

# Characterizing Human and Zero-Shot GPT-3.5 Object-Similarity Judgments

Anonymous NAACL submission

## Abstract

Recent advancements in large language models' (LLMs) capabilities have yielded few-shot, human-comparable performance on a range of tasks. At the same time, researchers expend significant effort and resources gathering human annotations. At some point, LLMs may be able to perform some simple annotation tasks, but studies of LLM annotation accuracy and consistency are sparse. In this paper, we characterize OpenAI's ChatGPT's judgment on a behavioral task for implicit object categorization. We characterize the embedding spaces of models trained on human vs. GPT responses and note similarities, but also systematic differences between them. We also find that augmenting a dataset of humans' responses with ChatGPT predictions causes models to diverge well before performance saturation.

## 1 Introduction

Large language models (LLMs) are capable of accomplishing a variety of language-oriented tasks in zero- or few-shot settings (Brown et al., 2020). Examples include common natural-language understanding and processing (NLU/P) tasks such as sentiment analysis and classification (Brown et al., 2020), language translation (Hendy et al., 2023), and named entity recognition (Ji, 2023); but also applied domains such as text tagging (Gilardi et al., 2023), multimodal tagging (Li et al., 2023), and text sample augmentation (Dai et al., 2023).

Current LLM performance indicates we may be able to use pre-trained high-resource LLMs to augment human annotations for tasks where data is sparse or compute resources are low (Møller et al., 2023). However, we do not currently know for which domains it is appropriate to augment human data with LLM-generated responses. This uncertainty stems from a poor understanding of how LLM and human annotation responses systematically differ. Thus, characterizing the ways in which

world knowledge manifests itself in the generations of LLMs is crucial for incorporating LLMs into annotation workflows.<sup>1</sup>

The domain of object-similarity judgment is a useful base-case for exploring the similarities and substitutability of LLM for human responses. On a human level, object-similarity judgment informs how we interact with objects (Desmarais et al., 2007), organize our world (Smith, 1981) and acquire new concepts from a young age (Markman and Hutchinson, 1984). Meanwhile, many corpus-based computational models, including deep transformer models that leverage corpora such as ChatGPT, leverage lexical co-occurrence relations to derive semantic meaning (i.e. the distributional hypothesis). Despite differences in process, these models' representations display correspondences with human judgment (Torabi Asr et al., 2018; Chandrasekaran and Mago, 2022).

In this paper we collect ChatGPT responses to an object similarity task introduced by Hebart et al. (2020). We reformat their image-based paradigm as a text completion task for ChatGPT.<sup>2</sup> Like Hebart et al. (2020), we also train a sparse embedding model that can predict object-similarity judgments. We annotate the dimensions of the embedding model to provide an interpretable characterization of the reasoning behind such judgements. We train a variety of models on different mixes of human and ChatGPT-derived responses and examine the effects of ChatGPT completions on the learned embedding spaces.

## 2 Methodology

**The Odd-One-Out (OOO) Task** To obtain object-similarity responses from ChatGPT, we use

<sup>1</sup>There is evidence that LLMs may already be incorporated into annotation workflows without researcher knowledge, as crowd workers are already using LLMs to speed up their own annotation tasks (Veselovsky et al., 2023).

<sup>2</sup>At the time of our experimentation, a multimodal ChatGPT was not widely available.

the *odd-one-out* (OOO) task, wherein participants indicate the least similar amongst three objects. For example, we might ask, “Which of these concepts is the odd one out: apple, banana, car?” and expect factors such as edibility to affect the response. The OOO task is well-established in the field of psychology for eliciting concept-relational preferences (Mirman et al., 2017; Valenti and Firestone, 2019)

**Human OOO Responses** Hebart et al. (2020) used an image-based OOO task to collect millions of object-similarity judgements. They did this in two rounds, first collecting 1.46M responses (Hebart et al., 2020), then creating a larger, 5M response dataset Hebart et al. (2022).<sup>3</sup> We used these two datasets to create two disjoint OOO response sets of equal size (1.46M). We refer to the first of these datasets as the *full human dataset* and the second as the *baseline dataset*.

**ChatGPT OOO responses** We now create a parallel GPT-only dataset with answers to the OOO questions from the full human dataset. We reformat the original prompt from (Hebart et al., 2020) to create a text completion task suitable for ChatGPT. We refer to these ChatGPT prompts and answers as the *full GPT dataset*.

For cost and task-efficacy reasons, we use OpenAI’s ChatGPT (GPT-3.5-Turbo-0613). Preliminary analysis revealed that smaller models (Falcon-7B, Alpaca-7B, Vicuna-7B) had difficulty answering odd-one-out questions in a coherent manner with simple prompting. Larger models, (e.g. Falcon-40B), produced coherent responses, but not at the scale afforded by ChatGPT’s API.

Transformer models such as ChatGPT incorporate word position for next-word prediction, and ChatGPT demonstrates a strong positional preference (see Appendix C). While humans situationally exhibit ordered preferences, we found a roughly uniform distribution for this task (see Appendix C). Thus, to collect position-neutral responses, we permuted the order of the three objects in the prompts to create six total questions (3!). We then use relative majority voting across the six questions to compute ChatGPT’s odd-one-out choice, breaking ties randomly. See [Supplemental Data: GPT Response Dataset](#) for API calls and a formatted table of all responses.

<sup>3</sup>These datasets were collected before ChatGPT existed, and thus are free of ChatGPT-derived responses.

	metallic	food-related	...	cylindrical
aardvark	$a_{1,1}$	$a_{1,2}$	...	$a_{1,49}$
abacus	$a_{2,1}$	$a_{2,2}$	...	$a_{2,49}$
⋮	⋮	⋮	⋮	⋮
zucchini	$a_{1854,1}$	$a_{1854,2}$	...	$a_{1854,49}$
Learned object-similarity embeddings				

Figure 1: An example embedding space with words as rows and *characterizing dimensions* as columns.

**Human-GPT Datasets** We wish to study the effect of replacing only some human responses with ChatGPT responses. Thus, we create <1.46M-count *partial human response sets* by taking proportions [0.125, 0.25, 0.375, 0.5, 0.625, 0.75, and 0.875] of the 1.46M full human-only response set. We then create a 1.46M-count *mixed GPT-human response set* for each partial human set by considering each unused human response and including the corresponding GPT response.

## 2.1 Model Details

The embedding model creates a vector representation of each object  $v_{\text{obj}}$ . The similarity between objects  $i$  and  $j$  is given by  $v_i \cdot v_j$ .

When considering an OOO question with objects  $i, j$  and  $k$ , we estimate object  $k$ ’s probability of being selected as the odd one out using the similarity between objects  $i$  and  $j$ :  $z_k = v_i \cdot v_j$ .  
 $P(k \text{ is the OOO}) = \sigma(\mathbf{z}_k) = e^{z_k} / (e^{z_i} + e^{z_j} + e^{z_k})$

**Model Training** To train each model, we use a cross-entropy loss with an  $\ell^1$ -norm penalty on the embedding to encourage sparsity. Hebart et al. (2020) found that training sparse models in this manner resulted in an embedding space with interpretable dimensions. We refer these dimensions as *characterizing dimensions*. An example embedding space matrix is shown in Figure 1.

Using a set of odd-one-out responses  $\mathcal{S}$ , we take the average cross-entropy loss,  $\frac{1}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} H(q, p)|_s$ . Here,  $H(q, p)|_s$  is the cross-entropy of the model prediction probability  $p$  for the odd-one-out question  $s$  relative to the entry  $q$  in the actual one-hot response vector. We incorporate an  $\ell^1$ -norm penalty on the embedding space to encourage sparsity, weighted by a hyperparameter  $\lambda$ . Elaborated loss details are given in Appendix D.

For training, we assume concavity of validation accuracy on the choice of  $\lambda$  and perform a two-tiered four-fold grid-search over 90–10 train–test

dataset splits: we start with  $\lambda = 0.0078$  and take steps of 0.0016 to find a coarse maximum, then take steps of 0.004 around that coarse maximum to establish a finer maximum. We train for a fixed 1000 epochs for each model, mirroring the setup of Zheng et al. (2019) to ensure convergence. Further specifics are given in Appendix E.

We train ten models each on the full human, full GPT, and baseline human sets; and four each on the partial human and mixed human–GPT datasets to produce *full human*, *full GPT*, *partial human*, *mixed human–GPT*, and *baseline models*.

## 2.2 All-GPT Model Characterization

To better understand the basis for ChatGPT responses to OOO questions, we manually annotated each dimension of the full GPT embedding space as in Hebart et al. (2020). Annotators were presented with images of objects at pre-determined intervals along a dimension’s range (e.g., Appendix F). Six respondents gave up to three descriptors for each dimension. We iteratively generated aggregate labels for each annotation until none were ungrouped, then chose the aggregate labels that covered the most participants. We call this the *labelled GPT model*, and we compare it to a previous *labelled human model* produced with the full human dataset from Hebart et al. (2020).

The labels for the nine dimensions with the highest means are given in Figure 2, while those for the 39 dimensions with max value above 0.1 are given in Appendix G; see Supplemental Data: Survey Responses for raw responses and coding.

**Labelled Correlations** We compute the correlations of each of these ChatGPT-derived dimensions with dimensions from the labelled human model. The correlations of the top 9 dimensions (by column mean) from each labelled model are shown in Figure 2; the full 39-by-49 correlation matrix, as well as correlation matrices ordered by maximal correlation matching, appear in Appendix H.

We also perform PCA and UMAP (McInnes et al., 2018) on the labelled dimensions, which are displayed in Appendix J.

## 2.3 GPT–Human Response Substitutability

To determine the impact of augmenting human responses with GPT responses, we compare embedding spaces trained on datasets with varying amounts of each. For this comparison, we use representational similarity analysis (RSA) (Kriegesko-

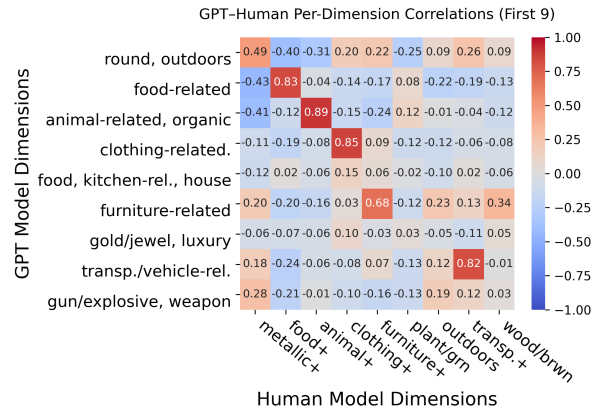


Figure 2: Correlation heatmap between the dimensions of the labelled GPT model and the labelled human model (all-dimension version located in Appendix H).

rte, 2008) with a linear kernel.

Given two embeddings  $\mathbf{X}_1$  and  $\mathbf{X}_2$ , we obtain their respective Gram matrices  $\text{sim}(\mathbf{X}) = \mathbf{X}^\top \mathbf{X}$ . These are the representational similarity matrices, or RSMs, of each space. Then, we calculate the Pearson correlation between the upper triangle of each RSM. The result is the *RSA correlation*, and we report an *RSA score*, the average RSA correlation of a model with the baseline human models.

**GPT Response Substitution** Given a full human dataset, if we replace some of the human responses with GPT responses, how does that affect the RSA score? Here we are comparing the purple pluses with the large red circle in Figure 4. To examine the effects of mixing GPT completion-driven responses into a human dataset, we compute the RSA scores of the mixed human–GPT embeddings. These results are given in Figure 4. A table of these values can be found in Appendix K. Even though the corresponding dataset size was larger, the mixed GPT–human embeddings each had lower RSA scores than the corresponding partial human embeddings. The scores trend downward in a sigmoid fashion as the proportion of human data decreases, with the most noticeable effects happening after .25 of the human data has been replaced.

**GPT Response Augmentation** Now let’s compare models trained on the same amount of human data, but differing amounts of GPT augmentation. In contrast to the previous paragraph, in this situation we are comparing models trained on datasets of differing size. Comparing these models tells us whether adding GPT data hindered, facilitated, or neutrally impacted the final model’s ability to represent human similarity judgment. To make this comparison, consider the small red circles and

Figure 3: Maximal correlations of the labelled human characterizing dimensions with any dimension of a full GPT model (over 8 such full GPT models). For how this differs from using full human embeddings, see Appendix I.

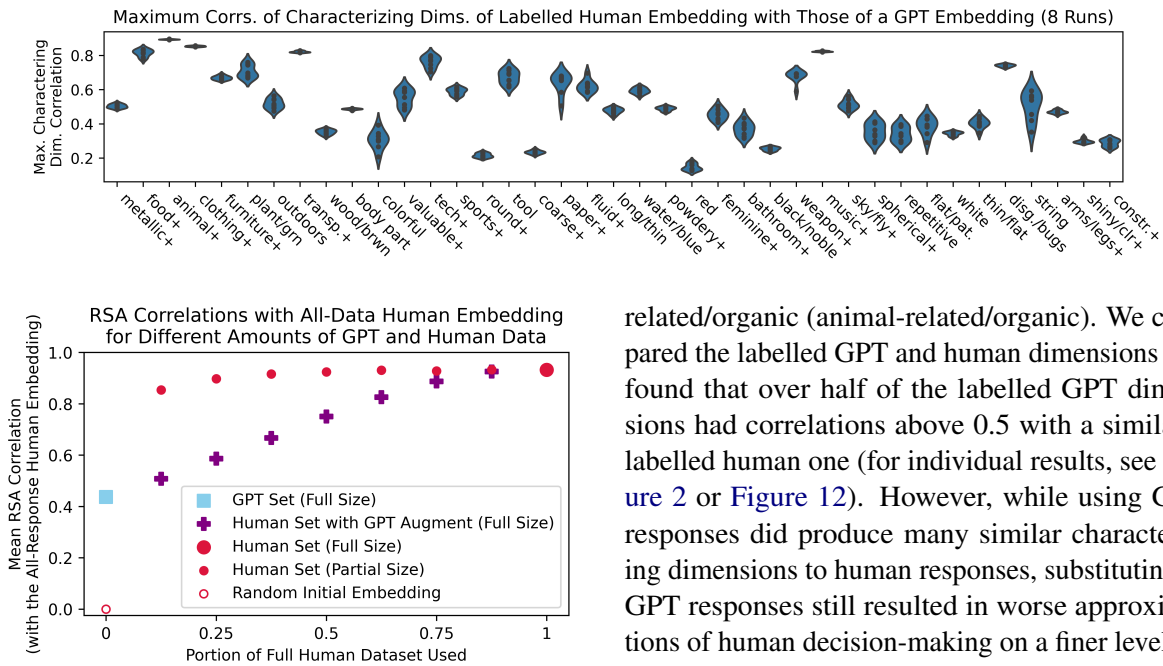


Figure 4: Average RSA scores for full GPT (blue), mixed GPT-human (purple), and full human (large red) models. Also plotted are the scores for the smaller, partial human (small red) models. The x-axis is the proportion of the original human dataset in each model’s training set. The RSA score for a no-data, random embedding (small hollow red) is given for comparison.

the corresponding purple pluses in Figure 4. We find that for all tested ratios, augmenting with GPT data results in lower RSA scores, even though the dataset size has increased.

**Individual Dimension Capturing** Finally, we explore the correspondence of dimensions from the labelled all-human model to those of the full GPT embeddings. To do this, for each labelled human dimension, we find the maximally correlated dimension in each GPT model, then plot those correlation values in Figure 3. Additional information is given in Appendix I.

### 3 Conclusions and Future Work

Our work illustrates GPT’s judgment in an odd-one-out similarity task, provides 39 judgment-characterizing dimensions with human annotations, and compares those dimensions with those derived from a human-only model. Notably, many GPT (and human) dimensions have similar, shared-word-or-synonym labelling, such as food-related (food-related/eating-related/kitchen-related) and animal-

related/organic (animal-related/organic). We compared the labelled GPT and human dimensions and found that over half of the labelled GPT dimensions had correlations above 0.5 with a similarly labelled human one (for individual results, see Figure 2 or Figure 12). However, while using GPT responses did produce many similar characterizing dimensions to human responses, substituting in GPT responses still resulted in worse approximations of human decision-making on a finer level, as demonstrated by Figure 4. Some of this is likely attributable to modality differences between the image and text questions, as some of the dimensions least captured by the model are color-oriented, such as “wood/brownish”, “red”, or “colorful”, as shown in Figure 3 and Appendix I.

More surprisingly, even when we have relatively little human data, adding GPT responses did not improve the trained model’s RSA score. This is shown by gaps between the RSA scores of the partial human models and the GPT-augmented mixed models in Figure 4. This is ostensibly in contrast with previous studies that show LLMs having human-comparable performance on a wide variety of tasks, but it is worth noting that human-level performance is different than human behavior.

In conclusion, our work provides characterizations of GPT object-similarity judgment and indicates utility in using LLM completions for no-resource environments as a high-level proxy for human judgment. It also, however, indicates disutility in mixing or even augmenting human responses should crowdsourced collection of human responses be a possibility, and that caution should be warranted about otherwise human-looking GPT responses infiltrating a dataset.

Our choice of LLM for our experiment was constrained by the sizes of (effective) current models, computing resources, and modality. As image-capable and more powerful models appear, future work should replicate this experiment using them.

269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309



## 310 Limitations

311 Our work uses text-only prompts, while the hu- 359  
312 man experiment uses images. The objects of the 360  
313 THINGS dataset were chosen to be highly image- 361  
314 able, but this nonetheless almost certainly played 362  
315 a role in shaping what ChatGPT found salient in 363  
316 the object-comparison task. At time of writing, 364  
317 GPT-4’s vision API had not seen full release. 365

318 Our prompts presented ChatGPT with objects 366  
319 in an ordered fashion that it heavily utilized (see 367  
320 [Appendix C](#)). To remedy this, we used aggregate 368  
321 responses on permuted prompts. However, humans 369  
322 may have used the ordering of questions (or re- 370  
323 sponses from previous questions) in ways our setup 371  
324 did not account for. 372

325 We used OpenAI’s GPT-3.5. It is possible cer- 373  
326 tain aspects of our characterization are specific to it. 374  
327 In particular, we anecdotally observed that smaller 375  
328 models had difficulty completing the odd-one-out 376  
329 task as far as we could understand; other models 377  
330 likely exhibit more or less similar behavior to hu- 378  
331 mans as well. 379

332 During the survey, multiple respondents men- 380  
333 tioned that the percentile structure made it difficult 381  
334 to discern continuous meaning across the entire 382  
335 dimension scale. This may be because the dimen- 383  
336 sions only hold palpable information at higher lev- 384  
337 els. Regardless, the common strategy employed 385  
338 was to look at the top and bottom objects rather 386  
339 than the ones in the middle. Our percentiles were 387  
340 chosen to align with previous work, but regardless, 388  
341 other methods may elucidate more nuances than 389  
342 our prompt and coding schema did. 390

343 Finally, GPT-3.5 is largely English-trained, and 391  
344 future work may wish to consider examining mod- 392  
345 els trained on data for other languages. 393

346 Our work serves as one data point for understand- 394  
347 ing LLMs. This should be sufficient for giving in- 395  
348 sight into related work, but (especially given the 396  
349 quickly-arriving ubiquity of LLMs and potential 397  
350 for harm; see [Ethics](#)), it is not in isolation nearly 398  
351 sufficient for determining whether LLMs should be 399  
352 used in real-world applications. 400

## 353 Ethics

### 354 Risks

355 Our model illuminates ChatGPT’s behavior in a 401  
356 direct odd-one-out task, and some of the charac- 402  
357 terizing dimensions have strong correlation with 403  
358 previously obtained dimensions that characterize 404  
405

human object-similarity judgment. There is a po- 359  
tential to misinterpret this as meaning ChatGPT 360  
uses these dimensions in the same way humans do 361  
or that these dimensions apply to all tasks ChatGPT 362  
performs. 363

## 364 3.1 Resources

Response-collection was performed using OpenAI’s 365  
GPT-3.5-Turbo-0613 endpoint. The 4,385,040 366  
responses took one week for OpenAI’s 367  
systems to process at a total cost of \$722 USD. 368  
Training was done with NVIDIA P100 GPUs, tak- 369  
ing about 16 hours per model. 370

## 371 Licensing and Artifacts

Our GPT odd-one-out response dataset and model 372  
are available under a CC-BY version 4 licence at 373  
[Supplementary Materials: Odd-One-Out GPT Re- 374  
sponse Set and Model](#). The intended use of our 375  
dataset is general-purpose, so long as it is not harm- 376  
ful. 377

We use the [THINGS images](#) under the public- 378  
domain terms under which it was released. We use 379  
the [THINGS odd-one-out](#) dataset under the terms 380  
of the CC-BY-4.0 license under which it was re- 381  
leased (see bibliography for citation). Its intended 382  
use is to further research (as per the Things Initia- 383  
tive’s [website](#)([Hebart et al., 2019](#))). 384

We use Pandas ([pandas development team 385  
\(2020\)](#); [Wes McKinney \(2010\)](#)) under its BSD 3 li- 386  
cense. We use Scikit-Learn ([Pedregosa et al., 2011](#)) 387  
under another BSD 3 licence. We use SciPy ([Virta- 388  
nen et al., 2020](#)) under the terms of a [similar licence](#). 389  
We use Matplotlib ([Hunter, 2007](#)) under a BSD-like 390  
licence. Finally, we also use PyTorch ([Paszke et al., 391  
2019](#)). We satisfy the licensing terms of it, along 392  
with the previous software packages, by not redis- 393  
tributing the source code. These software packages’ 394  
intended use is scientific and general-purpose ap- 395  
plication, and we satisfy both those criteria. 396

We also use representational similarity analysis 397  
(RSA) ([Kriegeskorte, 2008](#)) and uniform manifold 398  
approximation and projection (UMAP) ([McInnes 399  
et al., 2018](#)). Kriegeskorte and McInnes both likely 400  
intended others to use their algorithms for general 401  
research. 402

We use ChatGPT for some code generation under 403  
its commercial terms. At no point do we provide 404  
sensitive or copyrighted information to it. 405

## Response Collection

All respondents were members of the same research team. However, as responses were collected with respondents' choice of identifying keyword (initials were suggested), that identification was removed from any public release. This minimal information was necessary as respondents were informed they could have their responses deleted should they desire. All responses were part of the research team; no formal recruitment was done. For the same reason, no compensation was given. Respondents knew ahead of time what this project was for, but details were given in the instructions as well.

The instructions given can be found in [Supplemental Data: Response Form](#).

All responses were from graduate students and postdocs at a leading university. The country of origin of the respondents was diverse (only two respondents were from the same country), and all were fluent in English, although for half, it was not a first language.

## References

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language Models are Few-Shot Learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Dhivya Chandrasekaran and Vijay Mago. 2022. [Evolution of Semantic Similarity—A Survey](#). *ACM Computing Surveys*, 54(2):1–37.

Haixing Dai, Zhengliang Liu, Wenxiong Liao, Xiaoke Huang, Yihan Cao, Zihao Wu, Lin Zhao, Shaochen Xu, Wei Liu, Ninghao Liu, Sheng Li, Dajiang Zhu, Hongmin Cai, Lichao Sun, Quanzheng Li, Dinggang Shen, Tianming Liu, and Xiang Li. 2023. [AugGPT: Leveraging ChatGPT for Text Data Augmentation](#). ArXiv:2302.13007 [cs].

Geneviève Desmarais, Maria Cristina Pensa, Mike J. Dixon, and Eric A. Roy. 2007. [The importance of object similarity in the production and identification of actions associated with objects](#). *Journal of the International Neuropsychological Society*, 13(6):1021–1034.

Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. [ChatGPT Outperforms Crowd-Workers for Text-Annotation Tasks](#). ArXiv:2303.15056 [cs].

Martin N. Hebart, Adam H. Dickter, Alexis Kidder, Wan Y. Kwok, Anna Corriveau, Caitlin Van Wicklin, and Chris I. Baker. 2019. [THINGS: A database of 1,854 object concepts and more than 26,000 naturalistic object images](#). *PLOS ONE*, 14(10):e0223792.

Martin N. Hebart, Charles Y. Zheng, Francisco Pereira, and Chris I. Baker. 2020. [Revealing the multidimensional mental representations of natural objects underlying human similarity judgements](#). *Nature Human Behaviour*, 4(11):1173–1185.

M.N. Hebart, O. Contier, L. Teichmann, A.H. Rockter, C.Y. Zheng, A. Kidder, A. Corriveau, M. Vaziri-Pashkam, and C.I. Baker. 2022. [Things-data: A multimodal collection of large-scale datasets for investigating object representations in human brain and behavior](#).

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. [How Good Are GPT Models at Machine Translation? A Comprehensive Evaluation](#). ArXiv:2302.09210 [cs].

J. D. Hunter. 2007. [Matplotlib: A 2d graphics environment](#). *Computing in Science & Engineering*, 9(3):90–95.

Bin Ji. 2023. [VicunaNER: Zero/Few-shot Named Entity Recognition using Vicuna](#). ArXiv:2305.03253 [cs].

Nikolaus Kriegeskorte. 2008. [Representational similarity analysis – connecting the branches of systems neuroscience](#). *Frontiers in Systems Neuroscience*.

Chen Li, Yixiao Ge, Jiayong Mao, Dian Li, and Ying Shan. 2023. [TagGPT: Large Language Models are Zero-shot Multimodal Taggers](#). ArXiv:2304.03022 [cs].

Ellen M. Markman and Jean E. Hutchinson. 1984. [Children's sensitivity to constraints on word meaning: Taxonomic versus thematic relations](#). *Cognitive Psychology*, 16(1):1–27.

Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. 2018. [UMAP: Uniform Manifold Approximation and Projection](#). *Journal of Open Source Software*, 3(29):861.

Daniel Mirman, Jon-Frederick Landrigan, and Allison E. Britt. 2017. [Taxonomic and thematic semantic systems](#). *Psychological Bulletin*, 143(5):499–520.

Anders Giovanni Møller, Jacob Aarup Dalsgaard, Arianna Pera, and Luca Maria Aiello. 2023. [Is a prompt and a few samples all you need? Using GPT-4 for data augmentation in low-resource classification tasks](#). ArXiv:2304.13861 [physics].

512 The pandas development team. 2020. [pandas-](#)  
513 [dev/pandas: Pandas](#).

514 Adam Paszke, Sam Gross, Francisco Massa, Adam  
515 Lerer, James Bradbury, Gregory Chanan, Trevor  
516 Killeen, Zeming Lin, Natalia Gimelshein, Luca  
517 Antiga, Alban Desmaison, Andreas Kopf, Edward  
518 Yang, Zachary DeVito, Martin Raison, Alykhan Te-  
519 jani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,  
520 Junjie Bai, and Soumith Chintala. 2019. [Pytorch:](#)  
521 [An imperative style, high-performance deep learning](#)  
522 [library](#). In *Advances in Neural Information Process-*  
523 *ing Systems 32*, pages 8024–8035. Curran Associates,  
524 Inc.

525 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel,  
526 B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer,  
527 R. Weiss, V. Dubourg, J. Vanderplas, A. Passos,  
528 D. Cournapeau, M. Brucher, M. Perrot, and E. Duch-  
529 esnay. 2011. Scikit-learn: Machine learning in  
530 Python. *Journal of Machine Learning Research*,  
531 12:2825–2830.

532 Linda B Smith. 1981. [Importance of the Overall Simi-](#)  
533 [larity of Objects for Adults’ and Children’s Classifica-](#)  
534 [tions](#). *Journal of Experimental Psychology: Human*  
535 *Perception and Performance*, 7(4):811–824.

536 Fatemeh Torabi Asr, Robert Zinkov, and Michael Jones.  
537 2018. [Querying Word Embeddings for Similarity and](#)  
538 [Relatedness](#). In *Proceedings of the 2018 Conference*  
539 *of the North American Chapter of the Association for*  
540 *Computational Linguistics: Human Language Tech-*  
541 *nologies, Volume 1 (Long Papers)*, pages 675–684,  
542 New Orleans, Louisiana. Association for Computa-  
543 tional Linguistics.

544 J.J. Valenti and Chaz Firestone. 2019. [Finding the “odd](#)  
545 [one out”](#): Memory color effects and the logic of  
546 appearance. *Cognition*, 191:103934.

547 Veniamin Veselovsky, Manoel Horta Ribeiro, and  
548 Robert West. 2023. [Artificial Artificial Artificial](#)  
549 [Intelligence: Crowd Workers Widely Use](#)  
550 [Large Language Models for Text Production Tasks](#).  
551 ArXiv:2306.07899 [cs].

552 Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt  
553 Haberland, Tyler Reddy, David Cournapeau, Ev-  
554 geni Burovski, Pearu Peterson, Warren Weckesser,  
555 Jonathan Bright, Stéfan J. van der Walt, Matthew  
556 Brett, Joshua Wilson, K. Jarrod Millman, Nikolay  
557 Mayorov, Andrew R. J. Nelson, Eric Jones, Robert  
558 Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng,  
559 Eric W. Moore, Jake VanderPlas, Denis Laxalde,  
560 Josef Perktold, Robert Cimrman, Ian Henriksen, E. A.  
561 Quintero, Charles R. Harris, Anne M. Archibald, An-  
562 tônio H. Ribeiro, Fabian Pedregosa, Paul van Mul-  
563 bregt, and SciPy 1.0 Contributors. 2020. [SciPy 1.0:](#)  
564 [Fundamental Algorithms for Scientific Computing in](#)  
565 [Python](#). *Nature Methods*, 17:261–272.

566 Wes McKinney. 2010. [Data Structures for Statistical](#)  
567 [Computing in Python](#). In *Proceedings of the 9th*  
568 *Python in Science Conference*, pages 56 – 61.

Charles Y. Zheng, Francisco Pereira, Chris I. Baker,  
569 and Martin N. Hebart. 2019. [Revealing interpretable](#)  
570 [object representations from human behavior](#). 571

## A Response-Order Counts 572

573 For a set of triplets, each object is either ordered  
574 first, second, or third in their presentation to a re-  
575 spondent. Below are the holistic choice rates for  
576 each in the odd-one-out task for for ChatGPT ([Fig-](#)  
577 [ure 5](#)), for humans ([Figure 6](#)), and for ChatGPT  
578 aggregated ([Figure 7](#)).

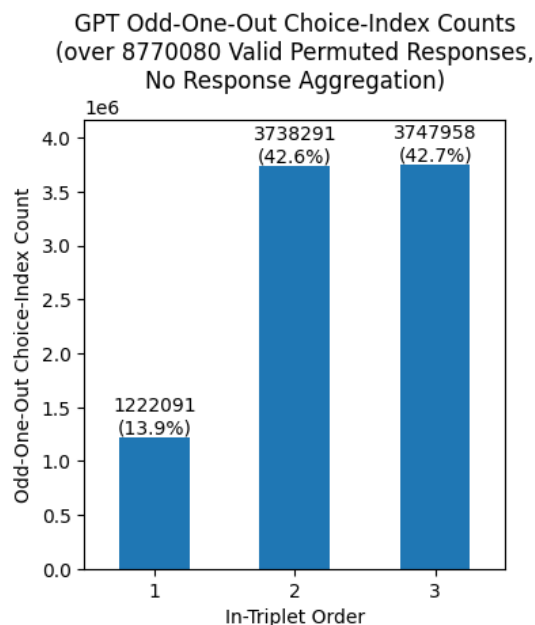


Figure 5: Counts of order-within-triplet responses for raw ChatGPT calls. For example, given a prompt asking about ‘apple’, ‘banana’, and ‘car’, in that order, and a response of ‘car’, this would be a response with an index of 3. These are unbalanced, so we resort to permuting them; see [section 2, Human–GPT Datasets](#) for details of this.

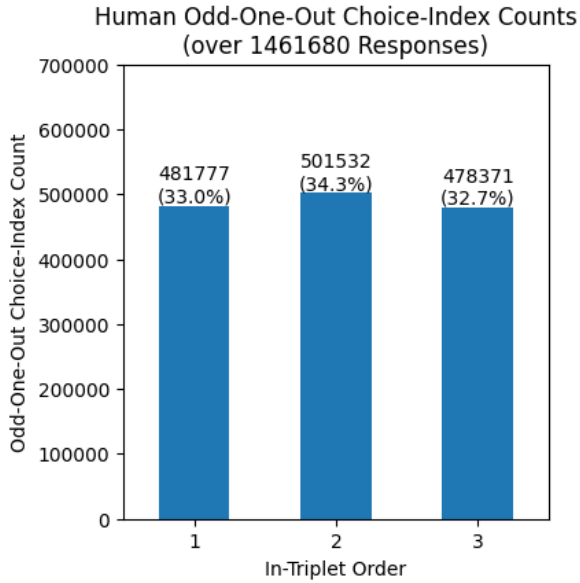


Figure 6: Counts of order-within-triplet responses for adult respondents on the dataset. For example, given a prompt asking about ‘apple’, ‘banana’, and ‘car’, in that order, and a response of ‘car’, this would be a response with an index of 3. These responses are from (Hebart et al., 2020).

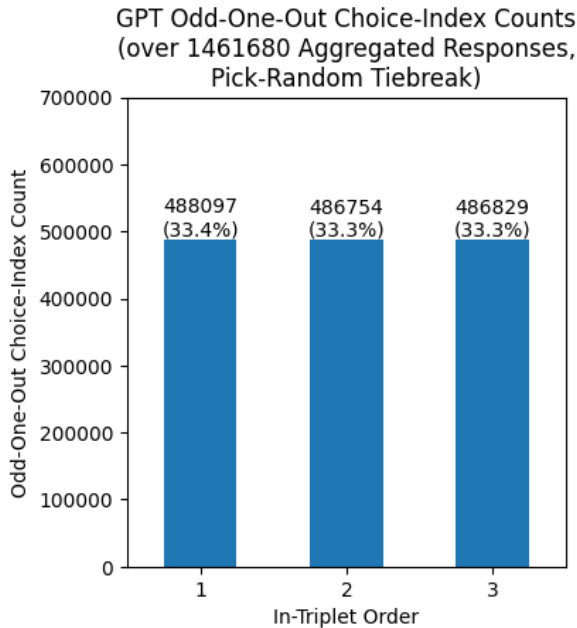


Figure 7: Counts of order-within-triplet responses for aggregated ChatGPT calls. For example, given ‘apple’, ‘banana’, and ‘car’, if the relative majority vote was ‘banana’, this would be a response of index 2. In the case of tiebreaks, in actuality the earliest tiebreaking indexed response was chosen; this is easier to reproduce and works out to be equivalent to choosing randomly due to the orders of the objects within the questions being completely random. See section 2, Human-GPT Datasets for permutation details.

## B Odd-One-Out Prompt

579

The prompts we provided to GPT were of the following form:

580

581

```
<|im_start|>system
Which of the objects are more similar to
each other? Say the object that
doesn't match. Format your choice as
[[object]]<|im_end|>
<|im_start|>user
{object1}, {object2}, {object3}.<|im_end|>
```

582

583

584

585

586

587

588

589

This was intended to be as close to the language used by (Hebart et al., 2020) as possible. Their instruction example is as follows:

590

591

592

```
The three pictures show {object1}, {
object2}, and {object3}. Which are
more similar to each other? Click on
the picture that doesn't match.
```

593

594

595

596

## C Permuted Response Distribution

597

For a given set of three objects, ChatGPT may answer differently when the objects’ order is permuted in the prompt. The rates of agreement of these individual permutations with the accepted aggregate response are given in Figure 8.

598

599

600

601

602

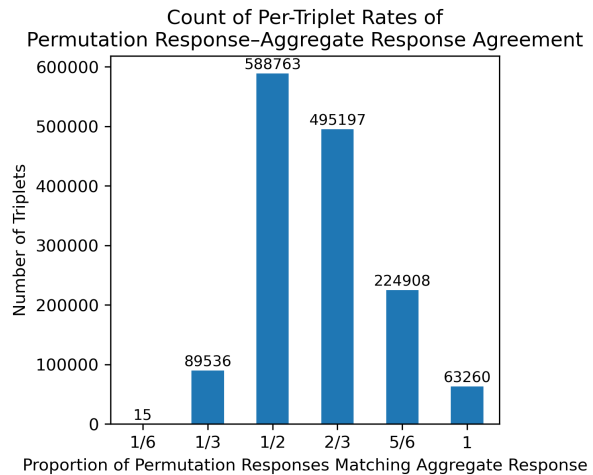


Figure 8: Distribution of the rate of agreement of model permutation responses with the aggregate model response (see section 2, Human-GPT Datasets for permuting details). 1.00 denotes that all 6 permutations of an odd-one-out triplet resulted in the same response;  $\frac{2}{3}$  indicate that 4 of 6 permutations resulted in the same response.  $\frac{1}{2}$  and  $\frac{2}{3}$  indicate possible ties, which were broken by choosing the first response at a tying index. Due to the questions being random ordered, consistently doing this is equivalent to choosing randomly between the options with the most votes.



## D Model Loss

The cross-entropy loss used by the model in training is given here.

$$\begin{aligned}
 H(q, p) & \text{object set is } \{i, j, k\}, \\
 & \text{ } k \text{ is the odd-one-out} \\
 &= \sum_{c \in \{i, j, k\}} q_{c \text{ is the odd-one-out}} \cdot \ln(p_{c \text{ is the odd-one-out}}) \\
 &= -\ln(p_{c_{\text{odd-one-out}}}) \\
 &= -\ln(\sigma(\mathbf{z})_c) = -\ln \frac{e^{z_k}}{e^{z_k} + e^{z_j} + e^{z_i}}
 \end{aligned}$$

where

- $H$  is the cross-entropy loss function
- $i, j, k$  denote the three objects of a triplet, where  $k$  is the true odd-one-out
- $z_c$  where  $c \in \{i, j, k\}$  and  $z_c$  represents the dot product between the vectors of the pair of objects  $\{i, j, k\} \setminus \{c\}$
- $\mathbf{z} = \{z_i, z_j, z_k\}$
- $\sigma$  is the softmax function
- $q$  is the probability of an object being the odd one out (so 100% for the identified odd-one-out, 0% for any other object)
- $p$  is the estimated probability the model gives that a given object is the odd-one-out

For the  $\ell^1$ -norm penalty, we flatten the embedding matrix and take the  $\ell^1$  norm of the resulting vector. We weight this norm by  $\lambda/\text{num\_items}$  and add it to the cross-entropy loss to obtain our full loss.

## E Grid Search Specifics

For a given training set, we perform a grid search: we take steps of 0.0016 over the range  $\lambda \in \{0.0078..0.027\}$  to find a maximum, expanding the search radius if necessary. We then perform ( $k = 4$ )-fold cross-validation ( $(k = 10)$ -fold for the full all-GPT set) in steps of 0.004 to find the optimal lambda in the region around that local maximum. We train on a 90% split for a fixed 1000 epochs for each model, mirroring the setup of Zheng et al. (2019) to ensure convergence. The per-epoch performance and final validation accuracies for the grid-search folds of the all-GPT model are given in Figure 9. The final validation accuracies for those  $\lambda$ s are given in Figure 10, illustrating the degree of local concavity.

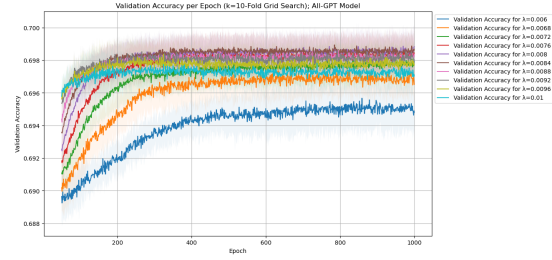


Figure 9: Per-epoch grid search for lambda for the full 1.46-million response all-GPT model.  $\lambda = 0.08$  is the highest-scoring performer.

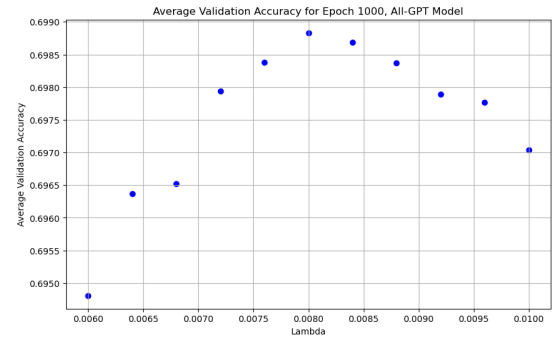


Figure 10: Average validation accuracy from the grid search lambdas at 1000 epochs.  $\lambda = 0.08$  is the highest performer.

## F Dimension Scales

For each dimension, we produced scales with objects whose values spanned the dimension, as in Figure 11.

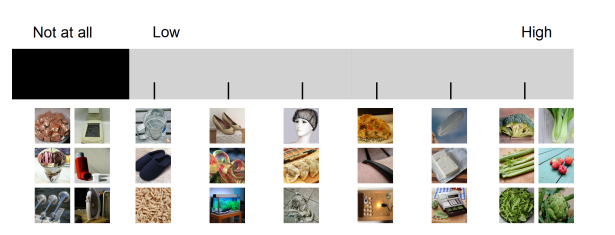


Figure 11: Scale produced for Dimension 12 for annotations

Namely, we made images as seen in Figure 11. The six images on the left have Dimension 12 values at the 0<sup>th</sup>, 1<sup>st</sup>, 5<sup>th</sup>, 10<sup>th</sup>, 15<sup>th</sup>, and 20<sup>th</sup> percentiles for the dimension. The images at the next tick have dimension values at the 33<sup>rd</sup> percentile, and thereafter the images at each successive tick are at a percentile 13.333 more. This continues until the last tick, denoting the 100<sup>th</sup> percentile, where the six top-scoring images are shown.

656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702

## G Dimension Labels

The aggregated dimension names for the 39 largest dimensions of the all-GPT model are given in [Table 1](#).

## H Correlation Heatmaps

The full heatmap of the 39 all-GPT model largest dimensions’ correlations with those of the full human-only model are given in [Figure 12](#). To illustrate the closest dimensions between the labelled GPT and labelled human embeddings, we have done a bipartite max-correlation-as-weight matching between the dimensions of our GPT embedding and the dimensions of the Hebart human embedding in [Figure 13](#) (and vice-versa in [Figure 14](#)).

## I Dimension Reproducibility

To gauge the reproducibility of labelled GPT embeddings, for each labelled dimension, we looked at eight of the other runs of the full GPT models and found, for each one, the maximal column correlation with the labelled dimension. The distribution of these maximal column correlations is given in [\(\)](#). We likewise did this for the labelled human embedding dimensions across full human models in [\(\)](#).

We were also interested in seeing to what extent GPT models captured labelled human dimensions (and vice-versa). [Figure 3](#) gives the maximal correlations of the labelled human embedding dimensions with the columns of random full GPT models. [\(\)](#) gives the maximal correlations of the labelled GPT embedding dimensions with the columns of random full human models.

Finally, we can determine how much worse GPT models were at producing the labelled human dimensions (and vice-versa) by subtracting the maximal correlations of the labelled human dimensions with GPT dimensions from the maximal correlations of the labelled human dimensions with random human dimensions, which are shown in [\(\)](#), and by

## J Dimension UMAP and PCA

We performed Uniform Manifold Approximation and Projection (UMAP) and Principal Component Analysis (PCA) on the aggregated human and GPT embedding dimensions for insight into the dimensions’ spatial relationships. These are given in [Figure 15](#).

## K Mixed Human–GPT RSA

[Table 2](#) gives a table of RSA correlations of the embeddings trained on mixed human–GPT datasets with the baseline human embeddings. Each value is an average across averages. The “Dot RSM” column represents using a dot-product kernel to take a representational similarity matrix for the correlation between the mixed–dataset model and all-human model RSMS, while the “Cos RDM Corr” column represents using cosine similarity to produce RDMs in lieu of those RSMS.

703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713

Dimension Ordering	Aggregate Dimension Label	Dimension Ordering	Aggregate Dimension Label
1	round, outdoors	21	alive/nature/plant-related
2	food-related	22	boats/water-related
3	animal-related, organic	23	box/container-related
4	clothing-related	24	sports-related
5	food, kitchen-related, house	25	small, (flying) insect-related
6	furniture-related	26	music-related
7	gold/jewel, luxury, ostentatious	27	vehicle-related, outdoors
8	transportation/vehicle-related	28	fruit-related
9	gun/explosive, weapon	29	aquatic/sea-related
10	electronics-related	30	crafts, push item through hole
11	(melee) weapon, long/thin	31	wound/rolled, thread-related
12	edible/vegetable-related	32	round, colorful, sports
13	tool-related	33	sanitation, garbage-related
14	(sharp) tools	34	medical (equipment/tools)
15	delicious/sweet liquid/food	35	toy-related
16	(metallic) housing hardware-related	36	vertical, elevated
17	earth/rock-related	37	industrial/mechanical
18	candy/sweet, food	38	paper/literacy-related
19	textiles	39	temperature/temperature-change related
20	container, tableware-related		

Table 1: Aggregate labels for the characterizing dimensions of the labelled GPT model. Labels obtained via the process described in [subsection 2.2](#).

Dataset	Proportion Human Data	Dot RSM Corr.	Cos RDM Corr.
Random Embedding	0	0	-0.01
Full GPT	0	0.437	0.438
Partial Human	0.125	0.853	0.638
Partial Human	0.25	0.897	0.710
Partial Human	0.375	0.916	0.752
Partial Human	0.5	0.924	0.772
Partial Human	0.625	0.930	0.797
Partial Human	0.75	0.928	0.763
Partial Human	0.875	0.933	0.808
Mixed	0.125	0.507	0.502
Mixed	0.25	0.585	0.566
Mixed	0.375	0.667	0.613
Mixed	0.5	0.750	0.680
Mixed	0.625	0.826	0.723
Mixed	0.75	0.887	0.774
Mixed	0.875	0.926	0.809
Full Human	1	0.933	0.808

Table 2: A table of RSA scores for different proportions of GPT data over 10 folds

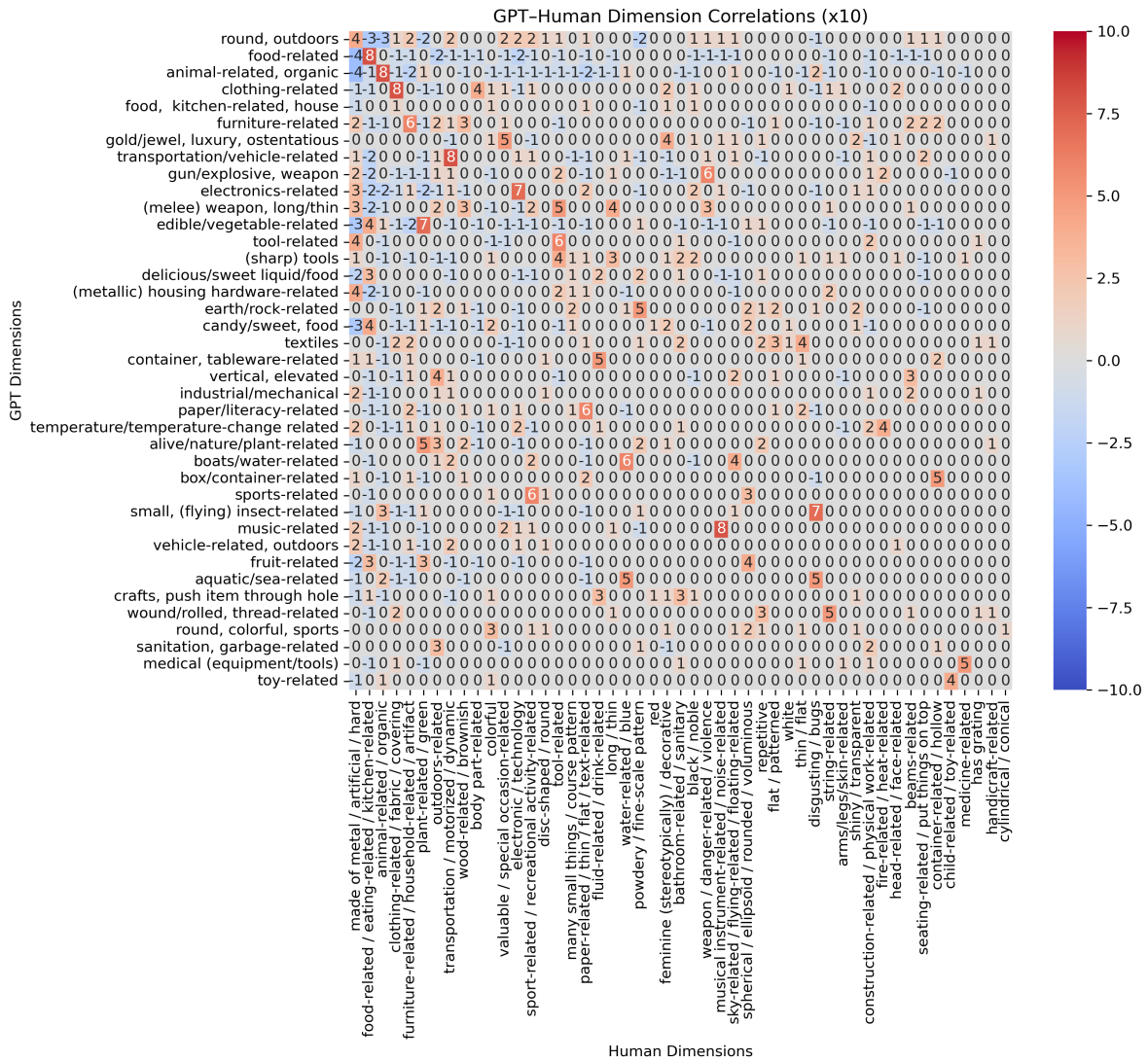


Figure 12: Full correlation heatmap between the dimensions of the all-GPT model and the human model, with aggregate labels on left. Dimensions are ordered by the mean value over objects. Correlations are multiplied by 10 and rounded to the nearest integer for text-size reasons.



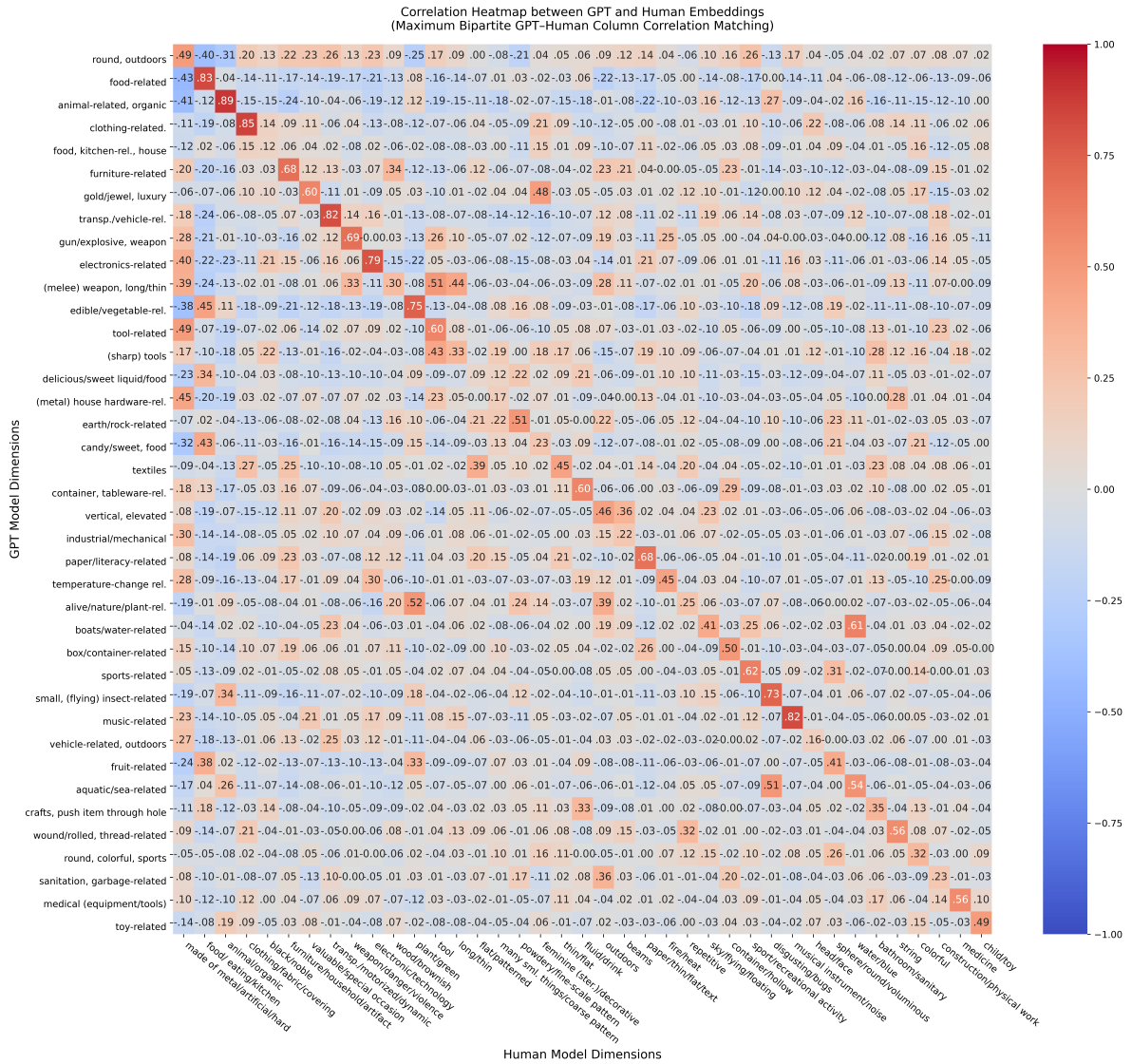


Figure 13: Correlation heatmap between each labelled GPT embedding dimension and the closest labelled human embedding dimension under bipartite max-correlation matching. GPT dimensions ordered by their mean value over all objects.

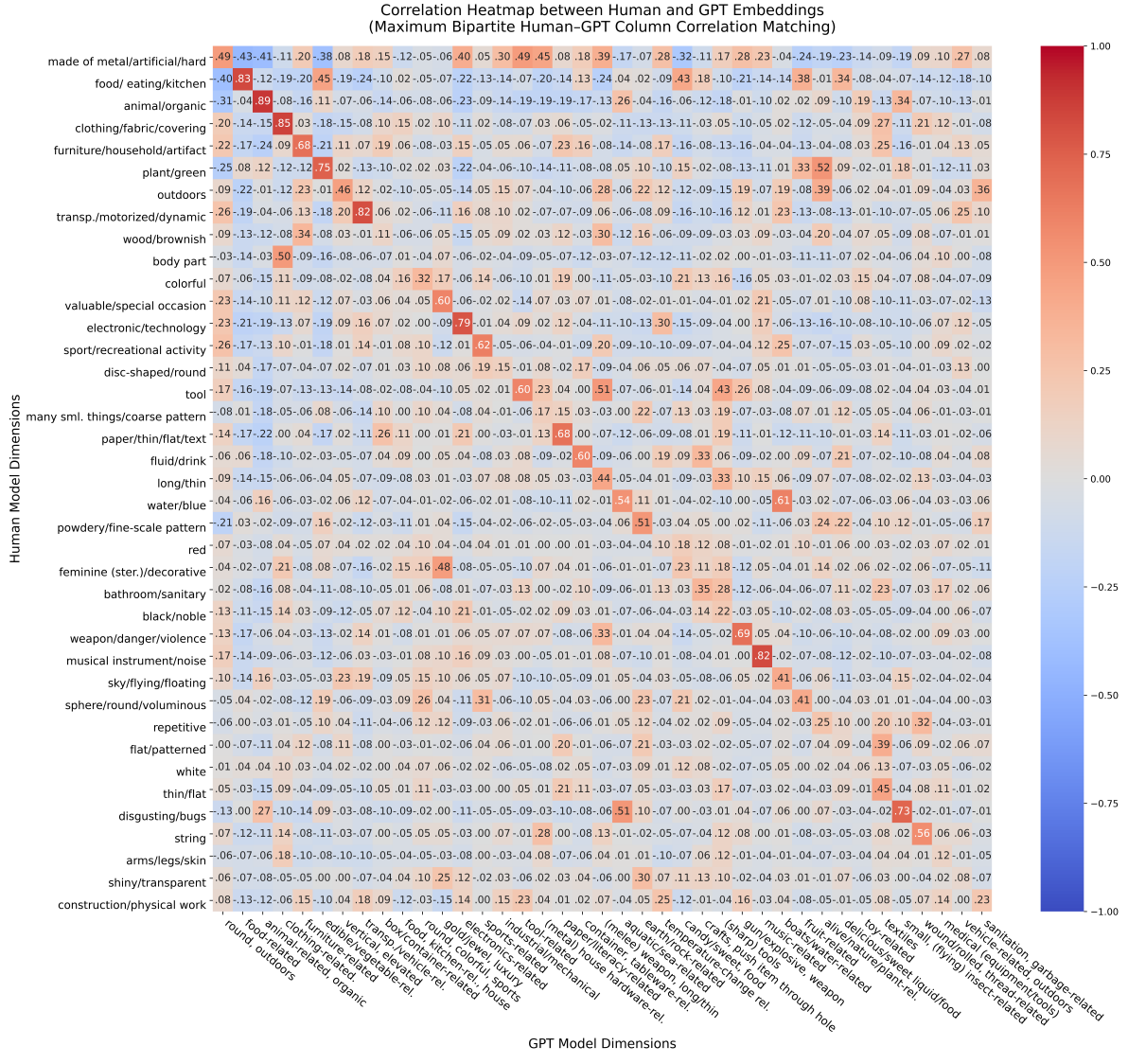


Figure 14: Correlation heatmap between each labelled human embedding dimension and the closest labelled GPT embedding dimension under bipartite max-correlation matching. Human dimensions ordered by their mean value over all objects.

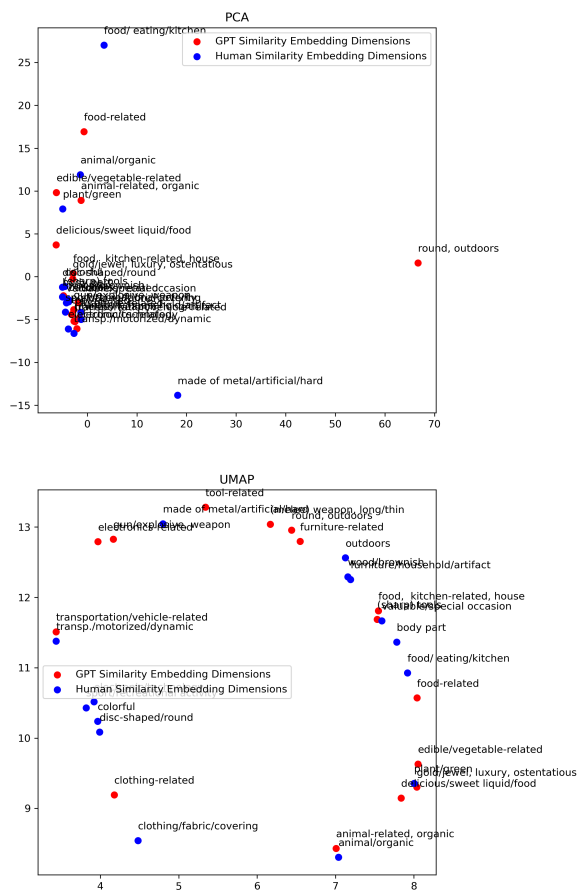


Figure 15: UMAP and PCA performed on the aggregated GPT-only and human-only embeddings' dimensions.