
CHARTING AND NAVIGATING THE SPACE OF SOLUTIONS FOR RECURRENT NEURAL NETWORKS

SUPPLEMENTARY MATERIAL

Anonymous Author(s)

Affiliation

Address

email

1 2D RNN

2 1.1 Training methods

3 The four parameters, $W \in \mathbb{R}^{2 \times 2}$, were iid sampled from a uniform distribution $\mathcal{U}(-1.5, 1.5)$. We
4 implemented the continuous time dynamics of the RNN in PyTorch[1] using the package torchdiffeq[2,
5 3].

6 Network parameters were optimized using gradient descent with no momentum and a learning rate
7 of 0.03. During each step of learning, the network dynamics were simulated for a single trajectory
8 (Equation 1 in main text), and the loss \mathcal{L} (Eq 2 in main text) was used to compute the gradient. We then
9 normalised the gradient by its Frobenius norm and scaled it by $\sqrt{\mathcal{L}}$ before updating the parameters.
10 This is a heuristic choice to motivate convergence of learning even when gradients are small for some
11 pathological initializations. Networks were trained for 2000 epochs, and only RNNs with $\mathcal{L} < 10^{-5}$
12 were accepted as solutions.

13 1.2 Different nonlinearity

14 To examine the effect of training hyperparameters on the space of solutions, we used $\phi := \text{ReLU}$
15 instead of $\phi := \tanh$ that was used in the main text. We find that this choice indeed leads to different
16 solution types. Specifically, ReLU RNNs did not converge to limit-cycle solutions. Some converged
17 to non-zero fixed points, accompanied by a saddle point at the origin, as in the main text. In addition,
18 two other solution types arised in this setting. A stable origin with large transient amplification, as in
19 the yellow curve of Figure 1), and a diverging trajectory, shown by the dark curve in the same figure.
20 The effect of the different non-linearity is also seen in the distribution of trace and determinants of
21 the solutions (Figure 2), where limit cycles are absent for ReLU, and stable solutions (bottom-right
22 quadrant) are absent for tanh.

23 2 Timing task

24 2.1 Training process

25 2.1.1 Network architecture

26 We studied three different RNN architectures and their exact equations are all summarized below.
27 The trained parameters are the weights W and biases b . The function $\sigma(z) = (1 + \exp(-z))^{-1}$ is
28 the sigmoid function, $h_t \in \mathbb{R}^N$ and $u_t \in \{0, 1\}^2$ are the state and the input at time t .

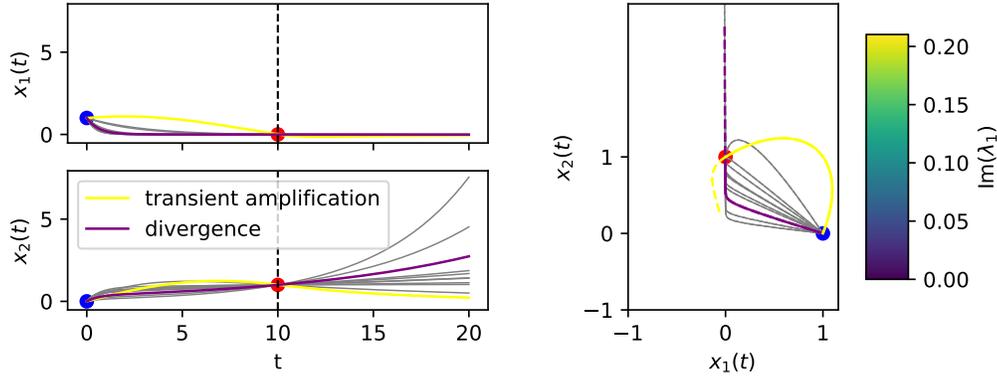


Figure 1: Trajectories of several 2D RNN solutions for $\phi := \text{ReLU}$.

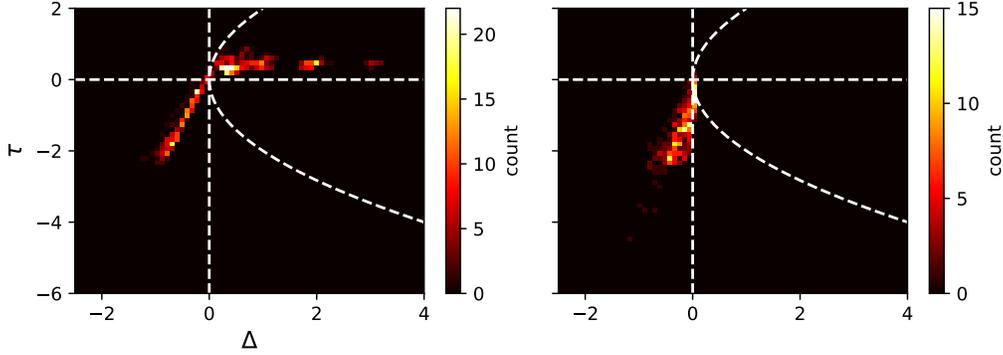


Figure 2: A 2D histogram of trace τ and determinant Δ of 2D RNNs solutions for the task specified in the main text. $\phi := \tanh$ (left) and $\phi := \text{ReLU}$ (right).

Vanilla [4]

$$h_t = \tanh(W_{ih}u_t + b_{ih} + W_{hh}h_{t-1} + b_{hh}) \quad (1)$$

GRU [5]

$$r_t = \sigma(W_{ir}u_t + b_{ir} + W_{hr}h_{t-1} + b_{hr}) \quad (2)$$

$$z_t = \sigma(W_{iz}u_t + b_{iz} + W_{hz}h_{t-1} + b_{hz}) \quad (3)$$

$$n_t = \tanh(W_{in}u_t + b_{in} + r_t * (W_{hn}h_{t-1} + b_{hn})) \quad (4)$$

$$h_t = (1 - z_t) * n_t + z_t * h_{t-1} \quad (5)$$

LSTM [6]

$$i_t = \sigma(W_{ii}u_t + b_{ii} + W_{hi}h_{t-1} + b_{hi}) \quad (6)$$

$$f_t = \sigma(W_{if}u_t + b_{if} + W_{hf}h_{t-1} + b_{hf}) \quad (7)$$

$$g_t = \tanh(W_{ig}u_t + b_{ig} + r_t * (W_{hg}h_{t-1} + b_{hg})) \quad (8)$$

$$o_t = \sigma(W_{io}u_t + b_{io} + W_{ho}h_{t-1} + b_{ho}) \quad (9)$$

$$c_t = f_t * c_{t-1} + i_t * g_t \quad (10)$$

$$h_t = o_t * \tanh(c_t) \quad (11)$$

29 The units had $N = 20, \dots, 50$ hidden neurons and the output of the network at every time-step is an
30 affine readout of the internal state. h_0 was always initialized to zero.

31 **2.1.2 Task and trial structure**

32 Each trial was comprised of seven consecutive epochs, as demonstrated in Figure 3. The *Ready* pulse
 33 was given after 20 – 30 steps. Both inputs and the required output were binary sequences with ones
 34 during each pulse (10 steps long) and zero elsewhere. When working with intervals from the range
 35 $[t_s^{min}, t_s^{max}]$, the length of all trials was set to $2 * t_s^{max} + 100$. This allowed the network time to
 relax back to rest for at least 70 steps after emitting a *Go* pulse. The training set always included 512

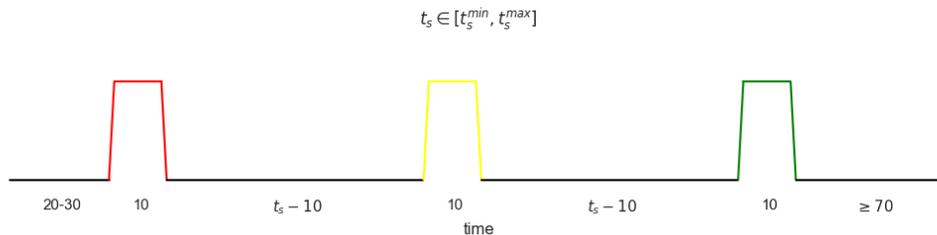


Figure 3: Ready-Set-Go timeline

36 random trials so, on average, every interval was included more than 5 times.
 37

38 **2.1.3 Training protocol**

39 All networks were trained using *Adam* [7] for 10000 epochs with a batch size of 64 and a decaying
 40 learning rate starting from $1e - 3$ up until $1e - 4$. Unless stated otherwise, the training set was
 41 comprised of 512 trials and their order was shuffled at the beginning of each epoch. We estimated the
 42 network’s performance with mean squared error (MSE), and training was halted when the minimal
 43 threshold of 10^{-5} was achieved over the training set.

44 **2.2 Feature extraction**

45 The feature extraction process in this work can be divided into two. First, we describe how we
 46 extracted numerical features from the neural activity during training. Later, we describe how we
 47 extracted topological features by analyzing the network dynamics outside the training set. Often in
 48 our analysis we quantified the neural velocity in phase space. For that purpose, we used the scalar
 49 function

$$q(h_t) = \|h_{t+1} - h_t\|_2, \quad t \in \mathbb{N} \tag{12}$$

50 that was introduced in [8]. This function estimates the distance that the network crosses in a single
 51 time-step, given its current location.

52 **2.3 Neural features**

53 As described in the main text, we extracted various features from the neural activity during the
 54 training set. These were related to the major dynamical objects and task epochs. Below, we explain
 55 and define the features according to the relevant epoch.

56 **Ready-Set features** As seen in the main text (FIGURE WITH CIRCULAR TRAJECTORY), the
 57 shape of the Ready-Set trajectory can indicate whether the network will eventually converge to a
 58 limit cycle. We thus considered the minimal and maximal curvature, the speed at its end, and the
 59 ratio between its initial and final speed. All these features were measured on a logarithmic scale.

60 **Set-Go features** The Set-Go manifold was defined by collecting the network states corresponding
 61 to trials of all delays ($t_s \in [30, 120]$), and using the time points from 10 after the Set pulse until 10
 62 before the Go pulse. Because this is a two-dimensional manifold (time by trials), we calculated the
 63 aspect ratio as follows. The nominator was the cumulative length of the trajectory corresponding
 64 to the initial states across all trials. The denominator was the length of the full trajectory of the
 65 longest trial ($t_s = 120$). Similarly, we extracted the aspect ratio with respect to the final states of the
 66 Set-Go manifold. We also measured how this ratio changes as a function of time in the following

67 manner.. Later, we calculated how the length of the Set-Go trajectory changes as a function of the
68 time interval that is being encoded, by fitting a linear regressor to the mapping $t_s \rightarrow \|\text{Set-Go}(t_s)\|_2^2$
69 and extracting its slope as a feature. To account for whether time-coding is concentrated in single
70 neurons or distributed across the population, we measured the following quantity. For each neuron,
71 we considered all points on the Set-Go manifold. We used linear regression to map the activity of the
72 neuron to the time remaining until the Go pulse. The minimal error across all neurons was used as a
73 feature.

74 **Ready-Set & Set-Go features** Here, we focused on the relationship between the trajectory
75 $\text{Ready-Set}(t_s^{max})$ and $\text{Set-Go}(t_s^{max})$. We extracted as features the Pearson correlation and the angle
76 between them, the ratio between their speeds, and the width of their separating hyper-plane obtained
77 from Linear SVM.

78 Figure 4 shows the density-histogram of each feature for each architecture.

79 2.4 Topological

80 As we discussed in the main text, defining what is a solution is not trivial. In the context of dynamical
81 systems, obtaining a qualitative description of the phase space is often enough. However, this
82 description requires full knowledge of the dynamical objects, which is often inaccessible. Particularly,
83 the transient nature of the Ready-Set-Go timing task renders the dynamical landscape less structured
84 and more difficult to analyze with classical dynamical systems tools. Therefore, to understand the
85 mechanism of the network we will suffice in identifying the key areas of the dynamics and the
86 transitions between them. In the RSG task, following the Ready cue the network entire activity
87 falls into one of the following categories: Ready-Set, Set-Go, Go, and its final configuration. These
88 state-space regions are not mutually exclusive, thus each different realization of them may indicate
89 a different algorithm. We derived a set of binary features (Figure 5) to partition the networks into
90 solution sets with distinct topological structures. To measure whether the Ready-Set and the Set-Go
91 epochs are the same objects (Fig 5 A) we fitted a linear SVM classifier to separate the neural states
92 that constitute these two epochs. If the classification was not successful, we considered them the
93 same. To measure whether the network transitions to the Go epoch after the Ready (Fig 5 E) we let it
94 evolve from there for 200 steps and measured its output. We initially let the network evolve after the
95 Ready pulse for 5000 steps and saved the final state and the neural velocity at that state. We applied a
96 threshold of 10^{-6} to determine whether it converges to a fixed point (Fig 5 D). To measure whether
97 a transition to Go occurs from any state (5 B) or whether its activity is periodic (Fig 5 C) we let it
98 evolve spontaneously for 1000 steps from that saved state and measured its output and periodicity
99 respectively. Using these features, we divided the space of solutions into six distinct clusters.

100 2.5 Different views of same object

101 To see whether the neural data from the training set contains information about the topology of the
102 networks, we evaluated the ability of the neural features to predict the topological classification
103 we described earlier. This was done by a Cross-Validation procedure that included 50 repetitions
104 of fitting a Decision Tree classifier to a randomly selected 70% of the data, and then evaluating
105 the *kappa-cohen* score and the confusion matrix of the classification on the remaining validation
106 set. For each architecture separately and combined, the mean and the standard errors of these two
107 measurements across all repetitions are shown in 6.

108 References

- 109 [1] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,
110 N. Gimeshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani,
111 S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style,
112 high-performance deep learning library,” in *Advances in Neural Information Processing Systems*
113 32, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds.
114 Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: [http://papers.neurips.cc/
115 paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
- 116 [2] R. T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud, “Neural ordinary differential
117 equations,” *Advances in Neural Information Processing Systems*, 2018.

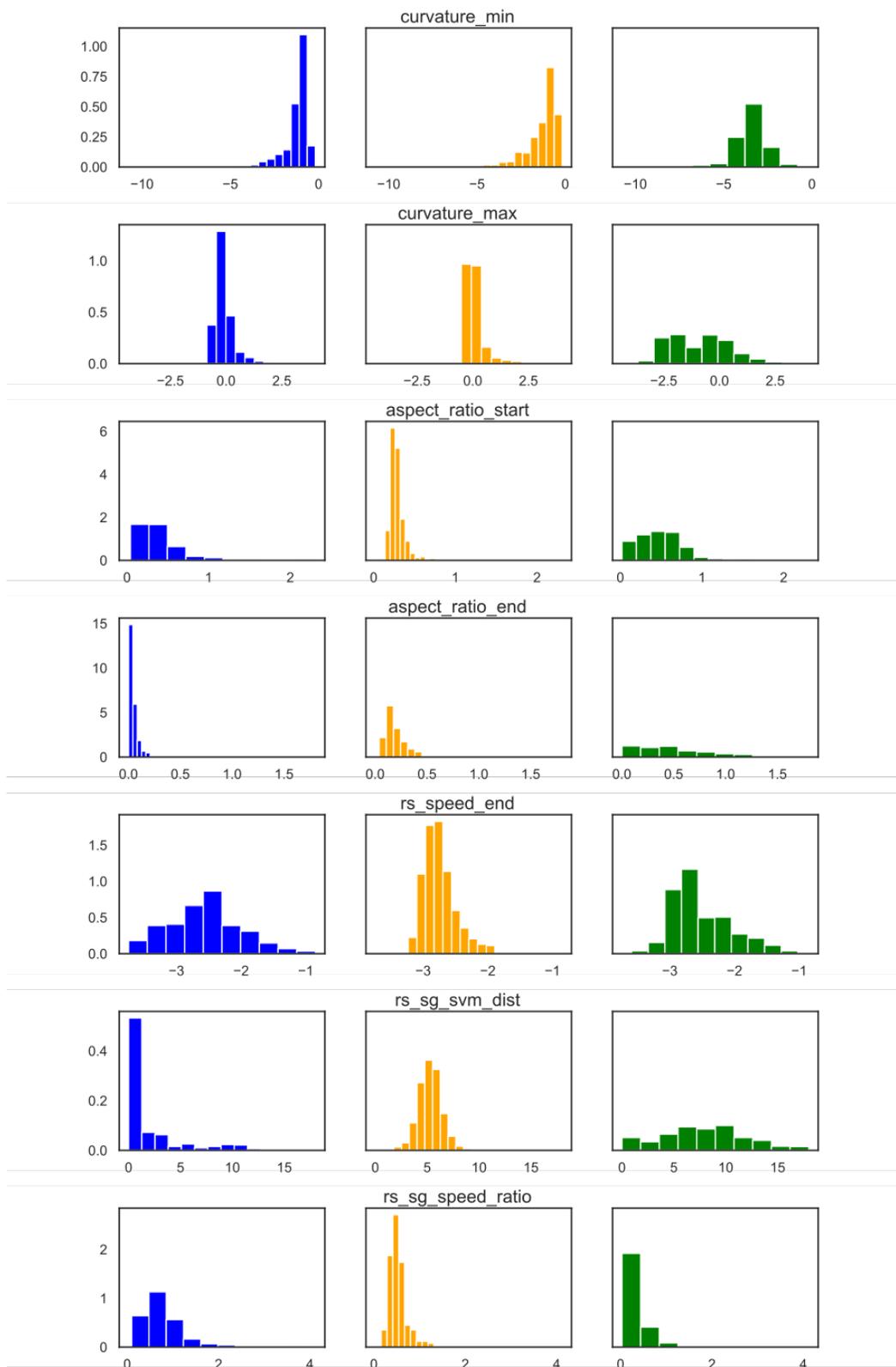


Figure 4: The histogram of every feature, shown for each architecture.

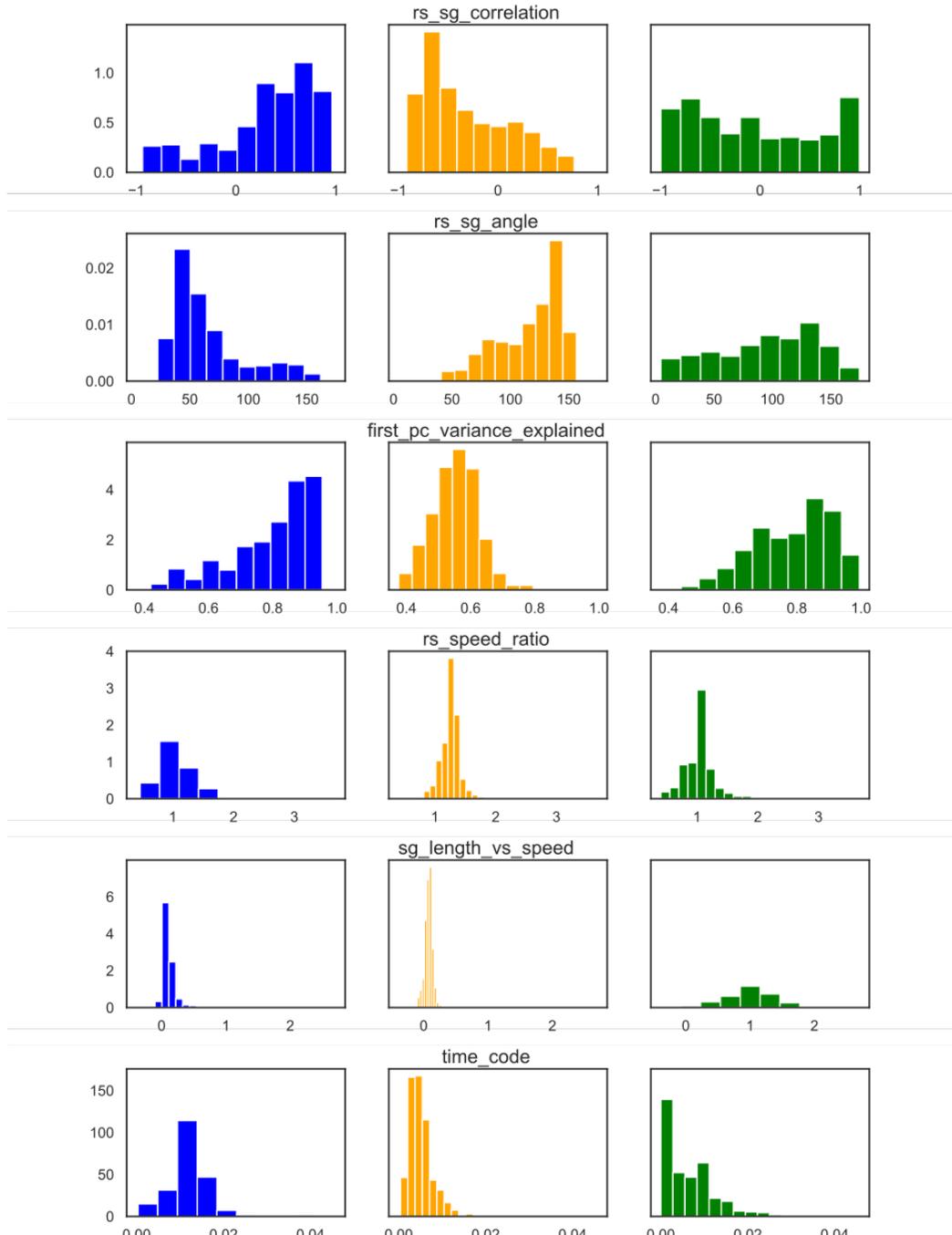


Figure 4: The histogram of every feature, shown for each architecture. (cont.)

- 118 [3] R. T. Q. Chen, B. Amos, and M. Nickel, “Learning neural event functions for ordinary differential equations,” *International Conference on Learning Representations*, 2021.
119
- 120 [4] J. L. Elman, “Finding Structure in Time,” vol. 14, no. 2, pp. 179–211. [Online]. Available:
121 http://doi.wiley.com/10.1207/s15516709cog1402_1
- 122 [5] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and
123 Y. Bengio, “Learning Phrase Representations using RNN Encoder–Decoder for Statistical
124 Machine Translation,” in *Proceedings of the 2014 Conference on Empirical Methods in Natural
125 Language Processing (EMNLP)*. Association for Computational Linguistics, pp. 1724–1734.

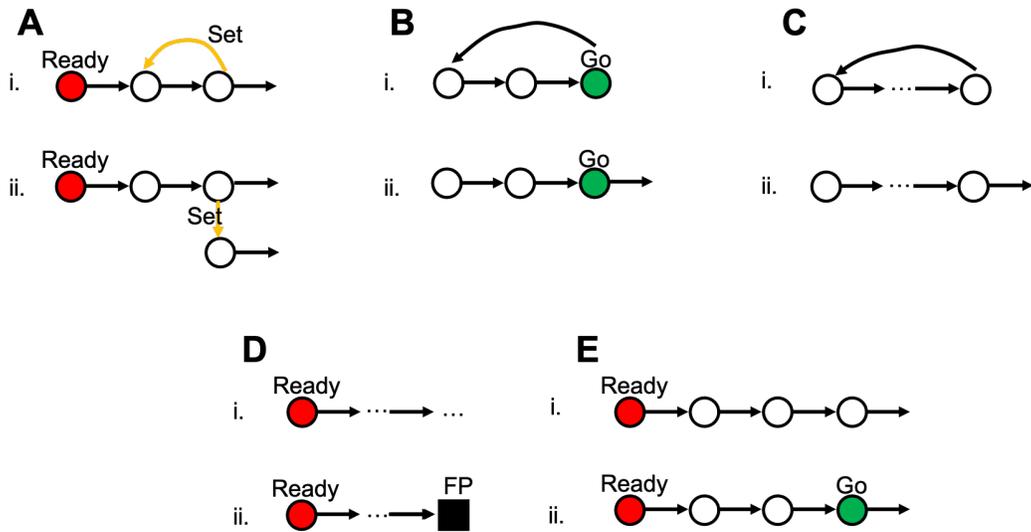


Figure 5: The five topological properties targeted by the binary features. Each feature separates alternatives i and ii. **A** Whether the Ready-Set epoch is identical to the Set-Go (i) or not (ii). **B** Whether the final state transitions to the Go epoch autonomously (i) or not (ii). **C** Whether the network enters a limit-cycle (i) or not (ii). **D** Whether the network reaches a fixed points autonomously from the Ready-Set epoch (i) or not (ii). **E** Whether the Ready-Set epoch transitions to the Go epoch autonomously (i) or not (ii).

126 [Online]. Available: <http://aclweb.org/anthology/D14-1179>

127 [6] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," vol. 9, no. 8, pp. 1735–1780.

128 [Online]. Available: <https://www.mitpressjournals.org/doi/abs/10.1162/neco.1997.9.8.1735>

129 [7] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. [Online]. Available:

130 <http://arxiv.org/abs/1412.6980>

131 [8] D. Sussillo and O. Barak, "Opening the Black Box: Low-Dimensional Dynamics in

132 High-Dimensional Recurrent Neural Networks," vol. 25, no. 3, pp. 626–649. [Online]. Available:

133 https://www.mitpressjournals.org/doi/abs/10.1162/NECO_a_00409

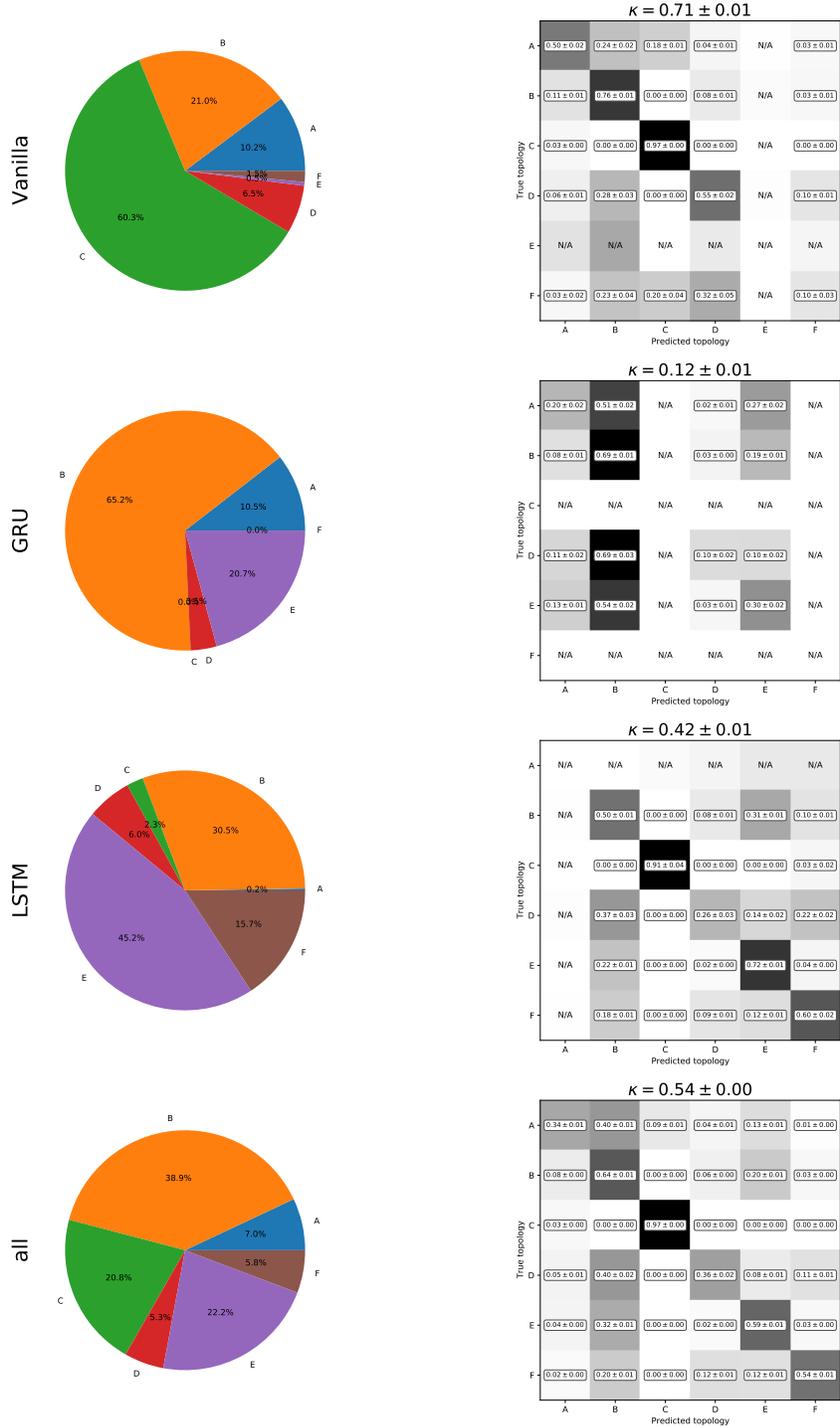


Figure 6: Left: The distribution of the different topologies for each architecture, and for all networks combined. Right: Confusion matrices (mean and standard errors) obtained from 50 repetitions of a Decision-Tree classifier. Cells corresponding to underrepresented topologies are shown as N/A. The value of Cohen's κ is shown for each matrix.