**Generating maximally disagreeing triplets.** To generate triplets that maximally differentiated the human-only and AI-enhanced norms, we computed cosine dissimilarity matrices for each set (Figure 4B), Procrustes-aligned them to minimize disparity, and identified concepts with the largest discrepancies in their distances to other concepts. For example, in NOVA space, ‘accordion’ was closer to ‘flute’ than to ‘geyser’, while the reverse was true in the human-only space (Figure 5A). We constructed 1,424 triplets where the two matrices produced divergent predictions, with each of the 786 concepts serving as the target approximately twice. The critical question was which matrix’s predictions would align more closely with human similarity judgments.

*Participants* 31 participants were recruited from the UW-Madison psychology subject pool. Participants completed the task online for course credit. Each participant provided informed consent in compliance with the UW-Madison IRB.

*Stimuli and Procedure.* The stimuli were the set of 1,424 triplets described above. Data were collected online via jsPsych (De Leeuw, 2015). On each trial, a randomly selected triplet was displayed, with participants indicating which of two options was more similar to the target concept using a mouse click. All triplets were judged by each participant<sup>1</sup>.

**Results.** Human similarity judgments agreed with predictions of the AI-enhanced norms for 86.20% of triplets, a result unlikely to arise by chance ( $p < 0.001$ , binomial test). Human judgments agreed with predictions of the FastText

embeddings on 60.40% of trials: reliably better than chance ( $p < 0.001$ , binomial test), but significantly worse than the AI-enhanced embeddings (paired  $t$ -test,  $t(1,423) = 18.37$ ,  $p < 0.001$ ). Thus the richer structure evident in the AI-enhanced feature norms appears to better express human-discerned semantic similarity structure than to norms derived from humans alone or from word-embeddings.

## Discussion

We presented a new approach for generating AI-enhanced semantic norms along with an accompanying dataset, NOVA. We first conducted controlled experiments evaluating LLM feature verification performance against a reliable subset of human norm judgments, using the results to find an optimal model, prompting strategy, and verification strategy. We then applied the best-performing approach to generate an AI-enhanced large-scale norm dataset spanning over 750 concepts and over 8,000 features. Concepts in the resulting NOVA dataset showed much higher feature density and a greater degree of feature overlap relative to the raw human-generated matrix. This overlap of features did not come at the cost of category selectivity, with concepts being reasonably organized into meaningful clusters. Finally, we used a triadic comparison task to show that NOVA vectors more accurately predicted human similarity judgments than did vectors based on human norms alone or word embeddings computed from natural language. The result suggests that AI-enhanced norms express semantic structure more similar to that discerned by human participants.

Taken together, our work addresses longstanding limitations in semantic norm generation by creating a dataset that includes a representative set of concepts and features, with AI-based feature verification validated against human judgments. The feature density of the AI-enhanced norms reveals semantic similarity structure richer than previous norm datasets, unlocking the potential to better understand both the cognitive and the neural bases of semantic memory (Clarke & Tyler, 2014; Cox et al., 2024; Fernandino et al., 2022; Rogers & McClelland, 2004) and to guide the development of future computational neurocognitive models (Dilkina et al., 2008; Giallanza et al., 2024; Riordan & Jones, 2011; Saxe et al., 2019; Suresh et al., 2024). Lastly, the present work highlights the promise of integrating large-language models into workflows for cognitive science research in a controlled and verifiable manner and provides a replicable framework for future endeavors in this domain (Dillion et al., 2023; Mukherjee et al., 2023, 2024; Suresh, Mukherjee, & Rogers, 2023; Trott, 2024).

<sup>1</sup>there was data loss of a few trials for some participants due to technical issues.

## Acknowledgments

We thank members of the Knowledge and Concepts Lab at UW-Madison and the Neuroscience of Cognitive Control Lab at Princeton for helpful discussion and feedback. This work was supported by a Vannevar Bush Faculty Fellowship (VBFF) administered through ONR to JDC and Multi University Research Initiative (MURI) award W911NF2110317 to TTR (co-PI).

All code and materials will be available at:  
<https://github.com/Knowledge-and-Concepts-Lab/llm-norms-cogsci2025>

## References

- Abdin, M., Aneja, J., Behl, H., Bubeck, S., Eldan, R., Gunasekar, S., Harrison, M., Hewett, R. J., Javaheripi, M., Kauffmann, P., et al. (2024). Phi-4 technical report. *arXiv preprint arXiv:2412.08905*.
- Ashcraft, M. H. (1978). Property norms for typical and atypical items from 17 categories: A description and discussion. *Memory & Cognition*, 6, 227–232.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *Advances in neural information processing systems*, 33, 1877–1901.
- Buchanan, E. M., Valentine, K. D., & Maxwell, N. P. (2019). English semantic feature production norms: An extended database of 4436 concepts. *Behavior Research Methods*, 51, 1849–1863.
- Clarke, A., & Tyler, L. K. (2014). Object-specific semantic coding in human perirhinal cortex. *Journal of Neuroscience*, 34(14), 4766–4775.
- Collins, A. M., & Loftus, E. F. (1975). A spreading-activation theory of semantic processing. *Psychological review*, 82(6), 407.
- Cox, C. R., Rogers, T. T., Shimotake, A., Kikuchi, T., Kunieda, T., Miyamoto, S., Takahashi, R., Matsumoto, R., Ikeda, A., & Lambon Ralph, M. A. (2024). Representational similarity learning reveals a graded multidimensional semantic space in the human anterior temporal cortex. *Imaging Neuroscience*, 2, 1–22.
- Cree, G. S., & McRae, K. (2003). Analyzing the factors underlying the structure and computation of the meaning of chipmunk, cherry, chisel, cheese, and cello (and many other such concrete nouns). *Journal of experimental psychology: general*, 132(2), 163.
- De Deyne, S., Verheyen, S., Ameel, E., Vanpaemel, W., Dry, M. J., Voorspoels, W., & Storms, G. (2008). Exemplar by feature applicability matrices and other dutch normative data for semantic concepts. *Behavior research methods*, 40, 1030–1048.
- De Leeuw, J. R. (2015). Jspsych: A javascript library for creating behavioral experiments in a web browser. *Behavior research methods*, 47, 1–12.
- Devereux, B. J., Tyler, L. K., Geertzen, J., & Randall, B. (n.d.). The centre for speech, language and the brain (cslb) concept property norms. *Behavior research methods*, 46, 1119–1127.
- Dilkina, K., McClelland, J. L., & Plaut, D. C. (2008). A single-system account of semantic and lexical deficits in five semantic dementia patients. *Cognitive Neuropsychology*, 25(2), 136–164.
- Dillion, D., Tandon, N., Gu, Y., & Gray, K. (2023). Can ai language models replace human participants? *Trends in Cognitive Sciences*.
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. (2024). The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Farah, M. J., & McClelland, J. L. (2013). A computational model of semantic memory impairment: Modality specificity and emergent category specificity (journal of experimental psychology: General, 120 (4), 339–357). *Exploring Cognition: Damaged Brains and Neural Networks*, 79–110.
- Fernandino, L., Tong, J.-Q., Conant, L. L., Humphries, C. J., & Binder, J. R. (2022). Decoding the information structure underlying the neural representation of concepts. *Proceedings of the National Academy of Sciences*, 119(6), e2108091119.
- Garrard, P., Ralph, M. A. L., Watson, P. C., Powis, J., Patterson, K., & Hodges, J. R. (2001). Longitudinal profiles of semantic impairment for living and nonliving concepts in dementia of alzheimer’s type. *Journal of Cognitive Neuroscience*, 13(7), 892–909.
- Giallanza, T., Campbell, D., Cohen, J. D., & Rogers, T. T. (2024). An integrated model of semantics and control. *Psychological Review*.
- Hansen, H., & Hebart, M. N. (2022). Semantic features of object concepts generated with gpt-3. *arXiv preprint arXiv:2202.03753*.
- Hebart, M. N., Contier, O., Teichmann, L., Rockter, A. H., Zheng, C. Y., Kidder, A., Corriveau, A., Vaziri-Pashkam, M., & Baker, C. I. (2023). Things-data, a multimodal collection of large-scale datasets for investigating object representations in human brain and behavior. *Elife*, 12, e82580.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., & Steinhardt, J. (2020). Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang, H., Chen, Q., Peng, W., Feng, X., Qin, B., et al. (2024). A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*.
- Jamieson, K. G., Jain, L., Fernandez, C., Glattard, N. J., & Nowak, R. D. (2015). Next: A system for real-world

- development, evaluation, and application of active learning. *NIPS*, 2656–2664.
- Kumar, A. A. (2021). Semantic memory: A review of methods, models, and current challenges. *Psychonomic Bulletin & Review*, 28(1), 40–80.
- McRae, K., Cree, G. S., Seidenberg, M. S., & McNorgan, C. (2005). Semantic feature production norms for a large set of living and nonliving things. *Behavior research methods*, 37(4), 547.
- Mehrer, J., Spoerer, C. J., Jones, E. C., Kriegeskorte, N., & Kietzmann, T. C. (2021). An ecologically motivated image dataset for deep learning yields better models of human vision. *Proceedings of the National Academy of Sciences*, 118(8), e2011417118.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 3111–3119.
- Mukherjee, K., Rogers, T. T., & Schloss, K. B. (2024). Large language models estimate fine-grained human color-concept associations. *arXiv preprint arXiv:2406.17781*.
- Mukherjee, K., Suresh, S., & Rogers, T. T. (2023). Human-machine cooperation for semantic feature listing. *arXiv preprint arXiv:2304.05012*.
- OLMo, T., Walsh, P., Soldaini, L., Groeneveld, D., Lo, K., Arora, S., Bhagia, A., Gu, Y., Huang, S., Jordan, M., et al. (2024). 2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*.
- Pennington, J., Socher, R., & Manning, C. D. (2014). Glove: Global vectors for word representation. *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 1532–1543.
- Riordan, B., & Jones, M. N. (2011). Redundancy in perceptual and linguistic experience: Comparing feature-based and distributional models of semantic representation. *Topics in Cognitive Science*, 3(2), 303–345.
- Rogers, T. T., & McClelland, J. L. (2004). *Semantic cognition: A parallel distributed processing approach*. MIT press.
- Rosch, E. (1975). Cognitive representations of semantic categories. *Journal of experimental psychology: General*, 104(3), 192.
- Ruts, W., De Deyne, S., Ameel, E., Vanpaemel, W., Verbeemen, T., & Storms, G. (2004). Dutch norm data for 13 semantic categories and 338 exemplars. *Behavior Research Methods, Instruments, & Computers*, 36(3), 506–515.
- Saxe, A. M., McClelland, J. L., & Ganguli, S. (2019). A mathematical theory of semantic development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23), 11537–11546.
- Sievert, S., Nowak, R., & Rogers, T. T. (2023). Efficiently learning relative similarity embeddings with crowdsourcing. *Journal of open source software*, 8(84).
- Suresh, S., Huang, W.-C., Mukherjee, K., & Rogers, T. T. (2024). Categories vs semantic features: What shape the similarities people discern in photographs of objects? *ICLR 2024 Workshop on Representational Alignment*.
- Suresh, S., Mukherjee, K., & Rogers, T. T. (2023). Semantic feature verification in flan-t5. *arXiv preprint arXiv:2304.05591*.
- Suresh, S., Mukherjee, K., Yu, X., Huang, W.-C., Padua, L., & Rogers, T. (2023). Conceptual structure coheres in human cognition but not in large language models. *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 722–738.
- Team, G., Riviere, M., Pathak, S., Sessa, P. G., Hardin, C., Bhupatiraju, S., Hussenot, L., Mesnard, T., Shahriari, B., Ramé, A., et al. (2024). Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.
- Trott, S. (2024). Can large language models help augment english psycholinguistic datasets? *Behavior Research Methods*, 1–19.
- Wei, J., Bosma, M., Zhao, V. Y., Guu, K., Yu, A. W., Lester, B., Du, N., Dai, A. M., & Le, Q. V. (2021). Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.