
RETHINKING ‘LANGUAGE-ALIGNMENT’ IN HUMAN VISUAL CORTEX WITH SYNTAX MANIPULATION AND WORD MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent success predicting human ventral visual system responses to images from large language model (LLM) representations of image captions has sparked renewed interest in the possibility that high-level visual representations are aligned to language. Here, we further explore this possibility using image-caption pairs from the Natural Scenes fMRI Dataset, examining how well language-only representations of image captions predict image-evoked human visual cortical responses, compared to predictions based on vision model responses to the images themselves. As in recent work, we find that unimodal language models predict brain responses in human visual cortex as well as unimodal vision models. However, we find that the predictive power of large language models rests almost entirely on their ability to capture information about the nouns present in image descriptions, with little to no role for syntactic structure or semantic compositionality in predicting neural responses to static natural scenes. We propose that the convergence between language-model and vision-model representations and those of high-level visual cortex arises not from direct interaction between vision and language, but instead from common reference to real-world entities, and the prediction of brain data whose principal variance is defined by common objects in common, non-compositional contexts.

1 INTRODUCTION

In order to support object-recognition and other visually-grounded cognitive tasks, the visual system must encode representations that abstract beyond identity-preserving transformations —such as variation in viewpoint, scale, and lighting —while simultaneously providing the discriminative power needed to tell apart thousands of objects, agents, and actions. These high-level visual representations must then interface with higher-level cognitive processes, including language. A fundamental question in cognitive science and visual cognitive neuroscience is how vision and language representations become aligned at this interface, and to what extent visual representations are altered by the handoff from perceptual processing of the sensory input to linguistic abstraction. Such questions have driven foundational research in visual cognitive neuroscience (Huth et al., 2012; Konkle and Oliva, 2012; Weiner and Grill-Spector, 2013; Devereux et al., 2013; Bracci and de Beeck, 2016), but have recently experienced new life with the emergence of multimodal (vision-language) deep neural network models (e.g. CLIP (Radford et al., 2021)) that have produced state-of-the-art results in both canonical computer vision tasks (e.g. image categorization) and in prediction of visual cortical activity evoked by natural images (Wang et al., 2023; Conwell et al., 2023).

The rich history of the debate about what happens at the interface between vision and language ultimately means that ‘language-alignment’ in the context of cognitive (neuro)science means different things to different people. For some, it evokes classic debates about the nature of how the language we speak shapes what we see (Hussein, 2012; Lupyan et al., 2020); for others, it evokes almost the opposite, exposing how the statistics of our perceptual ecology come to be reflected in the ways we communicate about it (Marjeh et al., 2023). One of the most substantive theoretical claims from recent research is that a primary function of high-level visual cortex may actually be generating full ‘semantic scene descriptions’ (Doerig et al., 2022) (i.e. a language-like description of the scene). This claim is predicated on two key findings: one, that purely linguistic embeddings from large language models (e.g. GUSE or Google’s Universal Sentence Encoder) applied to image captions are capable of predicting the majority of explainable variance in image-evoked high-level visual cortical activity;

054 two) that natural language descriptions to previously unseen images may be decoded with reasonable
055 accuracy from an embedding model fit directly to image-evoked brain activity. This latter finding
056 adds to a rapidly growing list of works that deploy models (seemingly with great success) to reading
057 out the contents and structure of mental life from neural response patterns (Takagi and Nishimoto,
058 2022; Luo et al., 2023; Tang et al., 2023a).

059 One overarching issue with work of this nature, however, is that it remains unclear the extent to
060 which the results are driven by information in the brain itself versus the priors of our models. In using
061 language models to predict visual brain activity, especially, there remains substantial ambiguity as to
062 which aspects of these models capture the structure and content of high-level visual representations,
063 and whether language-alignment *per se* improves the correspondence between language-model
064 representations and high-level visual cortex.

065 In this work, we address these questions by predicting visual responses to the large-scale human fMRI
066 Natural Scenes Dataset (NSD) (Allen et al., 2022) using a variety of unimodal language, unimodal
067 vision, and multimodal (language-aligned) vision models. Like others, we find that unimodal language
068 models are indeed capable of predicting image-evoked brain activity as accurately as unimodal vision
069 models. We also find, however, that the predictive power of large language models (in this dataset)
070 reduces almost entirely to a basis set of simple nouns in no syntactic order, with little to no role
071 for other parts of speech, or compositional semantics. Applying this intuition to recently proposed
072 ‘relative representation’ (anchor point embedding) techniques (c.f. Moschella et al., 2022; Maiorca
073 et al., 2024; Norelli et al., 2024), we show that we can even ‘hand-engineer’ a set of 62 simple words
074 whose relative ‘word’ coordinates in a multimodal foundation model (CLIP) effectively explain the
075 same amount of variance in the image-evoked brain data as the underlying image features themselves.
076 Taken together, these results suggest language model predictivity of visual cortical activity in the
077 Natural Scenes Dataset may have little to do with language *per se*, and far more to do with the
078 recovery of ‘grounded information’ from co-occurrence statistics. This adds as well to a growing
079 consensus that ‘vision’ and ‘language’ – at least as learned by modern artificial intelligence algorithms
080 – are in some sense *already aligned* (Pavlick, 2023; Huh et al., 2024), even in the absence of explicit
081 cross-modal learning.

082 2 RESULTS

083 Our main experimental assay consists of using features extracted from vision-only, language-only,
084 or hybrid vision-language hierarchical deep learning models, and shallow word-vectorizing models,
085 to predict the representational geometries of voxel responses in the early visual (EVC) and occipi-
086 totemporal cortices (OTC) of 4 subjects viewing 1000 MS-COCO images from the 7T fMRI Natural
087 Scenes Dataset (Allen et al., 2022). We split these 1000 images into a training and test set of 500
088 images each. Language descriptions of these images come in the form of captions (5 per image,
089 provided as part of the COCO metadata).

092 We employ two metrics of model-to-brain comparison: classical (unweighted) and voxelwise-
093 encoding representational similarity analysis (cRSA and eRSA, respectively) (Kriegeskorte et al.,
094 2008a; Kaniuth and Hebart, 2021; Konkle and Alvarez, 2022). cRSA considers all of the features
095 from a given model equally in computing an image-wise representational similarity matrix (RSM),
096 which is then directly compared with the target EVC or OTC RSM. eRSA involves first fitting
097 voxelwise encoding models from features maps with a combination of sparse random projection
098 (for dimensionality reduction) and ridge regression (with cross-validated lambda hyperparameters
099 for each voxel) The predicted responses from these encoding models are then used to generate a
100 reweighted RSM, which we then compare to the target EVC or OTC RSM¹. For DNN models, we
101 compute scores across all layers on the training set, and select the most predictive layer for assessment
102 on a held-out test set. All scores we report are the generalization scores of each RSA metric on this
103 held-out test set, with no contamination from selection procedures used on the training set (including
104 layer selection and voxel-encoding hyperparameters).

105 ¹An important methodological note here is that these ‘predicted responses’ can be directly converted into
106 ‘voxel-wise encoding scores’: that is, the correlation between predicted and actual responses *per voxel*. In this
107 work, we choose to report only the cRSA and eRSA scores so that they may be directly compared – both in
terms of their noise ceilings (GSN) and units (dissimilarity in $1 - r_{Pearson}$).

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

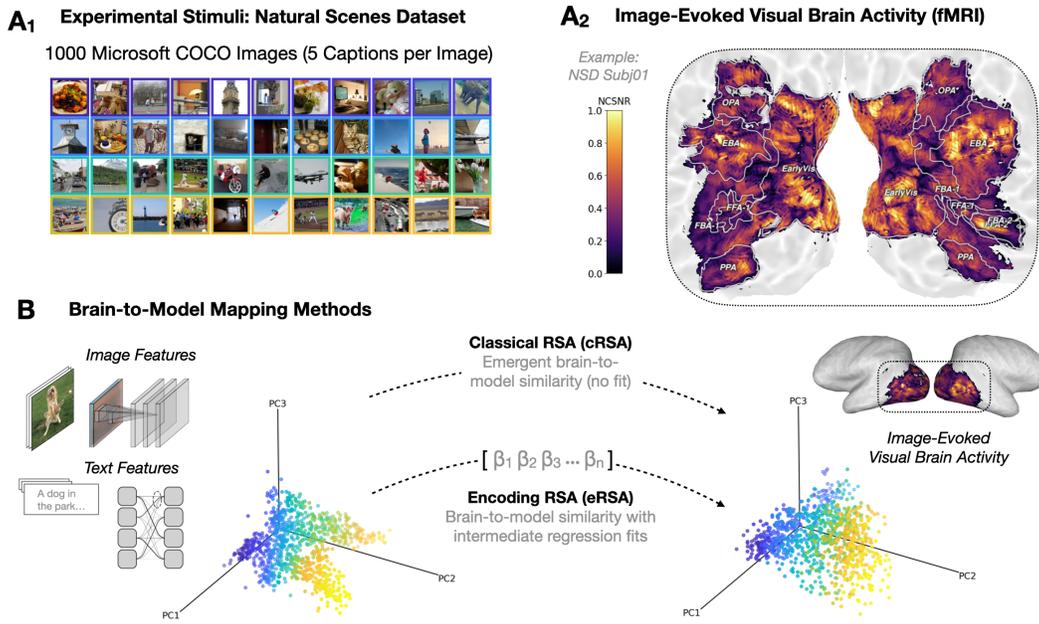


Figure 1: Overview of our experiments. **A** Our primary dataset consists of brain activity in the ventral stream of 4 subjects from the Natural Scenes Dataset viewing 1000 images from the Microsoft COCO (Lin et al., 2014) dataset. Each of these Microsoft COCO images is associated with 5-6 captions provided by human annotators. **B** To predict this activity, we first extract either image or text features from a deep neural network model, using either the images themselves, or the first 5 COCO captions associated with each image. To map these features to brain activity, we use one of two forms of representational similarity analysis (RSA): classical RSA, which enforces fully emergent similarity between a given model’s feature space and that of the brain, or, a voxelwise-encoding RSA, which involves first linearly re-weighting model features to predict each of the voxels in a target ROI, and then using held-out test images to generate a model-predicted representational similarity matrix (RSM) for comparison to the true brain RSM.

Finally, we compare these scores to subject- and ROI-specific noise ceilings that are estimated using the Generative Modeling of Signal and Noise (GSN) toolbox Kay et al. (2024). In brief, GSN attempts to model the distinct contributions of signal and noise covariance structure to the observed brain measurements, to then estimate the maximum degree to which a computational model can predict the RSM of a given pool of voxels. These noise ceilings are $r_{Pearson} = 0.69$, $CI_{95} = [0.63, 0.74]$ in EVC and $0.80 [0.73, 0.85]$ in OTC. An overview of our methodology is available in Figure 1. More details on all our experimental methods, including brain data preprocessing, voxel selection, model selection, feature extraction, brain mapping metrics, and noise ceiling calculations, may be found in the Methods Appendix.

A summary of results for all model comparisons and manipulations is displayed in Table 1. Unless otherwise noted, we use the following convention in the reporting of summary statistics: arithmetic mean [lower 95%, upper 95% bootstrapped confidence interval].

2.1 VISION-ONLY VERSUS LANGUAGE-ONLY MODELS

As a primary point of comparison, we consider the relative difference in brain-predictivity of the unimodal models (vision-only versus language-only). The main question we are asking in this analysis is whether these two drastically different model types (the latter of which learns in the absence of visual input) are nonetheless comparable in their ability to predict responses in visual cortex.

The vision-only models in this comparison consist entirely of self-supervised (visual) contrastive learning models, whose learned representations are the product of a training procedure that operates

Table 1: Comparison of Model Types in Predicting OTC Activity

Analysis Name	Model Type	OTC-Predictive Accuracy [\pm 95% BCI]			
		cRSA Score		eRSA Score	
Vision vs Language	Vision-only (mean)	0.338	[0.329, 0.348]	0.682	[0.662, 0.700]
	Language-only (mean)	0.277	[0.257, 0.297]	0.662	[0.650, 0.675]
	Vision-only (max)	0.384	[0.362, 0.392]	0.712	[0.703, 0.721]
	Language-only (max)	0.437	[0.414, 0.466]	0.689	[0.673, 0.694]
Word-Level Models	CountVec (trigrams)	0.210	[0.192, 0.221]	0.584	[0.556, 0.601]
	GLOVE (all words)	0.320	[0.300, 0.340]	0.650	[0.610, 0.700]
	GLOVE (nouns only)	0.310	[0.290, 0.330]	0.630	[0.590, 0.680]
Anchor Point Embeds	CLIP-Vision (768-D)	0.322	[0.311, 0.340]	0.668	[0.624, 0.734]
	62-word (hypothesis)	0.320	[0.305, 0.337]	0.671	[0.635, 0.708]
	62-word (random sample)	0.211	[0.202, 0.221]	0.556	[0.523, 0.590]

only over individual image instances, and involves no explicit semantic labels (e.g., they do not rely on the one-hot category encoding vectors that define category distinctions in supervised object recognition models).

The language-only models in this comparison consist entirely of transformer-based deep neural network models (Vaswani et al., 2017) trained on one of two tasks: masked language modeling (the prediction of a *masked* token removed at random from an input sequence of tokenized words) (Devlin et al., 2018) or causal language modeling (the prediction of the *next* word following an input sequence of tokenized words) (Radford et al., 2018). The primary operations in these models consist of multi-head attention computations over the tokenized words (plus positional embeddings) of a given sentence. Their learned representations can thus (in both theory and practice) capture transition probabilities and long-range dependencies between different words and concepts.

These model classes learn from different data modalities, under very different task constraints. Remarkably, despite these differences, we find the OTC predictivity of these two model types to be comparable (see Figure 2A): for the eRSA metric, these sets are not significantly different in their average brain predictivity (vision-only mean $r_{Pearson} = 0.682$ [0.662, 0.700], language-only mean $r_{Pearson} = 0.662$ [0.65, 0.675]; $p = 0.08$, $g_{Hedges} = -0.92$). For the more stringent cRSA metric, the scores of the vision-only models are significantly higher on average than those of the language-only models ($r_{Pearson} = 0.338$ [0.329, 0.348], mean $r_{Pearson} = 0.277$ [0.257, 0.297], respectively; $p = 0.031$, $g_{Hedges} = -0.915$). Worth noting, however, is that there is a high degree of variability in cRSA scores amongst the language-only models, which obscures a striking standout: SBERT-Mini-LM6, a language-only model whose cRSA score is the highest of all the models we survey (mean $r_{Pearson} = 0.437$ [0.414, 0.466]; the next highest-scoring model in cRSA yields $r_{Pearson} = 0.380$ [0.369, 0.393]). In other words, the highest-ranking (unweighted) model of the representational geometry of high-level visual cortex in this survey is not a visual model at all, and learns entirely without visual input.

We next compare the prediction of pure-vision and pure-language models in prediction of early visual cortex. Here, we find that the same language-only models that perform comparably with vision-only models in prediction of OTC perform uniformly worse (and by a large margin) in prediction of EVC. The mean accuracy of vision-only models in EVC is $r_{Pearson} = 0.295$ [0.28, 0.31] in cRSA and 0.48 [0.46, 0.5] in eRSA. The mean accuracy of language-only models is $r_{Pearson} = 0.087$ [0.081, 0.095] in cRSA and 0.149 [0.14, 0.16] in eRSA. This difference is significant and substantial in both metrics ($p = 4.11e^{-21}$, $g_{Hedges} = -11.2$ in cRSA; $p = 7.78e^{-19}$, $g_{Hedges} = -32.5$ in eRSA). Thus, despite their relative parity in high-level visual cortex, we find that language models yield very little in terms of the representational structure necessary to predict responses in early visual cortex, where voxels are tuned to lower-level image attributes such as oriented lines, edges, and textures.

A summary of these results is available in Figure 2 and Table 1 (Analysis Name: Vision versus Language). The results of two follow-up analyses (one that compares the randomly initialized versions of these models to assess for differences in architectural inductive bias, and another that uses

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323



Figure 3: **Syntax Manipulation & Simpler NLP Modeling:** **A** Examples of the 3 main syntactical manipulations we apply to each caption: full word shuffle (sentence scrambling); bag of nouns; bag of verbs. **B** Schematic of the sentence and word embeddings models we use to ‘decompose’ the LLM performance: a representative LLM (SBERT-Mini-LM) from our survey of unimodal language models; a word-count vectorized model; and a word embedding model (GLOVE). **C** Results (change in brain predictivity) of the syntax manipulations applied to the input of the 3 different NLP models we assess. Crossbars are the performance of a given model under a certain syntax manipulation (opaque crossbars are the eRSA metric scores; translucent crossbars are the cRSA metric scores)

The lowest performing model in this set of experiments consists of word count vectors computed over the verbs extracted from each image caption (mean $r_{Pearson} = 0.135$ [0.13, 0.15] in cRSA; 0.32 [0.27, 0.38] in eRSA). The highest performing word-level model consists of the averaged GLOVE word embeddings computed over all individual words in each caption (mean $r_{Pearson} = 0.32$ [0.30, 0.34] in cRSA; 0.65 [0.61, 0.7] in eRSA). Note already that the performance of this model (which involves no nonlinear hierarchies, and no explicit syntax) is comparable with that of the average LLM (mean $r_{Pearson} = 0.277$ [0.257, 0.297] in CRSA; 0.662 [0.65, 0.675] in eRSA). The second highest-performing word-level model is the average GLOVE embeddings for nouns only (mean $r_{Pearson} = 0.31$ [0.29, 0.33] in CRSA; 0.63 [0.59, 0.68] in eRSA). Thus, embeddings from a shallow (log bi-linear) word model (GLOVE), given only nouns, are capable of predicting OTC activity as accurately as LLM embeddings derived from full sentences. Further evidence for the representational significance of nouns may be seen in the performance of the word count model computed over nouns only. This model (unlike GLOVE) does not leverage co-occurrence statistics, yet still accounts for

the majority of explainable OTC variance in eRSA (mean $r_{Pearson} = 0.491$ [0.438, 0.564]) and is competitive with mean LLM performance in cRSA ($r_{Pearson} = 0.257$ [0.246, 0.272]).

In sum, this set of experiments suggests that the performance of the language models in predicting visual responses is accounted for almost entirely by the covariance structure instantiated by the nouns of the COCO captions. What exact role the co-occurrence statistics learned by a model like GLOVE are playing in these covariance structures remains unclear. What *is* clear is that whatever representations the LLMs are yielding in their prediction of high-level visual cortex, these representations need not be any more sophisticated than a single affine transformation of token-vectorized (embedded) nouns.

2.3 ‘HANDCRAFTED’ WORD MODELS BY WAY OF ANCHOR POINT ANALYSIS

Given the relatively high predictive power of the word-level models computed over human-provided captions, we next attempted to construct a simple, hypothesis-driven ‘word model’ that could effectively capture the variance in occipitotemporal cortex. This process involved two steps. The first step required an element of conjecture to determine which words might sufficiently capture the range of representations evoked by the images in our target stimulus set. (For a real-world example of a similar process, consider the prompt-engineering used in zero-shot evaluations of language-aligned models such as CLIP (Radford et al., 2021)). To score our word model in its prediction of the brain, we used a variant of relative representation analysis (Moschella et al., 2022) (sometimes, and in this work, referred to as anchor point embedding analysis).

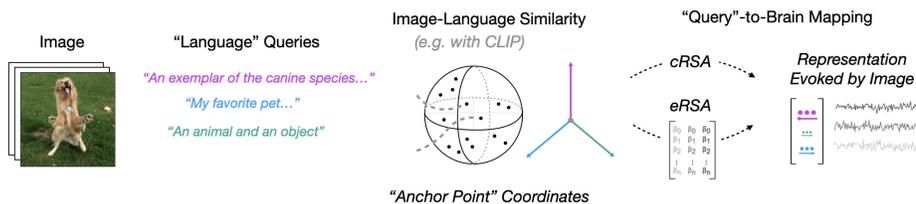
Words we included in our hypothesis-driven word model included global descriptors (adjectives applied to whole image, e.g. ‘high-resolution’ or ‘colorful’); agents and objects (e.g. ‘man’, ‘woman’, ‘child’, ‘animal’, ‘vehicle’, ‘food’); places (e.g. ‘desk’, ‘beach’, ‘snow’, ‘desert’); and times-of-day (e.g. ‘morning’, ‘night’). To derive brain-predictivity scores from these prompts, we first generated embeddings for each prompt using CLIP-ResNet50’s language encoder, which is a modified RoBERTa architecture. (Note that the use of CLIP is somewhat arbitrary, since this same analysis can be done with *any* algorithm that provides some form of image-text similarity score; see (Maiorca et al., 2024; Norelli et al., 2024)). We then generated the embeddings for each of our target image stimuli, and computed the cosine similarity between each image embedding and all the embeddings associated with our prompts. Finally, we aggregated the resultant image-text similarity matrix together (without softmax) as its own ‘feature set’, scoring this feature set’s brain-predictivity in the same way we scored the feature sets derived from all the models above. As a control, we generated 1000 random samples of unigrams and bigrams from the Brown NLTK corpus (with stimulus counts matching our ‘hypothesized’ word space), and scored the brain predictivity of these samples using the same encoding pipeline described above.

Strikingly, we found this CLIP-mediated anchor point analysis to perform remarkably well in predicting OTC activity: After some iterative selection (using nested cross-validation scoring on our training set of 500 images), we distilled a 62 word model (consisting only of adjectives and nouns) capable of describing as much variance in the occipitotemporal cortical responses as the underlying CLIP image embeddings (a 768-D vector) from which they were derived: 0.32 [0.305, 0.337] versus 0.322 [0.311, 0.34] in cRSA and 0.671 [0.635, 0.708] versus 0.668 [0.624, 0.734] in eRSA. The mean of the 1000 N=62 random word samples was noticeably lower, but not altogether poor: 0.556 [0.523, 0.59]. A summary of these results is displayed in Figure 4 and Table 1 (Analysis Name: Anchor Point Embeds).

While it might seem odd (or undermining of our hypothesized word space) that these randomly sampled words predicted brain activity so accurately on average, this finding actually strengthens the point this analysis intended to make. Specifically, what sometimes appears to be dense, compositionally complex, or richly structured information in the embedding spaces of large multimodal foundation models can often be reduced to far simpler basis sets. These simpler sets work so long as they provide sufficient coverage of the (representational) variance we are trying to explain. Since the Natural Scenes Dataset is composed of COCO images, its major axes of variance capture the objects and scene attributes found in each of its images. Models that sufficiently capture word co-occurrence appear to predict these visual brain responses as long as their set of natural language queries covers the relevant part of the representational space evoked by the objects in the images.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

A (CLIP-Based) Anchor Point Embedding Analysis



B 'Handcrafted' Visual Word Model Experiment

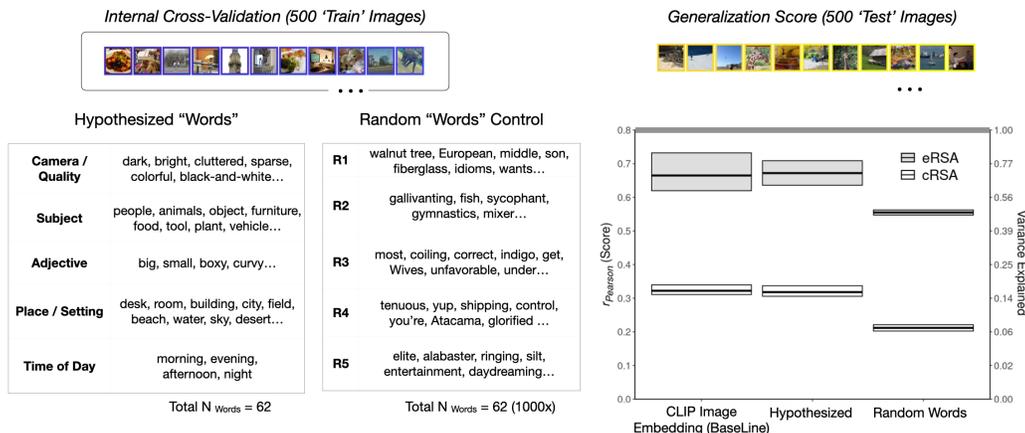


Figure 4: **Anchor Point "Word" Model Analysis.** In **A** we show a schematic of our analysis using arbitrary natural language queries embedded in CLIP to compute relative representations (Moschella et al., 2022), which are anchor points we subsequently use to predict image-evoked brain activity. This procedure involves first proposing candidate queries; then, computing the text embeddings for each of these queries; then, computing the text-to-image similarity of each image to each of the queries; and finally, using these text-to-image similarities as the basis of the mapping procedure to brains (cRSA or eRSA). In **B** we show the results of an experiment in which we contrast the predictivity of a 'hypothesis-driven' 62-Word model (derived from a mix of common words describing diverse object and scene attributes and CLIP-style prompt injections) with CLIP's image embeddings (768-D) and the mean of 1000 random 62-word selections. Note, that all 'development' of the anchor point word models is done via internal cross-validation on a training set of 500 images. The scores reported in the figure are the scores of the final models on a held-out test set of 500 images (and their associated captions).

3 DISCUSSION

In this work, we demonstrate that vision and language models predict visual cortex responses equally well, indicating that they may be converging on shared representational structure dominated by objects and agents (i.e. nouns). While significant future work will be necessary to further characterize this structure, our analyses suggest that this structure is obtainable either through purely visual bottom-up filtering operations over natural images, or, by performing multihead attention operations over the transition probabilities between words in the context of sentences. We further find both that attention and sentential context is largely unnecessary to match the predictive power of language in high-level vision. The 'or' of the first conclusion here is particularly important, as it argues against the idea that representations in high-level visual cortex are *predominantly* semantic and abstract in nature. To make this exclusion explicit, consider the evidence we would need for the opposite conclusion: a statistically significant case in which the unimodal language models or language-aligned vision models *outperformed* the unimodal vision models in prediction of OTC. Our sampling of the (recently SOTA) models from each of these categories does not provide such evidence.

We are not the first to demonstrate many of these trends. Long before the advent of sentence-processing language transformers, for example, Carlson et al. (2014) observed that the representational

432 geometry in occipitotemporal cortex is well captured by the semantic similarity features from earlier
433 NLP models such as WordNet (Fellbaum, 2010). Similar observations were made by Huth et al.
434 (2012), who used WordNet graphs to map a continuous semantic space of object and action categories
435 that bridged visual and nonvisual cortex alike.

436 What does the surprising effectiveness of word-level language models in predicting high-level visual
437 activity mean for the larger question of language alignment in the human visual system? Much
438 of this answer depends on the weight we attribute to a number of limitations that otherwise scope
439 or circumscribe our interpretation of the results we obtain from the particular models, metrics of
440 representational alignment, and neural dataset that scopes our analysis.

441 **Limitations of the Data (ROI Subset):** A first potential limitation is our particular treatment of
442 data from high-level visual cortex in this study (deriving a target RDM from a single broad mask of
443 occipitotemporal cortex). Popham et al. (2021) have previously suggested that language alignment in
444 high-level visual cortex may be limited to a narrower subset of neural real estate in the most anterior
445 portions of the ventral stream than those we have sampled here. Recent work from this same group
446 has found that multimodal (language-aligned vision) models may be particularly adept at predicting
447 the activity in these areas (Tang et al., 2023b). Others have analyzed data from the Natural Scenes
448 Dataset at the voxel level, finding that models trained with language feedback can account for modest
449 unique variance within certain brain areas such as the extrastriate body area (EBA) and some parts of
450 the fusiform face area (FFA), relative an otherwise identical model that does not include language
451 feedback (Wang et al., 2023).

452 **Limitations of the Data (Probe Stimuli):** The Natural Scenes Dataset is arguably the current best
453 dataset we have available for this kind of analysis in terms of reliability and stimulus-density – but it
454 is far less optimal in terms of the particular stimuli it relies on. COCO images and their associated
455 captions (solicited, by design, to be as visually grounded as possible) are almost certainly *not* the
456 optimal stimuli for assessing the presence of linguistically interesting structure in high-level visual
457 cortex. For this, we need far more targeted datasets evoking abstractions that language as a tool (c.f.
458 Fedorenko et al., 2024) is particularly well-suited for.

459 **Limitations of the Models:** Perhaps overlooked in discussions about the limitations of the brain
460 data and probe stimuli are the limitations of the candidate models themselves. The search for
461 ‘language-like’ structure in the visual brain – through comparison to pure-language or language-
462 aligned vision models – assumes such structure exists in the models (c.f. Doerig et al., 2022). However,
463 an increasingly large body of empirical work challenges even this core assumption. For instance,
464 CLIP and related models, including text-to-image systems like Stable Diffusion (Rombach et al.,
465 2022), often lack basic compositional structure, such as distinguishing ‘a spoon in a cup’ from ‘a
466 cup on a spoon’ (Conwell and Ullman, 2022; Yamada et al., 2023). These models frequently behave
467 as ‘bags-of-words’ (Yuksekgonul et al., 2023), representing concepts in ways easily learned by
468 simpler vision models. Even unimodal language models (LLMs) face similar criticism. Probes of
469 entailment, negation, and counting suggest these models often fail to meet linguistic standards of
470 finiteness, discreteness, syntactic well-formedness, and the ability to produce the new, but meaningful
471 constructions that syntactic well-formedness allows (Thrush et al., 2022; Press et al., 2022; Bertolini
472 et al., 2022; Hauser et al., 2002).

473 **(Cautiously) General Conclusions** The limitations above do potentially put somewhat strict and
474 finite limits on broader conclusions we might make based on this work. Nevertheless, given the
475 NSD’s current centrality in the landscape of visual cognitive neuroscience datasets, and pending
476 the emergence of datasets more deeply enriched with stimuli that elicit more complex linguistic
477 structure, we do (for now) interpret our results as follows: Language alignment in the visual system
478 is not an organizing force, but rather, a byproduct of the visual system’s goal to capture statistically
479 salient or ecologically relevant features of the visual world. Language models may align well with
480 the high-level visual cortex because vision organizes the world into units that map easily to learned
481 languages. These units, many of which correspond to distinct words, often have distinct visual
482 features. In datasets like NSD, which comprehensively sample natural image statistics, these visual
483 features account for most of the variance in visually evoked brain responses. Therefore, explicit
484 language alignment may do little to improve predictions of the high-level visual cortex because
485 language is already aligned with vision. Word-level models, even without sentence structure, appear
to already capture the primary dimensions of this alignment.

486 Concurrent trends in machine learning also provide (preliminary) mechanistic evidence supporting this
487 idea. Contrastive self-supervised learning models (Chen et al., 2020; Zbontar et al., 2021; Goyal et al.,
488 2021a; Konkle and Alvarez, 2022), trained without any semantic supervision, readily learn features
489 that support downstream object recognition with a single linear transformation. In language-alignment
490 research, some have noted that algorithms such as CLIP (which backpropagates its alignment loss
491 across the entirety of its vision and language encoders) may be inefficient because they do not leverage
492 pre-existing structure in vision and language modalities: Alternatives to CLIP, such as LiT and
493 DeCLIP, for example, maintain many of CLIP’s advantages (e.g., zero-shot classification, robustness,
494 guided conditional sampling) with largely frozen visual backbones pretrained via unimodal self-
495 supervision (Li et al., 2022; Zhai et al., 2022). The success of these models suggests that whatever
496 representational restructuring the language alignment task is doing, it need not be as deep or extensive
497 as we think. Zooming out even further, the relative parity of vision and language models in predicting
498 high-level visual cortical activity, may in some ways be a direct mirror of the datasets we’ve used to
499 train machine vision models since the earliest days of deep learning. The canonical 1000 categories
500 and 1.2 million images of ImageNet1K dataset (Deng et al., 2009) is a subset of a larger 14 million
501 image dataset whose labels are ‘synsets’ from WordNet. From AlexNet onwards, then, our most
502 popular visual models have often been trained on an image set whose primary dimensions of variance
503 are defined by language.

504 This latter point underscores the profound challenge of disentangling vision from language and may
505 explain why debates about which modality shapes the other are difficult to resolve conclusively
506 (Konkle and Oliva, 2012; Grill-Spector and Weiner, 2014; Bracci and de Beeck, 2016; Long et al.,
507 2018; de Beeck et al., 2023). Another trend we observe in the machine-learning literature is that large,
508 well-trained models increasingly converge on similar representations, even without direct cross-modal
509 learning (Pavlick, 2023; Huh et al., 2024). And while this convergence is illuminating, it highlights
510 the need for more precise diagnostic tests to separate model representations in domains where they
511 should be better specialized.

511 **Future Directions** Coming back to the question of dataset, then, we believe the data we really need is
512 data that pushes perception to its natural limit, and therefore necessitates the involvement of language
513 for understanding. Already, we are beginning to see work that uses the inherent ‘abstractability’ of
514 language (i.e. more or less grounded descriptions of the same perceptual stimulus, predicated as well
515 on factors like mutual familiarity and active co-reference) to build stimulus sets specifically targeted at
516 teasing apart vision from language. For example, recent work by (Shoham et al., 2024) using a custom
517 dataset and iEEG recordings accords well with what we might expect based on the modeling we have
518 done in this work: Language models given far more abstract descriptions of visual stimuli predict
519 high-level visual brain data far worse than vision models given the visual stimuli directly. We consider
520 this work an excellent step in the right direction, if not yet still the fullest picture we might obtain
521 with similarly targeted datasets. ‘Abstractability’ need only be one tool in a diverse toolkit of probes
522 that evoke the signature functional and representational structures we consider relatively more or even
523 uniquely ‘linguistic’. These structures include but are not limited to: compositional inversions and
524 role-filler distinctions (i.e. the difference of transmitted meaning in ‘man bites dog’ versus ‘dog bites
525 man’) (Frankland, 2015; Quilty-Dunn et al., 2023); external reference (i.e. the ability to represent
526 ‘cat’ in the sentence ‘the orange cat that I saw yesterday’); conceptual abstractions (e.g. justice); and
527 (more generally) any kind of representational invariance that no combination of feed-forward visual
528 filtering operations could feasibly produce (e.g. the representation of the written word ‘apple’ and
529 a picture of an apple as the same) (c.f. Quiroga et al., 2005). Language tends to excel in domains
530 requiring relational computations or logical operators, while vision excels in domains where granular
531 or holistic distinctions (e.g., texture) are not easily expressible in natural language. Combining
532 diagnostic tests that stress this difference with psychophysical approaches used by neuroscientists
533 to probe for multimodality outside the visual cortex could help us determine when, where, and how
534 truly “sentence-like” abstractions—those that can only exist through language—emerge in neural
535 learning systems grounded in sensory perception.

536 REFERENCES

537
538 Alexander G Huth, Shinji Nishimoto, An T Vu, and Jack L Gallant. A continuous semantic space
539 describes the representation of thousands of object and action categories across the human brain.
Neuron, 76(6):1210–1224, 2012.

540 Talia Konkle and Aude Oliva. A real-world size organization of object responses in occipitotemporal
541 cortex. *Neuron*, 74(6):1114–1124, 2012.

542 Kevin S Weiner and Kalanit Grill-Spector. Neural representations of faces and limbs neighbor in
543 human high-level visual cortex: evidence for a new organization principle. *Psychological research*,
544 77:74–97, 2013.

546 Barry J Devereux, Alex Clarke, Andreas Marouchos, and Lorraine K Tyler. Representational similarity
547 analysis reveals commonalities and differences in the semantic processing of words and objects.
548 *Journal of Neuroscience*, 33(48):18906–18916, 2013.

549 Stefania Bracci and Hans Op de Beeck. Dissociations and associations between shape and category
550 representations in the two visual pathways. *Journal of Neuroscience*, 36(2):432–444, 2016.

552 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
553 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
554 models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.

556 Aria Y Wang, Kendrick Kay, Thomas Naselaris, Michael J Tarr, and Leila Wehbe. Better models of
557 human high-level visual cortex emerge from natural language supervision with a large and diverse
558 dataset. *Nature Machine Intelligence*, 5(12):1415–1426, 2023.

559 Colin Conwell, Jacob S. Prince, Kendrick N. Kay, George A. Alvarez, and Talia Konkle. What can
560 1.8 billion regressions tell us about the pressures shaping high-level visual representation in brains
561 and machines? *BioRxiv*, 2023. doi: 10.1101/2022.03.28.485868. Publisher: Cold Spring Harbor
562 Laboratory.

563 Basel Al-Sheikh Hussein. The sapir-whorf hypothesis today. *Theory and Practice in Language
564 Studies*, 2(3):642–646, 2012.

566 Gary Lupyan, Rasha Abdel Rahman, Lera Boroditsky, and Andy Clark. Effects of language on visual
567 perception. *Trends in cognitive sciences*, 24(11):930–944, 2020.

568 Raja Marjeh, Ilia Sucholutsky, Pol van Rijn, Nori Jacoby, and Thomas L Griffiths. Large language
569 models predict human sensory judgments across six modalities. *arXiv preprint arXiv:2302.01308*,
570 2023.

572 Adrien Doerig, Tim C Kietzmann, Emily Allen, Yihan Wu, Thomas Naselaris, Kendrick Kay,
573 and Ian Charest. Semantic scene descriptions as an objective of human vision. *arXiv preprint
574 arXiv:2209.11737*, 2022.

575 Yu Takagi and Shinji Nishimoto. High-resolution image reconstruction with latent diffusion models
576 from human brain activity. *biorxiv*. CVPR, 2022.

578 Andrew F Luo, Margaret M Henderson, Michael J Tarr, and Leila Wehbe. Brainscuba: Fine-grained
579 natural language captions of visual cortex selectivity. *arXiv preprint arXiv:2310.04420*, 2023.

580 Jerry Tang, Amanda LeBel, Shailee Jain, and Alexander G Huth. Semantic reconstruction of
581 continuous language from non-invasive brain recordings. *Nature Neuroscience*, 26(5):858–866,
582 2023a.

584 Emily J Allen, Ghislain St-Yves, Yihan Wu, Jesse L Breedlove, Jacob S Prince, Logan T Dowdle,
585 Matthias Nau, Brad Caron, Franco Pestilli, Ian Charest, and others. A massive 7T fMRI dataset
586 to bridge cognitive neuroscience and artificial intelligence. *Nature neuroscience*, 25(1):116–126,
587 2022. doi: 10.1038/s41593-021-00962-x. Publisher: Nature Publishing Group US New York.

588 Luca Moschella, Valentino Maiorca, Marco Fumero, Antonio Norelli, Francesco Locatello, and
589 Emanuele Rodolà. Relative representations enable zero-shot latent space communication. *arXiv
590 preprint arXiv:2209.15430*, 2022.

591
592 Valentino Maiorca, Luca Moschella, Antonio Norelli, Marco Fumero, Francesco Locatello, and
593 Emanuele Rodolà. Latent space translation via semantic alignment. *Advances in Neural Information
Processing Systems*, 36, 2024.

-
- 594 Antonio Norelli, Marco Fumero, Valentino Maiorca, Luca Moschella, Emanuele Rodola, and
595 Francesco Locatello. Asif: Coupled data turns unimodal models to multimodal without training.
596 *Advances in Neural Information Processing Systems*, 36, 2024.
- 597 Ellie Pavlick. Symbols and grounding in large language models. *Philosophical Transactions of the*
598 *Royal Society A*, 381(2251):20220041, 2023.
- 600 Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. The platonic representation
601 hypothesis. *arXiv preprint arXiv:2405.07987*, 2024.
- 602 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
603 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European*
604 *conference on computer vision*, pages 740–755. Springer, 2014.
- 605 Nikolaus Kriegeskorte, Marieke Mur, and Peter A Bandettini. Representational similarity analysis-
606 connecting the branches of systems neuroscience. *Frontiers in systems neuroscience*, 2:4, 2008a.
- 607 Philipp Kaniuth and Martin N Hebart. Feature-reweighted rsa: A method for improving the fit
608 between computational models, brains, and behavior. *bioRxiv*, 2021.
- 609 Talia Konkle and George A Alvarez. A self-supervised domain-general learning framework for
610 human ventral stream representation. *Nature Communications*, 13(1):1–12, 2022.
- 611 Kendrick Kay, Jacob S Prince, Thomas Gebhart, Greta Tuckute, Jingyang Zhou, Thomas Naselaris,
612 and Heiko Schutt. Disentangling signal and noise in neural responses through generative modeling.
613 *bioRxiv*, pages 2024–04, 2024.
- 614 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
615 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
616 *systems*, 30, 2017.
- 617 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
618 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 619 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language
620 understanding by generative pre-training. 2018.
- 621 Thomas A Carlson, Ryan A Simmons, Nikolaus Kriegeskorte, and L Robert Slevc. The emergence
622 of semantic meaning in the ventral temporal pathway. *Journal of cognitive neuroscience*, 26(1):
623 120–131, 2014.
- 624 Christiane Fellbaum. Wordnet. In *Theory and applications of ontology: computer applications*, pages
625 231–243. Springer, 2010.
- 626 Sara F Popham, Alexander G Huth, Natalia Y Bilenko, Fatma Deniz, James S Gao, Anwar O Nunez-
627 Elizalde, and Jack L Gallant. Visual and linguistic semantic representations are aligned at the
628 border of human visual cortex. *Nature neuroscience*, 24(11):1628–1636, 2021.
- 629 Jerry Tang, Meng Du, Vy A Vo, Vasudev Lal, and Alexander G Huth. Brain encoding models based on
630 multimodal transformers can transfer across language and vision. *arXiv preprint arXiv:2305.12248*,
631 2023b.
- 632 Evelina Fedorenko, Steven T Piantadosi, and Edward AF Gibson. Language is primarily a tool for
633 communication rather than thought. *Nature*, 630(8017):575–586, 2024.
- 634 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
635 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
636 *ence on computer vision and pattern recognition*, pages 10684–10695, 2022. arXiv: 2112.10752
637 [cs.CV].
- 638 Colin Conwell and Tomer Ullman. Testing relational understanding in text-guided image generation.
639 *arXiv preprint arXiv:2208.00005*, 2022.
- 640 Yutaro Yamada, Yihan Bao, Andrew K Lampinen, Jungo Kasai, and Ilker Yildirim. Evaluating spatial
641 understanding of large language models. *arXiv preprint arXiv:2310.14540*, 2023.

648 Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and
649 why vision-language models behave like bags-of-words, and what to do about it?, 2023. URL
650 <https://arxiv.org/abs/2210.01936>.
651

652 Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Can-
653 dace Ross. Winoground: Probing vision and language models for visio-linguistic compositionality.
654 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
655 5238–5248, 2022.

656 Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. Measuring
657 and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350*,
658 2022.

659 Lorenzo Bertolini, Julie Weeds, and David Weir. Testing large language models on composi-
660 tionality and inference with phrase-level adjective-noun entailment. In *Proceedings of the*
661 *29th International Conference on Computational Linguistics*, pages 4084–4100, Gyeongju, Re-
662 public of Korea, October 2022. International Committee on Computational Linguistics. URL
663 <https://aclanthology.org/2022.coling-1.359>.
664

665 Marc D Hauser, Noam Chomsky, and W Tecumseh Fitch. The faculty of language: what is it, who
666 has it, and how did it evolve? *science*, 298(5598):1569–1579, 2002. doi: 10.1126/science.298.
667 5598.1569. Publisher: American Association for the Advancement of Science.

668 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
669 contrastive learning of visual representations. In *International Conference on Machine Learning*,
670 pages 1597–1607. PMLR, 2020. arXiv preprint arXiv:2002.05709.
671

672 Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised
673 learning via redundancy reduction. *arXiv preprint arXiv:2103.03230*, 2021.

674 Priya Goyal, Mathilde Caron, Benjamin Lefaudeaux, Min Xu, Pengchao Wang, Vivek Pai, Mannat
675 Singh, Vitaliy Liptchinsky, Ishan Misra, Armand Joulin, and others. Self-supervised pretraining of
676 visual features in the wild. *arXiv preprint arXiv:2103.01988*, 2021a.

677

678 Yangguang Li, Feng Liang, Lichen Zhao, Yufeng Cui, Wanli Ouyang, Jing Shao, Fengwei Yu,
679 and Junjie Yan. Supervision exists everywhere: A data efficient contrastive language-image
680 pre-training paradigm. In *International Conference on Learning Representations*, 2022. URL
681 <https://openreview.net/forum?id=zqliJkNk3uN>.

682 Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas Steiner, Daniel Keysers, Alexander Kolesnikov,
683 and Lucas Beyer. Lit: Zero-shot transfer with locked-image text tuning. In *Proceedings of the*
684 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18123–18133, 2022.

685

686 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
687 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
688 pages 248–255. Ieee, 2009. doi: 10.1109/CVPR.2009.5206848.

689 Kalanit Grill-Spector and Kevin S Weiner. The functional architecture of the ventral temporal cortex
690 and its role in categorization. *Nature Reviews Neuroscience*, 15(8):536–548, 2014.

691

692 Bria Long, Chen-Ping Yu, and Talia Konkle. Mid-level visual features underlie the high-level
693 categorical organization of the ventral stream. *Proceedings of the National Academy of Sciences*,
694 115(38):E9015–E9024, 2018. doi: 10.1073/pnas.1719616115. Publisher: National Acad Sciences.

695

696 Hans Op de Beeck et al. Category trumps shape as an organizational principle of object space in the
697 human occipitotemporal cortex. *Journal of Neuroscience*, 43(16):2960–2972, 2023.

698

699 Adva Shoham, Rotem Broday-Dvir, Rafael Malach, and Galit Yovel. The organization of high-level
700 visual cortex is aligned with visual rather than abstract linguistic information. *bioRxiv*, pages
701 2024–11, 2024.

702

703 Steven Michael Frankland. *Man bites dog: The representation of structured meaning in left-mid*
superior temporal cortex. PhD thesis, Harvard University, 2015.

702 Jake Quilty-Dunn, Nicolas Porot, and Eric Mandelbaum. The best game in town: The reemergence of
703 the language-of-thought hypothesis across the cognitive sciences. *Behavioral and Brain Sciences*,
704 46:e261, 2023.

705
706 R Quian Quiroga, Leila Reddy, Gabriel Kreiman, Christof Koch, and Itzhak Fried. Invariant visual
707 representation by single neurons in the human brain. *Nature*, 435(7045):1102–1107, 2005. doi:
708 10.1038/nature03687. Publisher: Nature Publishing Group.

709 Priya Goyal, Quentin Duval, Jeremy Reizenstein, Matthew Leavitt, Min Xu, Benjamin Lefaudeux,
710 Mannat Singh, Vinicius Reis, Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Ishan Misra.
711 Vissl. <https://github.com/facebookresearch/vissl>, 2021b.

712
713 Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori,
714 Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali
715 Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL [https://doi.org/10.5281/
716 zenodo.5143773](https://doi.org/10.5281/zenodo.5143773). If you use this software, please cite it as below.

717 Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets
718 language-image pre-training. *arXiv preprint arXiv:2112.12750*, 2021.

719
720 Ross Wightman. Pytorch image models. [https://github.com/rwightman/
721 pytorch-image-models](https://github.com/rwightman/pytorch-image-models), 2019.

722 Nils Reimers and Iryna Gurevych. Making monolingual sentence embeddings multilingual using
723 knowledge distillation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural
724 Language Processing*. Association for Computational Linguistics, 11 2020. URL [https://
725 arxiv.org/abs/2004.09813](https://arxiv.org/abs/2004.09813).

726
727 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
728 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers:
729 State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.

730
731 Matthew Honnibal, Ines Montani, Sofie Van Landeghem, Adriane Boyd, et al. spacy: Industrial-
732 strength natural language processing in python. 2020.

733
734 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
735 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
736 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
737 12:2825–2830, 2011.

738
739 Jacob S Prince, Ian Charest, Jan W Kurzwawski, John A Pyles, Michael J Tarr, and Kendrick N Kay.
740 Improving the accuracy of single-trial fmri response estimates using glmsingle. *Elife*, 11:e77599,
741 2022.

742
743 Leyla Tarhan and Talia Konkle. Reliability-based voxel selection. *NeuroImage*, 207:116350, 2020.

744
745 Talia Konkle and George A Alvarez. Beyond category-supervision: instance-level contrastive learning
746 models predict human visual system responses to objects. *bioRxiv*, 2021.

747
748 Dimitris Achlioptas. Database-friendly random projections. In *Proceedings of the twentieth ACM
749 SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 274–281, 2001.
750 doi: 10.1145/375551.375608.

751
752 Nikolaus Kriegeskorte, Marieke Mur, Douglas A Ruff, Roozbeh Kiani, Jerzy Bodurka, Hossein
753 Esteky, Keiji Tanaka, and Peter A Bandettini. Matching categorical object representations in
754 inferior temporal cortex of man and monkey. *Neuron*, 60(6):1126–1141, 2008b.

755
756 Mark D Lescroart, Dustin E Stansbury, and Jack L Gallant. Fourier power, subjective distance, and
757 object categories all provide plausible models of bold responses in scene-selective visual areas.
758 *Frontiers in computational neuroscience*, 9:135, 2015.

759
760 Submitter Mark Lescroart. Unique variance found by variance partitioning is superior to total variance
761 explained as a model comparison metric. *Cognitive Computational Neuroscience*, 2017.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A APPENDIX

A.1 MATERIALS AND METHODS

MODEL SELECTION

The models we test in this experiment consist of vision-only, language-only, and vision-language deep neural networks. These networks are supplemented with word-level models that operate over either the tokenized inputs of individual words (i.e. GLOVE) or directly over individual words (i.e. the count-vectorizing models).

Our sample of vision-only models consist of N=11 purely self-supervised visual contrastive-learning models from the VISSL model zoo (Goyal et al., 2021b). Our sample of hybrid vision-language models consists of N=10 models from the OpenAI CLIP, OpenCLIP, SLIP, and PyTorch-Image-Models repository (Radford et al., 2021; Ilharco et al., 2021; Mu et al., 2021; Wightman, 2019). Our sample of large language models (plus GLOVE) consists of N=12 models from Hugging Face and the SBERT repositories (Reimers and Gurevych, 2020; Wolf et al., 2019). To derive CLIP similarity scores for linguistic prompts, we use OpenAI CLIP’s ResNet50 backbone (Radford et al., 2021). For part-of-speech extraction and word vectorizing operations, we use spaCy’s English (BERT-Based) Large TRF model (Honnibal et al., 2020) and scikit-learn’s CountVectorizer (Pedregosa et al., 2011), respectively.

HUMAN FMRI DATA

The Natural Scenes Dataset (Allen et al., 2022) contains measurements of 73,000 unique stimuli from the Microsoft Common Objects in Context (COCO) dataset (Lin et al., 2014) at high resolution (7T field strength, 1.33s TR, $1.8mm^3$ voxel size). In this analysis, we focus on the brain responses to 1000 COCO stimuli that overlapped between subjects, and limit analyses to the 4 subjects (subjects 01, 02, 05, 07) for whom all 3 image repetitions are available for the overlapping images. The 3 image repetitions were averaged to yield the final voxel-level response values in response to each stimulus. All responses were estimated using a custom GLM toolbox (“GLMsingle” (Prince et al., 2022)), which was applied during the preprocessing of NSD time-series data, featuring optimized denoising and regularization procedures, to accurately measure changes in brain activity in response to each experimental stimulus.

Voxel Selection Procedure

To achieve a reasonable signal-to-noise ratio (SNR) in our target data, we implement a reliability-based voxel selection procedure (Tarhan and Konkle, 2020) to subselect voxels containing stable structure in their responses. Specifically, we use the NCSNR (“noise ceiling signal-to-noise ratio”) metric computed for each voxel as part of the NSD metadata (Allen et al., 2022) for our reliability metric. In this analysis, we include only those voxels with $NCSNR > 0.2$.

After filtering voxels based on their NCSNR, we then filtered voxels based on region-of-interest (ROI). In our main analyses, we focus on voxels within the early visual and occipitotemporal cortices (EVC and OTC, respectively). For OTC, we first considered voxels within a liberal mask of the visual system (“nsdgeneral” ROI, see (Allen et al., 2022) for details). Next we selected the subset within either the mid-to-high ventral or mid-to-high lateral ROIs (“streams” ROIs). Then, we included all voxels from 11 category-selective ROIs (face, body, word, and scene ROIs, excluding RSC) with a t-contrast statistic $t > 1$; while many of these voxels were already contained in the streams ROIs, this ensures that these regions were included in the larger scale OTC sector. The number of OTC voxels included were 8,088 for subject 01, 7,528 for subject 02, 8,015 for subject 05, and 5,849 for subject 07, for a combined total of 29,480 voxels.

The EVC ROI encapsulates the ventral and dorsal aspects of areas V1, V2, and V3, as well as area hV4 (see (Allen et al., 2022) for details on ROI localization). To define the EVC ROI for each subject, we again first isolated voxels within the “nsdgeneral” ROI, and then selected for analyses any voxels that both fell within one of the early visual regions listed above, and that exceeded the NCSNR threshold of 0.2. This procedure yielded a total of 4,657 voxels for subject 01, 3,757 voxels for subject 02, 3,661 voxels for subject 05, and 3,251 voxels for subject 07.

Noise Ceilings

810 To contextualize model performance results, we estimated noise ceilings for each of target brain ROIs.
811 These noise ceilings indicate the maximum possible performance that can be achieved given the level
812 of measurement noise in the data. Importantly, our noise ceiling estimates refer to within-subject
813 representational dissimilarity matrices (RDMs), where noise reflects trial-to-trial variability in a
814 given subject. This stands in contrast to more conventional group-level representational dissimilarity
815 matrices (Kriegeskorte et al., 2008a), where noise reflects variability across subjects. To estimate
816 within-subject noise ceilings, we applied a novel method based on generative modeling of data’s
817 signal and noise characteristics (GSN; Kay et al., 2024).

818 This method estimates, for a given ROI, multivariate Gaussian distributions characterizing the signal
819 and the noise under the assumption that observed responses can be characterized as sums of samples
820 from the signal and noise distributions. A post-hoc scaling is then applied to the signal distribution
821 such that the signal and noise distributions generate accurate matches to the empirically observed
822 reliability of RDMs across independent splits of the experimental data. Noise ceilings are estimated
823 using Monte Carlo simulations in which a noiseless RDM (generated from the estimated signal
824 distribution) is correlated with RDMs constructed from noisy measurements (generated from the
825 estimated signal *and* noise distributions).

826 FEATURE MAPPING METHODS

827 **Feature Extraction Procedure**

828 For each of our candidate DNN models, we extract features in response to each of our probe stimuli
829 at each distinct layer of the network. At the end of our feature extraction procedure, for each
830 model and each model layer, we arrive at a feature matrix of dimensionality number-of-images x
831 number-of-flattened-features.

832 When using large language models, we obtain a single embedding from each model layer by averaging
833 the 5 individual embeddings provided for each image. For both BERT-based and GPT-based models,
834 we hook directly into the transformer layers to extract hidden representations. The tokenized captions
835 are passed as a single tensor dataset, with token sequences padded to the maximum length. Attention
836 masks are applied to ignore padding tokens. We apply no aggregation or pooling beyond that which is
837 instantiated by submodules or functions in the feedforward pass. While non-standard, this approach
838 allows us to apply a single, consistent across architectures, and was validated empirically with
839 cross-reference to the embeddings obtained from standard Huggingface transfer-learning pipelines.

840 When using word-vectorizing (count) models (which do not contain layers), we obtain a single
841 embedding by summing the word counts across the 5 image captions.

842 **Classical RSA (cRSA)**

843 To compute the classical representational similarity (cRSA) score (Kriegeskorte et al., 2008a) for a
844 single layer, we used the following procedure: First, we split the 1000 images into two sets of 500
845 (a training set, and a testing set). Using the training set of images, we compute the representational
846 similarity matrices (RSMs) of each model layer (500 x 500 x number-of-layers) using Pearson
847 correlation distance metric. We then compare each layer’s RSM to the brain RSM, also using Pearson
848 similarity, and identify the layer with the highest correlation as the model’s most brain-predictive
849 layer. Finally, using the held-out test set of 500 images, we compute that target layer’s RDM and
850 correlate it with the brain RDM. This score serves as the overall cRSA score for the target model.

851 **Voxelwise Encoding RSA (eRSA)**

852 To arrive at a voxelwise encoding representational similarity (eRSA) score (Konkle and Alvarez,
853 2021; Kaniuth and Hebart, 2021) for a single model, the overall procedure was similar to that of
854 cRSA, but with the addition of an intermediate encoding procedure wherein layerwise model features
855 were fit to each voxel’s response profile.

856 The first step in the encoding procedure is the dimensionality reduction of model feature maps. We
857 perform this step for two reasons: (a) the features extracted from various deep neural networks can
858 sometimes be massive (the first convolutional layer of VGG16, for example, yields a flattened feature
859 matrix with 3.2 million dimensions per image); (b) the same dimensionality reduction procedure
860 applied to all layers ensures that the explicit degrees of freedom across model layers is constant.
861 To reduce dimensionality, we apply the scikit-learn implementation of sparse random projection

864 (Pedregosa et al., 2011). This procedure relies on the Johnson-Lindenstrauss (JL) lemma (Achlioptas,
865 2001), which takes in a target number of samples and an epsilon distortion parameter, and returns
866 the number of random projections necessary to preserve the euclidean distance between any two
867 points up to a factor of $1 \pm \epsilon$. (Note that this is a general formula; no brain data enter into this
868 calculation). In our case, with the number of samples set to 1000 (the total number of images) and an
869 epsilon distortion of 0.1, the Johnson-Lindenstrauss procedure yields a target dimensionality of 5920
870 projections.

871 After computing this target dimensionality, we then proceed to compute the sparse random projection
872 for each layer of our target DNN. The sparse random projection matrix consists of zeros and sparse
873 ones of nearly orthogonal dimensions, and the layerwise feature maps are then projected onto this
874 matrix by taking the dot product between them. The output of the procedure is a reduced layerwise
875 feature space of size of 1000 images x 5920 dimensions with a preserved representational geometry.
876 Note that in cases where the number of features is less than the number of projections suggested
877 by the JL lemma, the original feature map is effectively upsampled through the random projection
878 matrix, again yielding a matrix of 1000 x 5920 dimensions.

879 We compute our encoding model for each voxel as a weighted combination of these 5920 dimensions,
880 using brain data from our training set of 500 images. (We note that while the number of dimensions
881 needed for only 500 images would be only $D=5326$ according to the JL lemma, adding extra
882 dimensions will only preserve the geometry with nominally less distortion than the epsilon provided,
883 and does not affect the results). The fitting procedure for each voxel leverages SciKit-Learn’s
884 cross-validated ridge regression function, a hyperefficient regression method that uses generalized
885 cross-validation to provide a LOOCV prediction per image (per output). This fit was computed over
886 a logarithmic range of alpha penalty parameters ($1e^{-1}$ to $1e^{-7}$), to identify each voxel’s optimal
887 alpha parameter. We modified the RidgeCV function in order to select the best alpha using Pearson
888 correlation as a score function (the same score function we use to evaluate the model at large), and
889 to parallelize a slow loop for efficiency. This yielded a set of encoding weights for each voxel
890 (number-of-voxels x 5920 reduced-feature-dimensions).

891 Next, with these encoding weights and the 500 training images, we compute the predicted response of
892 every voxel to each image, and compute the corresponding *predicted* RSM using Pearson correlation.
893 After computing each layer’s RSA similarity value via Pearson correlation between the layer-predicted
894 RDM and the target brain RSM, we again select the most predictive layer on the basis of results from
895 the training set and compute this layer’s RSA correspondence to the brain data using the held-out set
896 of 500 test images. This test score from each model’s most-predictive layer serves as the final eRSA
897 score for each model.

898 We emphasize that this method contrasts with popular practices in primate and mouse benchmarking,
899 which treat predictivity of unit-level univariate response profiles as the key measure. Because fMRI
900 affords more systematic spatial sampling over the cortex, rather than taking the aggregate of single
901 voxel fits as our key measure, we choose to treat the population representational geometry over each
902 ROI as our critical target for prediction. This multi-voxel similarity structure provides different kinds
903 of information about the format of population-level coding than do individual units (Kriegeskorte
904 et al., 2008b). Computing the eRSA metric does, however, yield individual voxelwise encoding
905 models, the individual predictive accuracies of which we register and have available in addition to the
906 cRSA and eRSA scores for future analysis.

907

908 A.2 SUPPLEMENTARY ANALYSES

909

910 Here, we report the results of two supplementary analyses in the comparison of pure-vision and pure-
911 language models. In the first, we assess the OTC-predictivity of the randomly-initialized architectures
912 corresponding to all our of unimodal (pure-vision or pure-language) models. In the second, we use
913 a variance partitioning analysis (Lescroart et al., 2015; Lescroart, 2017) to assess the amount of
914 brain-predictivity both shared between and unique to either modality.

915 Note that we have opted here to use the average voxelwise encoding scores that are implicit to our
916 eRSA analysis, and report here scores only in the most-relevant contrast of occipitotemporal cortex
917 (OTC). This is for multiple reasons: one) for simplicity; two) as a demonstration that the results
derived from these scores (implicit to the eRSA analysis) largely concord with the RSA metrics used

in the main analysis; and three) for better synchrony with variance partitioning methods (which lend themselves naturally to regression-based metrics, but less so to RSA metrics).

RANDOMLY-INITIALIZED MODELS

An important question in the comparison of pure-vision and pure-language models is the question of just how much we can attribute their divergent behaviors to differences in the modality of their training data alone. After all, pure-vision and pure-language models differ not only in the modality of their training data, but in the format of their inputs (tokenized strings versus pixels), their architectures (e.g. CNN versus ViT), and their training task (e.g. next-word prediction versus contrastive learning over various image augmentation regimes). While in our main analysis, we have attempted to try and abstract over these other differences by assessing a relatively diverse sample of each model type, another way we can probe the differences between them (at least at the level of architecture) is by running our same brain-prediction pipeline on the randomly-initialized versions of each.

In the case of this particular comparison between vision and language models, this comparison of trained versus randomly-initialized models has the added benefit of contributing directly to the logic of our interpretation. If, as we argue in the main body of the work, the ‘language’ in language-models does not yield altogether ‘language-unique’ or ‘human language-like’ representation (in the sense of having complicated structure that extends beyond the co-occurrence statistics of common nouns, and not higher-order compositional meaning), then randomly-initialized versions of the language models should not suffer as substantial a decrease in performance without training as vision models will.

Indeed, we find this to be the case. In line with the idea that there is really not that much “language” structure in the language models, randomly-initialized language models only suffered a relatively minor drop in accuracy compared to their pretrained counterparts, with mean voxelwise encoding scores of $r_{Pearson} = 0.329$ [0.327, 0.331] from trained weights and 0.276 [0.273, 0.279] for untrained (randomly-initialized) weights. We can contrast this with the far more substantive drop with randomly initialized vision models: with $r_{Pearson} = 0.341$ [0.319, 0.367] for trained models and 0.156 [0.115, 0.193] for untrained models.

VARIANCE PARTITIONING ANALYSIS

Variance partitioning (Lescroart et al., 2015; Lescroart, 2017) is an analysis technique that can be used to determine how much of the variance explained in a multivariate regression model is shared between or unique to the predictors. Here, we deploy this technique to partition the variance in OTC responses explained by our unimodal models – in this case, the most-brain-predictive models from each class (SEER-RegNet64GF for pure-vision; SBERT-MiniLM-L12-PML for pure-language). Variance partitioning in the case of our two predictors (pure-vision and pure-language) involves fitting a total of three regressions: two with each predictor alone, and one regression with the two predictors combined. In this case, the results of those regressions – in units both of $r_{Pearson}$ and ($r_{Pearson}^2$) – are as follows:

1. OTC ~ Language (Alone): 0.354 (0.125)
2. OTC ~ Vision (Alone): 0.383 (0.147)
3. OTC ~ Language + Vision: 0.402 (0.161)

This produces the following unique and shared variances (in units of $r_{Pearson}^2$):

1. Unique to Language: 0.014
2. Unique to Vision: 0.036
3. Shared Language-Vision: 0.11

As we can see, the vast majority of the variance in brain prediction is shared between the pure-vision and pure-language models, though pure-vision has a slight advantage over pure-language in terms of unique variance.