

---

# Population Transformer: Learning Population-level Representations of Intracranial Activity

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We present a self-supervised framework that learns population-level codes for  
2 intracranial neural recordings at scale, unlocking the benefits of representation  
3 learning for a key neuroscience recording modality. The Population Transformer  
4 (PopT) lowers the amount of data required for decoding experiments, while in-  
5 creasing accuracy, even on never-before-seen subjects and tasks. We address  
6 two key challenges in developing PopT: sparse electrode distribution and varying  
7 electrode location across patients. PopT stacks on top of pretrained representa-  
8 tions and enhances downstream tasks by enabling learned aggregation of multiple  
9 spatially-sparse data channels. Beyond decoding, we interpret the pretrained PopT  
10 and fine-tuned models to show how it can be used to provide neuroscience in-  
11 sights learned from massive amounts of data. We release a pretrained PopT to  
12 enable off-the-shelf improvements in multi-channel intracranial data decoding and  
13 interpretability.

## 14 1 Introduction

15 Building effective representations of neural recordings is an important tool in enabling neuroscience  
16 research. We are particularly interested in modeling intracranial recordings, which rely on probes  
17 placed within the brain to provide high temporal resolution recordings of local neural activity [1, 2].  
18 Because of its dispersed placement within the brain volume, intracranial recordings suffer from  
19 data sparsity. Moreover, there is often significant variability in probe placement across subjects  
20 [1, 2], leading to high variability in input channel meaning. Historically, constructing decoders from  
21 intracranial data has relied on supervised learning [3, 2, 4–6], but this requires experimenters to  
22 collect annotated data, which is scarce due to patient availability and labor-intensive labeling.

23 To improve decoding data-efficiency, self-supervised pretraining on unannotated data can be employed  
24 to first learn generic representations of the recordings. This means that the model does not have to  
25 use valuable annotated samples to learn how to do feature extraction before it can do classification,  
26 improving the reach of neuroscientific research.

27 In this paper, we are interested in developing generic representations of multi-channel intracranial  
28 recordings that enable efficient adaptation to a wide range of downstream decoding tasks. Prior work  
29 has shown how to pretrain subject-specific [7] or channel-specific [8] models of intracranial data, but  
30 such techniques ignore inter-channel relationships or commonalities that might exist across subjects.  
31 The most general approach would be to pretrain using data from multiple datasets, but would require  
32 tackling the aforementioned challenges of sparse electrode coverage and variable electrode placement  
33 between subjects.

34 We propose Population Transformer (PopT), a self-supervised pretraining approach on transformers  
35 [9] that learns subject-generic representations of arbitrary electrode ensembles (Figure 1). During

36 pretraining, we simultaneously optimize both a channel-level and ensemble-level objective, that  
37 requires the model to (1) build representations of channels in the context of other channels and (2)  
38 meaningfully distinguish temporal relationships between different ensembles of channels.

39 Our approach builds on top of pretrained single-channel embeddings, which has two key advantages.  
40 First, by separating the single-channel embedding and multi-channel-aggregation into different  
41 modules, we make our approach agnostic to the specific type of temporal embedding used, leaving  
42 room for future independent improvements, an approach that has been validated in video modeling  
43 [10]. Second, by taking advantage of pretrained channel embeddings, we make our population-level  
44 training lightweight, allowing for adoption in lower compute resource environments.

45 Empirically, we find that our pretrained PopT outperforms existing aggregation approaches, highlight-  
46 ing the usefulness of learning spatial relationships during pretraining. Moreover, we find that these  
47 benefits hold even for subjects not seen during pretraining, lending to its usefulness for new subject  
48 decoding. We also find that the pretrained PopT weights themselves reveal interpretable patterns for  
49 neuroscientific study.

50 Our main contributions are:

- 51 1. a lightweight, generic SSL framework, Population Transformer (PopT) that learns arbitrary  
52 joint representations of channel embeddings across unannotated datasets of neural activity.
- 53 2. a demonstration that a pretrained PopT benefits downstream performance, interpretability,  
54 and training efficiency in comparison to baseline aggregation approaches.
- 55 3. a pretrained off-the-shelf model that computes population-level representations of intracra-  
56 nial neural recordings.

## 57 2 Population Transformer Approach

58 We propose a self-supervised training scheme to learn a subject-generic model that handles arbitrary  
59 electrode configurations. Our loss function has two discriminative components: (1) **channel-wise**  
60 — the model identifies outlier channels swapped with activity from a different timepoint, requiring  
61 sensitivity to surrounding channel context; (2) **ensemble-wise** — the model determines if two channel  
62 ensembles occurred consecutively, requiring ensemble-level context awareness. This objective  
63 effectively simulates many in-silico brain ablations, training the model to learn connections between  
64 regions in the presence of these ablations.

65 A key aspect of our method is the fact that our objective is **discriminative**, rather than reconstructive,  
66 as is often the case in self-supervision [11, 8]. We found this to be necessary, because in practice, the  
67 temporal embeddings often have low effective dimension (see [8]), and reconstruction rewards the  
68 model for overfitting to “filler” dimensions in the feature vector (see Appendix G).

69 To make our model subject and configuration agnostic, we provide the 3D position of each electrode,  
70 providing a common position embedding across subjects. We also vary subset sizes during sampling  
71 to ensure the model handles ensembles of different sizes, accommodating neuroscience experiments  
72 with varying electrode counts and analysis levels. Additionally, we select disjoint subsets to prevent  
73 the model from solving tasks through trivial copying.

74 **Architecture** A schematic of our Population Transformer (PopT) approach is shown in Figure 1.  
75 Consider a given subject with  $N_c$  channels indexed by  $C = \{1, \dots, N_c\}$ . Activity from channel  $i$  at  
76 time  $t$  can be denoted by  $x_i^t$ . The PopT takes as input an interval of brain activity  $X = \{x_i^t | i \in C\}$   
77 from a given time  $t$  and a special [CLS] token. Per channel, each interval of brain activity is passed  
78 through a temporal embedding model  $T$ , in our case BrainBERT, to obtain a representation of each  
79 channel’s temporal context.

80 Before being inputted to the PopT, each channel’s 3D position is added to this embedding, so the  
81 final input is  $X_B = \{T(x) + pos(i) + \mathcal{N}(0, \sigma) | x \in X\}$ . Here, we add Gaussian fuzzing to prevent  
82 overfitting to a particular set of coordinates. Spatial location is given by the electrode’s Left, Posterior,  
83 and Inferior coordinates [12]; see [8] for details on how these were obtained. Each coordinate is  
84 encoded using sinusoidal position encoding [9]. And the three encodings are concatenated together  
85 to form the position embedding  $pos(i) = [e_{\text{left}}; e_{\text{post.}}; e_{\text{inf}}]$ .

86 The core of PopT consists of a transformer encoder stack (see Appendix C: Architectures). The output  
87 of the PopT are spatial-contextual embeddings of the channels  $Y = \{y_i\}$  as well as an embedding of

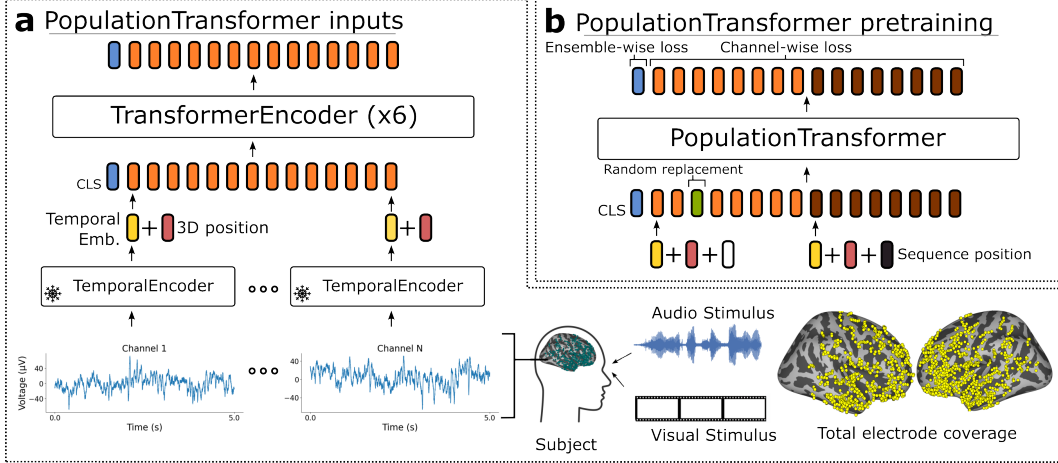


Figure 1: **Schematic of our approach.** The inputs to our model (a) are the combined neural activities from a collection of intracranial electrodes in a given time interval. These are passed to a frozen temporal embedding model, which produces a set of time-contextual embedding vectors (yellow). The 3D position of each electrode (red) is added to these vectors to produce the model inputs (orange). The PopT produces space-contextual embeddings for each electrode and a [CLS] token (blue), which can be fine-tuned for downstream tasks. During pretraining, (b) the PopT is trained on two objectives simultaneously: channel-wise loss and ensemble-wise loss. In channel-wise, PopT must determine whether an input channel has been replaced (green) with activity at a random other time that is inconsistent with the majority of inputs (orange or brown). In ensemble-wise, PopT determines whether two different ensembles (orange vs brown) represent consecutive or non-consecutive times.

88 the CLS token  $y_{cls}$ . During pretraining, the PopulationTransformer additionally is equipped with a  
 89 linear layer head for the [CLS] token output and separate linear layer heads for all other individual  
 90 token outputs. These produce the scalars  $\tilde{y}_{cls}$  and  $\tilde{y}_i$  and respectively, which are used in the objective  
 91 (see Figure 1a).

92 **Pretraining** Our pretraining objective has two components: **channel-wise** and **ensemble-wise**  
 93 losses (see Figure 1b). First, we describe our channel-wise discriminative learning. The model must  
 94 determine whether a channels activity has been swapped with activity from a random time. Precisely,  
 95 activity from each channel  $i$  is drawn from a time  $t_i$ . All channels are drawn from the same time  
 96  $t_i = T$ , and then 10% of the channels are randomly selected to have their activity replaced with  
 97 activity from the same channel, but taken from a random point in time  $t_i \neq T$ . Then, given the token  
 98 outputs of PopT, the channel-wise loss function  $\mathcal{L}_C$  is the binary cross entropy.

99 Next, we describe the ensemble-wise discrimination task. Two different subsets of channels  $S_A, S_B \subset$   
 100  $C$  are chosen with the condition that they be disjoint  $S_A \cap S_B = \emptyset$ . During pretraining, the  
 101 model receives the activities from these channels at separate times  $X_A = \{x_i^{t'} \mid i \in S_A\}$  and  
 102  $X_B = \{x_i^{t''} \mid i \in S_B\}$ . The objective of the task is then to determine whether these states  $X_A$  and  
 103  $X_B$  have occurred consecutively in time or are separated by some further, randomly selected interval.  
 104 Given the output of the classification head, the loss function  $\mathcal{L}_N$  is simply the binary cross entropy.  
 105 Then, our complete objective function is  $\mathcal{L} = \mathcal{L}_N + \mathcal{L}_C$ .

106 **Fine-tuning** During fine-tuning, the [CLS] intermediate representation,  $\tilde{y}_{cls}$  of the pretrained PopT  
 107 is passed through a single layer linear neural network to produce a scalar  $\hat{y}_{cls}$ . This scalar is the input  
 108 to binary cross entropy loss for our decoding tasks (see Section 3). After fine-tuning, we perform  
 109 interpretability analysis on [CLS] attention weights with techniques outlined in Appendix F.

### 110 3 Experiment Setup

111 **Data** We use the publicly available subject data from [8]. Data was collected from 10 subjects  
 112 (total 1,688 electrodes, with a mean of 167 electrodes per subject) who watched 26 movies while  
 113 intracranial probes recorded their brain activity. The movie transcripts were aligned to the brain

114 activity so that features such as volume, pitch, etc. could be associated with the corresponding sEEG  
 115 readings. 19 of the sessions are used for pretraining. 7 of the sessions are held-out for evaluation.

116 **Decoding** We evaluate the effectiveness of our pretrained PopT model by fine-tuning it on the  
 117 four downstream decoding task used in the evaluation of [8]. Two of the tasks are audio focused:  
 118 determining whether a word is spoken with a high or low pitch and determining whether a word is  
 119 spoken loudly or softly. And two of the tasks have a more linguistic focus: determining whether the  
 120 beginning of a sentence is occurring or determining whether any speech at all is occurring.

121 Our approach enables decoding on any arbitrary size of ensemble. We verify that our model is able to  
 122 leverage additional channels for improved decoding performance that scales the number of inputs. To  
 123 test this, we first order the electrodes by their individual linear decodability per task, and we increase  
 124 the number of channels available to the model at fine-tuning time.

125 **Baselines** We want to determine whether the information about spatial relationships learned during  
 126 pretraining was useful at fine-tuning time. For comparison, we concatenate the single-channel  
 127 temporal embeddings and train a linear (Linear) or non-linear (DeepNN) aggregator on the decoding  
 128 task. This sets a baseline for how much improvement is achievable from existing aggregation  
 129 approaches [13]. To determine whether our performance can be attributed to using a more powerful  
 130 architecture, we also fine-tune a PopT without pretraining, i.e. with randomly initialized weights.

## 131 4 Results

132 **Decoding Performance** We find that using a pretrained PopT significantly benefits downstream  
 133 decoding compared to baseline channel aggregation techniques (Table 1). Additionally, while scaling  
 134 performance with increasing number of channels is a challenging task for most baseline aggregation  
 135 approaches, a pretrained PopT is able to scale well with increasing ensemble sizes (Figure 2a).

136 To gain confidence on our method’s generalizability to channel encoders, we applied our framework  
 137 to two different channel encoders: (1) an sEEG temporal encoder (BrainBERT [8]) and (2) a general  
 138 time-series encoder (TOTEM [14]). We see significant boosts in performance with the pretrained  
 139 PopT in both cases when compared with baseline aggregation approaches, across all 4 auditory-  
 140 linguistic tasks (Table 1). These results suggest that our framework can generalize to benefit joint  
 141 aggregation of other single-channel embeddings and neural recording modalities.

	Pitch	Volume	Sent. Onset	Speech/Non-speech
BrainBERT: single channel	0.53 ± 0.05	0.58 ± 0.08	0.68 ± 0.04	0.66 ± 0.09
Linear + BrainBERT	0.59 ± 0.08	0.66 ± 0.08	0.70 ± 0.09	0.71 ± 0.11
Deep NN + BrainBERT	0.58 ± 0.08	0.67 ± 0.08	0.71 ± 0.10	0.72 ± 0.10
Non-pretrained PopT	0.53 ± 0.06	0.61 ± 0.13	0.74 ± 0.10	0.70 ± 0.08
Pretrained PopT	<b>0.69 ± 0.07</b>	<b>0.84 ± 0.06</b>	<b>0.86 ± 0.05</b>	<b>0.89 ± 0.07</b>
<hr/>				
TOTEM: single channel	0.53 ± 0.01	0.53 ± 0.02	0.69 ± 0.03	0.65 ± 0.04
Linear + TOTEM	0.55 ± 0.02	0.66 ± 0.03	0.79 ± 0.04	0.77 ± 0.05
Deep NN + TOTEM	0.57 ± 0.02	0.67 ± 0.03	0.78 ± 0.03	0.75 ± 0.05
Non-pretrained PopT	0.53 ± 0.02	0.64 ± 0.02	0.79 ± 0.03	0.77 ± 0.05
Pretrained PopT	<b>0.60 ± 0.02</b>	<b>0.73 ± 0.02</b>	<b>0.86 ± 0.03</b>	<b>0.84 ± 0.06</b>

Table 1: **Pretraining PopT is critical to downstream decoding performance.** We test on a variety of audio-linguistic decoding tasks (see Section 3) with either a single channel (row 1) or 90 channels (rows 2-5) as input. Two different temporal encoders are used: BrainBERT [8] (top section) and TOTEM [14] (bottom section). Shown are the ROC-AUC mean and standard error across subjects. We see that all aggregation approaches (rows 2-5) outperform single-channel decoding (row 1). Pretraining PopT and then fine-tuning it for downstream decoding results in significantly better performance (bold) compared to non-pretrained aggregation approaches (rows 2-4). This gain cannot be explained by simply providing more temporal embeddings, as evidenced by the performance of Linear and Deep NN (rows 2 and 3) that take the concatenated raw temporal embeddings as input. Neither can the gain be attributed to simply using a Transformer architecture, as is shown by a comparison with a non-pretrained PopT (row 4).

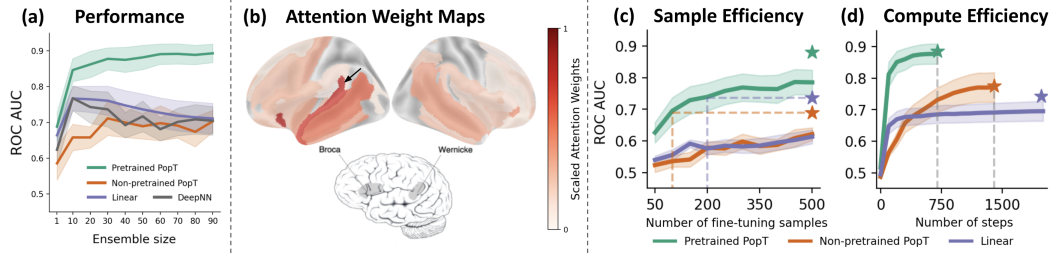


Figure 2: **(a) Pretrained PopT enables downstream performance scaling with ensemble size** Increasing channel ensemble size from 1 to 90 (x-axis), we see pretrained PopT (green) decoding performance (y-axis) not only beat non-pretrained approaches (orange, purple, grey), but also continually improve more with increasing channel count. Shaded bands show the standard error across subjects. **(b) Attention weights from a fine-tuned PopT identify candidate functional brain regions.** Candidate functional maps can be read from attention weights of a PopT fine-tuned on our decoding tasks. Note the weight placed on regions near Wernicke’s area (black arrow) for this Speech vs. Non-speech tuned model. Lower brain figure highlight regions related to auditory-linguistic processing such as language production area Broca’s area and language understanding Wernicke’s area (adapted from [15]). **(c) Pretrained PopT is more sample efficient when fine-tuning.** Varying the number of samples available to each model at train time (x-axis), we see how the pretrained PopT is highly sample efficient, requiring only a fraction of samples to reach the full performance level of non pretrained aggregation approaches (dashed lines). Bands show standard error across test subjects. Stars indicate performance of the model trained on the full fine-tuning dataset. **(d) Pretrained PopT is consistently more compute efficient when fine-tuning.** Number of steps required for each model to reach final performance during fine-tuning (dashed lines). We find that pretrained PopT consistently requires fewer than 750 steps (each step is training on a batch size of 256) to converge, in contrast to the 2k steps required for the non pretrained PopT. Linear aggregation can be similarly compute efficient, but occasionally benefits from more training steps depending on dataset size. Bands show standard error across test subjects. Stars indicate fully trained performance.

142 **Interpretability** To analyze what our massively pretrained + fine-tuned model for sEEG data may be  
 143 doing, we uncover the attention weights the model places on each input channel. We find agreement  
 144 in our model’s attention placement with brain regions typically involved in language processing (e.g.  
 145 Wernicke’s area), especially in the Speech vs. Non-speech downstream task (Figure 2b).

146 **Efficiency** To show that our technique is accessible to low data and compute regimes, we demonstrate  
 147 that a pretrained PopT reaches the same decoding performance as other baseline approaches with  
 148 an order of magnitude fewer samples and steps (Figure 2c and d). Pretraining PopT itself on more  
 149 unannotated data is also an order of magnitude more lightweight than pretraining existing end-to-end  
 150 temporal-spatial models (see Appendix B). By focusing on population-level learning and leveraging  
 151 the growing base of pretrained single-channel embedding techniques, our framework is efficient for  
 152 learning new decoding tasks and continual pretraining.

## 153 5 Conclusion

154 We presented a self-supervised learning scheme for learning effective representations of intracranial  
 155 activity from temporal embeddings. By decoupling temporal and spatial feature extraction, we are  
 156 able to leverage existing temporal embeddings to learn spatiotemporal representations efficiently  
 157 and with a smaller number of parameters. We showed that self-supervised pretraining imbues our  
 158 model with knowledge of spatial relationships between these embeddings and improved downstream  
 159 decoding that scales with the number of available channels. This scheme produces interpretable  
 160 weights from which attention weight maps can be read to help uncover learned relationships from  
 161 the massively pretrained framework. Finally, we release the pretrained weights for our PopT with  
 162 BrainBERT inputs as well as our code for plug-and-play pretraining with any temporal embedding.

## References

- 163
- 164 [1] Josef Parvizi and Sabine Kastner. Promises and limitations of human intracranial electroen-  
165 cephalography. *Nature neuroscience*, 21(4):474–483, 2018.
- 166 [2] Christian Herff, Dean J Krusienski, and Pieter Kubben. The potential of stereotactic-EEG for  
167 brain-computer interfaces: current progress and future directions. *Frontiers in neuroscience*, 14:  
168 483258, 2020.
- 169 [3] Sina Faezi, Rozhin Yasaei, and Mohammad Abdullah Al Faruque. Htnet: Transfer learning  
170 for golden chip-free hardware trojan detection. In *2021 Design, Automation & Test in Europe  
171 Conference & Exhibition (DATE)*, pages 1484–1489. IEEE, 2021.
- 172 [4] Stephanie Martin, Iñaki Iturrate, José del R Millán, Robert T Knight, and Brian N Pasley.  
173 Decoding inner speech using electrocorticography: Progress and challenges toward a speech  
174 prosthesis. *Frontiers in neuroscience*, 12:367292, 2018.
- 175 [5] Sean L Metzger, Kaylo T Littlejohn, Alexander B Silva, David A Moses, Margaret P Seaton,  
176 Ran Wang, Maximilian E Dougherty, Jessie R Liu, Peter Wu, Michael A Berger, et al. A  
177 high-performance neuroprosthesis for speech decoding and avatar control. *Nature*, 620(7976):  
178 1037–1046, 2023.
- 179 [6] Francis R Willett, Erin M Kunz, Chaofei Fan, Donald T Avansino, Guy H Wilson, Eun Young  
180 Choi, Foram Kamdar, Matthew F Glasser, Leigh R Hochberg, Shaul Druckmann, et al. A  
181 high-performance speech neuroprosthesis. *Nature*, 620(7976):1031–1036, 2023.
- 182 [7] Trung Le and Eli Shlizerman. Stndt: Modeling neural population activity with spatiotemporal  
183 transformers. *Advances in Neural Information Processing Systems*, 35:17926–17939, 2022.
- 184 [8] Christopher Wang, Vighnesh Subramaniam, Adam Uri Yaari, Gabriel Kreiman, Boris Katz,  
185 Ignacio Cases, and Andrei Barbu. Brainbert: Self-supervised representation learning for  
186 intracranial recordings. In *The Eleventh International Conference on Learning Representations*,  
187 2022.
- 188 [9] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
189 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information  
190 processing systems*, 30, 2017.
- 191 [10] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia  
192 Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international  
193 conference on computer vision*, pages 6836–6846, 2021.
- 194 [11] Andy T. Liu, Shang-Wen Li, and Hung-yi Lee. TERA: self-supervised learning of transformer  
195 encoder representation for speech. *IEEE ACM Trans. Audio Speech Lang. Process.*, 29:2351–  
196 2366, 2021. doi: 10.1109/TASLP.2021.3095662. URL [https://doi.org/10.1109/TASLP.  
197 2021.3095662](https://doi.org/10.1109/TASLP.2021.3095662).
- 198 [12] Graham Wideman. Orientation and voxel-order terminology: Ras, las, lpi, rpi, xyz and all that,  
199 2024. URL [http://www.grahamwideman.com/gw/brain/orientation/orientterms.  
200 htm](http://www.grahamwideman.com/gw/brain/orientation/orientterms.htm).
- 201 [13] Gaurav R Ghosal and Reza Abbasi-Asl. Multi-modal prototype learning for interpretable  
202 multivariable time series classification. *arXiv preprint arXiv:2106.09636*, 2021.
- 203 [14] Sabera Talukder, Yisong Yue, and Georgia Gkioxari. Totem: Tokenized time series embeddings  
204 for general time series analysis. *arXiv preprint arXiv:2402.16412*, 2024.
- 205 [15] What is aphasia? — types, causes and treatment, Mar 2017. URL [https://www.nidcd.nih.  
206 gov/health/aphasia](https://www.nidcd.nih.gov/health/aphasia).
- 207 [16] Ran Liu, Mehdi Azabou, Max Dabagia, Jingyun Xiao, and Eva Dyer. Seeing the forest and the  
208 tree: Building representations of both individual and collective dynamics with transformers.  
209 *Advances in neural information processing systems*, 35:2377–2391, 2022.

- 210 [17] Sabera J Talukder and Georgia Gkioxari. Time series modeling at scale: A universal representa-  
211 tion across tasks and domains. 2023.
- 212 [18] Hsiang-Yun Sherry Chien, Hanlin Goh, Christopher M Sandino, and Joseph Y Cheng. Maeeg:  
213 Masked auto-encoder for eeg representation learning. *arXiv preprint arXiv:2211.02625*, 2022.
- 214 [19] Demetres Kostas, Stephane Aroca-Ouellette, and Frank Rudzicz. Bendr: using transformers and  
215 a contrastive self-supervised learning task to learn from massive amounts of eeg data. *Frontiers*  
216 *in Human Neuroscience*, 15:653659, 2021.
- 217 [20] Ke Yi, Yansen Wang, Kan Ren, and Dongsheng Li. Learning topology-agnostic eeg represen-  
218 tations with geometry-aware modeling. In *Thirty-seventh Conference on Neural Information*  
219 *Processing Systems*, 2023.
- 220 [21] Armin Thomas, Christopher Ré, and Russell Poldrack. Self-supervised learning of brain  
221 dynamics from broad neuroimaging data. *Advances in Neural Information Processing Systems*,  
222 35:21255–21269, 2022.
- 223 [22] Xuan Kan, Wei Dai, Hejie Cui, Zilong Zhang, Ying Guo, and Carl Yang. Brain network  
224 transformer. *Advances in Neural Information Processing Systems*, 35:25586–25599, 2022.
- 225 [23] Josue Ortega Caro, Antonio Henrique Oliveira Fonseca, Christopher Averill, Syed A Rizvi,  
226 Matteo Rosati, James L Cross, Prateek Mittal, Emanuele Zappala, Daniel Levine, Rahul M  
227 Dhodapkar, et al. Brainlm: A foundation model for brain activity recordings. *bioRxiv*, pages  
228 2023–09, 2023.
- 229 [24] Antonis Antoniadis, Yiyi Yu, Joseph Canzano, William Wang, and Spencer LaVere Smith.  
230 Neuroformer: Multimodal and multitask generative pretraining for brain data. *arXiv preprint*  
231 *arXiv:2311.00136*, 2023.
- 232 [25] Daoze Zhang, Zhizhang Yuan, Yang Yang, Junru Chen, Jingjing Wang, and Yafeng Li. Brant:  
233 Foundation model for intracranial neural signal. *Advances in Neural Information Processing*  
234 *Systems*, 36, 2024.
- 235 [26] Chaoqi Yang, M Westover, and Jimeng Sun. Biot: Biosignal transformer for cross-data learning  
236 in the wild. *Advances in Neural Information Processing Systems*, 36, 2024.
- 237 [27] Weibang Jiang, Liming Zhao, and Bao liang Lu. Large brain model for learning generic  
238 representations with tremendous EEG data in BCI. In *The Twelfth International Conference on*  
239 *Learning Representations*, 2024. URL <https://openreview.net/forum?id=QzTpTRVtrP>.
- 240 [28] Joel Ye, Jennifer Collinger, Leila Wehbe, and Robert Gaunt. Neural data transformer 2: multi-  
241 context pretraining for neural spiking activity. *Advances in Neural Information Processing*  
242 *Systems*, 36, 2024.
- 243 [29] Donghong Cai, Junru Chen, Yang Yang, Teng Liu, and Yafeng Li. Mbrain: A multi-channel  
244 self-supervised learning framework for brain signals. In *Proceedings of the 29th ACM SIGKDD*  
245 *Conference on Knowledge Discovery and Data Mining*, KDD '23, page 130–141, New York,  
246 NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701030. doi: 10.1145/  
247 3580305.3599426. URL <https://doi.org/10.1145/3580305.3599426>.
- 248 [30] Mehdi Azabou, Vinam Arora, Venkataramana Ganesh, Ximeng Mao, Santosh Nachimuthu,  
249 Michael Mendelson, Blake Richards, Matthew Perich, Guillaume Lajoie, and Eva Dyer. A  
250 unified, scalable framework for neural population decoding. *Advances in Neural Information*  
251 *Processing Systems*, 36, 2024.
- 252 [31] Chethan Pandarinath, Daniel J O’Shea, Jasmine Collins, Rafal Jozefowicz, Sergey D Stavisky,  
253 Jonathan C Kao, Eric M Trautmann, Matthew T Kaufman, Stephen I Ryu, Leigh R Hochberg,  
254 et al. Inferring single-trial neural population dynamics using sequential auto-encoders. *Nature*  
255 *methods*, 15(10):805–815, 2018.
- 256 [32] Brianna M Karpowicz, Yahia H Ali, Lahiru N Wimalasena, Andrew R Sedler, Mohammad Reza  
257 Keshtkaran, Kevin Bodkin, Xuan Ma, Lee E Miller, and Chethan Pandarinath. Stabilizing  
258 brain-computer interfaces through alignment of latent dynamics. *BioRxiv*, pages 2022–04, 2022.

- 259 [33] Alan D Degenhart, William E Bishop, Emily R Oby, Elizabeth C Tyler-Kabara, Steven M Chase,  
260 Aaron P Batista, and Byron M Yu. Stabilization of a brain–computer interface via the alignment  
261 of low-dimensional spaces of neural activity. *Nature biomedical engineering*, 4(7):672–685,  
262 2020.
- 263 [34] Justin Jude, Matthew G Perich, Lee E Miller, and Matthias H Hennig. Robust alignment of  
264 cross-session recordings of neural population activity by behaviour via unsupervised domain  
265 adaptation. feb 2022. doi: 10.48550. *arXiv preprint arXiv:2202.06159*.
- 266 [35] Xuan Ma, Fabio Rizzoglio, Kevin L Bodkin, Eric Perreault, Lee E Miller, and Ann Kennedy.  
267 Using adversarial networks to extend brain computer interface decoding accuracy over time.  
268 *elife*, 12:e84296, 2023.
- 269 [36] Sabera Talukder, Jennifer J Sun, Matthew Leonard, Bingni W Brunton, and Yisong Yue. Deep  
270 neural imputation: A framework for recovering incomplete brain recordings. *arXiv preprint*  
271 *arXiv:2206.08094*, 2022.
- 272 [37] Xiang Zhang, Marko Zeman, Theodoros Tsiligkaridis, and Marinka Zitnik. Graph-guided  
273 network for irregularly sampled multivariate time series. *arXiv preprint arXiv:2110.05357*,  
274 2021.
- 275 [38] Geeling Chau, Yujin An, Ahamed Raffey Iqbal, Soon-Jo Chung, Yisong Yue, and Sabera  
276 Talukder. Generalizability under sensor failure: Tokenization+ transformers enable more robust  
277 latent spaces. *arXiv preprint arXiv:2402.18546*, 2024.
- 278 [39] Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli,  
279 Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for  
280 deep learning: Training BERT in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019.
- 281 [40] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*  
282 *arXiv:1711.05101*, 2017.
- 283 [41] ildoonet. ildoonet/pytorch-gradual-warmup-lr: Gradually-warmup learning rate scheduler for  
284 pytorch, 2024. URL <https://github.com/ildoonet/pytorch-gradual-warmup-lr>.
- 285 [42] Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. *arXiv preprint*  
286 *arXiv:2005.00928*, 2020.
- 287 [43] Christophe Destrieux, Bruce Fischl, Anders Dale, and Eric Halgren. Automatic parcellation  
288 of human cortical gyri and sulci using standard anatomical nomenclature. *Neuroimage*, 53(1):  
289 1–15, 2010.
- 290 [44] Nilearn, 2015. URL <https://nilearn.github.io/stable/index.html>.



291 **A Related Work**

292 **Self-supervised learning on neural data** Channel independent pretrained models are a popular  
 293 approach for neural spiking data [16], intracranial brain data [8, 17], and general time-series [14].  
 294 Additionally, in fixed-channel neural datasets, approaches exist for EEG [18–20], fMRI [21–23], and  
 295 calcium imaging [24] datasets. However, all of this work do not learn population-level interactions  
 296 across datasets with different recording layouts due to the single-channel focus or the ability to  
 297 assume fixed-channel setups. Several works pretrain spatial and temporal dimensions across datasets  
 298 with variable inputs [25–29], but most simultaneously learn the temporal embeddings with the spatial  
 299 modeling, which make them challenging to interpret and computationally expensive to train. As far  
 300 as we know, we are the first to study the problem of building pretrained channel aggregation models  
 301 on top of pre-existing temporal embeddings trained across datasets with variable sampling of input  
 302 channels, allowing for modeling of high quality (>2kHz sampling rate) intracranial data.

303 **Modeling across variable input channels** Modeling spatial representations on top of temporal  
 304 embeddings have been found to be beneficial for decoding [3, 7, 30], but prior works use supervised  
 305 labels, so do not leverage large amounts of unannotated data. The brain-computer-interface field has  
 306 been studying how to align latent spaces [31–35] which either still requires creating an alignment  
 307 matrix to learn across datasets or only provides post-training alignment mechanisms rather than  
 308 learning across datasets. Other approaches impute missing channels or learn latent spaces robust to  
 309 missing channels [36–38], but these are more suited for the occasional missing channel rather than  
 310 largely varying sensor layouts. We directly learn spatial-level representations using self-supervised  
 311 learning across datasets to leverage massive amounts of unannotated intracranial data.

312 **B Model and Compute Requirements**

	e5	e50	e90
PopT		20M	
Deep NN	3M	20M	36M
Linear	3.8k	38k	69k
Brant [25]		500M	
LaBraM [27]		350M	

Table 2: **Parameter counts.** Since PopT takes existing temporal embeddings as input, the number of parameters that must be trained is an order of magnitude less than recent end-to-end approaches.

	GPU count	GPU type	Time to train	TFLOPS
PopT	1	NVIDIA TITAN RTX (24GB)	2 days	2.1M
Brant [25]	4	NVIDIA Tesla A100 (80G)	2.8 days	18.8M
LaBraM [27]	8	NVIDIA Tesla A800 (40G)	–	–

Table 3: **Pretraining compute requirements** Based on published train times (none were given for LaBraM) it is evident that PopT has smaller hardware and shorter training time requirements.

## 313 C Architectures and training

314 **Pretrained PopT** The core Population Transformer consists of a transformer encoder stack with  
315 6 layers, 8 heads. All layers ( $N = 6$ ) in the encoder stack are set with the following parameters:  
316  $d_h = 512$ ,  $H = 8$ , and  $p_{\text{dropout}} = 0.1$ . We pretrain the PopT model with the LAMB optimizer [39]  
317 ( $lr = 1e - 4$ ), with a batch size of  $n_{\text{batch}} = 256$ , and train/val/test split of 0.98, 0.01, 0.01 of the data.  
318 We pretrain for 500,000 steps, and record the validation performance every 1,000 steps. Downstream  
319 evaluation takes place on the weights with the best validation performance. We use the intermediate  
320 representation at the [CLS] token  $d_h = 512$  and put a linear layer that outputs to  $d_{\text{out}} = 1$  for  
321 fine-tuning on downstream tasks. These parameters for pretraining were the same for any PopT that  
322 needed to be pretrained (hold-one-out subject, subject subsets, ablation studies).

323 **Non-pretrained PopT** The architecture for the non-pretrained PopT is the same as the pretrained  
324 PopT (above). However, no pretraining is done, and the weights are randomly initialized with the  
325 default initializations.

326 **Linear** The linear baseline consists of a single linear layer that outputs to  $d_{\text{out}} = 1$ . The inputs are  
327 flattened and concatenated BrainBERT embeddings  $d_{\text{emb}} = 756$  or TOTEM embeddings  $d_{\text{emb}} = 64$   
328 from a subset of channels  $S \subset C$ . Thus, the full input dimension is  $d_{\text{input}} = d_{\text{emb}} * |S|$ .

329 **Deep NN** The inputs are the same as above, but the decoding network now consists of 5 stacked  
330 linear layers, each with  $d_h = 512$  and a GeLU activation.

331 **Downstream Training** For both PopT models, we train with these parameters: AdamW optimizer  
332 [40],  $lr = 5e^{-4}$  where transformer weights are scaled down by a factor of 10 ( $lr_t = 5e^{-5}$ ),  
333  $n_{\text{batch}} = 256$ , a Ramp Up scheduler [41] with warmup 0.025 and Step LR gamma 0.99, reducing  
334 100 times within the 2000 total steps that we train for. For Linear and DeepNN models, we train with  
335 these parameters: AdamW optimizer [40],  $lr = 5e^{-4}$ ,  $n_{\text{batch}} = 256$ , a Ramp Up scheduler [41] with  
336 warmup 0.025 and Step LR gamma 0.95, reducing 25 times within the 17,000 total steps we train for.  
337 For all downstream decoding, we use a fixed train/val/test split of 0.8, 0.1, 0.1 of the data.

338 **Compute Resources** To run all our experiments (data processing, pretraining, evaluations, inter-  
339 pretability), one only needs 1 NVIDIA Titan RTXs (24GB GPU RAM). Pretraining PopT takes 2  
340 days on 1 GPU. Our downstream evaluations take a few minutes to run each. For the purposes of  
341 gathering all the results in the paper, we parallelized the experiments on roughly 8 GPUs.

## 342 D Decoding tasks

343 We follow the same task specification as in Wang et al. [8], with the modification that the pitch and  
344 volume examples are determined by percentile (see below) rather than standard deviation in order to  
345 obtain balanced classes.

346 **Pitch** The PopT receives an interval of activity and must determine if it corresponds with a high or  
347 low pitch word being spoken. For the duration of a given word, pitch was extracted using Librosa’s  
348 `piptrack` function over a Mel-spectrogram (sampling rate 48,000 Hz, FFT window length of 2048,  
349 hop length of 512, and 128 mel filters). For this task, for a given session, positive examples consist of  
350 words in the top-quartile of mean pitch and negative examples are the words in the bottom quartiles.

351 **Volume** The volume of a given word was computed as the average intensity of root-mean-square  
352 (RMS) (`rms` function, frame and hop lengths 2048 and 512 respectively). As before, positive examples  
353 are the words in the top-quartile of volume and negative examples are those in the bottom quartiles.

354 **Sentence onset** Negative examples are intervals of activity from 1s periods during which no speech is  
355 occurring in the movie. Positive examples are intervals of brain activity that correspond with hearing  
356 the first word of a sentence.

357 **Speech vs. Non-speech** Negative examples are as before. Positive examples are intervals of brain  
358 activity that correspond with dialogue being spoken in the stimuli movie.

359 **E Data**

Subj.	Age (yrs.)	# Elec- trodes	Movie	Recording time (hrs)	Held- out
1	19	91	Thor: Ragnarok	1.83	
			Fantastic Mr. Fox	1.75	
			The Martian	0.5	x
2	12	100	Venom	2.42	
			Spider-Man: Homecoming	2.42	
			Guardians of the Galaxy	2.5	
			Guardians of the Galaxy 2	3	
			Avengers: Infinity War	4.33	
			Black Panther	1.75	
			Aquaman	3.42	x
3	18	91	Cars 2	1.92	x
			Lord of the Rings 1	2.67	
			Lord of the Rings 2 (extended edition)	3.92	
4	9	135	Megamind	2.58	
			Toy Story	1.33	
			Coraline	1.83	x
5	11	205	Cars 2	1.75	x
			Megamind	1.77	
6	12	152	Incredibles	1.15	
			Shrek 3	1.68	x
			Megamind	2.43	
7	6	109	Fantastic Mr. Fox	1.5	
8	4.5	72	Sesame Street Episode	1.28	
9	16	102	Ant Man	2.28	
10	12	173	Cars 2	1.58	x
			Spider-Man: Far from Home	2.17	

Table 4: **Subject statistics** Subjects used in PopT training, and held-out downstream evaluation. Table taken from [8]. The number of uncorrupted, electrodes that can be Laplacian re-referenced are shown in the second column The average amount of recording data per subject is 4.3 (hrs).

360 **F Interpretation Methods**

361 **Candidate functional brain regions from attention weights** After fine-tuning our weights on a  
 362 decoding task, we can examine the attention weights of the [CLS] output for candidate functional  
 363 brain regions. We obtain a normalized Scaled Attention Weight metric (see next section) across all  
 364 subjects to be able to analyze candidate functional brain regions across sparsely sampled subject  
 365 datasets. The Scaled Attention Weight is computed from raw attention weights at the [CLS] token  
 366 passed through the attention rollout algorithm [42]. The resulting weights from each channel are then  
 367 grouped by brain region according to the Destrieux layout [43].

368 **Scaled Attention Weight** First, we obtain an attention weight matrix across all trials which includes  
 369 weights between all tokens. Then, we perform attention rollout [42] across layers to obtain the  
 370 contributions of each input channel by the last layer. We take the resulting last layer of rollout weights  
 371 for all channels, where the target is the [CLS] token, normalize within subject, and scale by ROC  
 372 AUC to obtain the Scaled Attention Weight per channel. Finally, we plot the 0.75 percentile weight  
 373 per region, as mapped by the Destrieux atlas [43] using Nilearn [44].

374 **G Ablation study**

375 **Ablation of loss components and position information** An ablation study confirms that both the  
 376 network-wise and channel-wise component of the pretraining objective contribute to the downstream  
 377 performance (Table 5). We also find that including the 3D position information for each channel is  
 378 critical for decoding. Additionally, we find that the discriminative nature of our loss is necessary for  
 379 decoding. Attempting to add an L1 reconstruction term to our pretraining objective results in poorer  
 380 performance, perhaps because the model learns to overfit on low-entropy features in the embedding.  
 381 Our discriminative loss requires the model to understand the embeddings in terms of how they can be  
 382 distinguished from one another, which leads the model to extract more informative representations.

	Pitch	Volume	Sent. Onset	Speech/Non-speech
PopT	<b>0.69 ± 0.07</b>	<b>0.84 ± 0.06</b>	<b>0.86 ± 0.05</b>	<b>0.89 ± 0.07</b>
PopT w/o group-wise loss	0.66 ± 0.07	0.83 ± 0.06	0.84 ± 0.04	0.88 ± 0.08
PopT w/o channel-wise loss	0.67 ± 0.06	0.81 ± 0.08	0.84 ± 0.06	0.87 ± 0.09
PopT w/o position encoding	0.59 ± 0.07	0.67 ± 0.10	0.75 ± 0.08	0.79 ± 0.08
PopT with reconstruction loss	0.60 ± 0.11	0.73 ± 0.11	0.81 ± 0.05	0.83 ± 0.09
PopT with L1 reconstruction only	0.56 ± 0.04	0.65 ± 0.08	0.73 ± 0.10	0.74 ± 0.10

Table 5: **PopT ablation study.** We individually ablate our losses and positional encodings during pretraining then decode on the resulting models. Shown are ROC-AUC mean and standard error across subjects. The best performing model across all decoding tasks uses all three of our proposed components, showing that they are all necessary. Removing our positional encoding during pretraining and fine-tuning drops the performance the most, indicating that position encoding is highly important for achieving good decoding. Additionally, we attempt adding a reconstruction component to the loss or purely using the L1 mask loss, but find that this leads to poorer performance (last two rows).