

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  

# NEURO-SYMBOLIC DECODING OF NEURAL ACTIVITY

Anonymous authors

Paper under double-blind review

## ABSTRACT

We propose NEURONA, a modular neuro-symbolic framework for fMRI decoding and concept grounding in neural activity. Leveraging image- and video-based fMRI question-answering datasets, NEURONA learns to decode interacting concepts from visual stimuli from patterns of fMRI signals, integrating symbolic reasoning and compositional execution with fMRI grounding across brain regions. We demonstrate that incorporating structure into the decoding pipeline improves both decoding accuracy and generalization performance. **NEURONA shows that modeling the compositional structure of concepts through hierarchical predicate-argument dependencies enables more precise decoding from fMRI, highlighting neuro-symbolic frameworks as promising tools for neural decoding.**

## 1 INTRODUCTION

A long-standing hypothesis in cognitive science, the Language of Thought (LoT) hypothesis (Fodor, 1975), proposes that human cognition operates over structured, symbolic representations that compose systematically. Rather than storing concepts as isolated units, the brain is thought to organize knowledge into compositional structures—such as predicates and their arguments—that enable flexible and generalizable reasoning. **To test whether such structures can improve neural decoding**, we study concept grounding in functional magnetic resonance imaging (fMRI), with the goal of predicting symbolic concepts (e.g., person, baseball-bat, and holding) from patterns of neural activity (Mitchell et al., 2008). **This alignment offers insights into whether incorporating structural priors can enable more accurate, precise, and generalizable neural decoding.**

There has been vast literature on concept grounding in the past decades, with several influential works studying how concepts are organized across the cortex (Mitchell et al., 2008; Palatucci et al., 2009; Huth et al., 2016; Pereira et al., 2018). Recent advances in machine learning has enabled growing efforts toward data-driven approaches to concept grounding. However, most large-scale fMRI decoding studies focus on isolated concepts or holistic stimulus reconstruction (Nishimoto et al., 2011; Naselaris et al., 2011; Chen et al., 2023a; Takagi & Nishimoto, 2023; Scotti et al., 2023; Chen et al., 2023b), **leaving open the question of how the to decode relational meaning between interacting visual concepts.** Specifically, we ask, does the decoding of relational concepts (e.g., holding) improve by accounting for the systematic combination of their constituent arguments (e.g., person and baseball-bat) across multiple brain regions?

To explore these questions, we leverage rich data from image- and video-based fMRI datasets, which naturally encode complex semantic and compositional structure. Naturalistic stimuli such as images and videos often involve multiple interacting concepts (e.g., a person holding a baseball-bat), **making them well-suited for probing how to decode entities and their relations from neural activity.** Hence, we build challenging fMRI-question-answering (fMRI-QA) datasets based on BOLD5000 (Chang et al., 2019) and CNNeuroMod (Gifford et al., 2024; Boyle et al., 2023), with the goal of learning concept grounding and **improving decoding accuracy based on computational hypotheses on composition.**

However, neither simple linear models nor purely end-to-end neural decoding models are sufficient for solving this task. Linear models lack the capacity to capture interactions between multiple interacting components, while large neural decoders (e.g., those with language model backbones) tend to encode stimuli holistically, without explicitly modeling modular concepts or their relationships. To address these limitations, we adopt a neuro-symbolic approach that integrates the compositionality of symbolic systems with the expressivity of neural networks for fMRI-QA: each query is decomposed

054 into a symbolic expression composing concepts, and neural activities in the brain are routed through  
 055 corresponding concept modules (implemented as neural networks) to answer the given query.  
 056

057 Specifically, we build upon the general neuro-symbolic framework of the Logic-Enhanced Foundation  
 058 Model (LEFT) (Hsu et al., 2023), and adapt it for the domain of fMRI-based question answering,  
 059 enabling the use of QA supervision to learn disentangled concept groundings. Crucially, we introduce  
 060 the incorporation of various *compositional priors* into the model, by defining candidate entities in  
 061 neural activity and specifying how they are composed based on the symbolic expressions. From  
 062 this paradigm, we propose **NEURONA**, a **NEURO**-symbolic framework for decoding in Neural  
 063 Activity, which integrate symbolic reasoning and compositional execution with fMRI grounding, and  
 064 significantly improves decoding accuracy compared to prior works.

065 With NEURONA, we find that incorporating structural priors to explicitly *guide* the concept grounding  
 066 process, such as enforcing hierarchical predicate-argument dependencies (e.g., grounding for  
 067 holding conditioned on the grounding of `baseball`-`bat`)—notably improves decoding accuracy  
 068 on fMRI-QA tasks. These priors guide the model to compose high-level relational concepts  
 069 from their constituent entity groundings, showing that relational meaning is better predicted across  
 070 multiple co-activated brain networks via its arguments, rather than localized to a single region or to  
 071 multiple regions without guidance.

072 Evaluating concept grounding in fMRIs is inherently challenging due to the lack of direct supervision:  
 073 there is no ground truth mapping from abstract concepts to specific brain regions. Instead, we  
 074 report experiments on BOLD5000 and CNNeuroMod fMRI-QA datasets, which demonstrate that  
 075 our neuro-symbolic framework significantly outperforms baseline neural decoding methods, and  
 076 importantly, exhibits strong generalization to unseen compositional queries. Notably, ablation studies  
 077 with NEURONA highlight the importance of encoding hierarchical structure: conditioning predicate  
 078 grounding modules on the regions associated with their subject and object arguments consistently  
 079 yields large performance gains in decoding of neural activity.

## 080 2 RELATED WORKS

081 **Visual decoding from fMRI.** Reconstructing visual content from fMRI signals has become a central  
 082 research focus of works in the field, with many approaches leveraging state-of-the-art generative  
 083 backbones for the task, following early studies (Thirion et al., 2006; Miyawaki et al., 2008; Kay  
 084 et al., 2008; Naselaris et al., 2009; Nishimoto et al., 2011). Takagi et al. demonstrated that a pre-  
 085 trained diffusion model can reconstruct high-resolution images from fMRI (Takagi & Nishimoto,  
 086 2023). MinD-Vis uses masked brain modeling with a diffusion model for semantically faithful image  
 087 generation (Chen et al., 2023a). MindEye projects fMRI into a CLIP embedding space and applies  
 088 a diffusion prior for pixel-level synthesis (Scotti et al., 2023). Extending to video, MinD-Video  
 089 and NeuroCLIP incorporate spatiotemporal masked modeling and keyframe-perception flow cues,  
 090 respectively, into diffusion-based reconstruction (Chen et al., 2023b; Gong et al., 2024). These visual  
 091 reconstruction works focus on recovering stimulus appearance from neural data; in contrast, our work  
 092 addresses a distinct goal of concept grounding: rather than generating pixel-level images or videos,  
 093 we aim to predict modular concepts and their relationships from neural activity.

094 **Concept grounding.** Several influential works have focused on how semantic information is  
 095 organized across the cortex. As a representative work, Huth et al. used voxel-wise encoding models  
 096 with natural narrative stimuli to construct a semantic atlas, showing that different semantic domains  
 097 selectively ground to distinct brain regions (Huth et al., 2016). Mitchell et al. predicted fMRI  
 098 patterns for concrete nouns using corpus-derived semantic features, showing generalization to unseen  
 099 words (Mitchell et al., 2008). SOC enables zero-shot decoding by mapping fMRI to semantic  
 100 codes and recognizing novel object categories (Palatucci et al., 2009). Pereira et al. introduced a  
 101 general decoder that maps fMRI into a shared semantic space, enabling generalization from limited  
 102 data (Pereira et al., 2018). Beyond semantic mapping, several studies also explored how concepts  
 103 are organized in the brain (Frankland & Greene, 2015; Eichenbaum, 2001). There is converging  
 104 evidence that certain brain regions support invariant concept and rule representations. For example,  
 105 the prefrontal cortex—spanning networks such as the dorsal attention and default mode networks—  
 106 and the medial temporal lobe have been implicated in abstract concept and rule processing (Quiroga  
 107 et al., 2005; Rey et al., 2015; Tian et al., 2024; Dijksterhuis et al., 2024). Additionally, single-neuron  
 108 recordings in the human prefrontal cortex have revealed neurons that encode abstract task rules

108 independently of sensory or motor details (Mian et al., 2014). In contrast to these prior works,  
 109 our approach emphasizes **compositional grounding in the neural decoding process**, with focus on  
 110 **functional instead of representational compositionality**, where we explicitly model not only individual  
 111 concepts, but also the relationships between them. **This modeling allows us to investigate how**  
 112 **relational concepts can be more accurately decoded through guided grounding across brain networks.**

113  
**fMRI-question answering.** Recent works have explored using fMRI data for question-answering  
 114 by integrating large vision-language models (VLMs). These methods typically map neural activity to  
 115 visual embeddings, then generate answers using pre-trained VLMs. For example, SDRecon (Takagi  
 116 & Nishimoto, 2023) projects fMRI signals into BLIP (Li et al., 2022) embeddings for captioning;  
 117 BrainCap (Ferrante et al., 2023) maps fMRI to GIT (Wang et al., 2022) features for visual description;  
 118 and UMBRAE (Xia et al., 2024) aligns fMRI to multimodal embeddings with subject-specific  
 119 tokenization and answers questions via LLaVA (Liu et al., 2023). These methods commonly use  
 120 BLEU scores (Papineni et al., 2002) to measure alignment with ground-truth text, but they do not  
 121 explicitly verify whether the predicted answer captures exact concepts or relational structure. In  
 122 contrast, our framework grounds fMRI signals to modular concepts before performing structured  
 123 reasoning, enabling precise, accurate, and generalizable question answering.

### 124 3 METHOD

#### 125 3.1 NEURO-SYMBOLIC FRAMEWORK

126 We introduce NEURONA as a neuro-symbolic framework for concept grounding and decoding in  
 127 neural activity. Neuro-symbolic models are a class of methods that decompose queries into symbolic  
 128 expressions containing concepts, and then differentiably execute those expressions over input data,  
 129 using learned concept grounding modules to perform a variety of downstream tasks (Yi et al., 2018;  
 130 Mao et al., 2019; Hsu et al., 2023; Mao et al., 2025). Each symbolic concept (e.g., person,  
 131 holding) is associated with a small neural network that maps entity-centric representations from  
 132 input data to a predicted semantic signal, enabling the learning of intermediate concept grounding  
 133 from weak supervision of reasoning tasks. Execution is conducted via differentiable functions  
 134 combining concept grounding modules, which enables end-to-end training.

135 In this work, we build our model, NEURONA, based on the Logic-Enhanced Foundation Model  
 136 (LEFT) (Hsu et al., 2023). LEFT is a general neuro-symbolic framework that unifies grounding  
 137 and reasoning via a differentiable executor for logic programs. It is designed to support concept  
 138 grounding across various visual domains (e.g., 2D images, 3D scenes), notably, where the relevant  
 139 entities are known a priori, such as objects in a room. In contrast, our setting introduces a unique  
 140 challenge: concept grounding from fMRI signals, where the relevant neural “entities” (i.e., brain  
 141 regions) are not predefined, and must instead be inferred as part of the learning process. **This makes**  
 142 **our setting significantly more challenging than prior works: we are testing hypotheses for what these**  
 143 **entities should be, as well as learning how they best compose to precisely decode neural activity.**

144 Concretely, we frame concept grounding in neural activity as a weakly supervised fMRI-question  
 145 answering (fMRI-QA) task. Each input consists of two elements (See Figure 1). The first is a visual  
 146 stimuli  $v$  that is parsed into a symbolic query that describe possible concepts in  $v$ . The second is  
 147 an fMRI recording  $f \in \mathbb{R}^{N \times T}$  of the neural activity recorded when the stimuli was viewed, where  
 148  $N$  is the number of fine-grained brain networks, and  $T$  is the number of time steps. The goal is to  
 149 predict an answer  $a$ , either as a Boolean label (e.g., True/False) or as a classification over a concept  
 150 vocabulary. Crucially, no supervision is provided on which brain regions are relevant to each concept.

151 To enable modular concept grounding to neural activity, we first map the fine-grained networks  $N$  to  
 152  $P$  functional networks defined by an atlas. We experiment with Yeo-7, Yeo-17 (Yeo et al., 2011),  
 153 DiFuMo-64, DiFuMo-128 (Dadi et al., 2020), and Schaefer-100 (Schaefer et al., 2018) atlases, where  
 154  $P$  is 7, 17, 64, 128, 100 respectively. While all atlases yield consistently strong performance with  
 155 NEURONA, we report decoding accuracy with the Yeo-17 atlas in the main text, as performance is  
 156 highest. Experiment results across all atlases can be found in Appendix A. Hence, we have  $P = 17$   
 157 resulting network-specific fMRI signals, which we then encode to form a unified set  $\mathcal{E}$  of embeddings  
 158  $\{e_1, \dots, e_P\}$ , to form parcellation embeddings. From these base embeddings, we propose hypotheses  
 159 of candidate entities from which concepts can be mapped to. Here, neural entities precisely are

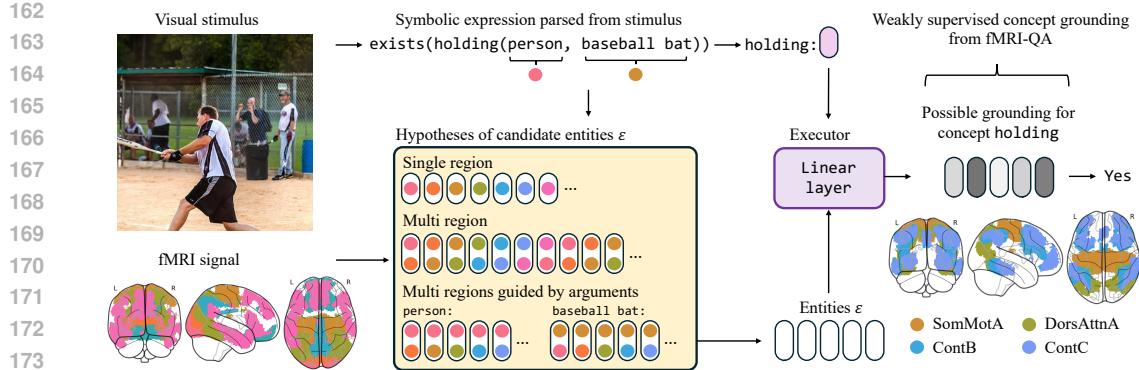


Figure 1: NEURONA is a neuro-symbolic framework for neural decoding that parses visual stimuli into symbolic expressions and fMRI signals into candidate entities, optionally guided by the concepts in the expression. These composed concepts are then grounded to the entities via a linear layer for question answering; the predicted answer provides weak supervision for the grounding process.

region-specific fMRI signals derived from functional parcellations (e.g., Yeo-17, DiFuMo-128 atlases) that serve as possible groundings for symbolic concepts.

In NEURONA, each concept is associated with a grounding module that maps between the concepts and such candidate neural entity embeddings—intuitively, selecting the brain regions that best predict the modular concept. Then, NEURONA uses a differentiable executor to compose these concept groundings according to the structure of the symbolic expression, yielding a final output  $a$  of either a Boolean value or a distribution over concept labels. Together, these components allow NEURONA to accurately decode concepts from fMRI, and test hypotheses on how to guide the grounding process to patterns of neural activity without direct supervision. In the following sections, we describe our method for grounding concepts in fMRI data, our hypotheses, and NEURONA’s training objective.

### 3.2 GROUNDING CONCEPTS ONTO NEURAL ACTIVITIES

In concept grounding, we aim to identify which neural entities predict modular concepts or compositions of concepts, over  $C$  total concepts. Let us consider a base set of candidate neural entities, which we extract by aggregating fMRI signals into functional parcellations. This yields a set of entity embeddings  $\mathcal{E} = [e_1, \dots, e_P]$  for each network  $p$ , where each entity is embedded to dimension  $d$ .

To model unary concepts (e.g., subject and objects such as `person` and `baseball-bat`), we formulate grounding as a  $C$ -way classification problem over  $\mathcal{E}$  using a linear classifier with weight matrix  $\mathcal{W}_{\text{unary}} \in \mathbb{R}^{d \times C}$  and bias  $\mathbf{b}_{\text{unary}} \in \mathbb{R}^C$ . For each entity  $e_p$ , the predicted logits are computed via a linear layer  $\mathbf{z}_p = \mathcal{W}_{\text{unary}}^T e_p + \mathbf{b}_{\text{unary}} \in \mathbb{R}^C$ . The overall predicted logits  $\mathcal{Z}_{\text{unary}} \in \mathbb{R}^{P \times C}$  is the grounding probability of all concepts, and the *grounding score* for concept  $c$  across all entities is  $G_{\text{unary}}(c) = [\mathbf{z}_{1c}, \dots, \mathbf{z}_{Pc}]^\top \in \mathbb{R}^P$ .

To model high-arity relational concepts (e.g., predicates such as `holding`), we augment the neural entity embeddings  $\mathcal{E}$  with learnable embeddings  $\mathcal{E}_b \in \mathbb{R}^{P \times d}$  and concatenate them to form  $\mathcal{E}_c \in \mathbb{R}^{P \times 2d}$ , where  $\mathcal{E}_c = \mathcal{E} \oplus \mathcal{E}_b$  and  $\oplus$  denotes feature concatenation.  $\mathcal{E}_b$  provides features that represent each brain network. For each pair of entities  $(i, j)$ , we concatenate their embeddings  $e_{ci} \oplus e_{cj}$  and apply a learnable transformation  $\mathcal{W}_{\text{pair}} \in \mathbb{R}^{4d \times d}$  to obtain a pairwise representation  $\mathcal{E}'_{ij} \in \mathbb{R}^d$ . We then apply a linear classifier to compute the logits  $\mathbf{z}_{ij}$  for each pair. The overall logits of high-arity concept are  $\mathcal{Z}_{\text{binary}} \in \mathbb{R}^{P \times P \times C}$ , which represents the grounding probability of all concepts. The grounding score for a concept  $c$  is  $G_{\text{binary}}(c) = [\mathbf{z}_{ij,c}]_{1 \leq i, j \leq P} \in \mathbb{R}^{P \times P}$ .

These grounded concepts are then used (and optionally composed) to answer queries from our weakly supervised fMRI-question answering task. Let  $c_p$ ,  $c_s$ , and  $c_o$  be the concepts for predicate, subject, and object, respectively. For Boolean queries, we first optionally condition  $G_{\text{binary}}(c_p)$  on the unary groundings  $G_{\text{unary}}(c_s)$  and  $G_{\text{unary}}(c_o)$ , and aggregate the resulting scores. Then, we apply a sigmoid over the scores, and threshold the result to produce a binary decision in inference. For concept classification queries, given a concept vocabulary  $\mathcal{V}$  of size  $|\mathcal{V}|$ , we compute grounding scores for

each  $v \in \mathcal{V}$ . Since the grounding logits  $\mathcal{Z}_{\text{unary}}$  and  $\mathcal{Z}_{\text{binary}}$  have already been computed for all concepts, we can directly extract the grounding scores for any vocabulary concept  $v$ , treating them as the unary similarity  $\mathcal{S}_{\text{unary}}$  and relational similarity  $\mathcal{S}_{\text{binary}}$  between each neural candidate and  $v$ . Unary and relational similarity scores are computed as

$$S_{\text{unary}} = [z_{iv}]_{1 \leq i \leq P, v \in \mathcal{V}} \in \mathbb{R}^{P \times |\mathcal{V}|}, \quad S_{\text{binary}} = [z_{ij,v}]_{1 \leq i, j \leq P, v \in \mathcal{V}} \in \mathbb{R}^{P \times P \times |\mathcal{V}|}. \quad (1)$$

When subject and object groundings are available, we compute guided scores

$$S_{\text{unary}}^{\text{guided}} = G_{\text{unary}}(c_s)^\top S_{\text{unary}} + G_{\text{unary}}(c_o)^\top S_{\text{unary}} \in \mathbb{R}^{|\mathcal{V}|}, \quad (2)$$

$$S_{\text{binary}}^{\text{guided}} = G_{\text{unary}}(c_s)^\top (G_{\text{unary}}(c_o)^\top S_{\text{binary}}) \in \mathbb{R}^{|\mathcal{V}|}, \quad (3)$$

and combine them as  $S_{\text{final}}^{\text{final}} = S_{\text{unary}}^{\text{guided}} + S_{\text{binary}}^{\text{guided}} \in \mathbb{R}^{|\mathcal{V}|}$ , with the final predicted concept selected by  $\hat{v} = \arg \max_{v \in \mathcal{V}} S_v^{\text{final}}$ . This formulation enables unified, differentiable concept grounding for unary and high-arity concepts, over Boolean and vocabulary-based queries. [Additionally, NEURONA’s role-conditioned aggregation mechanism enables flexible hypothesis testing.](#)

### 3.3 TESTING HYPOTHESES OF GROUNDING STRUCTURES

Notably, we propose and evaluate five hypotheses on how concepts can be composed during grounding in brain networks. Due to NEURONA’s modular structure, we can conduct guided grounding via conditioning the representation of a relational concept on the activations of its constituent arguments, rather than modeling each concept independently, which we see in the latter three hypotheses.

**H1: Single-region localized.** Concepts are localized to a single brain network, where the grounding score is defined as  $G_{\text{H1}}(c) = G_{\text{unary}}(c)$ .

**H2: Multi-region co-activation.** Concepts are represented by co-activation across region pairs, where the grounding score is defined as  $G_{\text{H2}}(c_p) = G_{\text{binary}}(c_p)$ .

**H3: Subject-guided multi-region.** Predicate representations are guided by subject region activation, with the grounding score defined as  $G_{\text{H3}}(c_p) = G_{\text{binary}}(c_p) + G_{\text{unary}}(c_s)$ .

**H4: Object-guided multi-region.** Predicate representations are guided by object region activation, with the grounding score defined as  $G_{\text{H4}}(c_p) = G_{\text{binary}}(c_p) + G_{\text{unary}}(c_o)$ .

**H5: Full argument-guided multi-region.** Multi-region groundings are combined with argument-specific guidance. The predicate grounding score is computed as  $G_{\text{H5}}(c_p) = G(c_p) + G(c_s) + G(c_o)$ , where  $G(c)$  aggregates unary and binary grounding as  $G(c) = G_{\text{unary}}(c) + \frac{1}{P} \sum_{i=1}^P G_{\text{binary}}(c_i)$ . The subject and object scores are implemented as above.

Together, these hypotheses define different subspaces and ways of grounding  $\mathcal{E}$ , allowing us to test with NEURONA whether structure-guided grounding leads to better decoding. Full derivations and experimental details are provided in Appendix D.

### 3.4 TRAINING OBJECTIVE

Our model is trained in a weakly supervised fMRI-QA setting, where ground-truth answers are provided for each query, but no supervision is given on intermediate concept groundings, as none are available. We optimize a standard cross-entropy loss over the model’s predicted output distribution  $\mathcal{L}_{\text{CE}} = -\sum_{i=1}^K a_i \log(\hat{a}_i)$ . Here,  $\hat{a} \in \mathbb{R}^K$  is the predicted probability distribution over  $K$  possible answer classes, and  $a$  is the ground-truth label. The prediction  $\hat{a}$  is the model’s prediction for the full symbolic expression after composing the neural groundings of its component concepts.

## 4 DATASETS

To train NEURONA, we create fMRI-question-answering (fMRI-QA) datasets by adapting existing large-scale fMRI-vision datasets. We first extract structure defining multiple interacting concepts from the rich visual data. Each visual stimuli is processed with a pre-trained vision-language model, which outputs a scene graph that captures both object-level and relational semantics (e.g., unary

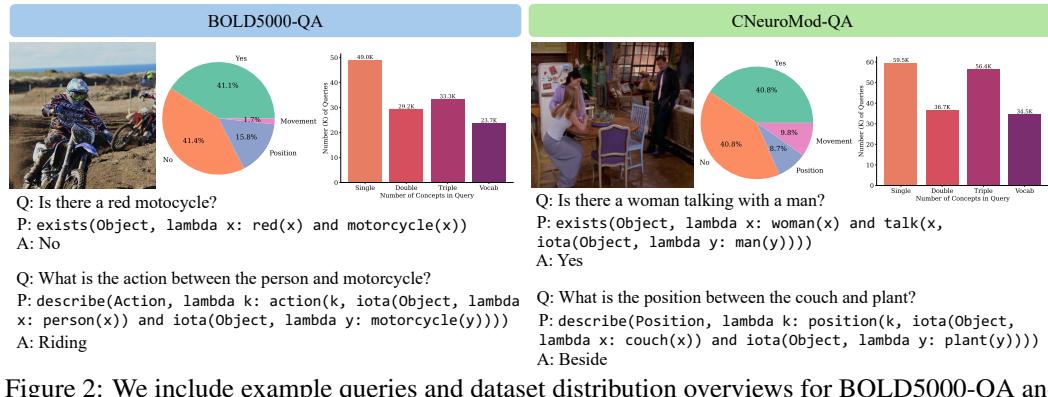


Figure 2: We include example queries and dataset distribution overviews for BOLD5000-QA and CNeuroMod-QA; both datasets span diverse queries and tasks.

objects such as `person`, `baseball-bat`, and high-arity predicates that capture their interaction such as `holding` (`person, baseball-bat`). We then convert the scene graph into a set of structured question-answer pairs, where each question corresponds to a symbolic query and each answer serves as a weak supervision signal for concept grounding. For example, given the above scene graph, we construct questions such as: “What is the relation between `person` and `baseball-bat`?” We additionally generate negative samples by randomly sampling concepts from other stimuli in the dataset. This process produces diverse fMRI-QA examples covering both unary entities and high-arity relations, and notably, with precise answers. Specifically, we apply our pipeline to two datasets: BOLD5000 (Chang et al., 2019) and CNeuroMod (Gifford et al., 2024; Boyle et al., 2023) (See Figure 2). More dataset details can be found in Appendix C.

**BOLD5000-QA.** The BOLD5000 dataset is a large-scale fMRI dataset collected while participants viewed naturalistic images. The dataset consists of approximately 5,000 distinct images from three datasets: Scene Images (Xiao et al., 2010), COCO (Lin et al., 2014), and ImageNet (Deng et al., 2009). Four participants viewed these images while undergoing whole-brain fMRI scanning. To create BOLD5000-QA, we follow the process above, and generate queries containing 4,258 unary and 135 relational concepts, with 133,146 train and 2,095 test examples.

**CNeuroMod-QA.** The CNeuroMod dataset is a large-scale fMRI dataset that includes recordings of participants watching full-length naturalistic videos. Specifically, we build upon the Friends dataset, and set season 1 to 5 to be the train samples, and unseen season 6 to be the test. To create CNeuroMod-QA, we sample video frames based on motion energy, defined as the absolute difference between consecutive frames. We then extract scene graphs for each sampled frame, aggregate the changes, and construct corresponding symbolic queries. The CNeuroMod-QA dataset includes 1,966 unary and 106 relational concepts, with 157,046 train and 30,059 test examples.

## 5 RESULTS

Our goal is precise neural decoding and consistent concept grounding in neural activity. We present quantitative metrics in Section 5.1, qualitative analyses in Section 5.2, and discussion in Section 5.3.

### 5.1 QUANTITATIVE PERFORMANCE

We test quantitative performance of NEURONA on both the BOLD5000-QA and CNeuroMod-QA datasets. We compare against baseline fMRI decoding methods, evaluate generalization to unseen compositional queries, conduct ablation studies to test hypotheses about grounding structures, and report quantitative consistency metrics of concept grounding. Results on performance across atlases, effect sizes between atlases, evaluation over all subjects, statistical tests across hypotheses, cross-dataset transfer experiments, network ablations, detailed concept accuracy, predicate argument binding, fine-grained generalization analyses, additional null models, and fMRI retrieval results are provided in Appendix A.

**Comparison to prior works.** We first evaluate NEURONA on the fMRI-question-answering (fMRI-QA) task, compared against existing decoding methods: a linear baseline, UMBRAE (Xia et al., 2024),

Table 1: We evaluate NEURONA on BOLD5000-QA (subject-CS11) and CNeuroMod-QA (subject-01), comparing its performance to prior fMRI language decoding models and a linear baseline.

| Method         | BOLD5000      |               |               |               | CNeuroMod     |               |               |               |
|----------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|                | Overall       | Action        | Position      | T/F           | Overall       | Action        | Position      | T/F           |
| Linear         | 0.4692        | 0.2069        | 0.1778        | 0.5260        | 0.4638        | 0.3043        | 0.1285        | 0.5192        |
| UMBRAE         | 0.4668        | 0.2069        | 0.1238        | 0.5328        | 0.4614        | 0.0549        | 0.1439        | 0.5442        |
| SDRecon        | 0.4711        | 0.2414        | 0.1937        | 0.5248        | 0.4430        | 0.1350        | 0.1481        | 0.5238        |
| BrainCap       | 0.4773        | 0.1937        | 0.1724        | 0.5551        | 0.4417        | 0.1257        | 0.1477        | 0.5112        |
| NEURONA (Ours) | <b>0.7041</b> | <b>0.6207</b> | <b>0.5079</b> | <b>0.7407</b> | <b>0.7046</b> | <b>0.6514</b> | <b>0.5746</b> | <b>0.7250</b> |

Table 2: Generalization results of NEURONA and prior work on unseen queries.

| Method         | BOLD5000      |               |               |               | CNeuroMod     |               |               |               |
|----------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|                | Overall       | Action        | Position      | T/F           | Overall       | Action        | Position      | T/F           |
| Linear         | 0.4587        | 0.0690        | 0.0794        | 0.5231        | 0.4143        | 0.0398        | 0.0323        | 0.5003        |
| UMBRAE         | 0.4162        | 0.1724        | 0.0540        | 0.4854        | 0.4306        | 0.1315        | 0.1022        | 0.5018        |
| SDRecon        | 0.4702        | 0.2414        | 0.1937        | 0.5237        | 0.4341        | 0.1236        | 0.1473        | 0.5008        |
| BrainCap       | 0.4754        | 0.1724        | 0.1937        | 0.5311        | 0.4347        | 0.1198        | 0.1402        | 0.5042        |
| NEURONA (Ours) | <b>0.6840</b> | <b>0.6207</b> | <b>0.4984</b> | <b>0.7184</b> | <b>0.6583</b> | <b>0.4365</b> | <b>0.5261</b> | <b>0.6991</b> |

SDRecon (Takagi & Nishimoto, 2023), and BrainCap (Ferrante et al., 2023). All models are trained on subject CSI1 (BOLD5000) and subject sub-01 (CNNeuroMod). As shown in Table 1, NEURONA consistently outperforms prior approaches across both the BOLD5000-QA and CNNeuroMod-QA datasets. Compared to prior linear or language model-based approaches, NEURONA achieves a 47% relative improvement on the top performing prior work. These results show that linear methods lack in expressivity, and purely end-to-end decoding pipelines struggle to learn fMRI embeddings that capture the detailed concepts required for fMRI-QA. In particular, NEURONA demonstrates significant gains on queries about actions and positions, which involve precise relational reasoning over subject and object roles. This highlights the strength of our neuro-symbolic framework, which explicitly grounds neural activity guided by hierarchical structures to answer compositional queries.

We further evaluate NEURONA’s ability to generalize to unseen compositions of concepts, which robustly tests whether learned concept groundings are semantically meaningful. Specifically, we construct evaluation splits where all training and testing queries are disjoint, with no overlapping combinations of entities and relations, ensuring that the model must generalize. For example, the training set may contain only queries such as `in_front_of(person, baseball-bat)`, while the test set includes novel compositions such as `in_front_of(baseball-bat, person)`. As shown in Table 2, NEURONA achieves the strongest performance across both datasets, substantially outperforming all baselines. Notably, prior methods suffer significant performance degradation, often falling to near-random levels when generalizing to unseen compositions. The language model-based baselines retain slightly better performance due to their use of pre-trained language models, which carry general linguistic priors. However, since they lack explicit concept grounding, their performance remains well below NEURONA. Overall, these results show that NEURONA’s neuro-symbolic framework learns meaningful decoding from neural activity and generalizes to new compositional queries for robust fMRI-QA.

**Ablations & hypotheses testing.** Notably, we conduct ablation studies on different hypotheses about the candidate neural entity space  $\mathcal{E}$ , and analyze how they affect decoding performance. We summarize results in Table 3 and Table 4: we compare single-region grounding, multi-region grounding without guidance, and forms of guided grounding based on subject and object arguments. We report standard deviation across 4 subjects in BOLD5000-QA and 3 subjects in CNNeuroMod-QA, and our analyses evaluate how different grounding assumptions affect QA accuracy. With NEURONA, we answer the following questions.

**DOES GROUNDING TO MULTIPLE REGIONS IMPROVE PERFORMANCE?** We first test whether allowing concepts to ground to combinations of brain regions improves decoding performance relative to grounding to a single region. As shown in Table 3 and Table 4, we observe that multi-region grounding alone does not yield substantial gains over single-region grounding in terms of overall accuracy. This is especially evident in vocabulary classification tasks, where both approaches tend to overfit to the most frequent vocabulary labels without capturing finer variations.

378 **Table 3: We ablate our hypotheses about the structure of concept grounding on BOLD5000-QA.**

| 379<br>380<br>381<br>382<br>383<br>384 | Method                                | 385<br>386<br>387<br>388<br>389<br>390<br>391<br>392<br>393<br>394<br>395<br>396<br>397<br>398<br>399<br>400<br>401<br>402<br>403<br>404<br>405<br>406<br>407<br>408<br>409<br>410<br>411<br>412<br>413<br>414<br>415<br>416<br>417<br>418<br>419<br>420<br>421<br>422<br>423<br>424<br>425<br>426<br>427<br>428<br>429<br>430<br>431 |                                       |                                       |     |  |
|--|---------------------------------------|---|---------------------------------------|---------------------------------------|-----|--|
|  |                                       | Overall   | Action                                | Position                              | T/F |  |
| Single region                          | 0.6451 $\pm$ 0.0161                   | 0.2973 $\pm$ 0.0361   | 0.2005 $\pm$ 0.0047                   | 0.7293 $\pm$ 0.0191                   |     |  |
| Multi-region (MR)                      | 0.6476 $\pm$ 0.0048                   | 0.2973 $\pm$ 0.0361   | 0.2005 $\pm$ 0.0047                   | 0.7324 $\pm$ 0.0056                   |     |  |
| Subject-guided MR                      | 0.6678 $\pm$ 0.0029                   | 0.2881 $\pm$ 0.0299   | 0.3469 $\pm$ 0.0134                   | 0.7308 $\pm$ 0.0024                   |     |  |
| Object-guided MR                       | 0.6733 $\pm$ 0.0051                   | 0.3892 $\pm$ 0.0752   | 0.3429 $\pm$ 0.0164                   | 0.7363 $\pm$ 0.0045                   |     |  |
| Full argument-guided MR                | <b>0.7102 <math>\pm</math> 0.0053</b> | <b>0.5965 <math>\pm</math> 0.0322</b>   | <b>0.5378 <math>\pm</math> 0.0135</b> | <b>0.7425 <math>\pm</math> 0.0057</b> |     |  |

385 **Table 4: We evaluate our hypotheses about compositional priors in neural activity on CNeuroMod.**

| 386<br>387<br>388<br>389<br>390<br>391<br>392 | Method                                | 387<br>388<br>389<br>390<br>391<br>392<br>393<br>394<br>395<br>396<br>397<br>398<br>399<br>400<br>401<br>402<br>403<br>404<br>405<br>406<br>407<br>408<br>409<br>410<br>411<br>412<br>413<br>414<br>415<br>416<br>417<br>418<br>419<br>420<br>421<br>422<br>423<br>424<br>425<br>426<br>427<br>428<br>429<br>430<br>431 |                                       |                                       |     |
|---|---------------------------------------|---|---------------------------------------|---------------------------------------|-----|
|   |                                       | Overall   | Action                                | Position                              | T/F |
| Single region                                 | 0.6429 $\pm$ 0.0013                   | 0.2165 $\pm$ 0.0193   | 0.2445 $\pm$ 0.0009                   | 0.7400 $\pm$ 0.0016                   |     |
| Multi-region (MR)                             | 0.6162 $\pm$ 0.0027                   | 0.2165 $\pm$ 0.0193   | 0.2445 $\pm$ 0.0009                   | 0.7042 $\pm$ 0.0023                   |     |
| Subject-guided MR                             | 0.6265 $\pm$ 0.0042                   | 0.2339 $\pm$ 0.0043   | 0.3722 $\pm$ 0.0045                   | 0.7008 $\pm$ 0.0051                   |     |
| Object-guided MR                              | 0.6872 $\pm$ 0.0040                   | 0.6320 $\pm$ 0.0019   | 0.4933 $\pm$ 0.0172                   | 0.7149 $\pm$ 0.0045                   |     |
| Full argument-guided MR                       | <b>0.7189 <math>\pm</math> 0.0009</b> | <b>0.6422 <math>\pm</math> 0.0072</b>   | <b>0.5931 <math>\pm</math> 0.0084</b> | <b>0.7417 <math>\pm</math> 0.0005</b> |     |

393 DOES GUIDED GROUNDING IMPROVE PERFORMANCE? Next, we evaluate whether guiding multi-  
394 region grounding based on the subject and object arguments of relational concepts improves per-  
395 formance. As shown, models that incorporate guided grounding significantly outperform both  
396 single-region and unguided multi-region baselines. We find that guiding NEURONA by object  
397 improves performance over by subject. In particular, grounding based on both subject and object  
398 regions achieves the highest accuracy, notably on action and position queries that require precise  
399 relational reasoning. With NEURONA, we demonstrate the importance of argument-conditioned  
400 composition for interpreting relational semantics in neural activity. Overall, these results demonstrate  
401 that while multi-region grounding provides a more flexible representation space, explicit structural  
402 guidance based on predicate-argument relationships is crucial. The consistent results on both natural  
403 image and video datasets validate NEURONA ability to conduct complex neural decoding.

404 **Concept grounding consistency.** As there is no ground truth to evaluate the reliability of inter-  
405 mediate concept grounding, we introduce a consistency metric to test whether the same concept  
406 grounds to consistent brain regions across different fMRI-QA instances. We calculate consistency  
407 as follows. Let a concept  $c$  appear in  $N$  QA examples, and let the predicted grounding in the  $i$ -th  
408 example be a set of brain regions  $B^{(i)}(c) \subseteq \{1, \dots, P\}$ , where  $P$  is the total number of regions,  
409 and each  $B^{(i)}(c)$  is the set of regions selected as the grounding for concept  $c$  in that example. We  
410 compute the frequency count for each region  $r$  as  $\text{Count}(r) = \sum_{i=1}^N \mathbf{1}[r \in B^{(i)}(c)]$ . Then, the score  
411 for concept  $c$  is defined as  $\text{Consistency}(c) = \frac{1}{|R|} \sum_{r \in R} \frac{\text{Count}(r)}{N}$ , where  $R$  is the set of all regions that  
412 appear in any grounding of  $c$ , and  $\frac{\text{Count}(r)}{N}$  is the fraction of times region  $r$  was selected.

413 Our proposed metric captures how concentrated the concepts groundings are: a score of 1.0 implies  
414 perfect consistency (all instances of the concept used the same region set), while lower scores  
415 indicate more variability in the grounding of concept  $c$ . We report these consistency scores over  
416 all concepts in BOLD5000 and CNeuroMod. As a baseline, we define a null model that randomly  
417 assigns each concept to a region subset via uniform sampling. In Table 5, we see that NEURONA  
418 significantly outperforms the null baseline. These results highlight that NEURONA not only improves  
419 QA accuracy but also grounds concepts in a structured and reproducible way across different stimuli.  
420 Consistency results over all other atlases can be found in Appendix A.

## 5.2 QUALITATIVE CONCEPT GROUNDING ANALYSES

424 Finally, we qualitatively examine how NEURONA  
425 learns to decode high-level relational concepts from  
426 neural activity using intermediate concept ground-  
427 ings, focusing on how this grounding varies with  
428 different subject-object pairs. Figure 3 shows rep-  
429 resentative examples from both BOLD5000-QA and  
430 CNeuroMod-QA, where we project the learned grounding scores onto network parcellations that  
431 define our neural entities. **We observe that the same relational predicate, such as `hold` or `look`, is**  
**best decoded from different brain networks depending on the object.** For example, in BOLD5000,

Table 5: Consistency of concept grounding  
between NEURONA and a null baseline.

|         | BOLD5000      | CNeuroMod     |
|---------|---------------|---------------|
| Null    | 0.5357        | 0.5358        |
| NEURONA | <b>0.8220</b> | <b>0.8700</b> |

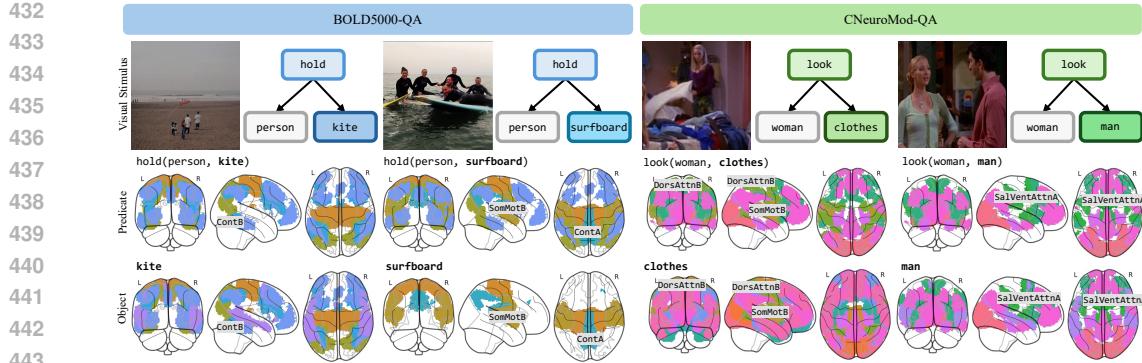


Figure 3: We show examples of learned concept grounding from NEURONA. On both BOLD5000-QA and CNeuroMod-QA, we see that predicate concepts ground to the brain regions that their constituent objects are grounded to, following hierarchical predicate-argument structure.

holding (person, kite) is best decoded using the Control B network, while holding (person, surfboard) is decoded using both the Somatomotor B and Control A networks. **This suggests that decoding accuracy improves when the grounding of a relational concept is modulated by its arguments, supporting the importance of argument-dependent composition in neural decoding.**

Interestingly, the model-inferred concept groundings are not confined to early visual areas, despite the task involving visual stimuli. Instead, objects such as baseball-bat and surfboard receive high grounding scores in motor-related regions, including the Somatomotor network. This qualitative pattern resembles prior findings that perceiving action-related objects is linked to motor and premotor areas (Martin, 2007; Gallese et al., 1996). Additionally, both hold and look are often decoded using prefrontal networks, including the Dorsal Attention and Salience/Ventral Attention networks, which have been associated with high-level cognitive control and abstract rule processing (Miller & Cohen, 2001; Quiroga et al., 2005; Tian et al., 2024). **We emphasize that these analyses reflect model-dependent decoding patterns rather than direct estimates of underlying encoding representations.** Additional concept grounding visualizations are provided in Appendix B.

### 5.3 DISCUSSION

Our findings demonstrate that incorporating compositional structure into the decoding pipeline significantly improves performance. While NEURONA does not establish representational compositionality in neural patterns, the substantial gains from modeling hierarchical predicate–argument structure provide proof-of-concept that compositional principles can inform neural decoding. However, our study also has several limitations. First, participants in our datasets engaged only in passive viewing, and the ground truth symbolic expressions were extracted through automated scene-graph parsing rather than participant-driven reasoning, limiting the cognitive conclusions that can be drawn. Second, we restrict candidate neural entities to predefined parcellations from established atlases (e.g., Yeo-7, Yeo-17, DiFuMo-64, DiFuMo-128, and Schaefer-100), which, while widely used, provide coarser representations of the brain from which we build upon. In the Appendix, we provide performance on BOLD5000 and CNeuroMod across atlases, effect sizes between atlases, and concept grounding consistency across atlases. These results demonstrate robust and modular concept grounding across different levels of spatial granularity, suggesting that our results are not tied to specific parcellation choices. However, we believe that scaling NEURONA to incorporate whole-brain voxel-level data is a promising next step.

## 6 CONCLUSION

We propose NEURONA, a neuro-symbolic framework for concept grounding and decoding in neural activity. By leveraging symbolic reasoning and compositional execution with fMRI grounding, NEURONA enables precise and generalizable decoding. Experiments on BOLD5000-QA and CNeuroMod-QA demonstrate that NEURONA outperforms baseline decoding methods and generalizes to novel compositions, with explicit grounding guidance significantly improving performance. Our findings show that incorporating predicate-argument structure improves decoding, and highlight neuro-symbolic modeling as a promising approach for interpreting and structuring fMRI decoding models.

486     **Reproducibility statement.** We refer readers to Section 3 for details on the grounding process and  
 487     Appendix D for train settings, and note that our work builds off the public codebase of LEFT. We will  
 488     release code upon acceptance. We also describe our dataset processing steps in detail in Appendix C.  
 489

490     REFERENCES  
 491

492     Julie Boyle, Basile Pinsard, Valentina Borghesani, Francois Paugam, Elizabeth DuPre, and Pierre  
 493     Bellec. The Courtois NeuroMod project: Quality Assessment of the Initial Data Release (2020).  
 494     In *2023 Conference on Cognitive Computational Neuroscience*, pp. 2023–1602, 2023.

495     Nadine Chang, John A Pyles, Austin Marcus, Abhinav Gupta, Michael J Tarr, and Elissa M Aminoff.  
 496     BOLD5000, A Public fMRI Dataset While Viewing 5000 Visual Images. *Scientific data*, 6(1):49,  
 497     2019.

498     Zijiao Chen, Jiaxin Qing, Tiange Xiang, Wan Lin Yue, and Juan Helen Zhou. Seeing Beyond  
 499     the Brain: Conditional Diffusion Model with Sparse Masked Modeling for Vision Decoding.  
 500     In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
 501     22710–22720, 2023a.

502     Zijiao Chen, Jiaxin Qing, and Juan Helen Zhou. Cinematic Mindscapes: High-quality Video  
 503     Reconstruction from Brain Activity. *Advances in Neural Information Processing Systems*, 36:  
 504     24841–24858, 2023b.

505     Kamalaker Dadi, Gaël Varoquaux, Antonia Machlouzarides-Shalit, Krzysztof J Gorgolewski, Demian  
 506     Wassermann, Bertrand Thirion, and Arthur Mensch. Fine-grain Atlases of Functional Modes for  
 507     fMRI Analysis. *NeuroImage*, 221:117126, 2020.

508     Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A Large-scale  
 509     Hierarchical Image Database. In *2009 IEEE conference on computer vision and pattern recognition*,  
 510     pp. 248–255. Ieee, 2009.

511     Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep  
 512     Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 conference  
 513     of the North American chapter of the association for computational linguistics: human language  
 514     technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.

515     Kingma Diederik. Adam: A Method for Stochastic Optimization. (*No Title*), 2014.

516     Doris E Dijksterhuis, Matthew W Self, Jessy K Posse, Judith C Peters, ECW van Straaten, Sander  
 517     Idema, Johannes C Baaijen, Sandra MA van der Salm, Erik J Aarnoutse, Nicole CE van Klink,  
 518     et al. Pronouns Reactivate Conceptual Representations in Human Hippocampal Neurons. *Science*,  
 519     385(6716):1478–1484, 2024.

520     Howard Eichenbaum. The Hippocampus and Declarative Memory: Cognitive Mechanisms and  
 521     Neural Codes. *Behavioural brain research*, 127(1-2):199–207, 2001.

522     Matteo Ferrante, Tommaso Boccato, Furkan Ozcelik, Rufin VanRullen, and Nicola Toschi. Multi-  
 523     modal Decoding of Human Brain activity into Images and Text. In *UniReps: the First Workshop  
 524     on Unifying Representations in Neural Models*, 2023.

525     Jerry Fodor. *The Language of Thought*. Harvard University Press, 1975.

526     Steven M Frankland and Joshua D Greene. An Architecture for Encoding Sentence Meaning in  
 527     Left Mid-superior Temporal Cortex. *Proceedings of the National Academy of Sciences*, 112(37):  
 528     11732–11737, 2015.

529     Vittorio Gallese, Luciano Fadiga, Leonardo Fogassi, and Giacomo Rizzolatti. Action Recognition in  
 530     the Premotor Cortex. *Brain*, 119(2):593–609, 1996.

531     Alessandro T Gifford, Domenic Bersch, Marie St-Laurent, Basile Pinsard, Julie Boyle, Lune Bellec,  
 532     Aude Oliva, Gemma Roig, and Radoslaw M Cichy. The Algonauts Project 2025 Challenge: How  
 533     the Human Brain Makes Sense of Multimodal Movies. *arXiv preprint arXiv:2501.00504*, 2024.

540 Zixuan Gong, Guangyin Bao, Qi Zhang, Zhongwei Wan, Duoqian Miao, Shoujin Wang, Lei Zhu,  
 541 Changwei Wang, Rongtao Xu, Liang Hu, et al. NeuroClips: Towards High-fidelity and Smooth  
 542 fMRI-to-Video Reconstruction. *Advances in Neural Information Processing Systems*, 37:51655–  
 543 51683, 2024.

544 Joy Hsu, Jiayuan Mao, Josh Tenenbaum, and Jiajun Wu. What’s Left? Concept Grounding with  
 545 Logic-Enhanced Foundation Models. *Advances in Neural Information Processing Systems*, 36:  
 546 38798–38814, 2023.

548 Alexander G Huth, Wendy A De Heer, Thomas L Griffiths, Frédéric E Theunissen, and Jack L  
 549 Gallant. Natural Speech Reveals the Semantic Maps that Tile Human Cerebral Cortex. *Nature*,  
 550 532(7600):453–458, 2016.

551 Kendrick N Kay, Thomas Naselaris, Ryan J Prenger, and Jack L Gallant. Identifying Natural Images  
 552 from Human Brain Activity. *Nature*, 452(7185):352–355, 2008.

554 Aaron Kucyi, Amy Daitch, Omri Raccah, Baotian Zhao, Chao Zhang, Michael Esterman, Michael  
 555 Zeineh, Casey H Halpern, Kai Zhang, Jianguo Zhang, et al. Electrophysiological Dynamics of  
 556 Antagonistic Brain Networks Reflect Attentional Fluctuations. *Nature communications*, 11(1):325,  
 557 2020.

558 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping Language-image Pre-  
 559 training for Unified Vision-language Understanding and Generation. In *International conference  
 560 on machine learning*, pp. 12888–12900. PMLR, 2022.

562 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
 563 Dollár, and C Lawrence Zitnick. Microsoft COCO: Common Objects in Context. In *Computer  
 564 vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6–12, 2014,  
 565 proceedings, part v 13*, pp. 740–755. Springer, 2014.

566 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual Instruction Tuning. *Advances in  
 567 neural information processing systems*, 36:34892–34916, 2023.

569 Jiayuan Mao, Chuang Gan, Pushmeet Kohli, Joshua B Tenenbaum, and Jiajun Wu. The Neuro-  
 570 Symbolic Concept Learner: Interpreting Scenes, Words, and Sentences from Natural Supervision.  
 571 *arXiv preprint arXiv:1904.12584*, 2019.

572 Jiayuan Mao, Joshua B Tenenbaum, and Jiajun Wu. Neuro-Symbolic Concepts. *arXiv preprint  
 573 arXiv:2505.06191*, 2025.

575 Alex Martin. The Representation of Object Concepts in the Brain. *Annu. Rev. Psychol.*, 58(1):25–45,  
 576 2007.

577 Vinod Menon. 20 Years of the Default Mode Network: A Review and Synthesis. *Neuron*, 111(16):  
 578 2469–2487, 2023.

580 Matthew K Mian, Sameer A Sheth, Shaun R Patel, Konstantinos Spiliopoulos, Emad N Eskandar,  
 581 and Ziv M Williams. Encoding of Rules by Neurons in the Human Dorsolateral Prefrontal Cortex.  
 582 *Cerebral cortex*, 24(3):807–816, 2014.

583 Earl K Miller and Jonathan D Cohen. An Integrative Theory of Prefrontal Cortex function. *Annual  
 584 review of neuroscience*, 24(1):167–202, 2001.

586 Tom M Mitchell, Svetlana V Shinkareva, Andrew Carlson, Kai-Min Chang, Vicente L Malave,  
 587 Robert A Mason, and Marcel Adam Just. Predicting Human Brain Activity Associated with the  
 588 Meanings of Nouns. *science*, 320(5880):1191–1195, 2008.

589 Yoichi Miyawaki, Hajime Uchida, Okito Yamashita, Masa-aki Sato, Yusuke Morito, Hiroki C Tanabe,  
 590 Norihiro Sadato, and Yukiyasu Kamitani. Visual Image Reconstruction from Human Brain Activity  
 591 Using a Combination of Multiscale Local Image Decoders. *Neuron*, 60(5):915–929, 2008.

593 Thomas Naselaris, Ryan J Prenger, Kendrick N Kay, Michael Oliver, and Jack L Gallant. Bayesian  
 Reconstruction of Natural Images from Human Brain Activity. *Neuron*, 63(6):902–915, 2009.

594 Thomas Naselaris, Kendrick N Kay, Shinji Nishimoto, and Jack L Gallant. Encoding and Decoding  
 595 in fMRI. *Neuroimage*, 56(2):400–410, 2011.  
 596

597 Shinji Nishimoto, An T Vu, Thomas Naselaris, Yuval Benjamini, Bin Yu, and Jack L Gallant.  
 598 Reconstructing Visual Experiences from Brain Activity Evoked by Natural Movies. *Current  
 599 biology*, 21(19):1641–1646, 2011.

600 Mark Palatucci, Dean Pomerleau, Geoffrey E Hinton, and Tom M Mitchell. Zero-shot Learning with  
 601 Semantic Output Codes. *Advances in neural information processing systems*, 22, 2009.  
 602

603 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: A Method for Automatic  
 604 Evaluation of Machine Translation. In *Proceedings of the 40th annual meeting of the Association  
 605 for Computational Linguistics*, pp. 311–318, 2002.

606 Francisco Pereira, Bin Lou, Brianna Pritchett, Samuel Ritter, Samuel J Gershman, Nancy Kanwisher,  
 607 Matthew Botvinick, and Evelina Fedorenko. Toward a Universal Decoder of Linguistic Meaning  
 608 from Brain Activation. *Nature communications*, 9(1):963, 2018.  
 609

610 R Quiroga, Leila Reddy, Gabriel Kreiman, Christof Koch, and Itzhak Fried. Invariant Visual  
 611 Representation by Single Neurons in the Human Brain. *Nature*, 435(7045):1102–1107, 2005.  
 612

613 Hernan G Rey, Matias J Ison, Carlos Pedreira, Antonio Valentin, Gonzalo Alarcon, Richard Selway,  
 614 Mark P Richardson, and Rodrigo Quiroga. Single-cell Recordings in the Human Medial  
 615 Temporal Lobe. *Journal of anatomy*, 227(4):394–408, 2015.

616 Alexander Schaefer, Ru Kong, Evan M Gordon, Timothy O Laumann, Xi-Nian Zuo, Avram J Holmes,  
 617 Simon B Eickhoff, and BT Thomas Yeo. Local-global Parcellation of the Human Cerebral Cortex  
 618 from Intrinsic Functional Connectivity MRI. *Cerebral cortex*, 28(9):3095–3114, 2018.  
 619

620 Paul Scotti, Atmadeep Banerjee, Jimmie Goode, Stepan Shabalin, Alex Nguyen, Aidan Dempster,  
 621 Nathalie Verlinde, Elad Yundler, David Weisberg, Kenneth Norman, et al. Reconstructing the  
 622 Mind’s Eye: fMRI-to-Image with Contrastive Learning and Diffusion Priors. *Advances in Neural  
 623 Information Processing Systems*, 36:24705–24728, 2023.

624 Yu Takagi and Shinji Nishimoto. High-Resolution Image Reconstruction with Latent Diffusion  
 625 Models from Human Brain Activity. In *Proceedings of the IEEE/CVF Conference on Computer  
 626 Vision and Pattern Recognition*, pp. 14453–14463, 2023.  
 627

628 Bertrand Thirion, Edouard Duchesnay, Edward Hubbard, Jessica Dubois, Jean-Baptiste Poline, Denis  
 629 Lebihan, and Stanislas Dehaene. Inverse Retinotopy: Inferring the Visual Content of Images from  
 630 Brain Activation Patterns. *Neuroimage*, 33(4):1104–1116, 2006.

631 Zhenghe Tian, Jingwen Chen, Cong Zhang, Bin Min, Bo Xu, and Liping Wang. Mental Programming  
 632 of Spatial Sequences in Working Memory in the Macaque Frontal Cortex. *Science*, 385(6716):  
 633 eadp6091, 2024.  
 634

635 Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu,  
 636 and Lijuan Wang. Git: A Generative Image-to-text Transformer for Vision and Language. *arXiv  
 637 preprint arXiv:2205.14100*, 2022.

638 Yanchen Wang, Adam Turnbull, Tiange Xiang, Yunlong Xu, Sa Zhou, Adnan Masoud, Shekoofeh  
 639 Azizi, Feng Vankee Lin, and Ehsan Adeli. Decoding visual experience and mapping semantics  
 640 through whole-brain analysis using fmri foundation models. *arXiv preprint arXiv:2411.07121*,  
 641 2024.

643 Weihao Xia, Raoul de Charette, Cengiz Oztireli, and Jing-Hao Xue. Umbrae: Unified Multimodal  
 644 Brain Decoding. In *European Conference on Computer Vision*, pp. 242–259. Springer, 2024.  
 645

646 Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. SUN Database:  
 647 Large-scale Scene Recognition from Abbey to Zoo. In *2010 IEEE computer society conference on  
 computer vision and pattern recognition*, pp. 3485–3492. IEEE, 2010.

648 BT Thomas Yeo, Fenna M Krienen, Jorge Sepulcre, Mert R Sabuncu, Danial Lashkari, Marisa  
649 Hollinshead, Joshua L Roffman, Jordan W Smoller, Lilla Zöllei, Jonathan R Polimeni, et al.  
650 The Organization of the Human Cerebral Cortex Estimated by Intrinsic Functional Connectivity.  
651 *Journal of neurophysiology*, 2011.

652 Kexin Yi, Jiajun Wu, Chuang Gan, Antonio Torralba, Pushmeet Kohli, and Josh Tenenbaum. Neural-  
653 symbolic VQA: Disentangling Reasoning from Vision and Language Understanding. *Advances in  
654 neural information processing systems*, 31, 2018.

655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

702

703

704

705

706

## Supplementary for: Neuro-Symbolic Decoding of Neural Activity

707

708

709

710

The appendix is organized into [five](#) main sections. Appendix A includes additional experiment results: performance across atlases, effect sizes between atlases, concept grounding consistency across atlases, [evaluation over all subjects, statistical tests across hypotheses](#), cross-dataset transfer experiments, [network ablations](#), detailed concept accuracy, [predicate argument binding](#), [fine-grained generalization analyses](#), [additional null models](#), and [fMRI retrieval results](#). Appendix B provides additional visualizations and analyses of NEURONA’s concept grounding performance. Appendix D describes the training procedure of NEURONA, the implementation of baseline methods, and the setup of our hypothesis ablation experiments. Appendix C presents more examples illustrating our fMRI-QA datasets and detail the data generation process. [Appendix E details our ethics statement](#). Here, we also note that we use large language models to make minor improvements to writing.

711

712

## A ADDITIONAL RESULTS

713

### A.1 PERFORMANCE ACROSS ATLASES

714

We report performance from NEURONA across atlases to show robustness of decoding. We map parcellated fMRI signals (1024 regions for BOLD5000, 1000 for CNNeuroMod) to multiple atlases, including Yeo-7, Yeo-17 (Yeo et al., 2011), DiFuMo-64, DiFuMo-128 (Dadi et al., 2020), and Schaefer-100 (Schaefer et al., 2018), then train NEURONA on these neural entities. Results in Table 6 and Table 7 show that NEURONA consistently learns across these atlases, and still significantly outperforms prior works in decoding accuracy. Yeo-17 yields the highest accuracy among all tested atlases, followed by DiFuMo-128 and Schaefer-100.

715

716

Table 6: NEURONA’s performance with coarse and fine-grained atlases on BOLD5000.

| BOLD5000   | Overall       | Action        | Position      | T/F           |
|------------|---------------|---------------|---------------|---------------|
| Yeo-7      | 0.6864        | 0.5517        | 0.4730        | 0.7270        |
| Yeo-17     | <b>0.7041</b> | <b>0.6207</b> | 0.5079        | <b>0.7407</b> |
| DiFuMo-64  | 0.6992        | 0.5517        | <b>0.5524</b> | 0.7282        |
| DiFuMo-128 | 0.7026        | 0.5862        | 0.5460        | 0.7327        |

717

718

Table 7: NEURONA’s performance with coarse and fine-grained atlases on CNNeuroMod.

| CNeuroMod    | Overall       | Action        | Position      | T/F           |
|--------------|---------------|---------------|---------------|---------------|
| Yeo-7        | 0.6969        | 0.6459        | 0.5577        | 0.7180        |
| Yeo-17       | <b>0.7046</b> | 0.6514        | <b>0.5746</b> | 0.7250        |
| Schaefer-100 | 0.7043        | <b>0.6549</b> | 0.5614        | <b>0.7258</b> |

719

720

### A.2 EFFECT SIZES BETWEEN ATLASES

721

722

To additionally evaluate the robustness of NEURONA to different brain parcellations, we compute Cohen’s d effect sizes between QA predictions from different atlases. For each atlas pair, we compute paired effect sizes using QA predictions across the test set. As seen in Table 8 and Table 9, effect sizes are consistently small, showing that NEURONA is robust to the choice of atlas and performs reliably across a range of parcellations.

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

## A.3 CONCEPT GROUNDING CONSISTENCY ACROSS ATLASES

In Table 10 and Table 11, we report additional consistency scores averaged over all concepts in BOLD5000 and CNeuroMod, under multiple atlas configurations: Yeo-7, Yeo-17, DiFuMo64, DiFuMo128, and Schaefer100. NEURONA achieves consistently high grounding consistency across all atlases, significantly above the null baseline. Unary concepts show higher consistency than relational ones, as expected due to their simpler structure. Across the atlases, all show consistent results, with DiFuMo-64 best for BOLD5000 and Yeo-17 best for CNeuroMod. NEURONA’s concept grounding is reproducible way across different stimuli and parcellations.

Table 8: Effect sizes between atlases in BOLD5000.

| BOLD5000   | Yeo-7 | Yeo-17 | Difumo-64 | Difumo-128 |
|------------|-------|--------|-----------|------------|
| Yeo-7      | -     | -0.017 | -0.133    | -0.012     |
| Yeo-17     | 0.017 | -      | -0.116    | 0.006      |
| Difumo-64  | 0.133 | 0.116  | -         | 0.121      |
| Difumo-128 | 0.012 | -0.006 | -0.121    | -          |

Table 9: Effect sizes between atlases in CNeuroMod.

| CNeuroMod    | Yeo-7  | Yeo-17 | Schaefer-100 |
|--------------|--------|--------|--------------|
| Yeo-7        | -      | 0.094  | 0.089        |
| Yeo-17       | -0.094 | -      | -0.006       |
| Schaefer-100 | -0.089 | 0.006  | -            |

Table 10: Concept grounding consistency in BOLD5000.

| BOLD5000           | Overall | Unary Concept | Relational Concept |
|--------------------|---------|---------------|--------------------|
| Yeo-7 Null         | 0.5738  | —             | —                  |
| Yeo-7 NEURONA      | 0.8207  | 0.8283        | 0.7343             |
| Yeo-17 Null        | 0.5357  | —             | —                  |
| Yeo-17 NEURONA     | 0.8220  | 0.8351        | 0.6646             |
| DiFuMo-64 Null     | 0.5075  | —             | —                  |
| DiFuMo-64 NEURONA  | 0.8462  | 0.8644        | 0.6064             |
| DiFuMo-128 Null    | 0.5039  | —             | —                  |
| DiFuMo-128 NEURONA | 0.8224  | 0.8380        | 0.6241             |

## A.4 EVALUATION OVER ALL SUBJECTS

We include evaluation over all subjects on both BOLD5000 (4 subjects) and CNeuroMod (3 subjects). For each dataset, we train an individual model for each subject and report the mean  $\pm$  standard deviation across subjects. We evaluate the in-distribution neural decoding performance in Table 12 and our generalization split in Table 13. We see that NEURONA continues to significantly outperform all prior works.

In addition, we evaluate cross-subject consistency in concept groundings. To evaluate whether concepts are grounded similarly across individuals, we aggregate grounding scores for each concept across all subjects and queries (for concept  $c$  appearing  $N$  times per subject, we obtain  $N \times S$  samples where  $S$  is the number of subjects). We then compute a cross-subject consistency metric that measures the similarity of the spatial grounding patterns for each concept across individuals. We evaluate NEURONA against a random assignment null model and a multinomial null model that preserves each concept’s distribution across regions. Our results in Table 14 show that NEURONA’s cross-subject grounding consistency scores, averaged across all concepts, are significantly higher than both null models ( $p < 0.001$ ), demonstrating that NEURONA’s concept groundings converge to more similar sets of regions across participants.

810

811

Table 11: Concept grounding consistency in CNeuroMod.

| CNeuroMod            | Overall | Unary Concept | Relational Concept |
|----------------------|---------|---------------|--------------------|
| Yeo-7 Null           | 0.5740  | –             | –                  |
| Yeo-7 NEURONA        | 0.8437  | 0.8564        | 0.7563             |
| Yeo-17 Null          | 0.5358  | –             | –                  |
| Yeo-17 NEURONA       | 0.8700  | 0.8967        | 0.6812             |
| Schaefer-100 Null    | 0.5029  | –             | –                  |
| Schaefer-100 NEURONA | 0.8346  | 0.8695        | 0.5838             |

820

821

Table 12: Results across subjects for BOLD5000 and CNeuroMod.

| BOLD5000  | Overall                               | Action                                | Position                              | T/F                                   |
|-----------|---------------------------------------|---------------------------------------|---------------------------------------|---------------------------------------|
| Linear    | $0.4702 \pm 0.0073$                   | $0.0677 \pm 0.0815$                   | $0.1817 \pm 0.0289$                   | $0.5280 \pm 0.0068$                   |
| UMBRAE    | $0.4948 \pm 0.0098$                   | $0.2488 \pm 0.0572$                   | $0.2092 \pm 0.0401$                   | $0.5401 \pm 0.0097$                   |
| SDRecon   | $0.4751 \pm 0.0064$                   | $0.2887 \pm 0.0499$                   | $0.2005 \pm 0.0047$                   | $0.5266 \pm 0.0075$                   |
| BrainCap  | $0.4842 \pm 0.0052$                   | $0.2412 \pm 0.0470$                   | $0.2005 \pm 0.0047$                   | $0.5383 \pm 0.0061$                   |
| NEURONA   | <b><math>0.7102 \pm 0.0053</math></b> | <b><math>0.5965 \pm 0.0322</math></b> | <b><math>0.5378 \pm 0.0135</math></b> | <b><math>0.7425 \pm 0.0057</math></b> |
| CNeuroMod | Overall                               | Action                                | Position                              | T/F                                   |
| Linear    | $0.4579 \pm 0.0142$                   | $0.2078 \pm 0.0400$                   | $0.1680 \pm 0.0338$                   | $0.5184 \pm 0.0153$                   |
| UMBRAE    | $0.4588 \pm 0.0425$                   | $0.1588 \pm 0.0102$                   | $0.1705 \pm 0.0220$                   | $0.5225 \pm 0.0248$                   |
| SDRecon   | $0.4368 \pm 0.0012$                   | $0.1428 \pm 0.0010$                   | $0.1554 \pm 0.0010$                   | $0.5012 \pm 0.0024$                   |
| BrainCap  | $0.4422 \pm 0.0026$                   | $0.1245 \pm 0.0099$                   | $0.1497 \pm 0.0074$                   | $0.5110 \pm 0.0024$                   |
| NEURONA   | <b><math>0.7189 \pm 0.0009</math></b> | <b><math>0.6422 \pm 0.0072</math></b> | <b><math>0.5931 \pm 0.0084</math></b> | <b><math>0.7417 \pm 0.0005</math></b> |

835

836

### A.5 STATISTICAL TESTS ACROSS HYPOTHESES

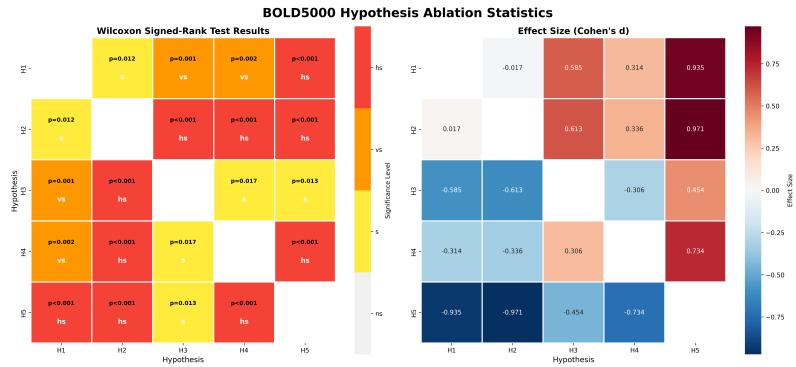
839

840

We conduct additional statistical analyses for our hypotheses ablations to more rigorously evaluate the effects of the five hypotheses across subjects. Specifically, we use the Wilcoxon signed-rank test to assess statistical significance and Cohen’s d to estimate effect size. For each subject, we collect accuracies across all metrics, and perform comparisons between hypotheses using these values. The BOLD5000 and CNeuroMod results are summarized in Figure 4 and Figure 5. Across both datasets, NEURONA’s full argument-guided multi-region hypothesis (H5) consistently shows statistically significant improvements over alternative hypotheses.

845

846



858

859

Figure 4: Statistical tests evaluating the effects of the five hypotheses across subjects in BOLD5000.

860

861

### A.6 CROSS-DATASET TRANSFER EXPERIMENTS

862

863

Here, we include cross-dataset generalization experiments by training our model on BOLD5000-QA and evaluating it on the CNeuroMod-QA test set. Since BOLD5000 spans a broader concept space,

864

865 Table 13: Generalization results across subjects for BOLD5000 and CNeuroMod.

| BOLD5000  | Overall                               | Action                                | Position                              | T/F                                   |
|-----------|---------------------------------------|---------------------------------------|---------------------------------------|---------------------------------------|
| Linear    | $0.4595 \pm 0.0093$                   | $0.0308 \pm 0.0312$                   | $0.1291 \pm 0.0545$                   | $0.5250 \pm 0.0059$                   |
| UMBRAE    | $0.4886 \pm 0.0068$                   | $0.1214 \pm 0.0217$                   | $0.0914 \pm 0.0212$                   | $0.5392 \pm 0.0080$                   |
| SDRecon   | $0.4752 \pm 0.0064$                   | $0.2973 \pm 0.0361$                   | $0.2005 \pm 0.0047$                   | $0.5266 \pm 0.0075$                   |
| BrainCap  | $0.4838 \pm 0.0055$                   | $0.2412 \pm 0.0470$                   | $0.1996 \pm 0.0040$                   | $0.5380 \pm 0.0064$                   |
| NEURONA   | <b><math>0.6812 \pm 0.0055</math></b> | <b><math>0.4952 \pm 0.0682</math></b> | <b><math>0.4696 \pm 0.0190</math></b> | <b><math>0.7217 \pm 0.0057</math></b> |
| CNeuroMod | Overall                               | Action                                | Position                              | T/F                                   |
| Linear    | $0.4206 \pm 0.0129$                   | $0.0462 \pm 0.0654$                   | $0.1045 \pm 0.0667$                   | $0.4986 \pm 0.0052$                   |
| UMBRAE    | $0.4487 \pm 0.0040$                   | $0.1107 \pm 0.0156$                   | $0.1386 \pm 0.0047$                   | $0.5247 \pm 0.0035$                   |
| SDRecon   | $0.4365 \pm 0.0014$                   | $0.1407 \pm 0.0008$                   | $0.1560 \pm 0.0012$                   | $0.5012 \pm 0.0019$                   |
| BrainCap  | $0.4382 \pm 0.0012$                   | $0.1217 \pm 0.0070$                   | $0.1455 \pm 0.0060$                   | $0.5069 \pm 0.0003$                   |
| NEURONA   | <b><math>0.6676 \pm 0.0119</math></b> | <b><math>0.2916 \pm 0.0878</math></b> | <b><math>0.5377 \pm 0.0215</math></b> | <b><math>0.7260 \pm 0.0072</math></b> |

872

873

874 Table 14: Cross-subject consistency of concept grounding.

|                    | BOLD5000            | CNeuroMod           |
|--------------------|---------------------|---------------------|
| Null (random)      | $0.5276 \pm 0.0629$ | $0.5267 \pm 0.0683$ |
| NULL (multinomial) | $0.6045 \pm 0.1685$ | $0.6481 \pm 0.1644$ |
| NEURONA            | $0.7000 \pm 0.2514$ | $0.7027 \pm 0.2165$ |

875

876

877

878 we selected overlapping queries across datasets. In the CNeuroMod test set, this includes 1,169  
 879 queries for the action task, 2,600 for the position task, and 23,038 for the T/F task (out of full test  
 880 set sizes of 2,912, 2,661, and 24,486, respectively).

881 In Table 15, we compare NEURONA to UMBRAE (Xia et al., 2024), the top performing baseline  
 882 model. NEURONA significantly outperforms UMBRAE across all queries, demonstrating stronger  
 883 cross-dataset robustness and generalization. Notably, while overall performance of NEURONA drops,  
 884 largely due to a performance gap on T/F queries, accuracy on action and position tasks remains high,  
 885 indicating some degree of cross-dataset transfer. This drop is expected, as our model is trained as a  
 886 subject-specific model and there is substantial variance across subjects. Additionally, the two datasets  
 887 differ in preprocessing pipelines: BOLD5000 uses the DiFuMo-1024 parcellation, while CNeuroMod  
 888 uses the Schaefer-1000 atlas. This difference requires us to apply padding to align the feature  
 889 dimensions when evaluating on CNeuroMod. Furthermore, some concepts in CNeuroMod, such as  
 890 `telephone`, occur infrequently in BOLD5000, which limits NEURONA’s ability to generalize to  
 891 them. Nonetheless, we find that NEURONA maintains strong performance on queries such as action  
 892 decoding, suggesting meaningful transfer of motor-related neural representations across datasets.

893

894 Table 15: Cross-dataset generalization results, where models are trained on BOLD5000 and tested on  
 895 CNeuroMod.

|        | Overall       | Action        | Position      | T/F           |
|--------|---------------|---------------|---------------|---------------|
| UMBRAE | 0.4494        | 0.0106        | 0.0485        | 0.5036        |
| Ours   | <b>0.5535</b> | <b>0.7237</b> | <b>0.5246</b> | <b>0.5481</b> |

896

897

898

### A.7 NETWORK ABLATIONS

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

To evaluate how each functional network contributes to decoding performance, we perform a network-level ablation in which NEURONA was trained using only a constrained subset of Yeo-7 networks at a time, isolating the contribution of each network. In Table 16, we see that the top performing networks are the Default Mode, Control, Dorsal Attention, and Visual networks, all of which have been closely linked to visual processing and high-level perceptual representations in prior works Menon (2023); Kucyi et al. (2020); Miyawaki et al. (2008).

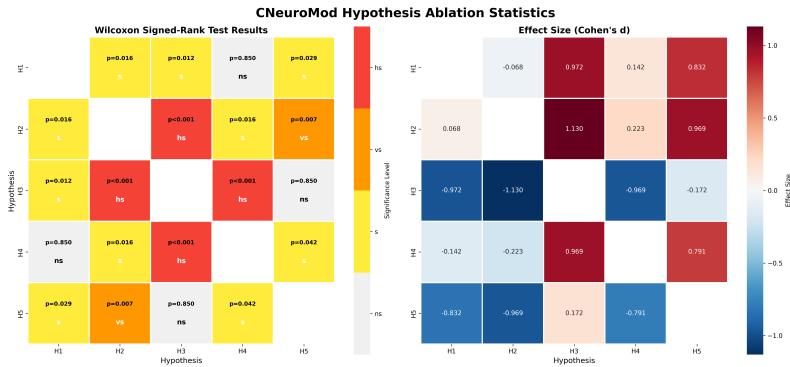


Figure 5: Statistical tests evaluating the effects of the five hypotheses across subjects in CNeuroMod.

Table 16: Network ablations to evaluate the contributions of each network to decoding performance.

| Selected network     | Overall             | Action              | Position            | T/F                 |
|----------------------|---------------------|---------------------|---------------------|---------------------|
| VisCent VisPeri      | $0.6750 \pm 0.0067$ | $0.3577 \pm 0.0459$ | $0.3765 \pm 0.0159$ | $0.7330 \pm 0.0072$ |
| SomMotA/B            | $0.6724 \pm 0.0069$ | $0.3480 \pm 0.0317$ | $0.3620 \pm 0.0131$ | $0.7326 \pm 0.0079$ |
| DorsAttnA/B          | $0.6753 \pm 0.0073$ | $0.3281 \pm 0.0786$ | $0.3850 \pm 0.0282$ | $0.7324 \pm 0.0076$ |
| SalVentAttnA/B       | $0.6709 \pm 0.0069$ | $0.2910 \pm 0.0655$ | $0.3839 \pm 0.0267$ | $0.7281 \pm 0.0060$ |
| LimbicA/B            | $0.6669 \pm 0.0101$ | $0.3379 \pm 0.0918$ | $0.3512 \pm 0.0603$ | $0.7282 \pm 0.0044$ |
| ContA/B/C            | $0.6785 \pm 0.0038$ | $0.3380 \pm 0.0808$ | $0.4129 \pm 0.0173$ | $0.7310 \pm 0.0035$ |
| DefaultA/B/C TempPar | $0.6804 \pm 0.0080$ | $0.3248 \pm 0.0752$ | $0.4357 \pm 0.0324$ | $0.7297 \pm 0.0066$ |
| All                  | $0.7102 \pm 0.0053$ | $0.5965 \pm 0.0322$ | $0.5378 \pm 0.0135$ | $0.7425 \pm 0.0057$ |

#### A.8 DETAILED CONCEPT ACCURACY

In Table 17 and Table 18, we report QA accuracy for unary and relational concepts separately across BOLD5000 and CNeuroMod, to analyze whether query structure affects performance. We see that that performance is generally stable across concept types, and across multiple atlases.

In Figure 6 and Figure 7, we illustrate confusion matrices for both BOLD5000 and CNeuroMod. We see that, as expected, while many concepts are reliably decoded (e.g., visually distinctive actions), some still cause confusion (e.g., spatial relations that are semantically similar).

Table 17: Accuracy breakdown between unary and relational concepts in BOLD5000.

| BOLD5000   | Overall | Unary | Relation |
|------------|---------|-------|----------|
| Yeo-7      | 0.727   | 0.717 | 0.751    |
| Yeo-17     | 0.740   | 0.732 | 0.760    |
| DiFuMo-64  | 0.728   | 0.727 | 0.728    |
| DiFuMo-128 | 0.732   | 0.733 | 0.730    |

#### A.9 PREDICATE ARGUMENT BINDING

We investigate whether unguided predicate representations contain decodable information about their arguments, by computing pairwise correlations between concept groundings across brain regions. For each relational query (e.g., `hold(person, baseball)`), we extract grounding score vectors across functional networks for the subject, object, unguided predicate, and guided predicate. We then compute the correlation matrix between these grounding patterns across all relational queries in both BOLD5000 (4 subjects) and CNeuroMod (3 subjects).

In Table 19, we report the mean correlations, and see minimal correlation between unguided predicate groundings and their subject ( $-0.0940$ ) or object ( $0.0048$ ) arguments in BOLD5000, with similarly

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

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

Probability Confusion Matrix for OVERALL Vocabulary Queries

| True          |      | Predicted |               |          |         |       |       |        |        |               |           |      |        |       |            |         |       |        |        |
|---------------|------|-----------|---------------|----------|---------|-------|-------|--------|--------|---------------|-----------|------|--------|-------|------------|---------|-------|--------|--------|
|               |      | _in -     | attached_to - | beside - | climb - | eat - | fly - | grow - | hold - | in_front_of - | next_to - | on - | ride - | sit - | surround - | under - | use - | walk - | wear - |
| _in -         | 54.2 | 12.5      | 16.7          | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 4.2    | 0.0           | 8.3       | 0.0  | 0.0    | 0.0   | 4.2        | 0.0     | 0.0   | 0.0    | 0.0    |
| attached_to - | 0.0  | 64.0      | 12.0          | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 4.0    | 0.0           | 12.0      | 0.0  | 0.0    | 8.0   | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| beside -      | 2.0  | 0.0       | 25.5          | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 45.1   | 0.0           | 7.8       | 0.0  | 0.0    | 9.8   | 9.8        | 0.0     | 0.0   | 0.0    | 0.0    |
| climb -       | 0.0  | 0.0       | 0.0           | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 0.0    | 0.0           | 0.0       | 0.0  | 0.0    | 100.0 | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| eat -         | 0.0  | 0.0       | 0.0           | 0.0      | 0.0     | 0.0   | 0.0   | 100.0  | 0.0    | 0.0           | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| fly -         | 0.0  | 0.0       | 0.0           | 0.0      | 0.0     | 100.0 | 0.0   | 0.0    | 0.0    | 0.0           | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| grow -        | 0.0  | 0.0       | 0.0           | 0.0      | 0.0     | 0.0   | 100.0 | 0.0    | 0.0    | 0.0           | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| hold -        | 0.0  | 0.0       | 0.0           | 0.0      | 0.0     | 0.0   | 0.0   | 100.0  | 0.0    | 0.0           | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| in_front_of - | 2.0  | 2.0       | 15.7          | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 58.8   | 0.0           | 7.8       | 0.0  | 0.0    | 3.9   | 9.8        | 0.0     | 0.0   | 0.0    | 0.0    |
| next_to -     | 7.1  | 10.7      | 32.1          | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 35.7   | 0.0           | 10.7      | 0.0  | 0.0    | 3.6   | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| on -          | 0.0  | 4.9       | 16.4          | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 4.9    | 0.0           | 70.5      | 0.0  | 0.0    | 1.6   | 1.6        | 0.0     | 0.0   | 0.0    | 0.0    |
| ride -        | 0.0  | 0.0       | 0.0           | 0.0      | 0.0     | 0.0   | 0.0   | 100.0  | 0.0    | 0.0           | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| sit -         | 0.0  | 0.0       | 0.0           | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 33.3   | 0.0           | 0.0       | 0.0  | 0.0    | 66.7  | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| surround -    | 7.1  | 7.1       | 35.7          | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 14.3   | 0.0           | 7.1       | 0.0  | 0.0    | 14.3  | 14.3       | 0.0     | 0.0   | 0.0    | 0.0    |
| under -       | 1.6  | 4.9       | 1.6           | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 14.8   | 0.0           | 0.0       | 0.0  | 0.0    | 0.0   | 77.0       | 0.0     | 0.0   | 0.0    | 0.0    |
| use -         | 0.0  | 0.0       | 0.0           | 0.0      | 0.0     | 0.0   | 0.0   | 100.0  | 0.0    | 0.0           | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| walk -        | 0.0  | 0.0       | 0.0           | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 50.0   | 0.0           | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 50.0    | 0.0   | 0.0    | 0.0    |
| wear -        | 0.0  | 0.0       | 0.0           | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 0.0    | 0.0           | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 100.0   | 0.0   | 0.0    | 0.0    |
| work -        | 0.0  | 0.0       | 0.0           | 0.0      | 0.0     | 0.0   | 0.0   | 100.0  | 0.0    | 0.0           | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |

Figure 6: Accuracy confusion matrix across concepts in BOLD5000.

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

Probability Confusion Matrix for OVERALL Vocabulary Queries

| True          |     | Predicted |               |          |         |       |       |        |        |               |           |      |        |       |            |         |       |        |        |
|---------------|-----|-----------|---------------|----------|---------|-------|-------|--------|--------|---------------|-----------|------|--------|-------|------------|---------|-------|--------|--------|
|               |     | _in -     | attached_to - | beside - | climb - | eat - | fly - | grow - | hold - | in_front_of - | next_to - | on - | ride - | sit - | surround - | under - | use - | walk - | wear - |
| _in -         | 1.3 | 0.0       | 30.1          | 0.0      | 0.0     | 42.9  | 0.0   | 0.0    | 0.0    | 21.7          | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 4.1   | 0.0    | 0.0    |
| attached_to - | 7.7 | 0.0       | 0.0           | 0.0      | 0.0     | 46.2  | 0.0   | 0.0    | 0.0    | 46.2          | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| beside -      | 1.7 | 0.0       | 12.4          | 0.0      | 0.0     | 49.4  | 0.0   | 0.0    | 0.0    | 20.0          | 0.0       | 0.0  | 0.0    | 0.0   | 0.3        | 0.0     | 0.0   | 16.2   | 0.0    |
| climb -       | 0.0 | 0.0       | 0.0           | 0.0      | 87.5    | 0.0   | 0.0   | 0.0    | 0.0    | 0.0           | 0.0       | 0.0  | 0.0    | 12.5  | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| eat -         | 0.0 | 0.0       | 0.0           | 0.0      | 66.9    | 0.0   | 0.0   | 0.0    | 10.1   | 0.0           | 0.0       | 0.0  | 0.0    | 18.2  | 0.0        | 8.8     | 0.0   | 0.0    | 0.0    |
| hold -        | 0.0 | 0.0       | 0.0           | 0.0      | 0.0     | 0.0   | 0.0   | 0.0    | 0.0    | 0.0           | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| in_front_of - | 2.1 | 0.0       | 8.3           | 0.0      | 0.0     | 53.6  | 0.0   | 0.0    | 0.0    | 20.8          | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 15.2   | 0.0    |
| next_to -     | 0.0 | 0.0       | 11.9          | 0.0      | 0.0     | 56.3  | 0.0   | 0.0    | 0.0    | 14.8          | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 17.0   | 0.0    |
| on -          | 1.2 | 0.0       | 10.1          | 0.0      | 0.0     | 17.7  | 0.0   | 0.0    | 0.0    | 66.3          | 0.0       | 0.0  | 0.0    | 0.0   | 2.0        | 0.0     | 0.0   | 2.7    | 0.0    |
| play -        | 0.0 | 0.0       | 0.0           | 0.0      | 42.9    | 0.0   | 0.0   | 0.0    | 0.0    | 0.0           | 14.3      | 0.0  | 21.4   | 0.0   | 21.4       | 0.0     | 0.0   | 0.0    | 0.0    |
| smile -       | 0.0 | 0.0       | 0.0           | 0.0      | 4.7     | 0.0   | 0.0   | 90.7   | 0.0    | 0.0           | 0.0       | 0.0  | 0.0    | 0.0   | 2.3        | 0.0     | 0.0   | 0.0    | 0.0    |
| stand -       | 0.0 | 0.0       | 0.0           | 0.0      | 8.4     | 0.0   | 0.0   | 0.7    | 0.0    | 0.0           | 0.0       | 18.4 | 0.0    | 72.2  | 0.0        | 0.4     | 0.0   | 0.0    | 0.0    |
| surround -    | 7.9 | 0.0       | 10.5          | 0.0      | 0.0     | 57.9  | 0.0   | 0.0    | 0.0    | 21.1          | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 2.6    | 0.0    |
| talk -        | 0.0 | 0.0       | 0.0           | 0.0      | 8.0     | 0.0   | 0.0   | 85.1   | 0.0    | 0.0           | 0.0       | 2.4  | 0.0    | 2.8   | 0.0        | 0.9     | 0.0   | 0.0    | 0.7    |
| touch -       | 0.0 | 0.0       | 0.0           | 0.0      | 0.0     | 27.4  | 0.0   | 0.0    | 32.3   | 0.0           | 0.0       | 0.0  | 14.5   | 0.0   | 19.4       | 0.0     | 6.5   | 0.0    | 0.0    |
| under -       | 4.2 | 0.0       | 0.7           | 0.0      | 0.0     | 27.1  | 0.0   | 0.0    | 0.0    | 3.2           | 0.0       | 0.0  | 0.0    | 0.0   | 0.0        | 0.0     | 0.0   | 64.8   | 0.0    |
| walk -        | 0.0 | 0.0       | 0.0           | 0.0      | 20.3    | 0.0   | 0.0   | 10.2   | 0.0    | 0.0           | 0.0       | 28.8 | 0.0    | 40.7  | 0.0        | 0.0     | 0.0   | 0.0    | 0.0    |
| wear -        | 0.0 | 0.0       | 0.0           | 0.0      | 0.0     | 62.5  | 0.0   | 0.0    | 0.0    | 0.0           | 4.2       | 0.0  | 4.2    | 0.0   | 29.2       | 0.0     | 0.0   | 0.0    | 0.0    |

Figure 7: Accuracy confusion matrix across concepts in CNeuroMod.

1013 low correlations in CNeuroMod (0.0230 for subject, 0.0019 for object). This suggests that when no structural guidance is provided, predicate representations do not bind to specific arguments.

1014 In contrast, guided predicate representations showed substantially higher correlations with both subjects (0.2020 in BOLD5000, 0.1095 in CNeuroMod) and objects (0.2975 in BOLD5000, 0.2555 in CNeuroMod). This significant increase in correlation indicates that NEURONA’s explicit guidance notably helps decode relations from predicate-argument interactions.

### A.10 FINE-GRAINED GENERALIZATION ANALYSES

1023 Our generalization experiments on both BOLD5000 (subject-CSI1) and CNeuroMod (subject-01) include argument swapping, predicate transfer, and role systematicity. In Table 20, we report fine-1024 grained systematicity tests across each of these settings. We see that NEURONA shows strong1025 performance on each category, with relatively small drops from all to unseen combinations.

1026

1027 Table 18: Accuracy breakdown between unary and relational concepts in CNeuroMod.

| CNeuroMod    | Overall | Unary | Relation |
|--------------|---------|-------|----------|
| Yeo-7        | 0.718   | 0.696 | 0.754    |
| Yeo-17       | 0.725   | 0.707 | 0.754    |
| Schaefer-100 | 0.725   | 0.708 | 0.754    |

1032

1033

1034 Table 19: Results of predicate argument binding.

|                                       | BOLD5000 | CNeuroMod |
|---------------------------------------|----------|-----------|
| Unguided-predicate / subject          | -0.0940  | 0.0230    |
| Unguided-predicate / object           | 0.0048   | 0.0019    |
| Unguided-predicate / guided-predicate | 0.0336   | 0.1488    |
| Guided-predicate / subject            | 0.2020   | 0.1095    |
| Guided-predicate / object             | 0.2975   | 0.2552    |

1041

1042

1043 Table 20: Fine-grained generalization analyses in BOLD5000 and CNeuroMod.

| Bold5000  | Argument swapping | Predicate transfer | Role systematicity | Other  | All (T/F) |
|-----------|-------------------|--------------------|--------------------|--------|-----------|
| Unseen    | 0.7451            | 0.75               | 0.7391             | 0.7096 | 0.7184    |
| All       | 0.7562            | 0.6667             | 0.8030             | 0.7393 | 0.7407    |
| CNeuroMod | Argument swapping | Predicate transfer | Role systematicity | Other  | All (T/F) |
| Unseen    | 0.7526            | 0.7333             | 0.7390             | 0.6842 | 0.6991    |
| All       | 0.7011            | 0.5088             | 0.7020             | 0.7307 | 0.7250    |

1052

1053

## A.11 ADDITIONAL NULL MODEL

1054

1055 For our consistency metric, we include a stronger null model that uses a multinomial distribution to  
 1056 randomly assign concept groundings across brain networks, preserving the exact total sample count  
 1057 for each concept while distributing samples uniformly across all networks. This ensures that the  
 1058 null model retains the same data structure as the observed data (i.e., same total grounding counts  
 1059 per concept) while randomizing only the network assignments. We run this null model 10 times and  
 1060 report the mean and standard deviation in Table 21. NEURONA similarly significantly outperforms  
 1061 the null model in grounding consistency.

1062

1063

Table 21: Consistency metric comparison with a multinomial null model.

|                    | BOLD5000            | CNeuroMod           |
|--------------------|---------------------|---------------------|
| Null (multinomial) | $0.6267 \pm 0.0067$ | $0.6424 \pm 0.0057$ |
| NEURONA            | $0.8475 \pm 0.0186$ | $0.8592 \pm 0.0084$ |

1068

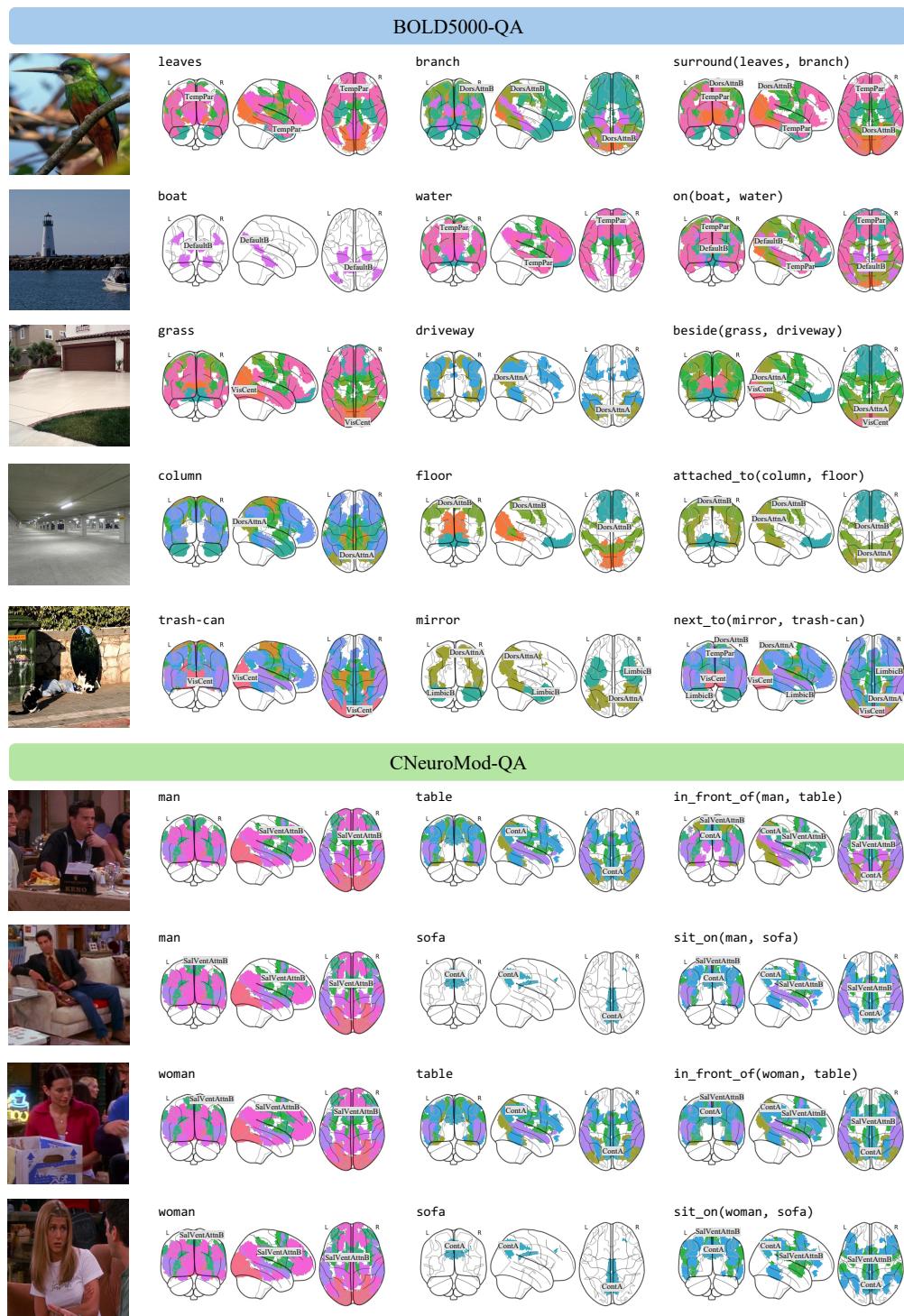
1069

## A.12 fMRI RETRIEVAL RESULTS

1070

1071

1072 Here, we detail results from an fMRI retrieval task, where we adapt NEURONA to retrieve the  
 1073 corresponding fMRI given a symbolic query. We use positive queries to ensure that the concept  
 1074 and fMRI are well aligned. Each concept is represented as a one-hot vector over the full concept  
 1075 vocabulary, from which a concept embedding is obtained using a small MLP. We then follow the  
 1076 structure of the symbolic expression by applying an aggregation operation to combine the specified  
 1077 concepts in the query. A parallel MLP encodes the fMRI input into an fMRI embedding, and we  
 1078 train the embedding spaces jointly using a CLIP-based contrastive loss. In this setting, we achieve a  
 1079 test top-1 retrieval accuracy of 0.1325 and a top-5 accuracy of 0.3012, substantially outperforming a  
 random-choice baseline (top-1: 0.0120, top-5: 0.0602).

1080 **B CONCEPT GROUNDING VISUALIZATIONS**  
10811082 In Figure 8, we present concept grounding examples from the BOLD5000 (Chang et al., 2019) and  
1083 CNeuroMod (Gifford et al., 2024; Boyle et al., 2023) datasets.  
10841132 Figure 8: We show examples of learned concept grounding by NEURONA on BOLD5000-QA and  
1133 CNeuroMod-QA, across subject, object, and predicate concepts.  
1134

1134  
1135

## C DATASETS

1136  
1137

### C.1 LICENSE FOR EXISTING DATASETS

1138  
1139  
1140

We train NEURONA on the BOLD5000 (Chang et al., 2019) and CNNeuroMod (Gifford et al., 2024; Boyle et al., 2023) datasets, which are both licensed under the Creative Commons 0 License. More information can be found on their websites: BOLD5000 and CNNeuroMod.

1141

### C.2 fMRI-QA DATASETS

1143

**BOLD5000.** We utilize the BOLD5000 dataset, which has been preprocessed and aligned with image stimuli following WAVE (Wang et al., 2024). The fMRI data has a shape of  $[5, 1024]$ , representing 5 TRs and 1024 brain regions. We use preprocessed, image-aligned fMRI data provided here. Each TR (repetition time) is 2 seconds, resulting in a chunk duration of 10 seconds ( $5 \times 2$ s). We account for a hemodynamic lag of 2 TRs. All four subject pairs are included, following the same train-test split as in previous studies.

1144

**CNeuroMod.** We use the CNNeuroMod dataset preprocessed by the Algonauts Challenge (Gifford et al., 2024; Boyle et al., 2023). The fMRI data has a TR of 1.49 seconds and a shape of  $[5, 1000]$ , representing 5 TRs and 1000 brain regions based on the Schaefer-1000 atlas (Schaefer et al., 2018). This yields a chunk duration of 7.45 seconds ( $5 \times 1.49$ s), with a hemodynamic lag of 3 TRs. For each chunk, we select the most motion-informative video frame by computing motion energy as the absolute difference between consecutive frames. The chunks are extracted from Friends episodes, with seasons 1–5 used for training and the unseen season 6 reserved for testing. We use 1,000 fMRI-video chunks per season, resulting in 5,000 training samples and 1,000 testing samples.

1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156

We provide additional examples from our datasets in Figure 9.

1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

1188

1189

1190

1191

1192

1193

1194



Q: Is there a dog?

A: Yes

## BOLD5000-QA



Q: Are there skis on the snow?

A: Yes



Q: Is there a bus on the road?

A: No



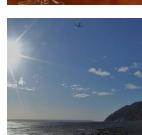
Q: What is the action between the person and the wine glass?

A: Hold



Q: What is the action between the sheep and the grass?

A: Eat



Q: Is there a moon?

A: No



Q: Is there a person riding the motorcycle?

A: No



Q: Is there a helmet?

A: Yes



Q: Is there a bear?

A: Yes



Q: Is there a woman sitting on the sofa?

A: No

1209

1210

1211

1212

1213

1214

1215

## CNeuroMod-QA



Q: Is there a bed?

A: Yes



Q: Is there a man in green?

A: Yes



Q: Is there a man sitting on the desk?

A: No



Q: What is the position between the door and the man?

A: Beside



Q: Is there a man looking at the painting?

A: No



Q: What is the position between the paper and the desk?

A: On



Q: What is the action between the man and the chair?

A: Sit



Q: Is there a woman?

A: Yes



Q: Is there a red shirt?

A: Yes



Q: What is the action between the man and the woman?

A: Talk

1240

1241

Figure 9: Examples of queries in BOLD5000-QA and CNeuroMod-QA.

1242 **D EXPERIMENT DETAILS**  
12431244 **D.1 TRAIN SETTINGS**  
12451246 We train and evaluate NEURONA on the specified training and test sets for both BOLD5000 (Chang  
1247 et al., 2019) and CNNeuroMod (Gifford et al., 2024; Boyle et al., 2023) datasets. Training is conducted  
1248 for 100 epochs using the Adam optimizer (Diederik, 2014), with learning rate 0.001 and batch size  
1249 32.  
12501251 **D.2 COMPUTE RESOURCES**  
12521253 Since NEURONA consists of a lightweight convolutional neural network for fMRI feature extraction  
1254 followed by a linear classifier for concept grounding and execution, its computational requirements  
1255 are minimal. All experiments are conducted on a single NVIDIA A100 GPU, with training completing  
1256 in approximately 30 minutes. Data loading is parallelized using 16 CPU workers, and the system  
1257 uses 64 GB of RAM.  
12581259 **D.3 BASELINE IMPLEMENTATIONS**  
12601261 We describe the implementation details of the baseline models compared in our study below. In all  
1262 methods, we treat the fMRI input as a sequence of length 5, with each time step as a token.  
12631264 **Linear** We tokenize the input query using the BERT tokenizer (Devlin et al., 2019) and pad all  
1265 sequences to a fixed length. The tokens are then encoded using a linear layer. We concatenate the  
1266 fMRI token sequence and the query token sequence, and apply a final linear classification layer to  
1267 predict the output (either a binary T/F answer or a vocabulary token).  
12681269 **SDRecon** We implement SDRecon (Takagi & Nishimoto, 2023) following the official repository\*.  
1270 A ridge regression model aligns fMRI features with image embeddings, which are then passed to a  
1271 VQA-GIT language model (Wang et al., 2022) to generate answers. We set the ridge regularization  
1272 parameter to  $\lambda = 20$ . A custom parser (described below) is used to map the generated language  
1273 response to a valid prediction.  
12741275 **BrainCap** BrainCap similarly uses a linear encoder to align fMRI features with visual embed-  
1276 dings (Ferrante et al., 2023). The aligned embeddings are passed to a BLIP language model (Li et al.,  
1277 2022) to generate answers. We apply the same parser as in SDRecon to extract final predictions from  
1278 the language output. We follow the implementation of the official repository†.  
12791280 **UMBRAE** UMBRAE leverages a transformer-based encoder to map fMRI features to image  
1281 embeddings (Xia et al., 2024). These embeddings are then passed to LLaVA (Liu et al., 2023) for  
1282 language-based prediction, followed by response parsing. We follow the implementation of the  
1283 official repository‡.  
12841285 We implement a rule-based parser to convert language model outputs into structured predictions.  
1286 The parser first cleans the text by removing punctuation, digits, and formatting inconsistencies. It  
1287 identifies binary answers (“yes” or “no”) when the query requires. For other queries, it extracts  
1288 the first valid word from a predefined vocabulary. If no valid word is found, it defaults to the most  
1289 common answers of “on” for spatial queries or “hold” otherwise. For all image-grounded baselines,  
1290 the ground-truth image embeddings are derived from the visual encoder of a vision-language model  
1291 for BOLD5000. We use the embedding of the first selected video frame as the ground truth for  
1292 CNNeuroMod.  
12931294 \*<https://github.com/yu-takagi/StableDiffusionReconstruction>  
1295 †<https://github.com/enomodnara/BrainCaptioning>  
1296 ‡<https://github.com/weihaox/UMBRAE>

1296 D.4 ABLATION DETAILS  
12971298 In this section, we provide the full definitions of the entities and grounding hypotheses introduced in  
1299 the main paper.  
13001301 D.4.1 ENTITY PROCESSING  
13021303 To enable concept grounding to neural activity  $f \in \mathbb{R}^{N \times T}$ , we first map the fine-grained networks  
1304  $N = 1024$  to  $P$  functional networks defined by the given atlas. This results in  $P$  network-specific  
1305 fMRI signals  $\{f_1, \dots, f_P\}$ , where each  $f_p \in \mathbb{R}^{m_p \times T}$  represents the aggregated signal from  $m_p$   
1306 fine-grained regions assigned to network  $p$ . Since the number of regions  $m_p$  vary across networks, we  
1307 apply a linear stitcher to project each  $f_p \in \mathbb{R}^{m_p \times T}$  to a fixed-dimensional representation  $e_p \in \mathbb{R}^{d \times T}$ ,  
1308 where  $d = 256$ , using network-specific linear projections  $W_p \in \mathbb{R}^{m_p \times d}$ , such that  $e_p = W_p^\top f_p$ . This  
1309 produces a unified set  $\mathcal{E}$  of  $P$  embeddings  $\{e_1, \dots, e_P\}$ , which are then processed by a small 1-D  
1310 convolutional encoder to form parcellation embeddings. From these base embeddings, we propose  
1311 hypotheses of candidate entities from which concepts can be grounded.  
13121313 D.4.2 GENERAL GROUNDING FORMULATION  
13141315 We define single-region (unary) grounding for a concept  $c$  as  $G_{\text{unary}}(c)$  and multi-region (binary)  
1316 grounding as  $G_{\text{binary}}(c)$ . We further define a fused grounding score combining unary and binary  
1317 components:  
1318

1319 
$$G(c) = G_{\text{unary}}(c) + \frac{1}{P} \sum_{i=1}^P G_{\text{binary}}(c_i). \quad (4)$$
  
1320

1321 D.4.3 HYPOTHESES DEFINITIONS  
13221323 Let  $c_p$ ,  $c_s$ , and  $c_o$  be the concepts for predicate (e.g., holding), subject (e.g., person), and object  
1324 (e.g., baseball–bat), respectively.  
13251326 **H1: Single-region grounding.** Concepts are grounded to a single brain region:  
1327

1328 
$$\begin{aligned} G_{\text{H1}}(c_p) &= G_{\text{unary}}(c_p), \\ G_{\text{H1}}(c_s) &= G_{\text{unary}}(c_s), \\ G_{\text{H1}}(c_o) &= G_{\text{unary}}(c_o). \end{aligned} \quad (5)$$
  
1329

1332 **H2: Multi-region co-activation.** Concepts are grounded through co-activation across region pairs:  
1333

1334 
$$\begin{aligned} G_{\text{H2}}(c_p) &= G_{\text{binary}}(c_p), \\ G_{\text{H2}}(c_s) &= G_{\text{unary}}(c_s), \\ G_{\text{H2}}(c_o) &= G_{\text{unary}}(c_o). \end{aligned} \quad (6)$$
  
1335

1338 **H3: Predicate conditioned on subject.** Predicate representations are guided by the activation of  
1339 the subject region:  
1340

1341 
$$\begin{aligned} G_{\text{H3}}(c_p) &= G_{\text{binary}}(c_p) + G_{\text{unary}}(c_s), \\ G_{\text{H3}}(c_s) &= G_{\text{unary}}(c_s), \\ G_{\text{H3}}(c_o) &= G_{\text{unary}}(c_o). \end{aligned} \quad (7)$$
  
1342

1345 **H4: Predicate conditioned on object.** Predicate representations are guided by the activation of the  
1346 object region:  
1347

1348 
$$\begin{aligned} G_{\text{H4}}(c_p) &= G_{\text{binary}}(c_p) + G_{\text{unary}}(c_o), \\ G_{\text{H4}}(c_s) &= G_{\text{unary}}(c_s), \\ G_{\text{H4}}(c_o) &= G_{\text{unary}}(c_o). \end{aligned} \quad (8)$$
  
1349

1350  
 1351 **H5: Full argument-guided grounding.** Our proposed method combines multi-region grounding  
 1352 with subject and object guidance. The grounding scores are defined as:  
 1353

$$\begin{aligned} G_{H5}(c_p) &= G(c_p) + G(c_s) + G(c_o), \\ G_{H5}(c_s) &= G(c_s), \\ G_{H5}(c_o) &= G(c_o). \end{aligned} \tag{9}$$

1355  
 1356 These formulations enable systematic evaluation of how structural priors affect downstream neural  
 1357 decoding accuracy.  
 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

1404 **E ETHICS STATEMENT**  
14051406 This work uses only publicly released fMRI datasets, and focuses on decoding concepts from visual  
1407 stimuli rather than personal traits or identity, which reduces the risk of individual fingerprinting.  
1408 However, we acknowledge the inherent risks in building architectures that enable neural decoding  
1409 from participants. We emphasize that our method is intended solely for scientific analysis, and should  
1410 not be used for identification or inference of sensitive personal attributes from neural activity.  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457