MULTIWAY MULTISLICE PHATE: VISUALIZING HID-DEN DYNAMICS OF RNNS THROUGH TRAINING

Anonymous authors

Paper under double-blind review

ABSTRACT

Recurrent neural networks (RNNs) are a widely used tool for sequential data analysis, however, they are still often seen as black boxes of computation. Visualizing the internal dynamics of RNNs is a critical step in understanding the functional principles of these networks and developing ideal model architectures and optimization strategies. Previous studies typically only emphasize the network representation post-training, overlooking their evolution process throughout training. Here, we present Multiway Multislice PHATE (MM-PHATE), a novel method for visualizing the evolution of RNNs' hidden states. MM-PHATE is a graphbased embedding using structured kernels across the multiple dimensions spanned by RNNs: time, training epoch, and units. We demonstrate on various datasets that MM-PHATE uniquely preserves hidden representation community structure among units and identifies information processing and compression phases during training. The embedding allows users to look under the hood of RNNs across training and provides an intuitive and comprehensive strategy to understanding the network's internal dynamics and draw conclusions, e.g., on why and how one model outperforms another or how a specific architecture might impact an RNN's learning ability.

026 027 028

029

025

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Recurrent neural networks (RNNs) are artificial neural networks (ANNs) designed for sequential data. Unlike feedforward neural networks (FNNs), which treat each input as independent, RNNs model sequences of inputs to produce one or a sequence of outputs. RNNs achieve this by processing each sequence element in order while retaining a memory of past inputs through memory units (e.g., LSTMs or GRUs) or through recurrent feedback connections (Lipton et al., 2015; Kaur & Mohta, 2019). This memory effectively passes information forward in time, updating the internal state of the RNN with each input to reflect the sequence's context. RNNs are thus particularly suited for tasks where the order and relationship between elements are crucial for understanding the whole sequence, e.g., neural time-series, postural action data, etc. (Hewamalage et al., 2021).

Since their initial rise in popularity in the 1990s, researchers have created various RNN variants 040 and training strategies to improve their training stability and ability to learn long-range time de-041 pendencies within data, such as Long-Short Term Memory (LSTM) (Hochreiter & Schmidhuber, 042 1996) networks and Structurally Constrained Recurrent Network (SCRN) (Mikolov et al., 2014) ad-043 dressing the vanishing gradient problem (Salehinejad et al., 2018). Moreover, RNNs have a natural 044 advantage in modeling irregular or incomplete sequential data, thanks to their flexibility in managing inputs of varying lengths. These intrinsic properties, as well as the developments of different architectures (Salehinejad et al., 2018), have led to RNNs showcasing exceptional performance across 046 numerous domains, such as natural language processing (NLP), neuroscience and biomedical signal 047 processing (Chen & Li, 2021; Barak, 2017; Khalifa et al., 2021). Moreover RNNs remain the state 048 of the art for neural decoding models used in intracortical brain-computer interfaces (Deo et al., 049 2024). 050

Despite their extensive use, RNN dynamics remain poorly understood. The opaque nature of their
 learned representations complicates the interpretation of their performance and robustness, hinders
 the selection of appropriate architectures and training parameters, and slows the development of
 more effective models. Significant efforts have been made to address the "black-box" nature of

FNNs (Wang et al., 2021; Yu et al., 2014; Li et al., 2018). However, similar advances have not been as prevalent for RNNs, which have primarily been analyzed from the standpoint of dynamical systems analyses after training (e.g., via fixed points) (Sussillo & Barak, 2013), or by comparing the performance of different RNN architectures at the network level and investigating the role of the various components within them (Chung et al., 2014; Greff et al., 2017). Thus, new methods that facilitate the interpretation of RNNs' latent representations and their evolution during training are clearly needed.

061 In the field of explainable deep learning, dimensionality reduction methods are one of the more 062 popular techniques (Karpathy et al., 2015; Hidaka & Kurita, 2017; Rauber et al., 2016; Gigante 063 et al., 2019; Hong et al., 2020; Holtz et al., 2022). These approaches are useful as they allow 064 researchers and engineers to readily visualize the structure of high-dimensional data, allowing them to gain intuition on structures and interactions within the data quickly. Nevertheless, they may 065 inadvertently obscure meaningful data relationships by either emphasizing local (e.g., t-SNE Maaten 066 & Hinton (2008)) or global (e.g., PCA Maćkiewicz & Ratajczak (1993) or Isomap Tenenbaum et al. 067 (2000)) structures, or possess sensitivity to noise and outliers (Moon et al., 2019). These challenges 068 are compounded in cases such as RNNs where structures along multiple dimensions (time, epoch, 069 unit) must all be preserved to produce summaries that capture the data's complexity. As a result, traditional dimensionality reduction techniques fall short in faithfully visualizing learning in RNNs. 071

To study learning in FNNs, Gigante et al. (Gigante et al., 2019) introduced Multislice PHATE (M-072 PHATE) to visualize the network's hidden state's representations during training. M-PHATE com-073 bines a multislice graph representation to model the community structure and temporal relationship 074 between hidden units over epochs, with Potential of Heat Diffusion for Affinity-based Transition 075 Embedding (PHATE) (Moon et al., 2019) for dimensionality reduction. Each slice in the multi-slice 076 graph captures the network's state at a specific epoch during training. This collection of graphs 077 represents the dynamical system that governs the evolution of the network's hidden states, where each hidden unit is connected to itself across epochs. While M-PHATE requires only the activa-079 tions of hidden units during training, the generated visualization captures key properties of network 080 performance, such as test error accuracy, without any held-out validation data.

081 Despite its numerous advantages, M-PHATE is designed to visualize the evolution of FNNs across epochs but fails to account for the recurrent and sequential nature of RNNs. Research has shown that 083 hidden states from all time-steps are crucial for RNNs' representations (Su & Shlizerman, 2020). 084 Therefore, to fully understand the evolution of these representations, it is essential to consider net-085 work dynamics across both time-steps and epochs. To address this need, we propose Multiway Multislice PHATE (MM-PHATE). Our method captures the latent dynamics of RNNs by visual-087 izing hidden states across both time-steps and training epochs, providing deeper insights into the 880 complex learning processes of RNNs. Our findings indicate that MM-PHATE retains considerably more dynamic details necessary for understanding RNNs' performance compared to PCA, t-SNE, 089 Isomap, and M-PHATE. 090

- 091 Our main contributions are as follows:
 - We present MM-PHATE, a novel framework for visualizing the hidden dynamics of RNNs across both time-steps and epochs simultaneously, providing a new perspective to RNN's leaning trajectory, learned representation, and model performance.
 - We demonstrate that MM-PHATE uniquely preserves the hidden representation community structure of hidden units throughout training by tracking each unit's learning trajectory and the correlations among their activations.
 - Applying MM-PHATE to RNN dynamics identifies phases of information processing and compression during learning, an observation that aligns with information bottleneck theory.

102 2 RELATED WORK

094

095 096

098 099

100

101

We group existing methods to interpret RNNs into two categories: 1) performance-oriented and 2) application-oriented post-training analysis.

107 Performance-oriented analyses include investigating the role of components within various RNN architectures, as well as comparing different architectures and training parameters based on their

network-level performance post-training. For example, Chung et al. (2014) conducted
 a performance evaluation comparing gated RNNs, particularly GRUs and LSTMs. Similarly, Greff
 et al. (2017) carried out a comprehensive analysis of LSTM components. While these studies interpret the performance of different RNNs at the network level, they did not delve deeply into analyzing
 their hidden states, thus providing limited insight into the nature of the representations learned by
 these networks.

114 In contrast, application-oriented analyses of trained networks often involve visualizing and inter-115 preting activation maps that depict the representations learned by hidden units after training, often 116 in the context of a specific task or network configuration. Numerous studies have applied this ap-117 proach within the field of NLP. For instance, Karpathy et al. overlaid activation heat maps over texts, 118 demonstrating that certain units developed interpretable representations, such as tracking text length, new line beginnings, and quote initiations (Karpathy et al., 2015). Li et al. used saliency heat maps 119 to identify words critical to the learned representations (Li et al., 2016). Strobelt et al. and Ming et 120 al. created interactive visualizations to correlate hidden state patterns with phrases in texts (Strobelt 121 et al., 2018; Ming et al., 2017). Beyond NLP applications, some studies explore other domains, 122 including speech recognition (Tang et al., 2017), earth sciences (Titos et al., 2022)), and medical 123 states (Kwon et al., 2019). While these studies provide intuitive insights into the representations 124 learned by the networks post-training, they are tailored for a particular task and may not generalize 125 to other tasks. Other studies tailored to understand general RNN properties include applying Proper 126 Orthogonal Decomposition (POD) to the internal states of encoder and decoder units in Seq2Seq 127 RNNs (Su & Shlizerman, 2020) or employing PCA to visualize the activation of recurrent units and 128 its link to generalization(Farrell et al., 2022).

While these application-oriented studies have identified critical aspects of RNNs' hidden representation, they largely overlook the network's learning trajectory over training epochs, a crucial process that must be understood to improve RNN-based machine learning systems. To the best of our knowledge, no existing visualizations of RNNs' hidden states interpret the hidden dynamics simultaneously across time-steps and epochs.

134 135

3 BACKGROUND

136 137

138 **PHATE:** PHATE is a recent visualization technique that can capture both the local and global structure of data using diffusion processes (Moon et al., 2019). The PHATE algorithm optimizes 139 the diffusion kernel (Coifman & Lafon, 2006) for the visualization of high dimensional data. Let 140 x_i be a point in a high-dimensional dataset. PHATE begins by computing the Euclidean distance 141 matrix E between all data points, where $E_{ij} = \|x_i - x_j\|_2$. These distances are then trans-142 formed into affinities using an adaptive α -decay kernel $K_{k,\alpha}(x_i, x_j) = \frac{1}{2} \exp\left(-\left(\frac{E_{ij}}{\epsilon_k(x_i)}\right)^{\alpha}\right) +$ 143 144 $\frac{1}{2}\exp\left(-\left(\frac{E_{ij}}{\epsilon_k(x_j)}\right)^{\alpha}\right)$, which adapts to the data density around each point and captures local infor-145 mation. The parameters $\epsilon_k(x_i)$ and $\epsilon_k(x_i)$ are the k-nearest-neighbor distance of x_i and x_j , and 146 α controls the decay rate. The affinities are then row-normalized to obtain the diffusion operator 147 $P = D^{-1}K_{k,\alpha}$ that represents the single-step transition probabilities between data points, where 148 D is a diagonal matrix whose entries are row sums of $K_{k,\alpha}$. PHATE calculates the information dis-149 tance between points based on their transition probabilities: dist_{ij} = $\sqrt{\|\log P_i^t - \log P_j^t\|^2}$, where 150 151 P^t captures the transition probabilities of a diffusion process on the data over t steps and i and j 152 are rows in the matrix. These distances are embedded into low dimensions using Multidimensional 153 Scaling (MDS) (Ramsay, 1966) for visualization. Local and global distances within the data's manifold are represented in PHATE by multistep diffusion probabilities. The diffusion probability of 154 each point can capture the local context surrounding said point, allowing the construction of pair-155 wise comparisons between all points (both neighboring and distant points) that represents the entire 156 global context. For further details, see (Moon et al., 2019). We will use PHATE to embed RNN 157 training dynamics, however we alter the initial graph construction to emphasize certain structures in 158 the data we wish to visualize. 159

160

161 **M-PHATE:** Gigante et al. (2019) model the evolution of the hidden units in a feedforward neural network and their community structure using a multislice graph. Let F be an FNN with a total of m

162 hidden units, and let $F^{(\tau)}$ be the representation of the network after being trained for $\tau \in \{1, ..., n\}$ 163 epochs on the training data X sampled from a larger dataset Π . The algorithm first calculates a 164 shared feature space using the normalized activations of all hidden units $i \in \{1, ..., m\}$ on the input 165 data, as a 3-dimensional tensor:

$$\boldsymbol{T}(\tau, i, k) = \frac{\boldsymbol{F}_i^{(\tau)}(\boldsymbol{Y}_k) - \frac{1}{p} \sum_{\ell} \boldsymbol{F}_i^{(\tau)}(\boldsymbol{Y}_{\ell})}{\sqrt{\operatorname{Var}_{\ell}[\boldsymbol{F}_i^{(\tau)}(\boldsymbol{Y}_{\ell})]}}$$

where $F_i^{(\tau)}(Y_k) : \mathbb{R}^d \to \mathbb{R}$ denotes the activation of the *i*-th hidden unit of F for the *k*-th sample 170 from input data Y. Here, Y is a subset of p samples from the d-dimensional training data X, with 171 an equal number of samples from each input class, and $p \ll |X|$. This activation tensor T is then 172 used to calculate intraslice affinities between pairs of hidden units within an epoch τ during the 173 training, as well as the interslice affinities between a hidden unit i and itself at different epochs: 174

$$\mathbf{K}_{\text{intraslice}}^{(\tau)}(i,j) = \exp\left(\frac{-\|\mathbf{T}(\tau,i) - \mathbf{T}(\tau,j)\|_2^{\alpha}}{\sigma_{(\tau,i)}^{\alpha}}\right), \mathbf{K}_{\text{interslice}}^{(i)}(\tau,\upsilon) = \exp\left(\frac{-\|\mathbf{T}(\tau,i) - \mathbf{T}(\upsilon,i)\|_2^{2}}{\epsilon^2}\right)$$

where α is the α -decay parameter, $\sigma_{(\tau,i)}$ is the intraslice bandwidth for unit i in epoch τ , and ϵ is the fixed interslice bandwidth. These matrices are combined to form an $nm \times nm$ multislice kernel matrix K, which is then symmetrized, row-normalized, and visualized using PHATE in 2D or 3D.

180 181 182

183

175

177

178

179

MULTIWAY MULTISLICE PHATE 4

M-PHATE was shown to be a powerful tool for visualizing FNNs. However, to effectively visualize the evolution of RNNs' hidden representations, we need to consider hidden state dynamics across 185 *time-steps* within the sequence and training epochs concurrently. In RNNs, the output from previous time-steps is fed as an input to current time-steps. This is useful in the treatment of sequences and 187 building a memory of the previous inputs into the network. The network iteratively updates a hidden 188 state h. At each time-step t, the next hidden state h_{t+1} is computed using the input x_t and the current 189 hidden state h_t . Importantly, the network uses the same weights W and biases b for each *time-step*. 190 Thus the output y_t at *time-step* t is $y_t = f(W \cdot h_t + b)$, where f is some activation function. 191

Let ${m R}^{(au)}$ be the representation of an m-unit RNN after being trained for $au \in \{1,\ldots,n\}$ epochs on 192 the training data $X \subset \Pi$. We denote $R_{i,w}^{(\tau)}(Y_k) : \mathbb{R}^d \to \mathbb{R}$ the activation of the *i*-th hidden unit of 193 R at time-step $w \in \{1, \ldots, s\}$ in epoch τ for the k-th sample of Y, where Y consists of p samples 194 from the training data X. We construct the 4-way tensor T using the hidden unit activations as a 195 196 shared feature space, which we use to calculate unit affinities across all epochs and time-steps. The tensor T is an $n \times s \times m \times p$ tensor containing the activations at each epoch $\tau \in \{1 \dots n\}$ and time-197 step $w \in \{1 \dots s\}$ of each hidden unit \mathbf{R}_i $(i \in \{1 \dots m\})$ with respect to each sample $\mathbf{Y}_k \subset \mathbf{X}$. To eliminate the variability in T due to the bias term b, we z-score the activation of each hidden unit at 199 time-step w and epoch τ : 200

203 204

$$\boldsymbol{T}(\tau,\omega,i,k) = \frac{\boldsymbol{R}_{i,\omega}^{(\tau)}(\boldsymbol{Y}_k) - \frac{1}{p} \sum_{\ell} \boldsymbol{R}_{i,\omega}^{(\tau)}(\boldsymbol{Y}_{\ell})}{\sqrt{Var_{\ell}[\boldsymbol{R}_{i,\omega}^{(\tau)}(\boldsymbol{Y}_{\ell})]}}.$$
(1)

We construct a kernel over T utilizing our prior knowledge of the temporal aspect of T to capture its 205 dynamics over epochs and time-steps. This constructed kernel, denoted K, represents the weighted 206 edges in the multislice graph of the hidden units. In this representation, each unit has two types of 207 connections: edges between the unit to itself across epochs and time steps and, within a fixed epoch 208 and time-step, edges between a unit and its community-the other units which have the most similar 209 representation. The edges are weighted by the similarity in activation pattern. We define K as a 210 $nsm \times nsm$ kernel matrix between all m hidden units at all s time-steps in all n training epochs. The 211 $((\tau - 1)sm + (\omega - 1)m + j)_{th}$ row or column of **K** refers to the j_{th} unit at time-step w in epoch 212 τ . We henceforth refer to the row as $K((\tau, \omega, j), :)$ and the column as $K(:, (\tau, \omega, j))$. In order 213 to capture the evolution of hidden units of R across time-steps and epochs, while preserving the unit's community structure, we construct a multiway multislice kernel matrix reflecting two types 214 of connections simultaneously. Given the α -decay parameter α , the intrastep bandwidth for unit *i* at 215 time-step w and epoch τ : $\sigma_{(\tau,\omega,i)}$, and the fixed interstep bandwidth ϵ , we define:

217 218

219 220

221 222

223

232

233

234 235

241

242

243

244

245

253 254

255

256

257

262

264

• Intrastep affinities between hidden units i and j at time-step ω in epoch τ :

$$\boldsymbol{K}_{\text{intrastep}}^{(\tau,\omega)}(i,j) = \exp(-\parallel \boldsymbol{T}(\tau,\omega,i) - \boldsymbol{T}(\tau,\omega,j) \parallel_2^{\alpha} / \sigma_{(\tau,\omega,i)}^{\alpha})$$

• Interstep affinities between a hidden unit *i* and itself at different time-steps and epochs:

$$\boldsymbol{K}_{\text{interstep}}^{(i)}((\tau,\omega),(\eta,\nu)) = \exp(-\parallel \boldsymbol{T}(\tau,\omega,i) - \boldsymbol{T}(\eta,\nu,i) \parallel_2^2 /\epsilon^2)$$

The bandwidth $\sigma_{(\tau,\omega,i)}$ of the α -decay kernel is set to be the distance of unit *i* at time-step *w* from epoch *n* to its *k*-th nearest neighbor across units at that time-step and epoch: $\sigma_{(\tau,\omega,i)} = d_k(\mathbf{T}(\tau,\omega,i), \mathbf{T}(\tau,\omega,:))$, where $d_k(z,Z)$ denotes the ℓ_2 distance from *z* to its *k*-th nearest neighbor in *Z*. We used k = 5 in all the results presented. The use of this adaptive bandwidth means the kernel is not symmetric and thus requires symmetrization. In the interstep affinities $\mathbf{K}_{interstep}^{(i)}$, we use a fixed-bandwidth Gaussian kernel $\epsilon = \frac{1}{nsm} \sum_{\tau=1}^{n} \sum_{w=1}^{s} \sum_{i=1}^{m} d_k(\mathbf{T}(\tau,\omega,i),\mathbf{T}(:,:,i))$, the average across all time-steps in all epochs and all units of the distance of unit *i* at time-step *t* to its k_{th} nearest neighbor among the set consisting of the same unit *i* at all steps.

The combined kernel matrix of these two matrices contains one row and column for each unit at each time-step in each epoch, such that the intrastep affinities form a block diagonal matrix and the interstep affinities form off-diagonal blocks composed of diagonal matrices (Fig. 1a, 1b).

$$\boldsymbol{K}((\tau,\omega,i),(\eta,\nu,j)) = \begin{cases} \boldsymbol{K}_{\text{intrastep}}^{(\tau,\omega)}(i,j) & \text{if } (\tau,\omega) = (\eta,\nu) \\ \boldsymbol{K}_{\text{interstep}}^{(i)}((\tau,\omega),(\eta,\nu)) & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

We symmetrize this final kernel as $\mathbf{K}' = \frac{1}{2}(\mathbf{K} + \mathbf{K}^T)$, and row-normalize it to obtain $\mathbf{P} = \mathbf{D}^{-1}\mathbf{K}'$, where \mathbf{D} is a diagonal matrix whose entries are row sums of \mathbf{K}' and which \mathbf{P} represents a random walk over all units at all time-steps in all epochs, where propagating from one state to another is conditional on the transition probabilities between time-step ω in epoch τ and time-step ν in epoch η . PHATE is applied to \mathbf{P} to visualize the tensor \mathbf{T} in two or three dimensions. This resulting visualization thus simultaneously captures information regarding the evolution of the units across both time-steps and epochs.



Figure 1: Example schematic of the multiway multislice graph (a) and kernel (b) used in MM-PHATE for RNNs. The intrastep kernels represent the similarities between the graph nodes at the same time-steps. The interstep kernels represent the similarities between the nodes and themselves at different time-steps and epochs.

5 Results

We demonstrate the ability of MM-PHATE to capture useful properties of RNN learning on two
datasets: 1) The Area2Bump dataset (Chowdhury et al., 2022) consisting of neural activity recorded
from Brodmann's area 2 of the somatosensory cortex while macaque monkeys performed a slightly
modified version of a standard center-out reaching task and 2) The Human Activity Recognition
(HAR) Using Smartphones dataset (Reyes-Ortiz et al., 2012), kinematic recordings of 30 subjects
performing daily living activities with a smartphone embedded with an inertial measurement unit.

270 5.1 NEURAL ACTIVITY 271

290 291 292

293

295

296 297

298 299

300

301

306

272 We begin with the Area2Bump dataset, which consists of spiking activity data from 273 macaques (Chowdhury et al., 2022). We trained an LSTM network comprised of a single layer of 20 units to classify the direction of arm-reaching movements, and applied MM-PHATE to vi-274 sualize the learning of the LSTM. In an example session where our network achieved a validation 275 accuracy of 74%, we applied the MM-PHATE visualization to the tensor T that consisted of the 276 network activations of all 20 units over 600 time-steps for each of 200 training epochs. In practice, 277 we sampled time-steps and epochs to reduce memory load. Each point in the visualization repre-278 sents a hidden unit at a given time-step in a given epoch (Fig. 2). Through this visualization, we 279 observed a smooth transition of the hidden states across both time-steps and epochs, reflecting the 280 dynamic changes in the network's internal representations throughout training. Here we compare 281 MM-PHATE to three other dimensionality reduction techniques: PCA, t-SNE, and Isomap, and we 282 compare to M-PHATE in the supplement (Fig. S1), using the same T tensor. We first flattened Talong the epoch, time-step, and unit axis, and embedded it with each of the dimensionality reduction 283 284 methods. Notably, the MM-PHATE visualization reveals a distinct split in representations during the later time-steps and epochs, highlighting unique learning patterns as the model converges-patterns 285 that are not discernible with other methods. PCA and Isomap, while showing a seemingly smooth 286 transition, do not show distinct differences between the early and late epochs, failing to capture how 287 the representation transforms during learning. On the other hand, t-SNE resulted in a visualization 288 that lacked smooth transitions across time-steps and epochs, possibly due to its sensitivity to noise. 289



Figure 2: Area2Bump: Visualization of a 20-unit LSTM network trained for 200 epochs. The visualizations are generated using MM-PHATE, PCA, t-SNE, and Isomap, from left to right, respectively. Points are colored based on epoch (top row) or time-step (bottom-row).

5.1.1 INTRA-STEP ENTROPY

307 To evaluate the effectiveness of our embedding in capturing meaningful structures within the net-308 work's hidden dynamics, we next analyzed the flow of information during training (Fig. 3). Specif-309 ically, we aimed to analyze the spread of points in the MM-PHATE space, which corresponds to 310 the diversity of the internal representations of the network. To quantify this property, we com-311 puted the entropy between hidden units in the embedded space at each time-step ω and compared 312 these intra-step entropies with the accuracy and loss metrics recorded at the end of each training 313 epoch. Conceptually, we model a general RNN trained on dataset X with label L as a Markov 314 chain $(L \to X \to R)$. This allows us to compute the mutual information between X and R as 315 $I(\mathbf{X}, \mathbf{R}) = H(\mathbf{R}) - H(\mathbf{R}|\mathbf{X})$, where $H(\mathbf{R})$ and $H(\mathbf{R}|\mathbf{X})$ are the marginal and conditional entropies. Given the deterministic nature of RNNs, $H(\mathbf{R}|\mathbf{X})$ equates to zero, indicating that $H(\mathbf{R})$ 316 at each time-step ω reflects the input information the network retains during training (Tishby & 317 Zaslavsky, 2015; Cheng et al., 2019). 318

Our analysis of the MM-PHATE embedding revealed a general increase in the entropy (Fig. 3) that aligns with the training epoch where the network begins to overfit (approximately epoch 100). We hypothesize that this increase may be due to memorization of noise or other nuisance variance in the training data. The observed changes in entropy at specific time-steps and epochs coincided with shifts in model performance throughout the training. This suggests that our algorithm successfully captures and retains dynamical information critical for the model's learning process. Notably, tran-



Figure 3: Intra-step entropy of all hidden units in embedding space at each time-step in each epoch, compared to training and validation accuracy (top) and losses (bottom), comparing embeddings of MM-PHATE, PCA, t-SNE, and Isomap.

sitions in validation accuracy are not always reflected in changes in the training loss curve, e.g., 340 fluctuation of validation accuracy curve before epoch 30 cannot be observed in the training loss 341 (Fig. 3). The MM-PHATE embedding, however, successfully identified relevant transitions coin-342 ciding with shifts in both performance measures, thus offering insights beyond those provided by 343 accuracy and loss functions alone. Moreover, the MM-PHATE embedding reveals details on in-344 formation processing through time-steps. This is not the case for the loss and accuracy metrics, 345 which are measured once per epoch. With MM-PHATE, we obtain intra-step entropy plots for each 346 time-step, and can thus compared dynamics across them. Notably, we observe diverging patterns 347 between the entropy of earlier and later time steps. In this network, earlier time steps' entropy increases and plateaus around epoch 100, while later time steps mimic a spiking activity around that 348 epoch, eventually returning back to the entropy value they had in earlier epochs. While the effects 349 of these diverging patterns remains uncertain, MM-PHATE clearly offers additional insights into the 350 model's dynamics. 351

352 We performed the same intra-step entropy analysis on other methods (Fig. 3). PCA failed to cap-353 ture the critical dynamics associated with early subtle changes in model performance, likely due to 354 its linear nature and emphasis on the data's global structures. Similarly, Isomap failed to capture significant dynamics related to performance changes before the onset of overfitting. Although both 355 PCA and Isomap exhibited a rise in entropy, this increase occurred after the loss had plateaued. 356 In contrast, MM-PHATE demonstrated an earlier rise in entropy, aligning with the transition into 357 the plateau and thus more accurately capturing the underlying dynamics. On the other hand, t-SNE 358 struggled to identify the increase in retained information corresponding to model overfitting and was 359 generally inadequate in capturing the global structural dynamics evident in other techniques. 360

Intra-step entropy reflects how much input information the network retains during training. As the 361 model is trained and its latent representations evolve, intra-step entropy should change accordingly. 362 MM-PHATE clearly shows these changes in the first 100 training epochs, aligning with changes 363 in validation accuracy, whereas traditional methods do not exhibit such responsiveness. Moreover, 364 as the model converges, the latent representations and intra-step entropy should stabilize. In later 365 epochs, only MM-PHATE demonstrates this desired stabilization, unlike PCA, t-SNE, and Isomap. 366 Moreover, each time step in the input contains different information, leading the network to learn in 367 a unique manner at each step. Consequently, intra-step entropy should vary across time steps. MM-368 PHATE effectively captures this variation, while PCA and Isomap show uniform entropy changes, 369 failing to reflect the unique learning dynamics across time steps.

370

324

325

326

327

328

330

331

332

333

334

338 339

 371
 5.1.2
 INTER-STEP ENTROPY

To further assess the quality of our embedding, we quantified the inter-step entropy of hidden unit activations across various time-steps ω (Fig. 4a). This metric measures the entropy between a unit's activations at different time-steps, providing insights into the temporal dynamics of each hidden unit. Our analysis identified distinct patterns among the hidden units. Notably, certain units exhibited significantly higher inter-step entropy, indicative of higher sensitivity to input changes over time and their potential role in capturing complex dependencies. These units peaked around epoch 80,



aligning with the peak in intra-step entropy at later time-steps. This suggests that these units are largely responsible for the increase in mutual information at these time-steps.

Figure 4: a) Inter-step entropy of each hidden unit across epochs, alongside model accuracy (left) and loss (right). b) Clusters of hidden units from the model's inter-step entropy trajectories across epochs. Left: Trajectories of the units in cluster 0 (12 units) and cluster 1 (8 units). Right: cluster 398 center trajectories across epochs. c-d) Confusion matrices of clusters on training and validation data. 399

400 To validate the insights from our entropy analyses and confirm the significance of the units' distinct learning behaviors with respect to the learned representation and model performance, we clustered 401 the hidden units into 2 groups, expecting one group to showcase the higher inter-step entropy. We 402 then analyzed each group's specific properties. For clustering, we used k-means clustering with the 403 Dynamic Time Warping (DTW) metric to group the hidden units into clusters based on their inter-404 step entropy trajectories (Fig. 4b). DTW determines an optimal match between two time series, 405 making it more suitable as a distance metric for our sequence data. Clustering revealed distinct 406 trends in inter-step entropy trajectories across epochs. Cluster 1 (8 units) corresponds to the units 407 previously identified as exhibiting higher inter-step entropy, and units in cluster 0 (12 units) showed 408 a sharp decrease in inter-step entropy around epoch 100. Following the clustering, we built new 409 networks comprised of the post-training weighted units contained in one particular cluster, forming 410 two sub-networks of the original network architecture to analyze their learning capabilities. By 411 computing confusion matrices for each sub-network and comparing their predictive performance on both training and validation data, we assessed the quality of their learned representation (Fig. 4c-d). 412 Interestingly, despite having fewer units, cluster 1 performed significantly better with both training 413 and validation data. This outcome aligns with our previous prediction, where cluster 1 was identified 414 as learning more complex dependencies from the input. This differential learning capability between 415 clusters demonstrates the utility of our visualization and clustering approach in revealing critical 416 differences in how information is processed and represented within the network. 417

Previous studies have indicated that the clustering property of hidden unit activation is crucial for 418 understanding the quality of the learned representation (Su & Shlizerman, 2020; Oliva & Lago-419 Fernández, 2021; Ming et al., 2017). Our findings affirm the importance of capturing each unit's 420 temporal dynamics, their community structure, as well as the evolution of these structures across 421 time-steps and epochs. 422

Comparisons with other dimensionality reduction techniques highlighted their limitations (Fig. S8). 423 PCA and Isomap failed to capture subtle dynamic variations, particularly in early epochs before the 424 onset of overfitting. t-SNE displayed a general trend corresponding to each performance transition, 425 however, the trajectories are noisy. More importantly, all these methods failed to differentiate indi-426 vidual units' learning behavior or to capture their community structure effectively. This emphasizes 427 the superior performance of MM-PHATE in maintaining the integrity of the hidden dynamics. 428

In the supplemental material section, we present results for different RNN architectures, namely 429 GRU (D.2) and Vanilla RNN (D.3), both with 20 hidden units. Using the Area2Bump dataset, we 430 additionally tried different LSTM sizes (10, 20, 30, 40, 50) (D.4). We observed consistent visualiza-431 tions when varying network size. Changing the learning rate helped confirm that the visualization

378

379

380

381

382

384 385

386

387

392

indeed reflects the model learning, i.e. the resulting change of entropies should always follow the change of model performance.

434 435 436

5.2 ANALYSIS WITH HUMAN ACTIVITY RECOGNITION MODEL

437 We next considered an action recognition dataset. We trained 30-unit LSTM network designed for 438 kinematics-based Human Activity Recognition (HAR) (Reyes-Ortiz et al., 2012). The model was 439 trained for 1000 epochs, achieving a final validation accuracy of 84% (Fig. 5). We applied MM-440 PHATE to the tensor of hidden unit activations and repeated our analysis as in the Area2Bump 441 dataset. Similar to the patterns observed in the Area2Bump LSTM network, the HAR network displayed an increase in intra-step entropy with the onset of overfitting. Particularly noteworthy was the 442 gradual increase in entropy across time-steps before epoch 300. However, after this epoch, entropy 443 at later time-steps significantly dropped, while entropy at early time-steps remained stable through-444 out the training epochs. This indicates a reduction in mutual information between the input and 445 the hidden states as the model processed more input over time. This behavior aligns with findings 446 by Farrell et al. (2022), who reported similar dynamics of information expansion and compression 447 across time-steps in trained RNNs. Their results further elucidated how gradient-based learning 448 mechanisms contribute to the development of robust representations by balancing these processes 449 of expansion and compression. 450

Further insights were gained from the inter-step entropy analysis, which indicated that the entropy of 451 many hidden units began to rise significantly around epoch 300, coinciding with a decrease in intra-452 step entropy and an improvement in model performance. As previously discussed, high inter-step 453 entropy indicates that the units are more sensitive to changes in input across time and may be crucial 454 for the model to learn complex dependencies. This observation also implies that despite the increase 455 in time-dependent information retained by many hidden units, the network successfully compresses 456 this information in its overall learned representation across epochs. Such compression seems ben-457 eficial to model performance and aligns with the principles in information bottleneck theory. This 458 theory states that deep network learning involves a fitting phase followed by a compression phase, during which useful information is distilled from the input to enhance generalization (Cheng et al., 459 2019; Tishby & Zaslavsky, 2015; Butakov et al., 2023). 460

461 Building on the discussion of information compression in the HAR model, it is pertinent to question 462 why the Area2Bump model does not exhibit a similar compression phase. Cheng et al. Cheng et al. 463 (2019) suggest that models with insufficient generalizability often fail to demonstrate a compression phase in practice, especially when applied to complex datasets with relatively simple network 464 architectures. Given that our models have comparable structures, we investigated the complexity of 465 the two datasets by analyzing the number of principal components (PCs) required to explain 95% of 466 the variance. We find that the Area2Bump data is significantly more complex than the HAR data, 467 necessitating 35 PCs compared to only 6 for the HAR data (Fig. 6). This difference in complexity 468 likely accounts for the absence of a noticeable compression phase in the Area2Bump model. 469



Figure 5: a) MM-PHATE visualization of a 30-unit LSTM network trained on HAR data, colored by epoch (left) and time-step (right). b) Intra-step entropy, c) and inter-step entropy.

Our analysis demonstrates that information compression occurs across both time-steps and epochs in
 RNNs, closely aligning with performance improvements. These results affirm the practical utility of
 the information bottleneck theory in RNNs and confirm that MM-PHATE effectively reveals detailed
 insights into the model's hidden representation and its evolution. Further research should investigate

9

480 481 482

479

470

the implications and interactions of various compression dynamics in RNNs, which could lead to
 more robust and generalizable network architectures.



Figure 6: PCA on the 2 datasets. The Area2Bump data requires 35 PCs (left) to cover 95% of the variance. The HAR data (right) requires 6 PCs to cover 95% of variance.

6 CONCLUSION

489

490

491

492

493

494

495

496

497 498 499

500

501 In this paper, we introduced MM-PHATE, a novel dimensionality reduction technique designed 502 to visualize the hidden dynamics of RNNs during training. MM-PHATE captures the evolution of these dynamics across both time-steps and epochs, offering insights beyond traditional metrics 504 like accuracy and loss curves, as well as other commonly used dimensionality reduction methods. 505 Our entropy-based analysis demonstrates that MM-PHATE can reveal distinct learning behaviors and the roles of hidden units in information flow, aligning with principles from the information 506 bottleneck theory. This approach is especially valuable in data-limited settings, as it does not rely 507 on external validation data. Moreover, we demonstrate the utility of analyzing the hidden state 508 dynamics throughout the model's learning trajectory, which provides deeper insights into the internal 509 learning processes and the evolving structure of the representation space, facilitating a more nuanced 510 understanding of how the model captures and processes information over time. 511

Despite its strengths, our approach rests on several assumptions. Notably, we assume continu-512 ity in time and over epochs, as encoded in the structured graph kernel. This assumption may be 513 violated in cases involving large learning rates, multiple significant restarts, or discontinuous acti-514 vation functions, which could result in less informative visualizations. While these scenarios are 515 not commonly explored in the literature, and PHATE has been proven to be robust against sub-516 sampling of data points, their impact should be considered (Moon et al., 2019). Additionally, while 517 MM-PHATE is based on a computationally efficient implementation of PHATE that has shown 518 better efficiency than traditional dimensionality reduction methods, the memory complexity of the 519 multiway-multislice kernel introduces a bottleneck in scalability for larger architectures or datasets. 520 Future work could address this by developing methods for sub-sampling the kernel or utilizing more 521 memory-efficient algorithms without compromising the overall structure (Holtz et al., 2022), such 522 as graph partitioning and merging of data points used by Kuchroo et al. (2022) in their Multiscale PHATE. Moreover, MM-PHATE does not currently leverage the internal structure of RNNs, such as 523 attention mechanisms present in transformers (Vaswani et al., 2017). Future research will explore 524 extending MM-PHATE to analyze transformers, which are increasingly dominant in sequential data 525 analysis. However, RNNs remain valuable across various fields (Deo et al., 2024), especially when 526 working with limited datasets where transformers may not be practical. Thus, while MM-PHATE 527 is built for RNNs, its future adaptations could provide deeper insights into transformer architectures 528 and their hidden dynamics. 529

530 531

532 533

ETHICS STATEMENT

This paper presents work whose goal is to advance the field of Machine Learning. The paper has no foreseeable societal impact.

534 535 536

538

REFERENCES

539 Omri Barak. Recurrent neural networks as versatile tools of neuroscience research. *Current Opinion in Neurobiology*, 46:1–6, October 2017. ISSN 09594388. doi: 10.1016/j.conb.2017.06.003.

- Ivan Butakov, Aleksander Tolmachev, Sofia Malanchuk, Anna Neopryatnaya, Alexey Frolov, and Kirill Andreev. Information bottleneck analysis of deep neural networks via lossy compression, May 2023. URL http://arxiv.org/abs/2305.08013.
- Yuexing Chen and Jiarun Li. Recurrent neural networks algorithms and applications. In 2021 2nd International Conference on Big Data & Artificial Intelligence & Software Engineering (ICBASE), pp. 38–43, Zhuhai, China, September 2021. IEEE. ISBN 978-1-66542-709-8. doi: 10.1109/ICBASE53849.2021.00015. URL https://ieeexplore.ieee.org/ document/9696155/.
- Hao Cheng, Dongze Lian, Shenghua Gao, and Yanlin Geng. Utilizing information bottleneck to
 evaluate the capability of deep neural networks for image classification. *Entropy*, 21(5):456, May 2019. ISSN 1099-4300. doi: 10.3390/e21050456.
- Raeed Chowdhury, Joshua Glaser, and Lee Miller. Data from: Area 2 of primary somatosensory cortex encodes kinematics of the whole arm, 2022. URL https://doi.org/10.5061/dryad.nk98sf7q7. Dataset.
- Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling, December 2014. URL http://arxiv.org/abs/1412.3555. arXiv:1412.3555 [cs].
- Ronald R. Coifman and Stéphane Lafon. Diffusion maps. *Applied and Computational Harmonic Analysis*, 21(1):5–30, July 2006. ISSN 10635203. doi: 10.1016/j.acha.2006.04.006.
- Darrel R. Deo, Francis R. Willett, Donald T. Avansino, Leigh R. Hochberg, Jaimie M. Henderson, and Krishna V. Shenoy. Brain control of bimanual movement enabled by recurrent neural networks. *Scientific Reports*, 14(1):1598, January 2024. ISSN 2045-2322. doi: 10.1038/s41598-024-51617-3. URL https://www.nature.com/articles/s41598-024-51617-3.
- Matthew Farrell, Stefano Recanatesi, Timothy Moore, Guillaume Lajoie, and Eric Shea-Brown.
 Gradient-based learning drives robust representations in recurrent neural networks by balancing compression and expansion. *Nature Machine Intelligence*, 4(6):564–573, June 2022. ISSN 2522-5839. doi: 10.1038/s42256-022-00498-0.
- Scott Gigante, Adam S. Charles, Smita Krishnaswamy, and Gal Mishne. Visualizing the
 PHATE of neural networks, August 2019. URL http://arxiv.org/abs/1908.02831.
 arXiv:1908.02831 [cs, stat].
- Klaus Greff, Rupesh Kumar Srivastava, Jan Koutník, Bas R. Steunebrink, and Jürgen Schmidhuber.
 LSTM: A search space odyssey. *IEEE Transactions on Neural Networks and Learning Systems*, 28(10):2222–2232, October 2017. ISSN 2162-237X, 2162-2388. doi: 10.1109/TNNLS.2016.
 2582924. arXiv:1503.04069 [cs].
- Hansika Hewamalage, Christoph Bergmeir, and Kasun Bandara. Recurrent neural networks for time series forecasting: Current status and future directions. *International Journal of Forecasting*, 37 (1):388–427, January 2021. ISSN 01692070. doi: 10.1016/j.ijforecast.2020.06.008.
- Akinori Hidaka and Takio Kurita. Consecutive dimensionality reduction by canonical correlation analysis for visualization of convolutional neural networks. In *Proceedings of the ISCIE international symposium on stochastic systems theory and its applications*, volume 2017, pp. 160–167. The ISCIE Symposium on Stochastic Systems Theory and Its Applications, 2017.
- Sepp Hochreiter and Jürgen Schmidhuber. Lstm can solve hard long time lag problems. Advances
 in neural information processing systems, 9, 1996.
- Chester Holtz, Gal Mishne, and Alexander Cloninger. Evaluating disentanglement in generative models without knowledge of latent factors. In *Topological, Algebraic and Geometric Learning Workshops 2022*, pp. 161–171. PMLR, 2022.
- 592 Chang Woo Hong, Changmin Lee, Kwangsuk Lee, Min-Seung Ko, Dae Eun Kim, and Kyeon Hur.
 593 Remaining useful life prognosis for turbofan engine using explainable deep neural networks with dimensionality reduction. *Sensors*, 20(22):6626, 2020.

595

596

631

632

633 634

635

636

641

Andrej Karpathy, Justin Johnson, and Li Fei-Fei. Visualizing and understanding recurrent networks, November 2015. URL http://arxiv.org/abs/1506.02078. arXiv:1506.02078 [cs].

- Manjot Kaur and Aakash Mohta. A review of deep learning with recurrent neural network. 597 In 2019 International Conference on Smart Systems and Inventive Technology (ICSSIT), pp. 598 460-465, Tirunelveli, India, November 2019. IEEE. ISBN 978-1-72812-119-2. doi: 10. 1109/ICSSIT46314.2019.8987837. URL https://ieeexplore.ieee.org/document/ 600 8987837/. 601
- 602 Yassin Khalifa, Danilo Mandic, and Ervin Sejdić. A review of hidden markov models and recurrent 603 neural networks for event detection and localization in biomedical signals. Information Fusion, 69:52-72, May 2021. ISSN 15662535. doi: 10.1016/j.inffus.2020.11.008. 604

605 Manik Kuchroo, Jessie Huang, Patrick Wong, Jean-Christophe Grenier, Dennis Shung, Alexander Tong, Carolina Lucas, Jon Klein, Daniel B. Burkhardt, Scott Gigante, Abhinav Godavarthi, 607 Bastian Rieck, Benjamin Israelow, Michael Simonov, Tianyang Mao, Ji Eun Oh, Julio Silva, 608 Takehiro Takahashi, Camila D. Odio, Arnau Casanovas-Massana, John Fournier, Yale IMPACT 609 Team, Abeer Obaid, Adam Moore, Alice Lu-Culligan, Allison Nelson, Anderson Brito, Angela 610 Nunez, Anjelica Martin, Anne L. Wyllie, Annie Watkins, Annsea Park, Arvind Venkataraman, 611 Bertie Geng, Chaney Kalinich, Chantal B. F. Vogels, Christina Harden, Codruta Todeasa, Cole Jensen, Daniel Kim, David McDonald, Denise Shepard, Edward Courchaine, Elizabeth B. White, 612 Eric Song, Erin Silva, Eriko Kudo, Giuseppe Deluliis, Haowei Wang, Harold Rahming, Hong-613 Jai Park, Irene Matos, Isabel M. Ott, Jessica Nouws, Jordan Valdez, Joseph Fauver, Joseph Lim, 614 Kadi-Ann Rose, Kelly Anastasio, Kristina Brower, Laura Glick, Lokesh Sharma, Lorenzo Se-615 wanan, Lynda Knaggs, Maksym Minasyan, Maria Batsu, Maria Tokuyama, M. Cate Muenker, 616 Mary Petrone, Maxine Kuang, Maura Nakahata, Melissa Campbell, Melissa Linehan, Michael H. 617 Askenase, Michael Simonov, Mikhail Smolgovsky, Nathan D. Grubaugh, Nicole Sonnert, Nida 618 Naushad, Pavithra Vijayakumar, Peiwen Lu, Rebecca Earnest, Rick Martinello, Roy Herbst, 619 Rupak Datta, Ryan Handoko, Santos Bermejo, Sarah Lapidus, Sarah Prophet, Sean Bickerton, 620 Sofia Velazquez, Subhasis Mohanty, Tara Alpert, Tyler Rice, Wade Schulz, William Khoury-621 Hanold, Xiaohua Peng, Yexin Yang, Yiyun Cao, Yvette Strong, Shelli Farhadian, Charles S. 622 Dela Cruz, Albert I. Ko, Matthew J. Hirn, F. Perry Wilson, Julie G. Hussin, Guy Wolf, Akiko Iwasaki, and Smita Krishnaswamy. Multiscale PHATE identifies multimodal signatures of covid-623 19. Nature Biotechnology, 40(5):681–691, May 2022. ISSN 1087-0156, 1546-1696. doi: 624 10.1038/s41587-021-01186-x. 625

- 626 Bum Chul Kwon, Min-Je Choi, Joanne Taery Kim, Edward Choi, Young Bin Kim, Soonwook Kwon, 627 Jimeng Sun, and Jaegul Choo. RetainVis: Visual analytics with interpretable and interactive 628 recurrent neural networks on electronic medical records. IEEE Transactions on Visualization and Computer Graphics, 25(1):299–309, January 2019. ISSN 1077-2626, 1941-0506, 2160-9306. doi: 10.1109/TVCG.2018.2865027. arXiv:1805.10724 [cs, stat]. 630
 - Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets, June 2018. URL https://arxiv.org/abs/1712.09913.
 - Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. Visualizing and understanding neural models in NLP, January 2016. URL http://arxiv.org/abs/1506.01066. arXiv:1506.01066 [cs].
- 637 Zachary C. Lipton, John Berkowitz, and Charles Elkan. A critical review of recurrent neural net-638 works for sequence learning, October 2015. URL http://arxiv.org/abs/1506.00019. 639 arXiv:1506.00019 [cs]. 640
- Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. Journal of Machine Learning Research, 9(86):2579–2605, 2008. 642
- 643 Andrzej Maćkiewicz and Waldemar Ratajczak. Principal components analysis (PCA). Com-644 puters & Geosciences, 19(3):303-342, 1993. ISSN 0098-3004. doi: https://doi.org/10.1016/ 645 0098-3004(93)90090-R. 646
- Tomas Mikolov, Armand Joulin, Sumit Chopra, Michael Mathieu, and Marc'Aurelio Ranzato. 647 Learning longer memory in recurrent neural networks. arXiv preprint arXiv:1412.7753, 2014.

662

677

685

687

689

690

691

692

693

694

696

- 648 Yao Ming, Shaozu Cao, Ruixiang Zhang, Zhen Li, Yuanzhe Chen, Yangqiu Song, and Huamin 649 Qu. Understanding hidden memories of recurrent neural networks. In 2017 IEEE Conference on 650 Visual Analytics Science and Technology (VAST), pp. 13–24, Phoenix, AZ, October 2017. IEEE. ISBN 978-1-5386-3163-8. doi: 10.1109/VAST.2017.8585721. URL https://ieeexplore. 651 652 ieee.org/document/8585721/.
- Kevin R. Moon, David Van Dijk, Zheng Wang, Scott Gigante, Daniel B. Burkhardt, William S. Chen, 654 Kristina Yim, Antonia Van Den Elzen, Matthew J. Hirn, Ronald R. Coifman, Natalia B. Ivanova, 655 Guy Wolf, and Smita Krishnaswamy. Visualizing structure and transitions in high-dimensional 656 biological data. Nature Biotechnology, 37(12):1482-1492, December 2019. ISSN 1087-0156, 657 1546-1696. doi: 10.1038/s41587-019-0336-3. 658
- 659 Christian Oliva and Luis F. Lago-Fernández. Stability of internal states in recurrent neural networks 660 trained on regular languages. *Neurocomputing*, 452:212–223, September 2021. ISSN 09252312. doi: 10.1016/j.neucom.2021.04.058. 661
- James O. Ramsay. Some statistical considerations in multidimensional scaling. ETS Research 663 Bulletin Series, 1966(1):i-96, 1966. doi: https://doi.org/10.1002/j.2333-8504.1966.tb00976.x. 664 URL https://onlinelibrary.wiley.com/doi/abs/10.1002/j.2333-8504. 665 1966.tb00976.x. 666
- 667 Paulo E Rauber, Samuel G Fadel, Alexandre X Falcao, and Alexandru C Telea. Visualizing the 668 hidden activity of artificial neural networks. IEEE transactions on visualization and computer graphics, 23(1):101–110, 2016. 669
- 670 Jorge Reyes-Ortiz, Davide Anguita, Alessandro Ghio, Luca Oneto, and Xavier Parra. Human 671 Activity Recognition Using Smartphones. UCI Machine Learning Repository, 2012. DOI: 672 https://doi.org/10.24432/C54S4K. 673
- 674 Hojjat Salehinejad, Sharan Sankar, Joseph Barfett, Errol Colak, and Shahrokh Valaee. Recent ad-675 vances in recurrent neural networks, February 2018. URL http://arxiv.org/abs/1801. 676 01078. arXiv:1801.01078 [cs].
- Hendrik Strobelt, Sebastian Gehrmann, Hanspeter Pfister, and Alexander M. Rush. LSTMVis: A 678 tool for visual analysis of hidden state dynamics in recurrent neural networks. IEEE Transactions 679 on Visualization and Computer Graphics, 24(1):667–676, January 2018. ISSN 1077-2626. doi: 680 10.1109/TVCG.2017.2744158. 681
- 682 Kun Su and Eli Shlizerman. Clustering and recognition of spatiotemporal features through inter-683 pretable embedding of sequence to sequence recurrent neural networks. Frontiers in Artificial Intelligence, 3:70, September 2020. ISSN 2624-8212. doi: 10.3389/frai.2020.00070. 684
- David Sussillo and Omri Barak. Opening the black box: Low-dimensional dynamics in high-686 dimensional recurrent neural networks. Neural Computation, 25(3):626-649, March 2013. ISSN 0899-7667. doi: 10.1162/NECO_a_00409. 688
 - Zhiyuan Tang, Ying Shi, Dong Wang, Yang Feng, and Shiyue Zhang. Memory visualization for gated recurrent neural networks in speech recognition, February 2017. URL http://arxiv. org/abs/1609.08789. arXiv:1609.08789 [cs].
 - Joshua B. Tenenbaum, Vin De Silva, and John C. Langford. A global geometric framework for nonlinear dimensionality reduction. Science, 290(5500):2319-2323, December 2000. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.290.5500.2319.
 - Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle, March 2015. URL http://arxiv.org/abs/1503.02406. arXiv:1503.02406 [cs].
- Manuel Titos, Luz Garcia, Milad Kowsari, and Carmen Benitez. Toward knowledge extraction 699 in classification of volcano-seismic events: Visualizing hidden states in recurrent neural net-700 works. IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, 15:2311–2325, 2022. ISSN 1939-1404, 2151-1535. doi: 10.1109/JSTARS.2022.3155967.

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information Processing Systems, volume 30, 2017. URL https://arxiv.org/pdf/1706.03762. pdf.
- Zijie J. Wang, Robert Turko, Omar Shaikh, Haekyu Park, Nilaksh Das, Fred Hohman, Minsuk Kahng, and Duen Horng Polo Chau. CNN explainer: Learning convolutional neural networks with interactive visualization. *IEEE Transactions on Visualization and Computer Graphics*, 27 (2):1396–1406, February 2021. ISSN 1077-2626, 1941-0506, 2160-9306. doi: 10.1109/TVCG. 2020.3030418.
- Wei Yu, Kuiyuan Yang, Yalong Bai, Hongxun Yao, and Yong Rui. Visualizing and comparing convolutional neural networks, December 2014. URL http://arxiv.org/abs/1412.6631.
 arXiv:1412.6631 [cs].
- 715 716

719 720

721

722

723

724

725

726

727

A DATASETS

- We employed two datasets.
- The Area2Bump dataset (Chowdhury et al., 2022) consists of neural activity recorded from Brodmann's area 2 of the somatosensory cortex while macaque monkeys performed a slightly modified version of a standard center-out reaching task. Dataset license: CC0 1.0. According to the authors, data was collected consistently "with the guide for the care and use of laboratory animals and approved by the institutional animal care and use committee of Northwestern University under protocol #IS00000367". This dataset, collected from monkeys thus does not contain personally identifiable information. Since it is neural data, we do not consider it to be offensive content.
- 728Data Statistics: The dataset includes 193 samples, split into 115 training samples and 78729testing samples. Each sample consists of 600 time steps with 65 features per time step. For730the training set, the mean values across features range from 0.0001 to 0.03, with standard731deviations between 0.001 and 0.02, and maximum values ranging from 0.003 to 0.12. In732the test set, the mean values range from 0.0008 to 0.03, standard deviations from 0.0007733to 0.018, and maximum values from 0.0007 to 0.13.
- 7.34 2. The Human Activity Recognition (HAR) Using Smartphones dataset (Reyes-Ortiz et al., 2012), kinematic recordings of 30 subjects performing daily living activities with a smartphone embedded with an inertial measurement unit. Dataset license: CC BY 4.0. Information relating to participant consent was not found relating to this dataset. This dataset does not reveal participant name or identifiable information. The kinematics contained in the dataset are not considered offensive content.
- Data Details: The dataset includes six activity classes (e.g., Walking, Sitting) based on accelerometer and gyroscope data sampled at 50 Hz. The sensor data has been pre-processed with noise filtering and separated into gravitational and body motion components. Data is windowed into 2.56-second segments (128 data points) with 50% overlap, resulting in 561dimensional feature vectors per window. The dataset is split into 70% training (21 subjects, 7352 samples) and 30% testing (9 subjects, 2947 samples).
- Data Statistics: For the training set, mean values across features range from -0.0008 to 0.8, with standard deviations between 0.1 and 0.41, and values spanning from -5.97 to 5.75. For the test set, mean values range from -0.013 to 0.8, with standard deviations between 0.095 and 0.41, and values spanning from -3.43 to 3.47.
- 749 750

751

B MODEL TRAINING

We used TensorFlow's Keras API for all models training and validation.

The network in Section 5.1 was trained as follows. The Area2Bump dataset was randomly split into training and validation subsets containing an equal number of samples for each input class with an 8 to 2 ratio. Additional samples that would make the training data uneven were added back to

the validation subset to make use of all samples. The network consists of an LSTM layer with 20 units. This is followed by a Flatten layer that converts the LSTM's output into a one-dimensional vector. Finally, a Dense layer with 8 units and a softmax activation function produces the output for the 8-class classification tasks. The network was trained with a batch size of 64. We used an Adam optimizer with a learning rate of $1e^{-4}$. During the training process, we recorded the activations from the LSTM layer into the activation tensor T for visualization.

The network in Section 5.2 was trained as follows. The HAR dataset was preprocessed and split by the authors into training and validation subsets according to the subjects with a 7 to 3 ratio. The network consists of an LSTM layer with 30 units and a Dense layer with 6 units and a softmax activation function produces the output for the 6-class classification tasks. The network was trained with a batch size of 32. We used an Adam optimizer with a learning rate of $2e^{-5}$. During the training process, we recorded the activations from the LSTM layer into the activation tensor T for visualization.

- 769
- 770

C IMPLEMENTATION OF VISUALIZATION METHODS

771 772 773

781

782 783

784

785

786 787

788

794

C.1 MM-PHATE

The first 29 epochs with every 10th epoch thereafter. We utilized the M-PHATE package to construct to ur multiway multislice graphs and for the application of PHATE.

C.2 PCA

PCA was performed using the "PCA" class from the "sklearn.decomposition" package to reduce the dimensionality of the time trace tensor T—recorded during training of the Area2Bump model—to three principal components.

C.3 T-SNE

⁷⁸⁹ In this analysis, we first conducted an initial dimensionality reduction on the same time trace ⁷⁹⁰ tensor T with PCA to 15 principle components, which explains 99.93% of the variance in the ⁷⁹¹ activations. Subsequently, t-SNE with the Barnes-Hut approximation was performed using the ⁷⁹² "sklearn.manifold.TSNE" class.

C.4 ISOMAP

Due to memory constraints, we only used a subset of the tensor T for Isomap computation. Specifically, epochs were sampled using an array combining the first 29 epochs with every 10th epoch thereafter, and intrinsic steps were sampled using a linear space from 0 to the end (600), resulting in 50 evenly spaced steps. We first conducted an initial dimensionality reduction on the sampled tensor with PCA to 15 principle components using "sklearn.decomposition.PCA" package. Then, we applied Isomap using "sklearn.manifold.Isomap" class to reduce the dimensionality to 3.

802 803

804

806

D ADDITIONAL EXPERIMENTS

D.1 M-PHATE

M-PHATE was applied to the same 20-unit LSTM network trained on the Area2Bump dataset. The
 algorithm only incorporated the final state, or time-step, from each epoch. While the resulting
 visualization captures smooth transitions across epochs, it omits critical information from earlier
 time-steps. This loss of temporal resolution obscures insights into how the network processes input

sequences over time—an essential aspect for understanding RNNs, for instance, how much sequential input information is retained by the network.



Figure S1: M-PHATE on Area2Bump LSTM: Visualization of a 20-unit LSTM network trained for 200 epochs. Each point represents a hidden unit at the last time-step in a given epoch throughout the entire training process. The points are colored based on epoch (left) or hidden unit (right).

D.2 AREA2BUMP WITH GRU

Here is the same analysis as section 5.1 using GRU (Fig. S2, S3, S4). Other parameters were kept the same.

From these figures, it is evident that PCA and t-SNE present similar visualizations of the hidden dynamics across different network architectures, while MM-PHATE distinctly captures the unique learning behaviors of each model. Consistent with Section 5.1, PCA displays an increasing intrastep entropy even after model accuracy has plateaued, and t-SNE produces a noisy visualization.
In contrast, MM-PHATE uniquely aligns its transitions well with the learning curve. Notably, the GRU model's representation appears more compact and organized compared to the LSTM model, potentially reflecting its superior performance and reduced overfitting.



Figure S2: Area2Bump GRU: Visualization of a 20-unit GRU network trained for 200 epochs. Each point represents a hidden unit at a specific time-step in a given epoch throughout the entire training process. The visualizations are generated using MM-PHATE, PCA, and t-SNE, from left to right, respectively. Points are colored based on epoch (top row) or time-step (bottom-row)



Figure S4: Area2Bump GRU: Inter-step entropy of all hidden units in embedding space of the Area2Bump model at each time-step in each epoch, compared to training and accuracies (top) and losses (bottom). From left to right, the dimensionality reduction metrics used are MM-PHATE, PCA, and t-SNE.

- 912 913
- 914
- 915
- 916
- 917

918 D.3 AREA2BUMP WITH VANILLA RNN

Here is the same analysis as section 5.1 using vanilla RNN (Fig. S5, S6, S7). Other parameters werekept the same.

From these figures, it is evident that regardless of the network architectures, PCA exhibits a revolving pattern with overly smooth transitions across epochs, while t-SNE produces a noisy visualization. In contrast, the MM-PHATE visualization reveals that the Vanilla RNN displays a more chaotic pattern compared to the LSTM and GRU models, which is likely associated with its reduced performance and increased overfitting. Furthermore, the intra-step entropies of MM-PHATE show reduced variation across time-steps, indicating that the model struggles to process the input data effectively to generate meaningful representations.



Figure S5: Area2Bump Vanilla: Visualization of a 20-unit Vanilla RNN trained for 200 epochs. Each point represents a hidden unit at a specific time-step in a given epoch throughout the entire training process. The visualizations are generated using MM-PHATE, PCA, and t-SNE, from left to right, respectively. Points are colored based on epoch (top row) or time-step (bottom-row)



Figure S6: Area2Bump Vanilla: Intra-step entropy of all hidden units in embedding space at each time-step in each epoch, compared to training and validation accuracy (top) and losses (bottom), comparing embeddings of MM-PHATE, PCA, and t-SNE.

D.4 AREA2BUMP WITH LSTM OF VARIOUS SIZES

Here we repeat the same MM-PHATE analysis as in section 5.1 using LSTM of various sizes (Fig. S9, S10, S11). Other parameters were kept the same. These results demonstrate that MM-PHATE consistently captures smooth yet distinct transitions across epochs and time steps, regard-less of the LSTM network size. The intra- and inter-step entropy analyses further reveal that these transitions closely correlate with performance changes throughout training. Specifically, we observe a general increase in intra-step entropy as models begin to overfit, suggesting that the networks increasingly memorize input information. In contrast, inter-step entropy shows a significant decline



Figure S8: Area2Bump: Inter-step entropy of all hidden units in embedding space of the Area2Bump model at each time-step in each epoch, compared to training and accuracies (top) and losses (bottom). From left to right, the dimensionality reduction metrics used are MM-PHATE, PCA, t-SNE, and Isomap.

1022

- 1022
- 1024
- 1025

as overfitting progresses, reflecting a loss of sensitivity to input changes over time. The loss curve shows worse overfitting as the network size increases, and the networks' inter-step entropy becomes less structured.



Figure S9: Area2Bump LSTM: MM-PHATE visualization of networks of size 10 to 50 (left to right). Each point represents a hidden unit at a specific time-step in a given epoch. Points are colored based on epoch (top) or timestep (bottom).



Figure S10: Area2Bump LSTM: Intra-step entropy of all hidden units in MM-PHATE embedding space at each time-step in each epoch, compared to accuracies (top) and losses (bottom). Network sizes range from 10 to 50 (left to right).





¹⁰⁸⁰ E COMPUTING INFRASTRUCTURE

All but the t-SNE computation was done on a 14-core laptop running Windows 11 Home with a NVIDIA GeForce RTX 3070 Ti Laptop graphics card and 40GB of RAM. The t-SNE was done on a single 95-core internal cluster running Ubuntu 18.04.6 LTS with 10 Quadro RTX 5000 graphics cards and 755GB of RAM.

F MATHEMATICAL NOTATIONS

1089		
1090	Notation	Definition
1091	x	Point in high dimensional space
1092		Euclidean distance matrix between all data points x
1093	\mathbf{R}	Animity kernel matrix k pearest neighbor distance of r
1094	$\epsilon_k(x_i)$	h -inclust-incignool distance of x_i
1095	α P	Diffusion operator
1096	D I	Diagonal matrix of row sums of K
1097	$\overset{oldsymbol{ extsf{P}}}{oldsymbol{P}^t}$	Transition probabilities of a diffusion process over t steps
1098	'n	Total number of epochs the network is trained for
1099	$oldsymbol{F}$	Feed-forward neural network
1100	m	Total number of hidden units in the network
1101	X	Training data, subset of Π
1102	Π	Larger dataset
1103	T	Activation tensor
1104	Y	Input data, subset of X with equal number of samples per class
1105	(-) p	Number of samples in Y
1106	$oldsymbol{K}_{ ext{intraslice}}^{(au)}(i,j)$	Intraslice affinities between pairs of hidden units within an epoch τ
1107	$oldsymbol{K}_{ ext{interslice}}^{(i)}(au, arcell)$	Interslice affinities between a hidden unit <i>i</i> and itself at different epochs
1108	$\sigma_{(au,i)}$	Intraslice bandwidth for unit <i>i</i> in epoch τ
1109	ϵ	Fixed interslice bandwidth
1110	au	Index for given epoch
1111	i, j	Index for given hidden unit
1112	h_t	RNN hidden state at time step t
1112	W	RNN weights
1114	0	KINN DIASES DNN output at time step t
1115	g_t	RNN activation function
1116	$\overset{j}{B}$	Recurrent neural network
1117		Index for given time step
1110	s	Total number of time steps in the RNN
1110	$\boldsymbol{K}^{(au,\omega)}_{i}$ (i,j)	Intrasten affinities between hidden units i and i at time-sten ω in enoch τ
1120	\mathbf{T} intrastep (i, j)	La sep a main des setween maden ands 7 and 7 at time step a mepoen 7
1120	$\mathbf{K}_{\text{interstep}}((\tau,\omega),(\eta,\nu))$	Interstep affinities between a nidden unit <i>i</i> and itself at different time-steps and epochs
1121	$\sigma_{(\tau,\omega,i)}$	Intrastep bandwidth for unit i at time-step w and epoch τ
1122	κ Γ	I abels of \mathbf{X}
1123	L	
1124		Table 1: Notations
1125		
1120		
1127		
1120		
1129		
1130		
1131		
1132		
1133		