
Learning Representations of Human Neural Activity via Contrastive Neural Forecasting

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Electrical patterns of brain activity form the basis of everything from perception
2 and movement to complex behaviors like decision-making and conscious thought.
3 While terabytes of human intracranial electroencephalography (iEEG) recordings
4 are openly available, deciphering and productively using them for downstream
5 use cases remains a challenging problem. We present Contrastive Neural Fore-
6 casting (CNF), a simple self-supervised framework for learning representations
7 of population-level neural activity across electrodes, time, and individuals from
8 unlabeled data at scale. The CNF objective requires the model to predict future
9 neural states in its latent space using cross-entropy over batched data samples.
10 Our objective is designed to resolve two major challenges in traditional MSE-
11 based autoencoding approaches. Forecasting in the latent space relieves the model
12 from overfitting to noise that is inherent in the data, and the cross-entropy loss
13 enables flexible capturing of high-dimensional, multimodal distributions under-
14 lying the evolution of neural dynamics. We validate the superior performance of
15 the contrastive objective on BrainBERT, and then train and open-source CNF-1,
16 a foundation model for human iEEG. We pretrain end-to-end directly from raw
17 voltage traces, without relying on handcrafted features or frequency band filtering.
18 While still closely followed by the linear baselines, which we found in many cases
19 score higher than other pretrained models, CNF-1 achieves state-of-the-art perfor-
20 mance on a suite of downstream decoding tasks. Surprisingly, and challenging
21 assumptions made in prior work, we obtain better performance by omitting the
22 spatial location of the electrodes from the embeddings, instead allowing the model
23 to learn its own channel-specific parameters. We show how CNF-1 can enable
24 novel approaches to extract neuroscientific insight from unlabeled data at scale. We
25 envision future clinical applications such as real-time functional region mapping
26 and model-guided electrical stimulation interventions in the operating room, as
27 well as next-generation brain-computer interfaces. Taken together, our work paves
28 the way for scalable brain foundation models trained entirely from observational
29 data.

30 **1 Introduction**

31 The human brain continuously processes rich, overlapping streams of information: from interpreting
32 speech and recognizing objects to reasoning about complex events [Schurz et al., 2014]. Despite
33 considerable progress in neuroscience over the past decades, building a comprehensive computational
34 model of the brain, where the brain state can be decoded, simulated, and interfaced with seamlessly,
35 remains a formidable challenge [Sejnowski et al., 2014]. Through invasive and non-invasive stimu-
36 lation interventions, these models could enable personalized treatments for neurological disorders
37 such as epilepsy [Herron et al., 2024, Morrell, 2011], and through superior decoding and encoding

38 capabilities enhance communication between brains and machines as well as between brains and
39 other brains [Pereira et al., 2018].

40 Foundation models have transformed fields like natural language processing and computer vision,
41 offering unparalleled performance across tasks and datasets [Bommasani et al., 2021, Brown et al.,
42 2020]. Yet, their potential in neuroscience remains untapped. Even if brain foundation models
43 (such as transformers pretrained on large volumes of data) do not provide simple and interpretable
44 models of brain function, their capacity for capturing complex and high-dimensional relationships
45 across brain areas, time, and individuals positions them as powerful tools for advancing medicine and
46 neuroscience [Parvizi and Kastner, 2012].

47 Human intracranial encephalography (iEEG) offers brain interfacing at an unprecedented combination
48 of spatial and temporal resolution. However, these raw voltage time series are noisy, high-dimensional,
49 highly nonlinear, and non-trivially dependent on physiological variables [Noury et al., 2016, Buzsáki
50 et al., 2012], which has been a major obstacle to gaining useful insight using this data and to creating
51 new tools and treatments relying on iEEG.

52 Our work introduces *Contrastive Neural Forecasting (CNF)*, a simple self-supervised framework
53 for learning population-level representations of human brain dynamics. We propose path toward
54 general-purpose brain models that excels by learning directly from raw voltage traces, without
55 requiring knowledge of electrode locations or prior domain-specific assumptions.

56 **Contributions** The key contributions of this work are:

- 57 • We propose a contrastive predictive learning objective tailored for neural time series data
58 that forecasts future neural states in latent space using a cross-entropy loss over real samples,
59 enabling flexible modeling of high-dimensional, multimodal distributions underlying neural
60 dynamics. We validate this objective by showing its superior performance when pretraining
61 BrainBERT [Wang et al., 2023].
- 62 • Our proposed objective enables the unification of representation learning from single channel
63 voltage traces with population activity over the whole brain by combining information from
64 many electrodes, without the need for pretrained channel feature extractors or spectrogram
65 encoders, handcrafted features, or frequency-based preprocessing.
- 66 • We propose the use of learned electrode embeddings for modeling of iEEG data. Our
67 method learns without access to spatial information about electrodes, recovering it during
68 pretraining in a purely data driven way, challenging the assumption that spatial embeddings
69 are necessary for accurate neural modeling.
- 70 • We introduce **CNF-1**, a foundation model trained end-to-end on raw iEEG voltage traces
71 across individuals and electrodes, achieving state-of-the-art performance on multiple decoding
72 benchmarks. We release the pretrained CNF-1 model and codebase to promote further
73 research on scalable, general-purpose brain foundation models (upon publication).

74 1.1 Related work

75 **Foundation Models for Neural Data.** Neuroformer introduced a multimodal, multitask generative
76 pretrained transformer tailored for systems neuroscience, capable of associating behavioral and neural
77 representations through joint training [Antoniades et al., 2024]. BrainBERT [Wang et al., 2023] learns
78 representations from single channels of intracranial EEG in a self-supervised manner by predicting
79 masked out spectrograms. Our work The Population Transformer (PopT, Chau et al. [2024]) and
80 [Zhang et al., 2023] extended pretrained embeddings from BrainBERT to enable decoding on the
81 population level. We provide two main improvements that raise the performance of the model: the
82 contrastive forecasting objective as opposed to the naïve MSE approach and learnable electrode
83 embeddings instead of positional coordinate embeddings. Learnable embeddings were introduced
84 by [Azabou et al., 2023] for the single unit modeling. We use the same learned embeddings and
85 adapt them for the continuous iEEG signal. NDT2 emphasized large-scale spatiotemporal pretraining
86 for neural spiking activity, facilitating adaptation to novel contexts in decoding tasks [Ye et al.,
87 2023]. Foundation models for neural data have been developed for single unit activity and fMRI [Liu
88 et al., 2022, Cai et al., 2023, Dong et al., 2024], however our work focuses specifically on human
89 intracranial EEG. For a review of brain foundation models, see Zhou et al. [2025].

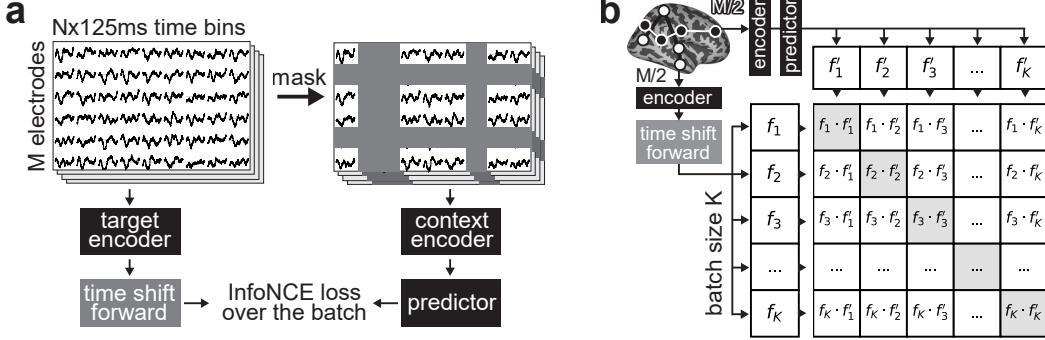


Figure 1: **Overview of the Contrastive Neural Forecasting approach.** (a) The neural activity timeseries is binned and split into two parallel streams: target and context. Each target timebin is encoded separately, and the whole timeseries is shifted forward by one timebin. The context stream is randomly masked in both the electrode and the time dimensions, and passed into the context encoder and then into the predictor. The target and context stream are compared using the InfoNCE objective, (b). The features generated by the model are compared with the corresponding electrode and timebin's features of every other item in the batch, and the InfoNCE objective requires that the corresponding pairs from the same item in the batch are close together in the feature space, and far away from the features of other items in the batch. All of the model components are trained end-to-end.

90 **Contrastive representation learning.** Contrastive Predictive Coding was first introduced by [van den
91 Oord et al., 2019] for representation learning, and successfully applied at scale by CLIP for language-
92 image pretraining [Radford et al., 2021]. We adapt these concepts from the ML literature to our
93 neural data setting. STNDT [Le and Shlizerman, 2022], in a single unit spiking modeling setting,
94 used contrastive learning as an auxiliary loss to further augment the data and constrain the model.
95 Similarly, [Vishnubhotla et al., 2023] use contrastive learning to learn representations for spike sorting
96 of single units. We too use contrastive learning, but directly for building models of continuous iEEG
97 signal.

98 2 The Contrastive Neural Forecasting approach

99 In this section, we overview the main components of our approach (Figure 1): predicting the future
100 latent representation of the signal and the contrastive forecasting objective based on the cross-entropy
101 loss. We assume that the neural data originates from a set of M channels (e.g., electrodes in the
102 brain), sampled at a constant rate to produce T data segments of length τ each, where τ is the desired
103 binning size for predictive modeling. Let's denote $x_t^{(m)}$ as the data sample from channel $m \leq M$
104 at time $t \leq T$. In the self-supervised learning setting, one is interested in modeling the distribution
105 of masked data conditioned on unmasked data. As an illustrative example, let's say we want to
106 autoregressively predict the joint distribution of the signal at time $t + 1$ using the previous T time
107 steps: $p(x_{t+1}^1 \dots x_{t+1}^M | x_1^1 \dots x_t^M)$. One approach, prevalent in the literature [Wang et al., 2023, Zhang
108 et al., 2023], is to define a parametrized predictive function F (i.e. a neural network), and train
109 the parameters to minimize the mean squared error (MSE) between the predicted and true masked
110 datapoints:

$$\mathcal{L}_{MSE} = \sum_i \|\hat{x}_{t+1}^i - x_{t+1}^i\|^2, \quad (1)$$

111 where $\hat{x}_{t+1}^i = F(i, x_1^1 \dots x_t^M)$ denotes the prediction of the model.

112 Despite the popularity of this approach, it has two flaws. First, in settings with an inherently high
113 level of noise, typical for recordings from intracranial electrodes, the MSE objective punishes the
114 model for poorly fitting the noise, encouraging overfitting to the noise pattern of the training dataset.
115 Further, these signals tend to be temporally autocorrelated.

116 The second flaw is in the implicit assumption that underlies the choice of MSE as the objective:
117 that the distribution of masked timepoints can be effectively captured with a unimodal Gaussian

118 centered around the mean which is equal to the prediction of the model (for an overview of this
119 equivalence, see Bishop [2006]). This assumption is not justified in our setting of interest. In practice,
120 the dynamical system is partially observed (the dimensionality of the signal, $M < 300$, is negligible
121 compared to the roughly 80 billion neurons in the brain), meaning that the observed input data put
122 very mild constraints on the multimodal distribution of the future evolution of the observed data.

123 To overcome these challenges, we introduce Contrastive Neural Forecasting. In CNF, the input data
124 is encoded with the context encoder $E_{context}$ and passed into a predictor P , and then compared to
125 the encoding of target data by a target encoder E_{target} . Specifically, we use the InfoNCE objective,
126 which pushes the predicted embedding $P(E_{context}(x_1^1 \dots x_t^M))$ to be close to the embedding of the
127 real target $E_{target}(x_{t+1}^1 \dots x_{t+1}^M)$ and far from the embeddings of other random timesamples in the
128 dataset. Formally, given a set of N random negative samples with timepoints t'_1, \dots, t'_N , the InfoNCE
129 loss is defined as:

$$\mathcal{L}_{InfoNCE} = -\log \frac{\exp(P(E_{context}(x_1^1 \dots x_t^M)) \cdot E_{target}(x_{t+1}^1 \dots x_{t+1}^M) / \tau)}{\sum_{j=1}^N \exp(P(E_{context}(x_1^1 \dots x_t^M)) \cdot E_{target}(x_{t'_j}^1 \dots x_{t'_j}^M) / \tau)}, \quad (2)$$

130 where \cdot denotes cosine similarity (dot product over the normalized features), and τ is a temperature
131 hyperparameter. In practice, this objective is efficiently implemented as the cross-entropy loss over
132 the batch dimension, meaning that negative samples for every item in the batch are taken from the
133 other items encountered in the same batch.

134 This formulation has three advantages. First, it automatically ensures that there is no incentive for the
135 model to encode noise in its latent space, where noise is defined as anything that is not helpful to
136 disambiguate the true future neural state from other random examples of neural states. Second, it
137 doesn't place assumptions (such as unimodality) on the distribution of the future timestep. Finally, it
138 turns the hard problem of modeling the high-dimensional, continuous distribution of the neural signal
139 into the "easy" problem of multi-class classification using the cross-entropy loss.

140 In the next sections, we describe the experiment setup and results that demonstrate these advantages.

141 3 Experimental Setup

142 **Data** To train and evaluate the performance of our objective, we use the publicly available Brain-
143 Treebank dataset [Wang et al., 2024]. The dataset consists of 43 hours of intracranial SEEG recordings
144 from 10 human subjects (ages 4–19) implanted with a total of 1,688 electrodes while passively watching
145 26 full-length Hollywood films. It includes aligned audio-visual and language annotations for
146 over 223,000 words across nearly 39,000 sentences, offering high temporal and spatial resolution
147 data suitable for multimodal neural decoding and large-scale modeling.

148 **Decoding evaluation tasks** We evaluate models on a suite of 14 standardized neural decoding
149 tasks spanning vision, audio, language, and multimodal domains, derived from the annotations in
150 the BrainTreebank dataset, such as audio volume, optical flow direction, face count, word onset,
151 LLM surprisal score, part-of-speech, speaker identity, etc. All of the tasks are formalized as binary
152 classification by thresholding the annotations. The models are tasked with classifying the task labels
153 from voltage traces of length 1 second aligned to each word onset. This decoding benchmark contains
154 labeled neural data from 12 recording sessions across 6 individuals. We evaluate the models by
155 fine-tuning on each task's training split, and testing on the non-intersecting test split that was taken
156 from a different recording session. For more details about the decoding tasks, see Appendix A.

157 **Models** We bin the neural data sampled at 2048 Hz into bins of 256 samples each (125 ms). The
158 target encoder E_{target} is a simple linear layer from the raw 256-dimensional feature vector into
159 the $d_{model} = 192$ dimensional latent space. For the initial validation experiments, we reimplement
160 a context encoder scaled-down version of the BrainBERT architecture [Wang et al., 2023] for
161 computational efficiency, which we call BrainBERT-mini. BrainBERT-mini is a transformer encoder
162 stack [Vaswani et al., 2023] with $N = 4$ layers and the hidden dimension size 192 with 12 attention
163 heads per layer. For CNF-1, the context encoder (Figure 2a) is a transformer with 4 layers, hidden
164 dimension d_{model} , that takes as input tokens which are 16 consecutive samples of the input data and
165 produces the latent representations of dimensionality d_{model} , and 4 attention heads per layer. The
166 outputs of the context encoder are concatenated to produce chunks of 256 samples for the next model

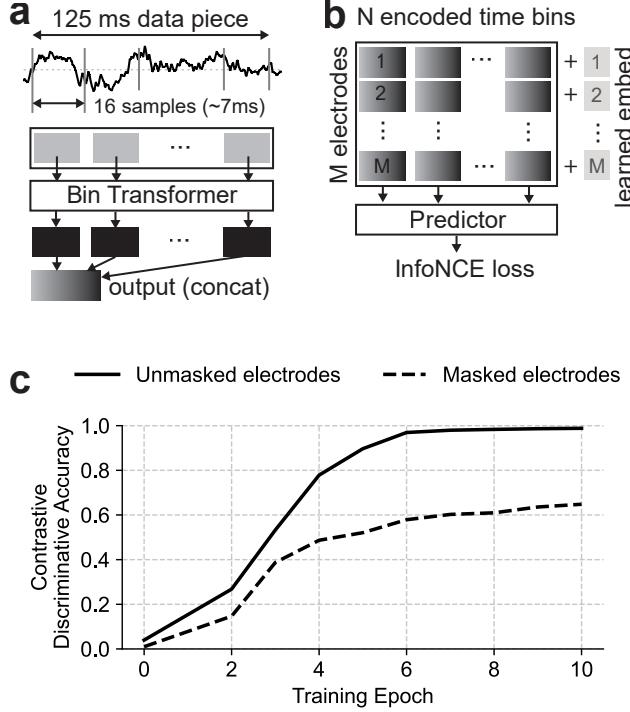


Figure 2: **Architecture and training dynamics of CNF-1.** Neural data sampled at 2048 Hz is binned into 125 ms segments (256 samples). The target encoder is a linear projection from 256 to a 192-dimensional latent space. The context encoder (a) is a 5-layer Transformer with 12 attention heads and hidden size 192, operating on tokens of 16 consecutive samples. The predictor (b) takes per-channel, per-timebin representations, adds learned electrode embeddings, and processes them through another 5-layer Transformer. Its outputs are compared with future target embeddings using the InfoNCE loss.

stage. The predictor (Figure 2b) takes in these representations for each channel $i \leq M$ and timebin $t \leq T = 1$ second, which are added to the learned per-channel electrode embeddings K_i (see next subsection), and again outputs the features of dimensionality d_{model} for every token by passing them through 5 layers of 12 attention heads each.

Pretraining For BrainBERT-mini experiments, we follow the approach of [Wang et al., 2023]. For the masking scheme, we set all data in $p = 10\%$ of timebins to 0, only passing the masked timebins into the objective (either MSE or our InfoNCE contrastive objective for this experiment). For CNF-1, the representations from its predictor output layer are then compared with the corresponding target embeddings *for the following future timebin*, using the InfoNCE objective as described in the previous section. For all models, we use a context of neural data of length 2 seconds. BrainBERT-mini is trained for 100 epochs on a small dataset containing just one subject’s session data. CNF-1 is trained for 10 epochs (CNF-1) on data from 20 sessions from all 10 subjects in the dataset. We train all networks with learning rate 0.003, the Muon optimizer [Jordan et al., 2024], and learning rate scheduling of 100 steps of warmup followed by linear decay to 0. The networks are trained on a single A100 GPU for 10 hours (CNF-1) or 2 hours (BrainBERT-mini).

Electrode embeddings To provide the Predictor transformer with the information about which channel each signal comes from in the brain, we allow the model to learn additional separate vectors of dimensionality d_{model} for each channel $1 \leq i \leq M$ in the dataset, implementing the technique used by [Azabou et al., 2023] for the single unit modeling. We contrast this approach with the prior iEEG-based approaches [Zhang et al., 2023, Chau et al., 2024] that provide coordinates of each channel via cosine positional embeddings. Empirically, we find higher performance with fully

188 learned embeddings, when not providing any spatial information of the electrodes into the model (see
189 Appendix B).

190 **Baselines and previous methods** We compare the performance of our model to six different
191 baselines and previous methods for feature learning on human intracranial EEG data:

- 192 • Linear regression from the raw voltage segments, aligned to the word onset.
- 193 • Linear regression from the spectrogram of the signal, normalized per frequency bin.
- 194 • Linear regression from the Fourier transform features (which include both magnitude and
195 phase information of the frequency bands).
- 196 • Population Transformer Chau et al. [2024], a previous state of the art in representation
197 learning from human intracranial EEG on the BrainTreebank dataset. We compare against
198 frozen PopT (only fine tuning the output linear layer and keeping the model weights frozen),
199 and a end-to-end finetuned PopT for each task.
- 200 • BrainBERT [Wang et al., 2023], a single-electrode representation extractor from iEEG on
201 the same BrainTreebank dataset. For this evaluation, the features from every electrode are
202 concatenated together before passing them into the linear regression layer to obtain the final
203 prediction of the task label.

204 For more detail on baselines and previous methods, see Appendix C.

205 4 Results

Training objective	Mean decoding AUROC (14 tasks)
MSE loss (voltage)	0.638 ± 0.008
MSE loss (spectrogram)	0.598 ± 0.009
Contrastive (voltage, latent space)	0.653 ± 0.010
Contrastive (spectrogram, latent space)	0.631 ± 0.011
Contrastive (voltage, data space)	0.664 ± 0.011
Contrastive (spectrogram, data space)	0.632 ± 0.010

Table 1: **In pretraining BrainBERT-mini, the contrastive objective performs better than the traditional MSE loss across 14 decoding tasks.** Trained for 100 epochs using the Muon optimizer. The weights of the models are frozen after pretraining with no labels, and a linear regression is applied on the features of the frozen models to obtain the AUROC (mean \pm SEM). The bolded entries indicate best performance (within one SEM of each other).

206 **Superior performance of the contrastive pretraining objective when training BrainBERT-mini**
207 First, we seek to validate the performance of the contrastive objective on an established architecture
208 from past literature - a scaled-down version of BrainBERT [Wang et al., 2023]. After the pretraining
209 phase with no labels concludes, we freeze the model and fine-tune only a single linear layer on top of
210 the model features on our downstream tasks of interest, in order to assess the quality of the generated
211 representations, with results shown in Table 1. Our findings demonstrate that the contrastive loss
212 performed better than its MSE counterpart in all experimental conditions (training on raw voltage
213 vs spectrogram of the signal). Furthermore, we find that pretraining on raw voltage instead of the
214 spectrogram of the signal, as often done in prior work, is beneficial for downstream performance
215 across many tasks. Taking insight from this smaller scale experiment, we next scale up our pretraining
216 to a larger chunk of the dataset, larger models, and expand to the population level information as
217 opposed to only single electrode with CNF-1.

218 **CNF-1 achieves state-of-the-art performance across the decoding tasks** We train CNF-1 on the
219 BrainTreebank dataset for 10 epochs and note that the model gets better at discriminating the true
220 next timestep from random samples over the course of training (Figure B). To assess the quality of
221 the representations learned by our model, and compare it to previously published models, we finetune

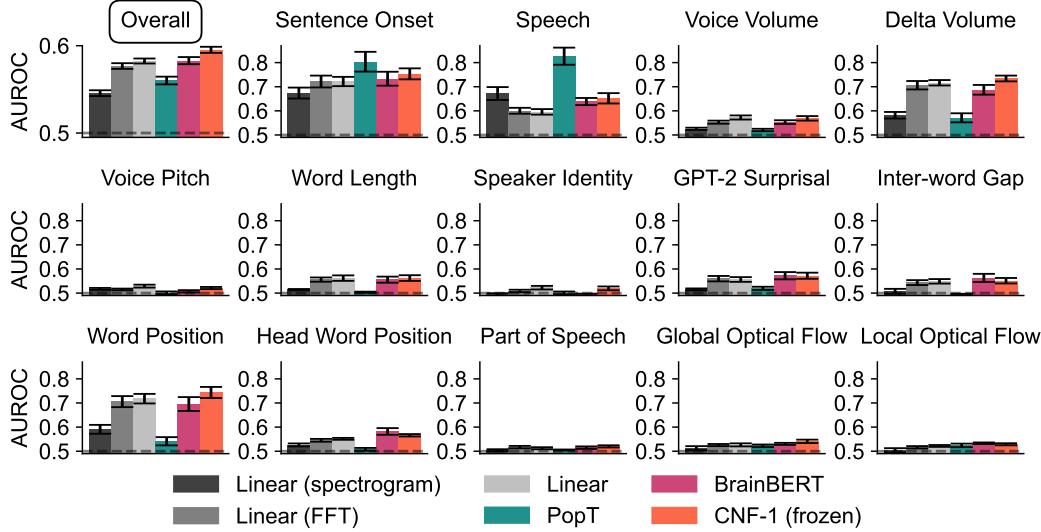


Figure 3: CNF-1 achieves state-of-the-art performance on a suite of benchmark decoding tasks, across subjects and sessions. We evaluate CNF-1, as well as baseline models and models from previous work, on 14 tasks that span language and visual domains, on 12 total recording sessions from 6 human subjects. We find that generally the linear baseline performs surprisingly well, at times outperforming pretrained models. CNF-1 (our model) outperforms all the considered pretrained models as well as all of the considered baselines overall.

the final linear layer that projects the features into a single output dimension on a suite of downstream decoding tasks (Figure 3). We find that generally the linear baseline performs surprisingly well, and surprisingly find that previous state of the art models, while outperforming all of the considered baseline methods on some tasks, fall behind them on others. we associate the lower performance of the BrainBERT and PopT models on some tasks with the low performance of the spectrogram regression baseline: spectrogram of the input signal is the base representation that both of these models use as input, and the baseline task decoding performance for some tasks (e.g. GPT-2 Surprisal, Word Length, Head Word Position etc) requires phase information as well, which is lost when taking the power spectrogram. In our experiments, we also found that pretraining often boosts decoding performance for some tasks (especially Onset and Speech) while decreasing the downstream performance on most other tasks.

While still close to the baselines, CNF-1 outperforms all the considered pretrained models and baselines (Figure 3, top left corner). While for PopT the performance peaks at some tasks and drops for others, CNF-1 shows a more uniform pattern of performance, suggesting that it contains representations that capture more aspects of the neural processing. The state-of-the-art performance of CNF-1 shows the potential and behind the Contrastive Neural Forecasting approach.

Investigating the learned electrode embeddings We now turn to what can be discovered in the data-driven way using our foundation model. An innovation from prior work is our entirely learned electrode embeddings, which replace the traditional coordinate positional embeddings.

We conjecture that over the course of training, the model may discover relationships between the input channels, and use the learned embedding parameters to store them across batches and employ them to improve the performance on the predictive objective. To test this hypothesis, we freeze the pretrained model and examine its learned electrode embeddings (example subject is shown in Figure 4). Across all pairs of electrodes, we find that the distance in embedding space of the model is strongly correlated with the physical distance between the electrodes in the brain (Figure 4a, $r = 0.400, p < 0.001$), despite the fact that spatial information was never available in pretraining. Furthermore, a t-SNE dimensionality analysis reveals spatially clustered groups of electrodes (Figure 4b) that are roughly corresponding to the gross anatomical and functional subdivision of the brain (Figure 4c). We note the consistent difference in the embeddings for the frontal, temporal and occipital lobe electrodes,

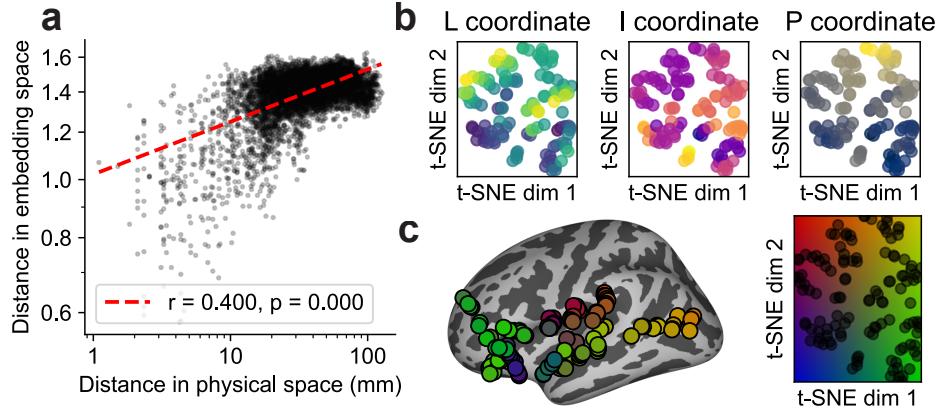


Figure 4: Learned electrode embeddings correlate with rough anatomical and functional brain regions, recovering them in a pure data-driven way. (a) Distance in the embedding space of the learned electrode embeddings is correlated with the distance in physical space in the brain, even though the spatial information was never made available during training. (b) Dimensionality reduction (t-SNE) reveals clustering of the electrode embeddings in the latent space, with the clusters generally grouping together according to the coordinate in the physical space. (c) Visualization of the t-SNE reduction result on an inflated map of the brain, which shows the anatomical locations of the embeddings. The results are shown for a representative Subject 3.

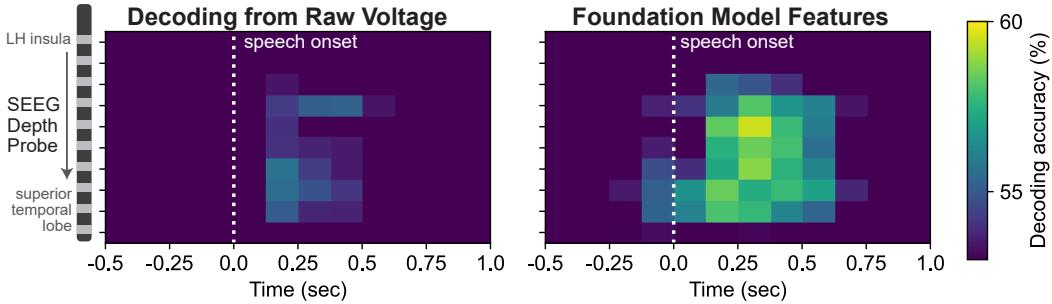


Figure 5: Foundation models enable effective functional mapping of brain regions. Validation of the model on an open human intracranial electroencephalography data (StereoEEG). For a given StereoEEG depth probe, we simulate a language mapping setting where the patient engages in an experimental task with two conditions: speech processing and non-speech. Then, we use either raw voltage (left) or features from our frozen foundation model (right) to decode GPT2 surprisal, as an indication of language processing, from every 125 ms timebin and every contact on the probe. The resulting model enabled stronger decoding of speech onset events than the raw voltage inputs. The probe spans multiple locations in the patient’s brain, enabling localization of the functional language processing region. Model: causal transformer, 5 layers, model hidden dimension 128, 4 attention heads per block, sampling rate 2048 Hz. Trained for 1000 steps with batch size 128 with Contrastive Neural Forecasting, using the Muon optimizer, learning rate 0.003, no weight decay.

251 as well as the language-selective parts of the superior temporal lobe, which suggests that the model
252 rediscovered the gross anatomical layout of the brain.

253 **Mapping function of brain regions with the foundation model features** Next, we validate a
254 practical application of our foundation model for functional brain mapping in a clinical setting. A
255 standard protocol in neurosurgery involves identifying brain regions involved in critical functions,
256 such as language processing, to guide tissue resection [Aron et al., 2021]. Traditionally, this relies on
257 visually inspecting raw intracranial recordings for stimulus-locked activity differences, which may be
258 subtle or ambiguous.

259 We propose using features extracted from our pretrained model to improve the clarity of this mapping.
260 In Figure 5, we compare the decodability of one of our features connected with language processing
261 (GPT-2 surprisal) across time and electrode contacts using raw voltage versus foundation model
262 representations. The model-derived features yield a much sharper spatial and temporal decoding
263 profile, revealing a more localized and time-locked peak in language-related activity along the probe.
264 This result demonstrates that we can enhance functional mapping by amplifying task-relevant signals
265 that may be difficult to detect in the raw data.

266 5 Discussion

267 Our results show that Contrastive Neural Forecasting (CNF) is a viable and scalable framework
268 for learning population-level representations of human intracranial neural activity directly from
269 raw voltage data. By forgoing handcrafted features and instead predicting future neural states in
270 latent space using a contrastive objective, CNF avoids several limitations of traditional approaches,
271 chiefly overfitting to noise and inflexibility in modeling multimodal dynamics. CNF-1 achieves
272 state-of-the-art decoding performance across a diverse suite of language and vision tasks.

273 One striking finding of our work is the effectiveness of learned electrode embeddings in the absence
274 of spatial coordinate information, reaffirming findings by Azabou et al. [2023]. Not only does this
275 challenge prevailing assumptions in neural modeling, but it also suggests that useful structural priors
276 can emerge from data alone when trained at scale, opening new opportunities for interpretability in
277 foundation models of brain activity. Future work will examine the possibility of delineating functional
278 and/or anatomical brain regions [Glasser et al., 2016] based solely on the activity statistics using
279 foundation models such as CNF-1.

280 **Limitations** Our work has several limitations and directions for future research. First, our model
281 outperforms linear baselines by only a small amount, and there is clearly room to grow. We anticipate
282 that training on datasets beyond the BrainTreebank, as well as incremental architecture and training
283 process improvements will greatly enhance the performance of our models. In addition, future work
284 may explore multimodal extensions that incorporate neural data with information about the sensory
285 inputs such as the viewed video. This can be achieved by incorporating CLIP representations of the
286 visual inputs and/or wav2vec or other audio representations (this data is available in datasets such as
287 BrainTreebank, but not used in this work).

288 More broadly, we view CNF as part of an emerging class of tools that treat the brain as a sequence-
289 generating system that is amenable to the same powerful modeling techniques that have revolutionized
290 NLP and vision. In this framing, iEEG signals become the neural analogue of text or pixels: high-
291 dimensional, temporally structured data with rich latent dynamics.

292 **Broader impacts** Importantly, our results support the broader vision of brain foundation models
293 (BFMs): pretraining once on large-scale observational recordings and reusing these representations
294 for a wide range of downstream clinical and scientific applications [Zhou et al., 2025]. For example,
295 we show that CNF-1 can enable real-time functional brain mapping, which could be used in clinical
296 settings such as operating rooms during brain resection surgeries [Richardson, 2022] to define
297 surgical and non-surgical targets. Thus, we anticipate this approach could accelerate workflows in
298 neurosurgery, diagnosis, and closed-loop brain-computer interfaces. The usecases of foundation
299 models should be strictly vetted to adhere to the ethical regulations, especially in the medical usecases
300 and when involved in decision making impacting human lives [Gordon and Seth, 2024].

301 Beyond the clinic, foundation models like CNF-1 offer exciting opportunities in basic neuroscience.
302 By unifying single-channel and population-level representations in a single model, CNF-1 can
303 help researchers probe the functional roles of specific brain regions, and simulate the evolution
304 of neural dynamics under different conditions, and generate new hypotheses to be tested in vivo
305 based on the findings in the foundation models (inception loops; Wang et al. [2025], Walker et al.
306 [2019]). Moreover, as brain foundation models grow larger and more expressive, they may serve as
307 computational proxies for in silico experimentation.

308 Taken together, CNF represents a step toward a general-purpose framework for modeling brain
309 dynamics, supporting the development of robust, scalable, and clinically useful brain-computer
310 interfaces and tools in neuroscience and medicine.

311 **Acknowledgements**

312 This work has been supported by ONR award N00014-19-1-2584, by NSF-CISE award IIS-2151077
313 under the Robust Intelligence program, by the ARO-MURI award W911NF-23-1-0277, by the Simons
314 Foundation SCGB program 1181110, and the K. Lisa Yang ICoN Center.

315 **References**

- 316 Antonis Antoniades, Yiyi Yu, Joseph Canzano, William Wang, and Spencer LaVere Smith. Neuro-
317 former: Multimodal and Multitask Generative Pretraining for Brain Data, March 2024.
- 318 Olivier Aron, Jacques Jonas, Sophie Colnat-Coulbois, and Louis Maillard. Language mapping using
319 stereo electroencephalography: A review and expert opinion. *Frontiers in Human Neuroscience*,
320 15:619521, 2021. doi: 10.3389/fnhum.2021.619521. URL <https://doi.org/10.3389/fnhum.2021.619521>.
- 322 Mehdi Azabou, Vinam Arora, Venkataramana Ganesh, Ximeng Mao, Santosh Nachimuthu, Michael J.
323 Mendelson, Blake Richards, Matthew G. Perich, Guillaume Lajoie, and Eva L. Dyer. A Unified,
324 Scalable Framework for Neural Population Decoding, October 2023.
- 325 Christopher M Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- 326 Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, et al.
327 On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- 328 Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
329 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
330 few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages
331 1877–1901, 2020.
- 332 György Buzsáki, Costas A Anastassiou, and Christof Koch. The origin of extracellular fields and
333 currents-eeg, ecog, lfp and spikes. *Nature Reviews Neuroscience*, 13(6):407–420, 2012.
- 334 Donghong Cai, Junru Chen, Yang Yang, Teng Liu, and Yafeng Li. MBraint: A Multi-channel
335 Self-Supervised Learning Framework for Brain Signals, June 2023.
- 336 Geeling Chau, Christopher Wang, Sabera Talukder, Vighnesh Subramaniam, Saraswati Soedarmadji,
337 Yisong Yue, Boris Katz, and Andrei Barbu. Population Transformer: Learning Population-level
338 Representations of Neural Activity, October 2024.
- 339 Zijian Dong, Ruilin Li, Yilei Wu, Thuan Tinh Nguyen, Joanna Su Xian Chong, Fang Ji, Nathanael
340 Ren Jie Tong, Christopher Li Hsian Chen, and Juan Helen Zhou. Brain-jepa: Brain dynamics
341 foundation model with gradient positioning and spatiotemporal masking, 2024. URL <https://arxiv.org/abs/2409.19407>.
- 343 Matthew F. Glasser, Timothy S. Coalson, Emma C. Robinson, et al. A multi-modal parcellation of
344 human cerebral cortex. *Nature*, 536:171–178, 2016. doi: 10.1038/nature18933. URL <https://doi.org/10.1038/nature18933>.
- 346 Emma C. Gordon and Anil K. Seth. Ethical considerations for the use of brain-computer interfaces
347 for cognitive enhancement. *PLOS Biology*, 22(10):e3002899, 2024. doi: 10.1371/journal.pbio.
348 3002899. URL <https://doi.org/10.1371/journal.pbio.3002899>.
- 349 Joel Herron, Vahe Kremen, John D. Simmeral, et al. The convergence of neuromodulation and
350 brain-computer interfaces. *Nature Reviews Bioengineering*, 2:628–630, 2024. doi: 10.1038/
351 s44222-024-00187-0. URL <https://doi.org/10.1038/s44222-024-00187-0>.
- 352 Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz Cesista, Laker Newhouse, and Jeremy
353 Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. URL <https://kellerjordan.github.io/posts/muon/>.

- 355 Tim Large, Yang Liu, Minyoung Huh, Hyojin Bahng, Phillip Isola, and Jeremy Bernstein. Scalable optimization in the modular norm. In A. Globerson, L. Mackey,
356 D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pages 73501–73548. Curran Associates,
357 Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/8629b0fff229b8a27efb1422e990605f-Paper-Conference.pdf.
- 360
- 361 Trung Le and Eli Shlizerman. STNDT: Modeling Neural Population Activity with a Spatiotemporal
362 Transformer, June 2022.
- 363 Ran Liu, Mehdi Azabou, Max Dabagia, Jingyun Xiao, and Eva L. Dyer. Seeing the forest and the tree:
364 Building representations of both individual and collective dynamics with transformers, October
365 2022.
- 366 Martha J Morrell. Responsive cortical stimulation for the treatment of medically intractable partial
367 epilepsy. *Neurology*, 77(13):1295–1304, 2011.
- 368 Norbert Noury, Joerg F Hipp, and Markus Siegel. Physiological processes non-linearly affect
369 electrophysiological recordings. *eNeuro*, 3(4):ENEURO.0191–16.2016, 2016.
- 370 Josef Parvizi and Sabine Kastner. Promises and limitations of human intracranial electroencephalogram-
371 raphy. *Nature Neuroscience*, 15(2):264–272, 2012.
- 372 Francisco Pereira, Brian Lou, Brian Pritchett, et al. Toward a universal decoder of linguistic meaning
373 from brain activation. *Nature Communications*, 9:963, 2018. doi: 10.1038/s41467-018-03068-4.
374 URL <https://doi.org/10.1038/s41467-018-03068-4>.
- 375 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
376 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
377 Learning transferable visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
- 378
- 379 R. Mark Richardson. Closed-loop brain stimulation and paradigm shifts in epilepsy surgery.
380 *Neurologic Clinics*, 40(2):355–373, 2022. doi: 10.1016/j.ncl.2021.12.002. URL <https://doi.org/10.1016/j.ncl.2021.12.002>.
- 381
- 382 Matthias Schurz, Joaquim Radua, Markus Aichhorn, Fabio Richlan, and Josef Perner. Fractionating
383 theory of mind: A meta-analysis of functional brain imaging studies. *Neuroscience & Biobehavioral
384 Reviews*, 42:9–34, 2014.
- 385 Terrence J Sejnowski, Patricia S Churchland, and J Anthony Movshon. Putting big data to good use
386 in neuroscience. *Nature Neuroscience*, 17(11):1440–1441, 2014.
- 387 Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive
388 coding, 2019. URL <https://arxiv.org/abs/1807.03748>.
- 389
- 390 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz
391 Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL <https://arxiv.org/abs/1706.03762>.
- 392 Ankit Vishnubhotla, Charlotte Loh, Liam Paninski, Akash Srivastava, and Cole Hurwitz. Towards
393 robust and generalizable representations of extracellular data using contrastive learning.
394 *bioRxiv*, 2023. doi: 10.1101/2023.10.30.564831. URL <https://doi.org/10.1101/2023.10.30.564831>. Preprint, not peer reviewed.
- 395
- 396 E. Y. Walker, F. H. Sinz, E. Cobos, et al. Inception loops discover what excites neurons most
397 using deep predictive models. *Nature Neuroscience*, 22:2060–2065, 2019. doi: 10.1038/s41593-019-0517-x. URL <https://doi.org/10.1038/s41593-019-0517-x>.
- 398
- 399 Christopher Wang, Vighnesh Subramaniam, Adam Uri Yaari, Gabriel Kreiman, Boris Katz, Ignacio
400 Cases, and Andrei Barbu. BrainBERT: Self-supervised representation learning for intracranial
401 recordings, February 2023.

- 402 Christopher Wang, Adam Uri Yaari, Aaditya K Singh, Vighnesh Subramaniam, Dana Rosenfarb, Jan
403 DeWitt, Pranav Misra, Joseph R. Madsen, Scellig Stone, Gabriel Kreiman, Boris Katz, Ignacio
404 Cases, and Andrei Barbu. Brain treebank: Large-scale intracranial recordings from naturalistic
405 language stimuli, 2024. URL <https://arxiv.org/abs/2411.08343>.
- 406 Evelyn Y. Wang, Patrick G. Fahey, Zheng Ding, et al. Foundation model of neural activity predicts
407 response to new stimulus types. *Nature*, 640:470–477, 2025. doi: 10.1038/s41586-025-08829-y.
408 URL <https://doi.org/10.1038/s41586-025-08829-y>.
- 409 Joel Ye, Jennifer L. Collinger, Leila Wehbe, and Robert Gaunt. Neural Data Transformer 2: Multi-
410 context Pretraining for Neural Spiking Activity, September 2023.
- 411 Daoze Zhang, Zhizhang Yuan, Yang Yang, Junru Chen, Jingjing Wang, and Yafeng Li. Brant: Foun-
412 dation Model for Intracranial Neural Signal. In *Thirty-Seventh Conference on Neural Information
413 Processing Systems*, November 2023.
- 414 Xinliang Zhou, Chenyu Liu, Zhisheng Chen, Kun Wang, Yi Ding, Ziyu Jia, and Qingsong Wen. Brain
415 foundation models: A survey on advancements in neural signal processing and brain discovery,
416 2025. URL <https://arxiv.org/abs/2503.00580>.

417 **NeurIPS Paper Checklist**

418 **1. Claims**

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

421 Answer: [Yes]

422 Justification: We describe our pretraining approach, and then validate and show its perfor-
423 mance and outline the applications of foundation models in the Results section.

424 Guidelines:

- 425 • The answer NA means that the abstract and introduction do not include the claims
426 made in the paper.
- 427 • The abstract and/or introduction should clearly state the claims made, including the
428 contributions made in the paper and important assumptions and limitations. A No or
429 NA answer to this question will not be perceived well by the reviewers.
- 430 • The claims made should match theoretical and experimental results, and reflect how
431 much the results can be expected to generalize to other settings.
- 432 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
433 are not attained by the paper.

434 **2. Limitations**

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

436 Answer: [Yes]

437 Justification: We included a limitations sections in our Discussion section.

438 Guidelines:

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

465 **3. Theory assumptions and proofs**

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

468 Answer: [NA]

469 Justification: This paper does not introduce any theory results or theorems, only showing
470 empirical results.

471 Guidelines:

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

482 4. Experimental result reproducibility

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

486 Answer: [Yes]

487 Justification: Our approach and architectures used for pretraining are described in the
488 Approach and Experimental Setup sections in detail, as well as in the included Appendix; in
489 addition, the code will be released upon publication.

490 Guidelines:

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

522 5. Open access to data and code

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

526 Answer: [Yes]

527 Justification: We release the Github repository with all of the code required to train the
528 models and reproduce the experiments, as well as the model weights, upon publication.

529 Guidelines:

- 530 • The answer NA means that paper does not include experiments requiring code.
- 531 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 532 • While we encourage the release of code and data, we understand that this might not be
533 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
534 including code, unless this is central to the contribution (e.g., for a new open-source
535 benchmark).
- 536 • The instructions should contain the exact command and environment needed to run to
537 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 538 • The authors should provide instructions on data access and preparation, including how
539 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 540 • The authors should provide scripts to reproduce all experimental results for the new
541 proposed method and baselines. If only a subset of experiments are reproducible, they
542 should state which ones are omitted from the script and why.
- 543 • At submission time, to preserve anonymity, the authors should release anonymized
544 versions (if applicable).
- 545 • Providing as much information as possible in supplemental material (appended to the
546 paper) is recommended, but including URLs to data and code is permitted.

547 6. Experimental setting/details

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

553 Answer: [Yes]

554 Justification: We provide information about our pretraining, architecture, and hyperparam-
555 eters, and optimizer in the Experimental Setup section in the paper.

556 Guidelines:

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

562 7. Experiment statistical significance

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

565 Answer: [Yes]

566 Justification: For our empirical results, we report standard error across cross-val folds, as
567 well as across subjects in case of evaluating on data from multiple subjects.

568 Guidelines:

- 569 • The answer NA means that the paper does not include experiments.
- 570 • The authors should answer “Yes” if the results are accompanied by error bars, confi-
571 dence intervals, or statistical significance tests, at least for the experiments that support
572 the main claims of the paper.

- 573 • The factors of variability that the error bars are capturing should be clearly stated (for
 574 example, train/test split, initialization, random drawing of some parameter, or overall
 575 run with given experimental conditions).
 576 • The method for calculating the error bars should be explained (closed form formula,
 577 call to a library function, bootstrap, etc.)
 578 • The assumptions made should be given (e.g., Normally distributed errors).
 579 • It should be clear whether the error bar is the standard deviation or the standard error
 580 of the mean.
 581 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 582 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 583 of Normality of errors is not verified.
 584 • For asymmetric distributions, the authors should be careful not to show in tables or
 585 figures symmetric error bars that would yield results that are out of range (e.g. negative
 586 error rates).
 587 • If error bars are reported in tables or plots, The authors should explain in the text how
 588 they were calculated and reference the corresponding figures or tables in the text.

589 **8. Experiments compute resources**

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

593 Answer: [Yes]

594 Justification: we specify the computational requirements when describing pretraining in our
 595 Experimental Setup section.

596 Guidelines:

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

605 **9. Code of ethics**

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

608 Answer: [Yes]

609 Justification: We adhere to the code of ethics.

610 Guidelines:

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

616 **10. Broader impacts**

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

619 Answer: [Yes]

620 Justification: We have included a broader impacts subsection in our Discussion section.

621 Guidelines:

- 622 • The answer NA means that there is no societal impact of the work performed.

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

11. Safeguards

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

Answer: [Yes]

Justification: Our paper describes ethical consideration with use of brain foundation models in the Discussion section.

Guidelines:

- 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.

12. Licenses for existing assets

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

Answer: [Yes]

Justification: The authors own all the assets in this paper, and credit with references whenever the openly available resources are used for the experiments or datasets.

Guidelines:

- 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.

- 675 • For scraped data from a particular source (e.g., website), the copyright and terms of
 676 service of that source should be provided.
 677 • If assets are released, the license, copyright information, and terms of use in the
 678 package should be provided. For popular datasets, paperswithcode.com/datasets
 679 has curated licenses for some datasets. Their licensing guide can help determine the
 680 license of a dataset.
 681 • For existing datasets that are re-packaged, both the original license and the license of
 682 the derived asset (if it has changed) should be provided.
 683 • If this information is not available online, the authors are encouraged to reach out to
 684 the asset's creators.

685 **13. New assets**

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

688 Answer: [Yes]

689 Justification: We document all of the code created for this paper in the comments and
 690 README files in the Github repository to be shared upon publication.

691 Guidelines:

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

700 **14. Crowdsourcing and research with human subjects**

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

704 Answer: [NA]

705 Justification: Our paper doesn't involve any crowdsourcing for our experiments and does
 706 not perform new experiments with human subjects.

707 Guidelines:

- 708 • The answer NA means that the paper does not involve crowdsourcing nor research with
 709 human subjects.
- 710 • Including this information in the supplemental material is fine, but if the main contribu-
 711 tion of the paper involves human subjects, then as much detail as possible should be
 712 included in the main paper.
- 713 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
 714 or other labor should be paid at least the minimum wage in the country of the data
 715 collector.

716 **15. Institutional review board (IRB) approvals or equivalent for research with human
 717 subjects**

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

722 Answer: [NA]

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

726 Guidelines:

- 727 • The answer NA means that the paper does not involve crowdsourcing nor research with
728 human subjects.
729 • Depending on the country in which research is conducted, IRB approval (or equivalent)
730 may be required for any human subjects research. If you obtained IRB approval, you
731 should clearly state this in the paper.
732 • We recognize that the procedures for this may vary significantly between institutions
733 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
734 guidelines for their institution.
735 • For initial submissions, do not include any information that would break anonymity (if
736 applicable), such as the institution conducting the review.

737 **16. Declaration of LLM usage**

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

742 Answer: [NA]

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

747 Guidelines:

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

752 A Decoding Tasks and Split Construction

753 **Decoding Tasks** The decoding benchmark used in this paper includes 14 decoding tasks derived
 754 from multimodal annotations in the BrainTreebank dataset [Wang et al., 2024]. These tasks span
 755 audio, vision, and language modalities. All tasks are cast as binary classification problems to ensure
 756 uniformity in evaluation across task types and models.

- 757 • **Scalar features** (e.g., GPT2 surprisal, pitch, volume): For each session, values are thresh-
 758 olded such that the top 25% of the distribution are labeled as the positive class, and the
 759 bottom 25% as the negative class. The middle 50% of values are excluded from training and
 760 evaluation to reduce ambiguity around class boundaries.
- 761 • **Categorical features** (e.g., part-of-speech, speaker identity): For each feature, a single
 762 target class is selected (typically the most frequent), and the task is defined as a one-vs-rest
 763 binary classification problem.

764 All tasks are aligned to word onsets. Neural data is segmented into 1-second windows starting at the
 765 onset of each word. Unless otherwise stated, all decoding experiments use these 1-second segments
 766 of neural activity as model inputs, and the corresponding BrainTreebank annotations as binary labels.
 767 See more information about the decoding tasks in the tables below.

Subj.	Age (yrs.)	# Elec- trodes	Movie	Recording time (hrs)	Used in benchmark
1	19	154	Thor: Ragnarok	1.83	x
			Fantastic Mr. Fox	1.75	
			The Martian	0.5	x
2	12	162	Venom	2.42	x
			Spider-Man: Homecoming	2.42	
			Guardians of the Galaxy	2.5	
			Guardians of the Galaxy 2	3	x
			Avengers: Infinity War	4.33	
			Black Panther	1.75	
			Aquaman	3.42	
3	18	134	Cars 2	1.92	x
			Lord of the Rings 1	2.67	x
			Lord of the Rings 2 (extended edition)	3.92	
4	12	188	Incredibles	1.15	
			Shrek 3	1.68	x
			Megamind	2.43	x
5	6	156	Fantastic Mr. Fox	1.5	
6	9	164	Megamind	2.58	
			Toy Story	1.33	
			Coraline	1.83	
7	11	246	Cars 2	1.75	x
			Megamind	1.77	x
8	4.5	162	Sesame Street Episode	1.28	
9	16	106	Ant Man	2.28	
10	12	216	Cars 2	1.58	x
			Spider-Man: Far from Home	2.17	x

Table S1: **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).

#	Feature	Description	Benchmark Task
1	global_flow (visual)	A camera motion proxy. The maximal average dense optical flow vector magnitude	Same as above
2	local_flow (visual)	A large displacement proxy. The maximal optical flow vector magnitude	Same as above
3	volume (auditory)	Average root mean squared watts of the audio	Binary classification: low (0%-25%) vs high (75%-100%)
4	pitch (auditory)	Average pitch of the audio	Same as above
5	delta_volume (auditory)	The difference in average RMS of the 500ms windows pre- and post-word onset	Same as above
6	speech (language)	Whether any speech is present in the given time interval	Binary classification
7	onset (language)	Whether a new sentence starts in the interval, or there is no speech at all	Binary classification
8	gpt2_surprisal (language)	Negative-log transformed GPT-2 word probability (given preceding 20s of language context)	Binary classification: low (0%-25%) vs high (75%-100%)
9	word_length (language)	Word length (ms)	Same as above
10	word_gap (language)	Difference between previous word offset and current word onset (ms)	Same as above
11	word_index (language)	The word index in its context sentence	2-way classification: 0 (the first word in the sentence), or other (1)
12	word_head_pos (language)	The relative position (left/right) of the word's dependency tree head	Binary classification
13	word_part_speech (language)	The word Universal Part-of-Speech (UPOS) tag	2-way classification: verb (0), or other (1)
14	speaker (multimodal)	The movie character that speaks the given word.	2-way classification: most frequent speaker (0), or other (1)

Table S2: **Extracted visual, auditory, and language features used to create the evaluations.** 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].

768 **Train/Test Split Construction** To probe model generalization under increasingly challenging
 769 conditions, we define the following split strategies:

- 770 • **Same Subject / Same Movie (SS/SM):** Training and testing data are drawn from the same
 771 subject and same movie (trial of recording). A contiguous 80/20 train-test split is applied,
 772 ensuring the training block precedes the test block to reduce temporal autocorrelation.
 773 Performance is computed via 5-fold cross-validation.
- 774 • **Same Subject / Different Movie (SS/DM):** Data is drawn from the same subject across two
 775 different movies. For the two movies selected for every subject for evaluation, both ways to
 776 split the pair into the train and test movie are used, and the resulting AUROC is averaged
 777 between the two splits.

778 BrainBERT-mini decoding experiments were run on the SS/SM split. CNF-1 (Contrastive Neural
 779 Forecasting) and Functional mapping analyses (i.e., the spatiotemporal decoding maps shown in
 780 Figure 5) evaluations were run on the SS/DM split.

781 We discard the data from electrodes which were labeled as corrupted by the BrainTreebank authors
 782 [Wang et al., 2024].

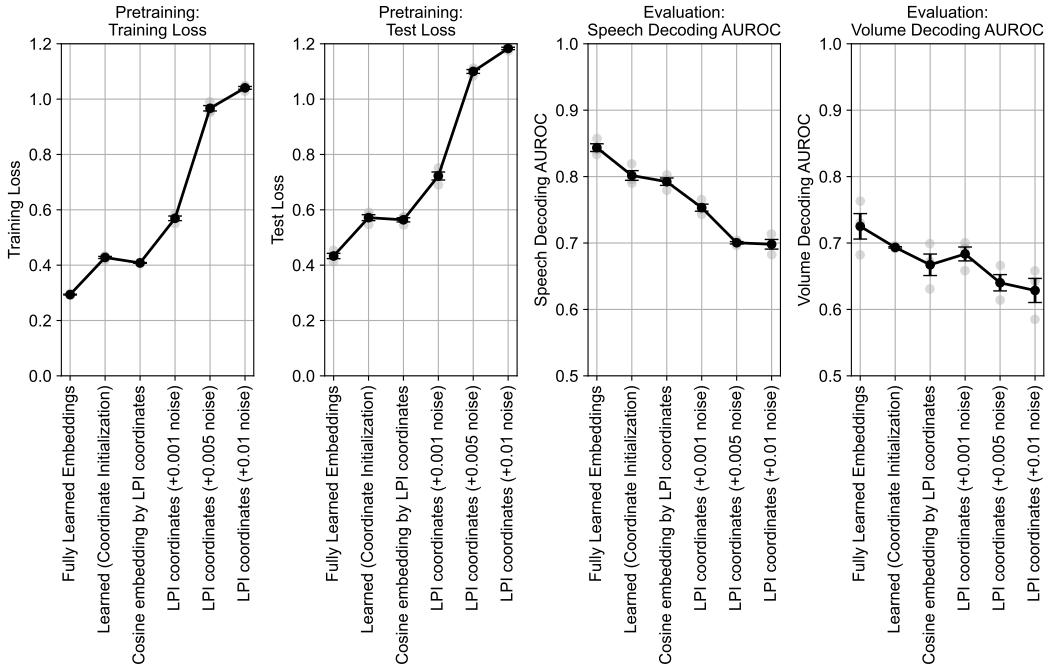


Figure S1: Fully learned electrode embeddings perform the best in Contrastive Neural Forecasting across both pretraining and evaluation. The two left graphs show the training and test loss during autoregressive pretraining, respectively. Fully learned embeddings outperform both traditional positional cosine embeddings with electrode LPI coordinates and the approach where the embeddings are initialized with the cosine embeddings but then are allowed to be updated during pretraining. Adding noise to the positional electrode embedding only increases the train and test pretraining error. The two right plots show the evaluation decoding AUROC (with frozen model weights), and demonstrate that the evaluation performance also decreases with increased pretraining loss. The error bars show the mean and s.e.m. across 3 random seeds. This experiment was performed in a model where the CNF objective was applied at the CLS token.

783 **Functional mapping experiment** To estimate the time course of information processing in the
 784 brain across space in a brain area (Figure 5), we used a sliding window of 125 ms across neural
 785 activity, in steps of 125 ms from -500 ms to $+1000$ ms relative to word onset. For each time bin
 786 and electrode, a separate linear decoder is trained for each task, either with raw voltage traces
 787 acting as features, or from the model features. The resulting decoding scores are averaged across
 788 cross-validation folds. The analysis was run specifically with subject 3, trial 0.

789 **B Details on the Model and Pretraining**

790 **Pretraining hyperparameters** All models were pretrained with learning rate 0.003 which we
 791 found works best across the range of different model sizes and architectures (which might be a feature
 792 of the Muon optimizer, for a discussion see Jordan et al. [2024], Large et al. [2024]). We use batch
 793 size 100, and for every electrode we use batch norm to normalize the input voltage traces across
 794 the batch and timesamples dimensions. We discard the data from electrodes which were labeled as
 795 corrupted by the BrainTreebank authors [Wang et al., 2024].

796 **Learned electrode embeddings** We found in our experiments that fully learned electrode embed-
 797 dings resulted in lower pretraining loss and higher decoding performance compared to the traditional
 798 approach from prior work [Chau et al., 2024, Zhang et al., 2023] which provides cosine positional
 799 embeddings from the electrode physical coordinates in 3D space (Supplementary Figure S1).

800 **C Comparison to baselines and previous methods**

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

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

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

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

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

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

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

827 **Population Transformer** For Population Transformer, we followed the implementation and used
828 the weights from [Chau et al., 2024]. The fine-tuning protocol is taken to be directly the same as in
829 the authors' original paper (including linear rate, number of epochs, a factor of 10 between learning
830 rates of the linear output layer vs the transformer blocks, etc). We found that frozen Population
831 Transformer's performance was almost always at chance and that pretraining through the whole
832 model was necessary to achieve comparable performance to other methods.