

UNCOVERING THE COMPUTATIONAL INGREDIENTS OF HUMAN-LIKE REPRESENTATIONS IN LLMs

Anonymous authors

Paper under double-blind review

ABSTRACT

The ability to translate diverse patterns of inputs into structured patterns of behavior has been thought to rest on both humans’ and machines’ ability to learn robust representations of relevant concepts. The rapid advancement of transformer-based large language models (LLMs) has led to a diversity of computational ingredients — architectures, fine tuning methods, and training datasets among others — but it remains unclear which of these ingredients are most crucial for building models that develop human-like representations. Further, most current LLM benchmarks are not suited to measuring representational alignment between humans and models, making existing benchmark scores unreliable for assessing if current LLMs are making progress towards becoming useful cognitive models. Here, we address these limitations by first evaluating a set of over 77 models that widely vary in their computational ingredients on a triplet similarity task, a method well established in the cognitive sciences for measuring human conceptual representations, using concepts from the THINGS database. Comparing human and model representations, we find that models that undergo instruction-finetuning and which have larger dimensionality of attention heads are among the most human aligned. We also find that factors such as choice of activation function, multimodal pretraining, and parameter size have limited bearing on alignment. Correlations between alignment scores and scores on existing benchmarks reveal that while some benchmarks (e.g., BigBenchHard) are better suited than others (e.g., MUSR) for capturing representational alignment, no existing benchmark is capable of fully accounting for the variance of alignment scores, demonstrating their insufficiency in capturing human-AI alignment. Taken together, our findings help highlight the computational ingredients most essential for advancing LLMs towards models of human conceptual representation and address a key benchmarking gap in LLM evaluation.

1 INTRODUCTION

The success of deep neural network models on diverse domains ranging from perception (Yamins et al., 2014; Yamins & DiCarlo, 2016), to robotic manipulation (Finn et al., 2016), to language understanding (Tuckute et al., 2024) has partly been attributed to their ability to learn powerful *representations* that can effectively aid in translating inputs into appropriate outputs. The continued development of a specific class of these systems — transformer-based large language models (LLMs) — and their capabilities on naturalistic human tasks might imply that these models, through training on large enough corpora of text (and often image) data and with the correct architectural ingredients, come to possess representations that are largely isomorphic to humans’ mental representations. This implication is critical for many domains of cognitive science that have long sought to characterize the nature of human conceptual representations and the operations performed on them that lead to naturalistic behavior. Indeed, there is a rich tradition of using deep neural networks as *computational cognitive models*: using their learned representations as proxies for human representations, and often mapping the evolution of learned representations to development of human conceptual representations (Jackson et al., 2022; Warstadt et al., 2025). While a patchwork of evidence in various domains has shown concordances in representations between language models and people’s behaviors and patterns in their neural activity, it remains unclear what properties of current LLMs are most predictive of strong human-model alignment. That

is, given the engineering-centric drive behind current LLM development, it is difficult to isolate what *computational ingredients* (model architecture, model size, instruction-tuning, activation functions, etc.) are important for giving rise to human-aligned conceptual representations in models for a large variety of concepts.

Identifying the computational ingredients underlying human-model alignment is critical for several reasons. While frontier models have achieved remarkable performance on diverse benchmarks spanning scientific reasoning (Rein et al., 2024; Suzgun et al., 2022), software engineering (Chen et al., 2021; Jimenez et al., 2023), and chart and humor understanding (Masry et al., 2022; Methani et al., 2020; Mukherjee et al., 2025; Zhou et al., 2025; Hessel et al., 2022) (domains thought to require ‘human-like’ thinking), we lack a unified theory explaining *why* model performance on one task might transfer to another and reasons for inconsistencies in benchmark rankings for the same model. Recent work has shown that optimizing for benchmark accuracy alone is insufficient for building models that fail in human-like ways and that are generally aligned with humans (Fel et al., 2022; Ying et al., 2025). We argue that prioritizing representational alignment (Sucholutsky et al., 2023; Sucholutsky & Griffiths, 2023b; Muttenthaler et al., 2024) offers a promising path forward. Measuring the similarity between model and human internal representations can provide a more unified account of model capabilities, guide the development of more robustly aligned systems (Collins et al., 2024; Muttenthaler et al., 2024; Peterson et al., 2018), and accelerate cognitive science by yielding more faithful computational models of the human mind.

2 RELATED WORK

Quantifying Human and Model Representations using Large-Scale Concept Datasets

Mapping the geometry human conceptual knowledge has long been a central goal of cognitive science (Rogers & McClelland, 2004; Shepard, 1980; Tversky, 1977; Rumelhart et al., 1986). Early attempts to do so relied on explicit feature ratings and pairwise similarity judgments, which scaled poorly limiting their utility in capturing the structure of the vast inventory of human concepts (Rosch, 1975; Nosofsky, 1984; Shepard, 1980). While structured probabilistic models have been useful in explaining aspects of concept learning, they too have been limited often being restricted to specific relational structures, limiting their scope and generalizability (Zhao et al., 2024; Wong et al., 2022; Kemp & Tenenbaum, 2008). Recent years have seen a resurgence of similarity-based measures of characterizing human conceptual knowledge, driven by scalable algorithms that convert human judgments into expressive high-dimensional *semantic embeddings* (King et al., 2019; Mukherjee & Rogers, 2025; Hebart et al., 2019; Suresh et al., 2023b; Peterson et al., 2018; Kriegeskorte & Mur, 2012; Cichy et al., 2019; Hebart et al., 2023; Muttenthaler et al., 2022b; Jamieson et al., 2015; Sievert et al., 2023) capable of expressing latent dimensions organizing human semantic knowledge. Specifically, triadic similarity judgment tasks — where participants select the most similar item from a triplet—are a particularly powerful paradigm (Tamuz et al., 2011; Wah et al., 2014; Kleindessner & von Luxburg, 2014). Embeddings derived from this measure have shown to be good predictors of both neural and behavioral patterns in humans across a variety of semantic tasks (Mukherjee & Rogers, 2025; Mukherjee et al.; Hebart et al., 2023; Colon & Rogers, 2023; Suresh et al., 2023b; Marjeh et al., 2022; Mukherjee et al., 2022; Mur et al., 2013; Suresh et al., 2024).

A complementary advancement in recent years that has galvanized research in representational alignment is the development of concept datasets that aim to capture the diversity of object concepts humans reason about (Mehrer et al., 2021; Hebart et al., 2019; Giallanza et al., 2024; Suresh et al., 2025). Datasets like THINGS (Hebart et al., 2019) and ECOSSET (Mehrer et al., 2021) consist of concept sets and associated images along with rich psychological and neural metadata (Hebart et al., 2023) that allow for rich comparisons between human and model representations along multiple levels of analysis — behavioral, neural, and computational.

Representational Alignment: Measurement and Downstream Task Performance The nascent field of representational alignment has developed formal tools for comparing internal representations across biological and artificial systems (Kriegeskorte et al., 2008; Sucholutsky et al., 2023; Barbosa et al., 2025; Oswal et al., 2016). Core measurement techniques include Representational Similarity Analysis (RSA) using correlation between similarity matrices derived from behavioral, neural, or neural network activation data (Kriegeskorte et al., 2008; Nili et al., 2014), Procrustes analysis for finding optimal linear alignments (Schönemann, 1966), and Centered Kernel Alignment (CKA)

108 for comparing representations invariant to linear transformations (Kornblith et al., 2019; Williams
109 et al., 2021). Recent work has further developed more sensitive metrics including shape-based
110 metrics (Williams et al., 2021), topological measures (Barannikov et al., 2021), information-
111 theoretic approaches (Bansal et al., 2021), and regularized regression-based methods (Oswal et al.
112 (2016) each of which iteratively address issues relating to model regularization, generalization,
113 and error-bounds. Growing evidence demonstrates that representational alignment with humans
114 correlates with practical benefits for model performance. Models with higher human alignment
115 show improved few-shot learning (Sucholutsky & Griffiths, 2023a; Huh et al., 2024), better out-
116 of-distribution generalization (Moschella et al., 2023; Norelli et al., 2023), enhanced adversarial
117 robustness (Engstrom et al., 2019), and more human-like error patterns (Geirhos et al., 2021;
118 Fel et al., 2022). These findings suggest that human-aligned representations capture meaningful
119 invariances that support robust intelligence.

120 **Computational Ingredients Driving Alignment** Understanding which design choices yield
121 human-aligned representations in foundation models can help uncover general principles
122 undergirding human intelligence. For instance, in vision, early work showed that optimizing models
123 for categorization (Yamins et al., 2014) and later for contrastive objectives (Konkle & Alvarez, 2022;
124 Zhuang et al., 2021) produced representations closely aligned with activity in the human visual
125 cortex, offering insight into the kinds of objective functions the brain may solve to construct visual
126 representations. Extending this line of inquiry, Muttenthaler et al. (2022a) evaluated over 75 models
127 spanning architectures, objectives, and datasets, and found that training data and objectives strongly
128 predicted human-model alignment, whereas architecture and scale had minimal impact. This finding
129 challenges scaling law accounts (Kaplan et al., 2020; Cherti et al., 2023) that predict generally larger
130 models should always perform better (at most tasks). Other key findings regarding the importance
131 of different computational ingredients include the finding that self-supervised contrastive methods
132 align more closely with human vision than non-contrastive approaches (Chen et al., 2020; Zbontar
133 et al., 2021; Caron et al., 2020); multimodal image–text contrastive training improves alignment
134 over vision-only training (Radford et al., 2021; Jia et al., 2021; Pham et al., 2022); training on
135 diverse, naturalistic datasets such as JFT-3B improves alignment beyond ImageNet (Zhai et al.,
136 2022; Dehghani et al., 2023); enforcing shape bias on models enhances alignment while texture
137 bias impairs it (Geirhos et al., 2019; Hermann et al., 2020); and intermediate layers often show
138 stronger perceptual alignment than final layers (Berardino et al., 2017; Kumar et al., 2022). Similar
139 approaches, dissecting the predictive power of different model building choices, have not been
140 brought to bear in investigations into conceptual structure in LLMs- leaving unclear what the relative
141 importance of different computational ingredients are when building foundation models of human
142 semantic knowledge.

143 **Language Models and Conceptual Alignment** Recent work has begun extending analyses of
144 representational alignment to language models (Suresh et al., 2023b; Marjeh et al., 2022). Early
145 studies showed that contextualized embeddings from models like BERT and GPT-2 could predict
146 human similarity judgments (Bommasani et al., 2020; Grand et al., 2022). More recent work
147 demonstrates that LLMs can achieve remarkably high alignment with human conceptual structure
148 (Marjeh et al., 2024a;b; Mukherjee et al., 2024; Suresh et al., 2025; 2023a; Mukherjee et al.,
149 2023b), with instruction-tuned models showing particular advantages (Tong et al., 2024; Gurnee
150 & Tegmark, 2024). Several studies have used triplet tasks specifically to probe LLM representations
151 (Nam et al., 2024; Studdiford et al., 2025; Rathi et al., 2024; Suresh et al., 2023b), finding that larger
152 models and those trained on more diverse data show better alignment. However, systematic analyses
153 disentangling the effects of scale, architecture, training data, and objectives remain limited. Work on
154 multimodal models suggests that vision-language training improves conceptual alignment beyond
155 either modality alone (Yuksekgonul et al., 2023; Thrush et al., 2022; Qin et al., 2025), though the
156 mechanisms remain unclear.

156 **Summary of key contributions.** We leverage the open-weight model ecosystem to construct a
157 suite of 75+ models spanning diverse computational ingredients. Using concepts and semantic
158 embeddings from THINGS dataset (Hebart et al., 2019) in conjunction with a triadic similarity
159 judgment paradigm (Jamieson et al., 2015; Hebart et al., 2023; Sievert et al., 2023), we obtain
160 both human and model-specific conceptual embeddings using analogous methods to ensure model-
161 fair comparisons (Firestone, 2020). We then compare these embeddings to human data and apply
statistical models to identify the ingredients most predictive of human-model alignment. Lastly, we

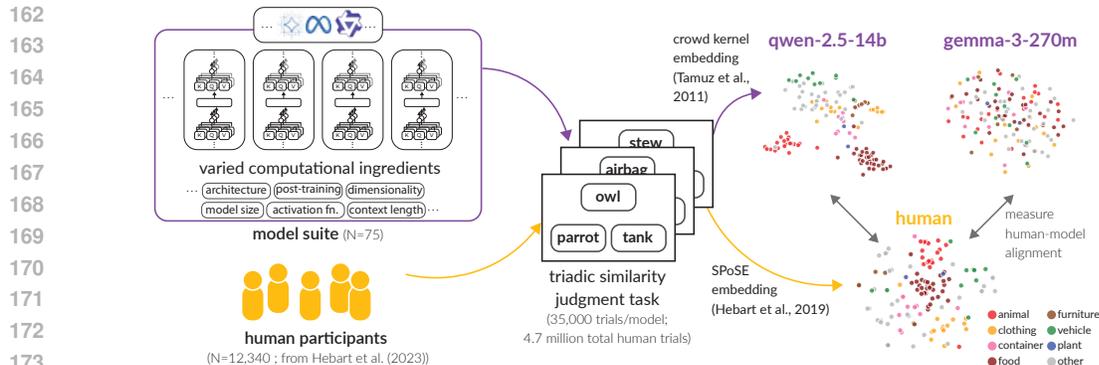


Figure 1: *Method for estimating human-model alignment.* We collected 35k triplet similarity judgments from each of the 77 models in our suite. We computed semantic embeddings based on these judgments and compared the representational geometry of model embeddings to human embeddings derived from Hebart et al. (2023).

relate alignment to performance on standard LLM benchmarks, highlighting points of convergence and divergence between representational alignment and task capabilities.

3 METHODS

Concept Dataset To investigate how various LLM characteristics influence human alignment, we carefully curated a set of test concepts that would provide robust and interpretable results. Starting from the full set of 1,854 THINGS object concepts (Hebart et al., 2019; 2023), we subsampled 128 concrete, real-world object concepts (e.g., “lion”, “banjo”, “car”) that spanned four main cognitive axes of variation: familiarity, artificiality, animacy, and size based on prior work (Mukherjee et al., 2023a). This sampling strategy ensured that we had adequate semantic diversity in our concept dataset (representative of the full THINGS database) while also ensuring that we could generate semantic embeddings for a large suite of models at scale in a tractable manner. The complete list of object concepts, along with detailed descriptions of our concept selection criteria and validation procedures, can be found in Appendix A.3.

Model Suite We evaluate a range of open-source models varying across several key computational ingredients: model architecture, primary activation functions used, pre-training modality, whether models were instruction fine-tuned or not, scale (both in terms of model parameters and training tokens), context length, dimensionality of attention heads/MLPs/residual streams. While in principle it is possible to offer even more granular analyses of computational ingredients, these factors constitute major architectural and training decisions when building foundation models, which we believe provide adequate resolution to investigate the drivers of human-model alignment.

We evaluated a set of 77 open-weight transformer models including models from popular families (Llama, Qwen, Gemma, etc.). The full set of models can be seen in Table 4. Notably, we included models that varied in size from 100M parameters to 42B parameters, ranged in the degree of post-training (from base to instruction fine-tuned), and architecture-style (decoder-only, encoder-decoder, and state space models (SSMs) (Voelker & Eliasmith, 2018)).

Triadic similarity judgment task The triadic similarity judgment task can be framed in two equivalent ways: (1) **anchored similarity judgments** where participants are presented with a triplet of items (images or words depicting different concepts) and select which of two option items is most similar to an anchored reference item, or **odd-one-out judgments** where they identify which of the three presented item is the odd-one-out. Both framings have been successfully employed to study conceptual representations in humans and models (Hebart et al., 2019; 2023; Mukherjee et al., 2023c; Colon & Rogers, 2023; Suresh et al., 2023b; Studdiford et al., 2025; Mukherjee et al.; Muttenthaler et al., 2022b). In this work, we leverage existing human similarity ratings data (Hebart et al., 2023) and corresponding embeddings, which were collected using the odd-one-out paradigm and embedded using the sparse positive similarity embedding (SPoSE) method, respectively. For model evaluations, we adopted the anchored similarity judgment paradigm for computational efficiency, deriving embeddings using ordinal embedding techniques (Jamieson et al.,

2015; Tamuz et al., 2011) following recent work (Suresh et al., 2023b; Mukherjee et al.; Colon & Rogers, 2023). Crucially, the ordinal embedding approach provides theoretical bounds on sample complexity (Jamieson et al., 2015), allowing us to determine the number of triplet comparisons needed for faithful representations a priori (see Section 3 for details). While these approaches differ in framing and embedding computation, they yield comparable representational structures. This methodological choice allows us to maximize the use of existing high-quality, validated human data while efficiently collecting model responses with known sampling requirements in a scalable manner, ensuring robust human-model comparisons across a large suite of models.

Human Semantic Embeddings Human semantic embeddings were obtained by Hebart et al. (2023) from the behavioral judgments of $N = 12,340$ online participants in a triplet odd-one-out task. On each trial, participants were shown three items from the THINGS database and were prompted to select the most dissimilar object ("Which is the odd one out?"). In total, 4.7 million behavioral triplet judgments were obtained for the full set of 1,854 concepts. To estimate an embedding space from the set of individual triplet judgments, Hebart et al. (2023) used the SPoSE algorithm (Hebart et al., 2023), which learns low-dimensional vectors for each object such that (1) pairs of objects judged as similar are placed closer in the embedding space, and (2) dimensions are constrained to be sparse in order to yield human-interpretable properties. The resulting space captured meaningful representational structure from the individual judgments of participants, revealing core dimensions along which human semantic representations are organized.

Model Semantic Embeddings While much work assessing representational similarity between models and humans focus on extracting activation vectors from models and computing their similarity to human behavior-based representations (using methods like RSA (Kriegeskorte et al., 2008)), we instead estimated model semantic embeddings from model similarity judgments using a triadic comparison task akin to those often used to estimate human embeddings. We did so for two principled reasons. Firstly, while computations in transformer layers and attention heads have been linked to human language processing (Tuckute et al., 2024), there is no general framework for assessing which parts (layers, heads) of models (across different model families) are likely to carry signature of human semantic knowledge, making it difficult to assess what constitutes a fair comparison between model activations and human representations. Second, due to the ability of modern LLMs to respond to structured tasks when prompted in natural language and given the relative simplicity of the triadic judgment task (often single token responses suffice), we felt that administering the same test to both models and humans constitutes a species-fair (Firestone, 2020) comparison approach. Further, this allows for a unified embedding method to be applied to both systems ensuring that any discrepancies in alignment are not attributable to vastly different embedding methods.

For each model in our suite, we collected 35,000 triadic judgment responses. Using the human semantic embeddings, we estimated that a rank 30 space explains over 95% of the variance in embeddings (see A.7). Based on known bounds (Sievert et al., 2023), we require $(n \cdot d \cdot \log_2(n))^1 \approx 26,900$ judgments on random sets of triplets to estimate a robust embedding. We collected 35,000 to account for noise in the judgments. Each triplet trial was randomly sampled from the $\binom{128}{3}$ possible triplet trials. We evaluated each model using its default huggingface sampling settings using a common system-level prompt — "You are a helpful assistant who gives responses to questions". We used the following prompt for each trial for each model — QUESTION: Which item is most similar to item_x: item_y or item_z? Respond only with the item name. If we were evaluating a base model (not instruction-tuned), we appended Answer: after the initial prompt. Specific instruction template formats were applied where appropriate. Each prompt was evaluated independently.

Using these similarity judgments, we fit an ordinal embedding algorithm (Tamuz et al., 2011; Sievert et al., 2023) that organized the semantic embedding space such that items that were judged to be similar more often were pulled closer together in this space. We fit 30D embeddings based on the dimensionality of human semantic embeddings. We also held out 20% of the data during the embedding fitting process and evaluated the quality of the embeddings by monitoring the crowd-kernel loss (Tamuz et al., 2011) to ensure the embeddings were of high quality.

¹ n is the number of items and d is the expected dimensionality of the space

4 RESULTS AND DISCUSSION

We first report on how predictive different computational ingredients were of human-model alignment, initially considering each ingredient separately via individual correlations or t -tests (Kim, 2015), then assessing relative contributions of each via mixed linear models (Lindstrom & Bates, 1988). Finally, we evaluate the relationship between human-model alignment and model performance on other key LLM benchmarks (Han et al., 2024; Rein et al., 2023; Srivastava et al., 2023; Zhou et al., 2024; Hendrycks et al., 2020). Our primary alignment measure was the variance in human semantic embeddings explained by model embeddings (R^2) after Procrustes alignment (Gower, 1975). Procrustes transforms find the optimal linear transformation (rotation, scaling, translation) between two vector spaces of equal dimensionality to minimize squared distances between corresponding points. Using human embedding variance (human SS) as the baseline, we compute $R^2 = 1 - \frac{\text{residual SSE}}{\text{human SSE}}$. Other alignment measures such as CKA (Kornblith et al., 2019) and RSA (Kriegeskorte et al., 2008) were highly correlated with Procrustes R^2 (Appendix A.7).

4.1 CONTRIBUTIONS OF MODEL COMPUTATIONAL INGREDIENTS TO ALIGNMENT

Instruction tuning leads to greater alignment. Models with instruction fine-tuning were significantly better aligned than those without ($t = 5.11, p < 0.001$; see Figure 2A). Qualitatively, instruction fine-tuning clusters the representations of semantic categories, such that within-category items become more proximal in representation space and between-category items become more distal (see Appendix 7). Thus, the same post-training techniques that enable models to be more adept at following prompt instructions also bring their representations closer into alignment with humans’.

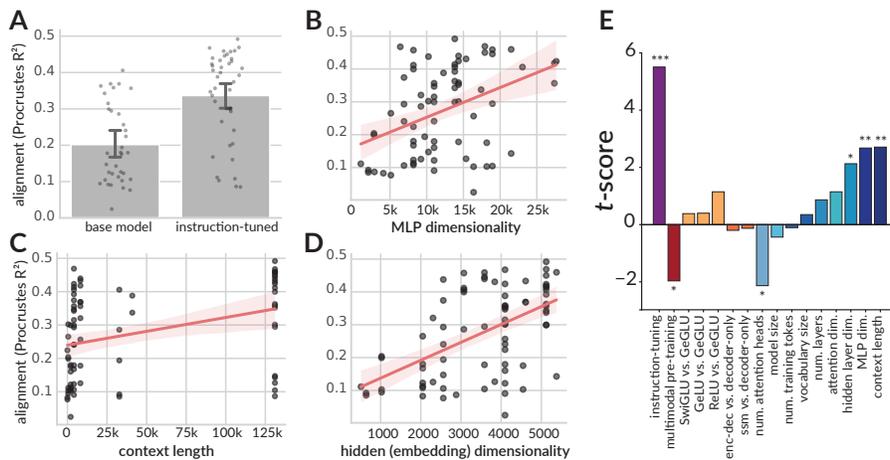


Figure 2: **Computational ingredients predictive of human model alignment.** A., B., C., and D. show the relationship between instruction-tuning, MLP dimensionality, context length, and embedding dimensionality on Procrustes R^2 . These were the four ingredients most predictive in the mixed-effects regression model. E. t -scores for each predictor, which highlight the relative contribution of each ingredient towards alignment when considered in a single statistical model.

Model architecture dimensionality corresponds to greater alignment. We found a strong positive relationship between the dimensionality of individual model architecture components — MLPs and attention heads — and the degree of model-human alignment. Specifically, the number of per-layer attention heads ($r = 0.52, p < 0.001$), attention matrix dimensionality ($r = 0.30, p = 0.026$), embeddings dimensionality ($r = 0.54, p < 0.001$), and MLP dimensionality ($r = 0.40, p < 0.001$) each correlated positively with human representational alignment. Figure 2B and D show the individual correlations between MLP and embedding dimensionality and model-human Procrustes R^2 (the attention dimensionality was not significant in multiple regression models reported in subsequent sections). Thus the greater expressivity granted by larger latent activation spaces and more attention heads predicts greater human alignment.

Multi-modal pretraining has no effect on alignment. While much of human concept learning is inherently multimodal (Rogers & McClelland, 2004; Patterson et al., 2007), there is limited work

324 showing whether vision-language pretraining (the most prominent class of multimodal training)
 325 leads to more aligned semantic representations (cf. (Qin et al., 2025)). In our model suite, however,
 326 multimodal (image) pre-training had no independent effect on human-model alignment relative to
 327 text-only models ($t = 0.34, p > 0.05$; Appendix Figure A.8).

328 **Larger model sizes showed greater human alignment.** For visual representations, Muttenthaler
 329 et al. (2022a) found that larger models need not yield more human-like representations. Neural
 330 scaling laws (Kaplan et al., 2020), however, might predict that larger and deeper models should
 331 be more performant, and perhaps more aligned. In accordance with this prediction, human-model
 332 alignment scaled approximately linearly with both model parameter count and the number of model
 333 layers ($r_{param} = 0.45, r_{layers} = 0.41, p < 0.001$; Appendix Figure A.8). Thus, larger scale
 334 predicts greater human-model alignment.

335 **Alignment increased with amount of training data.** Models exposed to more tokens in pretraining
 336 generally exhibited more human-like representations ($r = 0.33, p < 0.01$; Appendix Figure A.8).
 337 Although model training data also likely varied in quality and composition—thus limiting strong
 338 conclusions about which properties matter most—our results indicate that the *amount* of training data
 339 can predict alignment.

340 **Alignment scales with context length.** The triadic similarity task does not require a long
 341 context length, but we nevertheless found that models with a greater maximum context-length also
 342 demonstrated greater human alignment ($r = 0.35, p < 0.01$; Figure 2). This finding could be
 343 linked to post-training regimes that unlock long-context reasoning, which consequently lead models
 344 to perform well on short-context tasks as well.

345 **Choice of activation function does not affect alignment.** Activation functions have seen iterative
 346 development over the years with recent variants (e.g., SwiGLU (Shazeer, 2020)) being particularly
 347 well-suited to managing gradient flows during training and finetuning experiments leading to faster
 348 and more stable training regimes. But do choices in activation functions ground out in alignment?
 349 We found no advantage for any particular function in increasing human alignment ($F(3, 64) =$
 350 $1.36, p = .26$; Appendix Figure A.8).

351 **Larger vocabulary size does not increase alignment.** One might expect that larger model
 352 vocabularies unlock greater expressivity for models, leading them to be more human-aligned.
 353 Against this prediction, we found no relationship between model vocabulary size and human-
 354 alignment ($r = 0.10, p > 0.05$; Appendix Figure A.8).

355 **Which computational ingredients have the greatest relative contribution towards human-
 356 model alignment?** Thus far, we have investigated how different computational ingredients impact
 357 human-model alignment in isolation. In reality, many of these ingredients vary across models at the
 358 same time, leaving open the questions of (1) which ingredients are most important *relative to others*
 359 and (2) which ingredients account for unique variance after others are considered.

360 To answer these questions, we fit a mixed-effects multiple regression model (Lindstrom & Bates,
 361 1990) predicting Procrustes R^2 from *all* computational ingredients. Mixed-effects models combine
 362 the interpretability of OLS with random effects to capture item-level variance (Barr et al., 2013).
 363 We included random intercepts for model family (e.g., Llama, Qwen) to capture variance related to
 364 these families not grounded out in our set of ingredients, and fixed effects for all ingredients. The
 365 fitted model explained substantial variance in alignment ($R^2 = 0.78$; Fig. 2E). Instruction fine-
 366 tuning emerged as the strongest predictor ($\beta = 0.13, SE = 0.024, p < 0.001$). Dimensionality
 367 measures — MLP ($\beta = 0.049, SE = 0.018, p < 0.01$), embedding/hidden layers ($\beta = 0.048,$
 368 $SE = 0.022, p < 0.05$) — and context length ($\beta = 0.042, SE = 0.015, p < 0.01$) also accounted
 369 for significant unique variance. In contrast to the independent analyses, attention dimensionality
 370 was not significant, and more attention heads predicted worse alignment after accounting for other
 371 factors ($\beta = -0.045, p < 0.05$). Likewise, multimodal pretraining, which showed no independent
 372 effect on alignment, predicted reliably lower alignment after accounting for other factors ($\beta =$
 373 $-0.102, SE = 0.052, p < 0.05$). Thus, current multimodal training regimes may actually limit
 374 rather than improving human-model alignment. No other ingredient explained significant unique
 375 variance when others were included in the mode (see effect sizes in Figure 2E).

376 **A case study on the effects of post-training paradigms on alignment.** Modern LLM systems
 377 undergo various stages of post-training including *supervised fine-tuning* (SFT) (Ouyang et al., 2022),

378 *direct preference optimization* DPO (Rafailov et al., 2023), *reinforcement learning from human*
 379 *feedback* (RLHF) (Ouyang et al., 2022), and *instruction tuning* (IT) (Wei et al., 2021). It remains
 380 unclear what impact these post-training steps have on human-model alignment. Progress on this
 381 front has been limited due to a lack of open-weight models that have been checkpointed at the
 382 various possible stages. OLMo (Groeneveld et al., 2024) is a notable exception since it is available
 383 at four stages of post-training, including the base model (no post-training), SFT, DPO, and IT and is
 384 available at two model sizes (7B and 13B). Instella (Liu et al., 2025) and Llama-3.1 (Meta
 385 AI, 2024) also have similar checkpointed model variants available with Instella having an additional
 386 stage of pretraining checkpointed.

387 Using these models, we asked whether successive stages of post-training led to models’
 388 representations becoming more aligned with humans’. Figure 3 shows how alignment changes as a
 389 function of post-training. Generally, we found that the SFT and DPO stages led to improvements in
 390 alignment whereas the effect of RLVR was more muted.

391 4.2 RELATIONSHIP BETWEEN ALIGNMENT AND ESTABLISHED BENCHMARKS

392
 393 Prior work emphasizes that performance on a given model evaluation should “*translate to*
 394 *similar improvements on any other valid and reasonable evaluation data*” (Bowman & Dahl,
 395 2021). But these arguments were developed in the context of specific NLP tasks, leaving
 396 it unclear whether this notion extends to representational alignment. To understand the
 397 relationship between alignment and performance on standard LLM benchmarks we estimated the
 398 correlation between alignment, measured using three metrics (Procrustes R^2 , CKA, and RSA
 399 correlation), and model scores on five key benchmarks – MMLU, IFEval, BigBenchHard
 400 (BBH), GPQA, and MUSR. Figure 4A shows the relationship between alignment (Pro-
 401 crustes R^2) and scores on BBH, the benchmark most strongly correlated with alignment.
 402

403 While models that performed well on BBH
 404 generally had stronger alignment, we found
 405 some notable divergences especially for base
 406 models that scored lower on alignment than one
 407 might expect based on their BBH score. Figure
 408 4B shows the correlations between alignment
 409 and all the established benchmarks. While
 410 all three alignment metrics were in agreement,
 411 we found that the overall correlations varied
 412 ($r_{min} = 0.36, p < 0.05$; $r_{max} = 0.74, p <$
 413 0.001). We focus on Procrustes R^2 moving
 414 forward.

415 Performance on benchmarks that probe
 416 probe domain-general knowledge (BBH;
 417 $r = 0.74, p < 0.001$; MMLU;
 418 $r = 0.71, p < 0.001$) and instruction-following
 419 ability (IFEval; $r = 0.68, p < 0.001$) were
 420 more correlated with alignment than benchmarks testing more specific domain knowledge and
 421 reasoning ability (Math; $r = 0.61, p < 0.001$, MUSR; $r = 0.36, p < 0.05$).

422 This finding is sensible given that alignment with human semantic judgments requires representing
 423 semantic *knowledge* broadly in human-like ways (captured by BBH and MMLU) and *deploying* that
 424 knowledge in context-sensitive ways according to instructions (captured by IFEval). We also
 425 found that the correlation between BBH and IFEval was significantly weaker than the correlation
 426 between these evaluations and our measure of representational alignment ($r_{BBH,IFEval} = 0.44, p <$
 427 0.01), suggesting our alignment measure captures separate competencies unique to each benchmark.

428 5 CONCLUSION

429
 430 The rapid advancement of language models (both in scale and diversity) and the varied, often
 431 unpredictable, performance of models on standard benchmarks has left open the question “what

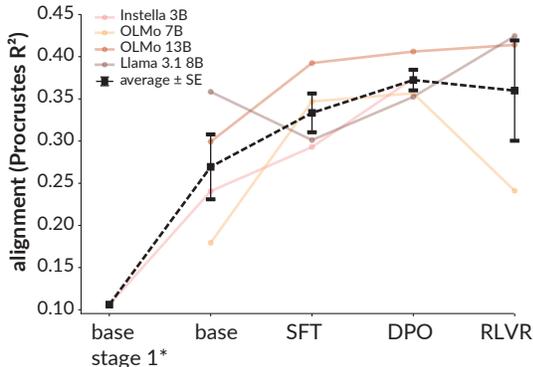


Figure 3: **Effects of post-training on alignment.** The black points show mean alignment (across models) at each post-training stage.*Instella only.

432 does it take to make human-like models”? A gold standard for ‘human-like’ would be to have
 433 models that not only perform tasks in human-like ways but also have internal representations and
 434 computations that are similar to that of humans. Here we attempt to identify the computational
 435 ingredients most predictive of human-model representational alignment.

436 We first leveraged a large-scale dataset of human triadic similarity judgments and associated
 437 semantic embeddings from the THINGS dataset and set them as the target for measuring alignment.
 438 Next, we carefully constructed a suite of 77 open-weight language models that varied in several
 439 key computational ingredients including model scale, training dataset size, dimensionality of model
 440 embeddings, and post-training regimes, among others. We estimated semantic embeddings for
 441 each model using a triadic similarity judgment task analogous to the human task and measured
 442 representational alignment between model and human embeddings as function of different model
 443 characteristics. We arrived at several key insights.

444 First, we found two convergent pieces
 445 of evidence that modern post-training
 446 techniques are central to human-model
 447 alignment: (1) instruction-finetuning was
 448 the strongest predictor of alignment even
 449 after controlling for variance explained
 450 by other ingredients (Figure 2 and
 451 (2) alignment tended to increase across
 452 different stages of post-training (Figure
 453 3). The importance of context length
 454 in predicting alignment further suggests
 455 that the methods that enable long context-
 456 reasoning also bring model representa-
 457 tions into closer alignment with humans.
 458 Second, we observed that architectural
 459 dimensionality (embedding layers, MLP)
 460 were relatively more important than
 461 overall model size and training corpus
 462 size in predicting alignment, especially
 463 in the mixed-effects model suggesting
 464 that representational capacity is key to
 465 developing models that align with human
 466 representations. Third, we found that
 467 while some existing benchmarks were
 468 correlated with our measure of alignment
 469 (e.g., *MMLU*, *BBH*), no single benchmark
 470 captured all the variance in representa-
 471 tional alignment. Notably, benchmarks
 472 emphasizing broad semantic knowledge
 473 and its flexible deployment showed the
 474 strongest correlations, consistent with the-
 475 ories of how humans represent semantic
 476 knowledge (Rogers & McClelland, 2004; Saxe et al., 2019).

475 **Limitations** To our knowledge, this is the first attempt to isolate the computational ingredients
 476 predictive of human–LLM conceptual alignment in language models, however several limitations
 477 remain that future work could address. First, while we evaluated many core computational
 478 ingredients, we could not assess the role of training dataset composition (e.g., text–code balance,
 479 language coverage) due to limited documentation for many models. Future work should look beyond
 480 token counts to clarify which dataset properties enable semantic alignment. Second, some classes
 481 of models were sparsely represented in our suite including SSMs, multimodal models, and mixture-
 482 of-expert models limiting the reliability of insights we can derive regarding ingredients that support
 483 them. Third, while similarity judgements are a well validated method and scale well to large concept
 484 sets, there may be other ways of probing semantic knowledge such as by observing how models
 485 deploy this knowledge in open-ended reasoning tasks (Wong et al., 2025; Mahowald et al., 2024).
 Future work can seek to study such tasks that tackle the *deployment* of semantic knowledge.

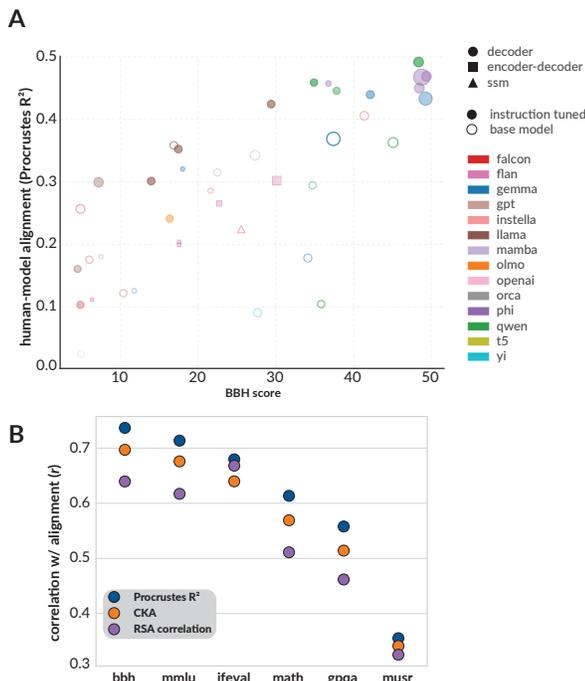


Figure 4: **Relationship between alignment and other LLM benchmarks.** **A.** BigBenchHard scores for each model in our suite vs. their Procrustes R^2 value w.r.t. human embeddings. **B.** Correlation between alignment and the six LLM benchmarks evaluated.

REFERENCES

- 486
487
488 01.AI Team. Yi: Open foundation models. arXiv:2403.04652, March 2024.
- 489 Yamini Bansal, Preetum Nakkiran, and Boaz Ravikumar. Revisiting model stitching to compare
490 neural representations. *Advances in neural information processing systems*, 34:225–236, 2021.
- 491
492 Serguei Barannikov, Ilya Trofimov, Nikita Balabin, and Evgeny Burnaev. Representation
493 topology divergence: A method for comparing neural network representations. *arXiv preprint*
494 *arXiv:2201.00058*, 2021.
- 495 Joao Barbosa, Amin Nejatbakhsh, Lyndon Duong, Sarah E Harvey, Scott L Brincat, Markus Siegel,
496 Earl K Miller, and Alex H Williams. Quantifying differences in neural population activity with
497 shape metrics. *bioRxiv*, pp. 2025–01, 2025.
- 498 Dale J Barr, Roger Levy, Christoph Scheepers, and Harry J Tily. Random effects structure for
499 confirmatory hypothesis testing: Keep it maximal. *Journal of memory and language*, 68(3):255–
500 278, 2013.
- 501
502 Alexander Berardino, Valero Laparra, Johannes Ballé, and Eero Simoncelli. Eigen-distortions of
503 hierarchical representations. *Advances in neural information processing systems*, 30, 2017.
- 504 S. Black, S. Biderman, et al. Gpt-neox-20b: An open-source autoregressive language model. In
505 *Proceedings of Machine Learning and Systems*, April 2022. [Online]. Available: [https://](https://aclanthology.org/2022.mlsys-1.72)
506 aclanthology.org/2022.mlsys-1.72.
- 507
508 Rishi Bommasani, Kelly Davis, and Claire Cardie. Interpreting pretrained contextualized
509 representations via reductions to static embeddings. *Proceedings of the 58th Annual Meeting*
510 *of the Association for Computational Linguistics*, pp. 4758–4781, 2020.
- 511 Samuel R. Bowman and George Dahl. What will it take to fix benchmarking in natural language
512 understanding? In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-
513 Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (eds.),
514 *Proceedings of the 2021 Conference of the North American Chapter of the Association for*
515 *Computational Linguistics: Human Language Technologies*, pp. 4843–4855, Online, June 2021.
516 Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.385. URL
517 <https://aclanthology.org/2021.naacl-main.385/>.
- 518 Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin.
519 Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural*
520 *information processing systems*, 33:9912–9924, 2020.
- 521 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared
522 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
523 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 524
525 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
526 contrastive learning of visual representations. *International conference on machine learning*, pp.
527 1597–1607, 2020.
- 528 Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade
529 Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for
530 contrastive language-image learning. In *Proceedings of the IEEE/CVF conference on computer*
531 *vision and pattern recognition*, pp. 2818–2829, 2023.
- 532 H. W. Chung, L. Hou, et al. Scaling instruction-finetuned language models. arXiv:2210.11416,
533 December 2022. [Online]. Available: <https://arxiv.org/abs/2210.11416>.
- 534
535 Radoslaw M Cichy, Nikolaus Kriegeskorte, Kamila M Jozwik, Jasper JF van den Bosch, and Ian
536 Charest. The spatiotemporal neural dynamics underlying perceived similarity for real-world
537 objects. *NeuroImage*, 194:12–24, 2019.
- 538 Katherine M Collins, Ilia Sucholutsky, Umang Bhatt, Kartik Chandra, Lionel Wong, Mina Lee,
539 Cedegao E Zhang, Tan Zhi-Xuan, Mark Ho, Vikash Mansinghka, et al. Building machines that
learn and think with people. *Nature human behaviour*, 8(10):1851–1863, 2024.

- 540 Y Ivette Colon and Timothy T Rogers. Yours and ours: Individual and group differences in semantic
541 organization from triplet judgements of faces. In *Proceedings of the annual meeting of the*
542 *cognitive science society*, volume 45, 2023.
- 543
- 544 T. Dao, A. Gu, M. Mazeika, et al. Mamba: Linear-time sequence modeling with selective state
545 spaces. arXiv:2312.00752, December 2023. [Online]. Available: <https://arxiv.org/abs/2312.00752>.
- 546
- 547 Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer,
548 et al. Scaling vision transformers to 22 billion parameters. *arXiv preprint arXiv:2302.05442*,
549 2023.
- 550
- 551 EleutherAI. Gpt-j-6b: 6 billion parameter open source gpt model. EleutherAI Blog, June 2021.
552 [Online]. Available: <https://6b.eleuther.ai>.
- 553
- 554 Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Brandon Tran, and Aleksander
555 Madry. Adversarial robustness as a prior for learned representations. *arXiv preprint*
556 *arXiv:1906.00945*, 2019.
- 557
- 558 Falcon-LLM Team. The falcon 3 family of open models. Published Dec 2024. [Online]. Available:
559 <https://huggingface.co/blog/falcon3>, December 2024.
- 560
- 561 Thomas Fel, Ivan Felipe, Drew Linsley, and Thomas Serre. Harmonizing the object recognition
562 strategies of deep neural networks with humans. *Advances in Neural Information Processing*
563 *Systems*, 35:9432–9446, 2022.
- 564
- 565 Chelsea Finn, Ian Goodfellow, and Sergey Levine. Unsupervised learning for physical interaction
566 through video prediction. *Advances in neural information processing systems*, 29, 2016.
- 567
- 568 Chaz Firestone. Performance vs. competence in human–machine comparisons. *Proceedings of the*
569 *National Academy of Sciences*, 117(43):26562–26571, 2020.
- 570
- 571 Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and
572 Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias
573 improves accuracy and robustness. *International Conference on Learning Representations*, 2019.
- 574
- 575 Robert Geirhos, Kantharaju Narayanappa, Benjamin Mitzkus, Tizian Thieringer, Matthias Bethge,
576 Felix A Wichmann, and Wieland Brendel. Partial success in closing the gap between human and
577 machine vision. *Advances in Neural Information Processing Systems*, 34:23885–23899, 2021.
- 578
- 579 Tyler Giallanza, Declan Campbell, Jonathan D Cohen, and Timothy T Rogers. An integrated model
580 of semantics and control. *Psychological Review*, 2024.
- 581
- 582 John C Gower. Generalized procrustes analysis. *Psychometrika*, 40(1):33–51, 1975.
- 583
- 584 Gabriel Grand, Idan Asher Blank, Francisco Pereira, and Evelina Fedorenko. Semantic projection
585 recovers rich human knowledge of multiple object features from word embeddings. *Nature*
586 *Human Behaviour*, 6(7):975–987, 2022.
- 587
- 588 D. Groeneveld, I. Beltagy, et al. Olmo: Accelerating the science of language models.
589 arXiv:2402.00838, June 2024.
- 590
- 591 S. Gunasekar, Y. Zhang, J. Aneja, et al. Textbooks are all you need. arXiv:2306.11644, October
592 2023.
- 593
- 594 Wes Gurnee and Max Tegmark. Language models represent space and time. *arXiv preprint*
595 *arXiv:2310.02207*, 2024.
- 596
- 597 Sishuo Han, Wenshan Cheng, Zhiruo Liu, Kaixin Shi, Jing Zhou, Xiang Ren, and Yizhou
598 Wang. MUSR: A multi-step unsupervised scientific reasoning benchmark. *arXiv preprint*
599 *arXiv:2405.08726*, 2024.

- 594 Martin N Hebart, Adam H Dickter, Alexis Kidder, Wan Y Kwok, Anna Corriveau, Caitlin
595 Van Wicklin, and Chris I Baker. Things: A database of 1,854 object concepts and more than
596 26,000 naturalistic object images. *PLoS ONE*, 14(10):e0223792, 2019. doi: 10.1371/journal.
597 pone.0223792.
- 598
599 Martin N Hebart, Oliver Contier, Lina Teichmann, Adam H Rockter, Charles Y Zheng, Alexis
600 Kidder, Anna Corriveau, Maryam Vaziri-Pashkam, and Chris I Baker. Things-data, a multimodal
601 collection of large-scale datasets for investigating object representations in human brain and
602 behavior. *eLife*, 12:e82580, 2023. doi: 10.7554/eLife.82580.
- 603 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
604 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*
605 *arXiv:2009.03300*, 2020.
- 606
607 Katherine Hermann, Ting Chen, and Simon Kornblith. The origins and prevalence of texture bias in
608 convolutional neural networks. *Advances in Neural Information Processing Systems*, 33:19000–
609 19015, 2020.
- 610 Jack Hessel, Ana Marasović, Jena D Hwang, Lillian Lee, Jeff Da, Rowan Zellers, Robert Mankoff,
611 and Yejin Choi. Do androids laugh at electric sheep? humor” understanding” benchmarks from
612 the new yorker caption contest. *arXiv preprint arXiv:2209.06293*, 2022.
- 613
614 Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. The platonic representation
615 hypothesis. *arXiv preprint arXiv:2405.07987*, 2024.
- 616
617 Rebecca L Jackson, Matthew A Lambon Ralph, and Timothy T Rogers. Late maturation of executive
618 control promotes conceptual development. *bioRxiv*, pp. 2022–03, 2022.
- 619
620 Lalit Jain, Kevin Jamieson, and Robert Nowak. Finite sample prediction and recovery bounds for
621 ordinal embedding, 2016. URL <https://arxiv.org/abs/1606.07081>.
- 622
623 Kevin G Jamieson, Lalit Jain, Chris Fernandez, Nicholas J Glattard, and Rob Nowak. Next: A
624 system for real-world development, evaluation, and application of active learning. *Advances in*
neural information processing systems, 28, 2015.
- 625
626 J. Ji, R. Kumar, et al. Gemma: Open models based on gemini
627 research and technology. Google Developers Blog, August 2024.
628 [Online]. Available: [https://developers.googleblog.com/en/
gemma-explained-overview-gemma-model-family-architectures/](https://developers.googleblog.com/en/gemma-explained-overview-gemma-model-family-architectures/).
- 629
630 Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan
631 Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning
632 with noisy text supervision. *International Conference on Machine Learning*, pp. 4904–4916,
633 2021.
- 634
635 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
636 Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint*
arXiv:2310.06770, 2023.
- 637
638 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,
639 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
640 models. *arXiv preprint arXiv:2001.08361*, 2020.
- 641
642 Charles Kemp and Joshua B Tenenbaum. The discovery of structural form. *Proceedings of the*
National Academy of Sciences, 105(31):10687–10692, 2008.
- 643
644 Tae Kyun Kim. T test as a parametric statistic. *Korean journal of anesthesiology*, 68(6):540–546,
645 2015.
- 646
647 Marcie L King, Iris IA Groen, Adam Steel, Dwight J Kravitz, and Chris I Baker. Similarity
judgments and cortical visual responses reflect different properties of object and scene categories
in naturalistic images. *NeuroImage*, 197:368–382, 2019.

- 648 Matthäus Kleindessner and Ulrike von Luxburg. Uniqueness of ordinal embedding. *Conference on*
649 *Learning Theory*, pp. 40–67, 2014.
- 650 Talia Konkle and George A Alvarez. A self-supervised domain-general learning framework for
651 human ventral stream representation. *Nature communications*, 13(1):491, 2022.
- 652 Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural
653 network representations revisited. In *International Conference on Machine Learning*, pp. 3519–
654 3529. PMLR, 2019.
- 655 Nikolaus Kriegeskorte and Marieke Mur. Inverse mds: Inferring dissimilarity structure from
656 multiple item arrangements. *Frontiers in Psychology*, 3:245, 2012.
- 657 Nikolaus Kriegeskorte, Marieke Mur, and Peter Bandettini. Representational similarity analysis–
658 connecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, 2:4, 2008.
- 659 Manoj Kumar, Neil Houlsby, Nal Kalchbrenner, and Ekin D Cubuk. Do better imagenet classifiers
660 assess perceptual similarity better? *Transactions on Machine Learning Research*, 2022.
- 661 Mary J Lindstrom and Douglas M Bates. Newton–raphson and em algorithms for linear mixed-
662 effects models for repeated-measures data. *Journal of the American Statistical Association*, 83
663 (404):1014–1022, 1988.
- 664 Mary J Lindstrom and Douglas M Bates. Nonlinear mixed effects models for repeated measures
665 data. *Biometrics*, pp. 673–687, 1990.
- 666 J. Liu, J. Wu, et al. Introducing instella: New state-of-the-art fully open 3b language models.
667 AMD ROCm Blog, March 2025. [Online]. Available: [https://rocm.blogs.amd.com/
668 artificial-intelligence/introducing-instella-3B/](https://rocm.blogs.amd.com/artificial-intelligence/introducing-instella-3B/).
- 669 Kyle Mahowald, Anna A. Ivanova, Idan A. Blank, Nancy Kanwisher, Joshua B. Tenenbaum, and
670 Evelina Fedorenko. Dissociating language and thought in large language models, 2024. URL
671 <https://arxiv.org/abs/2301.06627>.
- 672 Raja Marjieh, Ilia Sucholutsky, Theodore R Sumers, Nori Jacoby, and Thomas L Griffiths.
673 Predicting human similarity judgments using large language models. *arXiv preprint*
674 *arXiv:2202.04728*, 2022.
- 675 Raja Marjieh, Ilia Sucholutsky, Theodore R Sumers, Harin Lee, Thomas L Griffiths, and Nori
676 Jacoby. Do large language models predict human similarity judgments? *Cognitive Science*,
677 48(8):e13509, 2024a.
- 678 Raja Marjieh, Ilia Sucholutsky, Pol van Rijn, Nori Jacoby, and Thomas L Griffiths. Predicting human
679 similarity judgments using large language models. *arXiv preprint arXiv:2402.10910*, 2024b.
- 680 Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A
681 benchmark for question answering about charts with visual and logical reasoning. In *Findings*
682 *of the Association for Computational Linguistics: ACL 2022*, pp. 2263–2279, 2022.
- 683 Johannes Mehrer, Courtney J Spoerer, Emer C Jones, Nikolaus Kriegeskorte, and Tim C Kietzmann.
684 An ecologically motivated image dataset for deep learning yields better models of human vision.
685 *Proceedings of the National Academy of Sciences*, 118(8):e2011417118, 2021.
- 686 Meta AI. Llama 3.1 8b instruct (2024 release). Internal technical report, Meta, 2024.
- 687 Nitesh Methani, Pritha Ganguly, Mitesh M Khapra, and Pratyush Kumar. Plotqa: Reasoning over
688 scientific plots. In *Proceedings of the IEEE/CVF winter conference on applications of computer*
689 *vision*, pp. 1527–1536, 2020.
- 690 Luca Moschella, Valentino Maiorca, Marco Fumero, Antonio Norelli, Francesco Locatello, and
691 Emanuele Rodolà. Relative representations enable zero-shot latent space communication. *arXiv*
692 *preprint arXiv:2209.15430*, 2023.
- 693 Kushin Mukherjee and Timothy T Rogers. Using drawings and deep neural networks to characterize
694 the building blocks of human visual similarity. *Memory & Cognition*, 53(1):219–241, 2025.

- 702 Kushin Mukherjee, Holly Huey, Martin N Hebart, and Wilma A Bainbridge¹⁰. Drawings of things:
703 A large-scale drawing dataset of 1,854 object concepts. *PsyArXiv*.
704
- 705 Kushin Mukherjee, Karen Schloss, Laurent Lessard, Michael Gleicher, and Timothy Rogers. Color-
706 concept associations reveal an abstract conceptual space. *Journal of Vision*, 22(14):4408–4408,
707 2022.
- 708 Kushin Mukherjee, Holly Huey, Xuanchen Lu, Yael Vinker, Rio Aguina-Kang, Ariel Shamir, and
709 Judith Fan. Seva: Leveraging sketches to evaluate alignment between human and machine visual
710 abstraction. *Advances in Neural Information Processing Systems*, 36:67138–67155, 2023a.
- 711 Kushin Mukherjee, Siddharth Suresh, and Timothy T. Rogers. Human-machine cooperation for
712 semantic feature listing, 2023b. URL <https://arxiv.org/abs/2304.05012>. ICLR
713 2023 Tiny Papers.
714
- 715 Kushin Mukherjee, Timothy T. Rogers, and Karen B. Schloss. Large language models estimate fine-
716 grained human color-concept associations, 2024. URL [https://arxiv.org/abs/2406.](https://arxiv.org/abs/2406.17781)
717 [17781](https://arxiv.org/abs/2406.17781).
- 718 Kushin Mukherjee, Donghao Ren, Dominik Moritz, and Yannick Assogba. Encqa: Benchmarking
719 vision-language models on visual encodings for charts. *arXiv preprint arXiv:2508.04650*, 2025.
720
- 721 S. Mukherjee, A. Mitra, G. Jawahar, et al. Orca: Progressive learning from complex explanation
722 traces of gpt-4. *arXiv:2306.02707*, June 2023c.
- 723 Marieke Mur, Mirjam Meys, Jerzy Bodurka, Rainer Goebel, Peter A Bandettini, and Nikolaus
724 Kriegeskorte. Human object-similarity judgments reflect and transcend the primate-it object
725 representation. *Frontiers in psychology*, 4:128, 2013.
726
- 727 Lukas Muttenthaler, Jonas Dippel, Lorenz Linhardt, Robert A Vandermeulen, and Simon Kornblith.
728 Human alignment of neural network representations. *arXiv preprint arXiv:2211.01201*, 2022a.
- 729 Lukas Muttenthaler, Charles Y Zheng, Patrick McClure, Robert A Vandermeulen, Martin N Hebart,
730 and Francisco Pereira. Vice: Variational interpretable concept embeddings. *Advances in Neural*
731 *Information Processing Systems*, 35:33661–33675, 2022b.
- 732 Lukas Muttenthaler, Klaus Greff, Frieda Born, Bernhard Spitzer, Simon Kornblith, Michael C
733 Mozer, Klaus-Robert MÅžller, Thomas Unterthiner, and Andrew K Lampinen. Aligning machine
734 and human visual representations across abstraction levels. *arXiv preprint arXiv:2409.06509*,
735 2024.
736
- 737 Joosung Nam, Juyeon Lee, Dongkeun Chung, Juyoung Song, Sangwoo Yoon, Sungjin Jung,
738 Seungjong Min, and Yongjin Choi. Idealignbench: Measuring and improving value alignment
739 of large language models. *arXiv preprint arXiv:2412.03924*, 2024.
- 740 Hamed Nili, Cai Wingfield, Alexander Walther, Li Su, William Marslen-Wilson, and Nikolaus
741 Kriegeskorte. A toolbox for representational similarity analysis. *PLoS computational biology*,
742 10(4):e1003553, 2014.
743
- 744 Antonio Norelli, Marco Fumero, Valentino Maiorca, Luca Moschella, Emanuele Rodolà, and
745 Francesco Locatello. Asif: coupled data turns unimodal models to multimodal without training.
746 *arXiv preprint arXiv:2210.01738*, 2023.
- 747 Robert M Nosofsky. Choice, similarity, and the context theory of classification. *Journal of*
748 *Experimental Psychology: Learning, memory, and cognition*, 10(1):104, 1984.
749
- 750 Urvashi Oswal, Christopher Cox, Matthew Lambon-Ralph, Timothy Rogers, and Robert Nowak.
751 Representational similarity learning with application to brain networks. In *International*
752 *Conference on Machine Learning*, pp. 1041–1049. PMLR, 2016.
- 753 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
754 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
755 instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:
27730–27744, 2022.

- 756 Karalyn Patterson, Peter J Nestor, and Timothy T Rogers. Where do you know what you know? the
757 representation of semantic knowledge in the human brain. *Nature reviews neuroscience*, 8(12):
758 976–987, 2007.
- 759
760 Joshua C Peterson, Joshua T Abbott, and Thomas L Griffiths. Evaluating (and improving) the
761 correspondence between deep neural networks and human representations. *Cognitive science*, 42
762 (8):2648–2669, 2018.
- 763 Hieu Pham, Zihang Dai, Golnaz Ghiasi, Hanxiao Liu, Adams Wei Yu, Minh-Thang Luong,
764 Mingxing Tan, and Quoc V Le. Combined scaling for zero-shot transfer learning.
765 *Neurocomputing*, 555:126658, 2022.
- 766
767 Yulu Qin, Dheeraj Varghese, Adam Dahlgren Lindström, Lucia Donatelli, Kanishka Misra, and
768 Najoung Kim. Vision-and-language training helps deploy taxonomic knowledge but does not
769 fundamentally alter it. *arXiv preprint arXiv:2507.13328*, 2025.
- 770 Qwen Team. Qwen technical report (tongyi qianwen). Alibaba Cloud, August 2023. [Online].
771 Available: <https://github.com/QwenLM>.
- 772
773 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
774 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
775 models from natural language supervision. *International conference on machine learning*, pp.
776 8748–8763, 2021.
- 777 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea
778 Finn. Direct preference optimization: Your language model is secretly a reward model. In
779 *Advances in Neural Information Processing Systems*, volume 36, pp. 33499–33530, 2023.
- 780
781 C. Raffel, N. Shazeer, A. Roberts, et al. Exploring the limits of transfer learning with a unified
782 text-to-text transformer. *Journal of Machine Learning Research*, 21(140), 2020.
- 783 Vedant Rathi, Abhishek Kumar, et al. Can language models learn to skip steps? *arXiv preprint*
784 *arXiv:2411.01855*, 2024.
- 785
786 David Rein, Adarsh Raichur, Andre Csiszar, Ruben Gonzalez, Bomee Yuan, Vincent Zhong, Yi Tay,
787 Surbhi Narang, and Dani Yogatama. GPQA: A graduate-level google-proof q&a benchmark.
788 *arXiv preprint arXiv:2311.12022*, 2023.
- 789 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien
790 Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a
791 benchmark. In *First Conference on Language Modeling*, 2024.
- 792
793 Timothy T Rogers and James L McClelland. *Semantic cognition: A parallel distributed processing*
794 *approach*. MIT press, 2004.
- 795 Eleanor Rosch. Cognitive representations of semantic categories. *Journal of experimental*
796 *psychology: General*, 104(3):192, 1975.
- 797
798 David E Rumelhart, James L McClelland, PDP Research Group, et al. *Parallel distributed*
799 *processing, volume 1: Explorations in the microstructure of cognition: Foundations*. The MIT
800 press, 1986.
- 801 Andrew M Saxe, James L McClelland, and Surya Ganguli. A mathematical theory of semantic
802 development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116
803 (23):11537–11546, 2019.
- 804
805 Peter H Schönemann. A generalized solution of the orthogonal procrustes problem. *Psychometrika*,
806 31(1):1–10, 1966.
- 807
808 Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- 809
809 Roger N Shepard. Multidimensional scaling, tree-fitting, and clustering. *Science*, 210(4468):390–
398, 1980.

- 810 Scott Sievert, Robert Nowak, and Timothy T Rogers. Efficiently learning relative similarity
811 embeddings with crowdsourcing. *Journal of open source software*, 8(84), 2023.
812
- 813 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam
814 Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska,
815 Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W.
816 Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain,
817 Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders
818 Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La,
819 Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta,
820 Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum,
821 Arul Menezes, Arun Kirubakaran, Asher Mullokandov, Ashish Sabharwal, Austin Herrick,
822 Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph,
823 Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin
824 Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron
825 Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh,
826 Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites,
827 Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera,
828 Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette,
829 Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy,
830 Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito,
831 Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis
832 Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra,
833 Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina
834 Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth
835 Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie
836 Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii
837 Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé,
838 Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán
839 Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-
840 López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh
841 Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakura,
842 Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson
843 Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel,
844 James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema
845 Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova,
846 Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng
847 Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis,
848 Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph
849 Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua,
850 Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja
851 Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omondi, Kory Mathewson, Kristen Chiafullo,
852 Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo
853 Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency,
854 Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros Colón,
855 Luke Metz, Lütfi Kerem Şenel, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Maheen
856 Farooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria
857 Jose Ramírez Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast,
858 Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Melody
859 Arnaud, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy,
860 Michael Starritt, Michael Strube, Michał Swedrowski, Michele Bevilacqua, Michihiro Yasunaga,
861 Mihir Kale, Mike Cain, Mimeo Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal,
862 Moin Aminnaseri, Mor Geva, Mozhdah Gheini, Mukund Varma T, Nanyun Peng, Nathan A.
863 Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron,
864 Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar,
865 Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar
866 Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale
867 Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter
868 Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya

- 864 Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta
865 Rudolph, Raefer Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg,
866 Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan
867 Sikand, Roman Novak, Roman Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang,
868 Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib
869 Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel
870 Gruetter, Samuel R. Bowman, Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A.
871 Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann,
872 Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry
873 Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal,
874 Shyam Upadhyay, Shyamolima, Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi,
875 Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas
876 Dehaene, Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad,
877 Steven T. Piantadosi, Stuart M. Shieber, Summer Mishserghi, Svetlana Kiritchenko, Swaroop
878 Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo
879 Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick,
880 Timofei Kornev, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar
881 Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak,
882 Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus,
883 William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi
884 Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin
885 Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai,
886 Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. Beyond the
887 imitation game: Quantifying and extrapolating the capabilities of language models, 2023. URL
888 <https://arxiv.org/abs/2206.04615>.
- 889 Zach Studdiford, Timothy T Rogers, Siddharth Suresh, and Kushin Mukherjee. Evaluating steering
890 techniques using human similarity judgments. *arXiv preprint arXiv:2505.19333*, 2025.
- 891 Iliia Sucholutsky and Thomas L. Griffiths. Alignment with human representations supports robust
892 few-shot learning, 2023a. URL <https://arxiv.org/abs/2301.11990>.
- 893 Iliia Sucholutsky and Tom Griffiths. Alignment with human representations supports robust few-shot
894 learning. *Advances in Neural Information Processing Systems*, 36:73464–73479, 2023b.
- 895 Iliia Sucholutsky, Lukas Muttenthaler, Adrian Weller, Andi Peng, Andreea Bobu, Been Kim,
896 Bradley C Love, Erin J Achterberg, Christiane Fellbaum, Marek Grzes, et al. Getting aligned
897 on representational alignment. *arXiv preprint arXiv:2310.13018*, 2023.
- 898 Siddharth Suresh, Kushin Mukherjee, and Timothy T. Rogers. Semantic feature verification in
899 FLAN-T5, 2023a. URL <https://arxiv.org/abs/2304.05591>. ICLR 2023 Tiny
900 Papers.
- 901 Siddharth Suresh, Kushin Mukherjee, Xizheng Yu, Wei-Chun Huang, Lisa Padua, and Timothy
902 Rogers. Conceptual structure coheres in human cognition but not in large language models.
903 In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference*
904 *on Empirical Methods in Natural Language Processing*, pp. 722–738, Singapore, December
905 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.47. URL
906 <https://aclanthology.org/2023.emnlp-main.47/>.
- 907 Siddharth Suresh, Wei-Chun Huang, Kushin Mukherjee, and Timothy T. Rogers. Categories vs
908 semantic features: What shape the similarities people discern in photographs of objects? In *ICLR*
909 *2024 Workshop on Representational Alignment*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=iE5aXw3RFd)
910 [forum?id=iE5aXw3RFd](https://openreview.net/forum?id=iE5aXw3RFd).
- 911 Siddharth Suresh, Kushin Mukherjee, Tyler Giallanza, Xizheng Yu, Mia Patil, Jonathan D. Cohen,
912 and Timothy T. Rogers. Ai-enhanced semantic feature norms for 786 concepts, 2025. URL
913 <https://arxiv.org/abs/2505.10718>. Also presented at the ICLR 2025 Workshop on
914 Bidirectional Human-AI Alignment.
- 915 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,
916 Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big-bench tasks
917 and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*, 2022.

- 918 Omer Tamuz, Ce Liu, Serge Belongie, Ohad Shamir, and Adam Tauman Kalai. Adaptively learning
919 the crowd kernel. In *International Conference on Machine Learning*, pp. 673–680, 2011.
920
- 921 Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela,
922 and Candace Ross. Winoground: Probing vision and language models for visio-linguistic
923 compositionality. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
924 Recognition*, pp. 5238–5248, 2022.
- 925 Shengbang Tong, Erik Ferrara, Jie Wu, et al. Mass-producing failures of multimodal systems with
926 language models. *arXiv preprint arXiv:2306.12105*, 2024.
927
- 928 Greta Tuckute, Nancy Kanwisher, and Evelina Fedorenko. Language in brains, minds, and
929 machines. *Annual Review of Neuroscience*, 47(2024):277–301, 2024.
- 930 Amos Tversky. Features of similarity. *Psychological review*, 84(4):327, 1977.
931
- 932 Aaron R Voelker and Chris Eliasmith. Improving spiking dynamical networks: Accurate delays,
933 higher-order synapses, and time cells. *Neural computation*, 30(3):569–609, 2018.
- 934 Catherine Wah, Steve Branson, Pietro Perona, and Serge Belongie. Similarity comparisons for
935 interactive fine-grained categorization. *Proceedings of the IEEE Conference on Computer Vision
936 and Pattern Recognition*, pp. 859–866, 2014.
937
- 938 Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan Wilcox, Chengxu Zhuang, Juan Ciro, Rafael
939 Mosquera, Bhargavi Paranjape, Adina Williams, Tal Linzen, et al. Findings of the babylm
940 challenge: Sample-efficient pretraining on developmentally plausible corpora. *arXiv preprint
941 arXiv:2504.08165*, 2025.
- 942 Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Ankur Pasupat, Stella Wang, Quoc
943 Le, and Denny Zhou. Finetuned language models are zero-shot learners. *arXiv preprint
944 arXiv:2109.01652*, 2021.
945
- 946 Alex H Williams, Erin Kunz, Simon Kornblith, and Scott W Linderman. Generalized shape metrics
947 on neural representations. *Advances in Neural Information Processing Systems*, 34:4738–4750,
948 2021.
- 949 Catherine Wong, William P McCarthy, Gabriel Grand, Yoni Friedman, Joshua B Tenenbaum, Jacob
950 Andreas, Robert D Hawkins, and Judith E Fan. Identifying concept libraries from language about
951 object structure. *arXiv preprint arXiv:2205.05666*, 2022.
952
- 953 Lionel Wong, Katherine M Collins, Lance Ying, Cedegao E Zhang, Adrian Weller, Tobias
954 Gerstenberg, Timothy O’Donnell, Alexander K Lew, Jacob D Andreas, Joshua B Tenenbaum,
955 et al. Modeling open-world cognition as on-demand synthesis of probabilistic models. *arXiv
956 preprint arXiv:2507.12547*, 2025.
- 957 Daniel LK Yamins and James J DiCarlo. Using goal-driven deep learning models to understand
958 sensory cortex. *Nature neuroscience*, 19(3):356–365, 2016.
959
- 960 Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J
961 DiCarlo. Performance-optimized hierarchical models predict neural responses in higher visual
962 cortex. *Proceedings of the national academy of sciences*, 111(23):8619–8624, 2014.
- 963 Lance Ying, Katherine M Collins, Lionel Wong, Iliia Sucholutsky, Ryan Liu, Adrian Weller, Tianmin
964 Shu, Thomas L Griffiths, and Joshua B Tenenbaum. On benchmarking human-like intelligence in
965 machines. *arXiv preprint arXiv:2502.20502*, 2025.
- 966 Mert Yuksekogonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and
967 why vision-language models behave like bags-of-words, and what to do about it? *International
968 Conference on Learning Representations*, 2023.
969
- 970 Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised
971 learning via redundancy reduction. *International Conference on Machine Learning*, pp. 12310–
12320, 2021.

972 Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers.
973 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
974 12104–12113, 2022.

975
976 Bonan Zhao, Christopher G Lucas, and Neil R Bramley. A model of conceptual bootstrapping in
977 human cognition. *Nature Human Behaviour*, 8(1):125–136, 2024.

978
979 Kuan Lok Zhou, Jiayi Chen, Siddharth Suresh, Reuben Narad, Timothy T Rogers, Lalit K Jain,
980 Robert D Nowak, Bob Mankoff, and Jifan Zhang. Bridging the creativity understanding gap:
981 Small-scale human alignment enables expert-level humor ranking in llms. *arXiv preprint*
982 *arXiv:2502.20356*, 2025.

983
984 Yixin Zhou, Le Zhang, Chengyu Liu, Yinsong Yao, Yansong Feng, Yunzhe Wang, Linyi Zheng,
985 Chenghua Zhu, Jiayi Li, Binyuan Li, Libin Wang, Ruyi Chen, Jia Guo, and Lifeng Li. IFEval:
986 A new rigorous evaluation for instruction-following models. *arXiv preprint arXiv:2403.09886*,
987 2024.

988
989 Chengxu Zhuang, Siming Yan, Aran Nayebi, Martin Schrimpf, Michael C Frank, James J DiCarlo,
990 and Daniel LK Yamins. Unsupervised neural network models of the ventral visual stream.
991 *Proceedings of the National Academy of Sciences*, 118(3):e2014196118, 2021.

992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

1026 A APPENDIX

1027 A.1 ODD-ONE-OUT JUDGMENTS \rightarrow SPOSE ESTIMATION (HUMAN SEMANTIC

1028 EMBEDDINGS)

1029 To establish a ground-truth human semantic space, we employ the SPoSE (Sparse Positive Similarity
1030 Embedding) framework applied to odd-one-out behavioral data from the THINGS dataset (Hebart
1031 et al., 2023).

1032 *SPoSE Algorithm.* Let $\mathcal{X} = \{1, \dots, n\}$ index n items and let $X \in \mathbb{R}_{\geq 0}^{n \times d}$ be a nonnegative, low-
1033 dimensional embedding with rows $x_i^\top \in \mathbb{R}_{\geq 0}^d$. We define pairwise similarities using the inner
1034 product $s_{ij} = x_i^\top x_j$, yielding similarity matrix $S = XX^\top$. In an odd-one-out trial with unordered
1035 triplet $\{i, j, k\}$, participants select the *most similar pair* among $\{(i, j), (i, k), (j, k)\}$. SPoSE models
1036 this choice behavior with a multinomial logit over pairwise similarities:

$$1037 \Pr[(i, j) \text{ chosen} \mid \{i, j, k\}, X] = \frac{\exp(s_{ij})}{\exp(s_{ij}) + \exp(s_{ik}) + \exp(s_{jk})}. \quad (1)$$

1038 *Estimation and Implementation.* The SPoSE estimate \hat{X} is obtained via a sparsity-promoting MAP
1039 program that encourages both nonnegativity and interpretability:

$$1040 \hat{X} \in \arg \max_{X \geq 0} \left\{ \underbrace{\sum_{t=1}^T \log p(y_t \mid \{i_t, j_t, k_t\}, X)}_{\text{log-likelihood}} - \lambda \|X\|_1 \right\}, \quad (2)$$

1041 where y_t is the chosen pair on trial t , and $\lambda > 0$ controls sparsity.² We use the publicly released
1042 SPoSE embeddings as our human reference space.

1043 *Dimensionality Selection.* To determine the appropriate embedding dimension, we apply PCA to
1044 the released SPoSE matrix and retain the smallest d that explains at least 95% of the variance
1045 (Appendix A.7). For our corpus of $n = 128$ concepts, this yields $d = \hat{d}$ dimensions (*fill in*),
1046 which we use as the target dimensionality for model embeddings.

1047 A.2 ANCHORED SIMILARITY JUDGMENTS \rightarrow ORDINAL EMBEDDING (MODEL SEMANTIC

1048 EMBEDDINGS)

1049 To obtain comparable model-based semantic embeddings, we apply an ordinal embedding algorithm
1050 (Tamuz et al., 2011; Sievert et al., 2023) to model similarity judgments, creating semantic spaces
1051 where frequently co-judged similar items are positioned closer together.

1052 *Embedding Algorithm.* Each triplet $\{i, j, k\}$ from model judgments yields an ordinal constraint: “ i
1053 is closer to j than to k .” Let $\{x_i^*\}_{i=1}^n \subset \mathbb{R}^d$ denote the unknown true point locations we seek to
1054 estimate, with corresponding (squared) Euclidean distance matrix D^* . We model the probability of
1055 observing each ordinal constraint using a standard noisy triplet model with monotone link function:

$$1056 \Pr[y_{(i,j,k)} = 1] = f(D_{ik}^* - D_{ij}^*), \quad f(0) = \frac{1}{2}. \quad (3)$$

1057 *Estimation and Implementation.* The ordinal embedding algorithm minimizes crowd-kernel triplet
1058 loss (Tamuz et al., 2011) by optimizing Euclidean distances between item pairs. We minimize an
1059 empirical surrogate of the negative log-likelihood using a low-rank Gram matrix parameterization.
1060 To ensure reliable estimates, we reserve 20% of our triplet data for validation purposes and monitor
1061 crowd-kernel loss throughout the fitting procedure.

1062 *Sample Complexity Considerations.* When the true embedding has rank d and triplet judgments
1063 are sampled approximately uniformly at random, finite-sample theory provides the out-of-sample
1064 prediction error bound:

$$1065 \mathcal{E}_{\text{pred}} = \tilde{O}\left(\sqrt{\frac{dn \log n}{|S|}}\right), \quad \Rightarrow \quad |S| = \tilde{\Theta}(dn \log n). \quad (4)$$

1066 ²The nonnegativity constraint and elementwise ℓ_1 penalty are equivalent to an exponential prior on factors
1067 and empirically yield interpretable semantic dimensions.

Thus, $\tilde{\Theta}(nd \log n)$ triplet judgments suffice for accurate prediction of new comparisons, and at least $\Omega(dn \log n)$ ordinal comparisons are information-theoretically necessary (Jain et al., 2016). Guided by this theory and setting the model embedding dimension to match the human-derived \hat{d} , we collect

$$N_{\text{triplets}} = c d n \log n \quad (5)$$

triplet judgments per model, where c is chosen to balance statistical efficiency with computational constraints.

This parallel approach provides a principled framework for human-model semantic comparison. Both methods yield d -dimensional embeddings where geometric relationships reflect semantic similarities, with the SPoSE framework extracting interpretable human semantic structure and ordinal embedding converting model judgments into geometrically comparable representations with theoretically grounded sample complexity $\tilde{\Theta}(dn \log n)$.

Table 1: Model Alignment Across Training Stages

Model	Training Stage	Alignment (R ²)
Instella 3B	Base Stage 1	0.106
Instella 3B	Base	0.241
Instella 3B	SFT	0.293
Instella 3B	DPO	0.374
Llama 3.1 8B	Base	0.358
Llama 3.1 8B	SFT	0.301
Llama 3.1 8B	DPO	0.352
Llama 3.1 8B	RLVR	0.424
OLMo 13B	Base	0.299
OLMo 13B	SFT	0.392
OLMo 13B	DPO	0.406
OLMo 13B	RLVR	0.414
OLMo 7B	Base	0.179
OLMo 7B	SFT	0.347
OLMo 7B	DPO	0.357
OLMo 7B	RLVR	0.241
Average	Base Stage 1	0.106
Average	Base	0.269
Average	SFT	0.333
Average	DPO	0.372
Average	RLVR	0.360

A.3 OBJECT CONCEPTS

The following table presents the 128 concrete object concepts from the THINGS dataset used in our experiments, along with their semantic categories and definitions.

Table 2: Complete list of object concepts with categories and definitions

Concept	Category	Definition
Animal		
badger	animal	sturdy carnivorous burrowing mammal with strong claws; widely distributed in the northern hemisphere
bear	animal	massive plantigrade carnivorous or omnivorous mammals with long shaggy coats and strong claws
butterfly	animal	diurnal insect typically having a slender body with knobbed antennae and broad colorful wings

Table 2 continued from previous page

	Concept	Category	Definition
1134			
1135			
1136			
1137	chipmunk	animal	a burrowing ground squirrel of western America and Asia; has cheek pouches and a light and dark stripe running down the body
1138	cow	animal	female of domestic cattle: "‘moo-cow’ is a child’s term"
1139	crow	animal	black birds having a raucous call
1140	fly	animal	two-winged insects characterized by active flight
1141	giraffe	animal	tallest living quadruped; having a spotted coat and small horns and very long neck and legs; of savannahs of tropical Africa
1142			
1143	grasshopper	animal	terrestrial plant-eating insect with hind legs adapted for leaping
1144	hyena	animal	doglike nocturnal mammal of Africa and southern Asia that feeds chiefly on carrion
1145			
1146	iguana	animal	large herbivorous tropical American arboreal lizards with a spiny crest along the back; used as human food in Central America and South America
1147			
1148	lion	animal	large gregarious predatory feline of Africa and India having a tawny coat with a shaggy mane in the male
1149			
1150	mosquito	animal	two-winged insect whose female has a long proboscis to pierce the skin and suck the blood of humans and animals
1151			
1152	mouse	animal	a hand-operated electronic device that controls the coordinates of a cursor on your computer screen as you move it around on a pad; on the bottom of the device is a ball that rolls on the surface of the pad
1153			
1154			
1155	owl	animal	nocturnal bird of prey with hawk-like beak and claws and large head with front-facing eyes
1156			
1157	parrot	animal	usually brightly colored zygodactyl tropical birds with short hooked beaks and the ability to mimic sounds
1158			
1159	puffin	animal	any of two genera of northern seabirds having short necks and brightly colored compressed bills
1160			
1161	snail	animal	freshwater or marine or terrestrial gastropod mollusk usually having an external enclosing spiral shell
1162	snake	animal	limbless scaly elongate reptile; some are venomous
1163	tarantula	animal	large hairy tropical spider with fangs that can inflict painful but not highly venomous bites
1164			
1165			
1166	Clothing		
1167	bathrobe	clothing	a loose-fitting robe of towelling; worn after a bath or swim
1168	beanie	clothing	a small skullcap; formerly worn by schoolboys and college freshmen
1169	bow	clothing	a knot with two loops and loose ends; used to tie shoelaces
1170	bracelet	clothing	jewelry worn around the wrist for decoration
1171	button	clothing	a round fastener sewn to shirts and coats etc to fit through buttonholes
1172			
1173	chaps	clothing	(usually in the plural) leather leggings without a seat; joined by a belt; often have flared outer flaps; worn over trousers by cowboys to protect their legs
1174			
1175	hat	clothing	headdress that protects the head from bad weather; has shaped crown and usually a brim
1176			
1177	jeans	clothing	(usually plural) close-fitting trousers of heavy denim for manual work or casual wear
1178			
1179	kilt	clothing	a knee-length pleated tartan skirt worn by men as part of the traditional dress in the Highlands of northern Scotland
1180			
1181	kimono	clothing	a loose robe; imitated from robes originally worn by Japanese
1182	kneepad	clothing	protective garment consisting of a pad worn by football or baseball or hockey players
1183			
1184	tuxedo	clothing	semiformal evening dress for men
1185	uniform	clothing	clothing of distinctive design worn by members of a particular group as a means of identification
1186	veil	clothing	a garment that covers the head and face
1187	visor	clothing	a brim that projects to the front to shade the eyes

Table 2 continued from previous page

	Concept	Category	Definition
1188	Container		
1189	bag	container	a flexible container with a single opening
1190	cooker	container	a utensil for cooking
1191	doily	container	a small round piece of linen placed under a dish or bowl
1192	fishbowl	container	a transparent bowl in which small fish are kept
1193	honeypot	container	South African shrub whose flowers when open are cup-shaped resembling artichokes
1194	thermos	container	vacuum flask that preserves temperature of hot or cold drinks
1195	vase	container	an open jar of glass or porcelain used as an ornament or to hold flowers
1196	Food		
1197	appetizer	food	food or drink to stimulate the appetite (usually served before a meal or as the first course)
1198	applesauce	food	puree of stewed apples usually sweetened and spiced
1199	baklava	food	rich Middle Eastern cake made of thin layers of flaky pastry filled with nuts and honey
1200	beer	food	a general name for alcoholic beverages made by fermenting a cereal (or mixture of cereals) flavored with hops
1201	bok choy	food	Asiatic plant grown for its cluster of edible white stalks with dark green leaves
1202	bread	food	food made from dough of flour or meal and usually raised with yeast or baking powder and then baked
1203	crepe	food	small very thin pancake
1204	cupcake	food	small cake baked in a muffin tin
1205	dessert	food	a dish served as the last course of a meal
1206	dough	food	a flour mixture stiff enough to knead or roll
1207	enchilada	food	tortilla with meat filling baked in tomato sauce seasoned with chili
1208	grapefruit	food	large yellow fruit with somewhat acid juicy pulp; usual serving consists of a half
1209	gravy	food	a sauce made by adding stock, flour, or other ingredients to the juice and fat that drips from cooking meats
1210	hummus	food	a thick spread made from mashed chickpeas, tahini, lemon juice and garlic; used especially as a dip for pita; originated in the Middle East
1211	leek	food	plant having a large slender white bulb and flat overlapping dark green leaves; used in cooking; believed derived from the wild <i>Allium ampeloprasum</i>
1212	mango	food	large oval tropical fruit having smooth skin, juicy aromatic pulp, and a large hairy seed
1213	margarita	food	a cocktail made of tequila and triple sec with lime and lemon juice
1214	mashed potato	food	potato that has been peeled and boiled and then mashed
1215	milkshake	food	frothy drink of milk and flavoring and sometimes fruit or ice cream
1216	oatmeal	food	porridge made of rolled oats
1217	parfait	food	layers of ice cream and syrup and whipped cream
1218	pumpkin	food	usually large pulpy deep-yellow round fruit of the squash family maturing in late summer or early autumn
1219	quesadilla	food	a tortilla that is filled with cheese and heated
1220	quiche	food	a tart filled with rich unsweetened custard; often contains other ingredients (as cheese or ham or seafood or vegetables)
1221	ravioli	food	small circular or square cases of dough with savory fillings
1222	sea urchin	food	shallow-water echinoderms having soft bodies enclosed in thin spiny globular shells
1223	souffle	food	light fluffy dish of egg yolks and stiffly beaten egg whites mixed with e.g. cheese or fish or fruit
1224	stew	food	food prepared by stewing especially meat or fish with vegetables

Table 2 continued from previous page

Concept	Category	Definition
tortellini	food	small ring-shaped stuffed pasta
waffle	food	pancake batter baked in a waffle iron
wrap	food	a sandwich in which the filling is rolled up in a soft tortilla
Furniture		
bassinet	furniture	a basket (usually hooded) used as a baby’s bed
bathtub	furniture	a relatively large open container that you fill with water and use to wash the body
beanbag	furniture	a small cloth bag filled with dried beans; thrown in games
bed	furniture	a piece of furniture that provides a place to sleep
bench	furniture	a long seat for more than one person
coaster	furniture	a covering (plate or mat) that protects the surface of a table (i.e., from the condensation on a cold glass or bottle)
computer screen	furniture	a screen used to display the output of a computer to the user
cot	furniture	a small bed that folds up for storage or transport
couch	furniture	an upholstered seat for more than one person
crib	furniture	baby bed with high sides made of slats
Other		
album	other	a book of blank pages with pockets or envelopes; for organizing photographs or stamp collections etc
backscratcher	other	a long-handled scratcher for scratching your back
banjo	other	a stringed instrument of the guitar family that has long neck and circular body
baseball	other	a ball used in playing baseball
baseball glove	other	the handwear used by fielders in playing baseball
bassoon	other	a double-reed instrument; the tenor of the oboe family
beachball	other	large and light ball; for play at the seaside
blower	other	a device that produces a current of air
bongo	other	a small drum; played with the hands
cannonball	other	a solid projectile that in former times was fired from a cannon
canvas	other	an oil painting on canvas fabric
chainsaw	other	portable power saw; teeth linked to form an endless chain
cymbal	other	a percussion instrument consisting of a concave brass disk; makes a loud crashing sound when hit with a drumstick or when two are struck together
doorknob	other	a knob used to release the catch when opening a door (often called ‘doorhandle’ in Great Britain)
doorknocker	other	a device (usually metal and ornamental) attached by a hinge to a door
extinguisher	other	a manually operated device for extinguishing small fires
fire alarm	other	a device that makes a loud sound to warn people when there is a fire
guitar	other	a stringed instrument usually having six strings; played by strumming or plucking
hatchet	other	a small ax with a short handle used with one hand (usually to chop wood)
hula hoop	other	a large hoop spun around the body by gyrating the hips, for play or exercise.
knitting needle	other	needle consisting of a slender rod with pointed ends; usually used in pairs
knob	other	a round handle
lawnmower	other	garden tool for mowing grass on lawns
pinball	other	a game played on a sloping board; the object is to propel marbles against pins or into pockets
pocket watch	other	a watch that is carried in a small watch pocket

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Table 2 continued from previous page

Concept	Category	Definition
shuffleboard	other	a game in which players use long sticks to shove wooden disks onto the scoring area marked on a smooth surface
skeleton	other	the hard structure (bones and cartilages) that provides a frame for the body of an animal
swab	other	implement consisting of a small piece of cotton that is used to apply medication or cleanse a wound or obtain a specimen of a secretion
tennis ball	other	ball about the size of a fist used in playing tennis
toilet paper	other	a soft thin absorbent paper for use in toilets
treadmill	other	an exercise device consisting of an endless belt on which a person can walk or jog without changing place
trowel	other	a small hand tool with a handle and flat metal blade; used for scooping or spreading plaster or similar materials
volleyball	other	an inflated ball used in playing volleyball
Plant		
flower gourd	plant	a plant cultivated for its blooms or blossoms
	plant	any of numerous inedible fruits with hard rinds
Vehicle		
airbag	vehicle	a safety restraint in an automobile; the bag inflates on collision and prevents the driver or passenger from being thrown forward
carriage	vehicle	a vehicle with wheels drawn by one or more horses
exhaust pipe	vehicle	a pipe through which burned gases travel from the exhaust manifold to the muffler
hot-air balloon	vehicle	balloon for travel through the air in a basket suspended below a large bag of heated air
humvee	vehicle	a high mobility, multipurpose, military vehicle with four-wheel drive
missile	vehicle	a rocket carrying a warhead of conventional or nuclear explosives; may be ballistic or directed by remote control
odometer	vehicle	a meter that shows mileage traversed
taillight	vehicle	lamp (usually red) mounted at the rear of a motor vehicle
tank	vehicle	an enclosed armored military vehicle; has a cannon and moves on caterpillar treads
taxi	vehicle	a car driven by a person whose job is to take passengers where they want to go in exchange for money

A.4 MODEL SUITE

Table 3: List of all the models used

Model		Attention Type	Training Tokens	#Heads	Hidden Dim
Falcon3-3B-Base Team (2024)	Falcon-LLM	GQA	100B	12 (4 KV)	3072
tiuae/falcon-7b Team (2024)	Falcon-LLM	MQA	1.5T	71	4544
tiuae/falcon-7b-instruct LLM Team (2024)	Falcon-LLM	MQA		71	4544
Falcon3-10B-Base Team (2024)	Falcon-LLM	GQA	2T	12 (4 KV)	3072
tiuae/falcon-11b Team (2024)	Falcon-LLM	GQA	5T	32 (8 KV)	4096

Continued on next page

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Table 3 continued from previous page

Model	Attention Type	Training To-kens	#Heads	Hidden Dim
flan-t5-small Chung et al. (2022)	Multi-Head		8	512
flan-t5-large Chung et al. (2022)	Multi-Head		16	1024
flan-t5-xl Chung et al. (2022)	Multi-Head		32	1024
flan-t5-xxl Chung et al. (2022)	Multi-Head		128	1024
gemma-3-270m Ji et al. (2024)	GQA		8 (shared KV)	2048
gemma-3-270m-it Ji et al. (2024)	GQA		8 (shared KV)	2048
gemma-2-2b Ji et al. (2024)	MQA		8 (shared KV)	2048
gemma-2-2b-it Ji et al. (2024)	MQA		8 (shared KV)	2048
gemma-2-9b Ji et al. (2024)	MHA		16	4096
gemma-2-9b-it Ji et al. (2024)	MHA		16	4096
gemma-2-27b Ji et al. (2024)	MHA		16	6144
gemma-2-27b-it Ji et al. (2024)	MHA		16	6144
gemma-3-4b-it Ji et al. (2024)	GQA		12 (4 KV)	3072
gemma-3-4b-pt Ji et al. (2024)	GQA		12 (4 KV)	3072
gemma-3-12b-it Ji et al. (2024)	GQA		16 (4 KV)	4096
gemma-3-12b-pt Ji et al. (2024)	GQA		16 (4 KV)	4096
gemma-3-27b Ji et al. (2024)	GQA		16 (4 KV)	5120
gemma-3-27b-it Ji et al. (2024)	GQA		16 (4 KV)	5120
GPT-OSS-20B Black et al. (2022)	Multi-Head	~300B	64	6144
amd/Instella-3B Liu et al. (2025)	Multi-Head	4.15T	32	2560
amd/Instella-3B-Instruct Liu et al. (2025)	Multi-Head		32	2560
Llama-3.1-8B-Instruct Meta AI (2024)	MHA		32	4096
state-spaces/mamba-1.4b-hf Dao et al. (2023)	(SSM)			
allenai/OLMo-2-1124-7B Groeneveld et al. (2024)	MHA		32	4096
allenai/OLMo-2-1124-7B-Instruct Groeneveld et al. (2024)	MHA		32	4096
allenai/OLMo-2-1124-13B-Instruct Groeneveld et al. (2024)	MHA		40	5120
OLMo-2-1124-7B-DPO Groeneveld et al. (2024)	MHA		32	4096
OLMo-2-1124-7B-SFT Groeneveld et al. (2024)	MHA		32	4096
AMD-OLMo-1B Groeneveld et al. (2024)	MHA		16	2048
AMD-OLMo-1B-SFT-DPO Groeneveld et al. (2024)	MHA		16	2048
GPT-J-6B EleutherAI (2021)	Multi-Head	~300B	16	4096
orca-2-7b Mukherjee et al. (2023c)	Multi-Head		32	4096
orca-2-13b Mukherjee et al. (2023c)	Multi-Head		40	5120
phi-3.5-mini-instruct Gunasekar et al. (2023)	Multi-Head		32	2048
phi-3-mini-128k-instruct Gunasekar et al. (2023)	Multi-Head		32	2048
phi-3-medium-128k-instruct Gunasekar et al. (2023)	Multi-Head		40	5120
Phi-3-medium-4k-instruct Gunasekar et al. (2023)	Multi-Head		40	5120

Continued on next page

Table 3 continued from previous page

Model	Attention Type	Training Tokens	To-	#Heads	Hidden Dim
phi-1.5 Gunasekar et al. (2023)	Multi-Head	7B		32	2048
qwen-14b Qwen Team (2023)	GQA	3T		40 (10 KV)	5120
qwen-14b-chat Qwen Team (2023)	GQA			40 (10 KV)	5120
qwen-2-7b Qwen Team (2023)	GQA			32 (8 KV)	4096
qwen-2-7b-instruct Qwen Team (2023)	GQA			32 (8 KV)	4096
qwen-2.5-7b Qwen Team (2023)	GQA			32 (8 KV)	4096
qwen-2.5-7b-instruct Qwen Team (2023)	GQA			32 (8 KV)	4096
qwen-2.5-14b Qwen Team (2023)	GQA			40 (10 KV)	5120
qwen-2.5-14b-instruct Qwen Team (2023)	GQA			40 (10 KV)	5120
t5-3b Raffel et al. (2020)	Multi-Head	~1T		32	1024
t5-11b Raffel et al. (2020)	Multi-Head	~1T		128	1024
yi-9B 01.AI Team (2024)	Multi-Head	3.1T		32	6144
yi-9B (coder) 01.AI Team (2024)	Multi-Head	0.8T		32	6144

A.5 DETAILS FOR SPOSE (ODD-ONE-OUT) ESTIMATION

Given $S = XX^\top$ with $X \geq 0$, the triplet likelihood is the three-way softmax in equation 1. The MAP program equation 2 is optimized by minibatch stochastic gradients; λ is tuned by held-out triplets. The released THINGS-data SPOSE embeddings are derived from large-scale odd-one-out judgments and approach the behavioral noise ceiling in predicting left-out triplets, despite far fewer than the $O(n^3)$ triplets required to fully determine a dense similarity matrix (?).

A.6 ORDINAL EMBEDDING: RISK BOUNDS AND RECOVERY

Let G^* be the centered Gram matrix of $\{x_i^*\} \subset \mathbb{R}^d$ and let $\mathcal{L}(G)$ denote the expected logistic triplet loss induced by equation 3. If G^* has rank d and $|S|$ triplets are drawn uniformly at random, then with probability $1 - \delta$,

$$\mathcal{L}(\hat{G}) - \mathcal{L}(G^*) \leq C_1 \sqrt{\frac{dn \log n}{|S|}} + C_2 \sqrt{\frac{\log(1/\delta)}{|S|}}, \quad (6)$$

for universal constants C_1, C_2 depending on the loss Lipschitz constant (?). Moreover, writing C^* for the component of the EDM orthogonal to the linear operator’s kernel, one can recover the distance structure with

$$\frac{1}{n^2} \|\hat{C} - C^*\|_F^2 = \tilde{O}\left(\frac{dn \log n}{|S|}\right), \quad (7)$$

which implies equation ?? samples suffice up to logarithmic factors; $\Omega(dn \log n)$ lower bounds are known (?).

A.7 DIMENSIONALITY SELECTION FROM HUMAN SPOSE

We compute PCA on the released SPOSE embedding matrix (items \times dimensions), retain the top d components explaining at least 95% variance, and set $d = \hat{d}$ for model-side ordinal embeddings. In our corpus ($n = 128$), this yields $\hat{d} = 29$, and therefore $N_{\text{triplets}} = c \hat{d} n \log n = 31,288$ triplets per model where $c = 4$. We increase N_{triplets} to **35,000** to provide a margin of error.

A.8 MEASURE CORRELATIONS

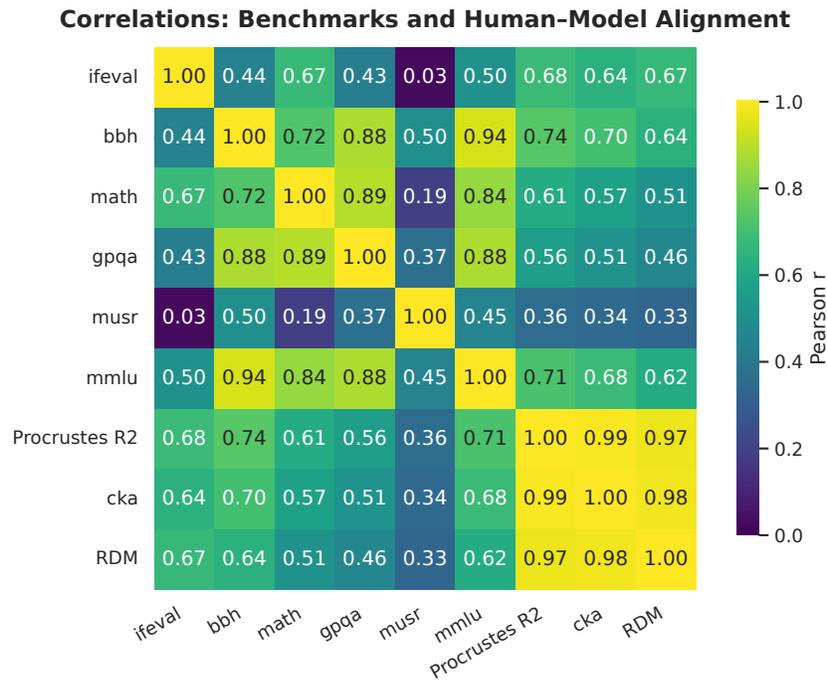


Figure 5: Correlations between measures of human-model alignment and model benchmarks, for models where benchmark scores are publicly available. (Procrustes R^2 , CKA, RSM)

A.9 RELATIONSHIPS BETWEEN COMPUTATIONAL INGREDIENT PROPERTIES AND HUMAN-MODEL ALIGNMENT

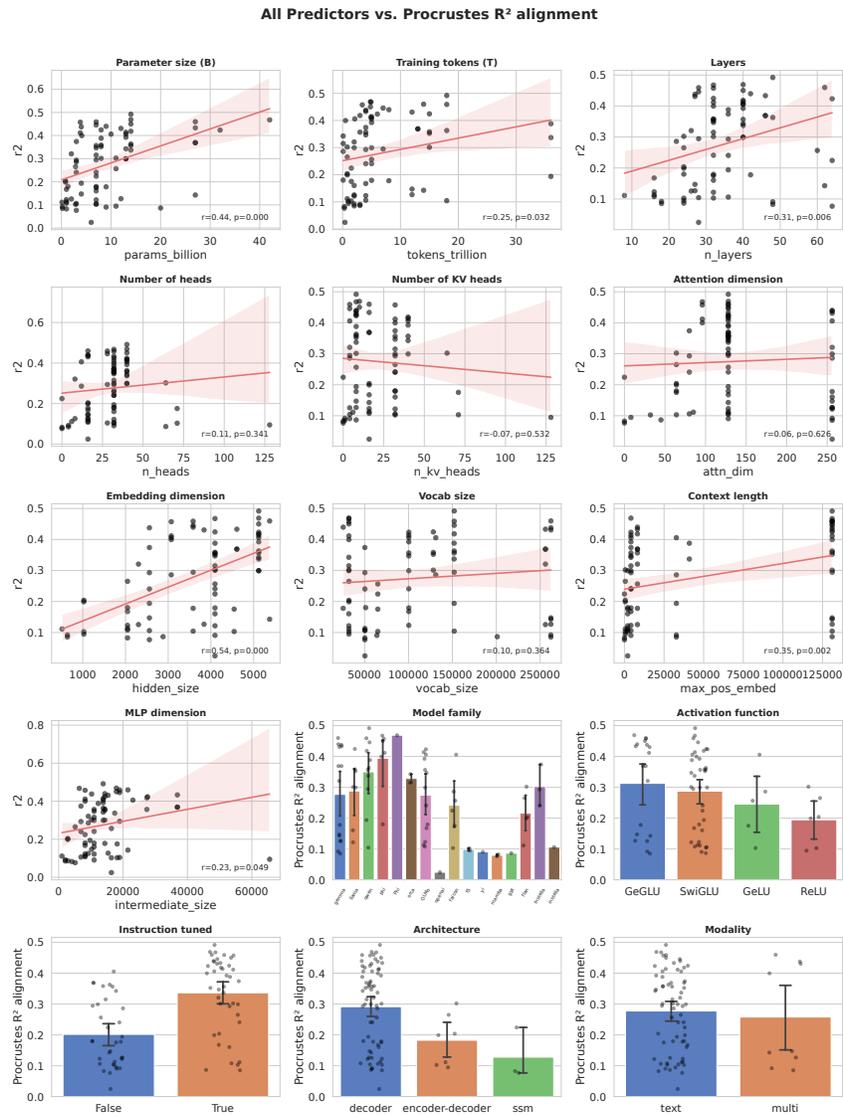


Figure 6: Categorical and continuous between all computational ingredients and human-model alignment. All significant effects in the full mixed linear-model are reported in the main text, Figure 2.

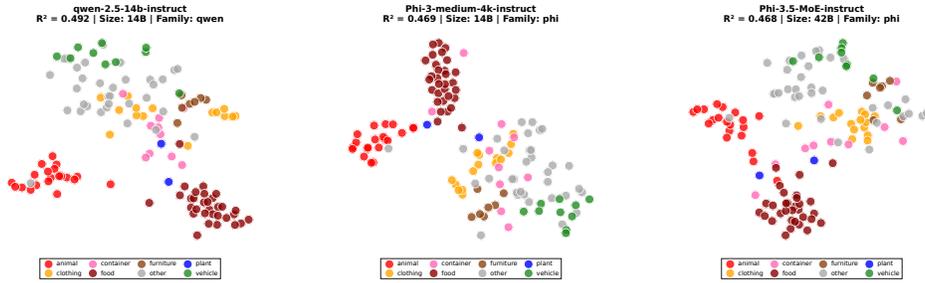
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

A.10 *t*-SNE VISUALIZATIONSTable 4: Complete ranking of models by R² score

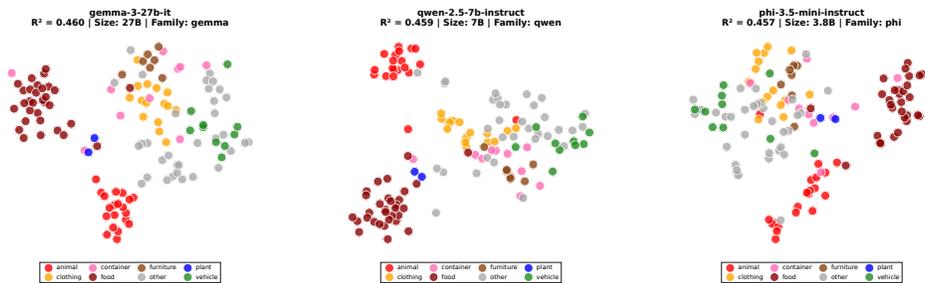
Rank	Model Name	R ² Score	Rank	Model Name	R ² Score
1	qwen25-14b-instruct	0.492	39	qwen2-7b	0.294
2	Phi-3-medium-4k-instruct	0.469	40	Instella-3B-SFT	0.293
3	Phi-3.5-MoE-instruct	0.468	41	Falcon3-3B-Base	0.286
4	gemma-3-27b-it	0.460	42	flan-t5-xl	0.265
5	qwen25-7b-instruct	0.459	43	falcon-11b	0.257
6	Phi-3.5-mini-instruct	0.457	44	OLMo-2-1124-7B-Instruct	0.241
7	phi-3-medium-128k-instruct	0.450	45	Instella-3B	0.241
8	qwen2-7b-instruct	0.446	46	Falcon3-Mamba-7B-Base	0.224
9	gemma-2-9b-it	0.440	47	flan-t5-large	0.199
10	gemma-3-4b-it	0.438	48	Qwen3-4B-Base	0.194
11	gemma-2-27b-it	0.433	49	phi-1.5	0.180
12	gemma-3-12b-it	0.431	50	OLMo-2-1124-7B	0.179
13	Llama-3.1-8B-Instruct	0.424	51	gemma-2-9b	0.178
14	OLMo-2-0325-32B-Instruct	0.423	52	falcon-7b	0.175
15	qwen-14b-chat	0.418	53	OLMo-2-0425-1B-SFT	0.168
16	OLMo-2-1124-13B-Instruct	0.414	54	llama2-7b-chat	0.160
17	phi-3-mini-128k-instruct	0.412	55	gemma-3-4b-pt	0.147
18	OLMo-2-1124-13B-DPO	0.406	56	gemma-3-27b-pt	0.143
19	Falcon3-10B-Base	0.406	57	gemma-3-12b-pt	0.127
20	Phi-3.5-vision-instruct	0.400	58	gemma-2-2b	0.126
21	OLMo-2-1124-13B-SFT	0.392	59	OLMo-2-0425-1B	0.123
22	qwen3-8b	0.388	60	llama2-7b	0.122
23	Instella-3B-Instruct	0.374	61	AMD-OLMo-1B	0.113
24	gemma-2-27b	0.369	62	flan-t5-small	0.111
25	qwen25-14b	0.363	63	AMD-OLMo-1B-SFT-DPO	0.108
26	Llama-3.1-Tulu-3-8B	0.358	64	Instella-3B-Stage1	0.106
27	qwen-14b	0.357	65	qwen25-7b	0.104
28	OLMo-2-1124-7B-DPO	0.357	66	falcon-7b-instruct	0.103
29	Llama-3.1-Tulu-3-8B-DPO	0.352	67	t5-3b	0.103
30	OLMo-2-1124-7B-SFT	0.347	68	t5-11b	0.095
31	orca-2-13b	0.343	69	gemma-3-270m	0.092
32	qwen3-14b	0.337	70	Yi-9B	0.091
33	gemma-2-2b-it	0.321	71	gpt-oss-20b	0.087
34	orca-2-7b	0.315	72	gemma-3-270m-it	0.085
35	flan-t5-xxl	0.302	73	mamba-1.4b-hf	0.083
36	Llama-3.1-Tulu-3-8B-SFT	0.301	74	mamba-2.8b-hf	0.077
37	llama2-13b-chat	0.299	75	gpt-j-6b	0.025
38	OLMo-2-1124-13B	0.299			

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

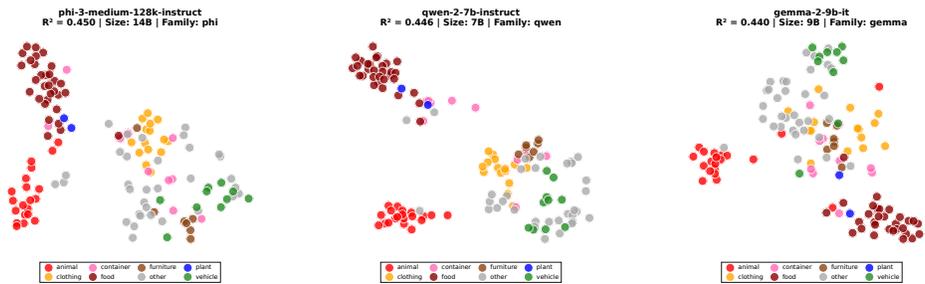
Figure 7: t-SNE visualizations of model embeddings (Part 1). R^2 scores are in parentheses.



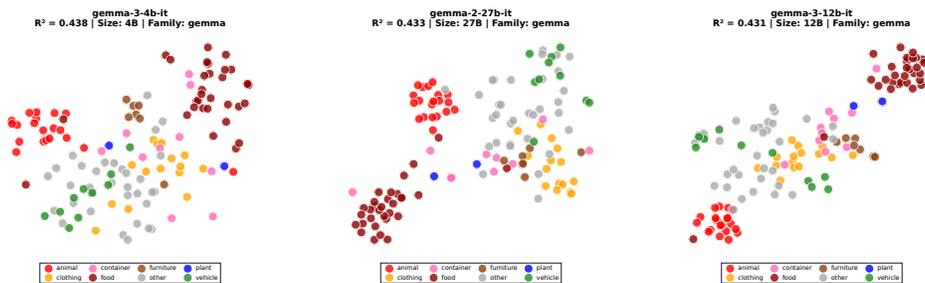
(1) Qwen2.5-14B-Inst (0.492) (2) Phi-3-Med-4k (0.469) (3) Phi-3.5-MoE (0.468)



(4) Gemma-3-27B-IT (0.460) (5) Qwen2.5-7B-Inst (0.459) (6) Phi-3.5-Mini (0.457)



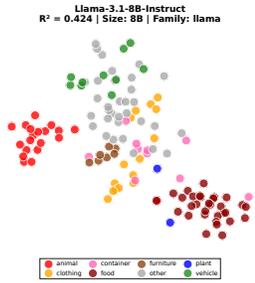
(7) Phi-3-Med-128k (0.450) (8) Qwen2-7B-Inst (0.446) (9) Gemma-2-9B-IT (0.440)



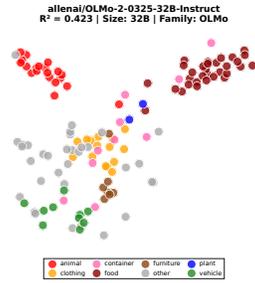
(10) Gemma-3-4B-IT (0.438) (11) Gemma-2-27B-IT (0.433) (12) Gemma-3-12B-IT (0.431)

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

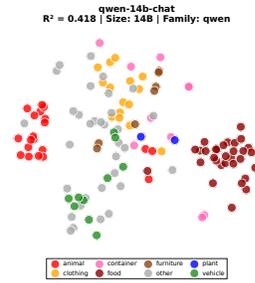
Figure 8: t-SNE visualizations of model embeddings (Part 2). R^2 scores are in parentheses.



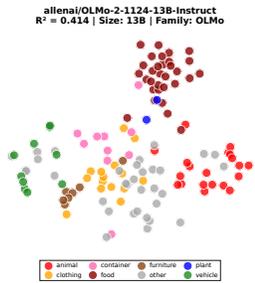
(13) Llama-3.1-8B (0.424)



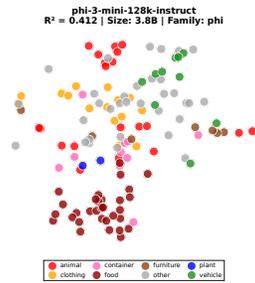
(14) OLMo-2-32B (0.423)



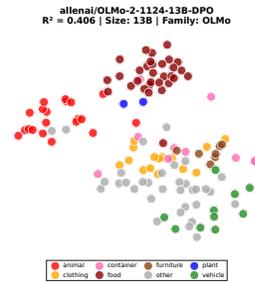
(15) Qwen-14B-Chat (0.418)



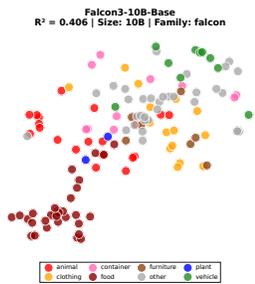
(16) OLMo-2-13B-Inst (0.414)



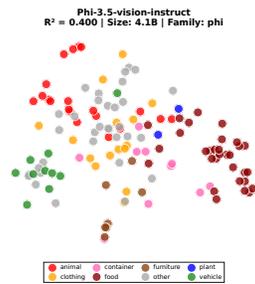
(17) Phi-3-Mini-128k (0.412)



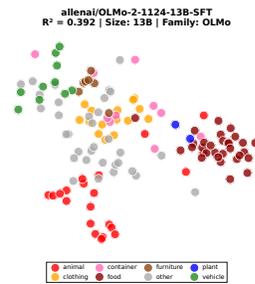
(18) OLMo-2-13B-DPO (0.406)



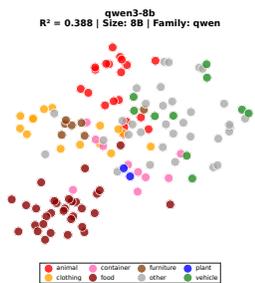
(19) Falcon3-10B (0.406)



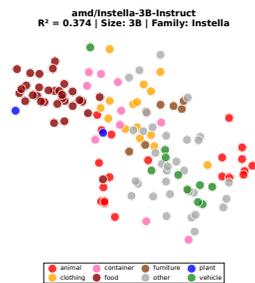
(20) Phi-3.5-Vision (0.400)



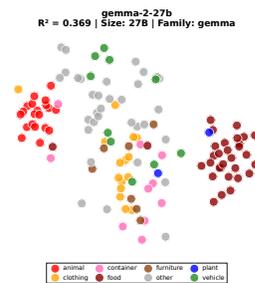
(21) OLMo-2-13B-SFT (0.392)



(22) Qwen3-8B (0.388)



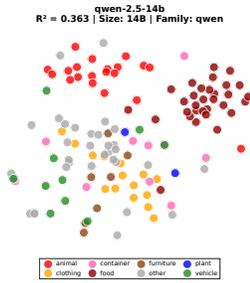
(23) Instella-3B-Inst (0.374)



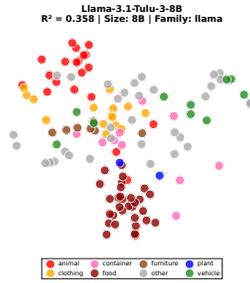
(24) Gemma-2-27B (0.369)

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

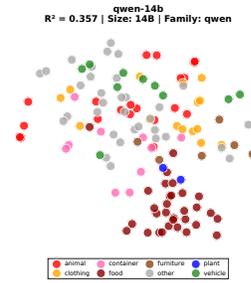
Figure 9: t-SNE visualizations of model embeddings (Part 3). R^2 scores are in parentheses.



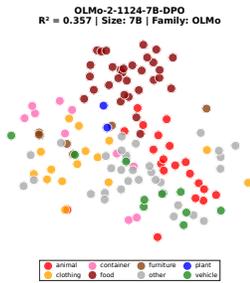
(25) Qwen2.5-14B (0.363)



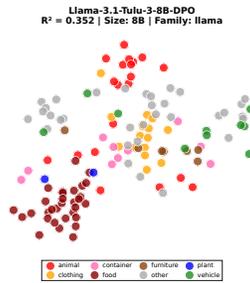
(26) Llama-3.1-Tulu (0.358)



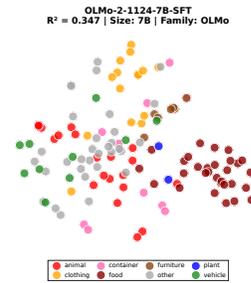
(27) Qwen-14B (0.357)



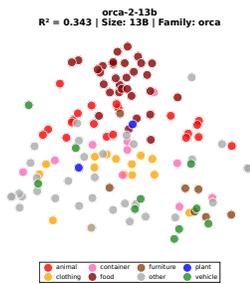
(28) OLMo-2-7B-DPO (0.357)



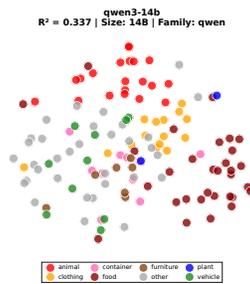
(29) Llama-3.1-Tulu-DPO (0.352)



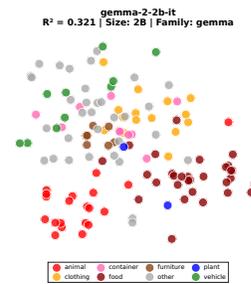
(30) OLMo-2-7B-SFT (0.347)



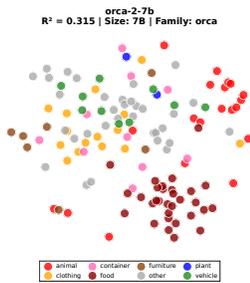
(31) Orca-2-13B (0.343)



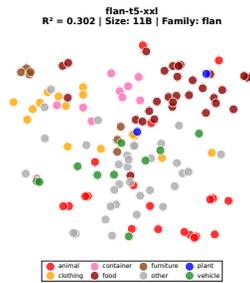
(32) Qwen3-14B (0.337)



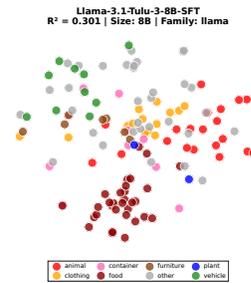
(33) Gemma-2-2B-IT (0.321)



(34) Orca-2-7B (0.315)



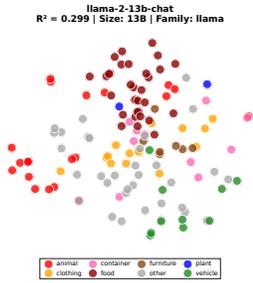
(35) FLAN-T5-XXL (0.302)



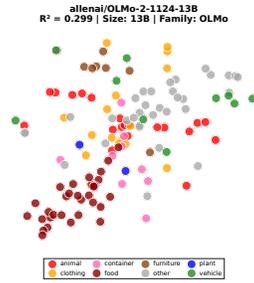
(36) Llama-3.1-Tulu-SFT (0.301)

Figure 10: t-SNE visualizations of model embeddings (Part 4). R^2 scores are in parentheses.

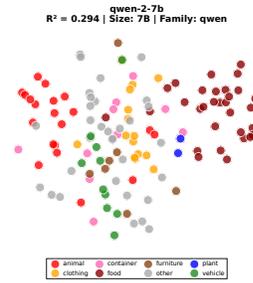
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835



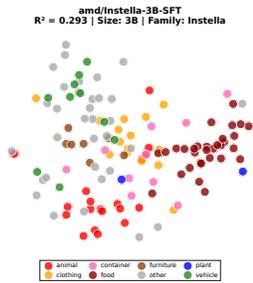
(37) Llama2-13B-Chat (0.299)



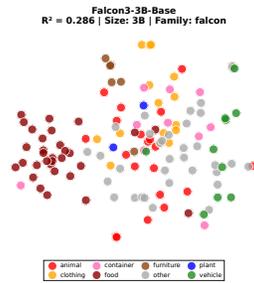
(38) OLMo-2-13B (0.299)



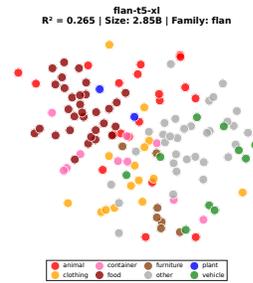
(39) Qwen2-7B (0.294)



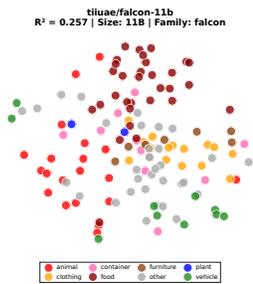
(40) Instella-3B-SFT (0.293)



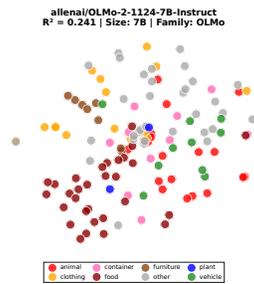
(41) Falcon3-3B (0.286)



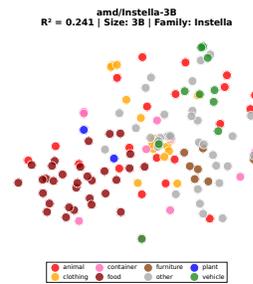
(42) FLAN-T5-XL (0.265)



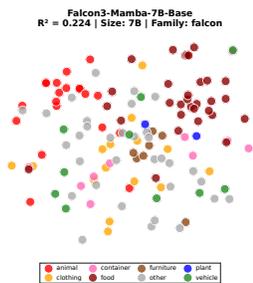
(43) Falcon-11B (0.257)



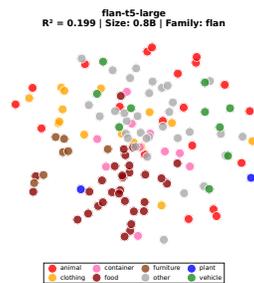
(44) OLMo-2-7B-Inst (0.241)



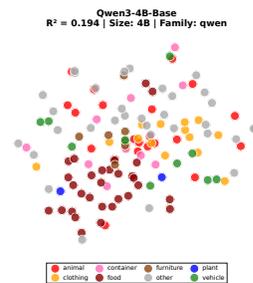
(45) Instella-3B (0.241)



(46) Falcon3-Mamba-7B (0.224)



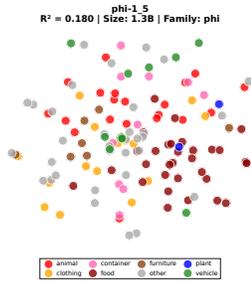
(47) FLAN-T5-Large (0.199)



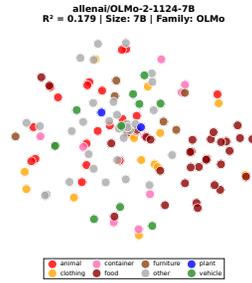
(48) Qwen3-4B (0.194)

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

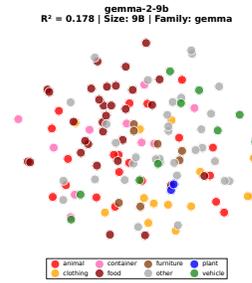
Figure 11: t-SNE visualizations of model embeddings (Part 5). R^2 scores are in parentheses.



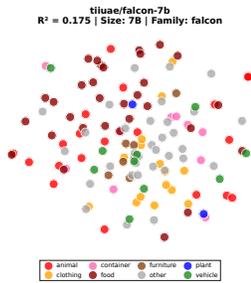
(49) Phi-1.5 (0.180)



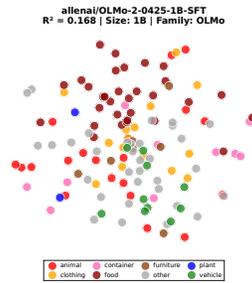
(50) OLMo-2-7B (0.179)



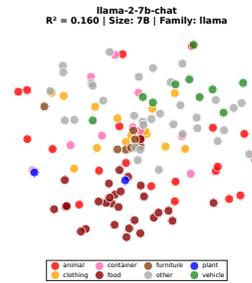
(51) Gemma-2-9B (0.178)



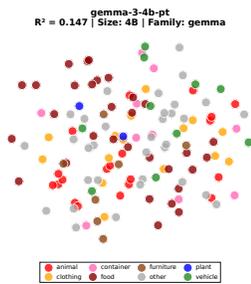
(52) Falcon-7B (0.175)



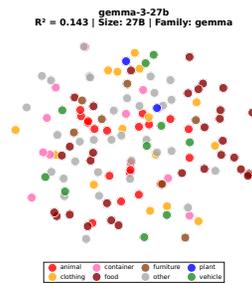
(53) OLMo-2-1B-SFT (0.168)



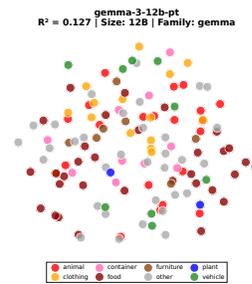
(54) Llama2-7B-Chat (0.160)



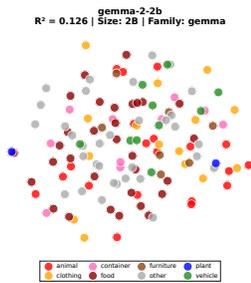
(55) Gemma-3-4B-PT (0.147)



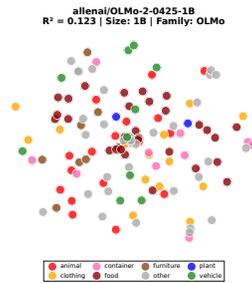
(56) Gemma-3-27B-PT (0.143)



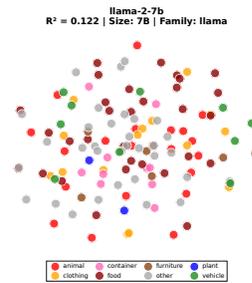
(57) Gemma-3-12B-PT (0.127)



(58) Gemma-2-2B (0.126)

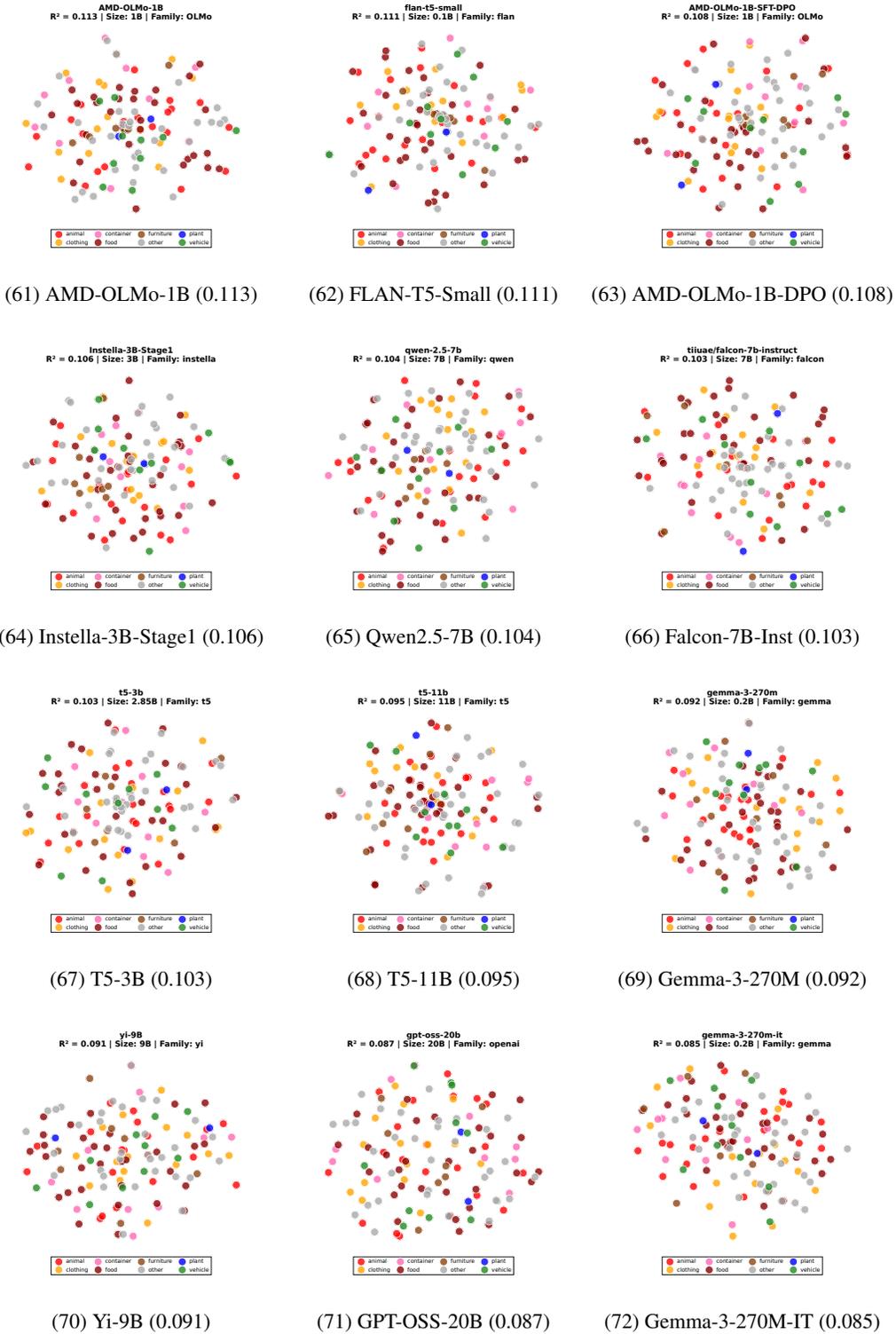


(59) OLMo-2-1B (0.123)



(60) Llama2-7B (0.122)

Figure 12: t-SNE visualizations of model embeddings (Part 6). R^2 scores are in parentheses.



1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

Figure 13: t-SNE visualizations of model embeddings (Part 7). R^2 scores are in parentheses.

