Know Thyself by Knowing Others: Learning Neuron Identity from Population Context

Vinam Arora

University of Pennsylvania Philadelphia, PA, USA vinam@upenn.edu

Mehdi Azabou

Columbia University New York, NY, USA ma4766@columbia.edu

Divyansha Lachi

University of Pennsylvania Philadelphia, PA, USA div11@upenn.edu

Blake Richards

McGill University, Mila Montréal, QC, Canada blake.richards@mila.quebec

Ian J. Knight

University of Pennsylvania Philadelphia, PA, USA ijknight@upenn.edu

Cole Hurwitz

Columbia University New York, NY, USA ch3676@columbia.edu

Joshua H. Siegle

Allen Institute for Neural Dynamics Seattle, WA, USA joshs@alleninstitute.org

Eva L. Dyer

University of Pennsylvania Philadelphia, PA, USA eva.dyer@upenn.edu

Abstract

Identifying the functional identity of individual neurons is essential for interpreting circuit dynamics, yet it remains a major challenge in large-scale in vivo recordings where anatomical and molecular labels are often unavailable. Here we introduce NuCLR, a self-supervised framework that learns context-aware representations of neuron identity by modeling each neuron's role within the broader population. NuCLR employs a spatio-temporal transformer that captures both within-neuron dynamics and across-neuron interactions. It is trained with a sample-wise contrastive objective that encourages temporally-stable and discriminative embeddings. Across multiple open-access datasets, NuCLR outperforms prior methods in both cell type and brain region classification. Critically, it exhibits strong zero-shot generalization to entirely new populations, without any retraining or access to stimulus labels. Furthermore, we demonstrate that our framework scales effectively with data size. Overall, our results demonstrate that modeling population context is crucial for understanding neuron identity and that rich signal for celltyping and neuron localization is present in neural activity alone. Code available at https://github.com/nerdslab/nuclr.

1 Introduction

Identifying the cell type and anatomical location of individual neurons is critical for understanding how neural circuits give rise to behavior [6, 35, 44], guiding closed-loop neurotechnologies [16], and tracking changes during learning or disease progression [5, 33]. While transcriptomic profiling, molecular labels, or morphological reconstructions can provide cell-type and region information post hoc [27, 32, 43], these techniques are resource intensive and rarely feasible in chronic or human recordings [14, 26, 18]. As a result, there is increasing interest in learning representations of neurons directly from their activity—representations that reflect functional identity and can support in vivo classification without external labels.

Recent work has shown that features intrinsic to each neuron—such as waveform shape, spike timing, or interspike intervals—can be used to train models that classify neurons into cell types or brain regions [28, 4, 42]. However, these approaches typically treat neurons independently, ignoring the broader population context in which each neuron operates. Yet a neuron's functional role is inherently relational, defined not just by its individual activity, but by its coordination and interaction with other neurons in the circuit. Some methods have attempted to include population-level structure [24, 40], but only via fixed, low-dimensional summaries of the surrounding activity. Additionally, the design of these methods prevents generalization to novel populations without self-supervised finetuning, which can be infeasible for in vivo clinical settings.

To address these limitations, we introduce NuCLR, a self-supervised framework that learns time-invariant, population-aware representations of individual neurons directly from neural activity. The core insight of NuCLR is that a neuron's functional identity can be inferred more reliably when considered within the context of the surrounding population and across time. To achieve this, we design a spatiotemporal transformer architecture that operates over the population of neurons recorded within a window of time, capturing both neuron-specific dynamics and across-neuron interactions. Each neuron's activity is treated as a sequence of temporal tokens, and the model uses attention mechanisms to integrate population-level structure without requiring fixed neuron ordering or session-specific alignment.

Training is guided by a sample-wise contrastive objective: representations of the same neuron observed at different times are treated as positives, while those from other neurons in the same population serve as negatives. This encourages the model to learn embeddings that are stable over time but discriminative across neurons, even under trial-to-trial variability and incomplete population observations. We additionally introduce neuron dropout during training, randomly masking neurons in each batch to promote robustness to changes in population composition.

Importantly, NuCLR is designed to generalize across sessions, animals, and experimental conditions. Unlike prior approaches that rely on stimulus labels, trial timing, or fixed population summaries, our method supports zero-shot classification in entirely new populations, enabling robust decoding of cell type or brain region identity without retraining or access to auxiliary metadata. This makes NuCLR well-suited for in vivo application where finetuning is infeasible and labeled data is sparse or unavailable.

Our core contributions are as follows:

- We present NuCLR, a self-supervised learning framework for learning neuron-level representations directly from raw neural population activity. The model captures population context and supports *zero-shot* generalization to unseen populations, without requiring retraining or behavioral labels.
- We show that NuCLR outperforms state-of-the-art supervised and self-supervised baselines on both cell type and brain region classification tasks. Our evaluations span multiple open-access datasets and test generalization across both transductive (within-session) and inductive (across-session/subject) settings.
- We investigate how performance scales with the amount of unlabeled pretraining data and label availability. Our results show that increasing the size of pretraining data consistently improves zero-shot decoding accuracy, highlighting the importance of large-scale, unsupervised representation learning for neuron identity inference.

2 Methods

The goal of this work is to design and train deep learning models that, given the activity of a population of neurons, produce a latent embedding for each neuron that captures its functional identity (i.e., cell type, anatomical location, functional tuning) within the population. Our approach is grounded in two core principles:

- 1. To correctly identify a neuron's cell type and anatomical location from activity alone, it is important not only to consider its individual activity but also the activity of the population, which provides critical contextual information.
- 2. The cell type and functional role in a circuit of a neuron remain stable over time, such that embeddings learned from population activity in one session should be transferable to new recordings.

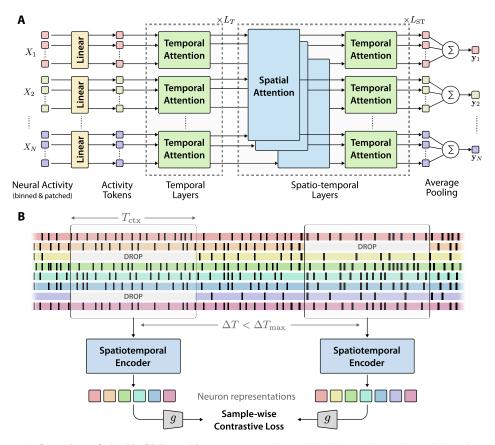


Figure 1: Overview of the NuCLR architecture. (A) The model takes as input the activity of a neural population over a fixed context window. Each neuron's activity is treated as a patched-token sequence and encoded across time using temporal transformer blocks. These temporally encoded tokens are then passed through spatio-temporal transformer layers that attends across neurons to incorporate population context. Finally, the model outputs one embedding for each neuron in the population. (B) The resulting representations are trained with a sample-wise contrastive objective, where positive pairs are derived from the same neuron across two windows of a neural recording, and negatives are drawn from other neurons in the same population. Neurons are also randomly dropped before encoding to build robustness to partial observations. The produced representations are encouraged to be stable across time and support downstream neuron-level classification tasks including zero-shot cell type and brain region decoding.

Designing models that can effectively incorporate rich population context is difficult, though, because the number of recorded neurons can vary widely across samples—across experiments, individuals, and sessions. This makes it non-trivial to define a consistent input structure or population-level operation. As well, recordings typically span experimental conditions, including varying stimuli and behaviors, which makes it hard for a model to generalize across recording sessions. These sources of variability can make it challenging to learn general embeddings that transfer to new recordings.

We address these sources of variability through a principled design of both the model architecture (Section 2.1) and the training objective (Section 2.2) that we describe below.

2.1 Model architecture

We aim to learn a function that maps the collective activity of a neural population to individual neuron-level embeddings that accurately reflect those neurons cell types and functional roles within the circuit. To support learning across recordings from different animals—each with a varying number of neurons—this function must be applicable to populations of arbitrary size. Moreover, we do not assume any relational bias between neurons, and treat them as a permutation-equivariant set where all neuron-neuron interactions are modeled. Finally, to enable embeddings that reflect a neuron's role within the circuit, the function should allow for information exchange among neurons, allowing

each one to contextualize its activity relative to the population. We formalize this as a set-to-set transformation:

$$\left\{\mathbf{y}_{1}, \mathbf{y}_{2}, \dots, \mathbf{y}_{N}\right\} = \mathcal{F}\left(\left\{X_{1}, X_{2}, \dots, X_{N}\right\}\right),\tag{1}$$

where X_n denotes the spike-train sequence of the n-th neuron, and y_n is its corresponding embedding. We implement \mathcal{F} using attention-based transformer blocks (Figure 1A), which naturally support set-to-set mappings through their permutation-equivariant structure [20].

Given spike trains of a population of N neurons over a temporal context window T_{ctx} , we begin by binning the spike train¹, and partitioning the bins into non-overlapping temporal patches of length T_{patch} . The resulting binned-and-patched activity of n-th neuron is a sequence of vectors $X_n = (\mathbf{x}_{n,1}, \mathbf{x}_{n,2}, \dots, \mathbf{x}_{n,P})$, where $P = T_{\text{ctx}}/T_{\text{patch}}$ is the number of patches. Each patch is linearly projected into a D-dimensional latent space, yielding a token sequence $Z_n^{(0)} = (\mathbf{z}_{n,1}^{(0)}, \dots, \mathbf{z}_{n,P}^{(0)})$ for each neuron.

These latent tokens are first processed independently per neuron using L_T layers of self-attention, which we refer to as **temporal transformer layers**. Each temporal layer \mathcal{T}_{temp} operates on the patch sequence of a single neuron:

$$\left(\mathbf{z}_{n,1}^{(l+1)}, \dots, \mathbf{z}_{n,P}^{(l+1)}\right) = \mathcal{T}_{\text{temp}}\left(\left(\mathbf{z}_{n,1}^{(l)}, \dots, \mathbf{z}_{n,P}^{(l)}\right)\right), \quad \forall n \in [N]. \tag{2}$$

Within these layers, the temporal structure is maintained by the use of rotary position embeddings [31, 1], which encode the relative timing of each patch without requiring absolute positional indices.

Following the stack of temporal layers, we apply $L_{\rm ST}$ layers of **spatio-temporal transformer layers**, which alternate between spatial and temporal attention (Figure 1A). In the spatial transformer layer, tokens at the same time index interact across the population via a shared transformer $\mathcal{T}_{\rm spatial}$:

$$\left\{\mathbf{z}_{1,p}^{(l+1)}, \dots, \mathbf{z}_{N,p}^{(l+1)}\right\} = \mathcal{T}_{\text{spatial}}\left(\left\{\mathbf{z}_{1,p}^{(l)}, \dots, \mathbf{z}_{N,p}^{(l)}\right\}\right) \quad \forall \ p \in [P]. \tag{3}$$

To maintain permutation-invariance across the population, no positional embeddings are used in the spatial blocks. The temporal blocks in these layers reuse the structure defined in Equation (2).

Finally, to obtain a fixed-dimensional embedding for each neuron, we apply mean pooling over the temporal axis:

$$\mathbf{y}_n = \frac{1}{P} \sum_{p=1}^{P} \mathbf{z}_{n,p}^{\text{final}} \quad \forall \ n \in [N].$$
 (4)

The core idea behind this architecture is to first build a temporally informed representation of each neuron based solely on its own activity. This is accomplished through the initial stack of temporal transformer layers, which process each neuron's token sequence independently. The resulting latent sequence captures the internal dynamics of each neuron over time ². Next, the spatial transformer layers allow neurons to exchange information at each timepoint, injecting population-level context into each neuron's representation. By alternating spatial and temporal layers, information received from other neurons at a given timepoint can be distributed to all other timepoints, creating more informed tokens to query the population again through the next spatial layer.

This architecture can also be readily adapted for calcium imaging data, with minimal modifications (see Appendix B.3). Additional details are provided in Appendix B.

What is considered a *population*? For electrophysiology datasets, we treat each probe insertion as a distinct population—both for the spatial transformer layers and for the contrastive loss, where negative pairs are limited to neurons recorded on the same insertion (i.e. a single Neuropixels [29] probe). In early experiments, we observed that allowing interactions across insertions led the model to cluster neurons based on probe identity rather than biologically meaningful properties. This is undesirable, as our goal is to produce embeddings that reflect intrinsic neuronal attributes such as

¹We use a fixed bin-size of 20ms, and also present a sweep over this value in Appendix A.1. We also attempted a spike-tokenization based approach, and discuss it in Appendix A.2.

²The initial temporal-only layers, while not essential for performance, contribute to improved computational efficiency. Please see Appendix A.5 for a longer discussion.

brain region and cell type—not experimental artifacts. Restricting spatial interactions to withininsertion populations also yields computational benefits, since the spatial transformer's compute cost scales quadratically with the number of neurons in a population.

2.2 Self-supervised training objective

Our goal is to learn the transformation \mathcal{F} in Equation (1) such that it produces neuron-level embeddings that capture biologically meaningful properties such as cell type and brain region. Since such labels are costly and difficult to obtain at scale, we adopt a self-supervised approach. NuCLR leverages natural spatial and temporal structure in neural population activity and applies a contrastive learning objective to train \mathcal{F} without requiring any explicit labels.

We begin by sampling two $T_{\rm ctx}$ -long windows (or "views") of neural population activity within $\Delta T_{\rm max}$ of each other (Figure 1B). Because neuron-level properties such as cell type or brain region remain stable over time, we can encode this invariance using a contrastive loss that encourages corresponding neurons across the two views to produce similar representations.

To promote robustness to partial observations and prevent overfitting to specific populations, we apply *neuron dropout* independently to each view—randomly removing *up to* 50% neurons per view. As a result, the number of neurons present in each view usually differs, and only a subset of neurons in one view have a corresponding neuron in the other view. We define the indices of neurons presented in both views as

$$\mathcal{M} = \{(n, m) \mid \text{neuron } n \text{ in view 1 corresponds to neuron } m \text{ in view 2}\},\tag{5}$$

which identifies all valid positive pairs for contrastive training.

Let $\tilde{\mathcal{X}}^1 = \{X_1^1, \dots, X_{N_1}^1\}$ and $\tilde{\mathcal{X}}^2 = \{X_1^2, \dots, X_{N_2}^2\}$ denote the two augmented views after neuron dropout. The encoder maps these views into corresponding sets of embeddings:

$$\mathcal{F}(\tilde{\mathcal{X}}^1) = \{ \mathbf{y}_1^1, \dots, \mathbf{y}_{N_1}^1 \} \text{ and } \mathcal{F}(\tilde{\mathcal{X}}^2) = \{ \mathbf{y}_1^2, \dots, \mathbf{y}_{N_2}^2 \}.$$
 (6)

Similar to SimCLR [9], these embeddings are passed through a projection head $g(\cdot)$ —a one-hidden-layer MLP—to obtain projected vectors $\mathbf{p}_n^v = g(\mathbf{y}_n^v)$, where $v \in \{1,2\}$ indexes the view. We then apply an InfoNCE-based contrastive loss over the set of valid positive pairs \mathcal{M} :

$$\mathcal{L}(\mathcal{F}, g \mid \tilde{\mathcal{X}}^{1}, \tilde{\mathcal{X}}^{2}) = \sum_{(n,m)\in\mathcal{M}} -\log\left(\frac{\exp\left(\langle \mathbf{p}_{n}^{1}, \mathbf{p}_{m}^{2} \rangle / \tau\right)}{\sum_{n'\neq n} \exp\left(\langle \mathbf{p}_{n}^{1}, \mathbf{p}_{n'}^{1} \rangle / \tau\right) + \sum_{(n,k)\notin\mathcal{M}} \exp\left(\langle \mathbf{p}_{n}^{1}, \mathbf{p}_{k}^{2} \rangle / \tau\right)}\right) + \text{symmetric term for view 2 to 1}$$

where $\langle \cdot, \cdot \rangle$ denote cosine-similarity, and τ is a temperature hyperparameter. This loss encourages embeddings of the same neuron to be close across views while pushing apart those of different neurons within the same population. As a result, the model learns neuron-level representations that are temporally stable and discriminative with respect to their functional roles.

A key distinction from standard SimCLR [9] lies in how we handle the minibatch setting. In NuCLR, the contrastive loss is computed *independently* within each view pair. That is, when training on minibatches containing samples from different animals or sessions, we do not treat neurons across samples as negatives. This design choice avoids an overabundance of *easy negatives* in the denominator of Equation (7), which can degrade contrastive learning performance [15]. As a result, the number of negatives per sample is limited to the neurons within a single recording (typically in the hundreds). This motivates our omission of positive-pair similarity terms from the denominator—following the decoupled contrastive loss (DCL) [41]—which has been shown to improve performance in regimes with few negatives.

In practice, for a minibatch containing B independently sampled view pairs (potentially from different recordings), the overall loss is computed as a weighted average across samples:

$$\mathcal{L}_{\text{batch}}(\mathcal{F}, g) = \frac{1}{\sum_{b=1}^{B} N_b} \sum_{b=1}^{B} N_b \cdot \mathcal{L}(\mathcal{F}, g | \tilde{\mathcal{X}}_b^1, \tilde{\mathcal{X}}_b^2), \tag{8}$$

where $\mathcal{X}_b^{1,2}$ is the *b*-th view pair, and N_b is the number of valid positive pairs (i.e., neurons present in both views after dropout). The inclusion of the N_b terms helps to deal with the imbalance otherwise created between views with very different numbers of valid positive pairs. We provide additional details about the model, hyperparameter choices, and other training details in Appendix B.

3 Results

We pre-train our spatio-temporal transformer (Section 2.1) using our self-supervised objective (Section 2.2) on a variety of open-access neural datasets spanning both electro- and optophysiology, and then evaluate the resulting neuron representations (averaged across the recordings) on two tasks: decoding the *cell type* and *brain region* of neurons from activity alone.

Evaluation strategies. For each task, we assess performance with a linear classifier across three generalization settings: (1) *Transductive*, where the testing populations are seen during self-supervised pretraining, and partial labels from these populations are used to train the classification head; (2) *Transductive zero-shot*, where the test populations are present during pretraining, but no labels from them are used when training the classifier; and (3) *Inductive zero-shot*, where the test populations are entirely unseen during pretraining, and no fine-tuning is performed on the encoder or classifier before evaluation. For all three settings, we use a linear head on the output embeddings for the classification probe.

Settings (1) and (2) emulate *offline* (or post-hoc) analysis of neural recordings, where pretraining is allowed on the test populations. Setting (1) reflects a scenario where a subset of neurons in a recording is labeled, and the goal is to infer labels for the rest. Setting (2) captures the case where the test population is collected but is not labeled, so it used for self-supervised pretraining but cannot be used to train the classification probe. Setting (3) corresponds to the *online* use-case, where a pretrained model must operate on completely new neural populations without any retuning—an "out-of-the-box" deployment scenario that reflects the intended practical use of neuron-type classifiers. This final setting is the most stringent and directly tests the model's ability to generalize in a truely out-of-the-box fashion.

Baseline methods. We compare NuCLR to several baselines: **NeuPRINT** [25], a population-context model that uses summary statistics of neighboring neurons and the recorded behavior; **NEMO** [42], a CLIP-style contrastive approach that encodes individual neurons using their activity autocorrelograms and waveform templates; and **LOLCAT** [28], a supervised model that focuses on temporal features of individual neurons and leverages trial structure in the recording. We refer the reader to Appendix E for further details on our implementation of these methods. Our main metric is the **macro F1-score**, which is robust to class imbalance and enables fair comparison across datasets with varying label distributions ³.

We note that not all methods can be tested across all settings: NeuPRINT learns neuron embeddings through back-propagation so it cannot be tested inductively without re-training the model. NEMO is designed only for electrophysiology data so it cannot be applied to Bugeon et al. which consists of calcium imaging recordings. Finally LOLCAT is a supervised method so the transductive zero-shot setting is not possible. The performance of baseline methods in such cases will be reported as "N/A".

3.1 Cell type decoding

To assess the quality of NuCLR's neuron representations, we evaluate performance on two datasets with ground-truth cell type annotations. Our primary goal is to test how well NuCLR supports generalization across animals and sessions, particularly in zero-shot settings where no labeled data from the target population is used for training.

We first use the Allen Brain Observatory Visual Coding (VC) Neuropixels dataset, which contains 58 sessions from unique mice. Of these, 16 sessions include optotagged inhibitory neurons labeled as one of 3 subclasses: Pvalb, Sst, or Vip. Since all labeled neurons within a session belong to the

³We attempted to compare with a recent method NeurPIR [40] but were not able to receive code from the authors to reproduce their method.

same subclass, we evaluate performance only in the two zero-shot settings ⁴. In the *transductive zero-shot* setting, all sessions are used during pretraining (including labeled ones), but the decoder is trained and tested on neurons from non-overlapping subsets of mice. In the more stringent *inductive zero-shot* setting, pretraining is restricted to unlabeled sessions only. Further details on our data splits and validation methodology are provided in Appendix D.

We also evaluate on the Bugeon et. al. [7] spatial transcriptomics dataset, which consists of calcium imaging recordings from 17 sessions across 4 mice, with cell types labeled as excitatory (E) or inhibitory (I), and inhibitory neurons further divided into 5 subclasses (Lamp5, Pvalb, Vip, Sncg, Sst). We test both binary classification (E vs. I) and five-way subclass classification.

As shown in Table 1, NuCLR achieves strong performance across both datasets and consistently outperforms all baselines in zero-shot settings. On the Allen VC dataset, NuCLR achieves a macro F1-score of 0.7218 in the transductive zero-shot setting and 0.7200 in the inductive setting—more than 0.29 F1 higher than the next best method (NEMO) in the inductive case. On the Bugeon dataset, NuCLR achieves a macro F1 of 0.701 on inductive zero-shot E vs. I classification and 0.444 on the more challenging five-way subclass task on a held-out subject. These results demonstrate that NuCLR produces stable, transferable neuron representations that generalize to previously unseen populations without retraining, enabling accurate zero-shot decoding of cell type identity across diverse datasets and experimental conditions.

Table 1: Macro F1-score for cell type classification across different generalization settings. Reported as mean \pm std. dev. across 5 training seeds. N/A indicates the method cannot operate in that evaluation setting.

Dataset	# Classes	Setting	NuCLR	NeuPRINT	NEMO	LOLCAT
Allen VC	3	Transductive zero-shot Inductive zero-shot	$\begin{array}{c} 0.7218 \pm 0.0113 \\ 0.7200 \pm 0.0267 \end{array}$	$0.4020 \pm 0.0238 \\ N/A$	$\begin{array}{c} 0.4256 \pm 0.0114 \\ 0.4194 \pm 0.0099 \end{array}$	$N/A = 0.4121 \pm 0.0800$
Bugeon et. al. (E vs. I)	2	Transductive Transductive zero-shot Inductive zero-shot	$\begin{array}{c} 0.8110 \pm 0.0035 \\ 0.6826 \pm 0.0293 \\ 0.6738 \pm 0.0563 \end{array}$	$\begin{array}{c} 0.6658 \pm 0.0090 \\ 0.6362 \pm 0.0073 \\ \text{N/A} \end{array}$	N/A N/A N/A	$\begin{array}{c} 0.7205 \pm 0.0127 \\ N/A \\ 0.7463 \pm 0.0095 \end{array}$
Bugeon et. al. (Subclass)	5	Transductive Transductive zero-shot Inductive zero-shot	$\begin{array}{c} 0.6101 \pm 0.0249 \\ 0.4014 \pm 0.0268 \\ 0.3938 \pm 0.0556 \end{array}$	$\begin{array}{c} 0.4952 \pm 0.0222 \\ 0.3529 \pm 0.0194 \\ \text{N/A} \end{array}$	N/A N/A N/A	$\begin{array}{c} 0.2900 \pm 0.0388 \\ N/A \\ 0.2418 \pm 0.0078 \end{array}$

Please refer to Appendix A.6 for confusion matrices. We present performance comparisons on an 11-class label set of Bugeon et. al. dataset in Appendix A.7. We also compared against POYO+ [2], a recent large-scale multi-task decoding method that learns neuron-level embeddings, and found that NuCLR outperforms these embeddings by a strong margin (Appendix A.4).

3.2 Brain region decoding

For brain region identification, we evaluate on two electrophysiological datasets with well-curated anatomical annotations: the International Brain Laboratory (IBL) Brain-wide Map [13] and the Steinmetz et. al. 2019 dataset [30]. The IBL dataset comprises recordings from 139 mice across approximately 700 probe insertions. Following the evaluation setup used by NEMO, we perform classification across 10 brain regions (Figure 2A). The Steinmetz et. al. dataset includes 39 recordings from 10 mice, with classification over 4 regions: HPF, MB, TH, and VIS. Zero-shot test populations correspond to entirely unseen experimental sessions (and thus unseen probe insertions) in the IBL dataset, and to unseen subjects in the Steinmetz et. al. dataset. Further details on our evaluation methodology and data folds are provided in Appendix D.

As shown in Table 2, NuCLR achieves the highest macro F1-scores across all evaluation settings. Notably, in the inductive zero-shot regime, NuCLR outperforms NEMO by 0.15 F1 on the IBL dataset. As we will see in Section 3.4, this performance gain can be majorly attributed to the presence of the spatial-attention layers in our model. We can also see in Figure 2A that the embeddings produced by NuCLR indeed are organized based on the brain region of neurons for the IBL dataset. We provide embedding visualization analysis for the remaining datasets in Appendix C.

⁴The transductive setting cannot allow for proper testing in this case, as a "perfect" classifier only has to infer which of the previously seen sessions does the neuron belong to and the cell-type corresponding to that session.

Table 2: Macro F1-scores for brain region classification across different generalization settings. Reported as mean \pm std. dev. across 5 training seeds. N/A indicates the method cannot operate in that evaluation setting.

Dataset	# Classes	Setting	NuCLR	NeuPRINT	NEMO	LOLCAT
IBL	10	Transductive Transductive zero-shot Inductive zero-shot	$\begin{array}{c} 0.6686 \pm 0.0034 \\ 0.5343 \pm 0.0115 \\ 0.5295 \pm 0.0040 \end{array}$	$\begin{array}{c} 0.2734 \pm 0.0153 \\ 0.2531 \pm 0.0137 \\ \text{N/A} \end{array}$	$\begin{array}{c} 0.4188 \pm 0.0041 \\ 0.3804 \pm 0.0011 \\ 0.3793 \pm 0.0011 \end{array}$	$0.2851 \pm 0.0008 \\ N/A \\ 0.2532 \pm 0.0016$
Steinmetz et. al.	4	Transductive Transductive zero-shot Inductive zero-shot	$\begin{array}{c} 0.9594 \pm 0.0027 \\ 0.7338 \pm 0.0226 \\ 0.5810 \pm 0.0110 \end{array}$	$0.4476 \pm 0.0166 \\ 0.4122 \pm 0.0326 \\ \text{N/A}$	$\begin{array}{c} 0.6989 \pm 0.0044 \\ 0.6681 \pm 0.0016 \\ 0.5595 \pm 0.0048 \end{array}$	$0.3205 \pm 0.0055 \\ N/A \\ 0.3191 \pm 0.0311$

3.3 Data scaling and label efficiency

We study how zero-shot performance scales with the amount of unlabeled pretraining data and the availability of labeled neurons for training the classification probe (Figure 2B). As one would expect, classification accuracy increases with more labeled data at a fixed pretraining data scale. We simulate varying supervision levels by subsampling the labeled neurons used to train the linear classifier (shown as "label ratio").

Interestingly, we find that for the Allen VC dataset, increasing the number of unlabeled pretraining sessions leads to a dramatic improvement in cell type classification performance. In fact, in some cases, doubling the amount of pretraining data is significantly more beneficial than doubling the number of labeled neurons, as annotated in Figure 2B. For the IBL dataset, we observe similar trends: brain region decoding performance continues to improve with additional pretraining data, although the gains are more modest. This is likely due to the already large size of IBL, suggesting that in this range (250–350 sessions), increasing labeled data may offer greater benefits.

Notably, even with only 25% of the labeled data, NuCLR significantly outperforms all other baselines with 100% of training examples on both IBL and Allen VC datasets. This demonstrates that our approach is both label-efficient and highly scalable, capable of leveraging large unlabeled datasets to improve downstream performance with minimal supervision.

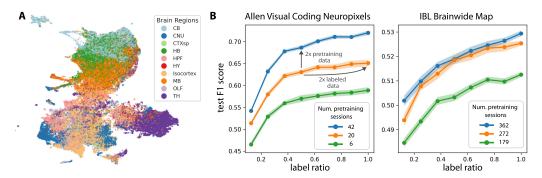


Figure 2: **Scaling the amount of pretraining and labeled data.** (A) A 2-D UMAP visualization of IBL neurons colored by brain region. (B) Inductive zero-shot cell type and brain region decoding performance improves with more pretraining data. In many cases, increasing the amount of unlabeled pretraining data is more effective than increasing the amount of labeled data.

3.4 Ablations

We performed a series of ablations to evaluate the importance of different architectural and training components in our model (Table 3). First, removing the spatial attention layers, which enable population-level interactions across neurons, resulted in the large performance drops on both Allen VC and IBL datasets. In this ablation, we *replaced* the spatial layers with temporal layers, maintaining the parameter count and depth of the model. This highlights the importance of spatial context: incorporating information from other neurons in the population provides complementary signals to a neuron's intrinsic activity, improving the model's ability to distinguish neuron types and functional roles.

We also evaluated the effect of neuron dropout, a regularization technique that randomly masks a subset of neurons during training. This component led to a clear improvement on the Allen Visual Coding dataset but had minimal effect on IBL. We attribute this to dataset scale: the Allen VC training set contains only 42 unique populations, making it more sensitive to regularization, whereas IBL includes over 600 populations and is less prone to overfitting.

Table 3: Ablation study on the Allen Visual Coding and IBL datasets. We test the impact of spatial attention layers and neuron dropout on the Macro-F1 score. Values are reported as mean \pm std. dev. across 5 seeds.

Model Variant	Allen VC	IBL
Full NuCLR w/o neuron dropout w/o spatial attention layers	$\begin{array}{c} 0.7200 \pm 0.0267 \\ 0.6181 \pm 0.0252 \\ 0.5550 \pm 0.0796 \end{array}$	$\begin{array}{c} 0.5295 \pm 0.0040 \\ 0.5248 \pm 0.0164 \\ 0.3573 \pm 0.0019 \end{array}$

4 Related Work

Cell type classification. Cell type classification seeks to assign neurons to meaningful biological or functional classes using structural, molecular, or physiological information [23, 43]. Transcriptomic approaches such as single-cell RNA sequencing [32] and spatial transcriptomics [7] provide high-resolution cell type labels, but these methods require extensive experimental infrastructure and are difficult to apply in chronic, large-scale, or in vivo recordings. Morphology-based approaches, such as self-supervised graph learning on neuronal reconstructions [38], and multimodal in vitro studies that integrate morphology, electrophysiology, and transcriptomics [11], provide complementary views of neuronal identity.

A major goal in recent work has been to infer functional cell types directly from large-scale physiological recordings, bypassing the need for extensive molecular and morphological profiling. Early approaches relied on hand-engineered features derived from electrophysiological signals, such as extracellular waveforms, autocorrelograms (ACGs), and peri-stimulus histograms as proxies for cell identity [10]. Subsequent methods incorporated stimulus-driven tuning properties, either by extracting features from responses to specific stimulus sets or by synthesizing maximally exciting or most discriminative inputs for classification [8, 37]. Large-scale functional characterization studies have further revealed substantial neuron-level diversity relevant to defining functional cell types, including work on tuning organization in primate V4 [39], chromatic feature detectors in retina [12], feature landscapes in visual cortex [34], and combinatorial codes in mouse V1 [36]. These studies show that functional cell types can be distinguished not only by static structural or molecular signatures, but also by patterns of tuning, feature selectivity, and population codes. While effective in controlled paradigms, stimulus-dependent approaches are inherently tied to the availability of specific stimuli and may fail to generalize across diverse experimental contexts.

Recent advances have shifted toward learning stimulus-agnostic, neuron-level embeddings that capture intrinsic activity dynamics. LOLCAT [28] learns trial-level representations and attends to subsets of trials to build a prediction of cell type over many trials, training in a supervised manner to classify individual neurons. NEMO [42] uses a CLIP-style contrastive loss between waveform and auto-correlogram views, while PhysMAP [19] and VAE-based models [4] combine multiple physiological signals into shared latent spaces.

NeuPRINT [24] and NeurPIR [40] aim to learn time-invariant or intrinsic representations for individual neurons from population dynamics using reconstruction and VICReg objectives [3], respectively. However, these methods include population structure only via fixed, low-dimensional summaries of the surrounding activity.

Channel-level transformer architectures and functional embeddings. Channel-level transformer architectures and related models such as POYO+ [2], POYO [1], NEDS [45], EIT [21], and STNDT [17] also generate neuron-level embeddings or tokens, but are not explicitly designed for neuron-level cell type or brain region readouts. POYO, POYO+, and EIT are trained primarily on supervised decoding tasks, while NEDS uses both encoding and decoding objectives. STNDT is trained with a masked modeling objective and uses a combination of neuron-level and population-level tokens to demonstrate strong performance on behavioral reach decoding tasks. We note that while POYO+

provides results for classification of different Cre-lines and brain regions, it is performed at the session level on averaged latent representations, rather than directly on embeddings of individual neurons.

5 Discussion

In this work, we introduced a self-supervised framework for learning population-aware, temporally-stable representations of individual neurons from spiking activity alone. By modeling both spatial and temporal relationships in population activity using a spatio-temporal transformer, and training with a sample-wise contrastive loss, our method produces biologically-meaningful representations that are generalizable across sessions and animals, and transferable to downstream decoding tasks. Importantly, our approach enables *zero-shot* decoding of cell types and brain regions on novel subjects, outperforming prior methods that rely on supervised labels or session-specific tuning.

A central insight of our model is that neuron identity (e.g. a neuron's cell type and/or brain region) is not only reflected in its intrinsic firing properties, but also in how it interacts with the surrounding population. Prior models either completely ignore this context or compress it into less informative low-dimensional features. Our results show that attention-based architectures can effectively leverage this population context, even when the number and identity of co-recorded neurons varies across recordings.

Our scaling analysis further demonstrates that the model continues to improve as more unlabeled data is introduced, and that it performs reliably even with limited labeled data. This data scaling trend makes NuCLR an ideal candidate for pretraining on multiple datasets and modalities for further improving the downstream classification performance.

While our approach demonstrates strong generalization across datasets and tasks, it has two notable limitations. First, when working with electrophysiological data, the model operates on binned neural activity and does not exploit the full temporal precision available in raw spike times, potentially discarding fine-grained timing information. Second, it requires modality-specific hyperparameters as calcium imaging and electrophysiology differ in timescales and signal properties, making it difficult to train a unified model across modalities.

More broadly, our work demonstrates the feasibility of building generalizable representations of neural identity from activity alone. This opens new opportunities for functional cell typing in the absence of ground truth labels, improving data harmonization across labs, and making progress towards achieving robust neuron classification in experimental or clinical settings. Future directions include integrating multi-modal inputs, extending to larger-scale populations, and applying these embeddings in closed-loop or adaptive neurotechnologies.

6 Acknowledgments and Disclosure of Funding

Thanks to Shivashriganesh P. Mahato for insightful discussions and feedback on this manuscript. This project was supported by NIH award 1R01EB029852-01, NSF award IIS-2146072 as well as generous gifts from the CIFAR Azrieli Global Scholars Program and The Hypothesis Fund.

References

- [1] M. Azabou, V. Arora, V. Ganesh, X. Mao, S. B. Nachimuthu, M. J. Mendelson, B. A. Richards, M. G. Perich, G. Lajoie, and E. L. Dyer. A unified, scalable framework for neural population decoding. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [2] M. Azabou, K. X. Pan, V. Arora, I. J. Knight, E. L. Dyer, and B. A. Richards. Multi-session, multi-task neural decoding from distinct cell-types and brain regions. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [3] A. Bardes, J. Ponce, and Y. LeCun. VICReg: Variance-invariance-covariance regularization for self-supervised learning. In *International Conference on Learning Representations*, 2022.
- [4] M. Beau, D. J. Herzfeld, F. Naveros, M. E. Hemelt, F. D'Agostino, M. Oostland, A. Sánchez-López, Y. Y. Chung, M. Maibach, H. N. Stabb, et al. A deep-learning strategy to identify cell types across species from high-density extracellular recordings. *Cell*, 2025.

- [5] H. Braak, K. del Tredici, C. Schultz, and E. Braak. Vulnerability of select neuronal types to alzheimer's disease. *Annals of the New York Academy of Sciences*, 924(1):53–61, 2000.
- [6] S. P. Brown and S. Hestrin. Cell-type identity: a key to unlocking the function of neocortical circuits. *Current opinion in neurobiology*, 19(4):415–421, 2009.
- [7] S. Bugeon, J. Duffield, M. Dipoppa, A. Ritoux, I. Prankerd, D. Nicoloutsopoulos, D. Orme, M. Shinn, H. Peng, H. Forrest, A. Viduolyte, C. B. Reddy, Y. Isogai, M. Carandini, and K. D. Harris. A transcriptomic axis predicts state modulation of cortical interneurons. *Nature*, 607(7918):330–338, July 2022.
- [8] M. F. Burg, T. Zenkel, M. Vystrčilová, J. Oesterle, L. Höfling, K. F. Willeke, J. Lause, S. Müller, P. G. Fahey, Z. Ding, K. Restivo, S. Sridhar, T. Gollisch, P. Berens, A. S. Tolias, T. Euler, M. Bethge, and A. S. Ecker. Most discriminative stimuli for functional cell type clustering. In *The Twelfth International Conference on Learning Representations*, 2024.
- [9] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PmLR, 2020.
- [10] S. E. J. de Vries, J. A. Lecoq, M. A. Buice, P. A. Groblewski, G. K. Ocker, M. Oliver, D. Feng, N. Cain, P. Ledochowitsch, D. Millman, and et al. A large-scale standardized physiological survey reveals functional organization of the mouse visual cortex. *Nature Neuroscience*, 23(1):138–151, Dec. 2019.
- [11] N. W. Gouwens, S. A. Sorensen, F. Baftizadeh, A. Budzillo, B. R. Lee, T. Jarsky, L. Alfiler, K. Baker, E. Barkan, K. Berry, and et al. Integrated morphoelectric and transcriptomic classification of cortical gabaergic cells. *Cell*, 183(4):935–953.e19, Nov. 2020.
- [12] L. Höfling, K. P. Szatko, C. Behrens, Y. Deng, Y. Qiu, D. A. Klindt, Z. Jessen, G. W. Schwartz, M. Bethge, P. Berens, K. Franke, A. S. Ecker, and T. Euler. A chromatic feature detector in the retina signals visual context changes. *eLife*, 13, Oct. 2024.
- [13] IBL, B. Benson, J. Benson, D. Birman, N. Bonacchi, M. Carandini, J. A. Catarino, G. A. Chapuis, A. K. Churchland, Y. Dan, et al. A brain-wide map of neural activity during complex behaviour. *bioRxiv*, pages 2023–07, 2023.
- [14] A. L. Juavinett, G. Bekheet, and A. K. Churchland. Chronically implanted neuropixels probes enable high-yield recordings in freely moving mice. *Elife*, 8:e47188, 2019.
- [15] Y. Kalantidis, M. B. Sariyildiz, N. Pion, P. Weinzaepfel, and D. Larlus. Hard negative mixing for contrastive learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 21798–21809. Curran Associates, Inc., 2020.
- [16] E. Krook-Magnuson, J. N. Gelinas, I. Soltesz, and G. Buzsáki. Neuroelectronics and biooptics: closed-loop technologies in neurological disorders. *JAMA neurology*, 72(7):823–829, 2015.
- [17] T. Le and E. Shlizerman. STNDT: Modeling neural population activity with spatiotemporal transformers. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [18] A. T. Lee, E. F. Chang, M. F. Paredes, and T. J. Nowakowski. Large-scale neurophysiology and single-cell profiling in human neuroscience. *Nature*, 630(8017):587–595, 2024.
- [19] E. K. Lee, A. E. Gül, G. Heller, A. Lakunina, S. Jaramillo, P. F. Przytycki, and C. Chandrasekaran. Physmap-interpretable in vivo neuronal cell type identification using multi-modal analysis of electrophysiological data. *BioRxiv*, pages 2024–02, 2024.
- [20] J. Lee, Y. Lee, J. Kim, A. Kosiorek, S. Choi, and Y. W. Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *Proceedings of the 36th International Conference on Machine Learning*, pages 3744–3753, 2019.

- [21] R. Liu, M. Azabou, M. Dabagia, J. Xiao, and E. L. Dyer. Seeing the forest and the tree: Building representations of both individual and collective dynamics with transformers. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [22] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.
- [23] R. H. Masland. Neuronal cell types. Current Biology, 14(13):R497–R500, 2004.
- [24] L. Mi, T. Le, T. He, E. Shlizerman, and U. Sümbül. Learning time-invariant representations for individual neurons from population dynamics. *Advances in Neural Information Processing* Systems, 36, 2023.
- [25] L. Mi, T. Le, T. He, E. Shlizerman, and U. Sümbül. Learning time-invariant representations for individual neurons from population dynamics. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 46007–46026. Curran Associates, Inc., 2023.
- [26] C. P. Mosher, Y. Wei, J. Kamiński, A. Nandi, A. N. Mamelak, C. A. Anastassiou, and U. Rutishauser. Cellular classes in the human brain revealed in vivo by heartbeat-related modulation of the extracellular action potential waveform. *Cell reports*, 30(10):3536–3551, 2020.
- [27] J.-F. Poulin, B. Tasic, J. Hjerling-Leffler, J. M. Trimarchi, and R. Awatramani. Disentangling neural cell diversity using single-cell transcriptomics. *Nature neuroscience*, 19(9):1131–1141, 2016.
- [28] A. Schneider, M. Azabou, L. McDougall-Vigier, D. F. Parks, S. Ensley, K. Bhaskaran-Nair, T. Nowakowski, E. L. Dyer, and K. B. Hengen. Transcriptomic cell type structures in vivo neuronal activity across multiple timescales. *Cell reports*, 42(4), 2023.
- [29] N. A. Steinmetz, C. Aydin, A. Lebedeva, M. Okun, M. Pachitariu, M. Bauza, M. Beau, J. Bhagat, C. Böhm, M. Broux, S. Chen, J. Colonell, et al. Neuropixels 2.0: A miniaturized high-density probe for stable, long-term brain recordings. *Science*, 372(6539), Apr. 2021.
- [30] N. A. Steinmetz, P. Zatka-Haas, M. Carandini, and K. D. Harris. Distributed coding of choice, action and engagement across the mouse brain. *Nature*, 576(7786):266–273, Nov. 2019.
- [31] J. Su, Y. Lu, S. Pan, A. Murtadha, B. Wen, and Y. Liu. Roformer: Enhanced transformer with rotary position embedding, 2023.
- [32] B. Tasic, Z. Yao, L. T. Graybuck, K. A. Smith, T. N. Nguyen, D. Bertagnolli, J. Goldy, E. Garren, M. N. Economo, S. Viswanathan, et al. Shared and distinct transcriptomic cell types across neocortical areas. *Nature*, 563(7729):72–78, 2018.
- [33] S. Temple. Advancing cell therapy for neurodegenerative diseases. Cell stem cell, 30(5):512–529, 2023.
- [34] R. Tong, R. da Silva, D. Lin, A. Ghosh, J. Wilsenach, E. Cianfarano, P. Bashivan, B. Richards, and S. Trenholm. The feature landscape of visual cortex. bioRxiv 2023.11.03.565500, Nov. 2023.
- [35] R. Tremblay, S. Lee, and B. Rudy. Gabaergic interneurons in the neocortex: from cellular properties to circuits. *Neuron*, 91(2):260–292, 2016.
- [36] I. Ustyuzhaninov, M. F. Burg, S. A. Cadena, J. Fu, T. Muhammad, K. Ponder, E. Froudarakis, Z. Ding, M. Bethge, A. S. Tolias, and A. S. Ecker. Digital twin reveals combinatorial code of non-linear computations in the mouse primary visual cortex. bioRxiv 2022.02.10.479884, Feb. 2022.
- [37] E. Y. Walker, F. H. Sinz, E. Cobos, T. Muhammad, E. Froudarakis, P. G. Fahey, A. S. Ecker, J. Reimer, X. Pitkow, and A. S. Tolias. Inception loops discover what excites neurons most using deep predictive models. *Nature Neuroscience*, 22(12):2060–2065, Nov. 2019.

- [38] M. A. Weis, L. Pede, T. Lüddecke, and A. S. Ecker. Self-supervised graph representation learning for neuronal morphologies. *Transactions on Machine Learning Research*, 2023.
- [39] K. F. Willeke, K. Restivo, K. Franke, A. F. Nix, S. A. Cadena, T. Shinn, C. Nealley, G. Rodriguez, S. Patel, A. S. Ecker, F. H. Sinz, and A. S. Tolias. Deep learning-driven characterization of single cell tuning in primate visual area v4 unveils topological organization. bioRxiv 2023.05.12.540591, May 2023.
- [40] W. Wu, C. Liao, Z. Deng, Z. Guo, and J. Wang. Neuron platonic intrinsic representation from dynamics using contrastive learning. *arXiv* preprint arXiv:2502.10425, 2025.
- [41] C.-H. Yeh, C.-Y. Hong, Y.-C. Hsu, T.-L. Liu, Y. Chen, and Y. LeCun. Decoupled contrastive learning. In S. Avidan, G. Brostow, M. Cissé, G. M. Farinella, and T. Hassner, editors, *Computer Vision ECCV 2022*, pages 668–684, Cham, 2022. Springer Nature Switzerland.
- [42] H. Yu, H. Lyu, Y. Xu, C. Windolf, E. K. Lee, F. Yang, A. M. Shelton, O. Winter, I. B. Laboratory, E. L. Dyer, C. Chandrasekaran, N. A. Steinmetz, L. Paninski, and C. L. Hurwitz. In vivo cell-type and brain region classification via multimodal contrastive learning. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [43] H. Zeng. What is a cell type and how to define it? Cell, 185(15):2739–2755, 2022.
- [44] H. Zeng and J. R. Sanes. Neuronal cell-type classification: challenges, opportunities and the path forward. *Nature Reviews Neuroscience*, 18(9):530–546, 2017.
- [45] Y. Zhang, Y. Wang, M. Azabou, A. Andre, Z. Wang, H. Lyu, I. B. Laboratory, E. L. Dyer, L. Paninski, and C. L. Hurwitz. Neural encoding and decoding at scale. In *Forty-second International Conference on Machine Learning*, 2025.

Appendix

A Additional Results

A.1 Effect of bin-size

The encoder used in this work assumes the neural activity is binned with a finite bin-size, which sets the minimum time resolution at which the encoder can views the data itself. A fair hypothesis would be that certain cell-types or neurons from certain brain-regions would be more (or less) identifiable at a certain bin-size setting.

We test this by performing a sweep over the bin-size, and measuring the class-wise and macro F1-scores. As we see in Figure 3, some kinds of neurons do indeed prefer certain bin-sizes, however, the overall effect of this hyperparameter is surprisingly small. In other words, the overall classifier performance is relatively robust to the choice of bin-size.

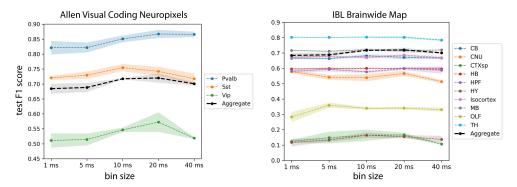


Figure 3: **Results of bin-size sweep.** Classification performance variation for different classes upon sweeping the bin-size hyperparameter of the encoder. "Aggregate" refers to the overall macro-F1 score. Error polygons represents standard error of mean (SEM) measured across 3 pretraining seeds.

A.2 Experiments with a spike-based encoder

We also explored a spike-tokenization strategy inspired by POYO [1]. Unlike POYO, our model aims to *produce* neuron-level embeddings at its outputs, and therefore cannot assume or learn neuron identities through separate "unit embeddings." Instead, we assign the same learnable embedding to all neurons, so that each spike token consists only of its timestamp and this shared embedding. Each neuron's spike train is processed using perceiver-style cross-attention layers, which are queried by learnable vectors placed at linearly spaced timestamps—using a spacing equal to the patch length in our main model. The shared embedding serves only to enable cross-attention over the spike sequence and does not encode any neuron-specific information. The resulting set of tokens is then treated as activity tokens (as shown in Figure 1A) and passed through the spatio-temporal transformer identically to our main architecture. A schematic of this spike-tokenization design is shown in Figure 4A.

We find that this spike-based encoder performs significantly worse on the Allen VC cell-type classification task, while achieving nearly comparable performance on the IBL brain-region classification task (Table 4). We also note that the training stability and convergence speed is much better in the binned-and-patched version of the model, as shown in Figure 4B. Alternative implementations of spike tokenization may yield improved results, however, we leave that for future work.

Table 4: **Spike-tokenization ablation.** Values are reported as mean \pm std. dev. across 5 seeds for the 20ms binning model (main NuCLR), and 3 seeds for the spike tokenization based model.

Model Variant	Allen VC	IBL
20ms binning (main NuCLR) Spike-tokenization based	$\begin{array}{c} 0.7200 \pm 0.0267 \\ 0.6665 \pm 0.0167 \end{array}$	$\begin{array}{c} 0.5295 \pm 0.0040 \\ 0.5268 \pm 0.0136 \end{array}$

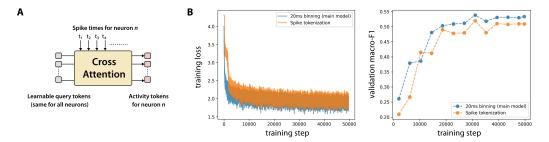


Figure 4: **Spike-tokenization based encoder.** (A) A spike-tokenization layer for the NuCLRencoder. (B) Pretraining-time metrics comparison for spike-tokenization version and binning version of NuCLR. These plots are for pretraining on the IBL dataset.

A.3 Probing intermediate layer representations

In all our experiments, we used a 6-layer implementation of NuCLR, with the first two layers being temporal-only layers followed by two spatio-temporal layers (each having 2 transformer blocks). In Figure 5 we measure the cell-type and brain-region decoding accuracies of the embeddings obtained at all 6 layers. The results show a general improvement of classification performance as we go deeper in the network, with major performance jumps being observed whenever the representation passes through spatial layers. This is another confirmation of the importance of spatial layers, as shown in Section 3.4.

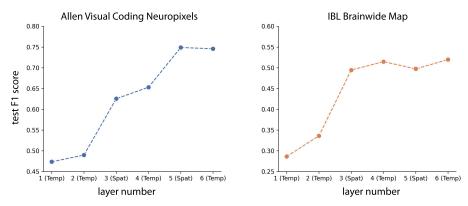


Figure 5: Classification performance of intermediate layer representations. Plots present the classification performance at the *output* of all layers in our 6 layer model for a single seed.

A.4 Comparison with POYO+ on Allen VC

POYO+ [2] is a multi-task behavior decoder model that can be trained simultaneously on multiple recordings. It has been shown to perform well on a Cre-line classification task on the Allen Brain Observatory calcium imaging dataset, using the *session-averaged* latent outputs of its encoder. This is possible since each recording (session) in that dataset has neurons with only one cell-type. This analysis method used in the POYO+ manuscript is not directly applicable here since we aim to produce neuron-level embeddings that can identify different types of neurons *within* a given population.

However, POYO+ learns a "unit embedding" for each neuron, which can, in theory, be used for the purpose of decoding neuron-level properties. We use this strategy to compare POYO+ with NuCLR. We train POYO+ on the Allen VC dataset for decoding the following behaviors: Drifting Gratings Orientation, Drifting Gratings Temporal Frequency, Gabor Orientation, Gabor Position, Natural Scenes, Running Speed, and Static Gratings Orientation.

Since POYO+ can only be tested under transductive evaluations, we choose the transductive zero-shot regime for our comparison, and find that NuCLR outperforms the learned unit embeddings of POYO+ by a strong margin (Table 5).

Table 5: Comparing POYO+ with NuCLR and other methods. Evaluation done in transductive zero-shot setting on the Allen VC dataset, with result being macro-F1 scores represented as mean \pm std. dev. over 5 seeds.

POYO+	NuCLR	NeuPRINT	NEMO
0.3521 ± 0.0233	$0.7218 \pm \textbf{0.0113}$	$0.3999 \pm \textbf{0.0312}$	$0.4256 \pm \textbf{0.0114}$

Pretraining POYO+ required an average of 7.5 hours on 4 NVidia B200s to converge (about 100 epochs), while NuCLR requires only about 1.5 hours on the same hardware. We used the main training configuration from the example provided by the authors⁵. We also report the behavior decoding performance achieved by POYO+ after training across five random seeds:

• Drifting Gratings Orientation Accuracy: 93.07% ± 1.76% (Chance: 12.5%)

• Drifting Gratings Temporal Frequency Accuracy: 93.40% ± 1.36% (Chance: 20%)

• Gabor Orientation Accuracy: 56.12% ± 1.18% (Chance: 25.0%)

• Gabor Position (2D) R^2 : 0.6888 ± 0.0487

• Natural Scenes Accuracy: $53.23\% \pm 10.18\%$ (Chance: 0.84%)

• Running Speed R^2 : 0.7681 ± 0.0115

• Static Gratings Orientation Accuracy: 75.99% ± 0.97% (Chance: 12.5%)

A.5 Ablating temporal-only attention layers

To assess the importance of dedicated temporal-only attention layers, we replace them with an equivalent number of spatio-temporal layers. Specifically, the first two temporal-only layers in the original model were substituted with one spatio-temporal layer to maintain a similar model capacity.

As shown in Table 6 the inclusion of dedicated temporal-only layers appears to be a relatively inconsequential design decision, as performance does not change significantly between the two model configurations. However, a benefit of retaining the temporal-only layers is computational efficiency. Given that the number of temporal patches (~ 10) is typically much smaller than the number of neurons (~ 100), a temporal layer has less computationally expensive than a using separate temporal layers reduces the overall computational complexity, as attention mechanisms scale quadratically with the number of tokens in a sequence.

Table 6: Ablation study for temporal-only layers. Values are reported as mean \pm std. dev. across 5 seeds for full NuCLR, and 3 seeds for the ablated model.

Model Variant	Allen VC	IBL
Full NuCLR	0.7200 ± 0.0267	0.5295 ± 0.0040
w/o temporal attention layers	$0.7184 \pm \textbf{0.0124}$	0.5259 ± 0.0089

A.6 Confusion matrices

Please refer to Appendix A.6 for confusion matrices. We present the confusion matrices achieved by NuCLR in Figure 6 for all datasets evaluated in the inductive zero-shot setting.

A.7 Additional results on Bugeon et. al. dataset

We evaluate the representations of NuCLR and baseline models on an 11-class label set of the Bugeon et. al. dataset. This label set consists of the classes: Lamp5-Chrna7, Lamp5-Npy, Lamp5-Tmem182, Pvalb-Tac1, Pvalb-Vipr2, Sncg-Pdzrn3, Sncg-Vip, Sst-Reln, Sst-Tac1, Vip-Cp, Vip-Reln. As seen in Table 7, NuCLR outperforms both NeuPRINT and LOLCAT baselines on this labelling of neurons.

⁵https://github.com/neuro-galaxy/torch_brain/tree/main/examples/poyo_plus

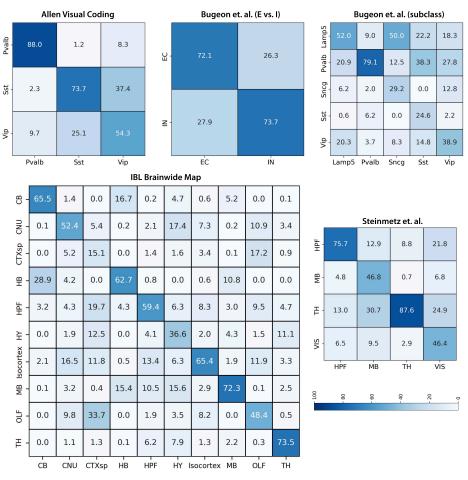


Figure 6: Confusion matrices achieved by NuCLR on inductive zero-shot evaluation.

Table 7: Macro F1-score for cell type classification on the 11-class label set of Bugeon et. al. Reported as mean \pm std. dev. across 5 training seeds. N/A indicates the method cannot operate in that evaluation setting.

Dataset	# Classes	Setting	NuCLR	NeuPRINT	LOLCAT
Bugeon et. al. (11-class)	11	Transductive Transductive zero-shot Inductive zero-shot		$\begin{array}{c} 0.2748 \pm 0.0131 \\ 0.1825 \pm 0.0158 \\ \text{N/A} \end{array}$	$\begin{array}{c} 0.1559 \pm 0.0318 \\ N/A \\ 0.1482 \pm 0.0256 \end{array}$

B Additional method details

B.1 View-pair sampling

We sample view pairs during pretraining in three steps:

- 1. **First view.** The first view is sampled using the RandomFixedWindowSampler found in torch_brain⁶. This sampler divides each recording into non-overlapping windows of length $T_{\rm ctx}$ and randomly samples one window (without replacement throughout an epoch). Across epochs, a shared random jitter to the window start times at the beginning of every epoch. This emulates uniform sampling while ensuring complete data coverage each epoch.
- 2. **Second view.** The second view is sampled relative to the first. Its start time is drawn uniformly from a range constrained by the $\Delta T_{\rm max}$ setting, ensuring the two views come from nearby temporal contexts.

⁶https://github.com/neuro-galaxy/torch_brain

3. **Neuron dropout.** To increase robustness to partial observations, we apply neuron dropout independently to each view. Given a view with N neurons, we first sample the number of neurons to drop from a uniform distribution: $N_{\text{drop}} \sim \mathcal{U}(0, 0.5N)$. We then randomly select N_{drop} neurons and exclude them from the view during that training step.

B.2 Model implementation details

Transformer. We use standard scaled-dot-product attention as implemented in xformers⁷, prenormalization using LayerNorm. Our feedforward network (FFN) uses the GEGLU activation function, with its hidden dimension being 4x the dimension of the tokens (*D*). In the temporal transformer layers, we use Rotary Embeddings to incorporate timing information, as described below.

Rotary Time Embeddings. We use rotary time embeddings following the formulation in Section A.1 of [1], which includes the use of value rotation in addition to query-key rotation. The only difference in our implementation lies in the choice of temporal scaling parameters. Because our input tokens are uniformly spaced with a stride of T_{patch} , we set $T_{\text{min}} = T_{\text{patch}}$ and $T_{\text{max}} = 8 \times T_{\text{ctx}}$.

B.3 Adapting to Calcium Imaging Data

To adapt NuCLR's spatio-temporal transformer to calcium imaging data, we replace the binning-and-patching step used for spike trains with a temporal patching operation applied directly to the calcium fluorescence time series (i.e. the $\Delta F/F$ signal). This modification affects only the input stage of the model; the architecture and training procedure remain unchanged. Unlike electrophysiology, calcium imaging does not involve physical probes or insertions. Therefore, we treat all simultaneously recorded neurons within a session as a single population—both for the spatial transformer layers and for contrastive loss computation.

B.4 Hyperparameters

We use the AdamW [22] optimizer with a linear learning rate warm-up over the first epoch, followed by cosine decay until end of training. All relevant hyperparameters for training NuCLR are listed in Table 8. These values were mainly selected via manual line searches on the IBL development set (for ephys data, Appendix D.2) and the Bugeon et al. development set (for calcium imaging data, Appendix D.4). Across all datsets, we pretrain the model for 50,000 steps and use the bfloat16 number format throughout. Pretraining takes approximately 3 hours on a machine with $4 \times NVidia H100$ GPUs.

Table 8: Key hyperparameters for NuCLR.

Parameter	Value for Ephys.	Value for Ca ⁺² (diff. only)
$T_{ m ctx}$	10s	30s
T_{patch}	1s	
$\Delta T_{ m max}$	30s	240s
Bin size	20ms	N/A
L_{T}	2	
$L_{ m ST}$	2	
D	256	
Num. attention heads	4	
Linear dropout	0.2	
Attention dropout	0.0	
Max. neuron dropout	50%	
Num. training steps	50,000	
Batch size	128	16
Max learning rate	5×10^{-4}	1.25×10^{-4}
Weight decay	0.01 (default)	
β_1	0.9 (default)	
β_2	0.999 (default)	

⁷https://github.com/facebookresearch/xformers

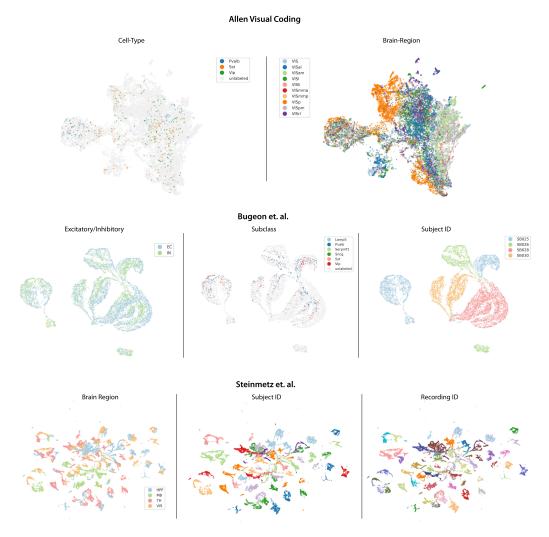


Figure 7: UMAP visualizations of NuCLR's embeddings colored by various properties.

C Embedding Visualizations

We present UMAP-based visualizations of NuCLR's output representations in Figure 7. For the Allen VC dataset, which includes a large number of subjects, the embeddings exhibit clear density modes that align most strongly with brain-region information.

In contrast, the Bugeon et al. and Steinmetz et al. datasets are relatively small, containing only 4 and 10 subjects respectively, with limited total number of recordings. For these datasets, NuCLR's embeddings cluster primarily by subject or session identity. However, *within* these clusters, we still observe meaningful density modes corresponding to cell-type and brain-region structure. Embeddings for the IBL dataset, shown in Figure 2A, reveal strong region-based organization without noticeable clustering by subject or session.

The subject- and session-specific clustering observed in smaller datasets (Bugeon et al. and Steinmetz et al.) may hinder data-driven discovery, as it suggests entanglement with recording-specific factors. While this effect appears only in small datasets, mitigating it remains an important direction for future work.

D Details on datasets and evaluation methodology

This section contains description of each dataset, and details our train-test splits for all evaluation settings. Additionally, our code-base⁸ includes all preprocessing, split preparation, and evaluation scripts.

D.1 Allen Visual Coding

The Allen Visual Coding (VC) dataset consists of 58 Neuropixels recordings, each from a unique mouse. Each recording spans approximately two hours, during which various visual stimuli are presented. Of these recordings, 16 contain optotagged neurons, with each labeled recording containing neurons from only one of the three inhibitory cell types: Pvalb, Vip, or Sst. Since reliable cell-type labels are only available for neurons in the visual cortex, we restrict all models—including NuCLR and baselines—to use only neurons from the VIS region during pretraining.

For zero-shot classifier evaluation, we follow a leave-one-subject-out strategy. Specifically, there are 16 test folds corresponding to the 16 labeled subjects. In each fold, one subject is held out for testing, and the remaining 15 are used for training. Within the training set, we perform 4-fold cross-validation (subject-wise) to select the best training epoch. The final test score is reported from a classifier trained on all 15 subjects using the selected epoch. We did not test classifier performance on this dataset in the transductive setting as a "perfect" classifier only has to infer which of the previously seen sessions does the neuron belong to, and the cell-type corresponding to that session.

D.2 IBL Brainwide Map

The IBL Brainwide Map dataset consists of 439 recordings from 139 unique mice, with each recording performed using one or two Neuropixels insertions. While performing zero-shot evaluations (both inductive and transductive), we construct train—test splits at the recording level, ensuring that no co-recorded neuronal populations appear in both sets. In total, 93 recordings are designated for the test set, and the remaining 346 are used for training and validation. Of the training set, 91 recordings are used as a development set for tuning hyperparameters specific to electrophysiology data. During classification, these same 91 recordings also serve as a validation set for selecting the best training epoch. Final performance metrics are reported using a model trained on all 346 training recordings, evaluated on the held-out test set.

For non-zero-shot evaluation, we create a neuron-wise stratified test-train split with a test size of 20%. The validation set for finding the best epoch is created from the train fold with a size of 20%.

D.3 Steinmetz et al.

The Steinmetz et al. dataset consists of 39 Neuropixels recordings from 10 unique mice, with each recording comprising 2 or 3 probe insertions. For non-zero-shot evaluation, we create train–test splits following the same procedure as described for the IBL dataset (Appendix D.2).

For transductive zero-shot evaluation, we use a 10-fold leave-one-subject-out strategy. In each fold, one subject is held out for testing, and the remaining nine are used for training. We select the best training epoch by performing a stratified 80/20 train–validation split within the training set. For inductive zero-shot evaluation, we hold out 3 subjects for testing and pretrain on the remaining 7. To select the best classifier epoch, we perform 4-fold subject-wise cross-validation within the training subjects.

D.4 Bugeon et al.

The Bugeon et al. dataset contains 17 spatial transcriptomic calcium imaging recordings from 4 unique mice. We train on all stimulus-specific sub-recordings within these sessions. Since this is the only optophysiology dataset in our evaluation, we adjusted several hyperparameters specifically for calcium activity. During development, we held out 4 recordings—one from each subject—as a

⁸https://github.com/nerdslab/nuclr

⁹We use the same set of labeled neurons as in the original LOLCAT publication [28].

validation set for tuning model hyperparameters. These recordings are never included in the test set when reporting final performance.

For non-zero-shot evaluation, we create neuron-wise train—test splits following the same protocol as in the IBL dataset (Appendix D.2). For transductive zero-shot evaluation, we adopt a 4-fold leave-one-subject-out strategy. In each fold, we perform 3-fold subject-wise cross-validation within the training subjects to select the best epoch. For inductive zero-shot evaluation, we hold out one subject (SB028) for testing and pretrain on the remaining three. As in the transductive case, we perform 3-fold subject-wise cross-validation on the training set to select the best model epoch.

E Implementation details for baseline models

E.1 NeuPRINT

Calcium Imaging (Bugeon et al.). For the Bugeon dataset, we used the publicly available NeuPRINT implementation and followed a transductive evaluation setup consistent with the original codebase. The model was first pretrained on the full dataset. For evaluation, the embedding table was reinitialized and optimized again using the self-supervised loss, while keeping the backbone encoder frozen. The resulting neuron embeddings were then used for downstream classification.

Electrophysiology Datasets. For electrophysiology datasets, we modified NeuPRINT's training loop and epoch structure to accommodate the larger scale and greater diversity of these datasets. In the original implementation, there is no true notion of an epoch—batches are drawn from individual sessions without ensuring full dataset coverage, and the sampling strategy is manually biased toward neurons with labeled cell types. This setup does not translate well to larger datasets such as Steinmetz et. al., IBL, and Allen VC.

We instead adopted the standard definition of an epoch, which is a full pass through all training data, and trained for 300 epochs. We also replaced the session-specific sampling strategy in the original implementation with uniform sampling across all sessions to ensure unbiased data coverage.

Furthermore, because these electrophysiology datasets are substantially larger than Bugeon et al. we did not re-initialize and relearn the embedding table during evaluation. The original NeuPRINT protocol is tailored for the Bugeon et al. dataset, where such re-initialization is computationally feasible. However, for large-scale datasets like Steinmetz, IBL, and Allen VC, this procedure becomes too computationally expensive. Therefore, we adopted a more scalable and simpler protocol that just uses learned embeddings from pretraining for downstream evaluation.

All hyperparameters were retained from the original NeuPRINT setup, except for the learning rate, which we set to 1×10^{-2} , and the backward context window, which was set to 10.

NeuPRINT relies on behavior features being provided along with the mean and standard deviations for the population activity. The behavioral features used for electrophysiology dataset were as follows:

- Steinmetz et al.: face motion, pupil area, and wheel velocity
- IBL: wheel velocity
- Allen VC: running speed

E.2 NEMO

We evaluated NEMO using the official implementation provided by the original authors upon request. When adapting the method to new datasets, we retained all hyperparameters from the original paper, with the exception of the waveform template dimensionality, which we adjusted to match the characteristics of each dataset.

E.3 LOLCAT

We evaluated LOLCAT using the official implementation provided by the original authors upon request. For all experiments (unless otherwise stated below), we used a batch size of 64, a learning rate of 1×10^{-3} , a weight decay of 1×10^{-5} , and a dropout rate of 0.5. The MLP hidden dimensions

were set to [64, 32, 16], and 4 attention heads were used. The snippet dropout rate during training was 0.45. We trained the models for 300 epochs with the optimizer outlined in the paper. These default parameters were chosen based on manual exploration of the "reduced" hyperparameter search range and explicit prescriptions in the LOLCAT paper.

- Allen VC: We used the same final hyperparameter selection as the paper with new results on our split of the data.
- IBL: We used the same hyperparameters as Allen VC.
- Steinmetz transductive: Batch size was increased to 1024, and the epochs were set to 1000.
- **Steinmetz inductive**: The minimum factor was set to 0.01, and we initialized the classes undersampling factors to [8 (HPF), 8 (MB), 0.01 (TH), 0.01 (VIS)].
- **Bugeon subclass**: The classes undersampling factors were initialized to [0.01 (Lamp5), 0.01 (Pvalb), 8 (Sncg), 8 (Sst), 8 (Vip)] and the minimum factor to 0.01.
- **Bugeon E vs. I**: Class undersampling factors were initialized to [0.01 (E), 8 (I)] and the minimum factor to 0.01.

NeurIPS Paper Checklist

1. Claims

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

Answer: [Yes]

Justification: All claims made have been confirmed in the results section Section 3, to the best of our ability.

Guidelines:

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

2. Limitations

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

Answer: [Yes]

Justification: We discuss the limitations of our work in Section 5.

Guidelines:

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

3. Theory assumptions and proofs

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

Answer: [NA]

Justification: No new theoretical results have been introduced in this work.

Guidelines:

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

4. Experimental result reproducibility

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

Answer: [Yes]

Justification: We have described our model architecture and training methodology clearly and fully in Section 2.

Guidelines:

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

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We will release our code along with rest of the supplementary material for this submission.

Guidelines:

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

6. Experimental setting/details

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

Answer: [Yes]

Justification: We describe all the training and hyperparamter details in Appendix B, and describe datasets and train-test split methodologies in Appendix D.

Guidelines:

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

7. Experiment statistical significance

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

Answer: [Yes]

Justification: We use SEM as our error reporting method, and mention it in all our results tables.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

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

Answer: [Yes]

Justification: We provide the details of computational resources required for reproducing the our experimets in Appendix B.

Guidelines:

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

9. Code of ethics

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

Answer: [Yes]

Justification: We have read through the guidelines and confirm that our research conforms to those guidelines.

Guidelines:

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

10. Broader impacts

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

Answer: [NA]

Justification: Our work is aimed at developing a fundamental tool for neuroscience research, advancing the field. While many potential and futuristic downstream societal consequences of advancing the field of neuroscience, we feel these do not need to be specifically highlighted in this work.

- 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: [NA]

Justification: This paper is about the design of a neuroscience research tool, which do not pose any such risks directly.

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: [NA]

Justification: The code provided in supplementary material is original work, and does not use existing assets.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

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

Answer: [Yes]

Justification: Attached code is well documented.

Guidelines:

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

14. Crowdsourcing and research with human subjects

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

Answer: [NA]

Justification: No studies with human subjects were involved in this research.

Guidelines:

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

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

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

Answer: [NA]

Justification: No studies with human subjects were conducted.

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

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

Answer: [NA]

Justification: LLMs were used only for writing, editing, or formatting purposes.

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