
Learning Useful Representations of Recurrent Neural Network Weight Matrices

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recurrent Neural Networks (RNNs) are general-purpose parallel-sequential com-
2 puters. The program of an RNN is its weight matrix. How to learn useful
3 representations of RNN weights that facilitate RNN analysis as well as down-
4 stream tasks? While the “mechanistic approach” directly looks at some RNN’s
5 weights to predict its behavior, the “functionalist approach” analyzes its overall
6 functionality—specifically, its input-output mapping. Our two novel functionalist
7 approaches extract information from RNN weights by ‘interrogating’ the RNN
8 through probing inputs. Our novel theoretical framework for the functionalist
9 approach demonstrates conditions under which it can generate rich representations
10 that help determine RNN behavior. RNN weight representations generated by
11 mechanistic and functionalist approaches are compared by evaluating them in two
12 downstream tasks. Our results show the superiority of functionalist methods.

13 1 Introduction

14 For decades, researchers have developed techniques for learning internal representations of complex
15 objects in deep neural networks (NNs). This expertise has significantly advanced the field by
16 enabling models to convert data into formats useful for solving problems. In particular, recurrent NNs
17 (RNNs) have been widely adopted due to their computational universality [18]. Low-dimensional
18 representations of the programs of RNNs (their weight matrices) are of great interest as they can
19 speed up the search for solutions to given problems. For instance, compressed representations of
20 RNN weight matrices have been used to evolve RNN parameters [10] for controlling a car from raw
21 video input [9], using Fourier-type transforms, e.g., the coefficient of the Discrete Cosine Transform
22 (DCT) [19], without using the capabilities of NNs to learn such representations. Recent work has
23 seen a rise of representation learning techniques for NN weights using NNs as encoders [20, 17, 1, 4].
24 However, there is a lack of methods for learning representations of RNNs. This paper introduces
25 novel techniques for learning them, using powerful NNs which may be RNNs themselves. Just like
26 representation learning in other fields, such as computer vision, facilitates solutions of specific tasks,
27 such techniques can facilitate learning, searching, and planning with RNNs.

28 2 Self-supervised Learning of Function Representations

29 We consider a Recurrent Neural Network (RNN), $f_\theta : \mathbb{R}^X \times \mathbb{R}^H \rightarrow \mathbb{R}^Y \times \mathbb{R}^H; (x, h_o) \mapsto (y, h_n)$,
30 parametrized by $\theta \in \Theta$, which maps an input x and hidden state h_o to an output y and a new hidden
31 state h_n . The RNN interacts with a potentially stochastic environment, \mathcal{E} , that maps an RNN’s output
32 y to a new input x . The environment may have its own hidden state η . By sequentially interacting

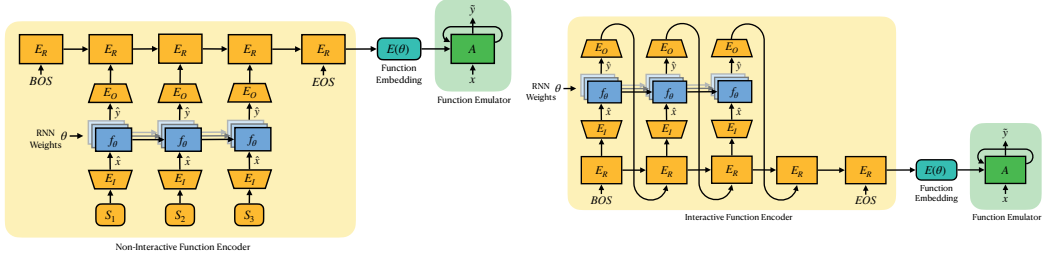


Figure 1: Left: Non-Interactive Encoder. Right: Interactive Encoder.

33 with the environment, the RNN produces a rollout defined by:

$$\begin{cases} x_t, \eta_t = \mathcal{E}(y_{t-1}, \eta_{t-1}) \\ y_t, h_t = f_\theta(x_t, h_{t-1}), \end{cases}$$

34 with fixed initial states y_0, η_0 and h_0 . For instance, f_θ might be an autoregressive generative model,
 35 with \mathcal{E} acting as a stochastic environment that receives a probability distribution over some language
 36 tokens, y_t —the output of f —, and produces a representation (e.g., a one-hot vector) of the new input
 37 token x_{t+1} . When the environment is stochastic, numerous rollouts can be generated for any $\theta \in \Theta$.
 38 A rollout sequence of a function f_θ in environment \mathcal{E} has the form $S_\theta = (x_1, y_1, x_2, y_2, \dots)$.

39 **Encoder and Emulator** Our primary objective is to propose, analyze, and train several methods
 40 for representing RNN weights. We define the Encoder $E_\phi : \Theta \rightarrow \mathbb{R}^M; \theta \mapsto z$, parametrized by
 41 $\phi \in \Phi$ as a function mapping the RNN parameters θ to a lower-dimensional representation z . To train
 42 the encoder E_ϕ , we consider an Emulator $A_\xi : \mathbb{R}^X \times \mathbb{R}^B \times \mathbb{R}^Z \rightarrow \mathbb{R}^Y \times \mathbb{R}^B; (x, b_o, z) \mapsto (\tilde{y}, b_n)$,
 43 parametrized by $\xi \in \Xi$. The Emulator is an RNN with hidden state b that learns to imitate different
 44 RNNs f_θ based on their function encoding $z = E(\theta)$.

45 **Dataset and Training** We consider a dataset $\mathcal{D} = \{(\theta_i, S_{\theta_i}) | i = 1, 2, \dots\}$ composed of tuples,
 46 each containing the parameters of a different RNN and a corresponding rollout sequence. We
 47 assume that all RNNs have the same initial state h_0 but have been trained on different tasks. Our
 48 self-supervised learning approach to training function representations is inspired by the work of [12]).
 49 The Encoder E_ϕ and the Emulator A_ξ are jointly trained by minimizing a loss function \mathcal{L} . This
 50 loss function measures the behavioral similarity between an RNN f_θ and the Emulator A_ξ , which
 51 is conditioned on the function representation $z = E_\phi(\theta)$ of θ as produced by the Encoder E_ϕ . Put
 52 simply, the Emulator utilizes the representations of a set of diverse RNNs f_θ to imitate their behavior:

$$\min_{\phi, \xi} \mathbb{E}_{(\theta, S) \sim \mathcal{D}} \sum_{(x_i, y_i) \in S} \mathcal{L}(A_\xi(x_i, b_{i-1}, E_\phi(\theta)), y_i). \quad (1)$$

53 In the case of continuous outputs y , the mean-squared error provides a suitable loss function. Con-
 54 versely, for categorical outputs, we employ the inverse Kullback-Leibler divergence.

55 2.1 RNN Encoders

56 In this section, we explore various mechanistic and functionalist methods for constructing RNN
 57 encoders. These approaches will be compared in our experimental section. Implementation details
 58 are described in Appendix B.

59 **Flattened Weights (Mechanistic)** Flattening the weights into a single vector presents the most
 60 straightforward method for encoding an RNN. While this technique has shown efficacy on a modest
 61 scale [3, 6], it faces challenges when applied to larger parameter vectors, especially in handling
 62 weight-space symmetries such as neuron permutations.

63 **Neural Functional (Mechanistic)** Fast Weight Programmers [14, 15, 13, 8] are neural networks
 64 that can process the gradients or weights of another neural network. A recent variant thereof, called
 65 Neural Functionals[21], has been used to learn representations of neural network weights that are

66 invariant to the permutation of hidden neurons. The architecture comprises layers that display
 67 equivariance to neuron permutations, followed by a final pooling operation that ensures the invariance
 68 property. Neural Functionals have been theoretically proven to be able to extract all information from
 69 the weights of a neural network [11]. However, their implementation to date has been confined to
 70 feedforward networks, such as MLPs and CNNs.

71 **Non-Interactive RNN Probing (Functionalist)** In the context of Reinforcement Learning and
 72 Markov Decision Processes, policy fingerprinting has emerged as an effective way to evaluate
 73 feedforward neural network policies [4, 5, 2]. In policy fingerprinting, a set of learnable probing
 74 inputs is given to the network. Based on the set of corresponding policy outputs, a function (policy)
 75 representation is produced. This approach can be adapted in a straightforward way for RNNs by
 76 learning whole probing sequences instead of probing inputs (see Figure 1, left). In the context of this
 77 paper, we refer to this approach as non-interactive RNN probing.

78 **Interactive RNN Probing (Functionalist)** The probing sequences for non-interactive RNN probing
 79 are static, i.e., at test time, the probing sequences do not depend on the specific RNN being evaluated.
 80 The alternative is to make the probing sequences dynamically dependent on the given RNN. Each
 81 item in the probing sequences should depend on the outputs of the given RNN to the previous items
 82 (Figure 1, right). This idea has been described previously to extract arbitrary information from a
 83 recurrent world model [16]. Our theoretical framework shows that this novel approach, which we
 84 call interactive RNN probing, is more powerful than non-interactive probing in certain cases.

85 2.2 A Theoretical Framework for the Functionalist Approach

86 Developing an abstract theoretical model based on the functionalist view of RNN weights provides
 87 fundamental insights into the potential and limitations of this approach. The functionalist perspective
 88 emphasizes the overall functionality, disregarding the specific underlying mechanisms of RNNs.
 89 Therefore, our abstraction adopts total Turing machines as a model of computation. Practically,
 90 the function encoder is trained using a dataset of functions. In contrast, the encoder maintains
 91 perpetual direct access to the dataset in our theoretical framework. Note that our framework does not
 92 incorporate the concept of generalization to unseen functions or networks. We detail our theoretical
 93 functionalist framework and explore the interactive and non-interactive approach.

94 Let D represent a set of n total computable functions $\{f_i : \mathbb{N} \rightarrow \mathbb{N} | i = 1, 2, \dots, n\}$. In other words,
 95 D comprises n Turing machines that halt on every input, with no pair being functionally equivalent.
 96 Let I_D denote another Turing machine, which we call the *Interrogator*. I_D has access to the function
 97 set D (e.g., the corresponding Turing numbers might be written somewhere on its tape). Moreover,
 98 I_D is given access to one function $f_C \in D$ as a black box. I_D can interact with f_C by providing an
 99 input $x \in \mathbb{N}$ and subsequently reading the corresponding output $f_C(x)$. The task of I_D is to identify
 100 which member of D corresponds to function f_C , while minimizing interactions with f_C . Specifically,
 101 I_D must return $i \in \{1, \dots, n\}$ such that $f_C = f_i$. This setup is depicted in Figure 3 (Appendix). The
 102 proofs of the following propositions can be found in Appendix A.

103 **Proposition 2.1.** *Any function f_C from a set D can be identified by an interrogator through at most*
 104 *$|D| - 1$ interactions.*

105 An Interrogator is called *interactive* if the value x_j of the j th probing input depends on
 106 $f_C(x_1), \dots, f_C(x_{j-1})$, i.e., the outputs corresponding to the previous probing inputs. This im-
 107 plies that the probing inputs generally depend on the specific function f_C given to I . Conversely, a
 108 *non-interactive* Interrogator can only provide a fixed set of probing inputs to f_C , and their values do
 109 not depend on the outputs of f_C . In the proof of Proposition 2.1, the probing inputs given to f_C do
 110 not dynamically depend on f_C . This means that the theorem holds for non-interactive Interrogators.
 111 A natural question arises: Can interactive Interrogators identify a function using fewer interactions?
 112 Although there are instances where they need exponentially fewer interactions, in the worst-case
 113 scenario, both methods necessitate an equivalent number of interactions:

114 **Proposition 2.2.** *The upper bound for probing interactions required to identify a function from a*
 115 *given function set D is $|D| - 1$ for both interactive and non-interactive Interrogators.*

116 **Proposition 2.3.** *There exist function sets for which an interactive Interrogator requires exponentially*
 117 *fewer probing interactions to identify a member than does a non-interactive one.*

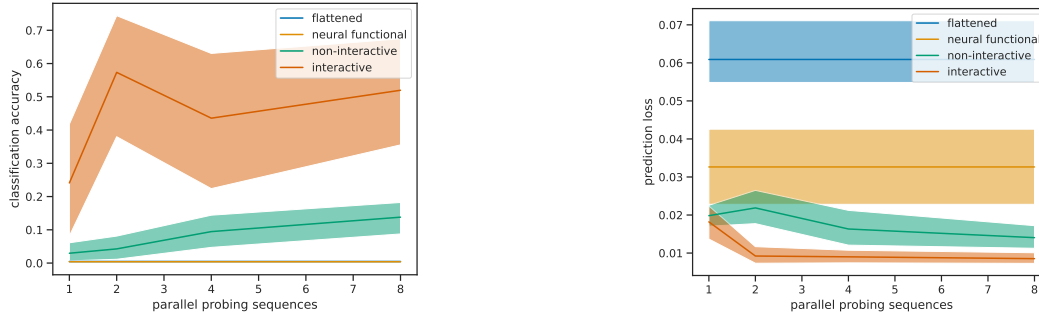


Figure 2: Left: Accuracy of a language classifier, trained using the generated function encodings. Right: Loss of a performance predictor, also trained on the generated function encodings, depicted on the test set. Plots are presented as a function of the number of parallel probing sequences (only relevant for interactive and non-interactive probing encoders). Both graphs display the mean and bootstrapped 95% confidence intervals, aggregated across 15 seeds.

118 3 Experiments

119 We empirically analyze various approaches to learning representations of RNNs, with a specific focus
 120 on LSTM [7] weights. Each LSTM in our dataset serves as an autoregressive generative model of a
 121 specific formal language. The set of formal languages is identical to the one constructed for the proof
 122 of Proposition 2.3. Each LSTM is trained on strings from a particular language using the standard
 123 language modelling objective. We split the dataset into training, validation, and out-of-distribution
 124 (OOD) test parts. The OOD split includes only tasks in which the relative frequencies of each token
 125 appearance are small (i.e., all tokens appear approximately the same number of times). The validation
 126 set is used for early stopping during training. All shown results are derived from the test set. The
 127 experiments employ the four types of function encoders described in Section 2.1. The encoders’
 128 hyperparameters are selected to ensure a comparable number of parameters among them. The training
 129 details remain consistent across all runs. All encoders are trained end-to-end together with an LSTM
 130 emulator to minimize the loss defined in Equation 1, utilizing the reverse Kullback-Leibler divergence
 131 as the loss function \mathcal{L} . Further details, along with additional experimental results, are given in
 132 Appendix C.

133 The objective is to ensure that the LSTM weight encodings z serve as generally useful representations.
 134 We verify this by training models for two downstream tasks using the fixed representations provided
 135 by the encoder E . The first task involves classifying the language on which an LSTM f_θ was trained,
 136 given its encoding $E(\theta)$. This classification is inherently challenging, considering the dataset contains
 137 a total of 216 different languages, and some networks are nearly untrained. The second task aims to
 138 predict the performance of f_θ , defined as the percentage of strings generated by f_θ belonging to the
 139 language that f_θ was trained on. We present the results for these tasks in Figure 2. A visualization
 140 of the learned embedding spaces can be found in Figure 8 in the Appendix. From the results, it is
 141 evident that the interactive probing encoder yields the most useful representations for both tasks.
 142 Having multiple probing sequences in parallel benefits both interactive and non-interactive encoders.
 143 The representations derived from the flattened weights and the neural functional encoder appear to
 144 contain no useful information for the language classifier. In predicting accuracy, representations
 145 from neural functionals outperform those based on flattened weights but fall short when compared to
 146 functionalist representations.

147 4 Conclusion and Future Work

148 We identified two classes of methods for learning RNN weight representations. Firstly, we adapted the
 149 Mechanistic Neural Functional approach to RNNs and, secondly, presented two novel Functionalist
 150 methods, theoretically demonstrating when their representations can be utilized to identify RNNs.
 151 Functionalist methods outperformed Mechanistic ones, learning more useful RNN weight representa-
 152 tions for two downstream tasks. Future work will explore the combination of both approaches and
 153 evaluate their performance on more challenging problems.

References

- 154
- 155 [1] E. Dupont, H. Kim, S. Eslami, D. Rezende, and D. Rosenbaum. From data to functa: Your data
156 point is a function and you can treat it like one. *arXiv preprint arXiv:2201.12204*, 2022.
- 157 [2] F. Faccio, V. Herrmann, A. Ramesh, L. Kirsch, and J. Schmidhuber. Goal-conditioned generators
158 of deep policies. *arXiv preprint arXiv:2207.01570*, 2022.
- 159 [3] F. Faccio, L. Kirsch, and J. Schmidhuber. Parameter-based value functions. *Preprint*
160 *arXiv:2006.09226*, 2020.
- 161 [4] F. Faccio, A. Ramesh, V. Herrmann, J. Harb, and J. Schmidhuber. General policy evaluation and
162 improvement by learning to identify few but crucial states. *arXiv preprint arXiv:2207.01566*,
163 2022.
- 164 [5] J. Harb, T. Schaul, D. Precup, and P.-L. Bacon. Policy evaluation networks. *arXiv preprint*
165 *arXiv:2002.11833*, 2020.
- 166 [6] V. Herrmann, L. Kirsch, and J. Schmidhuber. Learning one abstract bit at a time through
167 self-invented experiments encoded as neural networks. *arXiv preprint arXiv:2212.14374*, 2022.
- 168 [7] S. Hochreiter and J. Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–
169 1780, 1997.
- 170 [8] K. Irie, I. Schlag, R. Csordás, and J. Schmidhuber. Going beyond linear transformers with
171 recurrent fast weight programmers. *Advances in Neural Information Processing Systems*,
172 34:7703–7717, 2021.
- 173 [9] J. Koutník, G. Cuccu, J. Schmidhuber, and F. Gomez. Evolving large-scale neural networks for
174 vision-based reinforcement learning. In *Proceedings of the 15th annual conference on Genetic*
175 *and evolutionary computation*, pages 1061–1068, 2013.
- 176 [10] J. Koutník, F. Gomez, and J. Schmidhuber. Evolving neural networks in compressed weight
177 space. In *Proceedings of the 12th Annual Conference on Genetic and Evolutionary Computation*,
178 pages 619–626, 2010.
- 179 [11] A. Navon, A. Shamsian, I. Achituve, E. Fetaya, G. Chechik, and H. Maron. Equivariant
180 architectures for learning in deep weight spaces. *arXiv preprint arXiv:2301.12780*, 2023.
- 181 [12] R. Raileanu, M. Goldstein, A. Szlam, and R. Fergus. Fast adaptation via policy-dynamics value
182 functions. *arXiv preprint arXiv:2007.02879*, 2020.
- 183 [13] I. Schlag, K. Irie, and J. Schmidhuber. Linear transformers are secretly fast weight programmers.
184 In *International Conference on Machine Learning*, pages 9355–9366. PMLR, 2021.
- 185 [14] J. Schmidhuber. Learning to control fast-weight memories: An alternative to dynamic recurrent
186 networks. *Neural Computation*, 4(1):131–139, 1992.
- 187 [15] J. Schmidhuber. A ‘self-referential’ weight matrix. In *International Conference on Artificial*
188 *Neural Networks*, pages 446–450. Springer, 1993.
- 189 [16] J. Schmidhuber. On learning to think: Algorithmic information theory for novel combi-
190 nations of reinforcement learning controllers and recurrent neural world models. *Preprint*
191 *arXiv:1511.09249*, 2015.
- 192 [17] K. Schürholt, D. Kostadinov, and D. Borth. Hyper-representations: Self-supervised representa-
193 tion learning on neural network weights for model characteristic prediction. 2021.
- 194 [18] H. T. Siegelmann and E. D. Sontag. Turing computability with neural nets. *Applied Mathematics*
195 *Letters*, 4(6):77–80, 1991.
- 196 [19] R. K. Srivastava, J. Schmidhuber, and F. Gomez. Generalized compressed network search. In
197 *Proceedings of the fourteenth international conference on Genetic and evolutionary computation*
198 *conference companion, GECCO Companion ’12*, page 647–648, New York, NY, USA, 2012.
199 ACM, ACM.

- 200 [20] T. Unterthiner, D. Keysers, S. Gelly, O. Bousquet, and I. O. Tolstikhin. Predicting neural
201 network accuracy from weights. *ArXiv*, abs/2002.11448, 2020.
- 202 [21] A. Zhou, K. Yang, K. Burns, Y. Jiang, S. Sokota, J. Z. Kolter, and C. Finn. Permutation
203 equivariant neural functionals. *arXiv preprint arXiv:2302.14040*, 2023.

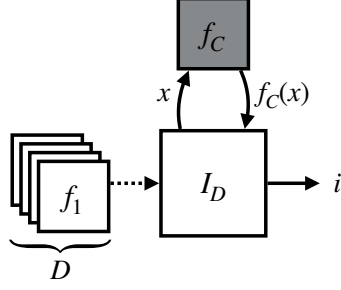


Figure 3: The Interrogator I_D has access to the set of functions D and can interact with one function f_C , which it has to identify.

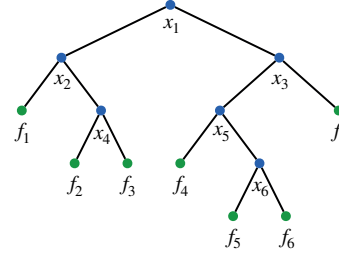


Figure 4: A binary tree constructed as described in the proof of Theorem 2.1. Giving the inputs x_j corresponding to all branching nodes to a function allows to uniquely identify it.

204 A Theoretical Results

205 **Lemma A.1.** *Given any subset $G \subseteq D$, $|G| \geq 2$, there exists an input x that can be computed for*
 206 *which $f_a(x) \neq f_b(x)$ with $f_a, f_b \in G$.*

207 *Proof.* This follows immediately from the fact that all functions in G are total computable and
 208 functionally distinct. \square

209 **Proposition 2.1.** *Any function f_C from a set D can be identified by an interrogator through at most*
 210 *$|D| - 1$ interactions.*

211 *Proof.* According to Lemma A.1, it is possible to split any set $G \subseteq D$, $|G| \geq 2$ into two nonempty,
 212 non-overlapping subsets: $G_a := \{f \in G \mid f(x_j) = f_a(x_j)\}$ and $G_b := \{f \in G \mid f(x_j) \neq f_a(x_j)\}$ for
 213 some $f_a \in G$ and $x_j \in \mathbb{N}$. Any resulting subset that has at least two members can be split again
 214 using the same procedure with a different probing input x_{j+1} . Starting from the full set D , it is
 215 possible to construct a binary tree (see Figure 4) where the leaves are subsets of D containing exactly
 216 one uniquely identified function. The branching (i.e., non-leaf) nodes correspond to the splitting
 217 operation, which involves observing the output of a specific probing input x_j .

218 The Interrogator can identify a given function $f_C \in D$ by providing it with all inputs x_j corresponding
 219 to the branching nodes in the binary tree and observing the outputs. Since any binary tree with n
 220 leaves has exactly $n - 1$ branching nodes, any function $f_C \in D$ can be identified using $|D| - 1$
 221 interactions. \square

222 Of course, there are ‘easy’ function sets in the sense that their members can be identified using much
 223 fewer interactions. Consider, for example, the set $\{n \mapsto i \mid 1 \leq i \leq L\}$. Here, only one (any)
 224 probing input is necessary, since the identity of the function can be directly read from the output.

225 **Proposition 2.2.** *The upper bound for probing interactions required to identify a function from a*
 226 *given function set D is $|D| - 1$ for both interactive and non-interactive Interrogators.*

227 *Proof.* It is easy to construct function sets D for which the members cannot be identified in less than
 228 $|D| - 1$ interactions, even by an interactive Interrogator.

229 One such function set is $\{\xi_i \mid 1 \leq i \leq L\}$ with $\xi_i : n \mapsto \begin{cases} 0 & \text{if } n = i, \\ n & \text{else} \end{cases}$. In the worst case, there is no
 230 way around trying all inputs $i, \dots, L - 1$. \square

231 **Proposition 2.3.** *There exist function sets for which an interactive Interrogator requires exponentially*
 232 *fewer probing interactions to identify a member than does a non-interactive one.*

233 *Proof.* We construct a concrete set of functions that an interactive Interrogator can identify exponen-
 234 tially faster than a non-interactive one. Consider the family of context-sensitive languages

$$L_{m_1, \dots, m_k} := \{a_1^{n+m_1} a_2^{n+m_2} \dots a_k^{n+m_k} \mid n \in \mathbb{N}\}, \quad (2)$$

235 with $m_1, \dots, m_k \in \mathbb{N}$ and a_1, \dots, a_k being the letters/tokens of the language. The parameters m_i
 236 define the relative number of times different tokens may appear. As an example, one member of the
 237 language $L_{3,1,2}$ is the string $a_1 a_1 a_1 a_1 a_2 a_2 a_3 a_3 a_3$.

238 Let $G_L := \{L_{m_1, \dots, m_k} \mid m_1, \dots, m_k \in \{1, \dots, M\}\}$, i.e., a set of such languages with different
 239 parameters m_i . G_L contains M^k languages. To each language L_{m_1, \dots, m_k} , we can assign a unique
 240 generative function g_{m_1, \dots, m_k} . This function, given a partial string from the language, returns a list
 241 of the allowed tokens for the next step. If the input string is not a partial string of the language,
 242 it returns the empty string ϵ . For example, $g_{3,1,2}(a_1 a_1 a_1) = (a_1, a_2)$, $g_{3,1,2}(a_1 a_1 a_1 a_1) = (a_2)$,
 243 and $g_{3,1,2}(a_1 a_1 a_2) = \epsilon$. Our function set D_L is a set of such generative functions, $D_L :=$
 244 $\{g_{m_1, \dots, m_k} \mid m_1, \dots, m_k \in \{1, \dots, M\}\}$.

245 For an interactive Interrogator, there is a simple strategy to identify a given function $g_C \in D_L$ using
 246 $M \cdot k$ interactions: The first input is the string $a_1^M a_2$. From there on, the Interrogator acts as an
 247 autoregressive generative model—it appends the allowed token returned by g_C to the string and uses
 248 it as the new input. Only one valid token will be returned by g_C for all probing input strings that are
 249 generated using this approach since the n is determined from the first input string. This is repeated
 250 until ϵ is returned, which is after a maximum of $(M - 1) \cdot k$ calls to g_C . The last probing input string
 251 will be $a_1^{r_1} \dots a_k^{r_k}$ with $r_1 = M$, from which the language can be inferred in the following way: Let
 252 $n := \min\{r_1, \dots, r_k\}$. The language g_C is generating is thus L_{r_1-n, \dots, r_k-n} .

253 The non-interactive Interrogator cannot use this strategy, since every probing input except the first
 254 depends on g_C 's output for the previous probing input. We can show that in the non-interactive
 255 setting, $(M - 1)^k$ calls to g_C are needed to identify it. Assuming $n = 0$, there are M^{k-1} unique
 256 prefixes for the first token a_k . Each of these prefixes is only allowed in M languages $L_{m_1, \dots, m_{k-1}, \cdot}$,
 257 namely the ones with specific m_1, \dots, m_{k-1} . Remember that g_C returns ϵ whenever it is given a
 258 substrings that is not part of its language. That means, to determine m_k , $M^{k-1}(M - 1)$ different
 259 inputs have to be given to f_C . It follows that in total, $\sum_{b=2}^{k-1} M^b(M - 1) = M^k - M^2$ inputs are
 260 needed to identify the exact language of g_C .

261 In short, to identify a function from the set D_L described above, an interactive Interrogator needs
 262 $O(Mk)$ probing inputs, whereas a non-interactive one needs $O(M^k)$. \square

263 B Implementation Details

264 **Flattened Weights** For the flattened weights Encoder, all parameters θ of the RNN to be encoded
 265 are flattened into a vector. This weight vector is given as input to a multi-layer perceptron (MLP)
 266 with ReLU nonlinearities, which outputs the RNN encoding z .

267 **Neural Functional** For the neural functional (NF) Encoder, we adapt the equivariant NF-layer
 268 [21] for LSTMs. To preserve both equivariance to neuron permutation and functional universality,
 269 the appropriate row- and column-wise feature extractors have to be added for input-to-hidden and
 270 hidden-to-hidden weights, considering rollouts across time and depth of the network.

271 **Non-interactive RNN Probing** A diagram of the non-interactive probing encoder is shown in
 272 Figure 1 (left). RNN probing Encoders have three main components: the core LSTM E_R , an input
 273 projection MLP E_I and an output projection MLP E_O .

274 For the the non-interactive Encoder, a learnable latent probing sequence (S_1, S_2, \dots, S_l) with a fixed
 275 length l is given to E_I . $E_I(S_i)$ is interpreted as either one a several parallel probing inputs \hat{x}_i and
 276 given to f_θ . The resulting probing outputs $\hat{y}_i := f_\theta(\hat{x}_i)$ are given to E_O (in the case of multiple
 277 parallel probing outputs, the values \hat{y}_i are concatenated). The sequence of probing output projections
 278 $(E_O(\hat{y}_1), \dots, E_O(\hat{y}_l))$ is given as input to E_R , preceded by a begin-of-sequence (BOS) and followed
 279 by an end-of-sequence (EOS) token. E_R 's output after the EOS token is transformed with a learned
 280 linear projection into the RNN representation z .

281 **Interactive RNN Probing** As can be seen in Figure 1 (right), the interactive probing encoder
 282 differs from the non-interactive one in one crucial aspect: Instead of having a learned but static latent
 283 probing sequence, the probing inputs at each step are based on the output of E_R from the current step,
 284 which in turn depends on the probing outputs of the previous step. This means that the interactive
 285 probing Encoder can dynamically adapt the probing sequences to the particular RNN f_θ that is being
 286 encoded.

287 **Emulator** The Emulator A_ξ is an LSTM network. The conditioning on the function representation
 288 z is done by adding a learned linear projection of z to the embedding of the begin-of-sequence token.

289 C Experimental Details

290 **LSTM Dataset** The set of languages is $\{L_{r,r+o_1,r+o_2,r+o_3} | o_1, o_2, o_3 \in \{-3, \dots, 2\} \text{ and } r =$
 291 $-\min\{o_1, o_2, o_3\}\}$, with L defined in Equation 2. This set contains $6^3 = 216$ uniquely identifiable
 292 languages. The training data for each LSTM are strings from a particular language of length ≤ 40 ,
 293 with an additional begin-of-sequence and end-of-sequence token.

294 The LSTMs trained for the dataset have two layers with a hidden size of 32, resulting in a total of
 295 13766 parameters. In total, 1000 such networks are trained, each on one of the 216 possible languages.
 296 For each LSTM, 10 snapshots (at steps 0, 100, 200, 500, 1000, 2000, 5000, 10000 and 20000) are
 297 saved during training. A snapshot consists of the LSTM’s current weights and 100 sequences, also of
 298 length 40, generated by it.

299 **Hyperparameters** Table 1 shows the hyperparameters shared by all four encoder types in the
 300 experiments. Hyperparameters specific to probing, flattened and neural functional encoders are shown
 301 in Tables 2, 3 and 4, respectively.

Hyperparameter	Value
A hidden size	256
A #layers	2
z size	16
batch size	64
optimizer	AdamW
learning rate	0.0001
weight decay	0.01
gradient clipping	0.1

Table 1: General hyperparameters

Hyperparameter	Value
E_R hidden size	256
E_R #layers	2
E_I hidden size	128
E_I #layers	1
E_O hidden size	128
E_O #layers	1
probing sequence length	22

Table 2: Hyperparameters for probing (inter-active and non-interactive) encoders

Hyperparameter	Value
hidden size	128
#layers	3

Table 3: Hyperparameters for flattened weights encoders

Hyperparameter	Value
#channels	32
#layers	4

Table 4: Hyperparameters for neural functional encoders

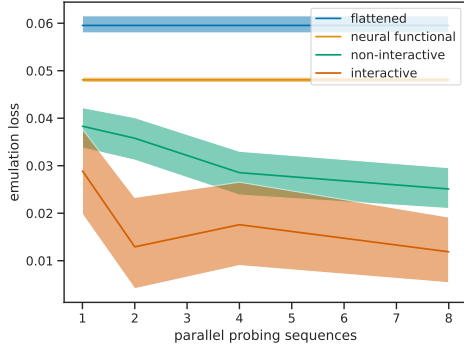


Figure 5: Emulation loss (Equation 1). Plotted as a function of the number of parallel probing sequences. Mean and bootstrapped 95% confidence intervals across 15 seeds.

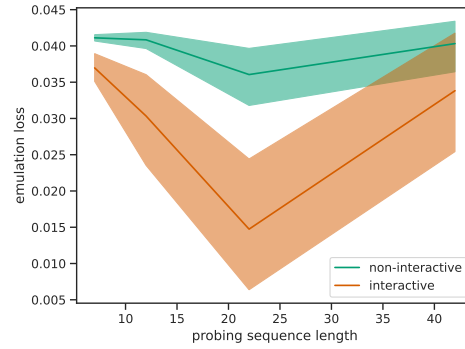


Figure 6: Emulation loss as a function of probing sequence length. Mean and bootstrapped 95% confidence intervals across 15 seeds.

302 **Additional Results** Figure 5 shows the test losses of the different Emulators (as defined in Equa-
 303 tion 1). The relative performance of the different encoders types is similar as for the downstream
 304 tasks shown Figure 2.

305 We also investigate the results for different lengths of the probing sequence for the interactive and
 306 non-interactive probing encoders. In principle, for an interactive probing encoder to correctly identify
 307 all languages from the dataset using the strategy explained in the proof of Proposition 2.3, a probing
 308 sequence of length at least 18 is needed. Figure 7 show the probing sequences of different interactive
 309 probing encoders generated for a specific LSTM that generates strings from the languages $L_{3,0,2,4}$. A
 310 sequence length of 7 is clearly too short, and no insightful probing sequence is learned. For sequence
 311 lengths 12, 22 and 42, the encoder learns to probe actual strings of different lengths from the language.
 312 Note that for this particular language, there exists a string of length < 12 , this is not the case for all
 313 languages used in the dataset. This can also be seen in Figure 6, where only the interactive encoder
 314 with a sequence length of 22 has good performance. A long probing sequence length leads to training
 315 instabilities.

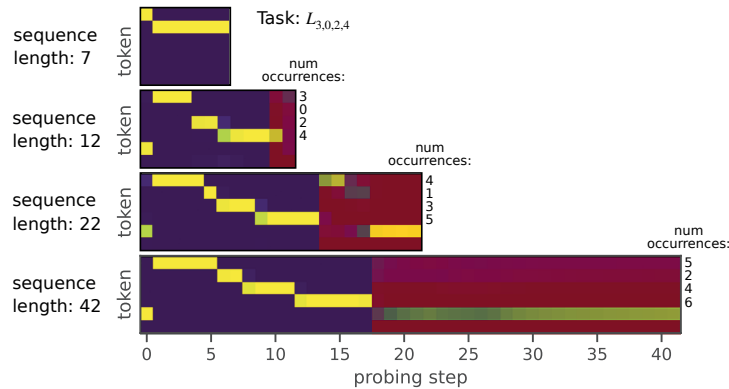


Figure 7: Probing sequences generated for $g_{2,0,1,3}$ by the best performing interactive function encoder with different probing sequence lengths. For the sequence lengths 12, 22 and 42, the encoder produces a insightful probing sequence, i.e. probing sequences that belong to the corresponding language.

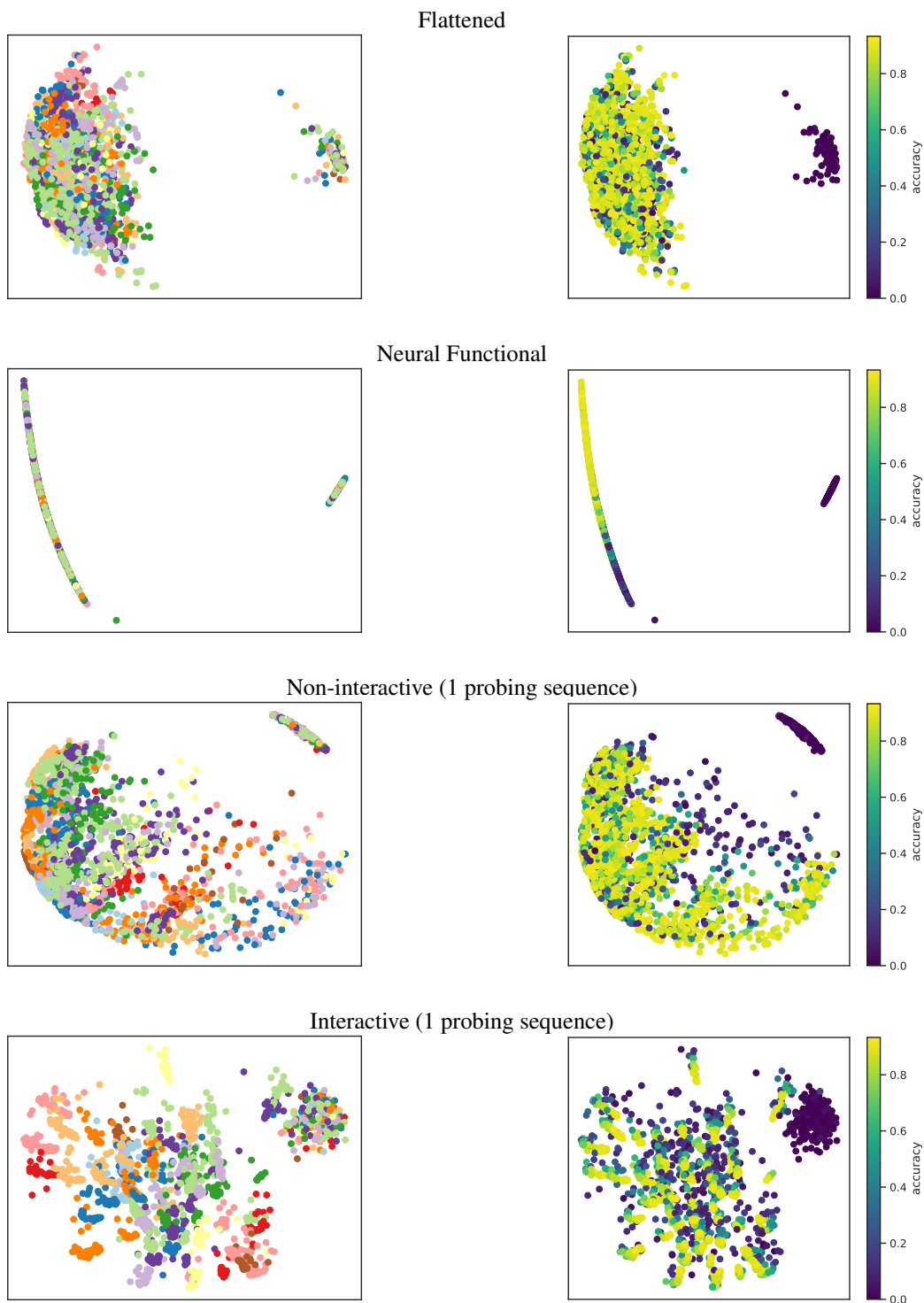


Figure 8: PCA of the function encodings generated by different encoders. Left column: Colored by language, Right column: Colored by performance (i.e., accuracy). The two columns correspond to the two down-stream tasks for which the results are shown in Figure 2.