INSTRUCTION-TUNED VIDEO-AUDIO MODELS ELUCIDATE FUNCTIONAL SPECIALIZATION IN THE BRAIN

Anonymous authors

000

001

002 003 004

010 011

012

013

014

016

017

018

019

021

025

026

028

029

031

034

037

038

040

041

042

043

044

045

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Recent voxel-wise multimodal brain encoding studies have shown that multimodal large language models (MLLMs) exhibit a higher degree of brain alignment compared to unimodal models in both unimodal and multimodal stimulus settings. More recently, instruction-tuned multimodal models have shown to generate taskspecific representations that align strongly with brain activity. However, prior work evaluating the brain alignment of MLLMs has primarily focused on unimodal settings or relied on non-instruction-tuned multimodal models for multimodal stimuli. To address this gap, we investigated brain alignment, that is, measuring the degree of predictivity of neural activity recorded while participants were watching naturalistic movies (video along with audio) with representations derived from MLLMs. We utilized instruction-specific embeddings from six video and two audio instruction-tuned MLLMs. Experiments on 13 video task-specific instructions show that instruction-tuned video MLLMs significantly outperform in-context learning multimodal models (by ~9%), non-instruction-tuned multimodal models (by $\sim 15\%$) and unimodal models (by $\sim 20\%$). Our evaluation of MLLMs for both video and audio tasks using language-guided instructions shows clear disentanglement in task-specific representations from MLLMs, leading to precise differentiation of multimodal functional processing in the brain. We also find that MLLM layers align hierarchically with the brain, with early sensory areas showing strong alignment with early layers, while higher-level visual and language regions align more with middle to late layers. These findings provide clear evidence for the role of task-specific instructions in improving the alignment between brain activity and MLLMs, and open new avenues for mapping joint information processing in both the systems.

1 Introduction

The alignment between internal representations of multimodal Transformer models and cortical activation patterns obtained from naturalistic stimuli has emerged as a key focus in the study of brainmodel correspondence. Recent research has demonstrated that multimodal models in brain encoding can be broadly categorized into two settings (see Appendix A Table 4): (i) multimodal models evaluated with unimodal stimuli (Doerig et al., 2022; Wang et al., 2023; Oota et al., 2022b; Popham et al., 2021; Tang et al., 2024; Oota et al., 2025a; Srijith et al., 2025), and (ii) multimodal models evaluated with multimodal stimuli (Nakagi et al., 2024; Subramaniam et al., 2024; Dong & Toneva, 2023a; Oota et al., 2025b; Sartzetaki et al., 2025). In the former setting, brain recordings are obtained from unimodal image stimuli, but representations from multimodal models, which integrate modalities such as vision and language, achieve a higher degree of brain alignment compared to vision-only models (Doerig et al., 2022; Wang et al., 2023; Oota et al., 2022b; Popham et al., 2021). This observation holds true to the new class of instruction-tuned multimodal large language models (MLLMs), especially when prompted with natural instructions (Oota et al., 2025a). In the latter setting, where brain recordings are obtained from multimodal stimuli (e.g., watching movies with both visual and auditory stimuli), studies show that multimodal models exhibit higher degree of brain alignment over unimodal models (Dong & Toneva, 2023a; Oota et al., 2025b). While prior studies have examined brain alignment with instruction-tuned MLLMs (IT-MLLMs), they have largely been limited to unimodal stimuli, or have used non-instruction-tuned models in the context of multimodal stimuli. In this work, we bridge this gap by systematically investigating IT-MLLMs in the presence of rich

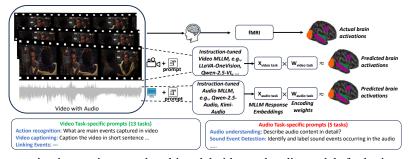


Figure 1: Leveraging instruction-tuned multimodal video and audio models for brain encoding with a diverse set of instructions. For the given movie clip, we can obtain different multimodal representations using instructions that ask the model to (i) generate the caption of the video, (ii) identify whether temporal events are present, (iii) determine the primary colors dominant in the video, etc. Using instruction-specific representations (X), we estimate the alignment using a simple linear function f (ridge regression), which maps MLLM representations to brain recordings. Here, W denotes voxelwise encoding model weights.

multimodal stimuli. Specifically, we assess how well representations elicited through naturalistic, task-specific instructions involving both video and audio align with brain activity across the cortical hierarchy, from early sensory regions to higher-order cognitive areas.

Several unimodal studies report that task-specific fine-tuned Transformer models better align with brain activity during text (Oota et al., 2022a; Aw & Toneva, 2023; Sun & Moens, 2023; Oota et al., 2024b), speech (Oota et al., 2023; Tuckute et al., 2023; Oota et al., 2024a), and vision (Wang et al., 2019; Conwell et al., 2022) processing, outperforming pretrained models in brain predictivity. However, these models are task-specific, limiting generalization, requiring separate data and training per task. Instruction-tuning (Xu et al., 2023; Dai et al., 2023; Liu et al., 2024) offers a scalable alternative, fine-tuning a single LLM across diverse NLP tasks and surpassing task-specific models (Taori et al., 2023; Touvron et al., 2023; Jiang et al., 2023; Abdin et al., 2024; Dubey et al., 2024), while showing stronger brain alignment (Sun et al., 2023; Sun & Moens, 2023; Loong Aw et al., 2024) (see Appendix B for more.) Building on this, recent work aligns IT-MLLMs with brain data for text (Benara et al., 2024) and images (Oota et al., 2025a), though limited to unimodal stimuli. Motivated by advances in multimodal MLLMs for video and audio tasks, we ask: Do instructiontuned video/audio MLLMs prompted with natural language yield better brain alignment than their pretrained in-context learning and non-instruction-tuned counterparts, while also distinguish taskspecific representations? To our knowledge, this is the first study to use such MLLMs to model fMRI responses across video and audio tasks (workflow in Fig. 1).

Using brain recordings from participants watching several popular movies with audio (Boyle et al., 2020), we investigate the brain alignment of IT-MLLMs. Specifically, we evaluate six video IT-MLLMs, two audio IT-MLLMs, two pretrained video MLLMs with in-context learning, two non-instruction-tuned multimodal models (video+audio), two unimodal models for video and one unimodal model for audio. These models are probed with 13 video task-specific instructions, and 5 audio task-specific instructions. Overall, this study addresses the following research questions: (1) How do different task-specific instructions influence the degree of brain alignment in instruction-tuned video and audio MLLMs? (2) Do instruction-tuned video MLLMs exhibit better brain alignment than their audio counterparts when exposed to multimodal stimuli? (3) Do IT-MLLMs produce functionally distinct representations that map onto different brain regions, offering a data-driven alternative to traditional experimental stimuli? (4) How do task instructions related to semantic categories (e.g., narrative understanding, spatial reasoning) explain differential activation across language, auditory, and visual brain regions?

To further quantify how IT-MLLMs capture shared and distinct neural processes across tasks, we use a variance partitioning approach. This analysis reveals the unique and overlapping contributions of individual task-specific representations to brain responses, enhancing our understanding of the brain's functional organization in processing multimodal information.

Our analysis of IT-MLLMs and brain alignment with multimodal stimuli reveals several key conclusions: (i) Video-based IT-MLLMs show significantly higher brain alignment than audio-based IT-MLLMs, pretrained in-context learning MLLMs, non-instruction-tuned multimodal models, as

well as unimodal video and audio models. This holds across the whole brain, as well as within language, visual and auditory regions. (ii) On the other hand, Audio MLLMs outperform both noninstruction-tuned multimodal and unimodal models only in the auditory cortex (AC) and middle frontal gyrus (MFG) language regions, while exhibiting comparable performance in other languagerelated areas. (iii) Surprisingly, both video and audio MLLMs generate task-specific representations based on task-instructions and effectively differentiate functional processing across brain regions. For example, audio understanding and captioning tasks show stronger alignment with language areas, while sound event detection aligns with the auditory cortex and temporal lobe. (iv) Grouping 13 video tasks into 5 semantic categories reveals strong alignment of MLLM representations with brain sub-regions in line with the existing literature. Tasks involving language and narrative understanding exhibit stronger alignment in language-related sub-regions such as angular gyrus and lateral temporal regions, consistent with prior findings on event structure representation in naturalistic stimuli (Baldassano et al., 2017). Similarly, spatial understanding tasks engage regions of the dorsal visual pathway, particularly the intraparietal sulcus and surrounding parietal cortex. Overall, our analysis reveals that IT-MLLMs capture both hierarchical and task-specific brain representations, making them powerful tools for studying functional specialization and bridging cognitive modeling with neuroscience. Our code is part of the supplementary material.

2 Dataset and Models

2.1 Brain Imaging Dataset

We experiment with Movie10 (Boyle et al., 2020), a multimodal naturalistic fMRI dataset, obtained from the Courtois NeuroMod databank. This dataset was collected while four human subjects (s1, s2, s3, s5; data for s4 and s6 is not public) passively watched four different movies: The Bourne supremacy (\sim 100 mins), The wolf of wall street (\sim 170 mins), Hidden figures (\sim 120 mins) and Life (\sim 50 mins). Among these, Hidden figures and Life are repeated twice, with the repeats used for testing and the remaining movies for training. We use Life movie for testing where we average the two repetitions to reduce noise. This is among the largest publicly available multimodal fMRI datasets by samples per participant, with 4024 TRs (Time Repetitions) of The Bourne supremacy and 6993 TRs of The wolf of wall street for training and 2013 TRs of Life for test. Train and test sets are totally disjoint. The fMRI data is collected every 1.49 seconds (= 1 TR).

The dataset is already preprocessed and projected onto the surface space ("fsaverage6"). We use the multimodal parcellation of the human cerebral cortex based on the Glasser Atlas (which consists of 180 regions of interest in each hemisphere) to report the ROI (region of interest) analysis for the brain maps (Glasser et al., 2016). This includes four visual processing regions (early visual cortex (EVC), object-related areas (LOC), face-related areas (OFA) and scene-related areas (PPA)), one early auditory area (AC), and eight language-relevant regions, encompassing broader language regions: angular gyrus (AG), anterior temporal lobe (ATL), posterior temporal lobe (PTL), inferior frontal gyrus (IFG), inferior frontal gyrus orbital (IFGOrb), middle frontal gyrus (MFG), posterior cingulate cortex (PCC) and dorsal medium prefrontal cortex (dmPFC), based on the Fedorenko lab's language parcels (Milton et al., 2021; Desai et al., 2023). We show the flatmap with these labeled ROIs in Appendix Fig. 6 and list the detailed sub-ROIs of these ROIs in Appendix C.

Estimating cross-subject prediction accuracy. To account for the intrinsic noise in biological measurements, we adapt Schrimpf et al. (2021)'s method to estimate the cross-subject prediction accuracy for a model's performance for the Movie10 fMRI dataset. Each subject $s \in ([1,4])$ is chosen as the prediction target and the other three are used to predict this target. We use a voxel-wise encoding model (see Section 3) to predict one participant's response from others. The detailed approach is described in Appendix D. Note that the estimated cross-subject prediction accuracy is based on the assumption of a perfect model, which might differ from real-world scenarios, yet offers valuable insights into model's performance. We present the cross-subject prediction accuracy across voxels for the Movie10 fMRI dataset for each of the four participants in Appendix D. The plots show that across all participants higher activity is observed in the language and visual regions with a max correlation up to 0.4 implying that data has low noise and low cross-subject variability.

2.2 Instruction-tuned Multimodal Models for Video and Audio

To investigate whether IT-MLLMs models, when prompted using natural language-guided instructions, align with the way humans process multimodal information in the brain, we consider six popular modern instruction-tuned video MLLMs (InstructBLIPVideo (Dai et al., 2023), Video-

Table 1: Pretrained MLLMs for video, audio vs. multimodal, unimodal models (IT: Instruction-tuned) (IC: In-context learning).

C. III-context learning).						
Model Name	IT	#Layers	Modality			
InstructBLIPVideo	1	33	Video+Text			
Video-LLaVA	/	33	Video+Text			
LLaVa-NeXT-Video	/	33	Video+Text			
Qwen-2.5-VL	1	29	Video+Text			
Videochat-R1	/	29	Video+Text			
LLaVA-OneVision	/	28	Video+Text			
Qwen-2.5-Audio	1	29	Audio+Text			
Kimi-Audio	1	29	Audio+Text			
Qwen-2.5-Omni (IC)	×	29	Video+Audio+Text			
InternVL (IC)	×	29	Video+Text			
VILA	×	29	Video+Audio			
TVLT	×	12	Video+Audio			
VideoMAE	×	24	Video			
TimeSFormer	×	12	Video			
AST	×	24	Audio			

Table 2: Instructions for various multimodal audio tasks.

Task	Description
Audio Understanding	Can you describe the audio content in detail?
Audio Comprehension	What are people doing in the audio?
Audio Captioning	Caption the audio in a short sentence.
Sound Event Detection	Identify and label the sound events occurring in the audio.
Speaker Identification	Who is speaking in the audio?

LLaVA (Lin et al., 2024), LLaVA-Next-Video (Zhang et al., 2024), Qwen-2.5-VL (Wang et al., 2024), Videochat-R1 (Li et al., 2025), LLaVA-OneVision (Li et al., 2025)) and two instruction-tuned audio MLLMs (Qwen-2.5-Audio (Chu et al., 2024), Kimi-Audio (Kimi Team, 2024)). We also experiment with two pretrained video MLLMs with in-context learning (Qwen-2.5-Omni (Xu et al., 2025) and InternVL (Chen et al., 2024)), two non-instruction-tuned multimodal (VILA (Lin et al., 2023) and TVLT (Tang et al., 2022)), two video unimodal models (VideoMAE (Tong et al., 2022) and TimeSFormer (Bertasius et al., 2021)), and one audio unimodal (AST (Baade et al., 2022)) model. Details for these models are reported in Table 1.

2.3 NATURAL LANGUAGE INSTRUCTIONS AND FEATURE EXTRACTION FROM IT-MLLMS

Video-specific tasks. To ensure the diversity of task-specific instructions while considering videos as input, we consider 13 instructions, as shown in Table 3, and extract the language-guided representations from multimodal instruction-tuned video models. This set of 13 tasks are inspired from VideoInstruct100K dataset (Maaz et al., 2024). We borrowed those tasks, which are generally applicable to any video regardless of the contents in the image frames. We provide a sample of generated outputs for all the six video MLLMs in Tables 5, 6, 7, 8, 9 and 10 in Appendix E.

To extract instruction-specific representations from multimodal instruction-tuned video models for the brain encoding task, we input a video and task instruction to obtain the embeddings for the language-guided instruction. For in-context learning models, a video is paired with a natural language prompt without instruction tuning. For TVLT and VILA, we input video and audio. For TimesFormer and VideoMAE we input video only. We perform zero-shot inference on these models. For all multimodal instruction-tuned video models, we use the pretrained Transformer weights, which generate hidden state representations at each layer. We then average these hidden state representations at layer level of output generated tokens to obtain final embedding at each layer for each video with respect to task instruction.

Audio-specific tasks. Similar to video tasks, we consider five natural instructions while considering audio as input, as shown in Table 2, and extract the language-guided representations from multimodal instruction-tuned audio model. We provide a sample of generated outputs for one of the instruction-tuned audio models across the five tasks in Tables 11 and 12 in Appendix E.

Similar to instruction-tuned video models, to extract instruction-specific representations from the multimodal instruction-tuned audio model for the brain encoding task, we input a audio and task instruction to obtain the embeddings for language-guided instruction. For AST we input audio only. We follow similar feature extraction method as video-tasks to extract audio task representations.

3 METHODOLOGY

Voxel-wise encoding model. We train banded ridge regression based voxel-wise encoding models (la Tour et al., 2022) to predict the fMRI brain activity associated with the stimulus representations obtained from 13 task-specific instructions from multimodal instruction-tuned video models. Banded ridge regression optimizes a different regularization hyperparameter per feature space, and decomposes the explained variance over feature spaces. This decomposition helps in identify-

Table 3: Instructions for various multimodal video tasks.

Task	Description
Action Recognition	What are the main events captured in the video?
Video Understanding	Can you describe the video content in detail?
Visual Question Answering	How many people are in the video, and what are they doing?
Video Captioning	Caption the video in a short sentence.
Object and Scene Recognition	What are the main objects and people visible in the video? Describe each one briefly.
Commonsense Reasoning	Why did the character take this action? What could have motivated them to do this?
Spatial Understanding	Where is this video taken from? What place/landmark is shown in the video?
Temporal Ordering	Step-by-step describe the activity shown in the video.
Video reasoning	What is unusual about this video?
Narrative Understanding	Summarize the main storyline of the movie. What is the central conflict, and how is it resolved?
Emotion and Sentiment Analysis	What emotions do the characters express during the video? How does the video make you feel overall?
Global Appearance	Describe changes in characters' appearances throughout the video, including any noticeable outfit changes.
Linking Events	Explain how an early event in the video influences later developments.

ing which task-specific instruction contributes most to the explainable variance in different brain regions. Overall, banded ridge regression helps to accurately identify the contribution of each task-specific instruction, leading to better prediction accuracy and better interpretability. We employ z-score thresholding separately for both input stimulus representations and brain recordings for training and test datasets. For each subject, we account for the delay in the hemodynamic response by modeling hemodynamic response function using a finite response filter (FIR) per voxel with 5 temporal delays (TRs) corresponding to \sim 7.5 seconds (Huth et al., 2022). Formally, at each time step t, we encode the stimuli as $X_t \in \mathbb{R}^D$ and brain region voxels $Y_t \in \mathbb{R}^V$, where D denotes the dimension of the concatenation of delayed 5 TRs, and V denotes the number of voxels. Overall, with N such TRs, we obtain N training examples. Detailed hyper-parameter settings are in Appendix F.

Evaluation metrics. We evaluate our models using Pearson Correlation (PC), which is a standard metric for evaluating brain alignment (Jain & Huth, 2018; Schrimpf et al., 2021; Goldstein et al., 2022). Let TR be #time repetitions in the test set. Let $Y = \{Y_i\}_{i=1}^{TR}$ and $\hat{Y} = \{\hat{Y}_i\}_{i=1}^{TR}$ denote actual and predicted value vectors for a single voxel. Thus, Y and $\hat{Y} \in \mathbb{R}^{TR}$. We use PC to compute the correlation function, $corr(Y, \hat{Y})$. The final measure of a model's performance is obtained by calculating Pearson's correlation between the model's predictions and neural recordings. To quantify the model predictions, the resulting model prediction correlations are divided by the estimated cross-subject prediction accuracy; and averaged across voxels, regions, and participants, resulting in a standardized measure of performance referred to as normalized brain alignment. For calculating normalized alignment, we select the voxels with cross-subject prediction accuracy ≥ 0.05 .

4 RESULTS

4.1 INSTRUCTION-TUNED VIDEO MLLMS REPRESENTATIONS ALIGN WELL WITH BRAIN ACTIVITY ACROSS WHOLE BRAIN, LANGUAGE, VISUAL AND AUDITORY REGIONS

First, we examine the brain alignment by measuring the degree of brain predictivity using representations extracted from instruction-tuned video MLLMs, focusing on whole brain, language, visual and auditory regions. For each instruction-tuned MLLM, we calculate the average normalized brain alignment across 13 tasks, multiple subjects, and best MLLM layer, using the Movie10 fMRI dataset. Similarly, for instruction-tuned Audio MLLMs, we calculate the average normalized brain alignment across five tasks, multiple subjects, and best MLLM layer. Additionally, we report the brain alignment performance of in-context learning video MLLMs, non-instruction-tuned multimodal models, unimodal video models, and unimodal audio model (AST). We treat the non-instruction-tuned multimodal models and unimodal models (audio and video) as the baselines when comparing against the IT-MLLMs.

Whole brain analysis. Fig. 2 (a) shows the results for whole brain analysis. We make the following observations: (i) At the whole-brain level, the Wilcoxon signed-rank test reveals that the differences in brain alignment between instruction-tuned video MLLMs and in-context learning models, the non-instruction-tuned multimodal and unimodal models are statistically significant. In particular, all instruction-tuned video MLLMs achieve over $\sim 9\%$ improvement in brain alignment compared to in-context learning models, and $\sim 15\%$ improvement compared to other baselines. This contrasts with prior findings on instruction-tuned image-based MLLMs, which demonstrated comparable performance to multimodal models when evaluated on unimodal image stimuli (Oota et al., 2025a),

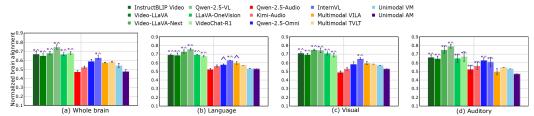


Figure 2: Average normalized brain alignment of instruction-tuned video MLLMs vs instruction-tuned audio MLLMs vs in-context learning video MLLMs vs multimodal and unimodal models across whole brain, language, visual and auditory regions. Error bars indicate the standard error of the mean across participants. * implies that instruction-tuned MLLM embeddings are significantly better than multimodal models and \land means that instruction-tuned MLLM embeddings are significantly better unimodal models with p ≤ 0.05 .

suggesting that instruction-tuned video MLLMs are more effective at capturing brain-relevant representations. (ii) Instruction-tuned audio MLLM embeddings show less alignment compared to non instruction-tuned multimodal and unimodal video models. These findings imply that instruction-tuned video MLLM models capture brain-relevant representations and contain additional information beyond the in-context learning, non-instruction-tuned multimodal and unimodal models.

Language, visual and auditory region analysis. We also present the average normalized brain alignment across language, visual and auditory regions in Fig. 2 (b, c & d). The results from Wilcoxon signed-rank test is consistent with whole-brain performance both in the language and visual regions i.e instruction-tuned video MLLMs embeddings exhibit significantly higher alignment in both language and visual regions compared to in-context learning video MLLMs, non-instruction-tuned multimodal, unimodal video, and audio models. On the other hand, instruction-tuned audio MLLM embeddings show significant alignment primarily in the auditory cortex and the middle frontal gyrus; when compared to non-instruction-tuned multimodal and unimodal models. Results for detailed language, visual and auditory sub-regions are shown in Fig. 8 and 9 in Appendix H.

These results suggest that instruction-tuned video MLLMs more effectively capture brain-relevant multimodal representations, particularly when processing naturalistic multimodal stimuli.

Additionally, we present contrast of brainmaps to display the average normalized brain alignment across voxels. Figs. 10 and 11 in Appendix I compare instruction-tuned video MLLMs with incontext learning video MLLMs (InternVL and Qwen-2.5-Omni, respectively). Figs. 12, 13, 14, 15, and 16 in Appendix J compare instruction-tuned video MLLMs with the non-instruction-tuned multimodal VILA and TVLT. The results show that instruction-tuned video MLLMs consistently achieve significantly higher alignment across all brain voxels. However, Figs. 17 & 18 in Appendix J reveal clear differences between audio MLLMs and multimodal models: the prediction performance of audio MLLMs lacks brain-relevant semantic information compared to multimodal models.

4.2 VIDEO AND AUDIO IT-MLLMS SUCCESSFULLY DIFFERENTIATE TASK-SPECIFIC INSTRUCTIONS

To investigate which instructions are more effective in predicting brain activity and whether IT-MLLMs differentiate task-specific representations and provide clear separation in brain regions, we analyze the voxels as follows. For each voxel, we select the instruction that results in the highest normalized brain alignment and apply the instruction-specific color code to the voxel.

Instruction-tuned video MLLMs. Fig. 3 (left) shows brain maps for Qwen-2.5-VL for video tasks for average normalized brain predictivity across subjects where the voxel color codes are projected onto the flattened cortical surface of the 'fsaverage' subject. The color-scheme corresponding to each instruction is also reported. We make the following observations: (i) Video understanding exhibits the strongest alignment across the whole brain. (ii) Tasks such as spatial understanding, narrative understanding, and visual question answering show higher alignment in language-related regions, including the angular gyrus, posterior temporal lobe, and visual regions. (iii) Higher-order language regions in the frontal cortex are predominantly identified by the video understanding task, with a smaller proportion of voxels also activated by video reasoning and temporal ordering tasks.

These findings suggest that instruction-tuned video MLLMs not only capture modality-specific representations (e.g., visual, linguistic), but also encode task-specific instructions involving semantic

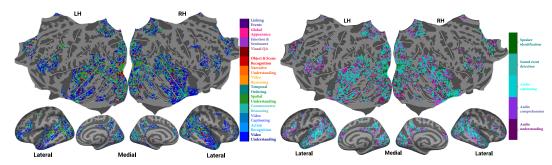


Figure 3: Each voxel is color-coded with the instruction that led to the highest normalized brain alignment. The color bar highlights color codes for each instruction. The voxels are projected onto the flattened cortical surface of the 'fsaverage' subject. (Left): video MLLM (Qwen-2.5-VL). (Right): audio MLLM (Qwen-2.5-Audio).

integration and event structure (like video understanding). This highlights that these models can encode complex neural patterns. We observe similar performance gains in other instruction-tuned video MLLMs, flatmaps showing task-specific encoding performance for average of subjects are shown in Figs. 19 and 20 in Appendix K.

Instruction-tuned audio MLLMs. Fig. 3 (right) shows brainmap for Instruction-tuned audio MLLM (Qwen-2.5-Audio) where the predictions are average across subjects. The voxel color codes are projected onto the flattened cortical surface of the 'fsaverage' subject. There is a clear distinction between different audio tasks. Audio captioning and sound detection are aligned with the auditory cortex (AC), while audio understanding activates higher-level regions like the inferior temporal (IT) cortex and inferior frontal gyrus (IFG). In contrast, speaker identification shows very sparse and scattered alignment, with some unexpected activation in the primary visual cortex (V1), suggesting it does not strongly reflect brain-relevant semantic processing. Fig. 21 in Appendix K shows similar brainmap for Kimi-Audio.

IT-MLLMs capture layer-wise cortical hierarchy. Inspired from previous literature (Namburi et al., 2023; Mitchell et al., 2022) which shows that Transformers process information differently across layers, we examine whether IT-MLLMs reflect the brain's hierarchy of information processing across layers by analyzing the voxels as follows. For each voxel, we select the layer that results in the highest normalized brain alignment and apply a color code for the 29/33 layers for each MLLM. Fig. 4 presents brain maps for the Qwen-2.5-VL & Qwen-2.5-Audio, where the voxels with their corresponding color codes are projected onto the flattened cortical surface of the 'fsaverage' subject. We make the following observations: (i) Early sensory areas-including early visual regions and early auditory cortex-are best aligned with the lower layers of the model, suggesting that shallow model representations capture low-level sensory features. (ii) High-level visual areas such as the lateral occipital complex (LOC) and parahippocampal place area (PPA), as well as language-related regions like the superior temporal sulcus and angular gyrus, show stronger alignment with the middle to deeper layers of the model. This reflects the model's progression toward more abstract and semantically rich representations. (iii) Notably, language-related areas such as the inferior frontal gyrus (IFG), anterior temporal lobe (ATL), and angular gyrus show strongest alignment with the deepest layers of the model. These results indicate that IT-MLLMs naturally develop a layered structure that maps well onto the brain's own representational hierarchy. Similar brain maps for the remaining models are provided in Fig. 22 in Appendix L.

4.3 REPRESENTATIONS FROM INSTRUCTION-TUNED VIDEO MLLMS FOR SEMANTIC TASK GROUPS REVEAL DISTINCT COGNITIVE AND NEURAL PROFILES

To further examine how instruction-tuned video MLLMs generate task-specific representations and reveal functional specialization in the brain, we group the 13 video tasks into 5 cognitively grounded categories: Perceptual visual processing, Cognitive reasoning and integration, Spatiotemporal understanding, Language and narrative understanding, and Social and affective understanding. Fig. 5 illustrates that this grouping captures meaningful distinctions.

Tasks in the **Language and narrative understanding** group show broader and denser cortical engagement, particularly across the temporal and parietal cortices, compared to visual and frontal regions. In particular, we observe strong activity in the bilateral temporal lobes for narrative under-

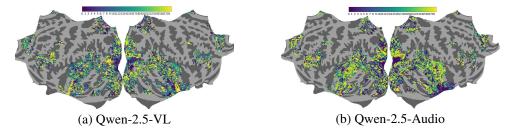


Figure 4: (a) Qwen-2.5-VL and (b) Qwen-2.5-Audio (layer-wise alignment): Each voxel is color coded with the MLLM layer number (out of 29) that led to the highest normalized brain alignment. The color bar highlights color codes for each layer. The voxels are projected onto the flattened cortical surface of average across subjects on 'fsaverage' surface.



Figure 5: Semantic Task Group Analysis: Each voxel is color coded with the task instruction that led to the highest normalized brain alignment. The color bar highlights color codes for each instruction. The voxels are projected onto the flattened cortical surface averaged across all subjects for video MLLM (Qwen-2.5-VL). While this plot shows brain maps for 3 groups, brain maps for remaining 2 task groups are in Fig. 23 in Appendix M.

standing, as well as in the angular gyrus, posterior superior temporal sulcus (pSTS), and posterior cingulate cortex (PCC) regions known to support multimodal integration, which is critical for narrative comprehension. This is aligned with previous work (Mar, 2011; Baldassano et al., 2017).

Spatiotemporal understanding. Temporal ordering elicits more widespread activation in the angular gyrus and posterior temporal lobe, whereas spatial understanding shows stronger engagement in the dorsal parietal cortex (part of the dorsal visual pathway) and anterior temporal lobe (Zacks et al., 2007; Baldassano et al., 2017). Additionally, we observe that early visual areas are more active during the spatial understanding task, whereas early auditory cortex shows higher activity in the temporal ordering task, likely due to its role in processing sound-based events (Belin et al., 2000). However, the brain does not strictly separate spatial and temporal processing. These representations often co-exist, particularly in narrative and event-based cognition.

Cognitive Reasoning. Commonsense reasoning elicits widespread activation in the temporal cortex, angular gyrus, and higher-order visual regions, reflecting its reliance on semantic processing and world knowledge. In contrast, video reasoning shows strong alignment with early visual areas (V1, V2, V3), indicating a greater dependence on visual perception and motion processing. Linking events tasks activate the early auditory cortex and show more distributed engagement of anterior temporal lobe (involved in word-level semantics), inferior frontal gyrus, and angular gyrus, highlighting the integration of temporal, linguistic, and episodic information necessary for narrative comprehension. These results show that different forms of higher-order reasoning highlights the brain's flexible organization for supporting diverse reasoning demands across modalities and timescales.

Similarly, we observe task-specific differences in brain regions for perceptual visual processing, and affective social processing (Appendix M). These patterns underscore the ability of IT-MLLMs to modulate their representations based on distinct cognitive demands reflected in the brain.

4.4 PARTITIONING EXPLAINED SHARED AND UNIQUE VARIANCE BETWEEN TASK-SPECIFIC INSTRUCTIONS

While the previous analysis reveals that task-specific instructions from MLLMs modulate their representations based on distinct cognitive demands, we further examine the representations of task-specific instructions to measure the overlap in brain variance explained by MLLMs. To accomplish this we use variance partitioning approach discussed in Appendix N.

Fig. 24 presents Venn diagrams for the whole brain, language and visual regions, depicting shared and unique variance across these regions between narrative understanding and other task instruc-

433

434

435

436

437

438

439 440 441

442 443

444

445

446

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479 480

481

482

483

484

485

tions. Similarly, we show analysis for all pairs from the 13 tasks in Table 13 in Appendix N. Across nearly all task pairs, the whole brain region consistently exhibits the highest shared variance. Tasks that are conceptually or functionally related exhibit high shared variance in all regions, indicating similar cognitive processing demands. Higher-level semantic and reasoning tasks (e.g., Narrative Understanding, Commonsense Reasoning, Temporal Ordering) show increased unique variance in the language network, indicating language-specific processing distinct from visual features. High visual load tasks (e.g., Action Recognition, Object and Scene Recognition, Global Appearance) contribute more uniquely in visual cortex, especially when paired with non-visual tasks.

5 DISCUSSION AND CONCLUSION

Using instruction-tuned representations from both video and audio MLLMs for various task-specific instructions, we evaluated how well these representations predict fMRI brain activity when participants viewed naturalistic movies (video included with audio). Additionally, we compared different video and audio MLLMs' representations, assessing their alignment with each instruction across whole brain, language, visual and auditory regions. We show that instruction-tuned video MLLMs exhibit significantly better brain alignment than audio MLLMs, vision-only, audio-only, and non-instruction-tuned multimodal models.

Our study on IT-MLLMs and their alignment with multimodal stimuli yields several key findings: (1) Although instruction-tuned video MLLMs demonstrate strong brain alignment across the whole brain (including language, visual, and auditory regions) audio MLLMs show effective alignment primarily in auditory and language-related areas such as the middle frontal gyrus (MFG). This highlights the potential of instruction-tuned audio MLLMs to capture different features relevant to auditory processing, providing information on the function of the auditory cortex similar to those observed in previous studies (Oota et al., 2024a; 2025b). However, their performance remains comparable to non-instruction-tuned multimodal models, indicating that further improvements are needed for instruction-tuned audio MLLMs to fully capture brain-relevant representations – an effort that aligns with recent work on inducing brain-relevant biases in model design (Moussa et al., 2025; Vattikonda et al., 2025). (2) The surprising effectiveness of task-specific instructions in predicting multimodal brain activity across different regions points out that both video and audio MLLMs generate distinct task-specific representations. These representations enable the models to effectively differentiate functional processing across brain regions, unlike prior work by Oota et al. (2025a), which did not observe such differentiation when using unimodal stimuli (e.g., static images). Specifically, certain audio instructions, such as audio captioning and audio understanding, show stronger alignment with language-related regions, while tasks such as sound event detection better align with the auditory cortex and temporal lobe. These findings imply that IT-MLLMs offer a powerful framework for designing controlled stimuli by a systematic manipulation of task goals through instructions, allowing researchers to isolate and examine task-specific brain responses using the same input. (3) By grouping task-specific instructions into functional categories, we find that narrative understanding consistently engages the bilateral temporal lobes, angular gyrus, and posterior cingulate cortex which are regions known for multimodal integration. Temporal ordering tasks elicit stronger responses in the angular gyrus and posterior temporal lobe, while spatial understanding activates the dorsal parietal cortex. These findings highlight the potential of instruction-tuned video MLLMs as powerful tools for probing functional specialization in the brain, offering a structured and interpretable framework for mapping high-level cognitive processes to specific neural substrates. (4) The observed correspondence between IT-MLLM layers and the brain's functional hierarchy suggests that these models inherently develop structured, brain-like representations, ranging from early sensory information processing in shallow layers to abstract semantic processing in deeper layers. This layered alignment not only enhances their interpretability but also highlights their potential as tools for investigating how the brain encodes and organizes complex, task-driven information.

Our findings also clearly show that despite the growing popularity of instruction-tuned video and audio MLLMs in handling generic task instructions, we are still far from fully interpreting how language-based instructions guide information flow through model layers and how fine-grained details are processed across layers to achieve brain-like representations. Future work should focus on leveraging the alignment strengths of these models using more fine-grained instruction-driven prompts, similar to controlled stimulus paradigms in neuroscience, to deepen our understanding of functional specialization in the brain. Lastly, we discuss limitations of our work in Appendix O.

REPRODUCIBILITY STATEMENT

Both the naturalistic stimuli (movies) and the fMRI recordings used in this study are publicly available, with preprocessing steps and experimental settings described in Section 2.1 and further detailed in Appendix C. Task-specific instruction representations from instruction-tuned video and audio MLLMs, as well as in-context learning video MLLMs, are described in Section 2.3. Implementation details of voxelwise brain encoding models and evaluation metrics are provided in Section 3, with hyperparameters listed in Appendix F. To facilitate reproducibility, we release anonymized source code for all models (instruction-tuned, in-context learning, multimodal, unimodal), brain encoding, and evaluation in the supplementary zip file.

REFERENCES

- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, August 2024.
- Khai Loong Aw and Mariya Toneva. Training language models to summarize narratives improves brain alignment. In *The Eleventh International Conference on Learning Representations*, 2023.
- Alan Baade, Puyuan Peng, and David Harwath. Mae-ast: Masked autoencoding audio spectrogram transformer. *Interspeech* 2022, 2022.
- Cordell M Baker, Joshua D Burks, Robert G Briggs, Andrew K Conner, Chad A Glenn, Kathleen N Taylor, Goksel Sali, Tressie M McCoy, James D Battiste, Daniel L O'Donoghue, et al. A connectomic atlas of the human cerebrum—chapter 7: the lateral parietal lobe. *Operative Neurosurgery*, 15(suppl_1):S295–S349, 2018.
- Christopher Baldassano, Janice Chen, Asieh Zadbood, Jonathan W Pillow, Uri Hasson, and Kenneth A Norman. Discovering event structure in continuous narrative perception and memory. *Neuron*, 95(3):709–721, 2017.
- Pascal Belin, Robert J Zatorre, Philippe Lafaille, Pierre Ahad, and Bruce Pike. Voice-selective areas in human auditory cortex. *Nature*, 403(6767):309–312, 2000.
- Vinamra Benara, Chandan Singh, John X Morris, Richard Antonello, Ion Stoica, Alexander G Huth, and Jianfeng Gao. Crafting interpretable embeddings by asking llms questions. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B (Methodological)*, 57(1):289–300, 1995.
- Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *Proceedings of the International Conference on Machine Learning (ICML)*, July 2021.
- Julie A. Boyle, Basile Pinsard, Amal Boukhdhir, Sylvie Belleville, Simona Brambatti, Jeni Chen, Julien Cohen-Adad, André Cyr, Adrian Fuente, Pierre Rainville, and Pierre Bellec. The courtois project on neuronal modelling first data release. In 26th OHBM annual meeting. Organization for Human Brain Mapping (OHBM), 2020. URL https://publications.polymtl.ca/50613/.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 24185–24198, 2024.
- Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen2-audio technical report. *arXiv* preprint arXiv:2407.10759, 2024.

- William Jay Conover. *Practical nonparametric statistics*, volume 350. john wiley & sons, 1999.
 - Colin Conwell, Jacob S Prince, Kendrick N Kay, George A Alvarez, and Talia Konkle. What can 1.8 billion regressions tell us about the pressures shaping high-level visual representation in brains and machines? *bioRxiv*, pp. 2022–03, 2022.
 - Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in Neural Information Processing Systems*, 2023.
 - Wendy A de Heer, Alexander G Huth, Thomas L Griffiths, Jack L Gallant, and Frédéric E Theunissen. The hierarchical cortical organization of human speech processing. *Journal of Neuroscience*, 37(27):6539–6557, 2017.
 - Fatma Deniz, Anwar O Nunez-Elizalde, Alexander G Huth, and Jack L Gallant. The representation of semantic information across human cerebral cortex during listening versus reading is invariant to stimulus modality. *Journal of Neuroscience*, 2019.
 - Rutvik H Desai, Usha Tadimeti, and Nicholas Riccardi. Proper and common names in the semantic system. *Brain Structure and Function*, 228(1):239–254, 2023.
 - Adrien Doerig, Tim C Kietzmann, Emily Allen, Yihan Wu, Thomas Naselaris, Kendrick Kay, and Ian Charest. Semantic scene descriptions as an objective of human vision. *arXiv* preprint *arXiv*:2209.11737, 2022.
 - Dota Tianai Dong and Mariya Toneva. Interpreting multimodal video transformers using brain recordings. In *ICLR 2023 Workshop on Multimodal Representation Learning: Perks and Pitfalls*, 2023a.
 - Dota Tianai Dong and Mariya Toneva. Vision-language integration in multimodal video transformers (partially) aligns with the brain. *arXiv preprint arXiv:2311.07766*, 2023b.
 - Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
 - Christopher R Genovese. A bayesian time-course model for functional magnetic resonance imaging data. *Journal of the American Statistical Association*, 95(451):691–703, 2000.
 - Matthew F Glasser, Timothy S Coalson, Emma C Robinson, Carl D Hacker, John Harwell, Essa Yacoub, Kamil Ugurbil, Jesper Andersson, Christian F Beckmann, Mark Jenkinson, et al. A multi-modal parcellation of human cerebral cortex. *Nature*, 536(7615):171–178, 2016.
 - Ariel Goldstein, Zaid Zada, Eliav Buchnik, Mariano Schain, Amy Price, Bobbi Aubrey, Samuel A Nastase, Amir Feder, Dotan Emanuel, Alon Cohen, et al. Shared computational principles for language processing in humans and deep language models. *Nature Neuroscience*, 25(3):369–380, 2022.
 - Alexander G Huth, Shinji Nishimoto, An T Vu, and T Dupre La Tour. Gallant lab natural short clips 3t fmri data. *G-Node doi*, 10, 2022.
 - Shailee Jain and Alexander Huth. Incorporating context into language encoding models for fmri. *Advances in Neural Information Processing Systems*, 31, 2018.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
 - Team Kimi Team. Kimi-audio technical report, 2024.
 - Tom Dupré la Tour, Michael Eickenberg, Anwar O Nunez-Elizalde, and Jack L Gallant. Feature-space selection with banded ridge regression. *NeuroImage*, 264:119728, 2022.

- Amanda LeBel, Shailee Jain, and Alexander G Huth. Voxelwise encoding models show that cerebellar language representations are highly conceptual. *Journal of Neuroscience*, 41(50):10341–10355, 2021.
 - Xinhao Li, Ziang Yan, Desen Meng, Lu Dong, Xiangyu Zeng, Yinan He, Yali Wang, Yu Qiao, Yi Wang, and Limin Wang. Videochat-r1: Enhancing spatio-temporal perception via reinforcement fine-tuning. *arXiv preprint arXiv:2504.06958*, 2025.
 - Bin Lin, Yang Ye, Bin Zhu, Jiaxi Cui, Munan Ning, Peng Jin, and Li Yuan. Video-llava: Learning united visual representation by alignment before projection. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 5971–5984, 2024.
 - Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad Shoeybi, and Song Han. Vila: On pre-training for visual language models, 2023.
 - Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in Neural Information Processing Systems*, 36, 2024.
 - Khai Loong Aw, Syrielle Montariol, Badr AlKhamissi, Martin Schrimpf, and Antoine Bosselut. Instruction-tuning aligns llms to the human brain. First Conference on Language Modeling, 2024.
 - Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL 2024)*, 2024.
 - Raymond A Mar. The neural bases of social cognition and story comprehension. *Annual Review of Psychology*, 62(1):103–134, 2011.
 - Camille K Milton, Vukshitha Dhanaraj, Isabella M Young, Hugh M Taylor, Peter J Nicholas, Robert G Briggs, Michael Y Bai, Rannulu D Fonseka, Jorge Hormovas, Yueh-Hsin Lin, et al. Parcellation-based anatomic model of the semantic network. *Brain and Behavior*, 11(4):e02065, 2021.
 - Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. Fast model editing at scale. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/pdf?id=0DcZxeWfOPt.
 - Omer Moussa, Dietrich Klakow, and Mariya Toneva. Improving semantic understanding in speech language models via brain-tuning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=KL8Sm4xRn7.
- Yuko Nakagi, Takuya Matsuyama, Naoko Koide-Majima, Hiroto Yamaguchi, Rieko Kubo, Shinji Nishimoto, and Yu Takagi. Unveiling multi-level and multi-modal semantic representations in the human brain using large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 20313–20338, 2024.
- Satya Sai Srinath Namburi, Makesh Sreedhar, Srinath Srinivasan, and Frederic Sala. The cost of compression: Investigating the impact of compression on parametric knowledge in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Singapore, December 2023. Association for Computational Linguistics. URL https://aclanthology.org/2023.findings-emnlp.349/.
- Subba Reddy Oota, Jashn Arora, Veeral Agarwal, Mounika Marreddy, Manish Gupta, and Bapi Surampudi. Neural language taskonomy: Which nlp tasks are the most predictive of fmri brain activity? In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 3220–3237, 2022a.
- Subba Reddy Oota, Jashn Arora, Vijay Rowtula, Manish Gupta, and Raju S Bapi. Visio-linguistic brain encoding. In *COLING*, pp. 116–133, 2022b.
- Subba Reddy Oota, Agarwal Veeral, Marreddy Mounika, Gupta Manish, and Raju Surampudi Bapi. Speech taskonomy: Which speech tasks are the most predictive of fmri brain activity? In *24th INTERSPEECH Conference*, 2023.

- Subba Reddy Oota, Emin Çelik, Fatma Deniz, and Mariya Toneva. Speech language models lack important brain-relevant semantics. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8503–8528. Association for Computational Linguistics, 2024a.
 - Subba Reddy Oota, Manish Gupta, and Mariya Toneva. Joint processing of linguistic properties in brains and language models. *Advances in Neural Information Processing Systems*, 36, 2024b.
 - Subba Reddy Oota, Akshett Rai Jindal, Ishani Mondal, Khushbu Pahwa, Satya Sai Srinath Namburi GNVV, Manish Shrivastava, Maneesh Kumar Singh, Bapi Raju Surampudi, and Manish Gupta. Correlating instruction-tuning (in multimodal models) with vision-language processing (in the brain). In *The Thirteenth International Conference on Learning Representations*, 2025a.
 - Subba Reddy Oota, Khushbu Pahwa, mounika marreddy, Maneesh Kumar Singh, Manish Gupta, and Bapi Raju Surampudi. Multi-modal brain encoding models for multi-modal stimuli. In *The Thirteenth International Conference on Learning Representations*, 2025b.
 - Sara F Popham, Alexander G Huth, Natalia Y Bilenko, Fatma Deniz, James S Gao, Anwar O Nunez-Elizalde, and Jack L Gallant. Visual and linguistic semantic representations are aligned at the border of human visual cortex. *Nature Neuroscience*, 24(11):1628–1636, 2021.
 - Aniketh Janardhan Reddy and Leila Wehbe. Can fmri reveal the representation of syntactic structure in the brain? *Advances in Neural Information Processing Systems*, 34:9843–9856, 2021.
 - Christina Sartzetaki, Gemma Roig, Cees GM Snoek, and Iris IA Groen. One hundred neural networks and brains watching videos: Lessons from alignment. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/pdf?id=LM4PYXBId5.
 - Martin Schrimpf, Idan Asher Blank, Greta Tuckute, Carina Kauf, Eghbal A Hosseini, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. The neural architecture of language: Integrative modeling converges on predictive processing. *Proceedings of the National Academy of Sciences*, 2021.
 - Padakanti Srijith, Khushbu Pahwa, Radhika Mamidi, Bapi Raju Surampudi, Manish Gupta, and SUBBA REDDY OOTA. Aligning text/speech representations from multimodal models with MEG brain activity during listening. In *The 2025 Conference on Empirical Methods in Natural Language Processing*, 2025. URL https://openreview.net/forum?id=111z0F1EZk.
 - V Subramaniam, C Wang, A Barbu, G Kreiman, and B Katz. Revealing vision-language integration in the brain with multimodal networks. In *International Conference on Machine Learning*. International Conference on Machine Learning (ICML), 2024.
 - Jingyuan Sun and Marie-Francine Moens. Fine-tuned vs. prompt-tuned supervised representations: which better account for brain language representations? In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, pp. 5197–5205, 2023.
 - Jingyuan Sun, Xiaohan Zhang, and Marie-Francine Moens. Tuning in to neural encoding: Linking human brain and artificial supervised representations of language. In *ECAI 2023*, pp. 2258–2265. IOS Press, 2023.
 - Jerry Tang, Meng Du, Vy Vo, Vasudev Lal, and Alexander Huth. Brain encoding models based on multimodal transformers can transfer across language and vision. *Advances in Neural Information Processing Systems*, 36, 2024.
 - Zineng Tang, Jaemin Cho, Yixin Nie, and Mohit Bansal. Tvlt: Textless vision-language transformer. *Advances in Neural Information Processing Systems*, 35:9617–9632, 2022.
 - Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. Stanford alpaca: An instruction-following llama model, 2023.

- Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *Advances in Neural Information Processing Systems*, 35:10078–10093, 2022.
 - Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
 - Greta Tuckute, Jenelle Feather, Dana Boebinger, and Josh H McDermott. Many but not all deep neural network audio models capture brain responses and exhibit correspondence between model stages and brain regions. *Plos Biology*, 21(12):e3002366, 2023.
 - Aditya R Vaidya, Shailee Jain, and Alexander Huth. Self-supervised models of audio effectively explain human cortical responses to speech. In *International Conference on Machine Learning*, pp. 21927–21944. PMLR, 2022.
 - Nishitha Vattikonda, Aditya R Vaidya, Richard J Antonello, and Alexander G Huth. Brainwavlm: Fine-tuning speech representations with brain responses to language. *arXiv* preprint arXiv:2502.08866, 2025.
 - Aria Wang, Michael Tarr, and Leila Wehbe. Neural taskonomy: Inferring the similarity of task-derived representations from brain activity. *Advances in Neural Information Processing Systems*, 32:15501–15511, 2019.
 - Aria Y Wang, Kendrick Kay, Thomas Naselaris, Michael J Tarr, and Leila Wehbe. Better models of human high-level visual cortex emerge from natural language supervision with a large and diverse dataset. *Nature Machine Intelligence*, 5(12):1415–1426, 2023.
 - Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
 - Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, et al. Qwen2. 5-omni technical report. *arXiv preprint arXiv:2503.20215*, 2025.
 - Zhiyang Xu, Ying Shen, and Lifu Huang. Multiinstruct: Improving multi-modal zero-shot learning via instruction tuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 11445–11465, 2023.
 - Jeffrey M Zacks, Nicole K Speer, Khena M Swallow, Todd S Braver, and Jeremy R Reynolds. Event perception: a mind-brain perspective. *Psychological Bulletin*, 133(2):273, 2007.
 - Yuanhan Zhang, Bo Li, haotian Liu, Yong jae Lee, Liangke Gui, Di Fu, Jiashi Feng, Ziwei Liu, and Chunyuan Li. Llava-next: A strong zero-shot video understanding model, April 2024. URL https://llava-vl.github.io/blog/2024-04-30-llava-next-video/.

Overview of Appendix Sections

- · Appendix A: Overview of multimodal model evaluation settings in brain encoding studies
- Appendix B: Related work
- Appendix C: Detailed sub-ROIs of language, visual and auditory regions
- Appendix D: Cross-subject prediction accuracy
- Appendix E: Model generated outputs across instructions
- Appendix F: Implementation details for reproducibility
- Appendix G: Statistical Significance
- Appendix H: Effectiveness of instruction-tuned video MLLMs vs audio MLLMs vs multimodal vs unimodal representations for various brain regions
- Appendix I: Contrasting Instruction-tuned video MLLMs with In-context learning video MLLMs

Appendix J: Contrasting Instruction-tuned video MLLMs with non-instruction-tuned mul-

Appendix K: Brain Maps for Task-specific instructions

Appendix L: Brain Maps showing Layer-wise Details for Video Instruction-based MLLMs

Appendix M: Details of Semantic Task Group Analysis

Appendix N: Details of explained variance partitioning

• Appendix O: Limitations

• Appendix P: LLM Usage

OVERVIEW OF MULTIMODAL MODEL EVALUATION SETTINGS IN BRAIN **ENCODING STUDIES**

Table 4: Overview of multimodal model evaluation settings in brain encoding studies.

Study	Model Type	Stimulus Modality	Brain	Dataset	Instruction-Tuned
			Data		
Doerig et al. (2022)	Vision-Language (CLIP)	Unimodal (Images)	fMRI	NSD	X
Wang et al. (2023)	Vision-Language (CLIP)	Unimodal (Images)	fMRI	NSD	Х
Oota et al. (2022b)	Vision-Language (CLIP, Vi-	Unimodal (Images)	fMRI	BOLD5000	Х
	sualBERT, LXMERT)				
Popham et al. (2021)	Vision-Only CNNs vs.	Unimodal (Silent Videos)	fMRI	Gallant lab short	Х
	Vision-Language			video clips	
Tang et al. (2022)	non-instruction-tuned multi-	Unimodal (Silent Videos),	fMRI	Gallant lab short	Х
	modal model (BridgeTower)	Unimodal (listening stories)		video clips	
Oota et al. (2025a)		Unimodal (Images)	fMRI	NSD	/
	age+Text MLLMs	, , ,			
Sartzetaki et al. (2025)	Image Recognition models,	Unimodal (Visual)	fMRI	Bold Moments	Х
	Action recognition models			Dataset	
Nakagi et al. (2024)	Language models (BERT,	Multimodal (Videos with	fMRI	8.3 hours of video	Х
	GPT-2, Lllama2, OPT)	audio)		dataset	
Subramaniam et al.	non-instruction-tuned multi-	Image frame-text pairs	SEEG	AMMT	Х
(2024)	modal models (SLIP-CLIP,	(Movies)			
	SimCLR, BLIP, BEIT)				
Dong & Toneva (2023a)	non-instruction-tuned mul-	Multimodal (Movies:	fMRI	Neuromod Friends	Х
	timodal models (Merlore-	Videos with audio)		dataset	
	serve)	·			
Oota et al. (2025b)	non-instruction-tuned multi-	Multimodal (Movies:	fMRI	Neuromod Movie10	Х
	modal models (TVLT and	Videos with audio)			
	ImageBind)				
Our study	instruction-tuned video and	Multimodal (Movies:	fMRI	Neuromod Movie10	✓
	audio MLLMs, in-context	Videos with audio)			
	learning video and audio				
	MLLMs				

RELATED WORK

Brain encoding using multimodal models. Our work is closely related to that of Conwell et al. (2022); Wang et al. (2023); Doerig et al. (2022); Tang et al. (2024); Nakagi et al. (2024); Dong & Toneva (2023b); Oota et al. (2025b), who proposed using multimodal model representations to study the contribution of brain alignment in unimodal and multimodal stimuli. The majority of brain encoding studies in using multimodal models focused on a single modality of input - vision alone (Conwell et al., 2022; Wang et al., 2023; Doerig et al., 2022; Wang et al., 2023; Tang et al., 2024; Nakagi et al., 2024). Recently, Dong & Toneva (2023b); Oota et al. (2022b) interpreted the effectiveness of multimodal Transformer language models in multimodal naturalistic stimuli. However, these studies focus on pretrained multimodal models which are not generic to tasks and lack the investigation of recent instruction-tuned models.

Task-based brain alignment. Our work is also closely related to that of Wang et al. (2019); Oota et al. (2022a); Aw & Toneva (2023); Sun et al. (2023) and Loong Aw et al. (2024), who propose using task-specific model representations to study the contribution of individual tasks to brain alignment. Wang et al. (2019) investigated 21 computer vision tasks to explore which vision tasks are more aligned with the brain while subjects engaged in viewing passive images. Similarly, Oota et al. (2022a) and Sun et al. (2023) explored 10 GLUE NLP tasks to study which NLP tasks are

more brain-aligned during reading and listening to stories. More recent work by Loong Aw et al. (2024) uses instruction-tuned LLMs to investigate the effect of natural language instruction model representations on brain alignment across layers for language comprehension. Further, Oota et al. (2025a) use IT-MLLMs (image+text), using natural language instructions across diverse vision tasks to analyze their alignment with brain activity across layers during visual processing. However, these studies primarily focused on unimodal stimuli and thus do not fully capture the capabilities of multimodal instruction-tuned models under multimodal conditions. We complement these works by examining the impact of a wide range of IT-MLLMs—spanning video and audio-based models with text-based prompts—on their alignment with brain activity from multimodal stimuli.

C DETAILED SUB-ROIS OF LANGUAGE, VISUAL AND AUDITORY REGIONS

The data covers seven brain regions of interest (ROIs) in the human brain with the following subdivisions: (i) early visual (EV: V1, V2, V3, V3B, and V4); (ii) object-related areas (LO1 and LO2); (iii) face-related areas (OFA), (iv) scene-related areas (PPA), (v) middle temporal (MT: MT, MST, LO3, FST and V3CD), (vi) late language regions, encompassing broader language regions: angular gyrus (AG: PFm, PGs, PGi, TPOJ2, TPOJ3), lateral temporal cortex (LTC: STSda, STSva, STGa, TE1a, TE2a, TGv, TGd, A5, STSdp, STSvp, PSL, STV, TPOJ1), inferior frontal gyrus (IFG: 44, 45, IFJa, IFSp) and middle frontal gyrus (MFG: 55b) (Baker et al., 2018; Milton et al., 2021; Desai et al., 2023).

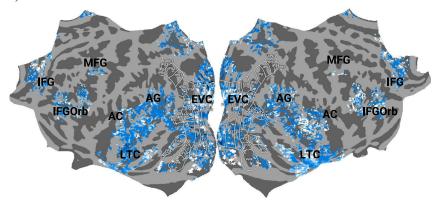


Figure 6: Flattened cortical surfaces for language-, visual- and auditory-selective regions displayed on the 'fsaverage' surface, used as the mask for all participants.

D CROSS-SUBJECT PREDICTION ACCURACY

We follow the method introduced by Schrimpf et al. (2021) to estimate how well brain activity in one individual can be predicted from others, using the Movie10 fMRI dataset. Starting with data from n participants (e.g., n=4), for each subject $s\in([1,4])$ is chosen as the prediction target and the other three are used to predict this target, we use a voxel-wise encoding model (see Sec. 3) to predict one participant's response from others. For every combination, one participant was randomly chosen as the target, and the model was trained to predict their brain responses using data from the remaining s-1 participants. This gave us an average prediction score (correlation) for each voxel at each participant. To extrapolate to infinitely many humans and thus to obtain the highest possible (most conservative) estimate, as suggested by Schrimpf et al. (2021), we fit the equation $v=v_0\times\left(1-e^{-\frac{x}{\tau_0}}\right)$ where x is each subsample's number of participants, v is each subsample's correlation score and v_0 and v0 are the fitted parameters. This fitting was performed for each sensor independently with 100 bootstraps each to estimate the variance where each bootstrap draws v0 and v1 with replacement. The final ceiling value was the median of the per-voxel ceilings v2.

Fig. 7 shows the estimated cross-subject prediction accuracy for all four participants for the naturalistic movie watching. Pearson correlation scores for each voxel in each subject are projected onto the subject's flattened cortical surface. The plots show that across all subjects higher activity is

observed in the language and visual regions with a max correlation up to 0.4 implying that data has low noise and low cross-subject variability.

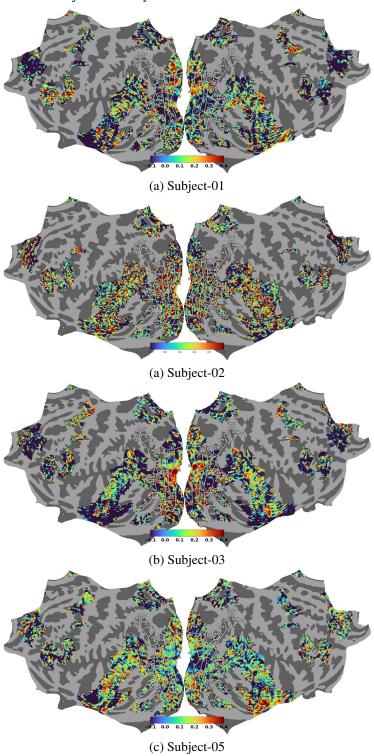


Figure 7: Estimated cross-subject prediction accuracy for all four participants for the naturalistic movie watching. Pearson correlation scores for each voxel in each subject are projected onto the subject's flattened cortical surface.

E MODEL GENERATED OUTPUTS ACROSS INSTRUCTIONS

Tables 5, 6, 7, 8, 9 and 10 show model generated outputs for a sample video from the Movie10 dataset using InstructBLIPVideo, Qwen-2.5-VL, Video-LLaVA, LLaVa-NeXT-Video, LLaVA-OneVision and VideoChat-R1 models, respectively. Similarly, Tables 11 and 12 show model generated outputs for a sample video from the Movie10 dataset using Qwen-2.5 Audio and Kimi-Audio models.

Table 5: Outputs from InstructBLIPVideo for a sample clip from the Movie10 dataset.

Frame 1	Frame 3	Frame 8	Frame 13	Frame 16
Instruction	Model Output			
Action Recognition	The main event capture versation with another other people, possibly wrapped around his ne	man. The man is wearing g in a bar or restaurant. The eck.	ting at a bar with a glass of liglasses and is sitting in a dinir scene also features a man wea	ng area, surrounded by aring a towel, which is
Video Understanding	man. The two men are bottles of beer visible people can also be see and interacting with the	e sitting at a table, with one in the background, placed en in the scene, possibly pa	n his hand, engrossed in a con of them holding a glass of wi on a table and around the b trons or other bar visitors, sc e scene is set in a dimly lit en en's conversation.	ne. There are multiple par area. A few other attered around the bar
Visual Question Answering	sitting at a table, and o	lrinking.	ne woman. They are having a	conversation at a bar,
Video Captioning		tting at a table talking with a		
Object and Scene Recognition			ng in a bar, talking to another a glass of beer, while the other	
Commonsense Reasoning	discussing something.	This action could have mot seeking to socialize, share	ngaging in conversation with ivated the man to take this sto his thoughts or ideas, or sin	ep for various reasons.
Spatial Understanding	Restaurant			
Temporal Ordering	The scene depicted in man. The man is wear or holding a drink. Th interest or sharing a st	ing glasses and is holding a ey appear to be having a ca- ory.	ng at a bar, engaging in a con glass in his hand, which sugg sual conversation, possibly di	ests that he is drinking scussing something of
Video reasoning	man in a bar, and they social setting for two environment, which ac	y are both holding glasses. men to be having a conversidds to the unusual nature of		ause it is not a typical ng in a dark, dimly lit
Narrative Understanding	glasses. The man in the of interest. It is not a movie's storyline.	e bar is having a conversation specified what the topic of	ng in a bar talking to another on with another man, possibly their conversation is, but it i	discussing something s likely related to the
Emotion and Sentiment Analysis	a social and relaxed a which creates a more	tmosphere. The man is weat intimate and cozy atmosphe	ring glasses, and the scene is	
Global Appearance	Man with glasses and			
Linking Events	possibly a friend, whil		ing glasses sitting at a bar, tall is early event sets the stage for eo.	

F IMPLEMENTATION DETAILS FOR REPRODUCIBILITY

All feature extraction experiments were conducted on a machine equipped with an NVIDIA A100 GPU with 80 GB of GPU RAM, partitioned into two devices of 40 GB each. The voxelwise encoding models were trained on NVIDIA GeForce RTX 3050 GPU with 4GB of GPU RAM. We used banded ridge-regression with the following parameters: MSE loss function; L2-decay (λ) varied from 10^{-1} to 10^3 ; the best λ was chosen by tuning on validation data that comprised a randomly chosen 10% subset from the train set used only for hyper-parameter tuning.

Table 6: Outputs from Qwen-2.5-VL for a sample video from the Movie10 dataset.

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

Frame 1	Frame 3	Frame 8	Frame 13	Frame 16	
Instruction	Model Output				
Action Recognition	The video appears to be a still from a movie or a scene from a film. It shows two individuals engaged in a conversation in what looks like a dimly lit bar or restaurant setting. The person on the left is wearing a suit and tie, holding a drink, and seems to be listening intently. The person on the right is gesturing with their hand, possibly emphasizing a point during the conversation. The atmosphere suggests a serious or intense discussion.				
Video Understanding	sation. The character on cigarette in his right han visible from behind. The	The video depicts a scene from a movie or TV show, featuring two characters engaged in a conversation. The character on the left is dressed in a formal gray suit with a patterned tie and is holding a cigarette in his right hand. He appears to be listening intently to the other character, who is partially visible from behind. The setting seems to be an indoor environment, possibly a bar or a similar establishment, as suggested by the background elements such as a mirror and what looks like a bar counter.			
Visual Question Answering	and appears to be engage from behind. The second	ed in a conversation or in I person is facing away frindoors, possibly in a b	sitting at a table, wearing a st teraction with another person rom the camera, and their attir ar or restaurant, given the bar	who is partially visible re includes a plaid shirt.	
Video Captioning	A man in a suit is engage	ed in a conversation with	another person in a dimly lit	bar setting.	
Object and Scene Recognition	a brown suit with a whit glass of what appears to	e shirt and a red paisley be a drink in his left har son. 2. The Other Perso	everal objects: 1. The Man in tie. He has dark hair styled ad. His expression seems serie on: Only the back of this indi	neatly and is holding a ous or contemplative as	
Commonsense Reasoning	is turned to the camera. or pointing at something Explanation: The chara	The character's hand is eg. This action could be cter might be trying to e	be in a conversation with anot xtended towards the other per motivated by several reason xplain something important o ld be used to emphasize a p	rson, possibly gesturing s: 1. Clarification or r clarify a point in their	
Spatial Understanding	released in 2013. The set of the characters' activiti	tting looks like a bar or a les in the financial world initively identify the exa	olf of Wall Street", directed by lounge, which is consistent vol. However, there are no specient location within the movie.	with the film's depiction fic landmarks visible in	
Temporal Ordering	The image depicts a sce is a step-by-step descrip environment, possibly a	ne from a movie where tion of the activity show bar or a restaurant, given 2. Characters: - The c	two characters are engaged i vn: 1. Setting: The setting a n the presence of a glass and character on the left is wearing	appears to be an indoor what looks like a menu	
Video reasoning	a conversation. The setti relaxed atmosphere. Wh gesturing with their hand	ing looks like a bar or a at might be considered u while holding a drink, v e conversation. However	ie or TV show, featuring two similar social environment, v musual is the character on the which could indicate they are e y, without additional context, it	with dim lighting and a right, who seems to be explaining something or	
Narrative Understanding	The movie depicted in th Leonardo DiCaprio as Jo himself. Main Storyline	e video is "The Wolf of ordan Belfort. The film is: The story follows Jord	Wall Street," directed by Martis based on the memoir of the an Belfort, a young stockbrok the financial world during the	e same name by Belfort ter who rises to become	
Emotion and Sentiment Analysis	The video depicts two individuals in what appears to be a dimly lit, possibly upscale setting, such as a bar or restaurant. The person on the left is dressed in a suit and tie, suggesting a formal or professional demeanor, while the person on the right is wearing a plaid shirt, which might indicate a more casual or relaxed attitude.				
Global Appearance	outfit changes. The man shirt. He has dark hair st	on the left appears to be yled neatly. The person of	significant changes in the cha wearing a brown suit with a p on the right is seen from behir changes in their attire or appea	patterned tie and a white and, wearing a plaid shirt	
Linking Events	DiCaprio's character, Jorappears to be a dimly lit l	dan Belfort, is engaged par or lounge, which is ty s. An early event that co	movie "The Wolf of Wall S in a conversation with another pical for scenes involving dis- uld influence later developmed d his current interlocutor.	er character. The setting cussions about business	

STATISTICAL SIGNIFICANCE

To determine if normalized predictivity scores are significantly higher than chance, we run a permutation test using blocks of 10 contiguous fMRI TRs (considering the slowness of hemodynamic response) rather than individual TRs. By permuting predictions 5000 times, we create an empirical

Table 7: Outputs from Video-LLaVA for a sample clip from the Movie 10 dataset

Frame 1	Frame 3	Frame 8	Frame 13	Frame 16
Instruction	Model Output			
Action Recognition	The man holding the b	eer is wearing a plaid shirt		, and the second
Video Understanding	drink and is engaged i a plaid shirt, giving hi and glasses on display, patrons.	in conversation with someo m a casual and relaxed den , indicating a lively and bust	drink in his hand. He appe ne else at the bar. The man neanor. The bar is well-lit, v ling atmosphere. In the back	is wearing glasses an with a variety of bottle ground, there are other
Visual Question Answering	and holding a beer glas	SS.	are sitting down, while the th	nird man is standing u
Video Captioning		laid shirt drinking beer at a		
Object and Scene Recognition			tting at a bar with a beverage	
	The bar is well-stocked blue plaid shirt, and a	d with various bottles and c beige blanket.	nen sitting at the bar, and a w ups. The man sitting at the b	ar is wearing glasses,
Commonsense Reasoning	motivation for the cha socialization, or simpl	racter to take this action co y enjoying a beverage duri	and appears to be in a bar of ould be due to various reasong leisure time. The charact with others in a social setting	ons, such as relaxation er may have chosen t
Spatial Understanding	holding a beer glass. T		e presence of multiple bottle an standing in front of a bar,	
Temporal Ordering	glasses. The scene is seen drinking from the	set in a dark room, and the beer glass, which is placed	a beer in his hand. He is we re are other people in the ba on the bar counter. The can g his relaxed and social atmo	ackground. The man in there captures the man'
Video reasoning	a bar and drinking a b	peer. The man is also wear ver, the most striking feature	a man in a plaid shirt and ging a beige scarf, which add e of this video is that the man	s a distinctive touch t
Narrative Understanding	to a bartender's refusa	al. The man's frustration ar at the bar. The conflict is ult	for a drink at a bar, but he ind disappointment lead to hi imately resolved when the m	m making a scene an
Emotion and Sentiment Analysis	drinks and conversatio	n. The man in the plaid shirts holding a beer glass. The	oman are seen laughing and a t is seen drinking beer and lo t video conveys a sense of re	ooking down, while th
Global Appearance	also seen wearing a dif The other man in the	ferent outfit in one of the sco	ard, is seen wearing a plaid a enes, which is not specified in and glasses. There are also her wearing a dark shirt.	the given information
Linking Events	The early event in the beer in his hand. This	video is the man in the blu scene introduces the setting . The man's interaction wit	e and white checkered shirts and the man's presence at the hother patrons, the woman,	he bar, setting the stag

distribution for chance performance, from which we estimate p-value of the actual performance. The choice of these specific permutation test configurations is based on established methodologies in previous research (Deniz et al., 2019; Reddy & Wehbe, 2021; Oota et al., 2024a). To estimate the statistical significance of performance differences, such as between the model's predictions and chance or residual predictions and chance, we utilized the Wilcoxon signed-rank test (Conover, 1999), applying it to the mean normalized predictivity for the participants. Finally, the Benjamini-Hochberg False Discovery Rate (FDR) correction for multiple comparisons (Benjamini & Hochberg, 1995) is used for all the tests (appropriate because fMRI data is considered to have positive dependence (Genovese, 2000)).

H EFFECTIVENESS OF INSTRUCTION-TUNED VIDEO MLLMS VS AUDIO MLLMS VS MULTIMODAL VS UNIMODAL REPRESENTATIONS FOR VARIOUS BRAIN REGIONS

Fig. 8 show average normalized brain alignment of instruction-tuned video MLLMs vs instruction-tuned audio MLLMs vs multimodal and unimodal models across several ROIs (AG, ATL, PTL, IFG,

Table 8: Outputs from LLaVA-NeXT-Video for a sample clip from the Movie10 dataset.

	Enw 2	Frame 9	France 12	Fig. 16
Frame 1	Frame 3	Frame 8	Frame 13	Frame 16
Instruction Action Recognition	is wearing a red and w tie. They are discussing	thite checkered shirt and gla ng a topic, possibly related	ation between two men at a lasses, while the other man is to a movie, as suggested by der is visible in the foreground	wearing a dark suit and the context. There are
Video Understanding	Bradley Cooper. In th	ne scene, Zach Galifianakis, oup of people. He is seen ho	e Hangoverfeaturing actors portraying a character nam lding a beer and appears to b	ed Chip Baskin, is sit-
Visual Question Answering	appears to be in a conv sitting nearby. There is visible, but it's unclear	ersation with someone. And a a bartender standing at the what their role is in the sce		ar, and there is a woman ground who is partially
Video Captioning	Two men are seen hav or disbelieving the ma		with the man in the plaid sh	nirt appearing confused
Object and Scene Recognition	a man in a suit standir man standing at a bar,	ng at a bar, a woman sitting and a man sitting at a table		nan standing at a bar, a
Commonsense Reasoning	be interpreted as a play	yful or humorous gesture, p is action might be motivate	ding up a beer bottle and mak ossibly to amuse others or to ed by a desire to entertain, b	express a lighthearted
Spatial Understanding	Robert De Niro. The fella's." This is a refe	scene is set in a bar, and the erence to the character Rob	ene in a movie, specifically be background includes a sign ert De Niro's character, what tent with the movie's depicti	n that says "The Good- o is a bartender in the
Temporal Ordering	the other is dressed in	a plaid shirt. They are botting suggests they might be	I in a conversation. One man th holding drinks and appear e at a social event, such as a	r to be having a casual
Video reasoning	appears to be engaged other man, who is not in the middle of a disc	in a conversation with son the focus of the scene. The ussion, but his face is not sh		mera is focused on the a beer and appears to be
Narrative Understanding	Burry, who predicts the	e financial crisis of 20008. T gage market, which Burry so	dy-drama film about the life the central conflict of the stores as unsustainable and warr	ry is the housing bubble
Emotion and Sentiment Analysis	even a bit of confusion be in a state of disbeli- drinking from his beer	 The video itself is quite here or shock as the other mare bottle. 	notions, from surprise and sh amorous, with the man in a s n, who is wearing a plaid shi	uit and tie appearing to rt and a hoodie, is seen
Global Appearance	in a more casual outfi	t with a plaid shirt. There a	s wearing a suit at the beginn are also other people in the s d a red shirt, and a woman ir	scene, including a man
Linking Events	a TV. The TV shows a himself. This early eve	news segment about a man ent in the video influences th	ting in a bar with a drink in f who has died, which is later e later developments because tablishing his presence in the	revealed to be the man e it sets up the narrative

MFG, IFGOrb, PCC and dmPFC) of language region. Fig. 9 show the same for visual, auditory and motor regions.

CONTRASTING INSTRUCTION-TUNED VIDEO MLLMs WITH IN-CONTEXT LEARNING VIDEO MLLMs

We present contrast of brainmaps to display the average normalized brain alignment across voxels for the instruction-tuned video MLLMs versus the in-context learning video MLLMs in Figures 10, and 11. The results show that instruction-tuned video MLLMs consistently achieve significantly higher alignment across all brain voxels.

Table 9: Outputs from LLaVA-OneVision Video for a sample clip from the Movie10 dataset.

Frame 1	Frame 3	Frame 8	Frame 13	Frame 16
Instruction	Model Output			
Action Recognition	The video captures a speaking while the ot		ndividuals at a bar, with one	e person gesturing and
Video Understanding	features a well-lit bar in the background, ad table. One person is d	counter adorned with various ding to the bustling ambiance	ng a warm and inviting atmosp is bottles and glasses, and seve e. In the foreground, two indi- ied and white patterned tie, whi is ded over their shoulders.	eral people can be seen ividuals are seated at a
Visual Question Answering	There are two people	in the scene, engaged in a co	onversation.	
Video Captioning	A man with curly has conversation.	r and glasses sits at a table	in a bar, holding a glass of b	beer, and engages in a
Object and Scene Recognition	bar or restaurant. One beige sweater over it,	man is dressed in a gray suit	are two men sitting at a table and tie, and the other is wear ound features a bar with vario	ring a plaid shirt with a
Commonsense Reasoning			point in the conversation, when need to convey a specific r	
Spatial Understanding		m a movie scene, and it show		
Temporal Ordering	The video takes place shelves. Two individe is dressed in a gray of the person on the rig shoulders and holding	e in a dimly lit bar with a li luals are seated at a table, e suit with a red and white pa ht, wearing a blue and red p g a glass of beer, listens atten	ively background filled with engaged in a conversation. I atterned tie and is gesturing volaid shirt with a yellow swe tively.	The person on the left with their hand, while eater draped over their
Video reasoning		that the scene remains static varieties a continuous conv	with minimal changes in the in versation or interaction.	ndividuals' expressions
Narrative Understanding	sweater being advised		par, with the character in the character. The conflict is resolution and posture.	
Emotion and Sentiment Analysis	creates a sense of intr	igue and emotional depth, dr	ing surprise, concern, and con awing the viewer into the cha	aracters' conversation.
Global Appearance			nd is wearing a blue plaid shir they are dressed in a dark suit	
Linking Events			ding a glass of beer and enga for the unfolding dialogue an	

J CONTRASTING INSTRUCTION-TUNED VIDEO MLLMS WITH NON-INSTRUCTION-TUNED MULTIMODAL

We present contrast of brainmaps to display the average normalized brain alignment across voxels for the instruction-tuned video MLLMs versus the non-instruction-tuned multimodal models VILA and TVLT in Figures 12, 13, 14, 15, and 16. The results show that instruction-tuned video MLLMs consistently achieve significantly higher alignment across all brain voxels. However, Figures 17 and 18 reveal clear differences between audio MLLMs and multimodal models: the prediction performance of audio MLLMs lacks brain-relevant semantic information compared to multimodal models.

K Brain Maps for Task-specific instructions

Figures 19 and 20 show brain maps for InstructBLIPVideo, Video-LLaVA, LLaVA-NeXT-Video, LLaVA-OneVision and VideoChat-R1 for video tasks for average normalized brain predictivity across subjects where the voxel color codes are projected onto the flattened cortical surface of the 'fsaverage' subject. The color-scheme corresponding to each instruction is also reported. We make the following observations: (i) Video understanding exhibits the strongest alignment across the whole brain. (ii) Tasks such as spatial understanding, narrative understanding, and visual question answering show higher alignment in language-related regions, including the angular gyrus, posterior temporal lobe, and visual regions. (iii) Higher-order language regions in the frontal cortex are predominantly identified by the video understanding task, with a smaller proportion of voxels also activated by video reasoning and temporal ordering tasks.

Table 10: Outputs from VideoChat-R1 for a sample clip from the Movie10 dataset.

1	1	8	8	
1	1	8	9	
1	1	9	0	
1	1	9	1	
1	1	9	2	
1	1	9	3	
1	1	9	4	
1	1	9	5	
1	1	9	6	
1	1	9	7	
1	1	9	8	
1	1	9	9	
	2			
	2			
	2			
	2			
	2			
	2			
1	2	0	6	
	2			
	2			
	2			
	2			
	2			
	2			
1	2	1	3	
	2			
	2			
	2			
	2			
	2			
	2			
	2			
	2			
	2			
	2			
	2			
	2			
	2			
	2			
1	2	2	8	
1	2	2	9	
	2			
1	2	3	1	

Frame 1	Frame 3	Frame 8	Frame 13	Frame 16	
nstruction	Model Output				
Action Recognition	The video captures a s individuals are engage over their shoulders, su in a suit and tie, indicate	d in conversation. One per- aggesting they might be col-	be a social setting, possibly son is wearing a plaid shirt a d or seeking comfort. The ot The background shows other z.	and has a blanket draped her individual is dressed	
Video Understanding	The video depicts a see Two individuals are see suit and tie, while the i	ene from what appears to be ated at a table, engaged in c individual on the right is dr	a movie or TV show, set in a onversation. The person on t essed in a plaid shirt with a y suggesting they might be er	the left is wearing a dark yellow scarf draped over	
Visual Question Answering	One person is wearing		be engaged in a conversation lanket draped over their sho casual social interaction.		
Video Captioning			surrounded by dim lighting a	nd patrons.	
Object and Scene Recognition	wearing a dark suit and shirt with a yellowish features a bar setting v	d appears to be listening at blanket draped over their	engaged in conversation. Tentively. The person on the shoulders and is holding a sses, indicating a social or che evening.	right is wearing a plaid drink. The background	
Commonsense Reasoning	looks like a bar or a s slightly, which might their conversation. The trying to show interest	imilar social setting. The indicate that they are liste e motivation for this action or engagement in the conv		s to be leaning forward nphasize a point during The character might be	
Spatial Understanding	counter with bottles an a social or dining estal	d glasses, suggesting it cou blishment. However, witho e exact location or landman	v set in a bar or restaurant. The ld be a scene from a film or sut additional context or recork shown in the video. It could be a scene from a film or sut additional context or recork shown in the video.	series that takes place in gnizable landmarks, it's	
Temporal Ordering	wearing a suit and app in a plaid shirt with a comfort. They both ha	bears to be listening attention blanket draped over their save drinks in front of them.	, engaged in conversation. The individual on the individual of the	right is dressed casually night be cold or seeking njoying a social setting.	
Video reasoning	setting. The unusual as a cartoonish character	spect here is the presence of , which seems out of place	ie or TV show, likely set in f a large, animated hand with in a real-world setting like edic or dramatic purposes w	h a glove that resembles a bar. This suggests it	
Narrative Understanding	seemingly engaged in	a conversation. The centra	e "The Hangover," where the al conflict in this movie revo out end up getting lost and ha	olves around a group of	
Emotion and Sentiment Analysis	friends who go on a bachelor party to Las Vegas but end up getting lost and having various mishaps. The video depicts two individuals engaged in what appears to be an intense conversation at a bar. The person on the right seems to be speaking with a sense of urgency or concern, as indicated by their open mouth and expressive eyes. The other individual is listening attentively, suggesting a moment of serious discussion. Overall, the scene conveys a sense of tension and engagement. The setting, with its dim lighting and crowded bar, adds to the atmosphere of a private yet public interaction.				
Global Appearance	tie, while the characte shoulders. There are n	r on the right is wearing a	The character on the left is plaid shirt with a yellow bleir appearances or outfits thr of the frames provided.	lanket draped over their	
Linking Events	In the video, two indivand possibly intimate significant. This early by establishing a found	iduals are seated at a bar, er interaction, which could in event of them sitting toge dation for their relationship narrative, the initial interact	ngaged in conversation. The supply that they are discussing ther and talking could influe or the context of their convetion might set up themes suc	g something personal or ence later developments ersation. For instance, if	

Fig. 21 shows brainmap for audio instruction-tuned MLLM (Kimi-Audio) where the predictions are average across subjects. Here, the voxel color codes are projected onto the flattened cortical surface of the 'fsaverage' subject. The figure shows a clear distinction between different audio tasks.

Table 11: Outputs from Qwen-2.5 Audio for a sample audio from the Movie10 dataset (Wolf of wallstreet).

Spectrogram

Instruction
Audio Understanding
Audio Comprehension
Audio Captioning
Sound Event Detection
Speaker Identification
A man speaking in English with an angry mood, saying 'And she grew up hot.
They are talking and laughing.
The audio is about someone's upbringing.
Sound Event Detection
Speaker Identification
A man speaking in English with an angry mood said 'And she grew up hot, you know?

Table 12: Outputs from Kimi-Audio for a sample audio from the Movie10 dataset (Wolf of wall-street)

Spectrogram									
Instruction	Model Output								
Audio Understanding	The audio contains a speech where kids are talking by the door in English with a female voice ar								
	neutral mood.								
Audio Comprehension	The people are talking and there is music playing in the background.								
Audio Captioning	A man is speaking loudly and aggressive.								
Sound Event Detection	There is no speech in this audio.								
Speaker Identification	The speaker in the audio is a Mandarin-speaking male, aged between 16-25 years old, with an angry								
	mood.								

L Brain Maps showing Layer-wise Details for Video Instruction-based MLLMs

To examine whether IT-MLLMs reflect the brain's hierarchy of information processing across layers, we analyze the voxels as follows. For each voxel, we select the layer that results in the highest normalized brain alignment and apply a color code for the 29/33 layers across the various MLLMs. Fig. 22 presents brain maps for four video MLLMs, where the voxels with their corresponding color codes are projected onto the flattened cortical surface of the 'fsaverage' subject.

M DETAILS OF SEMANTIC TASK GROUP ANALYSIS

To further examine how instruction-tuned video MLLMs generate task-specific representations and reveal functional specialization in the brain, we group the 13 video tasks into five cognitively grounded categories: Perceptual visual processing, Cognitive reasoning and integration, Spatiotemporal understanding, High-level language and narrative understanding, and Social and affective understanding. This categorization allows us to disentangle the functional specificity of brain regions engaged by different task types. The visualizations in Fig. 5 in Section 4.3 in the main paper and Fig. 23 illustrate that this grouping captures meaningful distinctions.

N DETAILS OF EXPLAINED VARIANCE PARTITIONING

Variance partitioning. To disentangle task-specific instruction representations from multimodal instruction-tuned models, we used a variance partitioning approach (de Heer et al., 2017; LeBel et al., 2021). This method measures the overlap in brain variance explained by different task-specific instruction representations. Specifically, variance partitioning separates the brain response variance

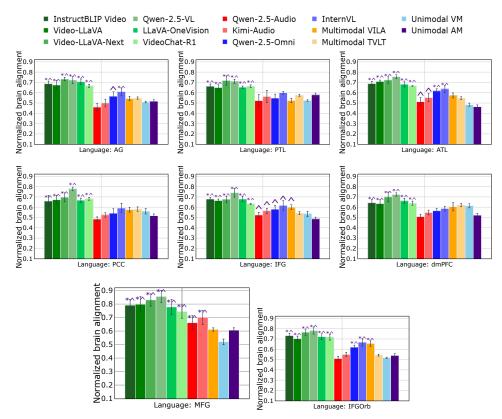


Figure 8: Average normalized brain alignment of instruction-tuned video MLLMs vs instruction-tuned audio MLLMs vs multimodal and unimodal models across several ROIs (AG, ATL, PTL, IFG, MFG, IFGOrb, PCC and dmPFC) of language region. Error bars indicate the standard error of the mean across participants. * implies that instruction-tuned MLLM embeddings are significantly better than multimodal models and \land means that instruction-tuned MLLM embeddings are significantly better unimodal models with p ≤ 0.05 .

that can be attributed to two models based on their unique and overlapping contributions (Vaidya et al., 2022; Deniz et al., 2019). To perform this, for every pair of instruction representations, we fit separate encoding models for each space as well as a joint encoding model, obtained by concatenating the features. Using set arithmetic, we can then derive the size of the intersection $(NBA)_v^{1\cap 2} = (NBA)_v^1 + (NBA)_v^2 - (NBA)_v^{1\cup 2}$, where NBA refers to normalized brain alignment, v refers to a specific voxel, $(NBA)_v^1$ denotes alignment of model 1, $(NBA)_v^2$ denotes alignment of model 2 and $(NBA)_v^{1\cup 2}$ denotes alignment of the joint model. Similarly, the unique contribution of model 1's feature space is computed as $(NBA)_v^{1\setminus 2} = (NBA)_v^1 - (NBA)_v^{1\cap 2}$.

Shared and Unique Variance between Narrative Understanding and Remaining Task Instructions

Fig. 24 shows the shared variance of the Narrative Understanding task with other video tasks for Qwen-2.5-VL.

Table 13 presents shared and unique variance explained by pairs of video tasks using brain-informed models across three neural regions: whole brain, visual cortex, and language network. The results are averaged across subjects and show how well representations from each task pair align with brain activity in specific regions.

Key Observations are as follows.

• Whole Brain Shows Dominant Shared Variance: Across nearly all task pairs, the whole brain region consistently exhibits the highest shared variance (often ¿80% in early task pairs). For example, the pair Action Recognition and Video Understanding (1–2) shows 90.69% shared variance, with very little unique variance from either task. This suggests

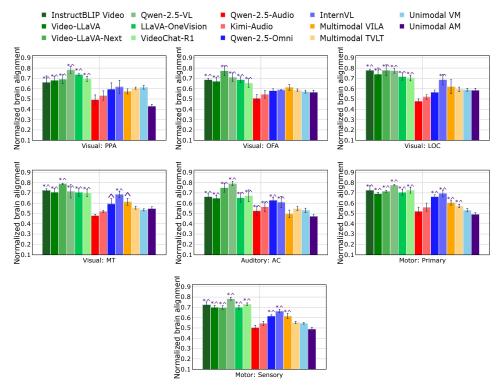


Figure 9: Average normalized brain alignment of instruction-tuned video MLLMs vs instruction-tuned audio MLLMs vs multimodal and unimodal models across several ROIs of visual cortex (PPA, OFA, LOC, MT), Auditory cortex (AC), and Motor Area (PMA and SMA). Error bars indicate the standard error of the mean across participants. * implies that instruction-tuned MLLM embeddings are significantly better than multimodal models and \land means that instruction-tuned MLLM embeddings are significantly better unimodal models with p ≤ 0.05 .

high redundancy and common processing across tasks when considering global brain activity.

- Visual and Language Regions Yield More Balanced Partitioning: In contrast, visual and language-selective voxels exhibit lower shared variance and comparatively higher unique contributions from individual tasks. For the same task pair (1–2), shared variance in visual is 72.05%, and in language it is 77.46%, with higher unique components (~10-14%). This suggests that fine-grained processing differences are more pronounced in modality-specific regions.
- Task Similarity Reflects in Shared Variance: Tasks that are conceptually or functionally related (e.g., Narrative Understanding-Linking Events (10-13) or Emotion and Sentiment Analysis-Linking Events (11-13)) exhibit high shared variance in all regions, indicating similar cognitive processing demands. Conversely, task pairs with less conceptual overlap (e.g., Object Recognition-Commonsense Reasoning (5-6) or Visual QA-Object Recognition (3-5)) show lower shared variance and higher unique variance, especially in language and visual regions.
- Language Regions Show Selectivity for High-Level Tasks: Higher-level semantic and reasoning tasks (e.g., Narrative Understanding, Commonsense Reasoning, Temporal Ordering) show increased unique variance in the language network, indicating language-specific processing distinct from visual features. For instance, pair 6-13 (Commonsense Reasoning-Linking Events) yields 16.75% unique variance for Linking Events in the language network.
- Visual Cortex Captures Scene and Action Differentiation: Tasks with high visual load (e.g., Action Recognition, Object and Scene Recognition, Global Appearance) contribute more uniquely in the visual cortex, especially when paired with non-visual tasks.

O LIMITATIONS

One possible limitation of our study lies in interpreting the differences in brain alignment between instruction-tuned video and audio MLLMs. The models we evaluate differ in several aspects, including the amount of training data and the specific objective functions used during training. To address this concern, we evaluated multiple models of each type, spanning a range of training objectives and dataset sizes, and found that our key results generalize within both video and audio MLLM categories. Still, it is possible that some of the differences in brain alignment may still be influenced by confounding factors related to model architecture, training objectives, or data scale. Future work should explore these questions using models that are more tightly controlled across these dimensions.

P LLM USAGE

We used OpenAI ChatGPT for grammar correction and language polishing.

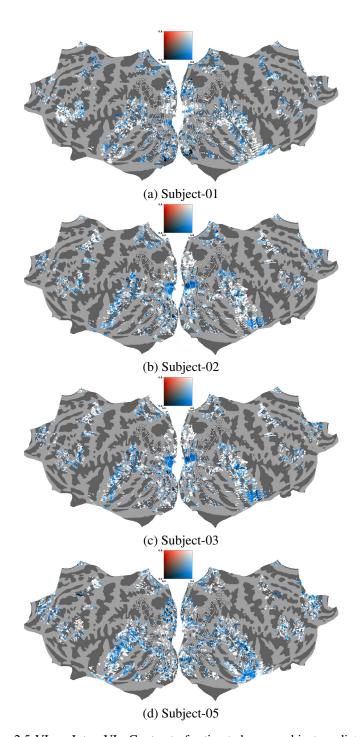


Figure 10: Qwen-2.5-VL vs.InternVL: Contrast of estimated cross-subject prediction accuracy for all participants for the naturalistic movie watching. Pearson correlation scores for each voxel in each subject are projected onto the subject's flattened cortical surface. Blue and Red voxels depict higher prediction accuracy estimates during instruction-tuned video MLLM and in-context learning video MLLM (InternVL), respectively. Voxels that have similar cross-subject prediction accuracy appear white. Here, middle frontal gyrus (MFG), inferior frontal gyrus (IFG), inferior frontal gyrus orbital (IFGOrb), angular gyrus (AG), and lateral temporal cortex (LTC) are late language regions, EVC denotes early visual cortex and AC denotes auditory cortex.

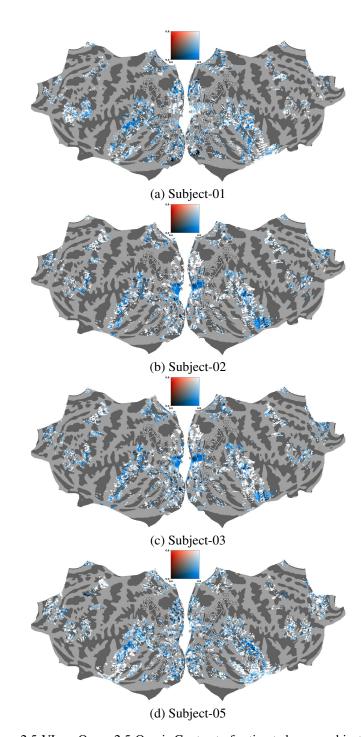


Figure 11: Qwen-2.5-VL vs.Qwen-2.5-Omni: Contrast of estimated cross-subject prediction accuracy for all participants for the naturalistic movie watching. Pearson correlation scores for each voxel in each subject are projected onto the subject's flattened cortical surface. Blue and Red voxels depict higher prediction accuracy estimates during instruction-tuned video MLLM and in-context learning video MLLM (Qwen-2.5-Omni), respectively. Voxels that have similar cross-subject prediction accuracy appear white. Here, middle frontal gyrus (MFG), inferior frontal gyrus (IFG), inferior frontal gyrus orbital (IFGOrb), angular gyrus (AG), and lateral temporal cortex (LTC) are late language regions, EVC denotes early visual cortex and AC denotes auditory cortex.

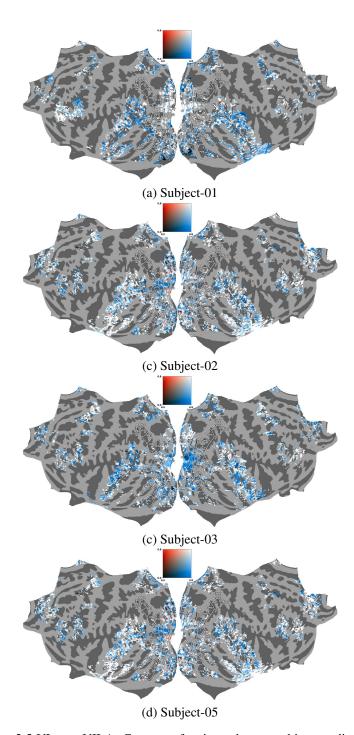


Figure 12: Qwen-2.5-VL vs. VILA: Contrast of estimated cross-subject prediction accuracy for all participants for the naturalistic movie watching. Pearson correlation scores for each voxel in each subject are projected onto the subject's flattened cortical surface. Blue and Red voxels depict higher prediction accuracy estimates during instruction-tuned video MLLM and multimodal VILA, respectively. Voxels that have similar cross-subject prediction accuracy appear white. Here, middle frontal gyrus (MFG), inferior frontal gyrus (IFG), inferior frontal gyrus orbital (IFGOrb), angular gyrus (AG), and lateral temporal cortex (LTC) are late language regions, EVC denotes early visual cortex and AC denotes auditory cortex.

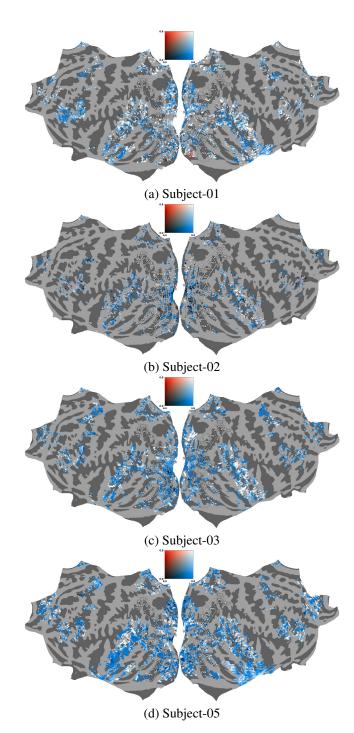


Figure 13: Qwen-2.5-VL vs. TVLT: Contrast of estimated cross-subject prediction accuracy for all participants for the naturalistic movie watching. Pearson correlation scores for each voxel in each subject are projected onto the subject's flattened cortical surface. Blue and Red voxels depict higher prediction accuracy estimates during instruction-tuned video MLLM and multimodal TVLT, respectively. Voxels that have similar cross-subject prediction accuracy appear white. Here, middle frontal gyrus (MFG), inferior frontal gyrus (IFG), inferior frontal gyrus orbital (IFGOrb), angular gyrus (AG), and lateral temporal cortex (LTC) are late language regions, EVC denotes early visual cortex and AC denotes auditory cortex.

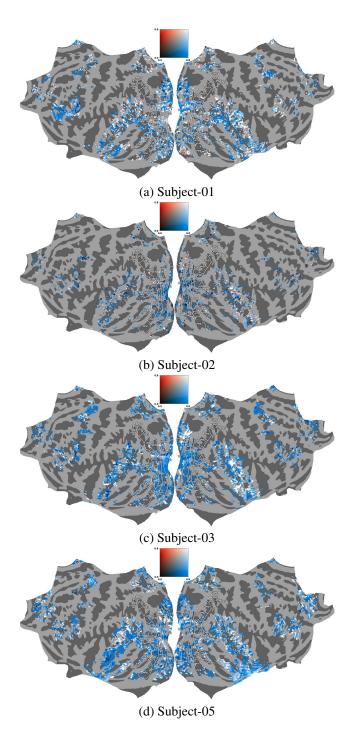


Figure 14: InstructBLIPVideo vs. TVLT: Contrast of estimated cross-subject prediction accuracy for all participants for the naturalistic movie watching. Pearson correlation scores for each voxel in each subject are projected onto the subject's flattened cortical surface. Blue and Red voxels depict higher prediction accuracy estimates during instruction-tuned video MLLM and multimodal TVLT, respectively. Voxels that have similar cross-subject prediction accuracy appear white.

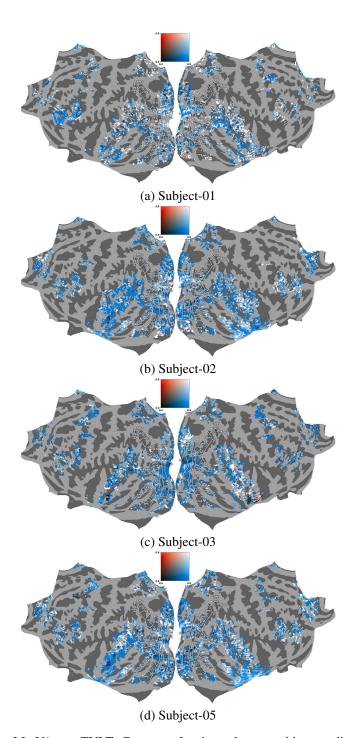


Figure 15: Video-LLaVA vs. TVLT: Contrast of estimated cross-subject prediction accuracy for all participants for the naturalistic movie watching. Pearson correlation scores for each voxel in each subject are projected onto the subject's flattened cortical surface. Blue and Red voxels depict higher prediction accuracy estimates during instruction-tuned video MLLM and multimodal TVLT, respectively. Voxels that have similar cross-subject prediction accuracy appear white.

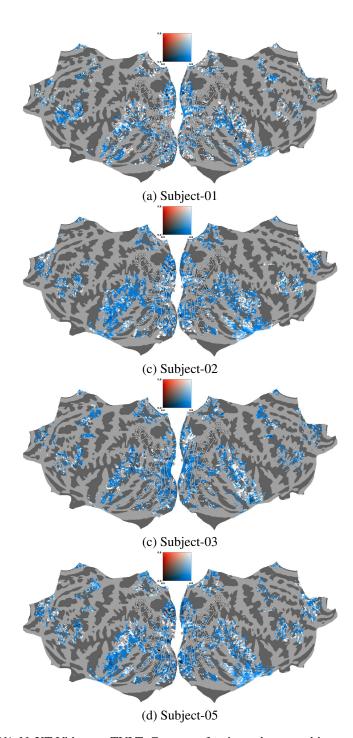


Figure 16: LLaVA-NeXT-Video vs. TVLT: Contrast of estimated cross-subject prediction accuracy for all participants for the naturalistic movie watching. Pearson correlation scores for each voxel in each subject are projected onto the subject's flattened cortical surface. Blue and Red voxels depict higher prediction accuracy estimates during instruction-tuned video MLLM and multimodal TVLT, respectively. Voxels that have similar cross-subject prediction accuracy appear white.

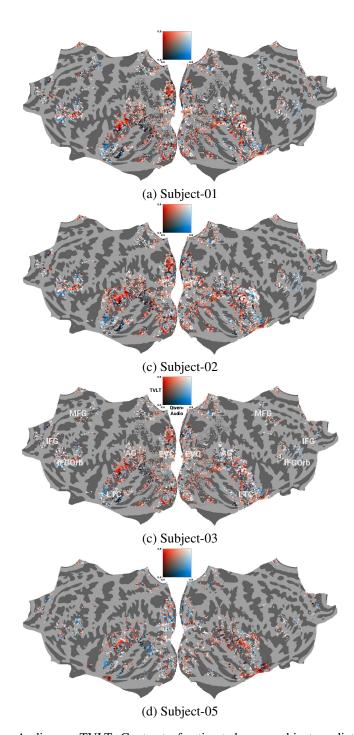


Figure 17: Qwen-Audio vs. TVLT: Contrast of estimated cross-subject prediction accuracy for all participants for the naturalistic movie watching. Pearson correlation scores for each voxel in each subject are projected onto the subject's flattened cortical surface. Blue and Red voxels depict higher prediction accuracy estimates during instruction-tuned audio MLLM and multimodal TVLT, respectively. Voxels that have similar cross-subject prediction accuracy appear white. Here, middle frontal gyrus (MFG), inferior frontal gyrus (IFG), inferior frontal gyrus orbital (IFGOrb), angular gyrus (AG), and lateral temporal cortex (LTC) are late language regions, EVC denotes early visual cortex and AC denotes auditory cortex.

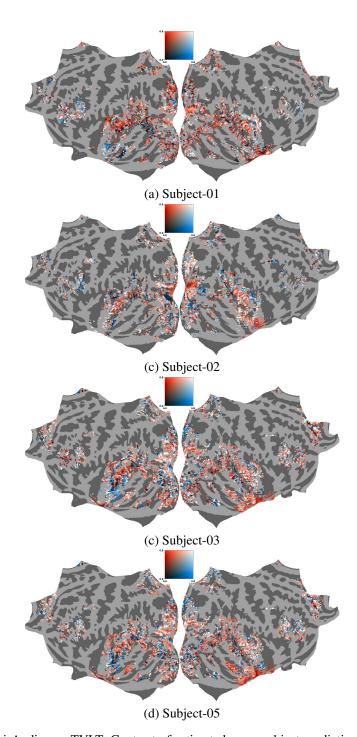


Figure 18: Kimi-Audio vs. TVLT: Contrast of estimated cross-subject prediction accuracy for all participants for the naturalistic movie watching. Pearson correlation scores for each voxel in each subject are projected onto the subject's flattened cortical surface. Blue and Red voxels depict higher prediction accuracy estimates during instruction-tuned audio MLLM and multimodal TVLT, respectively. Voxels that have similar cross-subject prediction accuracy appear white. Here, middle frontal gyrus (MFG), inferior frontal gyrus (IFG), inferior frontal gyrus orbital (IFGOrb), angular gyrus (AG), and lateral temporal cortex (LTC) are late language regions, EVC denotes early visual cortex and AC denotes auditory cortex.

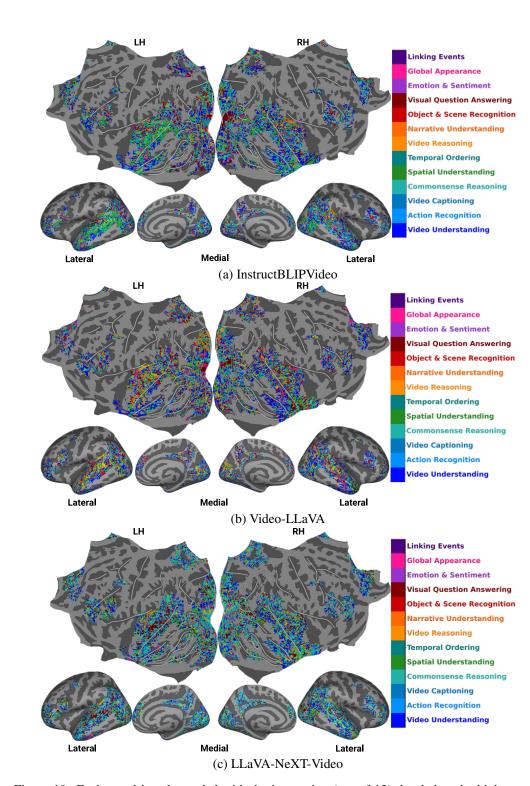


Figure 19: Each voxel is color coded with the instruction (out of 13) that led to the highest normalized brain alignment. The color bar highlights color codes for each instruction. The voxels are projected onto the flattened cortical surface averaged across all 4 subjects for 3 video MLLM (InstructBLIPVideo, Video-LLaVA and LLaVA-NeXT-Video).

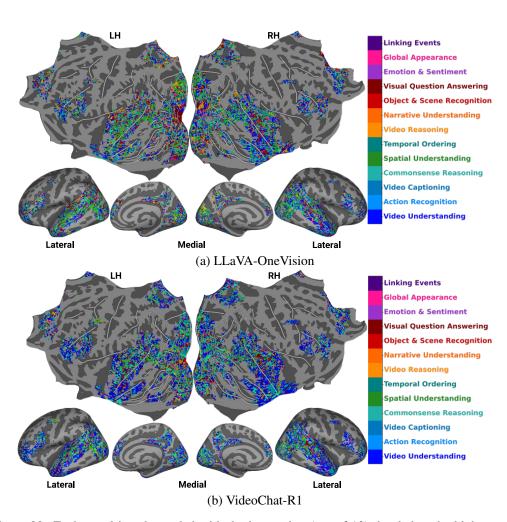


Figure 20: Each voxel is color coded with the instruction (out of 13) that led to the highest normalized brain alignment. The color bar highlights color codes for each instruction. The voxels are projected onto the flattened cortical surface averaged across all 4 subjects for 2 video MLLM (LLaVA-OneVision, VideoChat-R1).

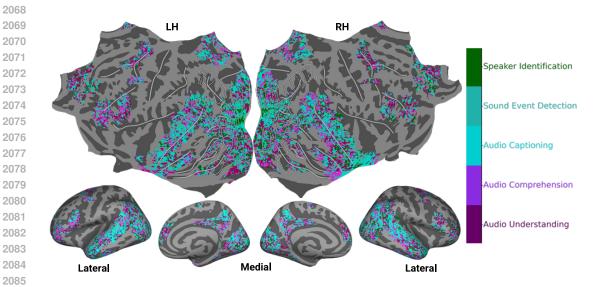


Figure 21: Kimi-Audio: Each voxel is color-coded with the instruction (out of 5) that led to the highest normalized brain alignment. The color bar highlights color codes for each instruction. The voxels are projected onto the flattened cortical surface of average across subjects on 'fsaverage' surface.

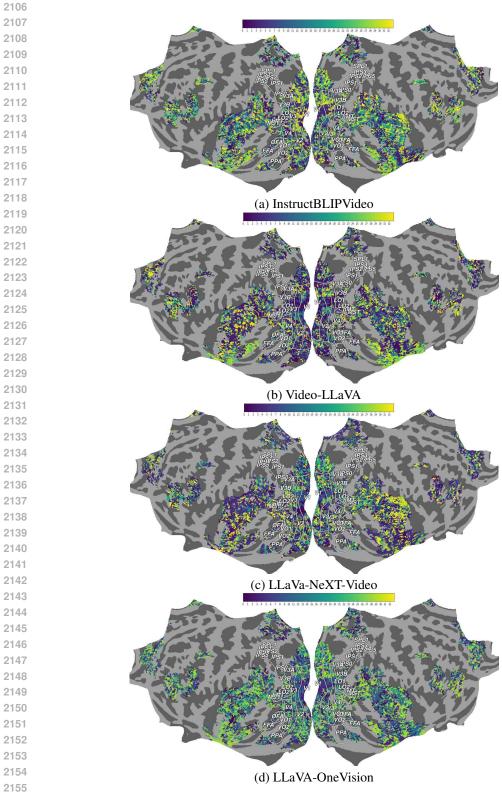


Figure 22: Each voxel is color coded with the video MLLM layer number (out of 33) that led to the highest normalized brain alignment. The color bar highlights color codes for each layer. The voxels are projected onto the flattened cortical surface of average across all 4 subjects on 'fsaverage' surface for four MLLMs.

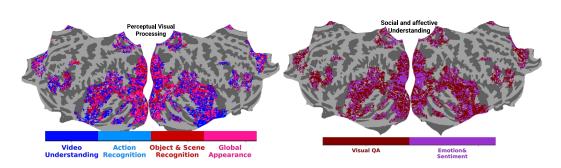


Figure 23: Semantic Task Group Analysis: Each voxel is color coded with the task instruction that led to the highest normalized brain alignment. The color bar highlights color codes for each instruction. The voxels are projected onto the flattened cortical surface averaged across all subjects for video MLLM (Qwen-2.5-VL). While this plot shows brain maps for 2 groups, brain maps for remaining 3 task groups are in Fig. 5 in Section 4.3 in the main paper.

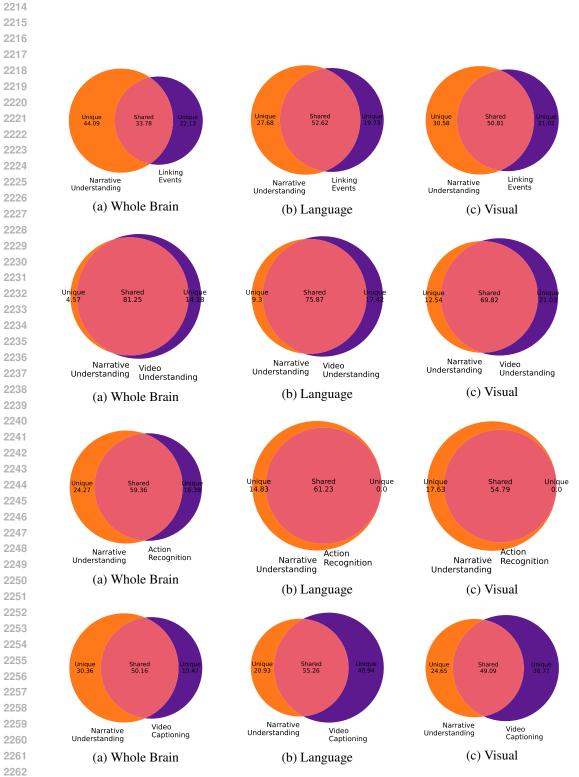


Figure 24: Shared and Unique Variance: Narrative Understanding vs. Linking Events Dark orange (left) shows variance unique to Narrative Understanding, indigo (right) shows variance unique to Linking Events, and the overlap indicates shared variance between both tasks.

T1-1	T1-2	Whole Brain			Visual Shared Uniq1 Uniq2			Language		
Task1	Task2	90.69	5.26	4.05		13.91	14.04	77.46	12.07	10.4
1	3	83.53	10.05	6.42	72.05 73.67	10.28	16.05	77.05	10.72	
1	4	84.51	9.65	5.84	71.87	13.82	14.31	75.97	12.27	11.
1	5	79.16	13.51	7.33	66.82	14.35	18.83	73.47	13.07	13.4
1	6	81.48	13.34	5.18	68.44	17.28	14.28	73.59	15.37	11.0
1	7	83.07	10.44	6.49	71.99	11.88	16.13	75.20	12.30	12.:
1	8	81.25	14.18	4.57	69.82	17.63	12.54	75.87	14.83	9.3
1	9	86.94	7.57	5.50	73.42	10.25	16.34	78.27	9.05	12.0
1	10	84.55	9.06	6.39	73.46	10.59	15.95	76.42	10.32	13.
1	11	85.44	8.51	6.05	74.92	11.12	13.96	76.56	10.96	12.4
1	12	82.46	11.66	5.88	72.88	12.75	14.37	76.02	12.50	11.
1	13	91.81	4.20	3.99	74.92	11.82	13.26	80.06	10.00	9.9
2	3	83.59	9.72	6.69	73.14	11.39	15.47	74.15	12.80	13.
2	4	86.25	7.40	6.36	73.32	13.52	13.16	74.41	12.14	13.
2	5	77.09	14.33	8.58	64.55	17.14	18.31	70.20	15.08	14.
2	6	79.86	13.99	6.15	69.43	17.86	12.71	73.10	14.96	11.
2	7	83.62	9.46	6.92	72.53	12.65	14.82	71.61	14.43	13.
2	8	81.30	13.10	5.60	67.98	18.96	13.05	72.05	16.07	11.
2	9	86.64	7.42	5.93	73.55	12.35	14.11	75.55	10.62	13.
2	10	85.25	7.97	6.78	72.98	12.28	14.73	73.28	12.51	14.
2 2	11	84.70 82.97	8.31 11.16	7.00	73.27 73.06	12.25	14.48	72.48 72.99	13.27	14.
2	13	91.78	3.66	5.88 4.55	74.89	12.59	12.54	78.19	9.77	12.
3	4	68.68	13.67	17.64	68.53	18.38	13.09	71.98	14.19	13.
3	5	50.07	24.61	25.32	52.60	24.08	23.32	60.68	17.79	21.
3	6	61.39	21.67	16.94	61.59	22.97	15.44	65.21	18.68	16.
3	7	65.21	17.99	16.80	64.73	20.33	14.94	66.85	17.80	15.
3	8	66.30	20.20	13.49	61.04	23.96	15.00	62.43	21.86	15.
3	9	70.23	13.71	16.06	70.07	16.68	13.25	72.20	12.52	15.
3	10	66.99	13.00	20.01	68.60	15.97	15.42	64.43	15.79	19.
3	11	68.07	14.39	17.54	66.84	17.50	15.66	66.97	16.85	16.
3	12	61.81	19.24	18.95	65.81	19.69	14.50	67.09	17.92	14.
3	13	83.92	6.44	9.64	71.83	16.87	11.31	76.76	12.86	
4	5	55.03	24.36	20.61	53.05	20.94	26.00	59.06	18.82	22.
4	6	61.72	25.66	12.62	59.66	24.72	15.62	63.75	21.99	14.
4	7	69.00	17.62	13.38	66.08	17.45	16.47	67.89	17.50	14.
4	8	63.88	21.85	14.27	60.24	23.59	16.17	65.25	19.95	14.
4	9	71.16	16.55	12.28	65.51	18.15	16.34	68.66	16.14	15.
4	10	66.37	18.11	15.53	63.85	17.11	19.04	57.73	20.94	21.
4	11	72.37	13.56	14.07	70.00	13.01	16.99	70.64	13.35	16.
4	12	66.38 86.69	18.76 6.09	7.23	64.80 71.23	17.67 16.28	17.53 12.49	67.94 76.56	17.21 13.87	9.:
5	6	50.09	27.24	22.63	51.63	27.81	20.56	58.56	23.05	18.
5	7	49.08	24.63	26.29	53.55	25.15	21.30	55.77	24.66	19.
5	8	47.03	27.55	25.43	53.22	28.86	17.93	53.88	26.92	19.
5	9	55.06	21.61	23.34	56.84	24.75	18.42	62.62	19.24	18
5	10	47.76	23.54	28.70	55.84	22.99	21.17	54.52	22.48	23.
5	11	52.17		25.25	57.44	22.32	20.24	57.94	22.48	
5	12	47.50	26.51	25.99	56.38	25.48	18.15	58.21	23.50	18.
5	13	79.36	6.98	13.67	66.31	16.96	16.74	71.80	12.91	15.
6	7	60.01	17.04	22.96	59.05	17.09	23.86	61.14	18.01	20.
6	8	54.31	21.48	24.22	57.44	21.55	21.01	62.62	18.13	
6	9	64.33	13.06	22.61	60.10	16.20	23.69	64.68	13.72	21.
6	10	57.84	16.91	25.25	61.41	14.59	24.00	61.01	16.15	
6	11	62.94	14.26	22.81	62.17	15.15	22.68	63.32	15.40	21.
6	12	55.82		24.54		17.37			18.93	
6	13	81.42	5.21	13.37	67.46	13.51	19.02	71.93	11.31	16.
7	8	58.19	23.15	18.65	60.58	23.47	15.95	61.00 71.25	20.86	18.
7	10	70.87 68.57	14.02 12.51	15.11 18.92	70.43 67.67	13.03	14.51	63.76	14.39	16. 21.
7	11	60.77	18.94	20.29	58.79	21.23	19.00	55.14	21.77	23.
7	12	66.57	17.86	15.57	67.97	17.05	14.98	67.18	17.38	15.
7	13	85.27	6.01	8.72	72.66	15.56	11.78	74.88	13.08	12.
8	9	62.84	15.99	21.18	63.11	15.66	21.22	68.03	13.67	18
8	10	60.10	17.38	22.52	59.39	16.80	23.81	60.46	16.80	22
8	11	60.31	14.63	25.07	61.67	13.24	25.09	61.38	15.64	22
8	12	60.04	18.69	21.28	62.31	17.41	20.28	65.74	16.70	17
8	13	81.06	5.66	13.27	68.01	14.38	17.61	74.50	11.65	13
9	10	69.21	14.34	16.44	68.83	12.98	18.19	67.69	15.88	16
9	11	70.80	13.15	16.05	69.96	14.08	15.96	70.82	14.04	15.
9	12	69.68	16.60	13.72	70.09	14.45	15.46	70.62	16.10	13
9	13	87.40	5.23	7.37	72.02	15.46	12.53	77.48	12.70	9.
10	11	68.63	16.35	15.02	67.96	16.43	15.61	64.85	19.12	16.
10	12	65.06	20.66	14.27	63.79	21.85	14.36	61.84	23.65	14.
		85.63	6.39	7.99	72.34	16.92	10.73	75.85	14.09	10.
10	13									
	13 12 13	61.95	22.51	15.54 7.58	65.60	19.55	14.85	63.80	21.51	14. 10.

Table 13: Variance partitioning for all the 13 video tasks averaged across all subjects for whole brain, visual and language regions with Qwen-2.5-VL model. Tasks are as follows: (1) Action Recognition (2) Video Understanding (3) Visual Question Answering (4) Video Captioning (5) Object and Scene Recognition (6) Commonsense Reasoning (7) Spatial Understanding (8) Temporal Ordering (9) Video reasoning (10) Narrative Understanding (11) Emotion and Sentiment Analysis (12) Global Appearance (13) Linking Events.