

---

# Neuroprobe: Evaluating Intracranial Brain Responses to Naturalistic Stimuli

---

Andrii Zahorodnii<sup>1,2\*</sup>

Bennett Stankovits<sup>1\*</sup>

Christopher Wang<sup>1\*</sup>

Charikleia Moraitaki<sup>1</sup>

Geeling Chau<sup>3</sup>

Ila R Fiete<sup>1,2</sup>

Boris Katz<sup>1</sup>

Andrei Barbu<sup>1</sup>

<sup>1</sup>MIT CSAIL, CBMM

<sup>2</sup>MIT McGovern Institute

<sup>3</sup>Caltech

## Abstract

Understanding the relationship between the various tasks the brain performs can shed light on its functional organization. We introduce a benchmark, Neuroprobe, which targets a wide range of multimodal tasks. Neuroprobe borrows several ideas from modern natural language processing: using large scale naturalistic datasets, probing at scale across tasks as a means to understand black box systems, and evaluating on large benchmarks that test many different skills. For artificial networks, probe analysis attempts to decode attributes from different layers. It is one of the main vehicles used to shed light on the relationship and dependencies between tasks and the algorithms that networks learn. While prior neuroscience benchmarks tend to focus on a single or a very small number of tasks, Neuroprobe uses a fixed set of subjects with a large amount of data across many annotated tasks, which will allow us to create an integrated picture. Furthermore, the results obtained from Neuroprobe evaluations can yield time-orderings between different tasks and recover the functional relationships between tasks that reveal properties of the algorithms the brain uses. The main remaining bottleneck to achieving these type of results is that decoding performance for many tasks is very poor. We demonstrate a few tasks both with simple linear decoders and neural foundation models, then introduce a large number of additional attributes that should, in principle, be decodable but are not. Neuroprobe gives us an opportunity to build higher accuracy decoders, better neural foundation models that are tested across many tasks, and to bring neuroscience closer to the methodology that has worked so well in natural language understanding, and to ultimately discover the functional organization of the brain across many tasks. We make our code publicly available <sup>2</sup> and will maintain a leaderboard <sup>3</sup> to track model progress upon publication.

## 1 Introduction

The human brain constantly engages in a variety of processing tasks simultaneously: parsing speech, interpreting visual scenes, and performing social reasoning (Schurz et al., 2014). However, a cohesive picture of how these computations are organized across time and regions in the brain remains poorly understood. While modern neuroscience offers glimpses into individual functions, a central challenge is that typical experiments isolate one or two tasks at a time, often using simplified stimuli and contrived lab settings (Nastase et al., 2020). A solution suggests itself from the field of machine learning interpretability, which has developed methods to reverse engineer neural network black

---

\*Equal contribution.

<sup>2</sup><https://github.com/azaho/neuroprobe>

<sup>3</sup><https://neuroprobe.dev>

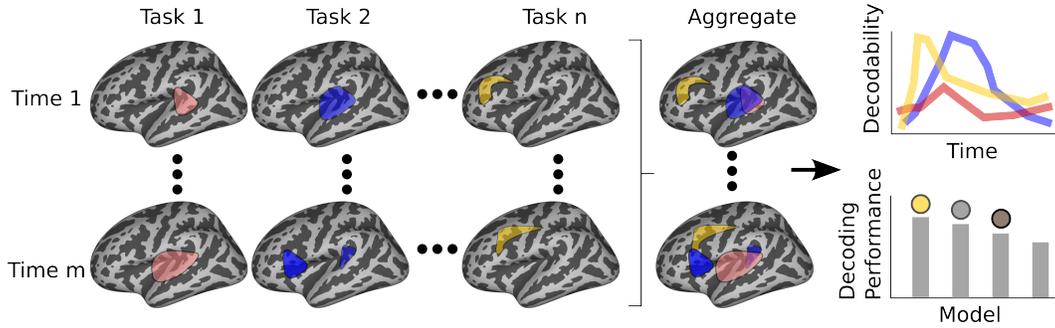


Figure 1: **Overview of Neuroprobe’s goals.** Neuroprobe consists of machine learning classification tasks derived from intracranial recordings aligned with annotated stimuli. By running a decoding analysis for each task, we can localize various aspects of multimodal language processing in the brain. Moreover, we can segment the neural recordings by time, repeat the decoding analyses across time bins, and discover a time evolution of each task. Previously, neuroscience experiments have been small, and focused on one task at a time. The results of our analyses can be combined to give a comprehensive picture of language processing in the brain. From this, two things can be achieved. First, we can derive neuroscience insights such as the relative timings for processing of certain tasks. Second, the tasks themselves can be used as a benchmark of neural decoding models.

34 boxes via probing experiments, e.g. Tenney et al. (2019); Alain & Bengio (2016). These methods are  
 35 powerful, but there is an obstacle in applying them to study the brain: decoding the contents of brain  
 36 activity remains a challenging task (Paninski & Cunningham, 2018). While intracranial data offers  
 37 high temporal and spatial resolution, the raw signals are noisy and high-dimensional. To these ends,  
 38 we introduce *Neuroprobe*, a benchmark that is designed both to be a setting in which neuroscience  
 39 probing experiment may be run *and* as a measure of progress to spur improvement of neural decoding  
 40 models.

41 Neuroprobe contains 19 decoding tasks that span vision and language, all on the same subjects  
 42 and the same neural recordings collected while subjects watched movies. Having many different  
 43 tasks on the same dataset allows one to derive constraints on the relationships between tasks, such  
 44 as: What is the temporal order between tasks across many subjects? Which tasks share neural real  
 45 estate? How does latency in one task influence latency in another task? These constraints can then  
 46 narrow the space of algorithms to regularize models of brain function. Unfortunately, as mentioned,  
 47 decoding today for many tasks is nowhere near accurate enough to systematically derive these kinds  
 48 of constraints. So, we develop a public leaderboard for hosting submissions to the Neuroprobe  
 49 benchmark. As submissions to the leaderboard increase, decoding accuracy will increase, in turn  
 50 raising our confidence in the spatial and temporal distribution of different tasks uncovered by the  
 51 probing experiments.

52 Meanwhile, on the modeling front, more and more foundation models are being developed for  
 53 neural recordings. There has been an explosion of neural foundation models as of late, including:  
 54 Neuroformer (Antoniades et al., 2024), BrainBERT (Wang et al., 2023), PopT (Chau et al., 2024),  
 55 STNDT (Le & Shlizerman, 2022), NDT2 (Ye et al., 2023), MBrain (Cai et al., 2023), Brant (Zhang  
 56 et al., 2023), MtM (Zhang et al., 2024b), and POYO (Azabou et al., 2023). Most of these models  
 57 are not tested on standardized decoding tasks. There are few cross-task decoding datasets at present  
 58 for testing new neural foundation models. This runs contrary to one of the main selling points of  
 59 foundation models for neuroscience, which is that they will improve decoding accuracy to enable  
 60 neuroscientists to run more experiments on a variety of tasks with less data. In addition, in a sense  
 61 the space of tasks determines the space of models considered, since only models that can show an  
 62 advantage are selected for and published. It is a long term problem for the community that larger  
 63 batteries of decoding tasks are not a common evaluation practice. This is already reflected in our  
 64 findings that state of the foundation models for neural recordings don’t make a massive different in  
 65 decoding performance for some tasks, and can even hurt it in a few cases (see section 4).

66 We have designed Neuroprobe to be usable by members of the ML community even if they have  
 67 no particular knowledge of neuroscience. Anyone can easily run models and contribute new ideas.

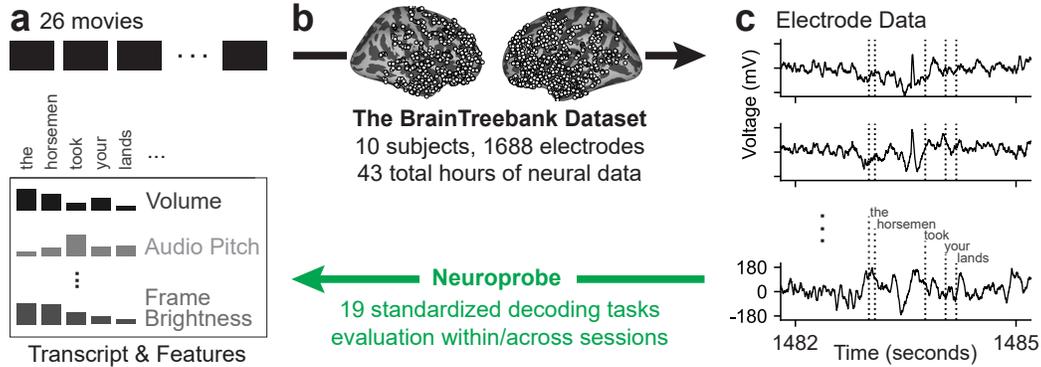


Figure 2: **From raw data to decoding tasks.** As part of the BrainTreebank dataset, 26 movies (a) are watched by 10 patients with stereoelectroencephalography electrodes implanted in various brain regions (b), and the local field potential from the implanted electrodes is recorded (c). Neuroprobe turns this dataset into an evaluation benchmark by segmenting the aligned data into various audio, language, and vision decoding tasks, such as, loudness and pitch of the audio, average pixel brightness, etc.

68 While Neuroprobe provides the analysis tools to interpret better decoding results. Lowering the  
 69 barriers to entry ensures that we have a healthier community and attracts many more researchers to  
 70 these problems.

71 Neuroprobe, see Figure 2, is derived from the Brain Treebank (Wang et al., 2024), which consists of  
 72 intracranial neural recordings aligned with the corresponding movie stimuli. The dataset contains  
 73 annotations from which we derive 19 decoding tasks, see Supplementary Table 1. We select the  
 74 BrainTreebank because it is at the scale at which modern NLP begins to operate and models being to  
 75 be understood (43 hours of recordings): comparable to datasets on low-resource languages.

76 In addition, we standardize a number of aspects of the benchmark. We select test/train splits in  
 77 different conditions: all the way from training and testing on the same subject and movie, to doing  
 78 cross-subject cross-movie decoding. We host a centralized website that aggregates results, both as a  
 79 whole and also by split-type and task, using a JSON schema to validate submissions.

80 Our contributions are:

- 81 1. A new large-scale multitask decoding benchmark: Neuroprobe.
- 82 2. Standardized splits and methods to rank neural foundation models and encourage their  
 83 development in a direction which benefits decoding tasks.
- 84 3. Results from a set of baselines and state-of-the-art models on Neuroprobe.
- 85 4. An early analysis of the timings and spatial distribution of different task processing pathways  
 86 in the brain.

87 In the long run we hope that Neuroprobe will both lead the way to an understanding of the general  
 88 architecture of the computations that the brain performs as well as bring the ML and neuroscience  
 89 communities into closer alignment by translating interesting neuroscience questions into questions  
 90 that are easily digested and then improved on by the ML community.

## 91 2 Related work

92 While there are many publicly available neural recordings that neural decoding models have been  
 93 developed on, neuroscience still suffers from a dearth of standardized, easy-to-run machine learning  
 94 benchmarks. This lack of defined decoding tasks, standardized train/test splits, and metrics make it  
 95 difficult to compare models.

96 **Neural recording datasets** The most recently developed models for neural data have relied on  
 97 several widely accessible datasets. For non-invasive EEG decoding, datasets from Zheng & Lu  
 98 (2015); Grootswagers et al. (2022); Bhattasali et al. (2020); Tangermann et al. (2012); Obeid &  
 99 Picone (2016); Broderick et al. (2018); Brennan & Hale (2019) have been used in the construction of

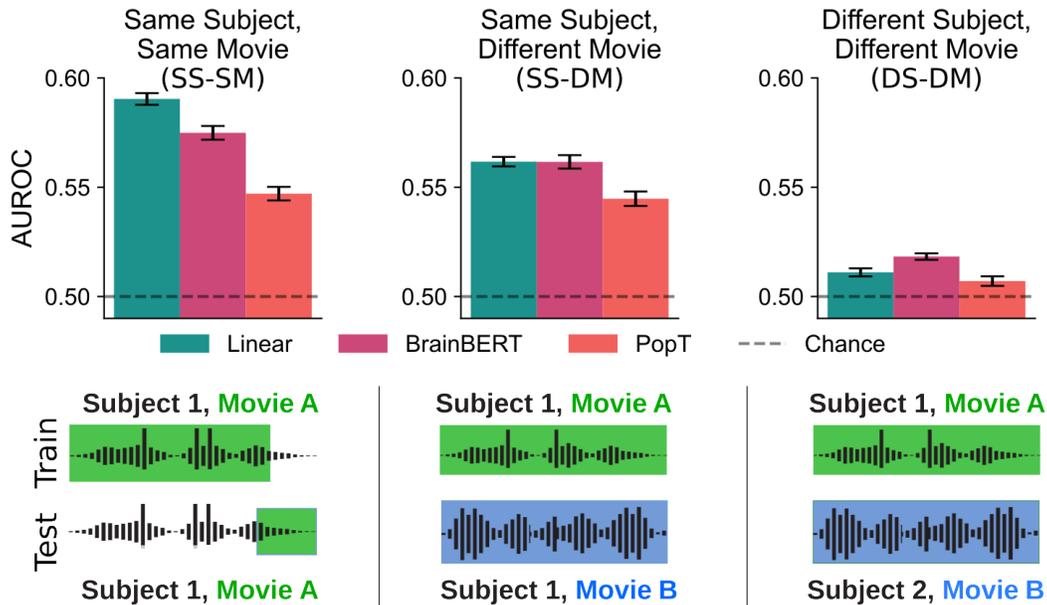


Figure 3: **Neuroprobe splits.** We perform analyses on three different types of splits. In *same subject/same movie* (SS-SM) we train on data from one subject and one movie segment, and evaluate on the same subject, but another segment of the same movie. Performance is measured via cross-validation. In *same subject/different movie* (SS-DM), we train on data from one subject and from one movie. Then, we evaluate on another movie. In *different subject/different movie* (DS-DM), we train on data from one subject and one movie and evaluate on data from an entirely different subject and movie. This is the most challenging split.

100 models such as those proposed by Jiang et al. (2024); Yang et al. (2023); Yuan et al. (2024); Défossez  
 101 et al. (2023). For fMRI decoding, (Wehbe et al., 2014; LeBel et al., 2023; Nastase et al., 2021; Li  
 102 et al., 2022; Allen et al., 2022) have led to models such as those proposed by Scotti et al. (2024);  
 103 Ozcelik & VanRullen (2023). For MEG decoding, Jan-Mathijs et al. (2019); Hebart et al. (2023)  
 104 have lead to models such as those proposed by Défossez et al. (2023); Benchetrit et al.. For neural  
 105 spike decoding Perich et al. (2025); Churchland et al. (2024); Manley et al. (2024); IBL (2024) have  
 106 lead to models such as those proposed by Azabou et al. (2023); Zhang et al. (2024a). For broadband  
 107 intracranial neural activity, datasets from (Peterson et al., 2022; Wang et al., 2024; Nejedly et al.,  
 108 2020) have fueled the development of models proposed by (Peterson et al., 2021; Wang et al., 2023;  
 109 Chau et al., 2024) However, these datasets do not provide rigorous splits or testing guidelines, so  
 110 each model is difficult to compare to others.

111 **Existing neural data benchmarks** There are a few benchmarks involving neural data. Some of the  
 112 earliest involve EEG BCI decoding (Tangermann et al., 2012), but are limited in data quality and scale  
 113 by today’s standards. The NaturalScenesDataset (Allen et al., 2022) is close to being a benchmark in  
 114 that they have splits, but it primarily benchmarks fMRI data, and focuses on visual processing. The  
 115 clinical-grade Temple University Hospital EEG dataset (Obeid & Picone, 2016) can also be used as a  
 116 benchmark, but it only contains EEG and has the labels are limited to seizure detection. Benchmarks  
 117 for neural spikes are proposed by Pei et al. (2021); Karpowicz et al. (2024); Willett et al. (2023);  
 118 Lueckmann et al. (2025), but these only contain spiking information rather than broadband signals  
 119 from ECoG or sEEG that capture more neural activity (Parvizi & Kastner, 2018). A benchmark like  
 120 Neuroprobe for high fidelity intracranial signals with corresponding challenging naturalistic language  
 121 stimuli is still needed to allow the field to progress forward in building better neural decoding models.

### 122 3 Approach

123 **Brain Treebank** Neuroprobe is an evaluation-only benchmark environment that uses the raw data  
 124 from the BrainTreebank (Wang et al., 2024), a publicly available dataset released under a CC

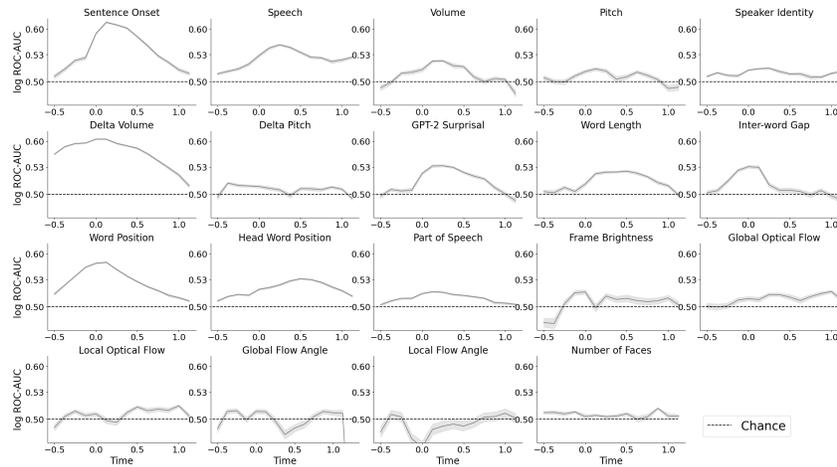


Figure 4: **Neuroprobe enables tracking of information processing in the brain across tasks.** A linear model is fit for a sliding 125ms window of activity. Here, we show the performance of the most decodable 100 electrodes per each task. Error bars show standard error across electrodes. Performance is plotted on a log scale to show trends for tasks that have lower decodability. The x-axis shows time, where  $t = 0$  corresponds with word onset. By plotting decoding performance across time, the time course of information availability for each task becomes visible. Audio-linguistic tasks, such as Speech vs. Non-speech, are most decodable closest to word onset.

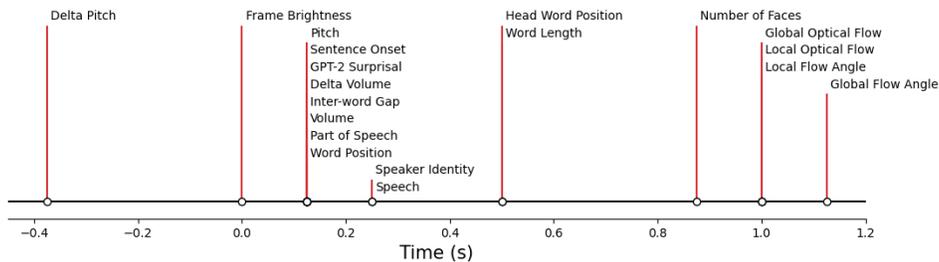


Figure 5: **Time-ranking of decodability** A simple method of finding relationships between tasks is to look at when each task is decodable. Consistencies in this order across subjects are an indication of a dependency between tasks. A shortcut to that, is to further restrict ourselves to when each task achieves maximum decodability. Note that we use a window of 125 ms which gives fairly coarse temporal localizations, which is why many tasks overlap. We can already observe some patterns from these results, even with poor decoding accuracy. Notably, *Word head position*, a semantic feature that pertains to the position of the dependency parse head, is decoded later than other language features. A caveat should be offered that these timings are dependent on the type of decoding analysis being performed. As different decoding methods are developed which solidify our ability to decode each task, it is certain that these ordering will change.

125 BY 4.0 license. The Brain Treebank is a large-scale dataset of intracranial electrophysiological  
 126 recordings (stereoelectroencephalography; sEEG) collected while 10 human subjects (5 male, 5  
 127 female, ages 4–19; Supplementary Table 3) watched 26 total Hollywood movies (Supplementary  
 128 Table 4). Electrode placements for each subject and their speech-selective responses are shown in  
 129 Supplementary Figure 10. Spanning 43 hours of neural activity, the dataset aligns recorded brain  
 130 signals with transcribed and manually corrected speech, word onsets, and universal dependency  
 131 parses across the 223,068 words in 38,572 sentences. This dataset enables the systematic evaluation  
 132 of computational models on multimodal neural decoding tasks.

133 **Decoding tasks** We use the movie annotations and the alignment with the corresponding neural data  
 134 to create a suite of 19 decoding tasks, spanning visual, audio, and language domains. For every task,  
 135 the neural data is the input and the annotation label is the target output, where we formalize all of the

136 tasks as binary classification by thresholding the labels. For example, for the GPT2 Surprisal task,  
137 the positive label corresponds to surprisal annotations above the 75%th percentile of the distribution  
138 within a session, and the negative label to the values below the 25%th percentile. For non-scalar labels  
139 (such as speaker identity or part of speech of the word) we pick a main target class (i.e. most frequent  
140 speaker, or Verb for the part of speech task), and formulate the task as one-versus-rest classification.  
141 See more details in Appendix A.

142 **Splits** The Neuroprobe evaluation takes place across three different types of splits. For the *same*  
143 *subject/same movie* (SS-SM) splits, train data and test data come from a single movie-viewing session.  
144 Decoding results are cross-validated with an 80-20 train-test split. Importantly, the indices for the  
145 cross-validation splits are not drawn from the whole movie uniformly, but rather the train examples  
146 are taken from a single contiguous block and the validation examples are taken from a separate block.  
147 This is done to prevent models from over-fitting to auto-correlation in the signal.

148 For the *same subject-different movie* SS-DM split, the train data consists of examples drawn from the  
149 longest movie viewed by a given patient, and the test data comes from the second longest movie.

150 For the *different subject-different movie* DS-DM split, the train data consists of data from a single  
151 session (trial 4), viewed by subject 2, chosen because this is the longest trial and the subject with  
152 the most electrodes in both hemispheres. Testing then consists of the average performance across  
153 selected sessions for all other subjects (see Appendix F). This split in particular presents a demanding  
154 test of model generalizability, especially since electrode placements vary widely between patients  
155 (see Figure 10).

156 **Experiments** In Neuroprobe, experiments can either be performed at the *single-electrode* level or  
157 the *population* level, i.e., using all electrodes in a given subject as model input. To give a sense of  
158 the types of neuroscience insights that can be derived from Neuroprobe, we perform a collection of  
159 single-electrode analyses across the SS-SM splits for all BrainTreebank sessions. In particular, for  
160 each task, we fit a linear classifier to do decoding over a fixed window of activity (250 ms). This  
161 window slides along a longer period, from 0.5s before word onset to 1.25s after word onset, with a  
162 stride of 125ms. This provides a picture of the time-course of decodability in the brain. Electrodes  
163 marked as corrupted in the original BrainTreebank dataset are excluded. See Section 4.

164 **Neuroprobe-Lite Benchmark** Outside of analyses described above, for the purposes of comparing  
165 models, running experiments over all sessions and electrodes is prohibitively expensive. To this end,  
166 we subset the data to create Neuroprobe-lite by selecting a smaller portion of subjects and sessions (6  
167 subjects, 2 trials each) for training and evaluation.

168 Furthermore, the total number of electrodes per subject is capped at 120. The electrodes in  
169 Neuroprobe-lite were chosen specifically to cover as much of the brain in each participant as possible.  
170 This was done by randomly taking a specified proportion of electrodes from every probe, to ensure  
171 that every probe is represented in the Neuroprobe-lite data features. This ensures that the input for  
172 each task is standardized matrix which has predictable memory and computational requirements. We  
173 maintain a public leaderboard which will display model performance on this benchmark, both on the  
174 single-electrode and population level; see Supplemental fig. 12.

175 **Models** To show the utility of the Neuroprobe tasks as a benchmark, we evaluate on a few baselines  
176 and models. For the purposes of benchmarking, all models are run on Neuroprobe-lite (see above).  
177 All inputs are given as a population, i.e., the data from all electrodes is provided as input, concatenated.  
178 The models we benchmark span the range of simple classifiers to large, pretrained models. These  
179 include three linear regression models, which take as input either the raw voltage time-series inputs,  
180 Fourier transform input, or Short-time Fourier transform (STFT) inputs. For pretrained models, we  
181 also train a regression on BrainBERT (Wang et al., 2023) inputs, and fine-tune a linear layer on  
182 top of a pretrained PopT (Chau et al., 2024), a pretrained transformer for encoding arbitrary sets of  
183 electrodes. More details on the models available in Appendix H.

184 **Metric calculations** The primary evaluation metric was the Area Under the Receiver Operating  
185 Characteristic curve (AUROC), aggregated across electrodes. We adjusted the aggregation strategy to  
186 be compatible with each model to obtain the different subjects-different movie DS/DM results shown  
187 in Figure 3. Before running our linear regressions, we preprocessed the neural data to represent  
188 activity in each cortical region (using averaging per subject/trial pair), as defined from the 34 regions  
189 by the Desikan-Killiany atlas. Similarly, we ran BrainBERT, with the same region averaging strategy.

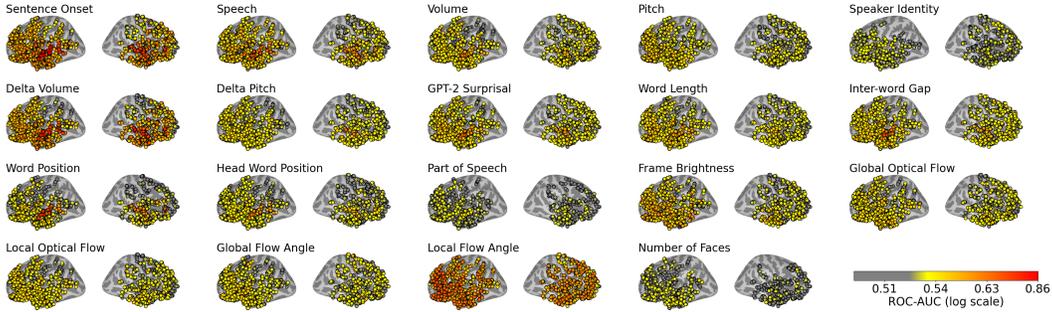


Figure 6: **Distribution of task processing throughout the brain** A linear decoder is trained on the *single-subject/single-movie* split. Color shows ROC-AUC on a logarithmic scale. Performance is computed by averaging over cross-validation folds ( $k = 5$ ) and movies and then taking a max over time bins. Language features like *Sentence Onset* and *GPT-2 Surprisal* are most decodable in the temporal and frontal lobes.

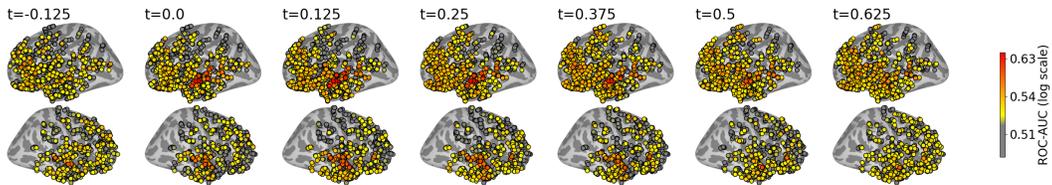


Figure 7: **Time evolution of surprisal decodability throughout the brain.** The decodability of features vary in both time and space. Close to word onset ( $t = 0$ ), surprisal is most decodable in the superior temporal gyrus. Time zero here refers to the onset of a given word. Most words are interior to sentences or to conversations. Since most modern surprisal metrics are contextualized, one can immediately predict surprisal even from the neural activity left over from prior words. As time progresses, surprisal becomes more decodable in the frontal areas. Full progressions for all tasks can be seen in Appendix L and in a movie at this url: [https://neuroprobe.dev/neuroprobe\\_time\\_course.mp4](https://neuroprobe.dev/neuroprobe_time_course.mp4).

190 For the PopulationTransformer we use all electrodes that can be bipolar-rereferenced and are in the  
 191 set of ‘clean’ electrodes (see (Chau et al., 2024)) for evaluation. No accommodation for the DS/DM  
 192 split was necessary for the PopulationTransformer, which is designed to handle subject-transfer.

## 193 4 Results

194 **Timing analysis** To investigate the time course of linguistic information processing in the brain, we  
 195 aligned neural data to word onsets and split it into narrow time-bins (width = 125ms), and train  
 196 a separate linear decoder on each bin for multiple tasks. Decodability is computed as the average  
 197 across cross-validation folds ( $k = 5$ ). For each task, we restrict our attention to the top 100 electrodes  
 198 with the highest decodability. Decoding performance as a function of time shows the course of  
 199 processing after the word onset ( $t = 0$ , Figure 4). Interestingly, the beginning of a new sentence can  
 200 be decoded with better-than-chance AUROC even before the word onset ( $\mu = 0.53$ ,  $\sigma_M = 0.0015$  at  
 201  $-250$ ms), hinting at the predictive nature of processing. Moreover, we can find a time-ranking of  
 202 features by looking at when decodability peaks for each feature (Figure 5). For example, we note  
 203 that the high-level semantic feature ‘word head position’ is decodable only later (decodability peaks  
 204 at  $t = 0.5s$  vs. volume and pitch at  $t = 0.125s$ ).

205 **Spatial analysis** By examining the linear decodability of features, a picture emerges of which features  
 206 modulate activity in which areas of the brain (Figure 6). Using the single electrode analysis, we find  
 207 that audio-linguistic tasks such as ‘sentence onset’, ‘speech vs. non-speech’, ‘delta volume’ are most  
 208 decodable in the superior temporal gyrus, especially close to Herschel’s and Wernicke’s area, with

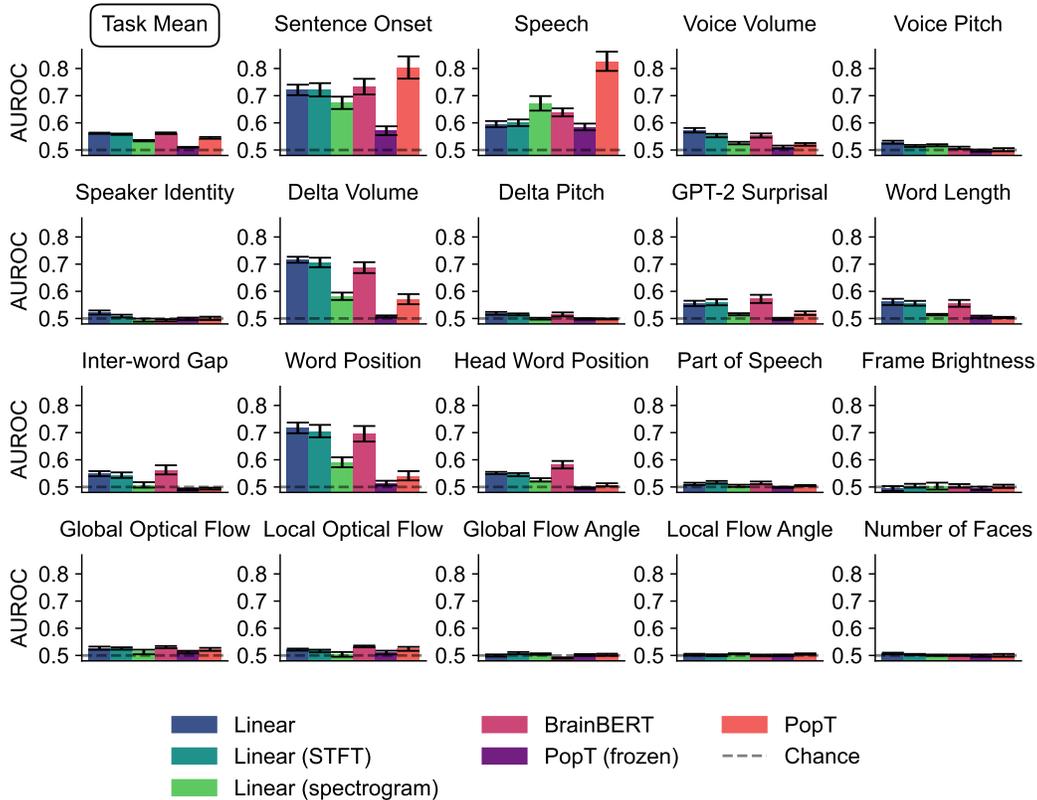


Figure 8: **Performance of baseline models on the 19 tasks of Neuroprobe.** Evaluation is done on the same subject, same trial (SS-ST), using 5-fold cross-validation. Normalized audio volume traces and the distribution of detected faces with corresponding word counts are shown in Supplementary Figures 9 and 11, respectively. The performance of four models is shown: (1) logistic regression either from raw voltage signal of all electrodes to the labels, or (2) from the spectrogram of the signal to the labels, as well as (3) BrainBERT (Wang et al., 2023) and (4) PopulationTransformer (Chau et al., 2024). Neural data was cut to include one second following each word onset. In case of multi-class classification, AUROC was computed using a one-vs-all strategy and averaged together. Performance on different trials for the same subject were averaged together. Error bars denote s.e.m. across all subjects. These results can be seen in tabular form in Appendix I.

209 average AUROCs of 0.61, 0.55, and 0.62, respectively in the gyrus of the temporal transverse. Here  
 210 region results are given with respect to the Destrieux atlas; see Appendix M.

211 **Spatio-Temporal analysis** We do a deep dive on the surprisal feature and show that after word onset,  
 212 it is most decodable in the temporal lobe (AUROC = 0.58 at  $t = 0$  in the transverse temporal),  
 213 but decodability spreads to the frontal lobe as time progresses (AUROC = 0.50 at  $t = -0.125$   
 214 and AUROC = 0.52 at  $t = 0.5$ ); see Figure 7. A movie of this for all tasks can be seen at  
 215 [https://neuroprobe.dev/neuroprobe\\_time\\_course.mp4](https://neuroprobe.dev/neuroprobe_time_course.mp4).

216 **Comparison of basic decoding methods on Neuroprobe.** We compare the performance of two  
 217 simple baseline models—logistic regression applied to raw voltage signals and logistic regression  
 218 applied to spectrogram features—across the 19 decoding tasks in Neuroprobe. Performance is  
 219 evaluated using area under the receiver operating characteristic curve (AUROC), with chance-level  
 220 performance ( $ROC = 0.5$ ) included for reference. We also compare with BrainBERT and PopT using  
 221 their publicly released off-the-shelf-weights. Because of this there may be some discrepancy due to the  
 222 fact that both models were trained on 5s intervals, whereas we train on 1s intervals across all models  
 223 for consistency. In general, linear decoding is very good (see Figure 3), achieving the best overall  
 224 performance on the SS/SM ( $0.590 \pm 0.003$ ) split, with the second best model being BrainBERT  
 225 ( $0.575 \pm 0.003$ ). On the SS/DM split, the linear baseline tied BrainBERT ( $0.562 \pm 0.002$  vs

226  $0.562 \pm 0.003$ , respectively), outperforming PopulationTransformer ( $0.545 \pm 0.003$ ). But BrainBERT  
227 performs the best on the difficult DS/DM split ( $0.518 \pm 0.001$ ) with the next best model being the  
228 linear baseline ( $0.511 \pm 0.002$ ).

229 Finally, for SS-SM, a breakdown by task can be seen in Figure 8. The PopulationTransformer, despite  
230 being pretrained, underperforms on many tasks, but achieves the highest performance on the Sentence  
231 Onset and Speech vs. Non-speech tasks.

## 232 5 Conclusion

233 Neuroprobe can be used in several ways by different communities: (1) Machine learning practitioners  
234 can contribute by improving decoding performance. (2) At the intersection of ML and neuroscience,  
235 Neuroprobe can be used to assess how good a given neural foundation model is at improving decoding  
236 accuracy. (3) Neuroscientists can use Neuroprobe to uncover relationships between different tasks  
237 that the brain executes which puts constraints on the kinds of algorithms our brains are using.

238 Using Neuroprobe, questions about processing in the brain become machine learning decoding tasks  
239 which can be rapidly iterated on. This will drive improvements both in decoding ability and the ability  
240 to draw neuroscience conclusions from large scale data. As we have seen in other fields, this can also  
241 lead to a virtuous cycle in which neuroscientists are encouraged to share more datasets to the effort.

242 Despite the weakness of current decoding models, Neuroprobe can still find interesting trends in both  
243 the spatial and temporal organization of tasks in the brain. As decoding models improve, the clarity  
244 of such findings will improve and their variance will decline. Each decoding task induces a map  
245 across the brain of when and where processing specific to that task is performed. By overlaying many  
246 of these maps, a functional picture of the brain emerges of which language, vision, and audio features  
247 modulate activity in each region. We see this approach as a way of answering the long-standing  
248 neuroscience question: What is the underlying circuit basis of language processing in the brain?

249 **Limitations** Our decoding results from the baselines we tested are low for a few tasks, such as  
250 speaker identity and pitch, and thus drawing any conclusions from their results is fraught. While  
251 our data offers unprecedented combination of scale and resolution, it is collected from a clinical  
252 population undergoing invasive monitoring, and results should not be overgeneralized. We only have  
253 10 subjects currently. This is because it is difficult to obtain this kind of data, which requires invasive  
254 surgery to implant electrodes. However, each subject has many sessions.

255 **Broader impacts** Neuroprobe provides a standardized benchmark for evaluating models of human  
256 brain activity, with potential applications in neuroscience, machine learning, and clinical technolo-  
257 gies such as brain-computer interfaces. By releasing our data, code, and leaderboard, we aim to  
258 democratize access to high-quality neural benchmarks and foster cross-disciplinary collaboration.

259 **Future work** Our framework is general enough to accommodate future annotations, allowing for  
260 investigations of low-level language processing, such as part of speech, or high-level semantic  
261 processing such as thematic roles or language model embeddings. We also seek, in near-term future  
262 work, to add to the library of tasks and datasets in Neuroprobe. As we continue to build out the  
263 benchmark, researchers will be able to study the question of how various tasks interact with each  
264 other.

## 265 6 Acknowledgements

266 This work was supported by the Center for Brains, Minds, and Machines, NSF STC award CCF-  
267 1231216, the NSF award 2124052, the MIT CSAIL Machine Learning Applications Initiative, the  
268 MIT-IBM Watson AI Lab, the CBMM-Siemens Graduate Fellowship, the DARPA Mathematics  
269 for the DIScovery of ALgorithms and Architectures (DIAL) program, the DARPA Knowledge  
270 Management at Scale and Speed (KMASS) program, the DARPA Machine Common Sense (MCS)  
271 program, the United States Air Force Research Laboratory and the Department of the Air Force  
272 Artificial Intelligence Accelerator under Cooperative Agreement Number FA8750-19-2-1000, and the  
273 Air Force Office of Scientific Research (AFOSR) under award number FA9550-21-1-0399. This work  
274 also has been supported by ONR award N00014-19-1-2584, by NSF-CISE award IIS-2151077 under  
275 the Robust Intelligence program, by the ARO-MURI award W911NF-23-1-0277, by the Simons  
276 Foundation SCGB program 1181110, the K. Lisa Yang ICoN Center, the Caltech Chen Institute,

277 and the Caltech Carver Mead New Adventures Fund. The views and conclusions contained in  
278 this document are those of the authors and should not be interpreted as representing the official  
279 policies, either expressed or implied, of the Department of the Air Force or the U.S. Government.  
280 The U.S. Government is authorized to reproduce and distribute reprints for Government purposes,  
281 notwithstanding any copyright notation herein.

## 282 **References**

- 283 International brain lab. <https://internationalbrainlab.org>, 2024. Accessed: 2024-11-23.
- 284 Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes.  
285 In *International Conference on Learning Representations (ICLR)*, 2016.
- 286 Emily J Allen, Ghislain St-Yves, Yihan Wu, Jesse L Breedlove, Jacob S Prince, Logan T Dowlle,  
287 Matthias Nau, Brad Caron, Franco Pestilli, Ian Charest, et al. A massive 7t fmri dataset to bridge  
288 cognitive neuroscience and artificial intelligence. *Nature neuroscience*, 25(1):116–126, 2022.
- 289 Antonis Antoniadis, Yiyi Yu, Joseph Canzano, William Wang, and Spencer LaVere Smith. Neuro-  
290 former: Multimodal and Multitask Generative Pretraining for Brain Data, March 2024.
- 291 Mehdi Azabou, Vinam Arora, Venkataramana Ganesh, Ximeng Mao, Santosh Nachimuthu, Michael J.  
292 Mendelson, Blake Richards, Matthew G. Perich, Guillaume Lajoie, and Eva L. Dyer. A Unified,  
293 Scalable Framework for Neural Population Decoding, October 2023.
- 294 Yohann Benchetrit, Hubert Banville, and Jean-Remi King. Brain decoding: toward real-time  
295 reconstruction of visual perception. october 2023. In URL <https://openreview.net/forum>.
- 296 Shohini Bhattachali, Jonathan Brennan, Wen-Ming Luh, Berta Franzluebbers, and John Hale. The alic  
297 datasets: fMRI & EEG observations of natural language comprehension. In Nicoletta Calzolari,  
298 Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi,  
299 Hitoshi Isahara, Bente Maegaard, Joseph Mariani, H el ene Mazo, Asuncion Moreno, Jan Odijk, and  
300 Stelios Piperidis (eds.), *Proceedings of the Twelfth Language Resources and Evaluation Conference*,  
301 pp. 120–125, Marseille, France, May 2020. European Language Resources Association. ISBN  
302 979-10-95546-34-4. URL <https://aclanthology.org/2020.lrec-1.15/>.
- 303 Jonathan R Brennan and John T Hale. Hierarchical structure guides rapid linguistic predictions  
304 during naturalistic listening. *PloS one*, 14(1):e0207741, 2019.
- 305 Michael P Broderick, Andrew J Anderson, Giovanni M Di Liberto, Michael J Crosse, and Edmund C  
306 Lalor. Electrophysiological correlates of semantic dissimilarity reflect the comprehension of  
307 natural, narrative speech. *Current Biology*, 28(5):803–809, 2018.
- 308 Donghong Cai, Junru Chen, Yang Yang, Teng Liu, and Yafeng Li. MBrain: A Multi-channel  
309 Self-Supervised Learning Framework for Brain Signals, June 2023.
- 310 Geeling Chau, Christopher Wang, Sabera Talukder, Vighnesh Subramaniam, Saraswati Soedarmadji,  
311 Yisong Yue, Boris Katz, and Andrei Barbu. Population Transformer: Learning Population-level  
312 Representations of Neural Activity, October 2024.
- 313 Mark Churchland, John P. Cunningham, Matthew T. Kaufman, Justin D. Foster, Paul Nuyujukian,  
314 Stephen I. Ryu, and Krishna V. Shenoy. Neural population dynamics during reaching. Data set,  
315 2024. URL <https://dandiarchive.org/dandiset/000070/draft>.
- 316 Alexandre D efosse, Charlotte Caucheteux, J er emy Rapin, Ori Kabeli, and Jean-R emi King. Decod-  
317 ing speech perception from non-invasive brain recordings. *Nature Machine Intelligence*, 5(10):  
318 1097–1107, 2023.
- 319 Tijl Grootswagers, Iris Zhou, Austin K. Robinson, et al. Human eeg recordings for 1,854 concepts  
320 presented in rapid serial visual presentation streams. *Scientific Data*, 9:3, 2022. doi: 10.1038/  
321 s41597-021-01102-7. URL <https://doi.org/10.1038/s41597-021-01102-7>.

- 322 Martin N Hebart, Oliver Contier, Lina Teichmann, Adam H Rockter, Charles Y Zheng, Alexis Kidder,  
323 Anna Corriveau, Maryam Vaziri-Pashkam, and Chris I Baker. Things-data, a multimodal collection  
324 of large-scale datasets for investigating object representations in human brain and behavior. *Elife*,  
325 12:e82580, 2023.
- 326 Schoffelen Jan-Mathijs, Robert Oostenveld, Lam Nietzsche HL, Uddén Julia, Hultén Annika, and  
327 Peter Hagoort. A 204-subject multimodal neuroimaging dataset to study language processing.  
328 *Scientific Data*, 6(1), 2019.
- 329 Wei-Bang Jiang, Li-Ming Zhao, and Bao-Liang Lu. Large brain model for learning generic represen-  
330 tations with tremendous eeg data in bci. *arXiv preprint arXiv:2405.18765*, 2024.
- 331 Brianna M Karpowicz, Joel Ye, Chaofei Fan, Pablo Tostado-Marcos, Fabio Rizzoglio, Clay Wash-  
332 ington, Thiago Scodeler, Diogo de Lucena, Samuel R Nason-Tomaszewski, Matthew J Mender,  
333 et al. Few-shot algorithms for consistent neural decoding (falcon) benchmark. *Advances in Neural*  
334 *Information Processing Systems*, 37:76578–76615, 2024.
- 335 Trung Le and Eli Shlizerman. STNDT: Modeling Neural Population Activity with a Spatiotemporal  
336 Transformer, June 2022.
- 337 Alexandre LeBel, Laura Wagner, Siddharth Jain, et al. A natural language fmri dataset for voxelwise  
338 encoding models. *Scientific Data*, 10:555, 2023. doi: 10.1038/s41597-023-02437-z. URL  
339 <https://doi.org/10.1038/s41597-023-02437-z>.
- 340 Jixing Li, Shohini Bhattasali, Shaolei Zhang, et al. Le petit prince multilingual naturalistic fmri  
341 corpus. *Scientific Data*, 9:530, 2022. doi: 10.1038/s41597-022-01625-7. URL <https://doi.org/10.1038/s41597-022-01625-7>.
- 343 Jan-Matthis Lueckmann, Alexander Immer, Alex Bo-Yuan Chen, Peter H Li, Mariela D Petkova, Nir-  
344 mala A Iyer, Luuk Willem Hesselink, Aparna Dev, Gudrun Ihrke, Woohyun Park, et al. Zapbench:  
345 A benchmark for whole-brain activity prediction in zebrafish. *arXiv preprint arXiv:2503.02618*,  
346 2025.
- 347 Jason Manley, Sihao Lu, Kevin Barber, Jeffrey Demas, Hyewon Kim, David Meyer, Fran-  
348 cisca Martínez Traub, and Alipasha Vaziri. Simultaneous, cortex-wide dynamics of up to 1  
349 million neurons reveal unbounded scaling of dimensionality with neuron number. *Neuron*, 112  
350 (10):1694–1709.e5, 2024. ISSN 0896-6273. doi: <https://doi.org/10.1016/j.neuron.2024.02.011>.  
351 URL <https://www.sciencedirect.com/science/article/pii/S0896627324001211>.
- 352 Samuel A. Nastase, Ariel Goldstein, and Uri Hasson. Keep it real: Rethinking the primacy of  
353 experimental control in cognitive neuroscience. *NeuroImage*, 222:117254, 2020. doi: 10.1016/  
354 j.neuroimage.2020.117254. URL <https://doi.org/10.1016/j.neuroimage.2020.117254>.  
355 Open access under CC license.
- 356 Samuel A. Nastase, Yung-Fang Liu, Harrison Hillman, et al. The “narratives” fmri dataset for  
357 evaluating models of naturalistic language comprehension. *Scientific Data*, 8:250, 2021. doi:  
358 10.1038/s41597-021-01033-3. URL <https://doi.org/10.1038/s41597-021-01033-3>.
- 359 Petr Nejedly, Vaclav Kremen, Vladimir Sladky, Jan Cimbálik, Petr Klimes, Filip Plesinger, Filip  
360 Mivalt, Vojtech Travnicek, Ivo Viscor, Martin Pail, et al. Multicenter intracranial eeg dataset for  
361 classification of graphoelements and artifactual signals. *Scientific data*, 7(1):179, 2020.
- 362 Iyad Obeid and Joseph Picone. The temple university hospital eeg data corpus. *Frontiers in*  
363 *neuroscience*, 10:196, 2016.
- 364 Furkan Ozcelik and Rufin VanRullen. Natural scene reconstruction from fmri signals using generative  
365 latent diffusion. *Scientific Reports*, 13(1):15666, 2023.
- 366 Liam Paninski and John P. Cunningham. Neural data science: Accelerating the experiment-analysis-  
367 theory cycle in large-scale neuroscience. *Current Opinion in Neurobiology*, 50:232–241, 2018.  
368 doi: 10.1016/j.conb.2018.04.007. URL <https://doi.org/10.1016/j.conb.2018.04.007>.  
369 Copyright © 2018 Elsevier Ltd. All rights reserved.

- 370 Josef Parvizi and Sabine Kastner. Promises and limitations of human intracranial electroencephalog-  
371 raphy. *Nature Neuroscience*, 21(4):474–483, 2018. doi: 10.1038/s41593-018-0108-2. URL  
372 <https://doi.org/10.1038/s41593-018-0108-2>.
- 373 Felix Pei, Joel Ye, David M. Zoltowski, Anqi Wu, Raed H. Chowdhury, Hansem Sohn, Joseph E.  
374 O’Doherty, Krishna V. Shenoy, Matthew T. Kaufman, Mark Churchland, Mehrdad Jazayeri, Lee E.  
375 Miller, Jonathan Pillow, Il Memming Park, Eva L. Dyer, and Chethan Pandarinath. Neural latents  
376 benchmark ’21: Evaluating latent variable models of neural population activity. In *Advances in*  
377 *Neural Information Processing Systems (NeurIPS), Track on Datasets and Benchmarks*, 2021.  
378 URL <https://arxiv.org/abs/2109.04463>.
- 379 Matthew G. Perich, Lee E. Miller, Mehdi Azabou, and Eva L. Dyer. Long-term recordings of  
380 motor and premotor cortical spiking activity during reaching in monkeys. Data set, 2025. URL  
381 <https://doi.org/10.48324/dandi.000688/0.250122.1735>.
- 382 Steven M Peterson, Zoe Steine-Hanson, Nathan Davis, Rajesh PN Rao, and Bingni W Brunton.  
383 Generalized neural decoders for transfer learning across participants and recording modalities.  
384 *Journal of Neural Engineering*, 18(2):026014, 2021.
- 385 Steven M Peterson, Satpreet H Singh, Benjamin Dichter, Michael Scheid, Rajesh PN Rao, and  
386 Bingni W Brunton. Ajile12: Long-term naturalistic human intracranial neural recordings and pose.  
387 *Scientific data*, 9(1):184, 2022.
- 388 Matthias Schurz, Joaquim Radua, Markus Aichhorn, Fabio Richlan, and Josef Perner. Fractionating  
389 theory of mind: A meta-analysis of functional brain imaging studies. *Neuroscience & Biobehavioral*  
390 *Reviews*, 42:9–34, 2014.
- 391 Paul S Scotti, Mihir Tripathy, Cesar Kadir Torrico Villanueva, Reese Kneeland, Tong Chen, Ashutosh  
392 Narang, Charan Santhirasegaran, Jonathan Xu, Thomas Naselaris, Kenneth A Norman, et al.  
393 Mindeye2: Shared-subject models enable fmri-to-image with 1 hour of data. *arXiv preprint*  
394 *arXiv:2403.11207*, 2024.
- 395 Michael Tangermann, Klaus-Robert Müller, Ad Aertsen, Niels Birbaumer, Christoph Braun, Clemens  
396 Brunner, Robert Leeb, Carsten Mehring, Kai J Miller, Gernot R Müller-Putz, et al. Review of the  
397 bci competition iv. *Frontiers in neuroscience*, 6:55, 2012.
- 398 Ian Tenney, Dipanjan Das, and Ellie Pavlick. Bert rediscovered the classical nlp pipeline. *arXiv*  
399 *preprint arXiv:1905.05950*, 2019.
- 400 Christopher Wang, Vighnesh Subramaniam, Adam Uri Yaari, Gabriel Kreiman, Boris Katz, Ignacio  
401 Cases, and Andrei Barbu. BrainBERT: Self-supervised representation learning for intracranial  
402 recordings, February 2023.
- 403 Christopher Wang, Adam Uri Yaari, Aaditya K Singh, Vighnesh Subramaniam, Dana Rosenfarb, Jan  
404 DeWitt, Pranav Misra, Joseph R. Madsen, Scellig Stone, Gabriel Kreiman, Boris Katz, Ignacio  
405 Cases, and Andrei Barbu. Brain treebank: Large-scale intracranial recordings from naturalistic  
406 language stimuli, 2024. URL <https://arxiv.org/abs/2411.08343>.
- 407 Leila Wehbe, Brian Murphy, Partha Talukdar, Alona Fyshe, Aaditya Ramdas, and Tom Mitchell.  
408 Simultaneously uncovering the patterns of brain regions involved in different story reading subpro-  
409 cesses. *PLOS ONE*, 9(11):e112575, November 2014. ISSN 1932-6203. doi: 10.1371/journal.pone.  
410 0112575. URL <http://dx.plos.org/10.1371/journal.pone.0112575>.
- 411 Francis R Willett, Erin M Kunz, Chaofei Fan, Donald T Avansino, Guy H Wilson, Eun Young  
412 Choi, Foram Kamdar, Matthew F Glasser, Leigh R Hochberg, Shaul Druckmann, et al. A high-  
413 performance speech neuroprosthesis. *Nature*, 620(7976):1031–1036, 2023.
- 414 Chaoqi Yang, M Westover, and Jimeng Sun. Biot: Biosignal transformer for cross-data learning in  
415 the wild. *Advances in Neural Information Processing Systems*, 36:78240–78260, 2023.
- 416 Joel Ye, Jennifer L. Collinger, Leila Wehbe, and Robert Gaunt. Neural Data Transformer 2: Multi-  
417 context Pretraining for Neural Spiking Activity, September 2023.

- 418 Zhizhang Yuan, Fanqi Shen, Meng Li, Yuguo Yu, Chenhao Tan, and Yang Yang. Brainwave: A brain  
419 signal foundation model for clinical applications. *arXiv preprint arXiv:2402.10251*, 2024.
- 420 Daoze Zhang, Zhizhang Yuan, Yang Yang, Junru Chen, Jingjing Wang, and Yafeng Li. Brant: Foun-  
421 dation Model for Intracranial Neural Signal. In *Thirty-Seventh Conference on Neural Information*  
422 *Processing Systems*, November 2023.
- 423 Yizi Zhang, Yanchen Wang, Donato Jiménez-Benetó, Zixuan Wang, Mehdi Azabou, Blake Richards,  
424 Renee Tung, Olivier Winter, Eva Dyer, Liam Paninski, et al. Towards a "universal translator" for  
425 neural dynamics at single-cell, single-spike resolution. *Advances in Neural Information Processing*  
426 *Systems*, 37:80495–80521, 2024a.
- 427 Yizi Zhang, Yanchen Wang, Donato Jimenez-Beneto, Zixuan Wang, Mehdi Azabou, Blake Richards,  
428 Olivier Winter, International Brain Laboratory, Eva Dyer, Liam Paninski, and Cole Hurwitz.  
429 Towards a "universal translator" for neural dynamics at single-cell, single-spike resolution, July  
430 2024b.
- 431 Wei-Long Zheng and Bao-Liang Lu. Investigating critical frequency bands and channels for eeg-  
432 based emotion recognition with deep neural networks. *IEEE Transactions on Autonomous Mental*  
433 *Development*, 7(3):162–175, 2015. doi: 10.1109/TAMD.2015.2431497.

## 434 **NeurIPS Paper Checklist**

### 435 **1. Claims**

436 Question: Do the main claims made in the abstract and introduction accurately reflect the  
437 paper's contributions and scope?

438 Answer: [\[Yes\]](#)

439 Justification: Yes, we outline the types insights that can be derived from our benchmark and  
440 then show preliminary neuroscience results that take steps to producing those insights in  
441 Section 4.

442 Guidelines:

- 443 • The answer NA means that the abstract and introduction do not include the claims  
444 made in the paper.
- 445 • The abstract and/or introduction should clearly state the claims made, including the  
446 contributions made in the paper and important assumptions and limitations. A No or  
447 NA answer to this question will not be perceived well by the reviewers.
- 448 • The claims made should match theoretical and experimental results, and reflect how  
449 much the results can be expected to generalize to other settings.
- 450 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
451 are not attained by the paper.

### 452 **2. Limitations**

453 Question: Does the paper discuss the limitations of the work performed by the authors?

454 Answer: [\[Yes\]](#)

455 Justification: Yes, we discuss this in the conclusion.

456 Guidelines:

- 457 • The answer NA means that the paper has no limitation while the answer No means that  
458 the paper has limitations, but those are not discussed in the paper.
- 459 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 460 • The paper should point out any strong assumptions and how robust the results are to  
461 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
462 model well-specification, asymptotic approximations only holding locally). The authors  
463 should reflect on how these assumptions might be violated in practice and what the  
464 implications would be.
- 465 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
466 only tested on a few datasets or with a few runs. In general, empirical results often  
467 depend on implicit assumptions, which should be articulated.
- 468 • The authors should reflect on the factors that influence the performance of the approach.  
469 For example, a facial recognition algorithm may perform poorly when image resolution  
470 is low or images are taken in low lighting. Or a speech-to-text system might not be  
471 used reliably to provide closed captions for online lectures because it fails to handle  
472 technical jargon.
- 473 • The authors should discuss the computational efficiency of the proposed algorithms  
474 and how they scale with dataset size.
- 475 • If applicable, the authors should discuss possible limitations of their approach to  
476 address problems of privacy and fairness.
- 477 • While the authors might fear that complete honesty about limitations might be used by  
478 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
479 limitations that aren't acknowledged in the paper. The authors should use their best  
480 judgment and recognize that individual actions in favor of transparency play an impor-  
481 tant role in developing norms that preserve the integrity of the community. Reviewers  
482 will be specifically instructed to not penalize honesty concerning limitations.

### 483 **3. Theory assumptions and proofs**

484 Question: For each theoretical result, does the paper provide the full set of assumptions and  
485 a complete (and correct) proof?

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

Answer: [NA]

Justification: We present a benchmark that only pertains to empirical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We release the code on github with a quickstart notebook as well as the scripts that produce all results and figures. The appendix contains specification of trials used in splits (see Appendix F).

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

540 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
541 tions to faithfully reproduce the main experimental results, as described in supplemental  
542 material?

543 Answer: [Yes]

544 Justification: Same as above. See item 4. We release our code on github and include a  
545 quickstart jupyter notebook as well as scripts to obtain our results.

546 Guidelines:

- 547 • The answer NA means that paper does not include experiments requiring code.
- 548 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
549 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 550 • While we encourage the release of code and data, we understand that this might not be  
551 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
552 including code, unless this is central to the contribution (e.g., for a new open-source  
553 benchmark).
- 554 • The instructions should contain the exact command and environment needed to run to  
555 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
556 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 557 • The authors should provide instructions on data access and preparation, including how  
558 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 559 • The authors should provide scripts to reproduce all experimental results for the new  
560 proposed method and baselines. If only a subset of experiments are reproducible, they  
561 should state which ones are omitted from the script and why.
- 562 • At submission time, to preserve anonymity, the authors should release anonymized  
563 versions (if applicable).
- 564 • Providing as much information as possible in supplemental material (appended to the  
565 paper) is recommended, but including URLs to data and code is permitted.

## 566 6. Experimental setting/details

567 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
568 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
569 results?

570 Answer: [Yes]

571 Justification: Hyperparameters are given in Appendix H and splits are specified in Ap-  
572 pendix F

573 Guidelines:

- 574 • The answer NA means that the paper does not include experiments.
- 575 • The experimental setting should be presented in the core of the paper to a level of detail  
576 that is necessary to appreciate the results and make sense of them.
- 577 • The full details can be provided either with the code, in appendix, or as supplemental  
578 material.

## 579 7. Experiment statistical significance

580 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
581 information about the statistical significance of the experiments?

582 Answer: [Yes]

583 For our empirical results, we report standard error across cross-val folds.

584 Guidelines:

- 585 • The answer NA means that the paper does not include experiments.
- 586 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
587 dence intervals, or statistical significance tests, at least for the experiments that support  
588 the main claims of the paper.
- 589 • The factors of variability that the error bars are capturing should be clearly stated (for  
590 example, train/test split, initialization, random drawing of some parameter, or overall  
591 run with given experimental conditions).

- 592 • The method for calculating the error bars should be explained (closed form formula,  
593 call to a library function, bootstrap, etc.)
- 594 • The assumptions made should be given (e.g., Normally distributed errors).
- 595 • It should be clear whether the error bar is the standard deviation or the standard error  
596 of the mean.
- 597 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
598 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
599 of Normality of errors is not verified.
- 600 • For asymmetric distributions, the authors should be careful not to show in tables or  
601 figures symmetric error bars that would yield results that are out of range (e.g. negative  
602 error rates).
- 603 • If error bars are reported in tables or plots, The authors should explain in the text how  
604 they were calculated and reference the corresponding figures or tables in the text.

## 605 8. Experiments compute resources

606 Question: For each experiment, does the paper provide sufficient information on the com-  
607 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
608 the experiments?

609 Answer: [Yes]

610 We discuss this in Appendix J.

611 Guidelines:

- 612 • The answer NA means that the paper does not include experiments.
- 613 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
614 or cloud provider, including relevant memory and storage.
- 615 • The paper should provide the amount of compute required for each of the individual  
616 experimental runs as well as estimate the total compute.
- 617 • The paper should disclose whether the full research project required more compute  
618 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
619 didn't make it into the paper).

## 620 9. Code of ethics

621 Question: Does the research conducted in the paper conform, in every respect, with the  
622 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

623 Answer: [Yes]

624 We adhere to the code of ethics.

625 Guidelines:

- 626 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 627 • If the authors answer No, they should explain the special circumstances that require a  
628 deviation from the Code of Ethics.
- 629 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
630 eration due to laws or regulations in their jurisdiction).

## 631 10. Broader impacts

632 Question: Does the paper discuss both potential positive societal impacts and negative  
633 societal impacts of the work performed?

634 Answer: [Yes]

635 We discuss this in the Conclusion.

636 Guidelines:

- 637 • The answer NA means that there is no societal impact of the work performed.
- 638 • If the authors answer NA or No, they should explain why their work has no societal  
639 impact or why the paper does not address societal impact.
- 640 • Examples of negative societal impacts include potential malicious or unintended uses  
641 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
642 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
643 groups), privacy considerations, and security considerations.

- 644
- 645
- 646
- 647
- 648
- 649
- 650
- 651
- 652
- 653
- 654
- 655
- 656
- 657
- 658
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 659 11. Safeguards

660 Question: Does the paper describe safeguards that have been put in place for responsible  
661 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
662 image generators, or scraped datasets)?

663 Answer: [NA]

664 Justification: We're using a public dataset only for evaluation purposes.

665 Guidelines:

- 666
- 667
- 668
- 669
- 670
- 671
- 672
- 673
- 674
- 675
- The answer NA means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 676 12. Licenses for existing assets

677 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
678 the paper, properly credited and are the license and terms of use explicitly mentioned and  
679 properly respected?

680 Answer: [Yes]

681 Justification: In the approach section we specify BrainTreebank's license.

682 Guidelines:

- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- The answer NA means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

696 • If this information is not available online, the authors are encouraged to reach out to  
697 the asset’s creators.

### 698 13. **New assets**

699 Question: Are new assets introduced in the paper well documented and is the documentation  
700 provided alongside the assets?

701 Answer: [Yes]

702 Justification: We are an evaluation-only benchmark. We make the code necessary for our  
703 benchmark public.

704 Guidelines:

- 705 • The answer NA means that the paper does not release new assets.
- 706 • Researchers should communicate the details of the dataset/code/model as part of their  
707 submissions via structured templates. This includes details about training, license,  
708 limitations, etc.
- 709 • The paper should discuss whether and how consent was obtained from people whose  
710 asset is used.
- 711 • At submission time, remember to anonymize your assets (if applicable). You can either  
712 create an anonymized URL or include an anonymized zip file.

### 713 14. **Crowdsourcing and research with human subjects**

714 Question: For crowdsourcing experiments and research with human subjects, does the paper  
715 include the full text of instructions given to participants and screenshots, if applicable, as  
716 well as details about compensation (if any)?

717 Answer: [NA]

718 Justification: We use a previously existing public dataset.

719 Guidelines:

- 720 • The answer NA means that the paper does not involve crowdsourcing nor research with  
721 human subjects.
- 722 • Including this information in the supplemental material is fine, but if the main contribu-  
723 tion of the paper involves human subjects, then as much detail as possible should be  
724 included in the main paper.
- 725 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
726 or other labor should be paid at least the minimum wage in the country of the data  
727 collector.

### 728 15. **Institutional review board (IRB) approvals or equivalent for research with human 729 subjects**

730 Question: Does the paper describe potential risks incurred by study participants, whether  
731 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
732 approvals (or an equivalent approval/review based on the requirements of your country or  
733 institution) were obtained?

734 Answer: [NA]

735 Justification: We use a public dataset that is openly published and available on the internet  
736 to construct our benchmark (BrainTreebank, <https://braintreebank.dev>). As such, we did not  
737 require any IRB approvals or equivalent to conduct our research.

738 Guidelines:

- 739 • The answer NA means that the paper does not involve crowdsourcing nor research with  
740 human subjects.
- 741 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
742 may be required for any human subjects research. If you obtained IRB approval, you  
743 should clearly state this in the paper.
- 744 • We recognize that the procedures for this may vary significantly between institutions  
745 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
746 guidelines for their institution.

747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763

- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

#### 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: We do not use LLMs as core components of our methods. One of our tasks is "GPT2 Surprisal", tasking the model with decoding the LLM negative log likelihood of the words in the dataset, however this feature was extracted from the sentences following standard protocol.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.

764 TODOs

- 765 • **TODO: chris, geeling** make an appendix with all the hyperparameters for PopT
- 766 • **TODO: bennet** write richer website description in appendix. Basically write up what will be  
767 displayed on the page. Put a new screenshot in.
- 768 • **DONE: chris** Put parcellation figure in appendix
- 769 • **DONE: chris** Put time series superposition figure in appendix
- 770 • **DONE: chris** Put time course for all features in appendix
- 771 • **DONE: chris** Make an appendix that has compute details of PopT
- 772 • **TODO: andrii** Make appendix H in tabular form.
- 773 • **TODO: andrii (only if you have time; low priority)** make a figure in the appendix for  
774 population level decoding over time.
- 775 • **TODO: chris / andrii** Fix table 2 to have corresponding info to the data.

#	Feature	Description	Benchmark Task
1	frame_brightness ( <i>visual</i> )	The mean brightness computed as the average HSV value over all pixels	Binary classification: low (percentiles 0%-25%) vs high (75%-100%)
2	global_flow ( <i>visual</i> )	A camera motion proxy. The maximal average dense optical flow vector magnitude	Same as above
3	local_flow ( <i>visual</i> )	A large displacement proxy. The maximal optical flow vector magnitude	Same as above
4	global_flow_angle ( <i>visual</i> )	As 2, averaged over orientation (degrees) and selected by maximal magnitude	2-way classification: Left vs Right (180 degree intervals)
5	local_flow_angle ( <i>visual</i> )	The orientation (degrees) of the largest local flow vector	Same as above
6	face_num ( <i>visual</i> )	The maximum number of faces per frame during the word	2-way classification: 0, or $\geq 1$
7	volume ( <i>auditory</i> )	Average root mean squared watts of the audio	Binary classification: low (0%-25%) vs high (75%-100%)
8	pitch ( <i>auditory</i> )	Average pitch of the audio	Same as above
9	delta_volume ( <i>auditory</i> )	The difference in average RMS of the 500ms windows pre- and post-word onset	Same as above
10	delta_pitch ( <i>auditory</i> )	The difference in average pitch of the 500ms windows pre- and post-word onset	Same as above
11	speech ( <i>language</i> )	Whether any speech is present in the given time interval	Binary classification
12	onset ( <i>language</i> )	Whether a new sentence starts in the interval, or there is no speech at all	Binary classification
13	gpt2_surprisal ( <i>language</i> )	Negative-log transformed GPT-2 word probability (given preceding 20s of language context)	Binary classification: low (0%-25%) vs high (75%-100%)
14	word_length ( <i>language</i> )	Word length (ms)	Same as above
15	word_gap ( <i>language</i> )	Difference between previous word offset and current word onset (ms)	Same as above
16	word_index ( <i>language</i> )	The word index in its context sentence	2-way classification: 0 (the first word in the sentence), or other (1)
17	word_head_pos ( <i>language</i> )	The relative position (left/right) of the word's dependency tree head	Binary classification
18	word_part_speech ( <i>language</i> )	The word Universal Part-of-Speech (UPOS) tag	2-way classification: verb (0), or other (1)
19	speaker ( <i>multimodal</i> )	The movie character that speaks the given word.	2-way classification: most frequent speaker (0), or other (1)

Table 1: **Extracted visual, auditory, and language features used to create the evaluations for Neuroprobe.** For all classification tasks, the classes were rebalanced. The difference between local and global flow is that global is the averaged optical flow, with the average being taken over all optical flow vectors on the screen, whereas local is the largest individual optical flow vector on the screen. The table is adapted from Chau et al. (2024).

777 **B Subject and movie information**

Subj.	Age (yrs.)	# Elec- trodes	Movie	Recording time (hrs)	Neuroprobe- Lite
1	19	154	Fantastic Mr. Fox	1.35	
			The Martian	2.43	x
			Thor: Ragnarok	1.77	x
2	12	162	Venom	1.54	x
			Spider-Man: Homecoming	2.05	
			Guardians of the Galaxy	1.90	
			Guardians of the Galaxy 2	2.13	x
			Avengers: Infinity War	2.30	
			Black Panther	1.42	
Aquaman	2.19				
3	18	134	Cars 2	1.64	x
			Lord of the Rings 1	2.25	x
			Lord of the Rings 2 (extended edition)	3.58	
4	12	188	Shrek 3	1.38	x
			Megamind	1.44	x
			Incredibles	0.85	
5	6	156	Fantastic Mr. Fox	1.35	
6	9	164	Megamind	0.68	
			Toy Story	1.29	
			Coraline	0.84	
7	11	246	Cars 2	1.64	x
			Megamind	1.44	x
8	4.5	162	Sesame Street Episode	0.94	
9	16	106	Ant Man	1.80	
10	12	216	Cars 2	1.33	x
			Spider-Man: Far from Home	1.93	x

Table 2: **Subject statistics** Subjects in the BrainTreebank dataset, and the trials used in the benchmark tasks. Table adapted from Wang et al. (2023). The second column shows the total number of electrodes. The average amount of recording data per subject is 4.3 (hrs).

Subj.	Age	Sex	Movies	Time (h)	# Sent.	# Words	# Lemmas	# Elec.	# Probes
1	19	M	7, 18, 19	5.6	4372	27424	4489	154	13
2	12	M	2, 3, 4, 8, 9, 17, 21	13.5	9870	57731	9164	162	47
3	18	F	5, 11, 12	7.5	5281	31596	4547	134	12
4	12	F	10, 13, 15	3.7	4056	23876	4017	188	15
5	6	M	7	1.35	1282	7908	1481	156	12
6	9	F	6, 13, 20	2.8	3789	20089	3349	164	12
7	11	F	5, 13	3.08	3523	19068	2828	246	18
8	4	M	14	0.94	860	3994	537	162	13
9	16	F	1	1.80	1558	9235	1480	106	12
10	12	M	5, 16	3.08	3981	22147	3004	216	17

Table 3: **All subjects language, electrodes and personal statistics.** Columns from left to right are the subject’s ID and information (age and gender), the IDs of the movies they watched (corresponding to Supplementary Table 4), the cumulative movie time (hours), number of sentences, number of words (tokens) and number of unique lemmas (canonical word forms), as well as the number of probes the subject had and their corresponding number of electrodes. Table adapted from Wang et al. (2024).

# Movie	Year	Length	Sent.	Words	Unique words	Nouns	Unique nouns	Verbs	Unique verbs
1 Antman	2015	7027	1558	9869	1944	1358	705	1545	580
2 Aquaman	2018	8601	1054	7233	1544	1069	520	1104	508
3 Avengers: Infinity War	2018	8961	1523	8529	1750	1083	607	1317	495
4 Black Panther	2018	8073	1254	7580	1606	1093	553	1209	508
5 Cars 2	2011	6377	2051	11407	2037	1572	724	1664	577
6 Coraline	2009	6036	997	5433	1232	784	409	805	348
7 Fantastic Mr. Fox	2009	5205	1282	8461	1864	1229	681	1227	484
8 Guardians of the Galaxy 1	2014	7251	1174	8295	1779	1096	603	1250	529
9 Guardians of the Galaxy 2	2017	8146	1290	9405	1824	1224	626	1370	532
10 Incredibles	2003	6926	1521	9430	1954	1226	652	1557	591
11 Lord of the Rings 1	2001	13699	1514	10566	1998	1473	679	1487	598
12 Lord of the Rings 2	2002	14131	1716	11041	2065	1588	743	1619	646
13 Megamind	2010	5735	1472	8891	1726	1172	602	1347	496
14 Sesame Street Ep. 3990	2016	3440	860	4220	787	717	231	706	217
15 Shrek the Third	2007	5568	1063	7226	1590	977	568	1071	422
16 Spiderman: Far From Home	2019	7764	1930	12189	1969	1459	668	1785	560
17 Spiderman: Homecoming	2017	8008	2196	12295	2066	1583	777	1808	572
18 The Martian	2015	9081	1570	11374	2192	1757	812	1677	622
19 Thor: Ragnarok	2017	7831	1583	9683	1789	1195	599	1419	548
20 Toy Story 1	1995	4863	1320	7216	1510	1019	548	1027	395
21 Venom	2018	6727	1379	7937	1513	897	507	1217	433

Table 4: **Language statistics for all movies.** Columns from left to right are the movie’s ID, name, year of production, length (seconds), number of sentences, number of words (tokens), number of unique words (types), number of nouns, number of unique nouns, number of verbs and number of unique verbs. Table adapted from Wang et al. (2024).

778 **C Composition of movies by volume**

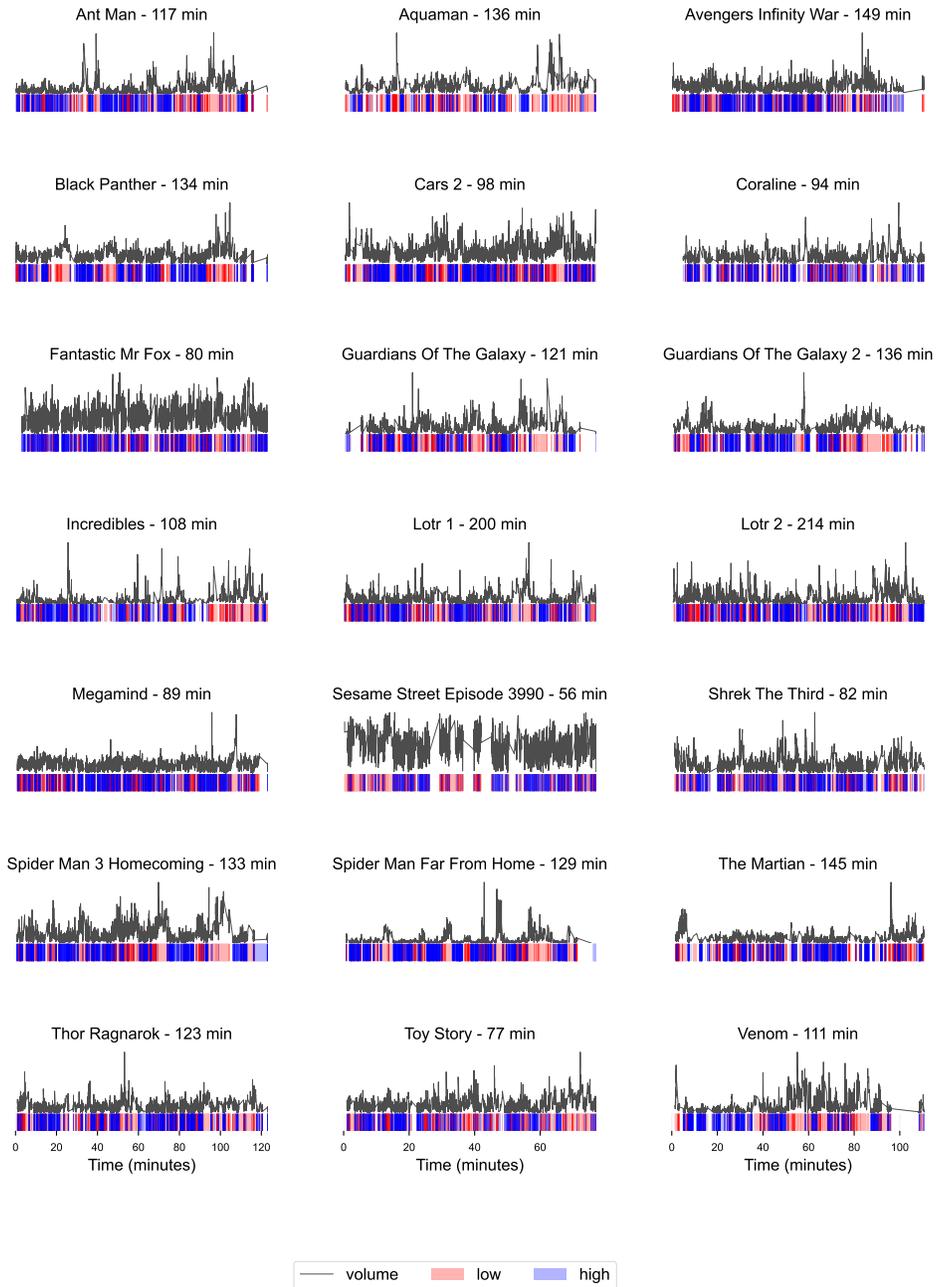


Figure 9: **Volume comparison across movies.** The black line shows the normalized audio volume over time for 18 feature-length films and one TV episode shown to subjects. Below each volume trace, colored bars indicate periods of relatively low (red) and high (blue) volume, defined as the bottom 25% and top 25% of volume values respectively.

779 **D Speech localization**

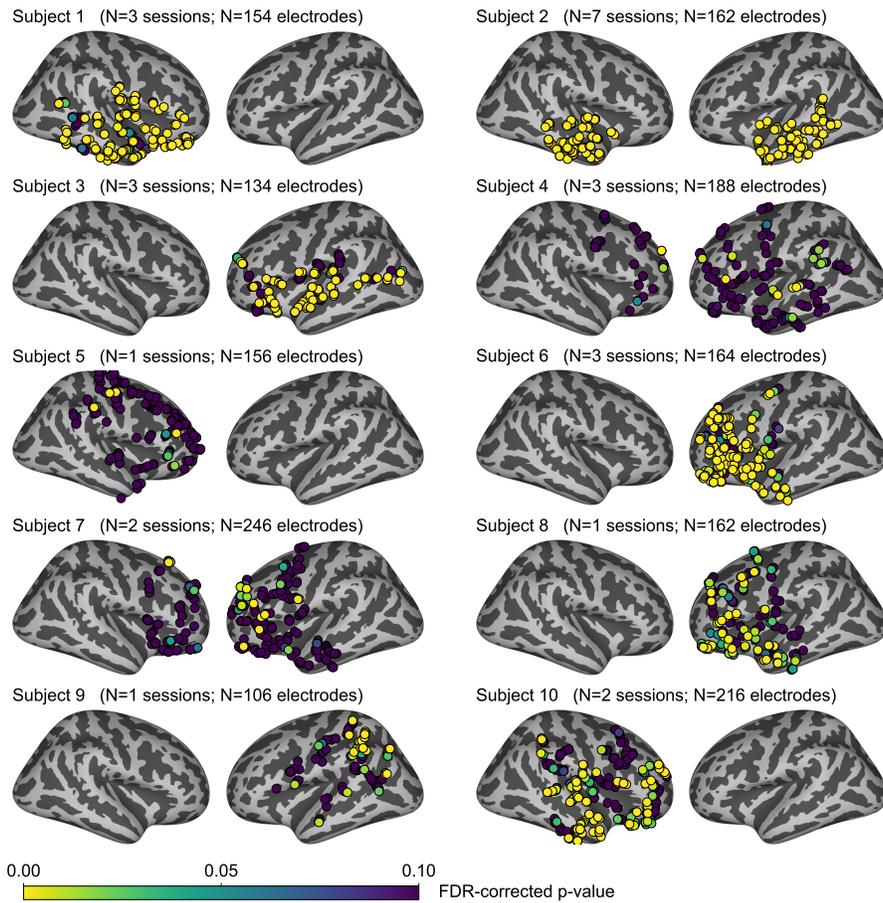


Figure 10: **Electrode locations and speech selectivity across subjects.** Brain reconstructions showing electrode placement and speech-selective responses for all 10 subjects. Each dot represents an electrode, colored by its FDR-corrected p-value from a speech vs. non-speech classification (color scale above, yellow indicating stronger selectivity). Left and right hemispheres are shown separately, with session counts and total electrodes noted. Speech selectivity was assessed by comparing high gamma power (70–300 Hz, dB) during the first 125 ms after word onset to non-speech intervals of equal duration. A two-sample t-test determined significance, with Benjamini-Hochberg correction applied for multiple comparisons.

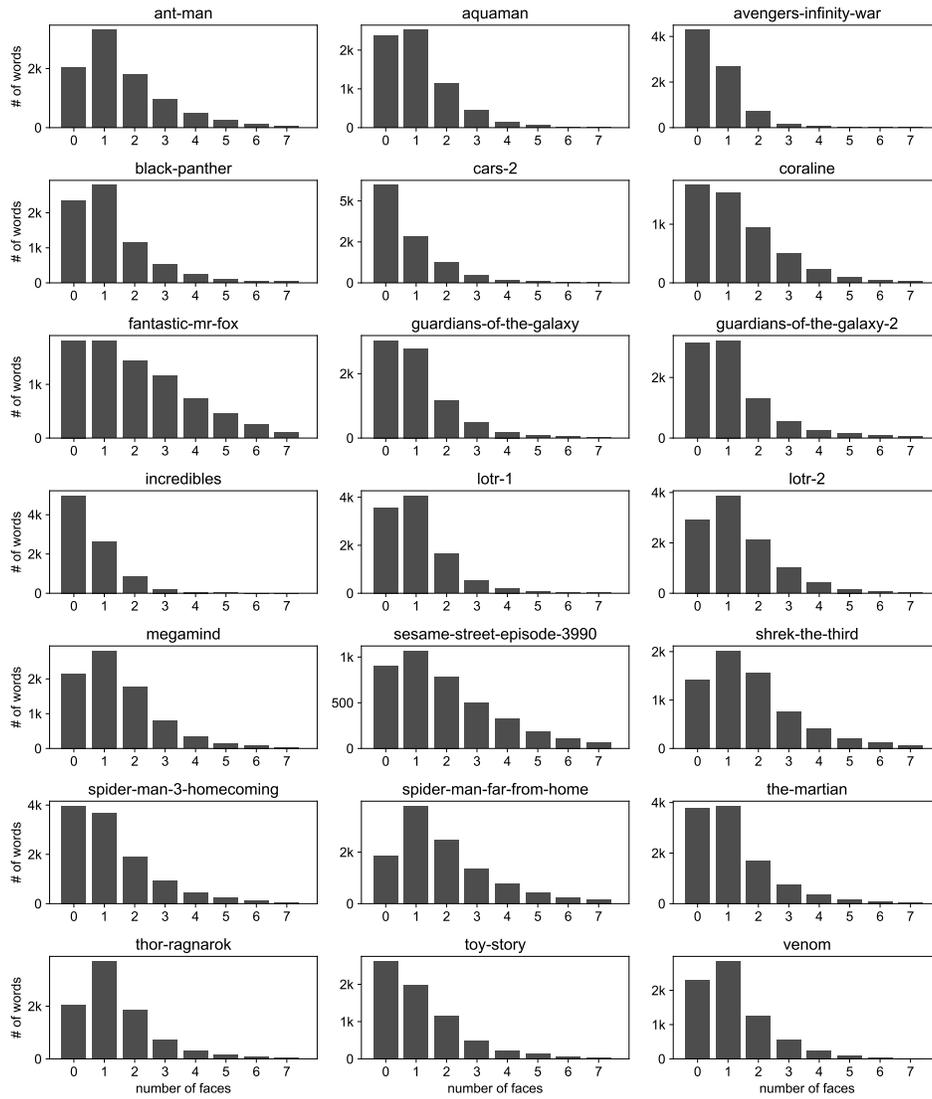


Figure 11: **Distribution of faces detected per frame across different movies.** Histograms show the number of words (y-axis) that occur during frames containing different numbers of faces (x-axis) for 18 feature-length films and one TV episode (Sesame Street)

781 **F Splits**

782 Neuroprobe includes 3 different types of splits.

783 **Same subject/Same trial**

784 **Same subject/Different movie** This is a slightly more difficult split. It ensures completely that no  
785 data-contamination due to auto-correlation has occurred.

786 **Different subject/Different movie** This is the most difficult split. It tests the model's ability to  
787 generalize between subjects *and* stimuli.

788 **TODO: describe what the splits are and which trials are in each split.**

## 789 G Neuroprobe-lite

790 The following subject-trial pairs are included in Neuroprobe Lite:

- 791 • Subject 1: Trials 1, 2
- 792 • Subject 2: Trials 0, 4
- 793 • Subject 3: Trials 0, 1
- 794 • Subject 4: Trials 0, 1
- 795 • Subject 7: Trials 0, 1
- 796 • Subject 10: Trials 0, 1

797 For every task, the number of datapoints was trimmed at 3500 datapoints (i.e. if a specific movie has  
798 more than 3500 annotations for any task, only the first 3500 are taken for the Lite benchmark). When  
799 selecting the subject/trial pairs for Neuroprobe Lite, we selected the trials that contained the most  
800 tasks which hit the 3500 datapoints limit.

## 801 H Models benchmarked

802 **Linear** For this evaluation, raw voltage traces sampled at 2048 Hz were taken from the BrainTree-  
803 bank data, then line noise was removed at  $60 \pm 5$  Hz and the 4 harmonics, and the resulting vectors of  
804 sampled features were fed as input to the linear regression. We found almost identical results when  
805 removing line noise or passing the data raw to the linear regression.

806 **Linear (STFT)** For this baseline evaluation, the features are the STFT of the raw signal with the  
807 following parameters (given that the sampling rate is 2048Hz):

- 808 • nperseg=256
- 809 • noverlap=0
- 810 • window=boxcar

811 After this step, the data turns into an array of arrays where first dimension is the time bin and the  
812 second dimension is the STFT result (a complex number); for the downstream regression, all of these  
813 features are concatenated together, with the real and imaginary parts of the complex features being  
814 split into two features each.

815 **Linear (spectrogram)** For this baseline evaluation, first the STFT of the raw voltage signal was  
816 taken as in the Linear (STFT) description, and then the absolute value of each complex number was  
817 taken to obtain the final real number features for each example.

818 **BrainBERT** For this evaluation, the BrainTreebank data was Laplacian rereferenced (as described  
819 in the original BrainBERT paper by Wang et al. (2023)), with line noise removed, and then passed into  
820 the BrainBERT model as provided by Wang et al. (2023). The output features were concatenated and  
821 used as input to the linear regression. For the electrodes which could not be Laplacian rereferenced,  
822 non-rereferenced data was inputted into BrainBERT. The BrainBERT model was frozen and only the  
823 final linear regression layer was fine tuned, in order to compare the quality of features generated by  
824 the foundation model.

825 For all linear regression, we used the sklearn package, class LinearRegression, with the tolerance  
826 parameter set as 0.001. In all cases, the features were first normalized using the sklearn StandardScaler.  
827 We found that it helps with convergence and often produces higher regression values for the baselines.

828 **PopulationTransformer Off-the-shelf** Population Transformer (PopT) is a SSL pretrained model for  
829 encoding arbitrary ensembles of iEEG electrode data for general downstream decoding (Chau et al.,  
830 2024). The model consists of a transformer backbone that learns functional and spatial relationships  
831 between input channels whose temporal activity is encoded. We use the publicly available weights  
832 which were pretrained on data from 10 iEEG subjects, using 5s BrainBERT temporal embeddings  
833 from individual channels. For Population Transformer, we followed the implementation and used

834 the weights from (Chau et al., 2024). The fine-tuning protocol is taken to be directly the same as in  
835 the authors’ original paper (including linear rate, number of epochs, a factor of 10 between learning  
836 rates of the linear output layer vs the transformer blocks, etc), but reduce the number of steps to  
837  $steps = 1000$ . We finetune Population Transformer in two conditions: either by only finetuning the  
838 final linear output layer while keeping the rest of the model weights frozen (the “frozen” condition),  
839 or finetuning through the whole model (the default PopT condition).

## 840 **I Benchmark results**

841 **TODO fill in with tabular form of fig. 8.**

## 842 **J Compute requirements**

843 Every Linear regression was run on a CPU-only instance, with 2 virtual CPU cores and 64GB RAM  
844 for the population level results and 2 CPU cores with 6GB RAM for the single electrode decoding  
845 results. For BrainBERT, the necessary resources also included a GPU with at least 9GB of memory  
846 along with 128GB of RAM and 2 CPU cores. For the PopulationTransformer, the fine-tuning was  
847 done on 2 GPUs (NVIDIA GeForce GTX TITAN X) with at least 12GB of GPU RAM.

848 **K Leaderboard**

TODO: describe leaderboard website

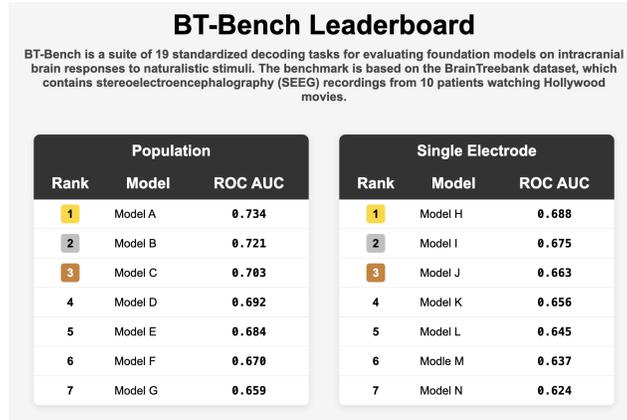


Figure 12: **The leaderboard for the task of classifying sentence onset.** The public webpage link will be made available upon publication. **TODO: revisit caption**

849

850 **L Time course of task decodability**

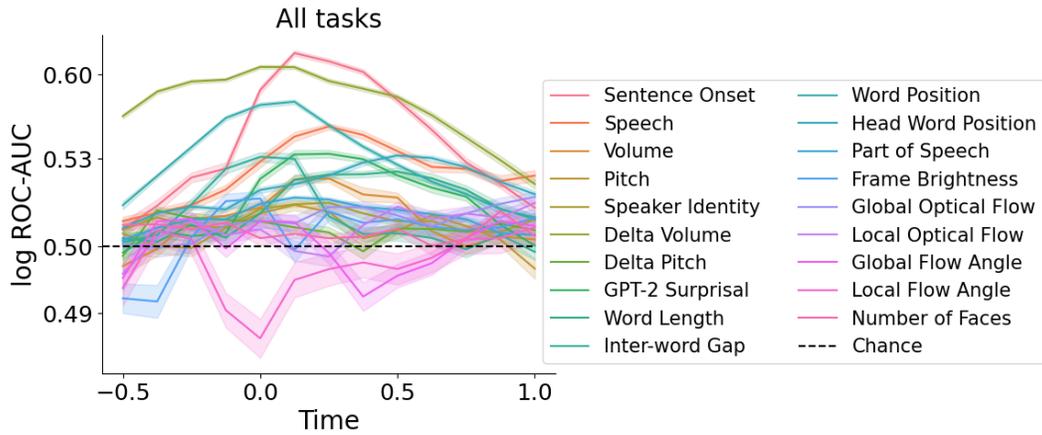


Figure 13: **TODO: revisit caption** All the plots from Figure 4 overlaid. Error bars show standard error from variability across all electrodes (from all subjects and all sessions).

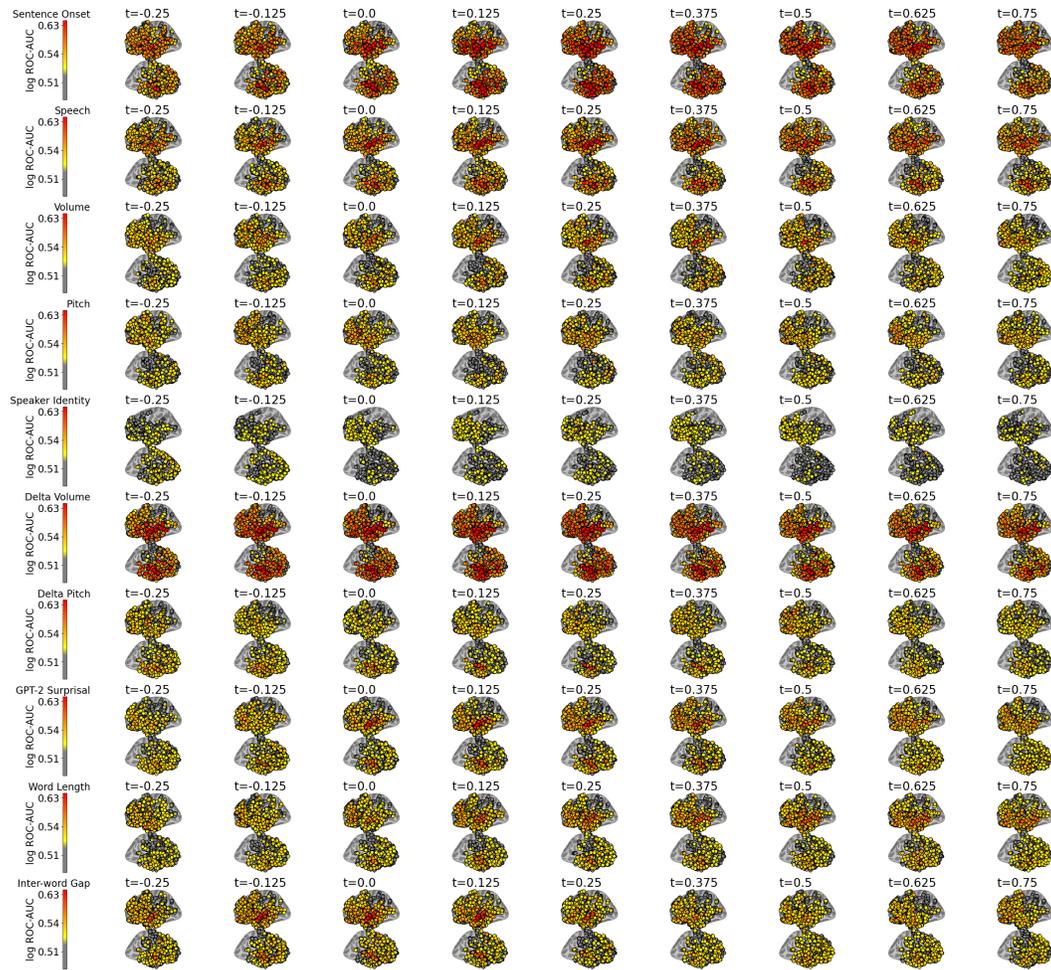


Figure 14: **TODO: revisit caption** Same as Figure 7 but for all features. Pt 1

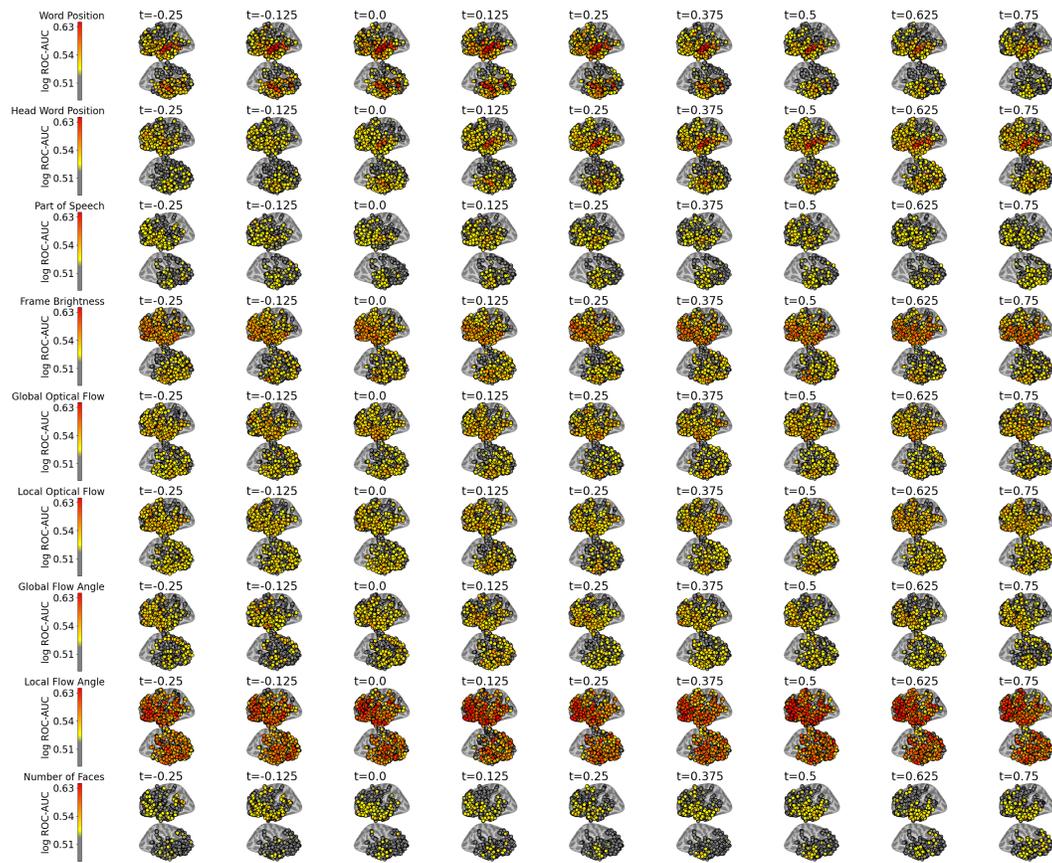


Figure 15: **TODO: revisit caption** Same as Figure 7 but for all features. Pt 2

851 **M Region analysis**

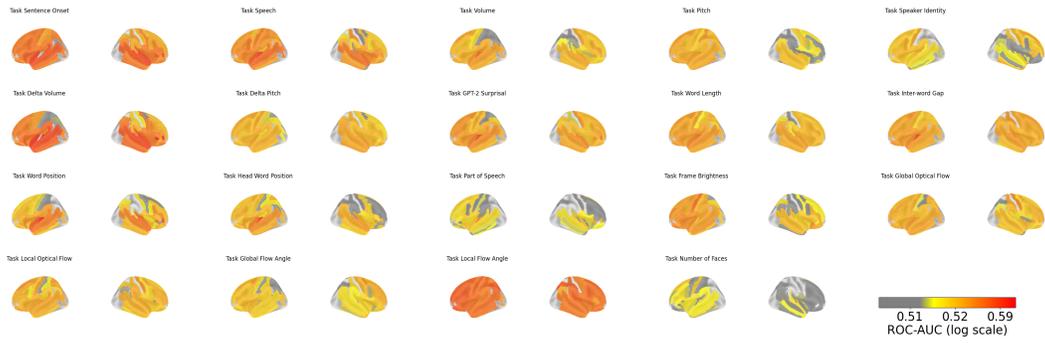


Figure 16: **TODO: revisit caption** top 10-th percentile of electrodes in each region are plotted  
**Make it top k=100?**