

# ON PROVABLE LENGTH AND COMPOSITIONAL GENERALIZATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Out-of-distribution generalization capabilities of sequence-to-sequence models can be studied from the lens of two crucial forms of generalization: length generalization – the ability to generalize to longer sequences than ones seen during training, and compositional generalization: the ability to generalize to token combinations not seen during training. In this work, we provide first provable guarantees on length and compositional generalization for common sequence-to-sequence models – deep sets, transformers, state space models, and recurrent neural nets – trained to minimize the prediction error. Taking a first principles perspective, we study the realizable case, i.e., the labeling function is realizable on the architecture. We show that *simple limited capacity* versions of these different architectures achieve both length and compositional generalization. In all our results across different architectures, we find that the learned representations are linearly related to the representations generated by the true labeling function.

## 1 INTRODUCTION

Large language models (LLMs), such as the GPT models (Achiam et al., 2023) and the Llama models (Touvron et al., 2023), have led to a paradigm shift in the development of future artificial intelligence (AI) systems. The accounts of their successes (Bubeck et al., 2023; Gunasekar et al., 2023) as well as their failures, particularly in reasoning and planning (Bubeck et al., 2023; Stechly et al., 2023; Valmeekam et al., 2023), continue to rise. The successes and failures of these models have sparked a debate about whether they actually learn general algorithms or if their success is primarily due to memorization and a superficial form of generalization (Dziri et al., 2024).

A model’s ability to perform well across different distribution shifts highlights its ability to learn general algorithms. For models with fixed-dimensional inputs, considerable efforts have led to methods with provable out-of-distribution (OOD) generalization guarantees (Rojas-Carulla et al., 2018; Rame et al., 2022; Chaudhuri et al., 2023; Wiedemer et al., 2023b; Eastwood et al., 2024). For sequence-to-sequence models, a large body of empirical works have investigated OOD generalization (Anil et al., 2022; Jelassi et al., 2023) but we lack efforts that study provable OOD generalization guarantees for these models. These provable guarantees provide a stepping stone towards explaining the success of the existing paradigm and also shine a light on where the existing paradigm fails.

OOD generalization capabilities of sequence-to-sequence models can be studied from the lens of two forms of generalization: length generalization – the ability to generalize to longer sequences than ones seen during training, and compositional generalization – the ability to generalize to token combinations not seen during training. While transformers (Vaswani et al., 2017) are the go-to sequence-to-sequence models for many applications, recently, alternative architectures based on state-space models, as noted by Gu et al. (2021), Orvieto et al. (2023b), and Gu & Dao (2023), have shown a lot of promise. This motivates us to study a range of natural sequence-to-sequence architectures, including deep sets (Zaheer et al., 2017), transformers, state space models (SSMs), and recurrent neural networks (RNNs). We focus on the realizable case, i.e., the labeling function is in the hypothesis class of the architecture. Further, in our theoretical analysis, we make certain simplifications to permit tractable analysis, for instance, we study RNNs with a limit on hidden state dimension. Our key contributions and insights are summarized below.

- Simple limited capacity versions of the different architectures namely deep sets, transformers, SSMs, and RNNs, provably achieve length and compositional generalization.
- In all our results across different architectures, we find that the learned representations are linearly related to the representations generated by the true labeling function, which is also termed *linear identification* (Khemakhem et al., 2020; Roeder et al., 2021).
- Through a range of experiments, we show the success in both forms of generalization, matching the predictions of the theory and even going beyond.

To the best of our knowledge, our provable guarantees for length and compositional generalization for sequence-to-sequence models are the first in the literature.

## 2 RELATED WORKS

**Length generalization** In the field of length generalization, many important empirical insights have been synthesized over the last few years. Shaw et al. (2018) discovered the drawbacks of absolute positional embeddings and suggested relative positional embeddings as an alternative. Subsequent empirical analyses, notably by Anil et al. (2022) and Jelassi et al. (2023), explored length generalization in different settings for transformer-based models. Key findings revealed that larger model sizes don’t necessarily enhance generalization, the utility of scratchpads varies, and the effectiveness of relative positional embeddings appeared task-dependent. In Kazemnejad et al. (2024), the authors did a comprehensive study of different positional embeddings and provided evidence to show that explicit use of positional encodings is perhaps not essential. In Delétang et al. (2022), the authors conducted experiments on tasks divided based on their placement in the Chomsky hierarchy and showed the importance of structured memory (stack, tape) in length generalization. In a recent work, Zhou et al. (2023) proposed RASP conjecture, which delineates the tasks where transformers excel or fall short in length generalization, emphasizing the necessity of task simplicity for the transformer and data diversity. Our work is inspired by the experimental findings of their work. While Zhou et al. (2023) provide empirical evidence for the conjecture, our work formalizes and proves simpler versions of the conjecture for a range of architectures.

On the theoretical side of length generalization, in Abbe et al. (2023), the authors showed an implicit bias of neural network training towards min-degree interpolators. This bias was used to explain the failures of length generalization on the parity task from Anil et al. (2022). In Xiao & Liu (2023), the authors leverage directed acyclic graphs (DAGs) to formulate the computation in reasoning tasks and characterize conditions under which there exist functions that permit length generalization. Our results crucially differ, we show a range of conditions under which models learned via standard expected risk minimization achieve length and compositional generalization.

**Compositional generalization** Compositionality has long been seen as a key piece to the puzzle of human-level intelligence (Fodor & Pylyshyn, 1988; Hinton, 1990; Plate et al., 1991; Montague, 1970). Compositionality is a large umbrella term associated with several aspects (Hupkes et al., 2020). In this work, we focus on systematicity, which evaluates a model’s capability to understand known parts and combine them in new contexts. The breadth of research on compositional generalization, encompassing studies like Lake & Baroni (2018); Loula et al. (2018); Gordon et al. (2019); Hupkes et al. (2020); Kim & Linzen (2020); Xu et al. (2022); Arora & Goyal (2023); Zhang et al. (2024), is too expansive to address comprehensively here, refer to these surveys (Lin et al., 2023; Sinha et al., 2024) for more detail.

In recent years, several works have taken first steps towards theoretical foundations of compositionality. We leverage the mathematical definition of compositionality from Wiedemer et al. (2023b), which focuses on generalization to the Cartesian product of the support of individual features. In Dong & Ma (2022), the authors analyze the conditions that provably guarantee generalization to the Cartesian product of the support of individual training features. Dong & Ma (2022) studied additive models, i.e., labeling function is additive over individual features. In (Wiedemer et al., 2023a), the authors focus on a more general model class than Dong & Ma (2022), where the labeling function is of the form  $f(x_1, \dots, x_n) = C(\psi_1(x_1), \dots, \psi_n(x_n))$ . However, to guarantee compositional generalization, Wiedemer et al. (2023b) require that the learner needs to know the exact function  $C$

that is used to generate the data. In our work, we do not make such an assumption, our data generation is dictated by the architecture in question, e.g., RNN, and we constrain the dimension of its hidden state. Lachapelle et al. (2023); Brady et al. (2023) extend these precursor results from Dong & Ma (2022) from the supervised setting to the unsupervised setting. In particular, Lachapelle et al. (2023); Brady et al. (2023) are inspired by the success of object-centric models and show additive decoder based autoencoders achieve compositional generalization.

### 3 PROVABLE LENGTH AND COMPOSITIONAL GENERALIZATION

We are given a dataset comprising of a sequence of inputs  $\{x_1, \dots, x_t\}$  and a corresponding sequence of labels  $\{y_1, \dots, y_t\}$ , where each  $x_i \in \mathbb{R}^n$  and  $y_i \in \mathbb{R}^m$ . Observe that this formulation includes both standard downstream tasks such as arithmetic tasks, e.g.,  $y_i = \sum_{j=1}^i x_j$ ,  $y_i = \prod_{j=1}^i x_j$  etc., as well as next-token prediction task, where  $\{y_1, \dots, y_t\} = \{x_2, \dots, x_{t+1}\}$ . We denote a sequence  $\{s_1, \dots, s_t\}$  as  $s_{\leq t}$ ,  $X_k$  is random variable for token at  $k^{\text{th}}$  position and its realization is  $x_k$ . Consider a sequence  $\{x_j\}_{j=1}^{\infty}$ , which is sampled from  $\mathbb{P}_X$ , and a subsequence of this sequence  $x_{\leq t} = \{x_j\}_{j=1}^t$ , whose distribution is denoted as  $\mathbb{P}_{X_{\leq t}}$ . The label  $y_t = f(x_{\leq t})$ , where  $f$  is the labeling function. The tuple of base distribution and the labeling function is denoted as  $\mathcal{P} = \{\mathbb{P}_X, f\}$  and the tuple of base distribution up to length  $t$  is denoted as  $\mathcal{P}(t) = \{\mathbb{P}_{X_{\leq t}}, f\}$ . The support of  $k^{\text{th}}$  token  $X_k$  in the sequence sampled from  $\mathbb{P}_X$  is denoted  $\text{supp}(X_k)$ . Given training sequences of length  $T$  from  $\mathcal{P}(T)$ , we are tasked to learn a model from the dataset that takes a sequence  $x_{\leq t}$  as input and predicts the true label  $y_t$  as well as possible. If the model succeeds to predict well on sequences that are longer than maximum training length  $T$ , then it is said to achieve length generalization (a more formal definition follows later). Further, if the model succeeds to predict well on sequences comprising of combination of tokens that are never seen under training distribution, then it is said to achieve compositional generalization (a more formal definition follows later.). We study both these generalization forms next.

**Learning via expected risk minimization** Consider a map  $h$  that accepts sequences of  $n$ -dimensional inputs to generate a  $m$ -dimensional output. We measure the loss of predictions of  $h$ , i.e.,  $h(x_{\leq t})$ , against true labels as  $\ell(h(x_{\leq t}), y_t)$ , where  $y_t$  is the true label for sequence  $x_{\leq t}$ . In what follows, we use the  $\ell_2$  loss. Given sequences sampled from  $\mathcal{P}(T)$ , the expected risk across all time instances up to maximum length  $T$  is defined as  $R(h; T) := \sum_{t=1}^T \mathbb{E}[\ell(h(x_{\leq t}), y_t)]$ . The learner aims to find an  $h^*$  that solves

$$h^* \in \arg \min_{h \in \mathcal{H}} R(h; T), \quad (1)$$

where  $\mathcal{H}$  is the hypothesis class of models. We seek to understand the properties of solutions to equation 1 through the lens of following questions.

When do common sequence-to-sequence models  $\mathcal{H}$  succeed at length & compositional generalization and when do they fail?

**Definition 1.** Consider the setting where a model is trained on sequences  $(x_{\leq t}, y_{\leq t})$  of length up to  $T$  drawn from  $\mathcal{P}(T)$ . If the model achieves zero error on sequences  $(x_{\leq t}, y_{\leq t})$  of length up to  $\tilde{T}$  drawn from  $\mathcal{P}(\tilde{T}), \forall \tilde{T} \geq 1$ , then it achieves length generalization w.r.t.  $\mathcal{P}$ .

In the above definition of length generalization, we simply ask if the model generalizes to longer sequences. We drop the phrase w.r.t  $\mathcal{P}$  hereafter to avoid repetition. We now define a test distribution that evaluates compositional generalization capabilities. We consider sequences of fixed length  $T$ . Define a uniform distribution  $\mathbb{Q}_{X_{\leq T}}$  such that the support of  $\mathbb{Q}_{X_{\leq T}}$  equals the Cartesian product of the support of each token  $X_k$  from  $\mathbb{P}_X$ , we write this joint support as  $\prod_{j=1}^T \text{supp}(X_j)$ . In this case as well, the labeling function continues to be  $f$ . Hence, we obtain the tuple  $\mathcal{Q}(T) = \{\mathbb{Q}_{X_{\leq T}}, f\}$ .

**Definition 2.** Consider the setting where a model is trained on sequences  $(x_{\leq t}, y_{\leq t})$  of length up to  $T$  drawn from  $\mathcal{P}(T)$ . If the model achieves zero error on sequences  $(x_{\leq t}, y_{\leq t})$  of length up to  $\tilde{T}$  drawn from  $\mathcal{Q}(\tilde{T})$ , then it achieves compositional generalization.

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

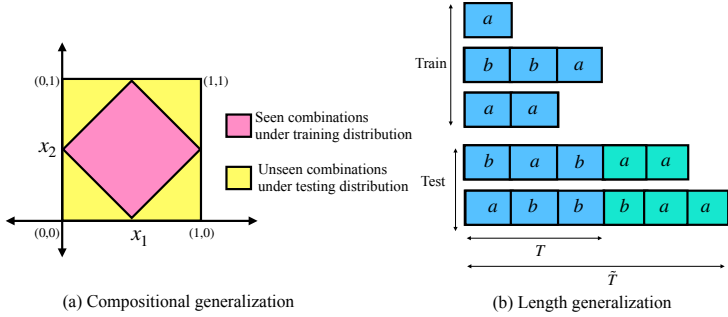


Figure 1: Illustrating support of train vs test distribution for (a) compositional generalization and (b) length generalization.

This definition of compositionality above is based on Wiedemer et al. (2023b); Brady et al. (2023). In this definition, we ask if the model generalizes to new combinations of seen tokens.

*Illustrative example* We teach the model multiplication on sequences of length 2, where each  $x_j$  is a scalar,  $y_i = \prod_{j=1}^i x_j$ . Say the support of the entire sequence drawn from  $\mathbb{P}_X$  is  $\{x \mid \|(x_1, x_2) - \frac{1}{2}\mathbf{1}\|_1 \leq \frac{1}{2}, x_k \in [0, 1], \forall k \geq 3\}$ . The support of training distribution  $\mathbb{P}_{X_{\leq 2}}$  is  $\{x \mid \|(x_1, x_2) - \frac{1}{2}\mathbf{1}\|_1 \leq \frac{1}{2}\}$  shown in the pink region in Figure 1a. In Figure 1a, we illustrate compositional generalization, the model is trained on pink region and asked to generalize to the yellow region. Further, if the model continues to correctly multiply on longer sequence lengths in  $\mathbb{P}_{X_{\leq \hat{T}}}$  for  $\hat{T} \geq T$ , then it achieves length generalization shown in Figure 1b.

**A preview of the technical challenges** Both notions of compositional generalization and length generalization introduced above involve testing on distributions whose support is not contained in the training distributions. The long line of work on distribution shifts (Sugiyama et al., 2007; David et al., 2010; Ben-David & Uner, 2014; Rojas-Carulla et al., 2018; Arjovsky et al., 2019; Ahuja et al., 2021) assume the support of test is contained in the support of train distribution. In recent years, there has been development of theory for distribution shifts under support mismatch Dong & Ma (2022); Abbe et al. (2023); Wiedemer et al. (2023b); Netanyahu et al. (2023); Shen & Meinhshausen (2023). Our work is closer to the latter line of work but it comes with its own *technical challenges*, which involve building new proofs different from the above line of work, as we study a new family of models, i.e., sequence-to-sequence models, and a new form of generalization, i.e., length generalization.

**RASP conjecture** Zhou et al. (2023) propose a conjecture backed by empirical evidence, which delineates the conditions that suffice for length generalization for transformers. The conjecture places three requirements – a) realizability: the task of interest is realizable on the transformer, b) simplicity: the task can be expressed as a short program in RASP-L language, c) diversity: the training data is sufficiently diverse such that there is no shorter program that achieves in-distribution generalization but not OOD generalization. We leverage assumptions similar to a) and b). We assume realizability, which means labeling function  $f$  is in the hypothesis class  $\mathcal{H}$ . As to simplicity, we consider hypothesis class  $\mathcal{H}$  with limited capacity, e.g., we study one block transformer, or RNNs with a limit on hidden state dimension. We emphasize that the third assumption c) on diversity from Zhou et al. (2023) is quite strong. In our setting, we do not invoke it and instead, we require that the support of test distribution is not larger than the Cartesian product of the marginal distribution of the tokens. We now move to proving simplified versions of this conjecture for different architectures.

### 3.1 DEEP SETS

Deep sets are a natural first choice of architecture to study here. These take sets as inputs and thus handle inputs of arbitrary lengths. These were introduced in Zaheer et al. (2017). Informally stated, Zaheer et al. (2017) show that a large family of permutation-invariant functions can be decomposed as  $\rho(\sum_{x \in \mathcal{X}} \phi(x))$ . Consider the examples of the sum operator or the multiplication operator, which take  $\{x_1, x_2, \dots, x_k\}$  as input, and return the sum  $y = \sum_{j=1}^k x_j$  or the product  $y = \prod_{j=1}^k x_j$ . These

operations are permutation invariant and can be expressed using the decomposition above. For the sum operator  $\rho$  and  $\phi$  are identity and for the multiplication operator  $\rho = \exp$  and  $\phi = \log$ . Consider another example from language. We construct a bag of words sentiment classifier, where  $\{x_1, x_2, \dots, x_i\}$  is the set of words that appear in the sentence,  $\phi(x_j)$  is the feature embedding for word  $j$ .  $\sum_{j \leq i} \phi(x_j)$  is the representation of the entire sentence which is passed to the final layer  $\rho$  that generates the sentiment label. In what follows, we aim to understand when such a classifier generalizes to sentences beyond training lengths and to new sentences comprised of unseen word combinations.

**Assumption 1.** *Each function in the hypothesis class  $\mathcal{H}$  takes a sequence  $\{x_1, \dots, x_i\}$  as input and outputs  $h(x_1, \dots, x_i) = \omega\left(\sum_{j \leq i} \psi(x_j)\right)$ , where  $\omega$  is a single layer perceptron with a continuously differentiable bijective activation (e.g., sigmoid) and  $\psi$  is a map that is differentiable.*

A simple mathematical example of a function from the above family when  $\psi(x_j) = [x_j, x_j^2, x_j^3]$  is a polynomial map of degree 3 and each  $x_j$  is a scalar  $-\sigma(a \sum_{j \leq i} x_j + b \sum_{j \leq i} x_j^2 + c \sum_{j \leq i} x_j^3)$ . In the assumption that follows, we assume that the support of the sequences is regular closed in the standard topology, i.e., the set is equal to the closure of its interior.

**Assumption 2.** *The joint support  $\text{supp}(X_{\leq i})$  is a regular closed set for all  $i \leq T$ .*

In most of our results in the main body, we invoke Assumption 2. This assumption is satisfied in many cases if the tokens are continuous random variables but it is not satisfied for discrete random variables. In the Appendix, we extend several of our key results to discrete tokens.

**Linear identification** Each architecture that we study in this work relies on a hidden representation that is passed on to a non-linearity to generate the label. Under the realizability condition for deep sets, the labeling function takes the form  $f(\mathcal{X}) = \rho(\sum_{x \in \mathcal{X}} \phi(x))$ , where  $\phi(x)$  is the hidden representation. If the learned deep set is denoted by  $\omega(\sum_{x \in \mathcal{X}} \psi(x))$ , then the learned hidden representation is  $\psi(x)$ . If  $\psi(x) = A\phi(x)$ , then the learned representation is said to *linearly identify* the true data generating representation  $\phi(x)$ . We borrow this definition from the representation identification literature (Khemakhem et al., 2020; Roeder et al., 2021).

**Theorem 1.** *If  $\mathcal{H}$  follows Assumption 1, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ ,  $\text{supp}(X_j) = [0, 1]^n$ ,  $\forall j \geq 1$ , and the regular closedness condition in Assumption 2 holds, then the model trained to minimize the risk in equation 1 with  $\ell_2$  loss generalizes to all sequences in the hypercube  $[0, 1]^{nt}$ ,  $\forall t \geq 1$  and thus achieves length and compositional generalization.*

The detailed proof is in Section C.1. In the above result, we require the support of the marginal distribution of each token to be  $[0, 1]^n$ . The support of  $T$  token length sequence under the joint training distribution can still be a much smaller subset of  $[0, 1]^{nT}$ , as illustrated in Figure 1a (and Figure 4 in the Appendix). Despite this the model generalizes to all sequences in  $[0, 1]^{nt}$  for all  $t$ . An important insight from the proof is that the hidden representation learned by the model is a linear transform of the true hidden representation, i.e., it achieves linear identification  $\psi = A\phi$  (Further details are in Corollary 1).

**Extensions of Theorem 1** In Theorem 8, we extend Theorem 1 to  $\omega$  from  $C^1$ -diffeomorphisms<sup>1</sup>. As a by product, we obtain length & compositional generalization for multiplication operator. In Theorem 7, we extend the above result to discrete tokens. Further, most results in this work translate to settings where we do not observe labels at all lengths from 1 to  $T$  (further discussion in Appendix).

**High capacity deep sets** In the above results, we operated with some constraints on the deep sets. In Theorem 1, we used limited capacity  $\omega$  that are represented via a single layer perceptron. In Theorem 8, we used  $\omega$  that are represented via  $C^1$ -diffeomorphisms, which implies the output dimension of  $\psi$  equals label dimension  $m$  and cannot be larger. What happens when we work with deep sets with arbitrary capacity, i.e., no constraints on  $\omega$  and  $\psi$ ? These models then express a large family of permutation invariant maps (Zaheer et al., 2017). Suppose  $\mathcal{H}$  is the class of all permutation invariant maps and the labeling function  $f \in \mathcal{H}$ . Consider a map  $h$  such that  $h = f$  for all sequences of length up to  $T$ , and  $h = f + c$  otherwise. Observe that  $h$  is permutation invariant and also belongs

<sup>1</sup> $C^1$ -diffeomorphism - a continuously differentiable map that has a continuously differentiable inverse.

to  $\mathcal{H}$ .  $h$  achieves zero generalization error on training sequences of length  $T$  but a non-zero error on longer sequences. Thus in the setting of high capacity deep sets, there exist solutions to equation 1, which do not achieve length generalization. We can construct the same argument for compositional generalization as well and say  $h = f$  on the training distribution (pink region) in Figure 1a and  $h = f + c$  on the testing distribution (yellow region) in Figure 1a. In order to show successful generalization (length or compositional) in Theorem 1, we require all solutions to risk minimization in equation 1 to match the predictions of true labeling function on data beyond the support of the training distribution. In order to show that high capacity models are not guaranteed to succeed, we focused on showing that there exists a solution to equation 1 that does not generalize beyond the support of training distribution. A more nuanced argument for failure should show that there exist solutions reachable via gradient descent that do not generalize. We leave a rigorous theoretical exploration of this to future work. However, we conduct experiments with high capacity models in the Appendix (Section D.3) to illustrate failures in high capacity regime.

### 3.2 TRANSFORMERS

Ever since their introduction in Vaswani et al. (2017), transformers have revolutionized all domains of AI. In this section, we seek to understand length generalization for these models. Transformer architectures are represented as alternating layers of attention and position-wise non-linearity. We drop layer norms for tractability. Following similar notation as previous section, we denote position-wise non-linearity as  $\rho$  and attention layer as  $\phi$ . We obtain the simplest form of causal transformer model as  $\rho\left(\sum_{j=1}^i \frac{1}{i} \cdot \phi(x_i, x_j)\right)$ . This decomposition captures linear attention, ReLU attention, sigmoid attention, ReLU squared attention, which were studied previously in Wortsman et al. (2023); Hua et al. (2022); Shen et al. (2023) and found to be quite effective in several settings. This decomposition does not capture softmax-based attention and developing provable length generalization guarantees for the same is an exciting future work. Other works (Bai et al., 2023) also replaced softmax with other non-linear attention for a more tractable analysis. We illustrate the sigmoid-based transformer from Wortsman et al. (2023) below. Let  $W_q \in \mathbb{R}^{k \times n}$ ,  $W_k \in \mathbb{R}^{k \times n}$ , and  $W_v \in \mathbb{R}^{k \times n}$  be the query, key and value matrices.  $\rho$  is parametrized via a multi-layer perceptron denoted as MLP.

$$q_i = W_q x_i, k_j = W_k x_j, v_j = W_v x_j, \phi(x_i, x_j) = \sigma\left(\frac{q_i^\top k_j}{\sqrt{d}}\right) v_j, \text{MLP}\left(\sum_{j=1}^i \frac{1}{i} \cdot \phi(x_i, x_j)\right). \quad (2)$$

In the above feedforward computation, the output of attention for the current query is computed and sent to the MLP to generate the label.

**Assumption 3.** Each function in the hypothesis class  $\mathcal{H}$  takes a sequence  $\{x_1, \dots, x_i\}$  as input and outputs  $h(x_1, \dots, x_i) = \omega\left(\sum_{j \leq i} \frac{1}{i} \cdot \psi(x_i, x_j)\right)$ , where  $\omega$  is a single layer perceptron with continuously differentiable bijective activation (e.g., sigmoid) and  $\psi$  is a map that is differentiable.

We denote the joint support of two tokens  $X_i, X_j$  as  $\text{supp}(X_i, X_j)$ .

**Theorem 2.** If  $\mathcal{H}$  follows Assumption 3, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ ,  $\text{supp}(X_i, X_j) = [0, 1]^{2n}$ ,  $\forall i \neq j$  and the regular closedness condition in Assumption 2 holds, then the model trained to minimize the risk in equation 1 (with  $T \geq 2$ ) with  $\ell_2$  loss generalizes to all sequences in the hypercube  $[0, 1]^{nt}$ ,  $\forall t \geq 1$  and thus achieves length and compositional generalization.

Similar to Theorem 1, we observe linear identification here too, i.e., learned attention representation denoted  $\psi$  is a linear transform of the true attention representation denoted  $\phi$ , i.e.,  $\psi(x_i, x_j) = C\phi(x_i, x_j)$ , (details in Section C.2, see Corollary 2). We now extend Theorem 2 from single layer perceptron  $\omega$  to  $C^1$ -diffeomorphism. We also extend Theorem 2 to discrete tokens in Theorem 10.

**Assumption 4.** Each function in  $\mathcal{H}$  takes  $\{x_1, \dots, x_i\}$  as input and outputs  $h(x_1, \dots, x_i) = \omega\left(\sum_{j=1}^{i-1} \frac{1}{i-1} \cdot \psi(x_i, x_j)\right)$ , where  $\omega$  is a  $C^1$ -diffeomorphism,  $\omega(0) = 0$ .

The reader would notice that the summation is up to  $i - 1$  and hence it computes attention scores w.r.t all other terms in the context except  $x_i$ . We conjecture that the theorem that we present next extends to the more general case where summation includes the  $i^{\text{th}}$  term.

**Assumption 5.** *The joint support  $\text{supp}(X_{\leq i})$  is a regular closed set for all  $i \leq T$ . The support of all pairs of tokens is equal, i.e.,  $\text{supp}(X_i, X_j) = [0, 1]^{2n}$ , where  $i \neq j$ ,  $i \geq 1, j \geq 1$ . The support of  $[\phi(X_1, X_2), \phi(X_1, X_3)]$  is  $\mathbb{R}^{2m}$ , where  $\phi$  is the embedding function for the labeling function  $\rho(\sum_{j \leq i} \phi(x_i, x_j))$ .*

**Theorem 3.** *If  $\mathcal{H}$  follows Assumption 4, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , and a further assumption on the support (Assumption 5) holds, then the model trained to minimize the risk in equation 1 (with  $T \geq 3$ ) with  $\ell_2$  loss generalizes to all sequences in  $[0, 1]^{nt}$ ,  $\forall t \geq 1$  and thus achieves length and compositional generalization.*

**Multiple attention heads and positional encoding** While the discussion in this section used a single attention head  $\phi$ , the results extend to multiple attention heads as shown in Section C.2. The model of transformers discussed so far uses the current query and compares it to keys from the past, it does not distinguish the keys based on their positions. For many arithmetic tasks such as computing the median, maximum etc., the positions of keys do not matter but for other downstream tasks such as sentiment classification, the position of the words can be important. In Section C.2, we adapt the architecture to incorporate relative positional encodings and show how some of the results extend. We modify the model as  $\rho(\sum_{j=1}^i \frac{1}{i} \cdot \phi_{i-j}(x_i, x_j))$ , where  $\phi_{i-j}(x_i, x_j)$  computes the query key inner product while taking the relative position  $i - j$  into account. We show that if  $\phi_{i-j} = 0$  for  $i - j > T_{\max}$ , i.e., two tokens sufficiently far apart do not impact the data generation, then length generalization and compositional generalization are achieved.

**High capacity transformers** In the above results, we operated with constraints on transformers, which limit their capacity. Similar to the setting of deep sets, observe that Assumption 3 constrains  $\omega$  to single layer perceptron, Assumption 4 constrains  $\omega$  to  $C^1$ -diffeomorphisms. What happens if we work with transformers with no constraint on  $\omega$  and  $\psi$ ? If  $\psi(x, y) = \psi(\tilde{x}, y), \forall x \neq \tilde{x}$ , then the decomposition for the causal transformer  $\omega(\sum_{j=1}^i \frac{1}{i} \cdot \psi(x_i, x_j))$  becomes  $\omega(\sum_{j=1}^i \frac{1}{i} \cdot \psi(x_j))$ , which is very similar to deep sets. In such a case, we can adapt arguments similar to that of arbitrary capacity deep sets and argue that there exist solutions to equation 1 that do not achieve length and compositional generalization. We now move to state-space models and RNNs.

### 3.3 STATE SPACE MODELS

In recent years, state space models Gu et al. (2021); Orvieto et al. (2023b) have emerged as a promising competitor to transformers. In (Orvieto et al., 2023a;b), the authors used the lens of linear recurrent layer followed by position-wise non-linearities as the main building block to understand these models. We illustrate the dynamics of these models to show the generation of  $x_{\leq t}$  and  $y_{\leq t}$  next. Given the current input  $x_t$ , we combine it linearly with the hidden state from the past to obtain the current hidden state. The hidden state is input to  $\rho$ , which generates the label as follows

$$\begin{aligned} h_1 &= Bx_1; & h_2 &= \Lambda h_1 + Bx_2; & \dots, & h_t &= \Lambda h_{t-1} + Bx_t, \\ y_1 &= \rho(h_1); & y_2 &= \rho(h_2); & \dots, & y_t &= \rho(h_t), \end{aligned} \quad (3)$$

where  $h_t \in \mathbb{R}^k$  is hidden state at point  $t$ ,  $\Lambda \in \mathbb{R}^{k \times k}$ ,  $B \in \mathbb{R}^{k \times n}$  and  $\rho : \mathbb{R}^k \rightarrow \mathbb{R}^m$ . We can succinctly write  $h_t = \sum_{j=0}^{t-1} \Lambda^j Bx_{t-j}$ .

**Assumption 6.** *Each function in the hypothesis class  $\mathcal{H}$  takes a sequence  $\{x_1, \dots, x_i\}$  as input and outputs  $h(x_1, \dots, x_i) = \omega(\sum_{j=0}^{i-1} \Lambda^j Bx_{i-j})$ , where  $\omega : \mathbb{R}^k \rightarrow \mathbb{R}^m$  is a  $C^1$ -diffeomorphism,  $B$  and  $\Lambda$  are square invertible. As a result,  $m = k = n$ .*

**Assumption 7.** *The joint support  $\text{supp}(X_{\leq i})$  is a regular closed set for all  $i \leq T$ . The support of  $X_1$  is  $\mathbb{R}^n$ . For some length  $2 \leq i \leq T$  there exists in sequences  $x_{\leq i}$  such that their concatenation forms a  $i \times i$  matrix of rank  $i$ .*

**Theorem 4.** *If  $\mathcal{H}$  follows Assumption 6, and the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , and a further condition on the support, i.e., Assumption 7, holds, then the model trained to minimize the risk in equation 1 with  $\ell_2$  loss ( $T \geq 2$ ) achieves length and compositional generalization.*

The proof is provided in Section C.3. Similar to previous theorems, the hidden state estimated by the learned model,  $\tilde{h}_t$ , and the true hidden state,  $h_t$ , bear a linear relationship (Corollary 4), i.e., linear identification is achieved. We extend Theorem 4 to discrete tokens in Theorem 12.

**High capacity SSMs** In the above result, we operated with certain constraints on SSMs, i.e., the input dimension, output dimension, and the hidden state dimension are equal. These constraints limit their capacity. What happens if we put no constraints on  $\Lambda$ ,  $B$  and  $\omega$ ? Orvieto et al. (2023a) showed that SSMs with appropriately large  $\Lambda$  and  $B$  matrices can approximate a sequence-to-sequence mapping up to some length with arbitrary precision. Consider the true labeling function  $f$  and another function  $h$ , which is equal to  $f$  for all sequences of length up to  $T$  and  $f + c$  for larger lengths. If we use such arbitrary capacity SSMs as our hypothesis class, then this hypothesis class contains both  $f$  and  $h$ . As a result,  $\tilde{h}$  is a solution to equation 1 and it does not achieve length generalization. We can extend the same argument to compositional generalization as well.

### 3.4 VANILLA RECURRENT NEURAL NETWORKS

Standard RNNs have a non-linear recurrence unlike the linear recurrence studied in the previous section. We use the same notation as the previous section and only add an activation for non-linear recurrence. We illustrate the dynamics to show the generation of  $x_{\leq t}$  and  $y_{\leq t}$  below.

$$\begin{aligned} h_1 &= \sigma(Bx_1); & h_2 &= \sigma(\Lambda h_1 + Bx_2); & \dots, & h_T &= \sigma(\Lambda h_{T-1} + Bx_T) \\ y_1 &= \rho(h_1); & y_2 &= \rho(h_2); & \dots, & y_T &= \rho(h_T), \end{aligned} \tag{4}$$

**Assumption 8.** Each function in the hypothesis class  $\mathcal{H}$  is a vanilla RNN of the form equation 4, where the position-wise non-linearity is a single layer perceptron  $\sigma \circ A$ , and  $\Lambda, B$  govern the hidden state dynamics (equation 4).  $A, \Lambda, B$  are square invertible matrices, and  $\sigma$  is the sigmoid activation.

**Theorem 5.** If  $\mathcal{H}$  follows Assumption 8, and the realizability condition holds, i.e.,  $f \in \mathcal{H}$  and regular closedness condition in Assumption 2 holds, then the model trained to minimize the risk in equation 1 with  $\ell_2$  loss (with  $T \geq 2$ ) achieves length and compositional generalization.

The hidden state estimated by the learned model, i.e.,  $\tilde{h}_t$ , and the true hidden state  $h_t$ , bear a linear relationship (See Corollary 5 in Section C.4 for details), where the linear relationship is a permutation map. We extend Theorem 5 to discrete tokens in Theorem 13.

**High capacity RNNs** In our result above, similar to previous sections we showed that limited capacity RNNs can achieve length and compositional generalization. How about RNNs with arbitrary capacity, i.e., no constraint on  $\Lambda, B$  and  $\rho$ ? These systems can approximate sequence-to-sequence models to arbitrary precision (Sontag, 1992; Gühring et al., 2020). Hence, we can use the same argument as previous sections to argue that if  $\mathcal{H}$  corresponds to RNNs with arbitrary capacity, then there exist solutions to equation 1 that do not achieve length and compositional generalization.

**Remark on proofs** Finally, we would like the reader to appreciate that our proofs follow different strategies in comparison to Wiedemer et al. (2023b); Dong & Ma (2022), due to the fact that we cater to sequence-to-sequence models. Consider the proofs in Wiedemer et al. (2023a), which reduce the solutions of equation 1 to solutions of set of ordinary differential equations, which under their assumptions are unique. That leads to exact identification in contrast to linear identification.

### 3.5 FINITE HYPOTHESIS CLASS

In the discussion so far, we have focused on different hypothesis class  $\mathcal{H}$  of infinite size. In this section, we focus on finite hypothesis class, i.e., the set  $\mathcal{H}$  has a finite size. We can construct such a finite hypothesis class for any architecture by restricting the parameter vectors (weights, biases etc.) to assume a finite set of values. Each possible parameter configuration denotes one distinct element in  $\mathcal{H}$ . Unlike the previous sections, we do not impose any further restrictions on  $\mathcal{H}$  other than the finite size. This allows us to consider arbitrary sequence to sequence models – RNNs, deep sets, transformers (e.g., with hard-coded positional encodings as in (Vaswani et al., 2017)) without restrictions on the depth and width as seen in the previous sections.



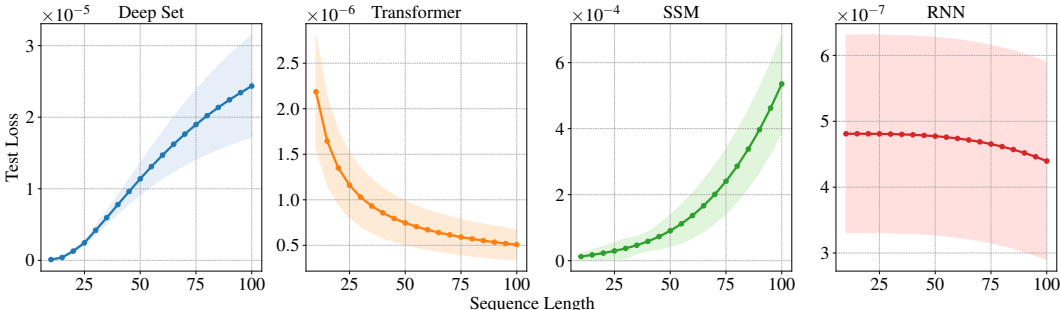


Figure 2: Length generalization: Test  $\ell_2$  loss on sequences of different lengths. The models are trained only on sequences of length up to  $T = 10$ . All models achieve small error values  $\approx 10^{-4} - 10^{-7}$  at all sequence lengths and thus length generalize. Since the error values are already quite small, the increasing or decreasing trends are not numerically significant.

**Theorem 6.** *If  $\mathcal{H}$  is a finite hypothesis class, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , then  $\exists T_0 < \infty$  such that the model trained to minimize the risk in equation 1 with  $\ell_2$  loss and  $T > T_0$  achieves length generalization.*

The above theorem states that for a finite hypothesis class, length generalization is provably achieved provided the training length is sufficiently large. Observe that the above theorem only focuses on length generalization and does not apply to compositional generalization. In the above result, the value of the threshold on  $T$ , i.e.,  $T_0$ , can be very large, and future work should consider quantifying bounds on  $T_0$ . In previous sections, where we had more structural restrictions on  $\mathcal{H}$ , the threshold on  $T$  was two.

## 4 EXPERIMENTS

We present the empirical evaluation of compositional and length generalization capabilities of the architectures from the previous section. All the experiments are carried out in the realizable case where  $f \in \mathcal{H}$ , i.e., depending on the architecture in question, we use a random instance of the architecture to generate the labels. We train a model  $h$  from the same architecture class to minimize the  $\ell_2$  loss between  $h$  and  $f$ . Under different scenarios, we ask if  $h$  achieves length generalization and compositional generalization. We also seek to understand the relationship between the hidden representations of  $h$  and hidden representations of  $f$ .

### 4.1 LENGTH GENERALIZATION

We sample sequences  $x_{\leq t}$  of varying length with a maximum length of  $T = 10$ . Each token  $x_i \sim \text{Uniform}[0, 1]^n$ , where  $n = 20$ . The sequences are then fed to the labeling  $f$ , which comes from the hypothesis class of the architecture, to generate the labels. We minimize the empirical risk version of equation 1 over the same hypothesis class with  $\ell_2$  loss. For evaluation, we present the  $\ell_2$  loss on the test datasets. We also evaluate  $R^2$  of linear regression between the learned hidden representations denoted  $\psi(x_i)$  and true hidden representations  $\phi(x_i)$  for all  $x_i \in x_{\leq t}$  from the test dataset sequences. This metric is often used to evaluate the claims of linear identification (Khemakhem et al., 2020), i.e., the higher this value, the closer the linear relationship. We present results averaged over five seeds for models with two hidden layer MLPs for  $\rho$  ( $\phi$  is two hidden layer MLP for deep sets). Figure 2 shows a very small test loss of models on increasing sequence lengths when only trained with sequences of up to length  $T = 10$ , which is in agreement with Theorem 1-5. Further, in Figure 3, we show an exemplar sequence from test set and how the trained transformer tracks it. Table 1 shows the average of  $R^2$  score of  $\psi(x_i), \phi(x_i)$  across different positions  $i$  at test time. These results demonstrate a linear relationship between learned and true hidden representations, which agrees with our theoretical claims. In Section D, we show that when realizability condition does not hold, i.e.,  $f \notin \mathcal{H}$ , then length generalization is not achieved. We also present *additional experiments with discrete tokens, failures in the high capacity settings*, and other experimental details in Section D.

Model	$R^2 (t = 20)$	$R^2 (t = 100)$	Model	Test Loss $\times 10^6$	$R^2$
Deep set	$0.97 \pm 0.01$	$0.97 \pm 0.01$	Deep set	$0.08 \pm 0.02$	$0.96 \pm 0.01$
Transformer	$0.99 \pm 0.01$	$0.99 \pm 0.01$	Transformer	$3.06 \pm 1.11$	$1.00 \pm 0.00$
SSM	$0.99 \pm 0.01$	$0.99 \pm 0.01$	SSM	$5.92 \pm 2.47$	$1.00 \pm 0.00$
RNN	$0.99 \pm 0.01$	$0.99 \pm 0.01$	RNN	$0.35 \pm 0.17$	$0.96 \pm 0.01$

Table 1: Average test  $R^2$  of true and learned hidden representations  $\psi(x_i), \phi(x_i)$  across all positions  $i$  at various lengths unseen during training. A strong linear relationship is observed for all models across lengths.

Table 2: Compositional generalization: Test  $\ell_2$  loss and  $R^2$  score for models with *two* hidden layers on sequences of length  $T = 10$ . A strong linear relationship is observed for all models for new sequences made of unseen token combinations.

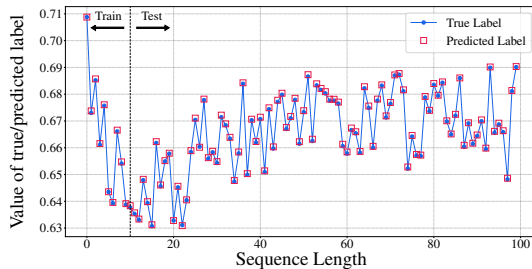


Figure 3: A transformer model with softmax attention with *two* hidden layer MLP for  $\omega$  trained on sequences of length up to  $T = 10$  length generalizes to sequences of length up to 100.

## 4.2 COMPOSITIONAL GENERALIZATION

For compositional generalization, we generate data following the illustration in Figure 1a. During training, we sample each component  $k$  of a token from  $\text{Uniform}[0, 1]$  and accept the sampled sequences that satisfy the following for all components  $i$ :  $-0.5 \leq \sum_{j=1}^T (x_j^k - 0.5) \leq 0.5 \quad \forall k$ , where  $x_j^k$  is the  $k^{\text{th}}$  component of token  $j$ . During testing, we sample  $x_{\leq t}$  from the complementary set of the training set, i.e., corners of hypercube  $[0, 1]^{nt}$ . We present the  $\ell_2$  loss on the test dataset, as well as the mean  $R^2$ , where the results are averaged over 5 seeds. The rest of the details are the same as the previous section, i.e.,  $T = 10$ ,  $n = 20$ ,  $\rho$  is a two hidden layer MLP ( $\phi$  is also a two hidden layer MLP for deep sets). Table 2 shows the test  $\ell_2$  loss and  $R^2$  scores for linear identification.

## 5 DISCUSSION AND LIMITATIONS

Our work is a step towards theoretical foundations of successes and failures of length and compositional generalization in sequence-to-sequence models. We prove simplified versions of the recently proposed RASP conjecture under weaker data diversity assumptions. In our analysis, we make certain simplifications, e.g., on the architectures considered, which motivates some of the important conjectures for future work. The main conjectures go as follows – a) **Conjecture 1:** Theorem 2 and 3 currently incorporate different non-linear attentions but not the softmax attention. We believe these guarantees on transformers extend to softmax attention given the experimental evidence in Section 4. b) **Conjecture 2:** Theorem 2 and 3 use one block of attention and one block of non-linearity. We believe that it is possible to extend these results to more expressive  $\mathcal{H}$ , e.g., with more alternating blocks. c) **Conjecture 3:** Our results focus on the generalization properties of all the possible solutions to risk minimization equation 1. However, in practice the optimization procedure may be biased towards a subset of those. Does accounting for the bias of optimization procedure give way to explaining the success of generalization in even higher capacity architectures?

## REFERENCES

- 540  
541  
542 Emmanuel Abbe, Samy Bengio, Aryo Lotfi, and Kevin Rizk. Generalization on the unseen, logic  
543 reasoning and degree curriculum. *arXiv preprint arXiv:2301.13105*, 2023.
- 544  
545 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
546 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
547 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 548  
549 Kartik Ahuja, Ethan Caballero, Dinghuai Zhang, Jean-Christophe Gagnon-Audet, Yoshua Bengio,  
550 Ioannis Mitliagkas, and Irina Rish. Invariance principle meets information bottleneck for out-of-  
551 distribution generalization. *Advances in Neural Information Processing Systems*, 34:3438–3450,  
2021.
- 552  
553 Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Ramasesh, Am-  
554 brose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length generalization  
555 in large language models. *Advances in Neural Information Processing Systems*, 35:38546–38556,  
556 2022.
- 557  
558 Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.  
559 *arXiv preprint arXiv:1907.02893*, 2019.
- 560  
561 Sanjeev Arora and Anirudh Goyal. A theory for emergence of complex skills in language models.  
562 *arXiv preprint arXiv:2307.15936*, 2023.
- 563  
564 Robert B Ash and Catherine A Doléans-Dade. *Probability and measure theory*. Academic press,  
2000.
- 565  
566 Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Prov-  
567 able in-context learning with in-context algorithm selection. *arXiv preprint arXiv:2306.04637*,  
2023.
- 568  
569 Shai Ben-David and Ruth Urner. Domain adaptation—can quantity compensate for quality? *Annals*  
570 *of Mathematics and Artificial Intelligence*, 70:185–202, 2014.
- 571  
572 Jack Brady, Roland S Zimmermann, Yash Sharma, Bernhard Schölkopf, Julius Von Kügelgen, and  
573 Wieland Brendel. Provably learning object-centric representations. In *International Conference*  
574 *on Machine Learning*, pp. 3038–3062. PMLR, 2023.
- 575  
576 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Ka-  
577 mar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general  
intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- 578  
579 Kamalika Chaudhuri, Kartik Ahuja, Martin Arjovsky, and David Lopez-Paz. Why does throwing  
580 away data improve worst-group error? In *International Conference on Machine Learning*, pp.  
581 4144–4188. PMLR, 2023.
- 582  
583 Shai Ben David, Tyler Lu, Teresa Luu, and Dávid Pál. Impossibility theorems for domain adaptation.  
584 In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*,  
pp. 129–136. JMLR Workshop and Conference Proceedings, 2010.
- 585  
586 Grégoire Delétang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot Catt,  
587 Chris Cundy, Marcus Hutter, Shane Legg, Joel Veness, et al. Neural networks and the chomsky  
588 hierarchy. *arXiv preprint arXiv:2207.02098*, 2022.
- 589  
590 Kefan Dong and Tengyu Ma. First steps toward understanding the extrapolation of nonlinear models  
591 to unseen domains. *arXiv preprint arXiv:2211.11719*, 2022.
- 592  
593 Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Sean  
Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, et al. Faith and fate: Limits of  
transformers on compositionality. *Advances in Neural Information Processing Systems*, 36, 2024.

- 594 Cian Eastwood, Shashank Singh, Andrei L Nicolicioiu, Marin Vlastelica Pogančić, Julius von  
595 Kügelgen, and Bernhard Schölkopf. Spuriousity didn’t kill the classifier: Using invariant pre-  
596 dictions to harness spurious features. *Advances in Neural Information Processing Systems*, 36,  
597 2024.
- 598 Jerry A Fodor and Zenon W Pylyshyn. Connectionism and cognitive architecture: A critical analy-  
599 sis. *Cognition*, 28(1-2):3–71, 1988.
- 600 Jonathan Gordon, David Lopez-Paz, Marco Baroni, and Diane Bouchacourt. Permutation equivari-  
601 ant models for compositional generalization in language. In *International Conference on Learning*  
602 *Representations*, 2019.
- 603 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*  
604 *preprint arXiv:2312.00752*, 2023.
- 605 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured  
606 state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- 607 Ingo Gühring, Mones Raslan, and Gitta Kutyniok. Expressivity of deep neural networks. *arXiv*  
608 *preprint arXiv:2007.04759*, 34, 2020.
- 609 Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth  
610 Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. Textbooks are  
611 all you need. *arXiv preprint arXiv:2306.11644*, 2023.
- 612 Geoffrey E Hinton. Mapping part-whole hierarchies into connectionist networks. *Artificial Intelli-*  
613 *gence*, 46(1-2):47–75, 1990.
- 614 Kaiying Hou, David Brandfonbrener, Sham Kakade, Samy Jelassi, and Eran Malach. Universal  
615 length generalization with turing programs. *arXiv preprint arXiv:2407.03310*, 2024.
- 616 Weizhe Hua, Zihang Dai, Hanxiao Liu, and Quoc Le. Transformer quality in linear time. In *Inter-*  
617 *national Conference on Machine Learning*, pp. 9099–9117. PMLR, 2022.
- 618 Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia Bruni. Compositionality decomposed:  
619 How do neural networks generalise? *Journal of Artificial Intelligence Research*, 67:757–795,  
620 2020.
- 621 Samy Jelassi, Stéphane d’Ascoli, Carles Domingo-Enrich, Yuhuai Wu, Yuanzhi Li, and François  
622 Charton. Length generalization in arithmetic transformers. *arXiv preprint arXiv:2306.15400*,  
623 2023.
- 624 Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Payel Das, and Siva  
625 Reddy. The impact of positional encoding on length generalization in transformers. *Advances*  
626 *in Neural Information Processing Systems*, 36, 2024.
- 627 Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoen-  
628 coders and nonlinear ica: A unifying framework. In *International Conference on Artificial Intel-*  
629 *ligence and Statistics*, pp. 2207–2217. PMLR, 2020.
- 630 Najoung Kim and Tal Linzen. Cogs: A compositional generalization challenge based on semantic  
631 interpretation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language*  
632 *Processing (EMNLP)*, pp. 9087–9105, 2020.
- 633 Sébastien Lachapelle, Divyat Mahajan, Ioannis Mitliagkas, and Simon Lacoste-Julien. Additive  
634 decoders for latent variables identification and cartesian-product extrapolation. *arXiv preprint*  
635 *arXiv:2307.02598*, 2023.
- 636 Brenden Lake and Marco Baroni. Generalization without systematicity: On the compositional skills  
637 of sequence-to-sequence recurrent networks. In *International conference on machine learning*,  
638 pp. 2873–2882. PMLR, 2018.
- 639 Baihan Lin, Djallel Bouneffouf, and Irina Rish. A survey on compositional generalization in appli-  
640 cations. *arXiv preprint arXiv:2302.01067*, 2023.

- 648 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-*  
649 *ence on Learning Representations*, 2019. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Bkg6RiCqY7)  
650 [Bkg6RiCqY7](https://openreview.net/forum?id=Bkg6RiCqY7).
- 651 Joao Loula, Marco Baroni, and Brenden M Lake. Rearranging the familiar: Testing compositional  
652 generalization in recurrent networks. *arXiv preprint arXiv:1807.07545*, 2018.
- 653 Boris Mityagin. The zero set of a real analytic function. *arXiv preprint arXiv:1512.07276*, 2015.
- 654 Richard Montague. Pragmatics and intensional logic. *Synthese*, 22(1):68–94, 1970.
- 655 Aviv Netanyahu, Abhishek Gupta, Max Simchowitz, Kaiqing Zhang, and Pulkit Agrawal. Learning  
656 to extrapolate: A transductive approach. *arXiv preprint arXiv:2304.14329*, 2023.
- 657 Antonio Orvieto, Soham De, Caglar Gulcehre, Razvan Pascanu, and Samuel L Smith. On the univer-  
658 sality of linear recurrences followed by nonlinear projections. *arXiv preprint arXiv:2307.11888*,  
659 2023a.
- 660 Antonio Orvieto, Samuel L Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pas-  
661 canu, and Soham De. Resurrecting recurrent neural networks for long sequences. *arXiv preprint*  
662 *arXiv:2303.06349*, 2023b.
- 663 Tony Plate et al. Holographic reduced representations: Convolution algebra for compositional dis-  
664 tributed representations. In *IJCAI*, pp. 30–35, 1991.
- 665 Alexandre Rame, Matthieu Kirchmeyer, Thibaud Rahier, Alain Rakotomamonjy, Patrick Gallinari,  
666 and Matthieu Cord. Diverse weight averaging for out-of-distribution generalization. *Advances in*  
667 *Neural Information Processing Systems*, 35:10821–10836, 2022.
- 668 Geoffrey Roeder, Luke Metz, and Durk Kingma. On linear identifiability of learned representations.  
669 In *International Conference on Machine Learning*, pp. 9030–9039. PMLR, 2021.
- 670 Mateo Rojas-Carulla, Bernhard Schölkopf, Richard Turner, and Jonas Peters. Invariant models for  
671 causal transfer learning. *Journal of Machine Learning Research*, 19(36):1–34, 2018.
- 672 Simon Schug, Seijin Kobayashi, Yassir Akram, Maciej Wołczyk, Alexandra Proca, Johannes  
673 Von Oswald, Razvan Pascanu, João Sacramento, and Angelika Steger. Discovering modular  
674 solutions that generalize compositionally. *arXiv preprint arXiv:2312.15001*, 2023.
- 675 Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representa-  
676 tions. *arXiv preprint arXiv:1803.02155*, 2018.
- 677 Kai Shen, Junliang Guo, Xu Tan, Siliang Tang, Rui Wang, and Jiang Bian. A study on relu and  
678 softmax in transformer. *arXiv preprint arXiv:2302.06461*, 2023.
- 679 Xinwei Shen and Nicolai Meinshausen. Engression: Extrapolation for nonlinear regression? *arXiv*  
680 *preprint arXiv:2307.00835*, 2023.
- 681 Sania Sinha, Tanawan Premsri, and Parisa Kordjamshidi. A survey on compositional learning of ai  
682 models: Theoretical and experimetal practices. *arXiv preprint arXiv:2406.08787*, 2024.
- 683 Eduardo D Sontag. Neural nets as systems models and controllers. In *Proc. Seventh Yale Workshop*  
684 *on Adaptive and Learning Systems*, volume 73, 1992.
- 685 Kaya Stechly, Matthew Marquez, and Subbarao Kambhampati. Gpt-4 doesn’t know it’s wrong: An  
686 analysis of iterative prompting for reasoning problems. *arXiv preprint arXiv:2310.12397*, 2023.
- 687 Masashi Sugiyama, Matthias Krauledat, and Klaus-Robert Müller. Covariate shift adaptation by  
688 importance weighted cross validation. *Journal of Machine Learning Research*, 8(5), 2007.
- 689 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
690 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
691 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

702 Karthik Valmeekam, Matthew Marquez, and Subbarao Kambhampati. Can large language models  
703 really improve by self-critiquing their own plans? *arXiv preprint arXiv:2310.08118*, 2023.  
704

705 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
706 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*  
707 *tion processing systems*, 30, 2017.

708 Thaddäus Wiedemer, Jack Brady, Alexander Panfilov, Attila Juhos, Matthias Bethge, and Wieland  
709 Brendel. Provable compositional generalization for object-centric learning. *arXiv preprint*  
710 *arXiv:2310.05327*, 2023a.

711 Thaddäus Wiedemer, Prasanna Mayilvahanan, Matthias Bethge, and Wieland Brendel. Composi-  
712 tional generalization from first principles. *arXiv preprint arXiv:2307.05596*, 2023b.  
713

714 Mitchell Wortsman, Jaehoon Lee, Justin Gilmer, and Simon Kornblith. Replacing softmax with relu  
715 in vision transformers. *arXiv preprint arXiv:2309.08586*, 2023.  
716

717 Changnan Xiao and Bing Liu. Conditions for length generalization in learning reasoning skills.  
718 *arXiv preprint arXiv:2311.16173*, 2023.

719 Zhenlin Xu, Marc Niethammer, and Colin A Raffel. Compositional generalization in unsupervised  
720 compositional representation learning: A study on disentanglement and emergent language. *Ad-*  
721 *vances in Neural Information Processing Systems*, 35:25074–25087, 2022.  
722

723 Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and  
724 Alexander J Smola. Deep sets. *Advances in neural information processing systems*, 30, 2017.

725 Min Zhang, Jianfeng He, Shuo Lei, Murong Yue, Linhan Wang, and Chang-Tien Lu. Can llm find  
726 the green circle? investigation and human-guided tool manipulation for compositional general-  
727 ization. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal*  
728 *Processing (ICASSP)*, pp. 11996–12000. IEEE, 2024.

729 Hattie Zhou, Arwen Bradley, Etai Littwin, Noam Razin, Omid Saremi, Josh Susskind, Samy Bengio,  
730 and Preetum Nakkiran. What algorithms can transformers learn? a study in length generalization.  
731 *arXiv preprint arXiv:2310.16028*, 2023.  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

756	APPENDIX	
757		
758	CONTENTS	
759		
760		
761	<b>A Illustration of the test support for compositional generalization</b>	<b>16</b>
762		
763	<b>B Supplement on Related Works</b>	<b>16</b>
764		
765	<b>C Proofs</b>	<b>16</b>
766		
767	C.1 Deep sets . . . . .	17
768	C.1.1 Extending Theorem 1 to discrete tokens . . . . .	20
769	C.1.2 Extending Theorem 1 to $\omega$ from $C^1$ -diffeomorphisms class . . . . .	21
770		
771	C.2 Transformers . . . . .	23
772	C.2.1 Extension of Theorem 2 to incorporate positional encodings . . . . .	25
773	C.2.2 Extending Theorem 2 to discrete tokens . . . . .	27
774	C.2.3 Extending Theorem 2 to $\omega$ from $C^1$ -diffeomorphisms . . . . .	28
775	C.2.4 Extending Theorem 3 to incorporate positional encodings . . . . .	30
776	C.2.5 Extending Theorem 3 to incorporate multiple attention heads . . . . .	31
777		
778	C.3 State space models . . . . .	32
779	C.3.1 Extending Theorem 4 to discrete tokens . . . . .	34
780		
781	C.4 Vanilla RNNs . . . . .	35
782	C.4.1 Extending Theorem 5 to discrete tokens . . . . .	38
783		
784	C.5 Finite Hypothesis Class . . . . .	41
785		
786		
787	<b>D Experiments</b>	<b>42</b>
788		
789	D.1 Length Generalization . . . . .	43
790	D.2 Compositional Generalization . . . . .	44
791	D.3 Failure Cases . . . . .	48
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

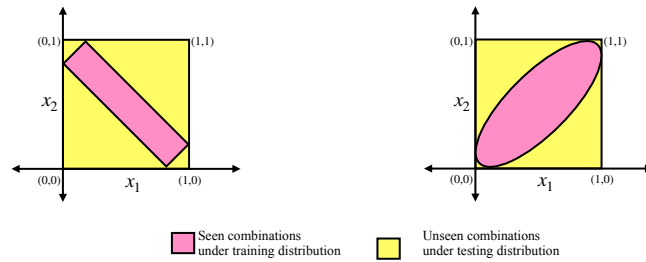


Figure 4: Illustration of observed support and its Cartesian product. These examples illustrate the support of the observed training data distribution can be much smaller than the Cartesian product of the support of the individual tokens.

## A ILLUSTRATION OF THE TEST SUPPORT FOR COMPOSITIONAL GENERALIZATION

The notion of compositional generalization we study requires us to evaluate the model on the Cartesian product of the support of individual token distributions. In Figure 4, we give some additional examples besides the ones shown in Figure 1a to illustrate the difference between the Cartesian product set and the observed support. These examples illustrate the support of the observed training data distribution can be much smaller than the Cartesian product of the support of the individual tokens.

## B SUPPLEMENT ON RELATED WORKS

We briefly discuss some other relevant works here, which could not be mentioned in the main body due to space constraints. In Schug et al. (2023), the authors exploit compositionality in the context of meta learning, where each task parameter is specified via a linear combination of some basis module parameters. They construct an approach that achieves provable compositional guarantees and outperforms meta-learning approaches such as MAML and ANIL. In a concurrent work Hou et al. (2024) propose an interesting scratch pad strategy inspired from the operation of Turing machines. They call this strategy Turing programs. The scratch pad emulates the operation of a Turing machine. The authors argue that there exist short RASP program ( $O(n)$  length) that can simulate the operation of a Turing machine for sufficiently long number of steps ( $O(\exp(n))$ ). Our current framework does not incorporate scratchpad strategies into it, and it is a promising future work to investigate provable length generalization guarantees with scratchpad.

## C PROOFS

In all the results that follow, we work with standard topology in  $\mathbb{R}^{nt}$ , where  $n$  is dimension of each token and  $t$  is the sequence length. We remind the reader of the definition of a regular closed set – if a set is equal to the closure of its interior, then it is said to be a regular closed set. In all the results that follow, we either work with continuous random variables for which the Radon-Nikodym derivative of  $X_{\leq t}$  is absolutely continuous w.r.t Lebesgue measure  $\forall t$  or we work with discrete random variables for which the Radon-Nikodym derivative of  $X_{\leq t}$  is absolutely continuous w.r.t counting measure  $\forall t$ .

**Lemma 1.** *Let  $\mathcal{X} \subseteq \mathbb{R}^n$ . If  $f : \mathcal{X} \rightarrow \mathbb{R}^m$  and  $g : \mathcal{X} \rightarrow \mathbb{R}^m$  are continuously differentiable functions that satisfy  $f(x) = g(x)$  almost everywhere in  $\mathcal{X}$ , where  $\mathcal{X}$  is a regular non-empty closed set, then  $f(x) = g(x), \forall x \in \mathcal{X}$  and  $\nabla f(x) = \nabla g(x), \forall x \in \mathcal{X}$ , where  $\nabla$  is the Jacobian w.r.t  $x$ .*



864 *Proof.* Let us consider the interior of  $\mathcal{X}$  and denote it as  $\mathcal{X}^{\text{int}}$ . We first argue that the two functions  
 865  $f$  and  $g$  are equal at all points in the interior. Suppose there exists a point  $x \in \mathcal{X}^{\text{int}}$  at which  
 866  $f(x) \neq g(x)$ . Consider a ball centered at  $x$  of radius  $r$  denoted as  $B(x, r) \subset \mathcal{X}^{\text{int}}$  (such a ball exists  
 867 as this point is in the interior of  $\mathcal{X}$ ). We argue that there exists at least one point  $x_1$  in this ball at  
 868 which  $f(x_1) = g(x_1)$ . If this were not the case, then the equality will not hold on the entire ball,  
 869 which would contradict the condition that the equality  $f(x) = g(x)$  can only be violated on a set  
 870 of measure zero. Note this condition holds true for all  $r > 0$ . Suppose the distance of  $x_1$  from  $x$  is  
 871  $r_1 \leq r$ . Consider another ball with radius  $r_2 < r_1$  and let  $x_2 \in B(x, r_2)$  where the equality holds.  
 872 By repeating this argument, we can construct a sequence  $\{x_k\}_{k \in \mathbb{N}}$  that converges to  $x$ , where  $\mathbb{N}$  is  
 873 the set of natural numbers. On this sequence, the following conditions hold.

$$874 \quad f(x_k) = g(x_k), \forall k \in \mathbb{N} \quad (5)$$

876 Further, from the continuity of  $f$  and  $g$  it follows that

$$877 \quad \lim_{k \rightarrow \infty} f(x_k) = f(x), \lim_{k \rightarrow \infty} g(x_k) = g(x) \quad (6)$$

879 Combining the above two conditions, we get that  $f(x) = g(x)$ . This leads to a contradiction since  
 880 we assumed that  $f(x) \neq g(x)$ . Thus there can be no such  $x$  in the interior at which  $f(x) \neq g(x)$ .  
 881 From this it follows that  $f(x) = g(x)$  for all  $x \in \mathcal{X}^{\text{int}}$ . Now let us consider the closure of  $\mathcal{X}^{\text{int}}$ ,  
 882 which is  $\mathcal{X}$  itself since it is a regular closed set. Every point  $x \in \mathcal{X}$  in the closure can be expressed  
 883 as limit of points in  $\mathcal{X}^{\text{int}}$ . Consider an  $x \in \mathcal{X}$  and from the definition of regular closed set it follows  
 884 that  $\lim_{k \rightarrow \infty} x_k = x$ , where  $x_k \in \mathcal{X}^{\text{int}}$ . We already know from the fact that  $f$  and  $g$  are equal in the  
 885 interior

$$886 \quad f(x_k) = g(x_k), \forall k \in \mathbb{N} \quad (7)$$

887 From the continuity of  $f$  and  $g$  it follows

$$888 \quad \lim_{k \rightarrow \infty} f(x_k) = f(x), \lim_{k \rightarrow \infty} g(x_k) = g(x) \quad (8)$$

891 Combining the above two we get that  $f(x) = g(x)$  for all  $x \in \mathcal{X}$ . After this we can use Lemma 6  
 892 from (Lachapelle et al., 2023) to conclude that  $\nabla f(x) = \nabla g(x), \forall x \in \mathcal{X}$ . We repeat their proof here  
 893 for completeness. For all points in the interior of  $\mathcal{X}$ , it follows that  $\nabla f(x) = \nabla g(x), \forall x \in \mathcal{X}^{\text{int}}$ .

894 Now consider any point  $x \in \mathcal{X}$ . Since  $\mathcal{X}$  is a regular closed set,  $\lim_{k \rightarrow \infty} x_k = x$ . Since each  $x_k$  is  
 895 in the interior of  $\mathcal{X}$  it follows that

$$896 \quad \nabla f(x_k) = \nabla g(x_k), \forall k \in \mathbb{N} \quad (9)$$

898 From the continuity of  $\nabla f$  and  $\nabla g$  it follows that

$$900 \quad \lim_{k \rightarrow \infty} \nabla f(x_k) = \nabla f(x), \lim_{k \rightarrow \infty} \nabla g(x_k) = \nabla g(x) \quad (10)$$

902 Combining the above conditions, we get that  $\nabla f(x) = \nabla g(x)$ . This completes the proof.

904  $\square$

## 905 C.1 DEEP SETS

907 In this section, we provide the proofs for length and compositional generalization for deep sets. We  
 908 first provide the proof for Theorem 1, followed by Corollary 1, where we establish linear identifica-  
 909 tion. We then present the discrete tokens counterpart to Theorem 1 in Theorem 7. In the next part of  
 910 this section, we extend Theorem 1 with  $\omega$  from  $C^1$ -diffeomorphism in Theorem 8.

911 We restate the theorems from the main body for convenience of the reader. In what follows, we  
 912 remind the reader that we denote the labeling function  $f(\mathcal{X}) = \rho(\sum_{x \in \mathcal{X}} \phi(x))$  and the function  
 913 learned is denoted as  $h(\mathcal{X}) = \omega(\sum_{x \in \mathcal{X}} \psi(x))$ .

914 **Theorem 1.** *If  $\mathcal{H}$  follows Assumption 1, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ ,  $\text{supp}(X_j) =$   
 915  $[0, 1]^n, \forall j \geq 1$ , and the regular closedness condition in Assumption 2 holds, then the model  
 916 trained to minimize the risk in equation 1 with  $\ell_2$  loss generalizes to all sequences in the hyper-  
 917 cube  $[0, 1]^{nt}, \forall t \geq 1$  and thus achieves length and compositional generalization.*

918 *Proof.* Consider any  $h$  that solves equation 1. Since  $\ell$  is  $\ell_2$  loss and realizability condition holds,  
 919  $f$  is one of the optimal solutions to equation 1. For all  $x_{\leq T} \in \text{supp}(X_{\leq T})$  except over a set of  
 920 measure zero the following condition holds

$$921 \quad h(x_{\leq T}) = f(x_{\leq T}). \quad (11)$$

922 The above follows from the fact that  $h$  solves equation 1, i.e.,  $\mathbb{E}[\|h - f\|^2] = 0$  and from The-  
 923 orem 1.6.6. (Ash & Doléans-Dade, 2000). Since  $\text{supp}(X_{\leq T})$  is regular closed,  $f, h$  are both  
 924 continuously differentiable, we can use Lemma 1, it follows that the above equality holds for all  
 925  $x_{\leq T} \in \text{supp}(X_{\leq T})$ . From realizability condition it follows that true  $f(x_{\leq T}) = \rho\left(\sum_{j \leq T} \phi(x_j)\right)$ .  
 926 We substitute the functional decomposition from Assumption 1 to get

$$927 \quad \omega\left(\sum_{j \leq T} \psi(x_j)\right) = \rho\left(\sum_{j \leq T} \phi(x_j)\right). \quad (12)$$

928  $\omega$  and  $\rho$  are both single layer perceptron with a bijective activation  $\sigma$ . We substitute the parametric  
 929 form of  $\omega$  and  $\rho$  to obtain

$$930 \quad \sigma\left(A \sum_{j \leq T} \psi(x_j)\right) = \sigma\left(B \sum_{j \leq T} \phi(x_j)\right) \implies A \sum_{j \leq T} \psi(x_j) = B \sum_{j \leq T} \phi(x_j). \quad (13)$$

931 The second equality in the above simplification follows from the fact that the activation  $\sigma$  is bijective,  
 932 the inputs to  $\sigma$  are equal. We take the derivative of the expressions above w.r.t  $x_r$  to get the following  
 933 condition and equate them (follows from Lemma 1). For all  $x_r \in \text{supp}(X_r)$ , i.e.,  $x_r \in [0, 1]^n$ ,

$$934 \quad \nabla_{x_r}\left(A \sum_{j \leq T} \psi(x_j)\right) = \nabla_{x_r}\left(B \sum_{j \leq T} \phi(x_j)\right). \quad (14)$$

935 We drop the subscript  $r$  to simplify the notation. Therefore, for all  $x \in [0, 1]^n$

$$936 \quad A \nabla_x \psi(x) = B \nabla_x \phi(x), \quad (15)$$

937 where  $\nabla_x \psi(x)$  is the Jacobian of  $\psi(x)$  w.r.t  $x$  and  $\nabla_x \phi(x)$  is the Jacobian of  $\phi(x)$  w.r.t  $x$ . We now  
 938 take the derivative w.r.t some component  $x^k$  of vector  $x = [x^1, \dots, x^n]$ . Denote the components  
 939 other than  $k$  as  $x^{-k} = x \setminus x^k$ . From the above condition, it follows that for all  $x \in [0, 1]^n$

$$940 \quad A \frac{\partial \psi(x)}{\partial x^k} = B \frac{\partial \phi(x)}{\partial x^k}. \quad (16)$$

941 Using fundamental theorem of calculus, we can integrate both sides for fixed  $x^{-k}$  and obtain the  
 942 following for all  $x^k \in [0, 1]$ ,

$$943 \quad A\psi(x^k, x^{-k}) = B\phi(x^k, x^{-k}) + C_k(x^{-k}) \implies A\psi(x) - B\phi(x) = C_k(x^{-k}). \quad (17)$$

944 The above condition is true of all  $k \in \{1, \dots, n\}$ . Hence, we can deduce that for all  $x \in [0, 1]^n$  and  
 945 for  $k \neq j$ , where  $j, k \in \{1, \dots, d\}$ ,

$$946 \quad A\psi(x) - B\phi(x) = C_k(x^{-k}) = C_j(x^{-j}). \quad (18)$$

947 Take the partial derivative of  $C_k(x^{-k})$  and  $C_j(x^{-j})$  w.r.t  $x^j$  to obtain, for all  $x^j \in [0, 1]$ ,

$$\frac{\partial C_k(x^{-k})}{\partial x^j} = \frac{\partial C_j(x^{-j})}{\partial x^j} = 0. \quad (19)$$

In the above simplification, we use the fact that  $\forall x^j \in [0, 1]$ ,  $\frac{\partial C_j(x^{-j})}{\partial x^j} = 0$ . Therefore,  $C_k(x^{-k})$  cannot depend on  $x^j$ . We can apply the same condition on all  $j \neq k$ . As a result,  $C_k(x^{-k})$  is a fixed constant vector denoted as  $C$ . We write this as

$$A\psi(x) = B\phi(x) + C. \quad (20)$$

Substitute the above into  $A \sum_{j \leq T} \psi(x_j) = B \sum_{j \leq T} \phi(x_j)$  to obtain

$$B \sum_{j \leq T} \phi(x_j) + CT = B \sum_{j \leq T} \phi(x_j) \implies C = 0. \quad (21)$$

Therefore, we get

$$\forall x \in [0, 1]^n, A\psi(x) = B\phi(x). \quad (22)$$

We now consider any sequence  $x_{\leq \tilde{T}}$  from  $[0, 1]^{n\tilde{T}}$ . The prediction made by  $h$  is

$$h(x_{\leq \tilde{T}}) = \sigma\left(A \sum_{j \leq \tilde{T}} \psi(x_j)\right) = \sigma\left(B \sum_{j \leq \tilde{T}} \phi(x_j)\right) = f(x_{\leq \tilde{T}}). \quad (23)$$

We use equation 22 in the simplification above. From the above, we can conclude that  $h$  continues to be optimal for distribution  $\mathbb{P}_{X_{\leq \tilde{T}}}$ .

□

**Corollary 1.** *If  $\mathcal{H}$  follows Assumption 1 with the condition that the output layer weight matrix is left invertible, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ ,  $\text{supp}(X_j) = [0, 1]^n$ ,  $\forall j \geq 1$ , and the regular closedness condition in Assumption 2 holds, then the model trained to minimize the risk in equation 1 with  $\ell_2$  loss achieves linear identification. Further, under the stated conditions linear identification is necessary for compositional and length generalization.*

*Proof.* We follow the exact same steps as in the previous proof of Theorem 1 up to equation 22. We restate equation 22 below.

$$\begin{aligned} \forall x \in [0, 1]^n, A\psi(x) &= B\phi(x) \\ \psi(x) &= A^{-1}B\phi(x) \end{aligned} \quad (24)$$

The above condition establishes linear identification, i.e., the learned model's representation for a token is a linear transform of the true model's representation. From the above, we can write that  $x_{\leq \tilde{T}}$  from  $[0, 1]^{n\tilde{T}}$

$$\sum_{j \leq \tilde{T}} \psi(x_j) = A^{-1}B \sum_{j \leq \tilde{T}} \phi(x_j) \quad (25)$$

The above shows linear relationship holds for the entire sequence as well.

Now let us turn to the part on necessity. From the proof of previous theorem, we know that

$$\forall x_{\leq T} \in \text{supp}(X_{\leq T}), \sigma\left(A \sum_{j \leq T} \psi(x_j)\right) = \sigma\left(B \sum_{j \leq T} \phi(x_j)\right) \implies \forall x \in [0, 1]^n, A\psi(x) = B\phi(x) \quad (26)$$

Thus if  $\forall x \in [0, 1]^n, A\psi(x) = B\phi(x)$  is not true, then  $\sigma\left(A \sum_{j \leq T} \psi(x_j)\right) = \sigma\left(B \sum_{j \leq T} \phi(x_j)\right)$  cannot be true either. Therefore, in the absence of linear identification neither length nor compositional generalization are achievable. □

**Remarks** A few remarks and observations from the proof are in order. Firstly, observe that we do not require  $\phi$  and  $\psi$  to have the same output dimension for the above proof to go through. Secondly, in Theorem 1, we observe all the labels from  $t = 1$  to  $T$ , i.e.,  $y_1$  to  $y_T$ . The result continues to hold if we only observe label at length  $T$ , i.e.,  $y_T$ . Finally, we make an observation in this result, which would apply to all the subsequent theorems. The definition of compositional generalization requires generalization to the Cartesian product over sequences of length  $T$ , where  $T$  is the training length. Since our model generalizes to the hypercube  $[0, 1]^{nt}$ ,  $\forall t$ , we achieve compositional generalization even beyond the training lengths.

### C.1.1 EXTENDING THEOREM 1 TO DISCRETE TOKENS

In our discussion, we have focused on settings where the support of each token has a non-empty interior (Assumption 2). In practice of language modeling, we use discrete tokens and hence Assumption 2 does not hold anymore. In this section, we discuss the adaptation of results for deepsets to setting when the the support of tokens is a finite set.

**Assumption 9.** *The marginal support of token for all positions is the same and denoted as  $\mathcal{X}$ . The joint support of first and second token is  $\mathcal{X} \times \mathcal{X}$ .*

**Theorem 7.** *If  $\mathcal{H}$  follows Assumption 1, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , and Assumption 9 holds, then the model trained to minimize the risk in equation 1 with  $\ell_2$  loss generalizes to all sequences in the hypercube  $[0, 1]^{nt}$ ,  $\forall t \geq 1$  and thus achieves length and compositional generalization.*

*Proof.* Consider any  $h$  that solves equation 1. Since  $\ell$  is  $\ell_2$  loss and realizability condition holds,  $f$  is one of the optimal solutions to equation 1. For all  $x_{\leq T} \in \text{supp}(X_{\leq T})$

$$h(x_{\leq T}) = f(x_{\leq T}). \quad (27)$$

The above follows from the fact that  $h$  solves equation 1, i.e.,  $\mathbb{E}[\|h - f\|^2] = 0$  and the fact that tokens are discrete random vectors. From realizability condition it follows that true  $f(x_{\leq T}) = \rho\left(\sum_{j \leq T} \phi(x_j)\right)$ . We substitute the functional decomposition from Assumption 1 to get

$$\omega\left(\sum_{j \leq T} \psi(x_j)\right) = \rho\left(\sum_{j \leq T} \phi(x_j)\right). \quad (28)$$

$\omega$  and  $\rho$  are both single layer perceptron with a bijective activation  $\sigma$ . We substitute the parametric form of  $\omega$  and  $\rho$  to obtain

$$\sigma\left(A \sum_{j \leq T} \psi(x_j)\right) = \sigma\left(B \sum_{j \leq T} \phi(x_j)\right) \implies A \sum_{j \leq T} \psi(x_j) = B \sum_{j \leq T} \phi(x_j). \quad (29)$$

The second equality in the above simplification follows from the fact that the activation  $\sigma$  is bijective, the inputs to  $\sigma$  are equal.

From Assumption 9, it follows that for all  $x_1, x_2 \in \mathcal{X} \times \mathcal{X}$

$$A\phi(x_1) + A\phi(x_2) = B\psi(x_1) + B\psi(x_2) \quad (30)$$

Set  $x_1 = x_2 = x$  (we can set this value due to Assumption 9) we get

$$\forall x \in \mathcal{X}, A\psi(x) = B\phi(x). \quad (31)$$

We now consider any sequence  $x_{\leq \bar{T}}$  from  $\mathcal{X}^{\bar{T}}$ . The prediction made by  $h$  is

$$h(x_{\leq \bar{T}}) = \sigma\left(A \sum_{j \leq \bar{T}} \psi(x_j)\right) = \sigma\left(B \sum_{j \leq \bar{T}} \phi(x_j)\right) = f(x_{\leq \bar{T}}). \quad (32)$$

We use equation 22 in the simplification above. From the above, we can conclude that  $h$  continues to be optimal for distribution  $\mathbb{P}_{X_{\leq T}}$ .

□

### C.1.2 EXTENDING THEOREM 1 TO $\omega$ FROM $C^1$ -DIFFEOMORPHISMS CLASS

**Assumption 10.** Each function in  $\mathcal{H}$  is expressed as  $h(x_1, \dots, x_i) = \omega(\sum_{j=1}^i \psi(x_j))$ , where  $\omega$  is a  $C^1$ -diffeomorphism.

**Assumption 11.** The joint support  $\text{supp}(X_{\leq i})$  is a regular closed set for all  $i \leq T$ . The support of all tokens is equal, i.e.,  $\text{supp}(X_j) = [0, 1]^n$ , where  $j \geq 1$ . The support of  $[\phi(X_1), \phi(X_2)]$  is  $\mathbb{R}^{2m}$ , where  $\phi$  is the embedding function for the labeling function  $f(\mathcal{X}) = \rho(\sum_{x \in \mathcal{X}} \phi(x))$ .

We provide a remark on the assumption and where it is used following the proof of the next theorem.

**Theorem 8.** If  $\mathcal{H}$  follows Assumption 10, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , and a further assumption on the support (Assumption 11) holds, then the model trained to minimize the risk in equation 1 (with  $T \geq 2$ ) with  $\ell_2$  loss generalizes to all sequences in  $[0, 1]^{nt}$ ,  $\forall t \geq 1$  and thus achieves length and compositional generalization.

*Proof.* We start with the same steps as earlier proofs and equate the prediction of  $h$  and  $f$ . We first use the fact  $h(x_{\leq i}) = f(x_{\leq i})$  everywhere in the support. For all  $x_{\leq i} \in \text{supp}(X_{\leq i})$

$$\begin{aligned} \omega\left(\sum_{j \leq i} \psi(x_j)\right) &= \rho\left(\sum_{j \leq i} \phi(x_j)\right) \implies \sum_{j \leq i} \psi(x_j) = \omega^{-1} \circ \rho\left(\sum_{j \leq i} \phi(x_j)\right) \implies \\ \sum_{j \leq i} \psi(x_j) &= a\left(\sum_{j \leq i} \phi(x_j)\right), \end{aligned} \quad (33)$$

where  $a = \omega^{-1} \circ \rho$ . In the above simplification, we used the parametric form for the true labeling function and the learned labeling function and use the invertibility of  $\omega$ . Let us consider the setting when  $i = 1$ . In that case summation involves only one term. Substitute  $x_1 = x$ . We obtain  $\forall x \in [0, 1]^n$ ,

$$\psi(x) = a(\phi(x)). \quad (34)$$

The above expression implies that  $\psi$  bijectively identifies  $\phi$ . Let us consider the setting when  $i = 2$ . Substitute  $x_1 = x$  and  $x_2 = y$ . We obtain

$$a(\phi(x)) + a(\phi(y)) = a(\phi(x) + \phi(y)). \quad (35)$$

We now use the that assumption  $[\phi(x), \phi(y)]$  spans  $\mathbb{R}^{2m}$ , where  $\phi(x)$  and  $\phi(y)$  individually span  $\mathbb{R}^m$ . Substitute  $\phi(x) = \alpha$  and  $\phi(y) = \beta$ . We obtain  $\forall \alpha \in \mathbb{R}^m, \forall \beta \in \mathbb{R}^m$

$$a(\alpha) + a(\beta) = a(\alpha + \beta). \quad (36)$$

Observe that  $a(0) = 0$  (substitute  $\alpha = \beta = 0$  in the above).

We use equation 36 to show that  $a$  is linear. To show that, we need to argue that  $a(c\alpha) = ca(\alpha)$  as we already know  $a$  satisfies additivity condition.

From the identity above, we want to show that equation 70  $a(p\alpha) = pa(\alpha)$ , where  $p$  is some integer.

Substitute  $\beta = -\alpha$  in  $a(\alpha + \beta) = a(\alpha) + a(\beta)$ . We obtain  $a(0) = a(\alpha) + a(-\alpha) \implies a(-\alpha) = -a(\alpha)$ . Suppose  $p$  is a positive integer. We simplify  $a(p\alpha)$  as follows  $a(\alpha + (p-1)\alpha) = a(\alpha) + a((p-1)\alpha)$ . Repeating this simplification, we get  $a(p\alpha) = pa(\alpha)$ . Suppose  $p$  is a negative integer. We can write  $a(p\alpha) = a(-p \times -\alpha) = -pa(-\alpha)$ . Since  $a(-\alpha) = -a(\alpha)$ , we get  $a(p\alpha) = pa(\alpha)$ .

Suppose  $c$  is some rational number, i.e.,  $c = p/q$ , where  $p$  and  $q$  are non-zero integers. We already know  $a(p\alpha) = pa(\alpha)$ . Further, we obtain

$$a(q\frac{1}{q}\alpha) = qa(\frac{1}{q}\alpha) \implies a(\frac{1}{q}\alpha) = \frac{1}{q}a(\alpha), \text{ where } q \text{ is some integer.}$$

Now combine these  $a(p/q\alpha) = pa(1/q\alpha) = \frac{p}{q}a(\alpha)$ . We have established the homogeneity condition for rationals.

We will now use the continuity of the function  $a$  and density of rationals to extend the claim for irrationals. Suppose  $c$  is some irrational. Define a sequence of rationals that approach  $c$  (this follows from the fact that rationals are dense in  $\mathbb{R}$ ).

$$a(c\alpha) = a(\lim_{n \rightarrow \infty} q_n\alpha) = \lim_{n \rightarrow \infty} a(q_n\alpha).$$

In the second equality above, we use the definition of continuity ( $a$  is continuous since composition of continuous functions is continuous). We can also use the property that we already showed for rationals to further simplify

$$\lim_{n \rightarrow \infty} a(q_n\alpha) = a(\alpha) \lim_{n \rightarrow \infty} q_n = ca(\alpha).$$

Observe that  $a : \mathbb{R}^m \rightarrow \mathbb{R}^m$  and for any  $\alpha, \beta \in \mathbb{R}^m$   $a(\alpha + \beta) = a(\alpha) + a(\beta)$  and  $a(c\alpha) = ca(\alpha)$ . From the definition of a linear map it follows that  $a$  is linear. As a result, we can write  $\forall x \in [0, 1]^n$

$$\psi(x) = A(\phi(x)) \tag{37}$$

Observe that  $a$  is invertible because both  $\rho$  and  $\omega$  are invertible. As a result, we know that  $A$  is an invertible matrix. From this we get

$$\phi(x) = A^{-1}\psi(x) = C(\psi(x)) \tag{38}$$

For all  $z \in \mathbb{R}^m$ , we obtain

$$a(z) = \rho^{-1} \circ \omega(z) = Cz \implies \omega(z) = \rho(Cz)$$

Let us consider any sequence  $x_{\leq \bar{T}} \in [0, 1]^{n\bar{T}}$ . We use the above conditions

$$\omega\left(\sum_{j \leq \bar{T}} \psi(x_j)\right) = \rho\left(C \sum_{j \leq \bar{T}} \psi(x_j)\right) = \rho\left(\sum_{j \leq \bar{T}} \phi(x_j)\right)$$

Thus we obtain length and compositional generalization. □

**Remark on Assumption 11** In Assumption 11, we require that the support of  $[\phi(X_1), \phi(X_2)]$  is  $\mathbb{R}^{2m}$ . This assumption is used in the proof in equation 36. We used this assumption to arrive at  $a(\alpha + \beta) = a(\alpha) + a(\beta), \forall \alpha, \beta \in \mathbb{R}^m$ . We then used continuity of  $a$  to conclude  $a$  is linear. Now suppose  $[\phi(X_1), \phi(X_2)]$  is some subset  $\mathcal{Z} \subseteq \mathbb{R}^{2m}$ . We believe that it is possible to extend the result to more general  $\mathcal{Z}$ , it might still be possible to arrive at  $a$  is linear. We leave this investigation to future work.

**Remark on expressivity under Assumption 10 and Assumption 11** Assumption 11 requires  $\omega$  is a  $C^1$ -diffeomorphism. Suppose the label is one dimensional, i.e.,  $m = 1$ . From Assumption 11 output dimension of  $\phi$  is restricted to be one dimensional. Consider the map  $h(x_1, \dots, x_i) = \rho(\sum_{j \leq i} \phi(x_j))$ . The output dimension of  $\phi$  is required to grow with sequence length to express all permutation invariant maps (See Theorem 7 in (Zaheer et al., 2017)). Thus by restricting the output dimension of  $\phi$  to one, we cannot express all the permutation invariant maps.

**Multiplication operator** Consider the multiplication operator  $y_i = \prod_{j=1}^i x_j$ , where each  $x_i > 0$ . Observe that we can rewrite this as  $y_i = \exp(\sum_{j=1}^i \log(x_j))$ . This operator is realizable on deep sets from hypothesis class described by Assumption 10 with  $\omega = \exp$  and  $\psi = \log$ . In Assumption 11, we require the support of  $[\phi(X_1), \phi(X_2)]$  to be  $\mathbb{R}^2$ . We let the support of  $X_1$  and  $X_2$  be  $(0, \infty)$ . In Assumption 11 we require that the support of each token was equal to  $[0, 1]$ . However, the proof of Theorem 8 still goes through even if support is  $(0, \infty)$ . Hence, we can use Theorem 8 to conclude that deep sets trained to predict the output of multiplication can multiply longer sequences and also multiply new token combinations.

## C.2 TRANSFORMERS

In this section, we provide the proofs for length and compositional generalization for transformers. We first provide the proof for Theorem 2, followed by Corollary 2, where we establish linear identification. We present an extension of Theorem 2 to incorporate positional encoding in Theorem 9. We then present the discrete tokens counterpart to Theorem 2 in Theorem 10. In the next part of this section, we extend Theorem 2 with  $\omega$  from  $C^1$ -diffeomorphism in Theorem 3. Theorem 11 adapts Theorem 3 to incorporate positional encodings.

We restate the theorems from the main body for convenience of the reader. In what follows, we want to remind the reader we denote the labeling function  $f(x_1, \dots, x_i) = \rho(\sum_{j \leq i} \phi(x_i, x_j))$  and the function learned is denoted as  $h(x_1, \dots, x_i) = \omega(\sum_{j \leq i} \psi(x_i, x_j))$ .

**Theorem 2.** *If  $\mathcal{H}$  follows Assumption 3, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ ,  $\text{supp}(X_i, X_j) = [0, 1]^{2n}$ ,  $\forall i \neq j$  and the regular closedness condition in Assumption 2 holds, then the model trained to minimize the risk in equation 1 (with  $T \geq 2$ ) with  $\ell_2$  loss generalizes to all sequences in the hypercube  $[0, 1]^{nt}$ ,  $\forall t \geq 1$  and thus achieves length and compositional generalization.*

*Proof.* Consider any  $h$  that solves equation 1. Since  $\ell$  is  $\ell_2$  loss and realizability condition holds,  $f$  is one of the optimal solutions to equation 1. For all  $i \leq T, x_{\leq i} \in \text{supp}(X_{\leq i})$  except over a set of measure zero the following condition holds

$$h(x_{\leq i}) = f(x_{\leq i}). \quad (39)$$

The above follows from the fact that  $h$  solves equation 1, i.e.,  $\mathbb{E}[\|h - f\|^2] = 0$  and from Theorem 1.6.6. (Ash & Doléans-Dade, 2000). Since  $\text{supp}(X_{\leq i})$  is regular closed,  $f, h$  are both continuously differentiable, we can use Lemma 1, it follows that the above equality holds for all  $x_{\leq i} \in \text{supp}(X_{\leq i})$ . From realizability condition it follows that true  $f(x_{\leq i}) = \rho(\sum_{k \leq i} \phi(x_i, x_k))$ . We substitute the parametric forms from Assumption 3 to get

$$\omega\left(\sum_{k \leq i} \frac{1}{i} \cdot \psi(x_i, x_k)\right) = \rho\left(\sum_{k \leq i} \frac{1}{i} \cdot \phi(x_i, x_k)\right). \quad (40)$$

Since  $\omega$  and  $\rho$  are single layer perceptron with bijective activation  $\sigma$ . We substitute the parametric form of  $\omega$  and  $\rho$  to obtain the following condition. For all  $x_{\leq i} \in \text{supp}(X_{\leq i})$ ,

$$\sigma\left(A \sum_{k \leq i} \frac{1}{i} \cdot \psi(x_i, x_k)\right) = \sigma\left(B \sum_{k \leq i} \frac{1}{i} \cdot \phi(x_i, x_k)\right) \implies A \sum_{k \leq i} \psi(x_i, x_k) = B \sum_{k \leq i} \phi(x_i, x_k). \quad (41)$$

The second equality follows from the fact that the activation  $\sigma$  is bijective and hence the inputs to  $\sigma$  are equal. We take the derivative of the expressions above w.r.t  $x_j$  to get the following (follows from Lemma 1). For  $j < i$  (there exists a  $j < i$  as  $T \geq 2$  and we can set  $i \geq 2$ ) and for all  $x_j \in \text{supp}(X_j)$ , i.e.,  $x_j \in [0, 1]^n$ ,

1242

1243

1244

1245

1246

1247

$$\begin{aligned} \nabla_{x_j} \left( A \sum_{k \leq i} \psi(x_i, x_k) \right) &= \nabla_{x_j} \left( B \sum_{k \leq i} \phi(x_i, x_k) \right) \implies \\ A \nabla_{x_j} \psi(x_i, x_j) &= B \nabla_{x_j} \phi(x_i, x_j), \end{aligned} \quad (42)$$

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

where  $\nabla_{x_j} \psi(x_i, x_j)$ ,  $\nabla_{x_j} \phi(x_i, x_j)$  are the Jacobians of  $\psi$  and  $\phi$  w.r.t  $x_j$  for a fixed  $x_i$ . Note that  $A \nabla_{x_j} \psi(x_i, x_j) = B \nabla_{x_j} \phi(x_i, x_j)$  holds for all  $x_i \in [0, 1]^n$ ,  $x_j \in [0, 1]^n$  (here we use the fact that joint support of every pair of tokens spans  $2n$  dimensional unit hypercube assumed in the Theorem 9). In this equality, we now consider the derivative w.r.t some component  $x_j^k$  of  $x_j$ . Denote the remaining components as  $x_j^{-k}$ . From the above condition it follows that for all  $x_i \in [0, 1]^n$ ,  $x_j \in [0, 1]^n$ ,

$$A \frac{\partial \psi(x_i, x_j)}{\partial x_j^k} = B \frac{\partial \phi(x_i, x_j)}{\partial x_j^k}. \quad (43)$$

1259

1260

1261

1262

Using fundamental theorem of calculus, we can integrate both sides for fixed  $x_j^{-k}$  and obtain the following for all  $x_j^k \in [0, 1]$ ,

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

$$\begin{aligned} A \psi(x_i, [x_j^k, x_j^{-k}]) &= B \phi(x_i, [x_j^k, x_j^{-k}]) + C_k(x_i, x_j^{-k}) = \\ A \psi(x_i, x_j) &= B \phi(x_i, x_j) + C_k(x_i, x_j^{-k}). \end{aligned} \quad (44)$$

The same condition is true of all  $k$ . Hence,  $\forall x_i \in [0, 1]^d, \forall x_j \in [0, 1]^d$  and for  $k \neq q$ , where  $q, k \in \{1, \dots, d\}$ ,

$$A \psi(x_i, x_j) - B \phi(x_i, x_j) = C_k(x_i, x_j^{-k}) = C_q(x_i, x_j^{-q}). \quad (45)$$

Take the partial derivative of both sides w.r.t  $x_j^q$  to obtain,  $\forall x_j^q \in [0, 1]$ ,

$$\frac{\partial C_k(x_i, x_j^{-k})}{\partial x_j^q} = \frac{\partial C_q(x_i, x_j^{-q})}{\partial x_j^q} = 0. \quad (46)$$

Therefore,  $C_k(x_i, x_j^{-k})$  cannot depend on  $x_j^q$ . We can apply the same condition on all  $q \neq k$ . As a result,  $C_k(x_i, x_j^{-k})$  is only a function of  $x_i$  denoted as  $C(x_i)$ . Therefore, for  $j < i$  and for all  $x_i \in [0, 1]^n, x_j \in [0, 1]^n$

$$A \psi(x_i, x_j) = B \phi(x_i, x_j) + C(x_i). \quad (47)$$

If we substitute  $x_i = x_j = x$ , then the above equality extends for  $i = j$  and thus we get

$$A \psi(x_i, x_i) = B \phi(x_i, x_i) + C(x_i). \quad (48)$$

Substitute the above equation 47, equation 48 into  $A \sum_{k \leq i} \psi(x_i, x_k) = B \sum_{k \leq i} \phi(x_i, x_k)$  to obtain

$$B \sum_{k \leq i} \phi(x_i, x_k) + (i)C(x_i) = B \sum_{k \leq i} \phi(x_i, x_k) \implies C(x_i) = 0. \quad (49)$$

Thus we obtain

$$\forall x_i \in [0, 1]^n, x_j \in [0, 1]^n \quad A \psi(x_i, x_j) = B \phi(x_i, x_j). \quad (50)$$

We now consider any sequence  $x_{\leq \bar{T}} \in [0, 1]^{n\bar{T}}$ . The prediction made by  $h$  is



$$h(x_{\leq \bar{T}}) = \sigma\left(A \sum_{j \leq \bar{T}} \psi(x_{\bar{T}}, x_j)\right) = \sigma\left(B \sum_{j \leq \bar{T}} \phi(x_{\bar{T}}, x_j)\right) = f(x_{\leq \bar{T}}) \quad (51)$$

We use equation 50 in the simplification above. From the above, we can conclude that  $h$  continues to be optimal for all sequences in  $[0, 1]^{n\bar{T}}$ .

□

**Corollary 2.** *If  $\mathcal{H}$  follows Assumption 3 with the condition that the output layer weight matrix is left invertible, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ ,  $\text{supp}(X_i, X_j) = [0, 1]^{2n}$ ,  $\forall i \neq j$  and the regular closedness condition in Assumption 2 holds, then the model trained to minimize the risk in equation 1 (with  $T \geq 2$ ) with  $\ell_2$  loss achieves linear identification. Further, linear identification is necessary for both length and compositional generalization.*

*Proof.* We follow the exact same steps as in the previous proof of Theorem 2 up to equation 50. We restate equation 50 below.

$$\forall x_i \in [0, 1]^n, x_j \in [0, 1]^n \quad A\psi(x_i, x_j) = B\phi(x_i, x_j) \quad (52)$$

$$\psi(x_i, x_j) = A^{-1}B\phi(x_i, x_j)$$

In the second step above, we use left invertibility of  $A$ . The above condition establishes linear identification, i.e., the learned model’s representation is a linear transform of the true model’s representation. From this we obtain that for any sequence  $x_{\leq \bar{T}} \in [0, 1]^{n\bar{T}}$

$$\sum_{j \leq \bar{T}} \psi(x_{\bar{T}}, x_j) = A^{-1}B\left(\sum_{j \leq \bar{T}} \phi(x_{\bar{T}}, x_j)\right) \quad (53)$$

The above establishes a linear relationship between the learned representation of the sequence and the representation of the sequence under the true model. Now let us turn to the part on necessity. From the proof of previous theorem, we know that

$$\omega\left(\sum_{k \leq i} \frac{1}{i} \cdot \psi(x_i, x_k)\right) = \rho\left(\sum_{k \leq i} \frac{1}{i} \cdot \phi(x_i, x_k)\right) \implies \forall x_i \in [0, 1]^n, x_j \in [0, 1]^n \quad A\psi(x_i, x_j) = B\phi(x_i, x_j) \quad (54)$$

Thus from the above it follows that in the absence of linear identification neither length nor compositional generalization are achievable. □

**On the absence of labels at all lengths from  $t = 1$  to  $t = T$**  A few important remarks are to follow. In the proof above, we do not require to observe all the labels from  $t = 1$  to  $t = T$ , where  $T \geq 2$ . The proof goes through provided we observe data at two different lengths.

### C.2.1 EXTENSION OF THEOREM 2 TO INCORPORATE POSITIONAL ENCODINGS

In what follows, we extend the above result (Theorem 2) to incorporate positional encoding. We start with extension of the hypothesis class to incorporate positional encoding.

**Assumption 12.** *Each function in the hypothesis class  $\mathcal{H}$  used by the learner is given as  $h(x_1, \dots, x_i) = \omega\left(\sum_{j \leq i} \frac{1}{i} \psi_{i-j}(x_i, x_j)\right)$ , where  $\omega$  is a single layer perceptron with continuously differentiable bijective activation (e.g., sigmoid) and each  $\psi_k$  is a map that is differentiable. Also,  $\psi_k = 0$  for  $k \geq T_{\max}$ , i.e., two tokens that are sufficiently far apart do not interact.*

In the above assumption, we incorporate relative positional encodings by making the function  $\psi_{i-j}$  depend on the relative positional difference between token  $x_i$  and token  $x_j$ . We would like to emphasize the reasons why we assume that the tokens that are sufficiently far apart do not interact. Suppose  $T_{\max} = \infty$ , which implies tokens at all positions interact. As a result, during training since

we only see sequences of finite length  $T$ , we will not see the effect of interactions of tokens that are separated at a distance larger than  $T$  on the data generation, which makes it impossible to learn anything about  $\phi_{i-j}$ , where  $i - j \geq T - 1$ .

In the theorem that follows, we show that we can achieve length and compositional generalization for the above hypothesis class.

**Theorem 9.** *If  $\mathcal{H}$  follows Assumption 12, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ ,  $\text{supp}(X_i, X_j) = [0, 1]^{2n}$ ,  $\forall i \neq j \in \{1, \dots, \infty\}$ , the regular closedness condition in Assumption 2 holds and  $T \geq T_{\max} \geq 2$ , then the model trained to minimize the risk in equation 1 with  $\ell_2$  loss generalizes to all sequences in the hypercube  $[0, 1]^{nt}$ ,  $\forall t$  and thus achieves length and compositional generalization.*

*Proof.* Consider any  $h$  that solves equation 1. Since  $\ell$  is  $\ell_2$  loss and realizability condition holds,  $f$  is one of the optimal solutions to equation 1. For all  $i \leq T$  and for all  $x_{\leq i} \in \text{supp}(X_{\leq i})$  except over a set of measure zero the following condition holds

$$h(x_{\leq i}) = f(x_{\leq i}). \quad (55)$$

The above follows from the fact that  $h$  solves equation 1, i.e.,  $\mathbb{E}[\|h - f\|^2] = 0$  and from Theorem 1.6.6. (Ash & Doléans-Dade, 2000). Since  $\text{supp}(X_{\leq i})$  is regular closed,  $f, h$  are both continuously differentiable, we can use Lemma 1, it follows that the above equality holds for all  $x_{\leq i} \in \text{supp}(X_{\leq i})$ . From realizability condition it follows that true  $f(x_{\leq i}) = \rho\left(\sum_{k \leq i} \phi_{i-k}(x_i, x_k)\right)$ . We substitute the parametric forms from Assumption 3 to get

$$\omega\left(\sum_{k \leq i} \frac{1}{i} \cdot \psi_{i-k}(x_i, x_k)\right) = \rho\left(\sum_{k \leq i} \frac{1}{i} \cdot \phi_{i-k}(x_i, x_k)\right). \quad (56)$$

Since  $\omega$  and  $\rho$  are single layer perceptron with bijective activation  $\sigma$ . We substitute the parametric form of  $\omega$  and  $\rho$  to obtain the following condition. For all  $x_{\leq i} \in \text{supp}(X_{\leq i})$ ,

$$\begin{aligned} \sigma\left(A \sum_{k \leq i} \frac{1}{i} \cdot \psi_{i-k}(x_i, x_k)\right) &= \sigma\left(B \sum_{k \leq i} \frac{1}{i} \cdot \phi_{i-k}(x_i, x_k)\right) \implies \\ A \sum_{k \leq i} \psi_{i-k}(x_i, x_k) &= B \sum_{k \leq i} \phi_{i-k}(x_i, x_k). \end{aligned} \quad (57)$$

The second equality follows from the fact that the activation  $\sigma$  is bijective and hence the inputs to  $\sigma$  are equal. We take the derivative of the expressions above w.r.t  $x_j$  to get the following (follows from Lemma 1). The equality holds true for all  $i \leq T$ .

From the above, we can use  $i = 1$  and obtain

$$A\psi_0(x_1, x_1) = B\phi_0(x_1, x_1), \forall x_1 \in [0, 1]^n.$$

From  $i = 2$ , we obtain

$$A\psi_0(x_2, x_2) + A\psi_1(x_2, x_1) = B\phi_0(x_2, x_2) + B\phi_1(x_2, x_1), \forall x_1 \in [0, 1]^n, x_2 \in [0, 1]^n$$

Combining the two conditions we get

$$A\psi_1(x_2, x_1) = B\phi_1(x_2, x_1), \forall x_1 \in [0, 1]^n, x_2 \in [0, 1]^n.$$

We can use this argument and arrive at

$$A\psi_{i-1}(x_i, x_1) = B\phi_{i-1}(x_i, x_1), \forall x_i \in [0, 1]^n, x_1 \in [0, 1]^n, \forall i \leq T.$$

1404 Thus we obtain

$$1405 \quad \forall i - j \leq T - 1, \forall x_i \in [0, 1]^n, x_j \in [0, 1]^n, \quad A\psi_{i-j}(x_i, x_j) = B\phi_{i-j}(x_i, x_j). \quad (58)$$

1407 From Assumption 12 and  $T \geq T_{\max}$ , we already know that

$$1409 \quad \forall i - j \geq T, \forall x_i \in [0, 1]^n, x_j \in [0, 1]^n, \quad A\psi_{i-j}(x_i, x_j) = B\phi_{i-j}(x_i, x_j) = 0. \quad (59)$$

1411 If  $A$  is left invertible, then the above condition implies that linear representation identification is  
1412 necessary for both compositional and length generalization.

1413 We now consider any sequence  $x_{\leq \bar{T}} \in [0, 1]^{n\bar{T}}$ . The prediction made by  $h$  is

$$1416 \quad h(x_{\leq \bar{T}}) = \sigma\left(A \sum_{j \leq \bar{T}} \psi_{\bar{T}-j}(x_{\bar{T}}, x_j)\right) = \sigma\left(B \sum_{j \leq \bar{T}} \phi_{\bar{T}-j}(x_{\bar{T}}, x_j)\right) = f(x_{\leq \bar{T}}) \quad (60)$$

1420 We use equation 50 in the simplification above. From the above, we can conclude that  $h$  continues  
1421 to be optimal for all sequences in  $[0, 1]^{n\bar{T}}$ .

1422  $\square$

## 1424 C.2.2 EXTENDING THEOREM 2 TO DISCRETE TOKENS

1426 In the above result we used Assumption 2. In practice of language modeling, we use discrete tokens  
1427 and hence Assumption 2 does not hold anymore. In this section, we discuss the adaptation of results  
1428 for transformers to setting when the the support of tokens is a finite set.

1429 **Assumption 13.** *The marginal support of token for all positions is the same and denoted as  $\mathcal{X}$ . The  
1430 joint support of first three tokens is  $\mathcal{X} \times \mathcal{X} \times \mathcal{X}$ .*

1431 **Theorem 10.** *If  $\mathcal{H}$  follows Assumption 3, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , and Assumption 13 holds, then the model trained to minimize the risk in equation 1 (with  $T \geq 2$ ) with  $\ell_2$  loss generalizes to all sequences in the hypercube  $[0, 1]^{nt}$ ,  $\forall t \geq 1$  and thus achieves length and compositional generalization.*

1436 *Proof.* Consider any  $h$  that solves equation 1. Since  $\ell$  is  $\ell_2$  loss and realizability condition holds,  
1437  $f$  is one of the optimal solutions to equation 1. For all  $i \leq T, x_{\leq i} \in \text{supp}(X_{\leq i})$  the following  
1438 condition holds

$$1440 \quad h(x_{\leq i}) = f(x_{\leq i}). \quad (61)$$

1441 The above follows from the fact that  $h$  solves equation 1, i.e.,  $\mathbb{E}[\|h - f\|^2] = 0$  and from the  
1442 fact that the tokens are discrete random vectors. From realizability condition it follows that true  
1443  $f(x_{\leq i}) = \rho\left(\sum_{k \leq i} \phi(x_i, x_k)\right)$ . We substitute the parametric forms from Assumption 3 to get

$$1445 \quad \omega\left(\sum_{k \leq i} \frac{1}{i} \cdot \psi(x_i, x_k)\right) = \rho\left(\sum_{k \leq i} \frac{1}{i} \cdot \phi(x_i, x_k)\right). \quad (62)$$

1448 Since  $\omega$  and  $\rho$  are single layer perceptron with bijective activation  $\sigma$ . We substitute the parametric  
1450 form of  $\omega$  and  $\rho$  to obtain the following condition. For all  $x_{\leq i} \in \text{supp}(X_{\leq i})$ ,

$$1452 \quad \sigma\left(A \sum_{k \leq i} \frac{1}{i} \cdot \psi(x_i, x_k)\right) = \sigma\left(B \sum_{k \leq i} \frac{1}{i} \cdot \phi(x_i, x_k)\right) \implies A \sum_{k \leq i} \psi(x_i, x_k) = B \sum_{k \leq i} \phi(x_i, x_k). \quad (63)$$

1456 The second equality follows from the fact that the activation  $\sigma$  is bijective and hence the inputs to  $\sigma$   
1457 are equal.

From Assumption 13, it follows that for all  $x_1, x_2, x_3 \in \mathcal{X} \times \mathcal{X} \times \mathcal{X}$

$$A\psi(x_3, x_1) + A\psi(x_3, x_2) = B\phi(x_3, x_1) + B\phi(x_3, x_2) \quad (64)$$

Set  $x_1 = x_2$  (we can do so owing to Assumption 13).

Thus we obtain

$$\forall x_i \in \mathcal{X}, x_j \in \mathcal{X} \quad A\psi(x_i, x_j) = B\phi(x_i, x_j). \quad (65)$$

We now consider any sequence  $x_{\leq \bar{T}} \in \mathcal{X}^{\bar{T}}$ . The prediction made by  $h$  is

$$h(x_{\leq \bar{T}}) = \sigma\left(A \sum_{j \leq \bar{T}} \psi(x_{\bar{T}}, x_j)\right) = \sigma\left(B \sum_{j \leq \bar{T}} \phi(x_{\bar{T}}, x_j)\right) = f(x_{\leq \bar{T}}) \quad (66)$$

We use equation 65 in the simplification above. From the above, we can conclude that  $h$  continues to be optimal for all sequences in  $[0, 1]^{n\bar{T}}$ .

□

### C.2.3 EXTENDING THEOREM 2 TO $\omega$ FROM $C^1$ -DIFFEOMORPHISMS

**Theorem 3.** *If  $\mathcal{H}$  follows Assumption 4, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , and a further assumption on the support (Assumption 5) holds, then the model trained to minimize the risk in equation 1 (with  $T \geq 3$ ) with  $\ell_2$  loss generalizes to all sequences in  $[0, 1]^{nt}, \forall t \geq 1$  and thus achieves length and compositional generalization.*

*Proof.* We start with the same steps as earlier proofs and equate the prediction of  $h$  and  $f$ . We first use the fact  $h(x_{\leq i}) = f(x_{\leq i}), \forall i \leq T$  almost everywhere in the support. We can use the continuity of  $h, f$  and regular closedness of the support to extend the equality to all points in the support (follows from the first part of Lemma 1) to obtain the following. For all  $x_{\leq i} \in \text{supp}(X_{\leq i})$

$$\begin{aligned} \omega\left(\sum_{j < i} \frac{1}{i-1} \cdot \psi(x_i, x_j)\right) &= \rho\left(\sum_{j < i} \frac{1}{i-1} \cdot \phi(x_i, x_j)\right) \implies \\ \sum_{j < i} \frac{1}{i-1} \psi(x_i, x_j) &= \omega^{-1} \circ \rho\left(\sum_{j < i} \frac{1}{i-1} \cdot \phi(x_i, x_j)\right) \implies \\ \sum_{j < i} \frac{1}{i-1} \psi(x_i, x_j) &= a\left(\sum_{j < i} \frac{1}{i-1} \phi(x_i, x_j)\right), \end{aligned} \quad (67)$$

where  $a = \omega^{-1} \circ \rho$ . In the above simplification, we used the parametric form for the true labeling function and the learned labeling function and use the invertibility of  $\omega$ . Let us consider the setting when  $i = 2$ . In that case summation involves only one term. Substitute  $x_1 = y$  and  $x_2 = x$ . We obtain  $\forall x \in [0, 1]^n, y \in [0, 1]^n$ ,

$$\psi(x, y) = a(\phi(x, y)). \quad (68)$$

The above expression implies that  $\psi$  bijectively identifies  $\phi$ . Let us consider the setting when  $i = 3$  (this is possible since  $T \geq 3$ ). We substitute  $x_3 = x, x_2 = y, x_1 = z$  and obtain

$$\frac{1}{2} \left[ a(\phi(x, y)) + a(\phi(x, z)) \right] = a\left(\frac{1}{2} (\phi(x, y) + \phi(x, z))\right). \quad (69)$$

Substitute  $\phi(x, y) = \alpha$  and  $\phi(x, z) = \beta$ . In the simplification that follows, we use the that assumption  $[\phi(x, y), \phi(x, z)]$  spans  $\mathbb{R}^{2m}$ , where  $\phi(x, y)$  and  $\phi(x, z)$  individually span  $\mathbb{R}^m$ .

1512

1513

1514

$$\frac{1}{2}(a(\alpha) + a(\beta)) = a\left(\frac{1}{2}(\alpha + \beta)\right). \quad (70)$$

1515

1516

Observe that  $a(0) = 0$  because  $\omega^{-1} \circ \rho(0) = 0$  because  $\omega^{-1}(0) = \rho(0) = 0$ .

1517

1518

1519

1520

1521

$$\begin{aligned} \frac{1}{2}(a(2\alpha) + a(0)) &= a\left(\frac{1}{2}(2\alpha + 0)\right) \\ a(2\alpha) &= 2a(\alpha) \end{aligned} \quad (71)$$

1522

1523

1524

Next, substitute  $\alpha$  with  $2\alpha$  and  $\beta$  with  $2\beta$  in equation 70 to obtain

1525

1526

1527

1528

$$\begin{aligned} \frac{1}{2}(a(2\alpha) + a(2\beta)) &= a\left(\frac{1}{2}(2\alpha + 2\beta)\right) \\ a(\alpha + \beta) &= a(\alpha) + a(\beta) \end{aligned} \quad (72)$$

1529

1530

We use equation 72 to show that  $a$  is linear. To show that, we need to argue that  $a(c\alpha) = ca(\alpha)$  as we already know  $a$  satisfies additivity condition.

1531

1532

Suppose  $c$  is some rational number, i.e.,  $c = p/q$ , where  $p$  and  $q$  are non-zero integers.

1533

From the identity it is clear that  $a(p\alpha) = pa(\alpha)$ , where  $p$  is some integer.

1534

1535

$a(q\frac{1}{q}\alpha) = qa(\frac{1}{q}\alpha) \implies a(\frac{1}{q}\alpha) = \frac{1}{q}a(\alpha)$ , where  $q$  is some integer.

1536

1537

Now combine these  $a(p/q\alpha) = pa(1/q\alpha) = \frac{p}{q}a(\alpha)$ . We have established the homogeneity condition for rationals.

1538

1539

1540

We will now use the continuity of the function  $a$  and density of rationals to extend the claim for irrationals. Suppose  $c$  is some irrational. Define a sequence of rationals that approach  $c$  (this follows from the fact that rationals are dense in  $\mathbb{R}$ ).

1541

1542

$a(c\alpha) = a(\lim_{n \rightarrow \infty} q_n\alpha) = \lim_{n \rightarrow \infty} a(q_n\alpha)$ .

1543

1544

1545

In the second equality above, we use the definition of continuity ( $a$  is continuous since composition of continuous functions is continuous). We can also use the property that we already showed for rationals to further simplify

1546

1547

$\lim_{n \rightarrow \infty} a(q_n\alpha) = a(\alpha) \lim_{n \rightarrow \infty} q_n = ca(\alpha)$ .

1548

1549

1550

Observe that  $a : \mathbb{R}^m \rightarrow \mathbb{R}^m$  and for any  $\alpha, \beta \in \mathbb{R}^m$   $a(\alpha + \beta) = a(\alpha) + a(\beta)$  and  $a(c\alpha) = ca(\alpha)$ . From the definition of a linear map it follows that  $a$  is linear. As a result, we can write  $\forall x \in [0, 1]^n, y \in [0, 1]^n$

1551

1552

$$\psi(x, y) = A(\phi(x, y)) \quad (73)$$

1553

1554

1555

Observe that  $a$  is invertible because both  $\rho$  and  $\omega$  are invertible. As a result, we know that  $A$  is an invertible matrix. From this we get

1556

1557

1558

1559

For all  $z \in \mathbb{R}^m$ , we obtain

1560

1561

1562

$$a(z) = \rho^{-1} \circ \omega(z) = Cz \implies \omega(z) = \rho(Cz)$$

1563

1564

1565

Let us consider any sequence  $x_{\leq \bar{T}} \in [0, 1]^{n\bar{T}}$ . We use the above conditions

$$\omega\left(\sum_{j < \bar{T}} \psi(x_{\bar{T}}, x_j)\right) = \rho\left(C \sum_{j < \bar{T}} \psi(x_{\bar{T}}, x_j)\right) = \rho\left(\sum_{j < \bar{T}} \phi(x_{\bar{T}}, x_j)\right).$$

Thus we obtain length and compositional generalization.  $\square$

**Corollary 3.** *If  $\mathcal{H}$  follows Assumption 4, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , and a further assumption on the support (Assumption 5) holds, then the model trained to minimize the risk in equation 1 (with  $T \geq 3$ ) with  $\ell_2$  loss achieves linear identification. Further, under the stated conditions linear identification is necessary for both length and compositional generalization.*

*Proof.* We follow the exact same steps as in the previous proof of Theorem 3 up to equation 74. We restate equation 74 below.

$$\forall x \in [0, 1]^n, y \in [0, 1]^n, \phi(x, y) = C(\psi(x, y)) \quad (75)$$

The above condition directly implies linear identification. We can use this to obtain that for any sequence  $x_{\leq \bar{T}} \in [0, 1]^{n\bar{T}}$

$$\sum_{j < \bar{T}} \psi(x_{\bar{T}}, x_j) = \sum_{j < \bar{T}} \phi(x_{\bar{T}}, x_j) \quad (76)$$

To show necessity of linear identification, from the proof of Theorem 3 observe that

$$\forall i \leq T, \forall x_{\leq i} \in \text{supp}(X_{\leq i}), \omega\left(\sum_{j < i} \frac{1}{i-1} \cdot \psi(x_i, x_j)\right) = \rho\left(\sum_{j < i} \frac{1}{i-1} \cdot \phi(x_i, x_j)\right) \implies \quad (77)$$

$$\forall x \in [0, 1]^n, y \in [0, 1]^n, \phi(x, y) = C(\psi(x, y))$$

Thus from the above it follows that in the absence of linear identification neither length nor compositional generalization are achievable.  $\square$

**On absence of labels at all lengths from 1 to  $T$**  We argue that the above proof can be adapted to the setting where we do not observe labels at all lengths from 1 to  $T$ . Suppose we only observe label at length  $T$ . Take equation equation 67 and substitute  $x_i = x$  and  $x_j = y$  for all  $j < i$  to obtain the same condition as equation equation 68. Suppose  $T$  is odd and larger than or equal to 3. Fix  $x_i = x, x_{2j-1} = y, \forall j \in \{1, \dots, (T-1)/2\}, x_{2j} = z, \forall j \in \{1, \dots, (T-1)/2\}$ . We obtain the same condition as equation equation 69. Rest of the proof can be adapted using a similar line of reasoning.

**Remark on Assumption 4** We require that the support of  $[\phi(X_1, X_2), \phi(X_1, X_3)]$  is  $\mathbb{R}^{2m}$ . This assumption is used in the proof in equation equation 72. We used this assumption to arrive at  $a(\alpha + \beta) = a(\alpha) + a(\beta), \forall \alpha, \beta \in \mathbb{R}^m$ . We then used continuity of  $a$  to conclude  $a$  is linear. Now suppose  $[\phi(X_1, X_2), \phi(X_1, X_3)]$  is some subset  $\mathcal{Z} \subseteq \mathbb{R}^{2m}$ . We believe that it is possible to extend the result to more general  $\mathcal{Z}$ , it might still be possible to arrive at  $a$  is linear. We leave this investigation to future work.

#### C.2.4 EXTENDING THEOREM 3 TO INCORPORATE POSITIONAL ENCODINGS

We next present the result when  $\omega$  is continuously differentiable and invertible.

**Assumption 14.** *Each function in the hypothesis class  $\mathcal{H}$  used by the learner is given as  $h(x_1, \dots, x_i) = \omega\left(\sum_{j \leq i} \psi_{i-j}(x_i, x_j)\right)$ , where  $\omega$  is a  $C^1$ -diffeomorphism. Also,  $\psi_{i-j} = 0$  for  $i - j > T_{\max} - 1$ , i.e., two tokens that are sufficiently far apart do not interact. For all  $k \leq T_{\max} - 1$  each  $x \in [0, 1]^n, \exists y \in [0, 1]^n$  we  $\psi_k(x, y) = 0$ .*

In the theorem that follows, we require the support of training distribution under consideration is already sufficiently diverse and hence we only seek to prove length generalization guarantees.

**Assumption 15.** *The joint support  $\text{supp}(X_{\leq T}) = [0, 1]^T$ . The support of  $[\phi_1(X_1, X_2), \phi_2(X_1, X_3)]$  is  $\mathbb{R}^{2k}$ , where  $\phi_{i-j}$  is the embedding function for the labeling function  $\rho\left(\sum_{j \leq i} \phi_{i-j}(x_i, x_j)\right)$ .*

**Theorem 11.** *If  $\mathcal{H}$  follows Assumption 14, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , Assumption 15 holds and  $T \geq T_{\max}$ , then the model trained to minimize the risk in equation 1 (with  $T \geq 2$ ) with  $\ell_2$  loss achieves length generalization.*

*Proof.* We start with the same steps as earlier proofs and equate the prediction of  $h$  and  $f$ . We first use the fact  $h(x_{\leq i}) = f(x_{\leq i})$  almost everywhere in the support. We can use the continuity of  $h, f$  and regular closedness of the support to extend the equality to all points in the support (follows from the first part of Lemma 1) to obtain the following. For all  $x_{\leq i} \in \text{supp}(X_{\leq i})$

$$\begin{aligned} \omega\left(\sum_{j<i} \frac{1}{i-1} \psi_{i-j}(x_i, x_j)\right) &= \rho\left(\sum_{j<i} \frac{1}{i-1} \phi_{i-j}(x_i, x_j)\right), \\ \sum_{j<i} \frac{1}{i-1} \psi_{i-j}(x_i, x_j) &= \omega^{-1} \circ \rho\left(\sum_{j<i} \frac{1}{i-1} \phi_{i-j}(x_i, x_j)\right), \\ \sum_{j<i} \frac{1}{i-1} \psi_{i-j}(x_i, x_j) &= a\left(\sum_{j<i} \frac{1}{i-1} \phi_{i-j}(x_i, x_j)\right), \end{aligned} \quad (78)$$

where  $a = \omega^{-1} \circ \rho$ . In the above simplification, we used the parametric form for the true labeling function and the learned labeling function. We also used the invertibility of  $\rho$ . Let us consider the setting when  $i = 2$ . In that case summation involves only one term. Substitute  $x_1 = y$  and  $x_2 = x$ . We obtain  $\forall x \in [0, 1]^n, y \in [0, 1]^n$ ,

$$\psi_1(x, y) = a(\phi_1(x, y)). \quad (79)$$

For  $i = 3$ , substitute  $x_1 = x, x_3 = z$  and set  $x_2 = y$  in such a way that  $\phi_1(x, y) = 0$  (follows from Assumption 14). Thus we obtain

$$\psi_2(x, y) = a(\phi_2(x, y)). \quad (80)$$

Similarly, we can obtain the following. For all  $k \leq T_{\max}$

$$\psi_k(x, y) = a(\phi_k(x, y)). \quad (81)$$

The above expression implies that  $\psi$  bijectively identifies  $\phi$ . Let us consider the setting when  $i = 3$  (this is possible since  $T \geq 3$ ). We substitute  $x_3 = x, x_2 = y, x_1 = z$  to give

$$\frac{1}{2}(a(\phi_1(x, y)) + a(\phi_2(x, z))) = a\left(\frac{1}{2}(\phi_1(x, y) + \phi_2(x, z))\right). \quad (82)$$

We now use the that assumption  $[\phi_1(x, y), \phi_2(x, z)]$  spans  $\mathbb{R}^{2k}$  and substitute  $\phi_1(x, y) = \alpha$  and  $\phi_2(x, z) = \beta$

$$\frac{1}{2}(a(\alpha) + a(\beta)) = a\left(\frac{1}{2}(\alpha + \beta)\right). \quad (83)$$

Rest of the proof follows the same strategy as proof of Theorem 3.  $\square$

### C.2.5 EXTENDING THEOREM 3 TO INCORPORATE MULTIPLE ATTENTION HEADS

Our choice of the architecture did not invoke multiple attention heads. If we include multiple attention heads, then also we can arrive at the same length generalization guarantees. The model class with two attention heads  $\psi_1, \psi_2$  can be stated as follows  $\omega\left(\sum_{j<i} A[\psi_1(x_i, x_j), \psi_2(x_i, x_j)]^\top\right)$ , where  $A$  combines the outputs of the attention heads linearly. Following the same steps of proof of Theorem 3, we obtain the following.

$$\begin{aligned}
& \omega\left(\sum_{j<i} A[\psi_1(x_i, x_j), \psi_2(x_i, x_j)]^\top\right) = \rho\left(\sum_{j<i} B[\phi_1(x_i, x_j), \phi_2(x_i, x_j)]^\top\right), \\
& \omega\left(\sum_{j<i} \tilde{\psi}(x_i, x_j)\right) = \rho\left(\sum_{j<i} \tilde{\phi}(x_i, x_j)\right), \\
& \sum_{j<i} \tilde{\psi}(x_i, x_j) = a\left(\sum_{j<i} \tilde{\phi}(x_i, x_j)\right),
\end{aligned} \tag{84}$$

where  $a = \omega^{-1} \circ \rho$ . In the above simplification, the RHS shows the labeling function and the RHS is the function that is learned. We can follow the same strategy as the proof of Theorem 3 for the rest of the proof. We set  $i = 2$  and obtain a condition similar to equation 68 and for  $i = 3$  we obtain a condition similar to equation 69. Following a similar proof technique, we obtain  $a$  is linear and the proof extends to multiple attention heads.

### C.3 STATE SPACE MODELS

In this section, we first provide the proof to Theorem 4. We then provide Corollary 4, where we describe how the learned representations linearly identify the true representations. In Theorem 12, we present the discrete tokens counterpart to Theorem 4.

**Theorem 4.** *If  $\mathcal{H}$  follows Assumption 6, and the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , and a further condition on the support, i.e., Assumption 7, holds, then the model trained to minimize the risk in equation 1 with  $\ell_2$  loss ( $T \geq 2$ ) achieves length and compositional generalization.*

*Proof.* We start with the same steps as earlier proofs and equate the prediction of  $h$  and  $f$ . We first use the fact  $h(x_{\leq i}) = f(x_{\leq i}), \forall i \leq T$  almost everywhere in the support. We can use the continuity of  $h, f$  and regular closedness of the support to extend the equality to all points in the support (from first part of Lemma 1) to obtain the following. For all  $x_{\leq i} \in \text{supp}(X_{\leq i})$ .

$$\begin{aligned}
f(x_{\leq i}) &= h(x_{\leq i}) = \\
& \rho\left(\sum_{j=0}^{i-1} \Lambda^j Bx_{i-j}\right) = \omega\left(\sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B}x_{i-j}\right) \implies \\
\omega^{-1} \circ \rho\left(\sum_{j=0}^{i-1} \Lambda^j Bx_{i-j}\right) &= \sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B}x_{i-j} = \\
c\left(\sum_{j=0}^{i-1} \Lambda^j Bx_{i-j}\right) &= \sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B}x_{i-j}
\end{aligned} \tag{85}$$

For  $i = 1, \forall x_1 \in \mathbb{R}^n, c(Bx_1) = \tilde{B}x_1$ . Substitute  $Bx_1 = x$ , we obtain  $\forall x \in \mathbb{R}^n, c(x) = \tilde{B}B^{-1}x = Cx$ , where we use the fact that  $Bx_1$  spans  $\mathbb{R}^n$  as  $B$  is invertible.

From linearity of  $c$ , we obtain

$$\omega^{-1} \circ \rho(z) = Cz \implies \rho(z) = \omega(Cz), \forall z \in \mathbb{R}^n \tag{86}$$



We use this linearity of  $c$  to simplify

$$\begin{aligned}
c\left(\sum_{j=0}^{i-1} \Lambda^j B x_{i-j}\right) &= \sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B} x_{i-j} \implies \\
C\left(\sum_{j=0}^{i-1} \Lambda^j B x_{i-j}\right) &= \sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B} x_{i-j} \implies \\
[CB, C\Lambda B, C\Lambda^2 B, \dots, C\Lambda^{i-1} B] \begin{bmatrix} x_i \\ x_{i-2} \\ \vdots \\ x_1 \end{bmatrix} &- [\tilde{B}, \tilde{\Lambda}\tilde{B}, \tilde{\Lambda}^2\tilde{B}, \dots, \tilde{\Lambda}^{i-1}\tilde{B}] \begin{bmatrix} x_i \\ x_{i-2} \\ \vdots \\ x_1 \end{bmatrix} = 0 \implies \\
\left[CB, C\Lambda B, C\Lambda^2 B, \dots, C\Lambda^{i-1} B\right] - [\tilde{B}, \tilde{\Lambda}\tilde{B}, \tilde{\Lambda}^2\tilde{B}, \dots, \tilde{\Lambda}^{i-1}\tilde{B}] \mathbf{X} &= 0, \tag{87}
\end{aligned}$$

where  $\mathbf{X} = \begin{bmatrix} x_i \\ x_{i-2} \\ \vdots \\ x_1 \end{bmatrix}$ .

Denote  $R = [CB, C\Lambda B, C\Lambda^2 B, \dots, C\Lambda^{i-1} B] - [\tilde{B}, \tilde{\Lambda}\tilde{B}, \tilde{\Lambda}^2\tilde{B}, \dots, \tilde{\Lambda}^{i-1}\tilde{B}]$ . We collect a set of points  $\mathbf{X}^+ = [\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(l)}]$  where  $l \geq ni$  and rank of  $\mathbf{X}^+ = ni$  (from Assumption 7). Since the matrix  $\mathbf{X}^+$  is full rank, we have

$$R\mathbf{X}^+ = 0 \implies R = 0.$$

This yields

$$CB = \tilde{B}, C\Lambda B = \tilde{\Lambda}\tilde{B}, \dots, C\Lambda^i B = \tilde{\Lambda}^i \tilde{B}. \tag{88}$$

Observe that from the second equality, we get  $\tilde{\Lambda} = C\Lambda C^{-1}$ . Given the parameters  $(\Lambda, B)$ , the set of parameters  $(\tilde{\Lambda}, \tilde{B})$  that solve the first two equalities are  $\{\tilde{B} \text{ is an arbitrary invertible matrix, } \tilde{\Lambda} = C\Lambda C^{-1}, \text{ where } C = \tilde{B}B^{-1}\}$ .

Take any solution of the first two equalities and compute

$$\tilde{\Lambda}^i \tilde{B} = C\Lambda^i C^{-1} \tilde{B} = C\Lambda^i B, \forall i \geq 1 \tag{89}$$

From equation 89 and equation 86, we obtain that for all  $x_{\leq i} \in \mathbb{R}^{ni}$

$$h(x_{\leq i}) = \omega\left(\sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B} x_{i-j}\right) = \omega\left(C \sum_{j=0}^{i-1} \Lambda^j B x_{i-j}\right) = \rho\left(\sum_{j=0}^{i-1} \Lambda^j B x_{i-j}\right) = f(x_{\leq i}) \tag{90}$$

This establishes both compositional and length generalization. □

**Corollary 4.** *If  $\mathcal{H}$  follows Assumption 6, and the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , and a further condition on the support, i.e., Assumption 7, holds, then the model trained to minimize the risk in equation 1 with  $\ell_2$  loss ( $T \geq 2$ ) achieves linear identification. Further, under the stated conditions linear identification is necessary for both length and compositional generalization.*

*Proof.* We follow the same steps as proof of Theorem 4 up to equation 89. From that we obtain that for all  $x_{\leq i} \in \mathbb{R}^{ni}$

$$\sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B} x_{i-j} = C\left(\sum_{j=0}^{i-1} \Lambda^j B x_{i-j}\right) \tag{91}$$

Recall that  $\sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B}x_{i-j} = \tilde{h}_j$  and  $\sum_{j=0}^{i-1} \Lambda^j Bx_{i-j} = h_j$ . From this it follows that  $\tilde{h}_j = Ch_j$ , which proves that learned hidden state are a linear transform of the hidden state underlying the labeling function. This establishes linear identification.

To show the necessity of linear identification, from the proof of Theorem 4 it follows that

$$\begin{aligned} \forall i \leq T, \forall x_{\leq i} \in \text{supp}(X_{\leq i}), \rho\left(\sum_{j=0}^{i-1} \Lambda^j Bx_{i-j}\right) = \omega\left(\sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B}x_{i-j}\right) &\implies \\ \forall i \geq 1, \forall x_{\leq i} \in \mathbb{R}^{ni}, \tilde{h}_j = Ch_j & \end{aligned} \quad (92)$$

If the latter condition in the above implication does not hold, then the former condition cannot hold. Hence, linear identification is necessary.

□

### C.3.1 EXTENDING THEOREM 4 TO DISCRETE TOKENS

In our discussion, we have focused on settings where the support of each token has a non-empty interior (Assumption 2). In practice of language modeling, we use discrete tokens and hence Assumption 2 does not hold anymore. In this section, we discuss the adaptation of results for SSMs to setting when the the support of tokens is a finite set.

**Assumption 16.** *Each function in the hypothesis class  $\mathcal{H}$  takes a sequence  $\{x_1, \dots, x_i\}$  as input and outputs  $h(x_1, \dots, x_i) = \omega\left(\sum_{j=0}^{i-1} \Lambda^j Bx_{i-j}\right)$ , where  $\omega : \mathbb{R}^k \rightarrow \mathbb{R}^m$  is a single layer perceptron denoted as  $\sigma \circ A$ .  $A$ ,  $B$  and  $\Lambda$  are square invertible. As a result,  $k = m = n$ .*

**Assumption 17.** *For some length  $2 \leq i \leq T$  an there exists in sequences  $x_{\leq i}$  such that their concatenation forms a  $in \times in$  matrix of rank  $in$ .*

**Theorem 12.** *If  $\mathcal{H}$  follows Assumption 16, and the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , and a further condition on the support, i.e., Assumption 17, holds, then the model trained to minimize the risk in equation 1 with  $\ell_2$  loss ( $T \geq 2$ ) achieves length and compositional generalization.*

*Proof.* We start with the same steps as earlier proofs and equate the prediction of  $h$  and  $f$ . We first use the fact  $h(x_{\leq i}) = f(x_{\leq i}), \forall i \leq T$  almost everywhere in the support. We can use the continuity of  $h, f$  and regular closedness of the support to extend the equality to all points in the support (from first part of Lemma 1) to obtain the following. For all  $x_{\leq i} \in \text{supp}(X_{\leq i})$ .

$$\begin{aligned} f(x_{\leq i}) &= h(x_{\leq i}) = \\ \rho\left(\sum_{j=0}^{i-1} \Lambda^j Bx_{i-j}\right) &= \omega\left(\sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B}x_{i-j}\right) \implies \\ \sigma\left(A \sum_{j=0}^{i-1} \Lambda^j Bx_{i-j}\right) &= \sigma\left(\tilde{A} \sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B}x_{i-j}\right) = \\ C\left(\sum_{j=0}^{i-1} \Lambda^j Bx_{i-j}\right) &= \sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B}x_{i-j}, \end{aligned} \quad (93)$$

where  $C = \tilde{A}^{-1}A$ .

We simplify the last identity in the above further.

$$\begin{aligned}
C\left(\sum_{j=0}^{i-1} \Lambda^j B x_{i-j}\right) &= \sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B} x_{i-j} \implies \\
[CB, C\Lambda B, C\Lambda^2 B, \dots, C\Lambda^{i-1} B] \begin{bmatrix} x_i \\ x_{i-2} \\ \vdots \\ x_1 \end{bmatrix} &- [\tilde{B}, \tilde{\Lambda}\tilde{B}, \tilde{\Lambda}^2\tilde{B}, \dots, \tilde{\Lambda}^{i-1}\tilde{B}] \begin{bmatrix} x_i \\ x_{i-2} \\ \vdots \\ x_1 \end{bmatrix} = 0 \implies \\
\left[ [CB, C\Lambda B, C\Lambda^2 B, \dots, C\Lambda^{i-1} B] - [\tilde{B}, \tilde{\Lambda}\tilde{B}, \tilde{\Lambda}^2\tilde{B}, \dots, \tilde{\Lambda}^{i-1}\tilde{B}] \right] \mathbf{X} &= 0,
\end{aligned} \tag{94}$$

where  $\mathbf{X} = \begin{bmatrix} x_i \\ x_{i-2} \\ \vdots \\ x_1 \end{bmatrix}$ .

Denote  $R = [CB, C\Lambda B, C\Lambda^2 B, \dots, C\Lambda^{i-1} B] - [\tilde{B}, \tilde{\Lambda}\tilde{B}, \tilde{\Lambda}^2\tilde{B}, \dots, \tilde{\Lambda}^{i-1}\tilde{B}]$ . We collect a set of points  $\mathbf{X}^+ = [\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(l)}]$  where  $l \geq ni$  and rank of  $\mathbf{X}^+ = ni$  (from Assumption 7). Since the matrix  $\mathbf{X}^+$  is full rank, we have

$$R\mathbf{X}^+ = 0 \implies R = 0.$$

This yields

$$CB = \tilde{B}, C\Lambda B = \tilde{\Lambda}\tilde{B}, \dots, C\Lambda^i B = \tilde{\Lambda}^i \tilde{B}. \tag{95}$$

Observe that from the second equality, we get  $\tilde{\Lambda} = C\Lambda C^{-1}$ . Given the parameters  $(\Lambda, B)$ , the set of parameters  $(\tilde{\Lambda}, \tilde{B})$  that solve the first two equalities are  $\{ \tilde{B} \text{ is an arbitrary invertible matrix, } \tilde{\Lambda} = C\Lambda C^{-1}, \text{ where } C = \tilde{B}B^{-1} \}$ .

Take any solution of the first two equalities and compute

$$\tilde{\Lambda}^i \tilde{B} = C\Lambda^i C^{-1} \tilde{B} = C\Lambda^i B, \forall i \geq 1 \tag{96}$$

From equation 96, we obtain that for all  $x_{\leq i} \in \mathbb{R}^{ni}$

$$h(x_{\leq i}) = \omega\left(\sum_{j=0}^{i-1} \tilde{\Lambda}^j \tilde{B} x_{i-j}\right) = \omega\left(C \sum_{j=0}^{i-1} \Lambda^j B x_{i-j}\right) = \rho\left(\sum_{j=0}^{i-1} \Lambda^j B x_{i-j}\right) = f(x_{\leq i}) \tag{97}$$

This establishes both compositional and length generalization. □

#### C.4 VANILLA RNNs

In this section, we discuss RNNs and present the proof of Theorem 5. We first build some lemmas in the form of Lemma 2 and 3 that are used to prove Theorem 5. In Corollary 5, we explain the learned hidden state are a permutation transform of the true hidden state and also show that its a necessary condition for length and compositional generalization. Finally, in Theorem 13, we present the discrete token counterpart to Theorem 5.

**Lemma 2.** *The  $k^{\text{th}}$  derivative of sigmoid function denoted  $\frac{\partial^k \sigma(s)}{\partial s^k}$  is not zero identically.*

*Proof.* The first derivative of the sigmoid function  $\frac{\partial \sigma(s)}{\partial s} = \sigma(s)(1 - \sigma(s))$ . We argue that the  $\frac{\partial^k \sigma(s)}{\partial s^k}$  is a polynomial in  $\sigma(s)$  with degree  $k + 1$ . Consider the base case of  $k = 1$ . This condition is true

as  $\frac{\partial \sigma(s)}{\partial s} = \sigma(s)(1 - \sigma(s))$ . Now let us assume that  $\frac{\partial^k \sigma(s)}{\partial s^k}$  is a polynomial of degree at most  $k + 1$  denoted as  $P_{k+1}(\sigma(s))$ . We simplify

$$\frac{\partial^k \sigma(s)}{\partial s^k} = P_{k+1}(\sigma(s)) = \sum_{j=1}^{k+1} a_j (\sigma(s))^j$$

We take another derivative of the term above as follows.

$$\frac{\partial^{k+1} \sigma(s)}{\partial s^{k+1}} = \frac{\partial P_{k+1}(\sigma(s))}{\partial s} = \sum_{j=1}^{k+1} a_j \frac{\partial (\sigma(s))^j}{\partial s} = \sum_{j=1}^{k+1} a_j j \sigma(s)^{j-1} (\sigma(s)(1 - \sigma(s)))$$

Observe that the  $\frac{\partial^{k+1} \sigma(s)}{\partial s^{k+1}}$  is also a polynomial in  $\sigma(s)$ . Observe that the degree  $k + 2$  term has one term with coefficient  $-a_{k+1} \cdot (k + 1)$ . Since  $a_{k+1} \neq 0$ , the coefficient of degree  $k + 2$ ,  $-a_{k+1} \cdot (k + 1)$ , is also non-zero. Since  $\frac{\partial^k \sigma(s)}{\partial s^k}$  is a polynomial in  $\sigma(s)$  with degree  $k + 1$  and hence, it cannot be zero identically.  $\square$

**Lemma 3.** Let  $x \in \mathbb{R}^n$  and  $A \in \mathbb{R}^{n \times n}$ . Suppose  $Ax = 0, \forall x \in \mathcal{X}$ , where  $\mathcal{X}$  has a non-empty interior. Under these conditions  $A = 0$ .

*Proof.* Since  $\mathcal{X}$  has a non-empty interior, we can construct a  $\ell_\infty$  ball centered on  $\theta$ , defined as follows  $\tilde{\mathcal{X}} = \{\theta + \sum_{j=1}^n \alpha_j e_j \mid \|\alpha\|_\infty \leq \alpha_{\max}\}$ , where  $e_j$  is a vector that is zero in all components and one on the  $j^{\text{th}}$  component. Suppose  $A$  was non-zero. One of the columns say  $a_j$  is non-zero. Consider two points in the ball  $\tilde{\mathcal{X}}$  such that  $j^{\text{th}}$  coefficients are non-zero but rest of the coefficients are zero. We denote the  $j^{\text{th}}$  components for the two components as  $\alpha_j$  and  $\tilde{\alpha}_j$ , where  $\alpha_j \neq \tilde{\alpha}_j$ . We now plug these two points into the condition that  $Ax = 0$

$$\begin{aligned} A(\theta + \alpha_j e_j) = 0 &\implies A\theta = -\alpha_j a_j, \\ A(\theta + \tilde{\alpha}_j e_j) = 0 &\implies A\theta = -\tilde{\alpha}_j a_j, \end{aligned} \tag{98}$$

We take a difference of the two steps above and obtain

$$(\alpha_j - \tilde{\alpha}_j) a_j = 0 \implies a_j = 0$$

This is a contradiction. Hence,  $A = 0$ .  $\square$

**Theorem 5.** If  $\mathcal{H}$  follows Assumption 8, and the realizability condition holds, i.e.,  $f \in \mathcal{H}$  and regular closedness condition in Assumption 2 holds, then the model trained to minimize the risk in equation 1 with  $\ell_2$  loss (with  $T \geq 2$ ) achieves length and compositional generalization.

*Proof.* We start with the same steps as earlier proofs and equate the prediction of  $h$  and  $f$  everywhere in the support of the training distribution (using first part of Lemma 1). We start with equating label at length 1, i.e.,  $y_1$ . For all  $x_1 \in \text{supp}(X_1)$

$$\begin{aligned} \sigma(A\sigma(Bx_1)) = \sigma(\tilde{A}\sigma(\tilde{B}x_1)) &\implies A\sigma(Bx_1) = \tilde{A}\sigma(\tilde{B}x_1) \implies \\ \sigma(B\tilde{B}^{-1}\tilde{B}x_1) = A^{-1}\tilde{A}\sigma(\tilde{B}x_1) \end{aligned} \tag{99}$$

Say  $y = \tilde{B}x_1$ ,  $A^{-1}\tilde{A} = U$ ,  $B\tilde{B}^{-1} = V$ . We substitute these expressions in the simplification below. We pick a  $y$  in the interior of  $\tilde{B} \cdot \text{supp}(X_1)$ .

$$\sigma(Vy) = U\sigma(y) \tag{100}$$

Take the first row of  $V$  and  $U$  as  $v^\top$  and  $u^\top$  to obtain

$$\sigma(v^\top y) = u^\top \sigma(y) \tag{101}$$

Suppose there is some non-zero component of  $v$  say  $i$  but the corresponding component is zero in  $u$ .

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

$$\frac{\partial \sigma(v_i y_i + v_{-i} y_{-i})}{\partial y_i} = \sigma'(v_i y_i + v_{-i} y_{-i}) v_i = \frac{\partial u_{-i}^\top \sigma(y_{-i})}{\partial y_i} = 0 \quad (102)$$

From the above we get  $\sigma'(v^\top y) = 0$ . But sigmoid is strictly monotonic on  $\mathbb{R}$ ,  $\sigma'(x) > 0, \forall x \in \mathbb{R}$  and  $v^\top y \in \mathbb{R}$ . Hence,  $\sigma'(v^\top y) = 0$  is not possible. Similarly, suppose some component is non-zero in  $u$  and zero in  $v$ .

$$\frac{\partial \sigma(v_{-i}^\top y_{-i})}{\partial y_i} = 0 = \frac{\partial (u_i \sigma(y_i) + u_{-i}^\top \sigma(y_{-i}))}{\partial y_i} = u_i \sigma'(y_i) \quad (103)$$

Since the derivative of  $\sigma$  cannot be zero, the above condition cannot be true.

From the above, we can deduce that both  $u$  and  $v$  have same non-zero components.

Let us start with the case where  $p \geq 2$  components of  $u, v$  are non-zero. Below we equate the partial derivative w.r.t all components of  $y$  that have non-zero component in  $u$  (since  $y$  is in the interior of the image of  $\tilde{B}x_1$ , we can equate these derivatives).

$$\begin{aligned} \sigma(v^\top y) &= u^\top \sigma(y), \\ \frac{\partial^p \sigma(s)}{\partial s^p} \Big|_{s=v^\top y} \left( \prod_{v_i \neq 0} v_i \right) &= 0 \implies \frac{\partial^p \sigma(s)}{\partial s^p} = 0. \end{aligned} \quad (104)$$

Since support  $X_1$  has a non-empty interior, the set of values  $v^\top y$  takes also has a non-empty interior in  $\mathbb{R}$ . Hence, the above equality is true over a set of values  $s$ , which have a non-empty interior. Since  $\sigma(s)$  is analytic,  $\frac{\partial^p \sigma(s)}{\partial s^p}$  is also analytic. From (Mityagin, 2015), it follows that  $\frac{\partial^p \sigma(s)}{\partial s^p} = 0$  everywhere. From Lemma 2, we know this condition cannot be true.

We are left with the case where  $u$  and  $v$  have one non-zero component each.

$$\frac{1}{1 + e^{-vy}} = \frac{u}{1 + e^{-y}} \implies 1 + e^{-y} = u + u e^{-vy}$$

In the simplification above, we take derivative w.r.t  $y$  to obtain  $e^{-(v-1)y} = 1/uv$ . We now again take derivative again w.r.t  $y$  to get  $v = 1$  and substitute it back to get  $u = 1$ . Note that no other row of  $U$  or  $V$  can have same non-zero element because that would make matrix non invertible. From this we deduce that  $U$  and  $V$  are permutation matrices. From  $\sigma(Vy) = U\sigma(y)$  it follows that  $U = V = \Pi$ . Thus  $B = \Pi\tilde{B}$  and  $\tilde{A} = \Pi\Pi$ .

Next, we equate predictions for  $y_2$  to the ground truth (label  $y_2$  exists as  $T \geq 2$ ). For all  $x_1 \in \text{supp}(X_1)$

$$\begin{aligned} \sigma(A\sigma(\Lambda\sigma(Bx_1) + Bx_2)) &= \sigma(\tilde{A}\sigma(\tilde{\Lambda}\sigma(\tilde{B}x_1) + \tilde{B}x_2)) \implies \\ A\sigma(\Lambda\sigma(Bx_1) + Bx_2) &= \tilde{A}\sigma(\tilde{\Lambda}\sigma(\tilde{B}x_1) + \tilde{B}x_2) \implies \\ \tilde{A}\sigma(\tilde{\Lambda}\sigma(\tilde{B}x_1) + \tilde{B}x_2) &= \Pi\Pi\sigma(\tilde{\Lambda}\Pi^\top\sigma(Bx_1) + \Pi^\top Bx_2) = A\sigma(\Pi\tilde{\Lambda}\Pi^\top\sigma(Bx_1) + Bx_2). \end{aligned} \quad (105)$$

We use the simplification in the second step to equate to LHS in the first step as follows.

$$\begin{aligned} A\sigma(\Pi\tilde{\Lambda}\Pi^\top\sigma(Bx_1) + Bx_2) &= A\sigma(\Lambda\sigma(Bx_1) + Bx_2) \\ \implies (\Pi\tilde{\Lambda}\Pi^\top - \Lambda)\sigma(Bx_1) &= 0. \end{aligned} \quad (106)$$

Since  $\sigma(Bx_1)$  spans a set that has a non-empty interior, we get that  $\tilde{\Lambda} = \Pi^\top \Lambda \Pi$  (from Lemma 3).

From the above conditions, we have arrived at  $\tilde{\Lambda} = \Pi^\top \Lambda \Pi$ ,  $\tilde{B} = \Pi^\top B$ ,  $\tilde{A} = \Pi\Pi$ .

We want to show that for all  $k \geq 1$

$$h_k = \Pi \tilde{h}_k, \quad (107)$$

where  $h_k = \sigma(\Lambda h_{k-1} + Bx_k)$  and  $\tilde{h}_k = \sigma(\tilde{\Lambda} \tilde{h}_{k-1} + \tilde{B}x_k)$  and  $h_0 = \tilde{h}_0 = 0$ . In other words, we define  $T_k$  as a mapping that takes  $x_{\leq k}$  as input and outputs  $h_k$ , i.e.,  $T_k(x_{\leq k}) = h_k$ . Similarly, we write  $\tilde{T}_k(x_{\leq k}) = \tilde{h}_k$ . We want to show

$$T_k = \Pi \tilde{T}_k, \forall k \quad (108)$$

We show the above by principle of induction. Let us consider the base case below. For all  $x_1 \in \mathbb{R}^n$

$$\tilde{A}\sigma(\tilde{B}x_1) = A\Pi\sigma(\Pi^\top Bx_1) = A\sigma(Bx_1) = Ah_1 \implies h_1 = \Pi\tilde{h}_1 \implies T_1(x_1) = \Pi\tilde{T}_1(x_1) \quad (109)$$

Suppose  $\forall j \leq k, T_j = \Pi\tilde{T}_j$ .

Having shown the base case and assumed the condition for  $j \leq k$ , we now consider the mapping  $\tilde{T}_{k+1}$

$$\Pi\tilde{T}_{k+1}(x_{\leq k+1}) = \Pi\sigma(\tilde{\Lambda}\tilde{h}_k + \tilde{B}x_{k+1}) = \Pi\sigma(\Pi^\top \Lambda\Pi\tilde{h}_k + \Pi^\top Bx_k) = \sigma(\Lambda h_k + Bx_k) = T_{k+1}(x_{\leq k+1}). \quad (110)$$

The prediction from the model  $(\tilde{A}, \tilde{\Lambda}, \tilde{B})$  at a time step  $k$  is denoted as  $\tilde{y}_k$  and it relates to  $\tilde{h}_k$  as follows  $\tilde{y}_k = \sigma(\tilde{A}\tilde{h}_k)$ . We use the above condition in equation equation 126 to arrive at the following result. For all  $x_{\leq k} \in \mathbb{R}^{nk}$

$$\tilde{y}_k = \sigma(\tilde{A}\tilde{h}_k) = \sigma(\tilde{A}\tilde{T}(x_{\leq k})) = \sigma(A\Pi\tilde{T}(x_{\leq k})) = \sigma(AT(x_{\leq k})) = y_k$$

This completes the proof.  $\square$

**Corollary 5.** *If  $\mathcal{H}$  follows Assumption 8, and the realizability condition holds, i.e.,  $f \in \mathcal{H}$  and regular closedness condition in Assumption 2 holds, then the model trained to minimize the risk in equation 1 with  $\ell_2$  loss (with  $T \geq 2$ ) achieves permutation identification. Further, under the stated conditions permutation identification is necessary for both length and compositional generalization.*

*Proof.* We follow the exact steps from the proof of Theorem 5 up to equation 128. From equation 128 it follows that for all  $x_{\leq k} \in \mathbb{R}^{nk}$

$$T_k(x_{\leq k}) = \Pi\tilde{T}(x_{\leq k}) \implies h_k = \Pi\tilde{h}_k \quad (111)$$

The above implies permutation identification. To show the necessity of permutation identification, from the proof of Theorem 5 observe that

$$\forall i \leq T, \forall x_{\leq i} \in \text{supp}(X_{\leq i}), \sigma(A\sigma(Bx_i + \Lambda h_{i-1})) = \sigma(\tilde{A}\sigma(\tilde{B}x_i + \tilde{\Lambda} \tilde{h}_{i-1})) \implies T_k(x_{\leq k}) = \Pi\tilde{T}(x_{\leq k}) \quad (112)$$

The latter condition implies permutation identification. If it does not hold, then the condition in LHS cannot hold and hence neither length nor compositional generalization can be achieved.  $\square$

#### C.4.1 EXTENDING THEOREM 5 TO DISCRETE TOKENS

In our discussion, we have focused on settings where the support of each token has a non-empty interior (Assumption 2). In practice of language modeling, we use discrete tokens and hence Assumption 2 does not hold anymore. In this section, we discuss the adaptation of results for vanilla RNNs to setting when the the support of tokens is a finite set.

Define  $\mathcal{S} = \{y = Bx \mid x \in \mathcal{X}\}$ , where  $\mathcal{X}$  is the marginal support of each token.

**Assumption 18.** a) For each component  $i$  of  $y$ ,  $\mathcal{S}$  contains two pairs where the first coordinate differs by the same amount. Mathematically stated, the two pairs are  $\left((y_i, y_{-i}), (y_i + \delta, y_{-i})\right)$  and  $\left((y'_i, y'_{-i}), (y'_i + \delta, y'_{-i})\right)$ .

b) For every pair of components  $i, j$  of  $y$ ,  $\mathcal{S}$  contains a point  $y$  that satisfies the following. There exists three points in  $\mathcal{S}$  such that they only differ in  $y_i, y_j$ , and form a rectangle,  $(y_i, y_j), (y'_i, y_j), (y_i, y'_j), (y'_i, y'_j)$ . Similarly, there exists another set of points where  $y'_i < y_i$  and  $y'_j < y_j$ .

**Theorem 13.** If  $\mathcal{H}$  follows Assumption 8, and the realizability condition holds, i.e.,  $f \in \mathcal{H}$  and regular closedness condition in Assumption 2 holds, then the model trained to minimize the risk in equation 1 with  $\ell_2$  loss (with  $T \geq 2$ ) achieves length and compositional generalization.

*Proof.* We start with the same steps as earlier proofs and equate the prediction of  $h$  and  $f$  everywhere in the support of the training distribution. We start with equating label at length 1, i.e.,  $y_1$ . For all  $x_1 \in \text{supp}(X_1)$

$$\begin{aligned} \sigma(A\sigma(Bx_1)) &= \sigma(\tilde{A}\sigma(\tilde{B}x_1)) \implies A\sigma(Bx_1) = \tilde{A}\sigma(\tilde{B}x_1) \implies \\ \tilde{A}^{-1}A\sigma(Bx_1) &= \sigma(\tilde{B}B^{-1}Bx_1) \end{aligned} \quad (113)$$

Say  $y = Bx_1$ ,  $\tilde{A}^{-1}A = U$ ,  $\tilde{B}B^{-1} = V$ . We substitute these expressions in the simplification below. We pick a  $y$  in the interior of  $\tilde{B} \cdot \text{supp}(X_1)$ .

$$\sigma(Vy) = U\sigma(y) \quad (114)$$

Take the first row of  $V$  and  $U$  as  $v^\top$  and  $u^\top$  to obtain

$$\sigma(v^\top y) = u^\top \sigma(y) \quad (115)$$

Say  $v_i \neq 0$  and  $u_i = 0$ . We consider a  $(y_i, y_{-i})$  and  $(y'_i, y_{-i})$  satisfying Assumption 18 a. We substitute these points in equation 115 and take the difference of the LHS and RHS in equation 115 to obtain.

$$\sigma(v_i y'_i + v_{-i} y_{-i}) - \sigma(v_i y_i + v_{-i} y_{-i}) = 0 \quad (116)$$

$\sigma$  is strictly monotonic and thus the above cannot be true. Similarly, we can rule out the case when  $u_i \neq 0$  and  $v_i = 0$ . Thus we can deduce that both  $u$  and  $v$  have same non-zero components.

Let us start with the case where  $p \geq 2$  components of  $u, v$  are non-zero. Without loss of generality say the first two components are among coordinates that are non-zero. Pick a  $y \in \mathcal{S}$  that satisfies Assumption 18 b. Suppose  $v^\top y \geq 0$ . We select the neighbors of  $y$  that form the rectangle such that each coordinate is greater than  $y$ . We substitute these points in equation 115 and the simplification procedure works as follows. Let

$$\begin{aligned} s_1 &= v_1 y'_1 + v_2 y'_2 + \cdots + v_n y_n, & s_3 &= v_1 y'_1 + v_2 y_2 + \cdots + v_n y_n \\ s_2 &= v_1 y_1 + v_2 y'_2 + \cdots + v_n y_n, & s_4 &= v_1 y_1 + v_2 y_2 + \cdots + v_n y_n \end{aligned} \quad (117)$$

Observe that  $s_1 > s_2 > s_4$  and  $s_1 > s_3 > s_4$ . It is possible that  $s_2 \geq s_3$  or  $s_3 > s_2$ . Suppose  $s_2 \geq s_3$ .

We can write

$$\begin{aligned} \sigma(s_1) &= u_1 \sigma(y'_1) + u_2 \sigma(y'_2) + \cdots + u_n \sigma(y_n), & \sigma(s_2) &= u_1 \sigma(y_1) + u_2 \sigma(y'_2) + \cdots + u_n \sigma(y_n) \\ \sigma(s_3) &= u_1 \sigma(y'_1) + u_2 \sigma(y_2) + \cdots + u_n \sigma(y_n), & \sigma(s_4) &= u_1 \sigma(y_1) + u_2 \sigma(y_2) + \cdots + u_n \sigma(y_n) \end{aligned} \quad (118)$$

We take a difference of the first two and the latter two, and subtract these differences to get

2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159

$$\left(\sigma(s_1) - \sigma(s_2)\right) - \left(\sigma(s_3) - \sigma(s_4)\right) = 0 \quad (119)$$

From mean value theorem, we get that  $\sigma'(\tilde{s}) = \sigma'(s^\dagger)$ , where  $\sigma'$  is the derivative of  $\sigma$ ,  $\tilde{s}$  is a value between  $s_1$  and  $s_2$ , and  $s^\dagger$  is a value between  $s_3$  and  $s_4$ . Since  $s_1 > s_2 > s_3 > s_4 > 0$ ,  $\tilde{s} > s^\dagger > 0$ . Since  $\sigma'$  strictly decreases on positive values, the above equality  $\sigma'(\tilde{s}) = \sigma'(s^\dagger)$  is not possible. Similarly, we can tackle the case  $v^\top y < 0$ .

We are left with the case where  $u$  and  $v$  have one non-zero component each. From Assumption 18a, we select two pairs that differ exactly in the non-zero component. We can resort to dealing with scalars as follows. We start with first pair  $(y, y + \delta)$ .

$$\begin{aligned} \sigma(vy) = u\sigma(y) &\implies \frac{1}{1 + e^{-vy}} = \frac{u}{1 + e^{-y}} \implies 1 - u = ue^{-vy} - e^{-y} \\ \sigma(v(y + \delta)) = u\sigma(y + \delta) &\implies 1 - u = ue^{-v(y+\delta)} - e^{-(y+\delta)} \end{aligned} \quad (120)$$

By equating the RHS in the above, we obtain

$$\frac{1 - e^{-\delta}}{1 - e^{-v\delta}} = ue^{-(v-1)y} \quad (121)$$

For the second pair  $(y', y' + \delta)$ , we obtain

$$\frac{1 - e^{-\delta}}{1 - e^{-v\delta}} = ue^{-(v-1)y'} \quad (122)$$

If we compare the RHS of equation 121 and equation 122, we obtain  $ue^{-(v-1)y} = ue^{-(v-1)y'}$ . Since  $u$  is non-zero, we obtain that  $v = 1$ . Substituting this into  $\sigma(vy) = u\sigma(y)$ , we also obtain  $u = 1$ .

Note that no other row of  $U$  or  $V$  can have same non-zero element because that would make matrix non invertible. From this we deduce that  $U$  and  $V$  are permutation matrices. From  $\sigma(Vy) = U\sigma(y)$  it follows that  $U = V = \Pi$ . Thus  $B = \Pi\tilde{B}$  and  $\tilde{A} = A\Pi$ .

Next, we equate predictions for  $y_2$  to the ground truth (label  $y_2$  exists as  $T \geq 2$ ). For all  $x_1 \in \text{supp}(X_1)$

$$\begin{aligned} \sigma(A\sigma(\Lambda\sigma(Bx_1) + Bx_2)) &= \sigma(\tilde{A}\sigma(\tilde{\Lambda}\sigma(\tilde{B}x_1) + \tilde{B}x_2)) \implies \\ A\sigma(\Lambda\sigma(Bx_1) + Bx_2) &= \tilde{A}\sigma(\tilde{\Lambda}\sigma(\tilde{B}x_1) + \tilde{B}x_2) \implies \\ \tilde{A}\sigma(\tilde{\Lambda}\sigma(\tilde{B}x_1) + \tilde{B}x_2) &= A\Pi\sigma(\tilde{\Lambda}\Pi^\top\sigma(Bx_1) + \Pi^\top Bx_2) = A\sigma(\Pi\tilde{\Lambda}\Pi^\top\sigma(Bx_1) + Bx_2). \end{aligned} \quad (123)$$

We use the simplification in the second step to equate to LHS in the first step as follows.

$$\begin{aligned} A\sigma(\Pi\tilde{\Lambda}\Pi^\top\sigma(Bx_1) + Bx_2) &= A\sigma(\Lambda\sigma(Bx_1) + Bx_2) \\ \implies (\Pi\tilde{\Lambda}\Pi^\top - \Lambda)\sigma(Bx_1) &= 0. \end{aligned} \quad (124)$$

Since  $\sigma(Bx_1)$  spans a set that has a non-empty interior, we get that  $\tilde{\Lambda} = \Pi^\top\Lambda\Pi$  (from Lemma 3).

From the above conditions, we have arrived at  $\tilde{\Lambda} = \Pi^\top\Lambda\Pi$ ,  $\tilde{B} = \Pi^\top B$ ,  $\tilde{A} = A\Pi$ .

We want to show that for all  $k \geq 1$

$$h_k = \Pi\tilde{h}_k, \quad (125)$$



where  $h_k = \sigma(\Lambda h_{k-1} + Bx_k)$  and  $\tilde{h}_k = \sigma(\tilde{\Lambda}\tilde{h}_{k-1} + \tilde{B}x_k)$  and  $h_0 = \tilde{h}_0 = 0$ . In other words, we define  $T_k$  as a mapping that takes  $x_{\leq k}$  as input and outputs  $h_k$ , i.e.,  $T_k(x_{\leq k}) = h_k$ . Similarly, we write  $\tilde{T}_k(x_{\leq k}) = \tilde{h}_k$ . We want to show

$$T_k = \Pi\tilde{T}_k, \forall k \quad (126)$$

We show the above by principle of induction. Let us consider the base case below. For all  $x_1 \in \mathbb{R}^n$

$$\tilde{A}\sigma(\tilde{B}x_1) = A\Pi\sigma(\Pi^\top Bx_1) = A\sigma(Bx_1) = Ah_1 \implies h_1 = \Pi\tilde{h}_1 \implies T_1(x_1) = \Pi\tilde{T}_1(x_1) \quad (127)$$

Suppose  $\forall j \leq k, T_j = \Pi\tilde{T}_j$ .

Having shown the base case and assumed the condition for  $j \leq k$ , we now consider the mapping  $\tilde{T}_{k+1}$

$$\Pi\tilde{T}_{k+1}(x_{\leq k+1}) = \Pi\sigma(\tilde{\Lambda}\tilde{h}_k + \tilde{B}x_{k+1}) = \Pi\sigma(\Pi^\top \Lambda\Pi\tilde{h}_k + \Pi^\top Bx_k) = \sigma(\Lambda h_k + Bx_k) = T_{k+1}(x_{\leq k+1}). \quad (128)$$

The prediction from the model  $(\tilde{A}, \tilde{\Lambda}, \tilde{B})$  at a time step  $k$  is denoted as  $\tilde{y}_k$  and it relates to  $\tilde{h}_k$  as follows  $\tilde{y}_k = \sigma(\tilde{A}\tilde{h}_k)$ . We use the above condition in equation 126 to arrive at the following result. For all  $x_{\leq k} \in \mathcal{X}^k$

$$\tilde{y}_k = \sigma(\tilde{A}\tilde{h}_k) = \sigma(\tilde{A}\tilde{T}(x_{\leq k})) = \sigma(A\Pi\tilde{T}(x_{\leq k})) = \sigma(AT(x_{\leq k})) = y_k$$

This completes the proof.  $\square$

## C.5 FINITE HYPOTHESIS CLASS

Before we present the proof of Theorem 6, we revisit some basics of convergence of sets. Consider a sequence of sets  $(A_n)$  which are a subset of  $\Omega$ , i.e.,  $A_n \subseteq \Omega$ . We define the lim inf first and then the lim sup.

$$\liminf_{n \rightarrow \infty} A_n = \bigcup_{n \geq 1} \bigcap_{j \geq n} A_j$$

$$\limsup_{n \rightarrow \infty} A_n = \bigcap_{n \geq 1} \bigcup_{j \geq n} A_j$$

The limit of this sequence of sets exists provided the lim inf and lim sup are equal, i.e.,

$$\lim_{n \rightarrow \infty} A_n = \bigcup_{n \geq 1} \bigcap_{j \geq n} A_j = \bigcap_{n \geq 1} \bigcup_{j \geq n} A_j$$

If the sequence comprises of non-increasing sets, i.e.,  $A_{n+1} \subseteq A_n$ , then the limit exists. For this non-increasing sequence observe that

$$\bigcap_{j \geq n} A_j = \bigcap_{j \geq 1} A_j$$

$$\bigcup_{j \geq n} A_j = A_n$$

We combine the above two observations to see both lim inf and lim sup are equal and thus the limit of non-increasing sets exists. There is another way to define the limit of sets using indicator functions that goes as follows.  $\mathbf{1}_A(\cdot)$  is the indicator function that checks if input belongs to the set

or not and takes the value of one if the input is in the set and zero otherwise. We define the limit using indicator functions as follows.

$$\lim_{n \rightarrow \infty} A_n = \{\omega \in \Omega, \lim_{n \rightarrow \infty} \mathbf{1}_{A_n}(\omega) = 1\},$$

where  $\mathbf{1}_{A_n}(\omega)$  is one if  $\omega \in A_n$  and zero otherwise. The limit of sequence of sets  $A_n$  exists if and only if  $\lim_{n \rightarrow \infty} \mathbf{1}_{A_n}(\omega)$  exists for all  $\omega \in \Omega$ .

**Theorem 6.** *If  $\mathcal{H}$  is a finite hypothesis class, the realizability condition holds, i.e.,  $f \in \mathcal{H}$ , then  $\exists T_0 < \infty$  such that the model trained to minimize the risk in equation 1 with  $\ell_2$  loss and  $T > T_0$  achieves length generalization.*

*Proof.* Let  $\mathcal{H}_T$  be the set of solutions to equation 1, where  $T$  is the maximum length of the sequence in the training distribution. Observe that each  $\mathcal{H}_T$  can take one of the possible values in the power set  $2^{\mathcal{H}}$ , i.e., the set of all the possible subsets of  $\mathcal{H}$ . Since the objective at length  $T$  in equation 1 evaluates the model at all lengths up to length  $T$  we obtain that  $\mathcal{H}_{T+1} \subseteq \mathcal{H}_T$ . Since the sequence  $\mathcal{H}_T$  indexed by  $T$  is non-increasing, the limit of the above sequence exists and is denoted as  $\mathcal{H}^*$ . From the indicator function definition of the limit, we can write  $\mathcal{H}^*$  as

$$\mathcal{H}^* = \{h \in \mathcal{H}, \lim_{T \rightarrow \infty} \mathbf{1}_{\mathcal{H}_T}(h) = 1\}$$

Since the limit of the sequence  $\mathcal{H}_T$  exists, for each  $h \in \mathcal{H}$ , the limit  $\lim_{T \rightarrow \infty} \mathbf{1}_{\mathcal{H}_T}(h)$  exists denoted as  $p(h)$ . Each element of this sequence  $\mathbf{1}_{\mathcal{H}_T}(h)$  indexed by  $T$  takes a value of one or zero. From the standard definition of limit, we know that for each  $\epsilon$ , there exists  $T(h, \epsilon)$  such that  $T > T(h, \epsilon)$ ,  $|\mathbf{1}_{\mathcal{H}_T}(h) - p(h)| < \epsilon$ . Both  $\mathbf{1}_{\mathcal{H}_T}(h)$  and  $p(h)$  can only take a value of 0 or 1 (for  $p(h)$  if there is any other value it takes, then the distance of sequence terms  $\mathbf{1}_{\mathcal{H}_T}(h)$  from  $p(h)$  will be bounded away from zero, which is not possible). If  $\epsilon < 1$ , then for all  $T > T(h, \epsilon)$ ,  $\mathbf{1}_{\mathcal{H}_T}(h) = p(h)$ .

Define  $T_0 = \sup_{h \in \mathcal{H}} T(h, \epsilon)$ . Since  $\mathcal{H}$  is finite,  $T_0 < \infty$ .

We can write the set  $\mathcal{H}_T$  as

$$\mathcal{H}_T = \{h \in \mathcal{H}, \mathbf{1}_{\mathcal{H}_T}(h) = 1\}$$

If  $T > T_0$ , then

$$\mathcal{H}_T = \{h \in \mathcal{H}, p(h) = 1\} = \{h \in \mathcal{H}, \lim_{T \rightarrow \infty} \mathbf{1}_{\mathcal{H}_T}(h) = 1\} = \mathcal{H}^*$$

We now argue that  $\mathcal{H}^*$  contains all length generalizing solutions. Since  $f \in \mathcal{H}_t$  for all  $t \geq 1$ ,  $f \in \mathcal{H}^*$ . Now let us suppose that there is a  $g \in \mathcal{H}^*$ , which does not length generalize. In other words, this  $g$  leads to a non-zero error for some finite length  $\tilde{T}$ . Thus  $g$  cannot be in  $\mathcal{H}_{\tilde{T}}$ . From the definition of limit, it follows that  $\mathcal{H}^* \subseteq \mathcal{H}_t$  for all  $t$ . This leads to contradiction. Hence, such a  $g$  cannot exist. Thus all the solutions in  $\mathcal{H}^*$  the set length generalize, which proves the claim.  $\square$

## D EXPERIMENTS

Here we provide additional experimental results as well as the training details.

**Model Architecture** In all the architectures, there are two types of non-linearities,  $\omega$  that generates the target label,  $\psi$  that operates on inputs (used in deep sets and transformers). We use MLPs to implement these non-linearities. We instantiate MLPs with  $l$  hidden layers, and the input, output, and hidden dimensions are all the same  $m = n = k$ . Recall that under the realizability assumption  $f \in \mathcal{H}$ . Therefore, we need to select the labeling function from  $\mathcal{H}$ . To do so, the weights of MLP

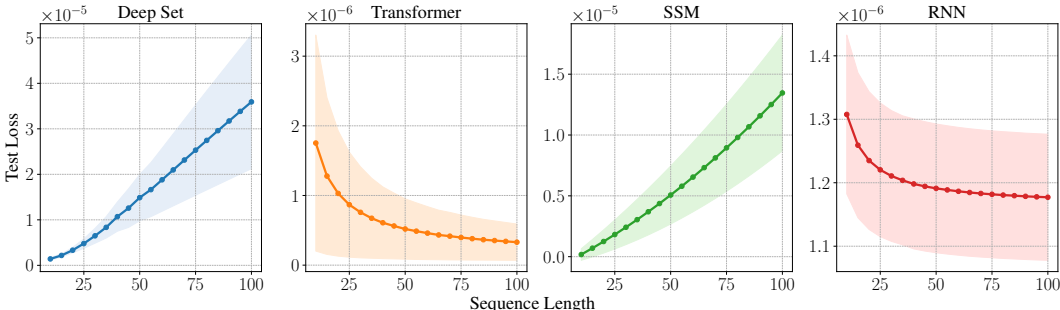


Figure 5: Length generalization: Test  $\ell_2$  loss on sequences of different lengths. The models are trained only on sequences of length up to  $T = 10$ . All models achieve small error values  $\approx 10^{-5} - 10^{-6}$  at all sequence lengths and thus length generalize. Since the error values are already quite small, the increasing or decreasing trends are not numerically significant.

are initialized according to  $\mathcal{N}(\mu, \sigma^2)$ , where  $\mu = 0.0, \sigma = 0.6$ . For RNNs and SSMs,  $A, B, \Lambda$  are initialized separately for the learner and true generating process as orthogonal matrices. All hidden layers, as well as the output layer are followed by a sigmoidal activation function.

**Training Details and Hyperparameter Selection** We train all models with AdamW optimizer (Loshchilov & Hutter, 2019) with a learning rate of  $10^{-3}$ , weight decay of 0.01,  $\epsilon = 10^{-8}, \beta_1 = 0.9, \beta_2 = 0.95$ . We reduce the learning rate by a factor of 0.8 if the validation loss is not improved more than  $10^{-6}$  for 1 epoch. This drop is followed by a cool-down period of 1 epoch, and the learning rate cannot decrease to lower than  $10^{-7}$ . For all datasets we use a streaming dataset where each epoch contains 100 batches of size 256 sampled online from the specified training and test distributions, and we train all models for 100 epochs. Therefore, the size of the training dataset is  $256 \times 10^4$  and the size of the testing dataset is  $256 \times 10^2$ . Since our models are generally small, running the experiments is rather inexpensive, and we carried out each experiment on 4 CPU cores using 20 GB of RAM. For inference, specially for SSM and RNN with very long sequences, we use RTX8000 GPUs.

#### D.1 LENGTH GENERALIZATION

In Figure 5, we present additional findings for length generalization capability of all architectures when both the learner and the generating process MLPs all consist of one hidden layers with input, output, and hidden size matching  $n = m = k = 20$ .

To complement Figure 3, in Figures 6, 7 we present the prediction behaviour of SSM and RNN architectures with two hidden layer MLPs for  $\omega$  trained on sequences output by two hidden layer MLPs for  $\rho$ .

Figures 8, 9, 10, 11 present the prediction behaviour of deep set, Transformer with softmax attention, SSM, and RNN architectures with one hidden layer in  $\rho$  (and one hidden layer MLPs for the learner  $\omega$ ). Training procedure remains the same. We can observe that all models length generalize.

Additionally, to support the theory on other types of attention, Figures 12, 13 demonstrate the loss and prediction of a Transformer with ReLU attention and one hidden layer MLPs for  $\omega, \psi$  trained on output sequences of a Transformer with ReLU attention and one hidden layer MLP for  $\rho, \phi$ . Similarly, all these models were trained to predict sequences of length up to  $T = 10$  output by a true labeling function  $f$  in their respective hypothesis classes  $\mathcal{H}$ , and were tested with sequences of length up to 100. As a reminder, the output tokens  $y_i \in \mathbb{R}^m$ , where  $m = 20$ , and the figures below show only one representative dimension for illustration. All models demonstrate strong length generalization capacity.

**Discrete Tokens** In Table 3 we present the results for successful length generalization of the different architectures when the inputs are discrete. We sample all components from  $[0, 1]$  interval and discretize the values to one of the 5 levels in  $[0.0, 0.2, 0.4, 0.6, 0.8]$ . Note that the small scale of values of loss at longer lengths indicate successful generalization. For a visual depiction of results,

Model	Test Loss $\times 10^6$ ( $t = 10$ )	Test Loss $\times 10^6$ ( $t = 90$ )
Deep set	$3.48 \pm 0.15$	$52.7 \pm 0.88$
Transformer	$1.72 \pm 0.27$	$48.8 \pm 2.44$
SSM	$0.2 \pm 0.0$	$4.06 \pm 0.0$
RNN	$0.22 \pm 0.0$	$1.3 \pm 0.0$

Table 3: Length generalization of different architectures when the input tokens are discrete. Models are trained in sequences of length up to  $T = 10$  and show successful generalization on much longer sequences.

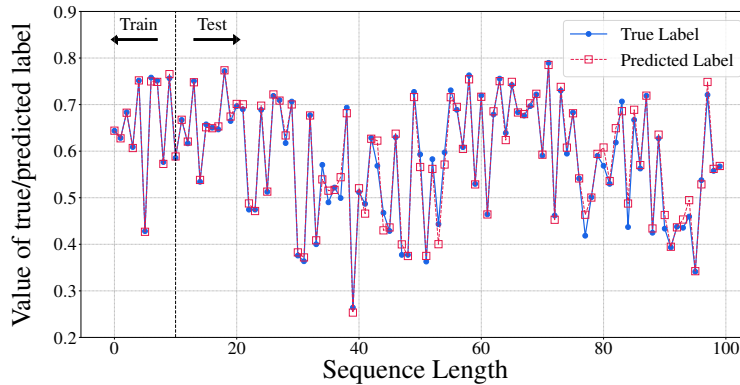


Figure 6: A SSM model with *two* hidden layer MLP for  $\omega$  trained on sequences of length up to  $T = 10$  length generalizes to sequences of length up to 100.

please see Fig. 14. Also note that in all architectures  $\rho$  in  $f$  and  $\omega$  in  $h$  are comprised of two hidden layers.

## D.2 COMPOSITIONAL GENERALIZATION

Here we present the prediction behavior of different architectures on the test sequences that consist of unseen token combinations during training. This helps us better interpret qualitatively how the model actually performs in following the true labels. Figures 16-19 show the prediction trajectories for different architectures. We can observe that not only do these models perform quite well on unseen sequences of length up to  $T = 10$ , but they also length generalize and continue to remain consistent with the true labels on unseen combinations at longer lengths than the training. Table 4 presents the test loss and  $R^2$  on the test set when the model is only trained on the red region in

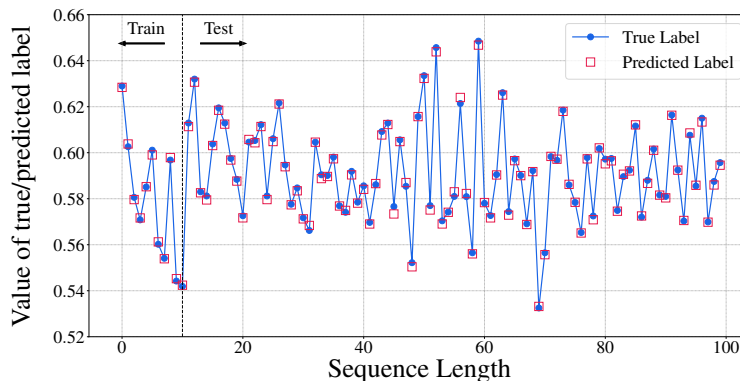


Figure 7: A RNN model with *two* hidden layer MLP for  $\omega$  trained on sequences of length up to  $T = 10$  length generalizes to sequences of length up to 100.

2376  
 2377  
 2378  
 2379  
 2380  
 2381  
 2382  
 2383  
 2384  
 2385  
 2386  
 2387  
 2388  
 2389  
 2390  
 2391  
 2392  
 2393  
 2394  
 2395  
 2396  
 2397  
 2398  
 2399  
 2400  
 2401  
 2402  
 2403  
 2404  
 2405  
 2406  
 2407  
 2408  
 2409  
 2410  
 2411  
 2412  
 2413  
 2414  
 2415  
 2416  
 2417  
 2418  
 2419  
 2420  
 2421  
 2422  
 2423  
 2424  
 2425  
 2426  
 2427  
 2428  
 2429

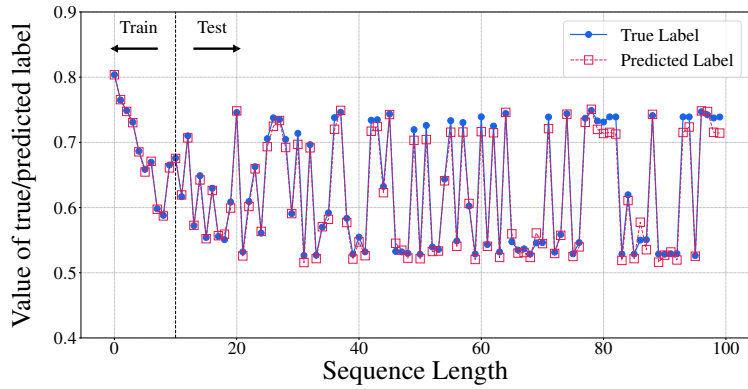


Figure 8: A deep set model with one hidden layer MLP for  $\psi, \omega$  trained on sequences of length up to  $T = 10$  shows perfect generalization to sequences of length up to 100.

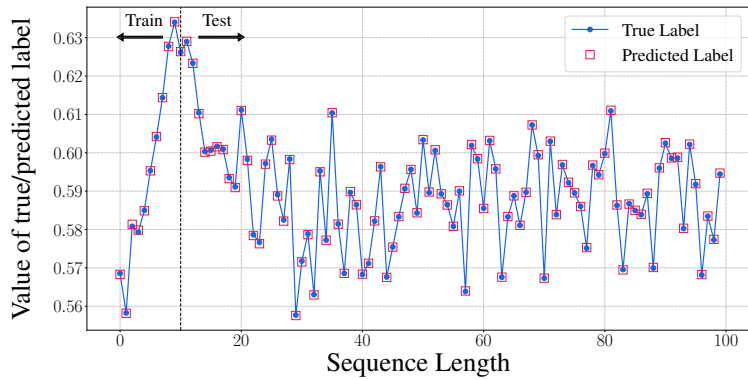


Figure 9: A Transformer model with softmax attention and one hidden layer MLP trained on sequences of length up to  $T = 10$  shows perfect generalization to sequences of length up to 100.

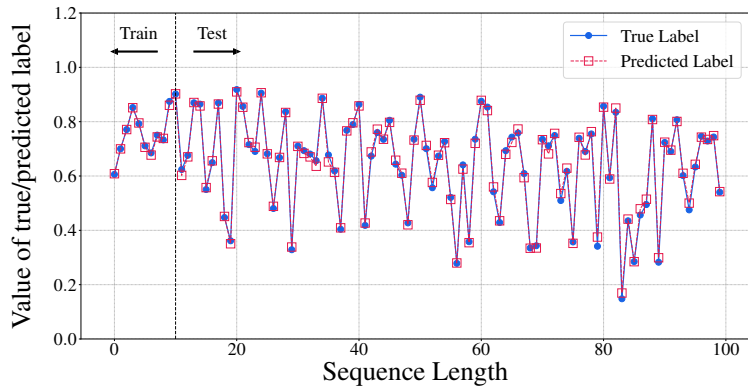


Figure 10: A SSM model with one hidden layer MLP for  $\omega$  trained on sequences of length up to  $T = 10$  length generalizes to sequences of length up to 100.

2430  
2431  
2432  
2433  
2434  
2435  
2436  
2437  
2438  
2439  
2440  
2441  
2442

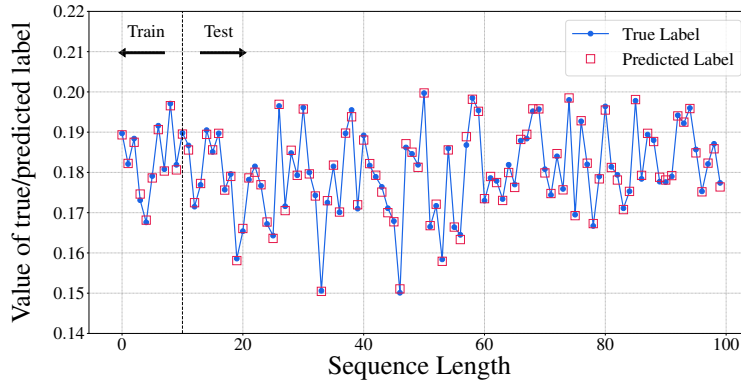


Figure 11: A RNN model with one hidden layer MLP for  $\omega$  trained on sequences of length up to  $T = 10$  length generalizes to sequences of length up to 100.

2446  
2447  
2448  
2449  
2450  
2451  
2452  
2453  
2454  
2455  
2456  
2457  
2458  
2459  
2460  
2461  
2462  
2463  
2464

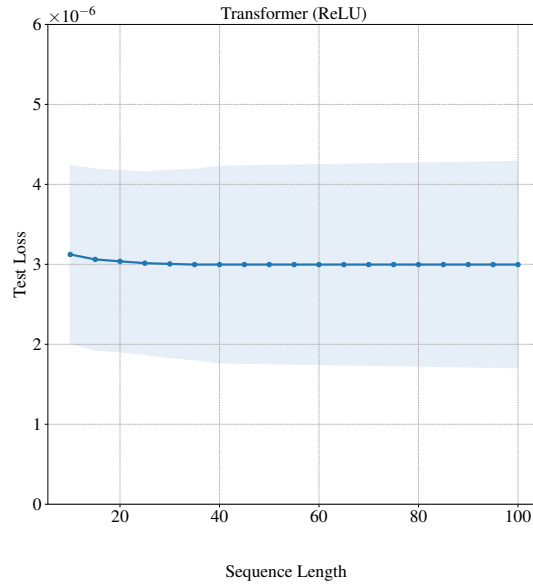


Figure 12: Test loss of a transformer model with ReLU attention and one hidden layer MLP for  $\omega, \psi$  trained on sequences of length up to  $T = 10$  length generalizes to sequences of length up to 100. The results are averaged over five seeds.

2465  
2466  
2467  
2468  
2469  
2470  
2471  
2472  
2473  
2474  
2475  
2476  
2477  
2478  
2479  
2480  
2481

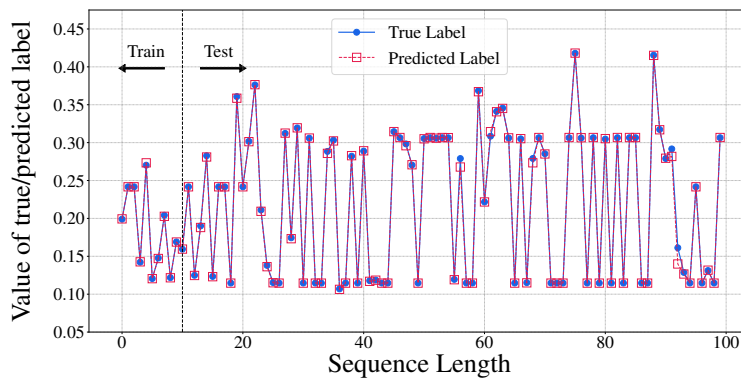


Figure 13: A transformer model with ReLU attention and one hidden layer MLP for  $\omega, \psi$  trained on sequences of length up to  $T = 10$  length generalizes to sequences of length up to 100.

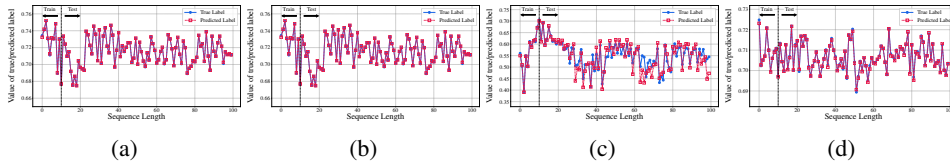


Figure 14: Successful length generalization of different architectures (with 2 hidden layers for  $\rho$ ) when the input is discrete. From the left: Deep set, Transformer, SSM, RNN.

Model	Test Loss $\times 10^6$	$R^2$
Deep set	$1.27 \pm 0.24$	$0.96 \pm 0.01$
Transformer	$4.50 \pm 3.28$	$1.00 \pm 0.00$
SSM	$11.00 \pm 10.92$	$1.00 \pm 0.00$
RNN	$1.22 \pm 0.12$	$0.99 \pm 0.00$

Table 4: Compositional generalization: Test  $\ell_2$  loss and  $R^2$  score for models with one hidden layers on sequences of length  $T = 10$ . A strong linear relationship is observed for all models for new sequences made of unseen token combinations.

Figure 15. All models generalize to unseen combination of tokens and the learned representations linearly identify the true hidden representations.

Figures 20, 21, 22, 23 present the prediction behaviour of deep set, transformer with softmax attention, SSM, and RNN architectures with two hidden layers in  $\rho$  (and two hidden layer MLPs for the learner  $\omega$ ) when trained on sequences of length up to  $T = 10$  sampled from the red region in Figure 15. We can observe that all models continue to generalize to unseen combinations beyond their training length.

**Discrete Tokens** Evaluating compositional generalization with discrete tokens introduces additional challenges. This is because we have to sample the training and test distribution according to Fig. 1 (and 15). There are multiple ways to achieve this but they become infeasible with long sequences of interest in practice:

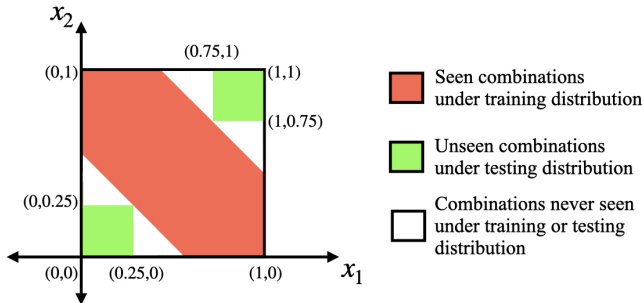
- We could continuously sample from the regions in Fig. 1 and 15 and then round up or down the components to one of the predefined set of values. However, with longer sequences this translates to sampling and then rounding values in a high-dimensional hyper-diamond where points are increasingly spread out toward the boundaries. Rounding up results in corners becoming part of the training samples, corrupting the test set. Rounding down will result in a training set in which the support of tokens no longer follows Fig. 1 (i.e., does not cover the discrete set of values predefined in  $[0, 1]$ ).
- We could instead sample continuously and then discretize based on finding the nearest neighbour of each point to the points in a discrete grid of values in  $\mathbb{R}^T$ . Having as few as 5 discrete levels renders this sampling procedure impossible for long sequences due to the complexity of finding nearest neighbours.
- Lastly, one could construct the set of discrete points in  $\mathbb{R}^T$  that satisfy the constraints in Fig. 1 and then sample from this set, however, this enumeration also proves infeasible as the search space grows exponentially.

Therefore, evaluating compositional generalization in the discrete case is not straightforward beyond very short sequences.

**Practical Considerations** For training and evaluating compositional aspect of generalization, we follow the sampling procedure described in Figure 1 with a slight modification that allows for faster sampling and easier training. This procedure is illustrated in Figure 15, and results in a more difficult testing strategy, as the test set spans a smaller area than the complement of the training set.

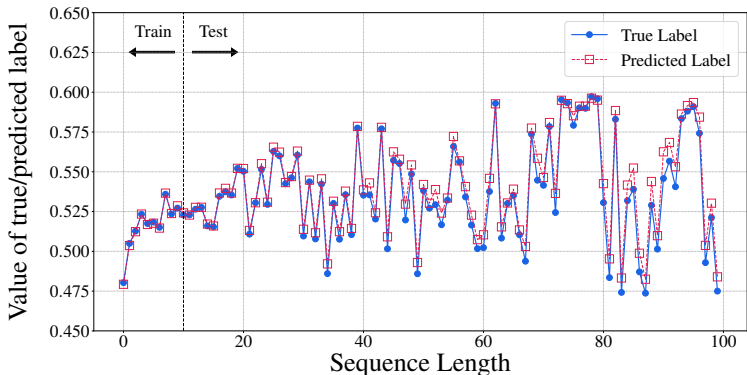
We opted for such a procedure because rejection sampling from the complement of the training set given in Figure 1 is extremely slow. In particular, given our batch size of 256, token dimension  $n =$

2538  
2539  
2540  
2541  
2542  
2543  
2544  
2545  
2546  
2547  
2548  
2549  
2550  
2551  
2552



2553  
2554 Figure 15: Illustrating the modified support of train vs test distribution for compositional generalization. This enables speed up in the sampling procedure, while keeping the challenging aspect of generalization to the corners.  
2555  
2556

2557  
2558  
2559  
2560  
2561  
2562  
2563  
2564  
2565  
2566  
2567  
2568  
2569



2570  
2571 Figure 16: A deep set model with one hidden layer MLP for  $\omega, \psi$  trained on sequences of length  
2572 up to  $T = 10$  sampled according to Figure 15 can generalize to unseen test sequences (Figure 15).  
2573 Additionally, the compositional generalization holds even beyond the training length.  
2574

2575  
2576  
2577  
2578  
2579  
2580  
2581  
2582  
2583  
2584  
2585  
2586

20, and having 100 batches per epoch, constructing the full test set requires finding  $256 \times 100 \times 20$  sequences of length  $t \leq T$  that are rejected by the original constraints. This becomes quite inefficient and slow specially in higher dimensions as the sum of the sequence along each component tends to concentrate more around  $t/2$ , therefore it becomes harder to find such sequences (the sum follows Irwin-Hall distribution since the components come from the Uniform distribution). To improve the speed of sampling the test dataset, we sample token dimensions  $x_i^k$  from the smaller corners shown in Figure 15 which allows for parallel sampling. These corners correspond to sampling  $x_i^k \sim \text{Uniform}[0, 1/2T]$  or  $x_i^k \sim \text{Uniform}[1/2 + 1/2T, 1]$ . This way we can sample token components independently and in parallel without having to reject any samples, since by construction no test sequence coincides with the training set. This procedure leaves a gap (see Figure 15) that will not be sampled neither during training nor testing.

2587  
2588  
2589  
2590  
2591

### D.3 FAILURE CASES

Although most of our focus has been on the success scenarios for length and compositional generalization, here we provide examples to show how a model might fail.



2592  
 2593  
 2594  
 2595  
 2596  
 2597  
 2598  
 2599  
 2600  
 2601  
 2602  
 2603  
 2604  
 2605  
 2606  
 2607  
 2608  
 2609  
 2610  
 2611  
 2612  
 2613  
 2614  
 2615  
 2616  
 2617  
 2618  
 2619  
 2620  
 2621  
 2622  
 2623  
 2624  
 2625  
 2626  
 2627  
 2628  
 2629  
 2630  
 2631  
 2632  
 2633  
 2634  
 2635  
 2636  
 2637  
 2638  
 2639  
 2640  
 2641  
 2642  
 2643  
 2644  
 2645

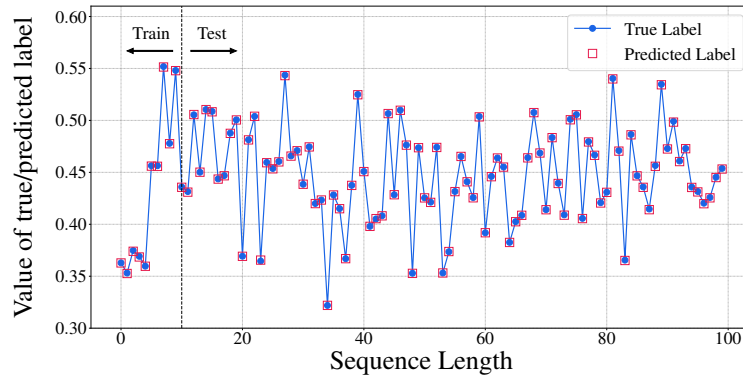


Figure 17: A Transformer model with softmax attention and one hidden layer MLP for  $\omega$  trained on sequences of length up to  $T = 10$  sampled according to Figure 15 can generalize to unseen test sequences (Figure 15). Additionally, the compositional generalization holds even beyond the training length.

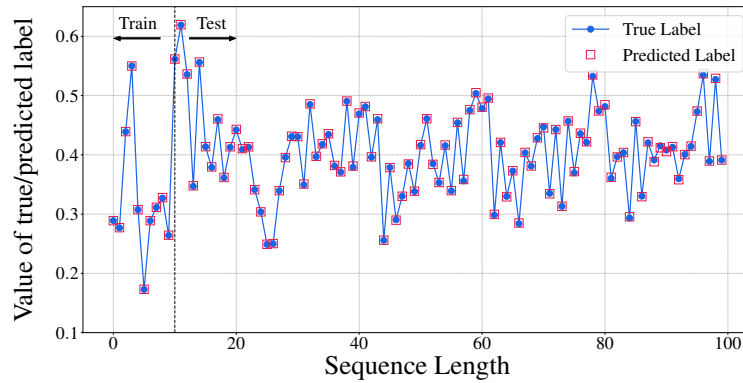


Figure 18: A SSM model with one hidden layer MLP for  $\omega$  trained on sequences of length up to  $T = 10$  sampled according to Figure 15 can generalize to unseen test sequences (Figure 15). Additionally, the compositional generalization holds even beyond the training length.

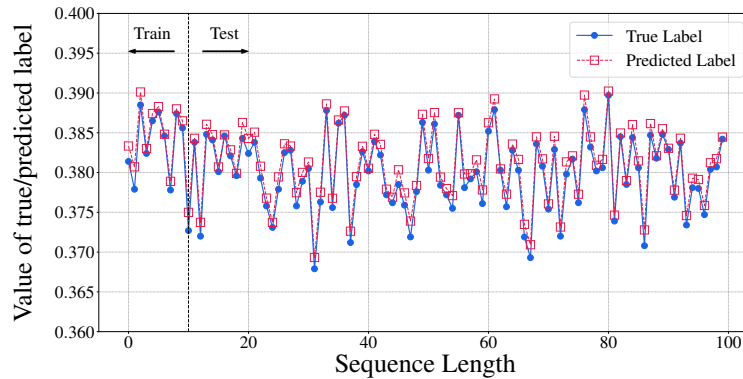
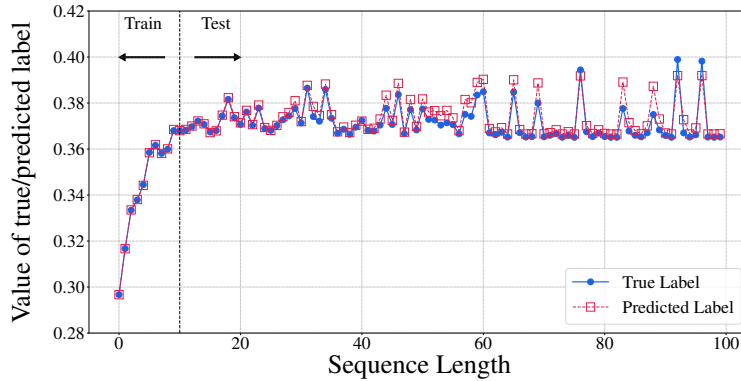


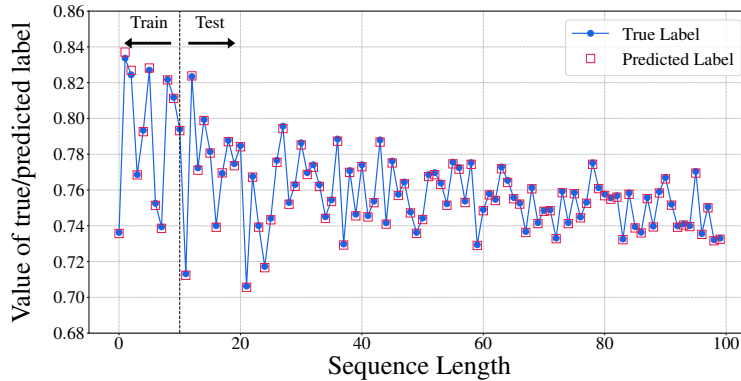
Figure 19: A RNN model with one hidden layer MLP for  $\omega$  trained on sequences of length up to  $T = 10$  sampled according to Figure 15 can generalize to unseen test sequences (Figure 15). Additionally, the compositional generalization holds even beyond the training length.

2646  
2647  
2648  
2649  
2650  
2651  
2652  
2653  
2654  
2655  
2656  
2657  
2658  
2659



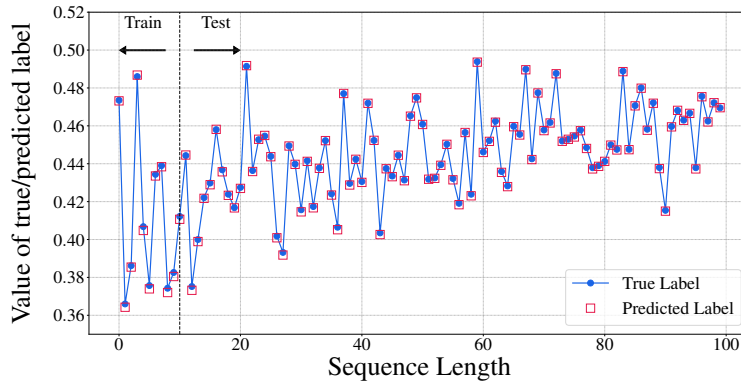
2660 Figure 20: A deep set model with *two* hidden layer MLP for  $\omega, \psi$  trained on sequences of length up  
2661 to  $T = 10$  sampled according to Figure 15 can generalize to unseen test sequences. Additionally,  
2662 the compositional generalization holds even beyond the training length.  
2663  
2664

2665  
2666  
2667  
2668  
2669  
2670  
2671  
2672  
2673  
2674  
2675  
2676  
2677



2678 Figure 21: A Transformer model with softmax attention and *two* hidden layer MLP for  $\omega$  trained  
2679 on sequences of length up to  $T = 10$  sampled according to Figure 15 can generalize to unseen test  
2680 sequences. Additionally, the compositional generalization holds even beyond the training length.  
2681  
2682

2683  
2684  
2685  
2686  
2687  
2688  
2689  
2690  
2691  
2692  
2693  
2694  
2695



2696 Figure 22: A SSM model with *two* hidden layer MLP for  $\omega$  trained on sequences of length up to  
2697  $T = 10$  sampled according to Figure 15 can generalize to unseen test sequences. Additionally, the  
2698 compositional generalization holds even beyond the training length.  
2699

2700  
2701  
2702  
2703  
2704  
2705  
2706  
2707  
2708  
2709  
2710  
2711  
2712  
2713  
2714  
2715  
2716  
2717  
2718  
2719  
2720  
2721  
2722  
2723  
2724  
2725  
2726  
2727  
2728  
2729  
2730  
2731  
2732  
2733  
2734  
2735  
2736  
2737  
2738  
2739  
2740  
2741  
2742  
2743  
2744  
2745  
2746  
2747  
2748  
2749  
2750  
2751  
2752  
2753

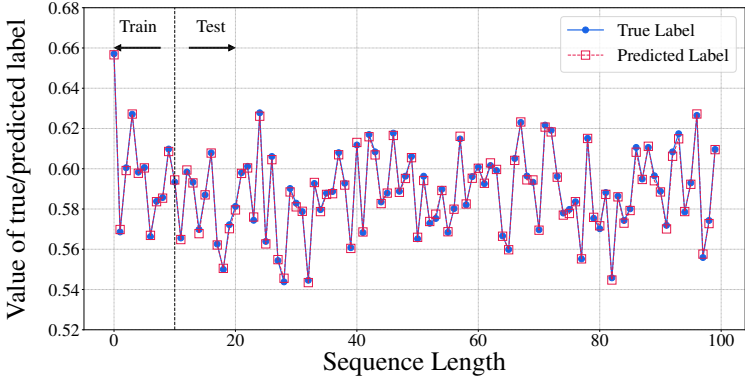


Figure 23: A RNN model with *two* hidden layer MLP for  $\omega$  trained on sequences of length up to  $T = 10$  sampled according to Figure 15 can generalize to unseen test sequences. Additionally, the compositional generalization holds even beyond the training length.

***f* is not realizable in  $\mathcal{H}$**  In Figures 25, 26, 27, 28, we present the predictions of learned models from different architectures initialized in their respective  $\mathcal{H}$  that does not contain the true  $f$ . In particular, we have the following for the different architectures:

- Deep set: The true labeling function  $f$  is a deep set with one hidden layer MLPs for  $\rho, \phi$ , but the learner uses  $h$ , a deep set model for which the MLPs  $\psi, \omega$  have no hidden layers.
- Transformer: The true labeling function  $f$  is a Transformer with 1 hidden layer in  $\rho$ , but the learner uses  $h$ , an RNN with 1 hidden layer in  $\omega$ .
- SSM: The true labeling function  $f$  is an SSM with 1 hidden layer in  $\rho$ , but the learner uses  $h$ , an RNN without any hidden layers in  $\omega$ .
- RNN: The true labeling function  $f$  is an RNN with 1 hidden layer in  $\rho$ , but the learner uses  $h$ , a Transformer with 1 hidden layer in  $\omega$ .

In each case, the learner is trained on sequences of length up to  $T = 10$  and its performance on the test set at longer lengths indicates whether generalization is possible or not. For a visual illustration of such failures beyond the training length, see Figures 25, 26, 27, 28. We can observe that the model can predict the test sequence well up to the length it has learned during training, but starts to diverge from the true labels beyond that. This demonstrates a failure case in which the realizability condition is violated.

***f* is realizable in a high capacity  $\mathcal{H}$**  For a given  $\mathcal{H}$ , if all solutions to 1 achieve length generalization or compositional generalization, then we can guarantee length or compositional generalization regardless of the training procedure. When the capacity of  $\mathcal{H}$  becomes very large, it continues to contain the right solutions but it starts to contain many incorrect solutions that match the true solution only on the support of training distribution. In such a case, there is no reason to presume that our learning procedure picks the right solution to 1 that also achieves length and compositional generalization. Figure 24 show experiments illustrating the above. We experiment with the following scenarios for deep sets and transformers:

- Deep set: We use the labeling function that takes the following form  $f = \rho(\sum_{i \leq t} \phi(x_i))$  for  $t \leq T$  and  $f = \rho(\sum_{i \leq t} \phi(x_i)) + c$  for  $t > T$  with  $c = 0.2, T = 5$ . We use 1 hidden layer MLPs for  $\rho, \phi$  (with no activation on the output of  $\rho$ ). We use 2 hidden layer MLPs for  $\omega, \psi$  for  $h$  so that it can express the above labeling function. The input, hidden, and output dimensions are all equal  $m = n = k = 20$  for  $f, h$ . We train on sequences of length longer than  $T$  to demonstrate this expressivity claim. When the model is trained on sequences of length less than  $T$ , due to the simplicity bias of the training procedure model learns  $\rho(\sum_{i \leq t} \phi(x_i))$  and uses it on longer sequences and hence fails.
- Transformer: We use the labeling function that takes the following form  $f = \rho(\sum_{j=1}^i \frac{1}{i} \phi(x_i, x_j))$  for  $t \leq T$  and  $f = \rho(\sum_{j=1}^i \frac{1}{i} \phi(x_i, x_j)) + c$  for  $t > T$  with

2754  
2755  
2756  
2757  
2758  
2759  
2760  
2761  
2762  
2763  
2764  
2765  
2766  
2767  
2768  
2769  
2770  
2771  
2772  
2773  
2774  
2775  
2776  
2777  
2778  
2779  
2780  
2781  
2782  
2783  
2784  
2785  
2786  
2787  
2788  
2789  
2790  
2791  
2792  
2793  
2794  
2795  
2796  
2797  
2798  
2799  
2800  
2801  
2802  
2803  
2804  
2805  
2806  
2807

Deep set	Loss ( $t < T_0$ )	Loss ( $t \geq T_0$ )
Fig 2-a	$0.001 \pm 10^{-4}$	$0.002 \pm 3 \times 10^{-4}$
Fig 2-b	$0.0007 \pm 10^{-4}$	$0.007 \pm 0.001$
Transformer	Loss ( $t < T_0$ )	Loss ( $t \geq T_0$ )
Fig 2-c	$0.0006 \pm 10^{-4}$	$0.006 \pm 3 \times 10^{-3}$
Fig 2-d	$10^{-5} \pm 10^{-6}$	$0.01 \pm 0.003$

Table 5: Length generalization of different architectures when the hypothesis class  $\mathcal{H}$  is highly expressive. For further details see Fig. 24

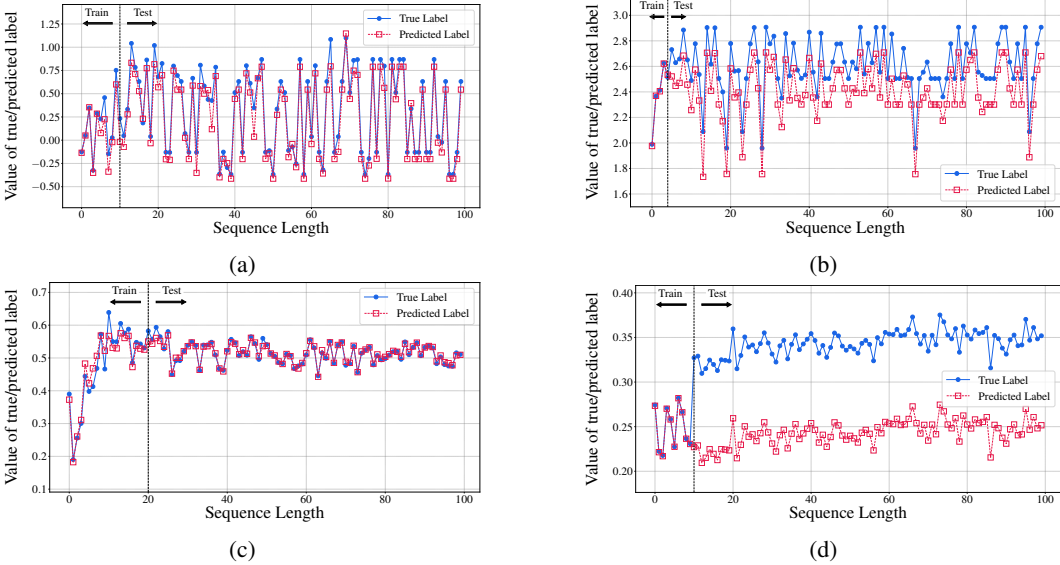
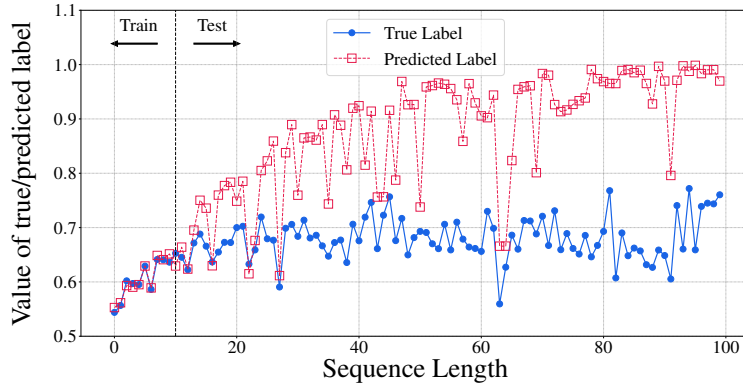


Figure 24: A failure case of length generalization under arbitrary expressive generative model with (a,b) Deep sets, (c,d) and Transformer. The generative function on both cases introduces an offset to sequences longer than some critical length ( $T_0$ ). The learner is once trained on sequences longer than  $T_0$  and successfully generalizes (a,c), and once is trained only on sequences shorter than  $T_0$  where the offset never appears, and hence fails to generalize beyond that.

$c = 0.1, T = 10$ . We use 1 hidden layer MLPs for  $\rho$  (with no activation on the output of  $\rho$ ). We use a Transformer with 3 hidden layer MLPs for  $\omega$  so that it can express the above labeling function. The input, hidden, and output dimensions are all equal  $m = n = k = 20$  for  $f, h$ . We train on sequences of length longer than  $T$  to demonstrate this expressivity claim. When the model is trained on sequences of length less than  $T$ , due to the simplicity bias of the training procedure model learns  $f = \rho(\sum_{j=1}^i \frac{1}{i} \phi(x_i, x_j))$  and uses it on longer sequences and hence fails.

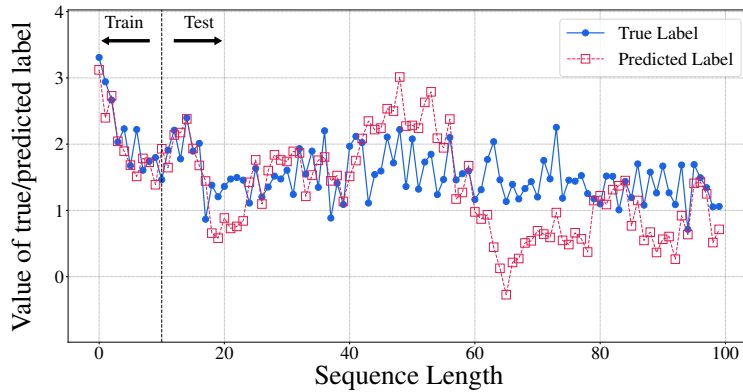
The failures of such degenerate solutions can be visualized in Figure 24 (right), where the predictions diverge from the true values when the model is only trained on sequences shorter than  $T_0$ . Figure 24 (left) shows that when the model is trained on sequences longer than  $T_0$ , it can successfully generalize to longer lengths. Table 5 further validates this observation numerically. It presents the test loss of each model at lengths shorter and longer than  $T_0$  under the two training schemes: a) When trained only on sequences of length shorter than  $T_0$  (rows corresponding to Fig 2-b and 2-d which result in failure due to degenerate solution), b) when trained on sequences of length longer than  $T_0$  (rows corresponding to Fig 2-a and 2-c which result in successful generalization).

2808  
2809  
2810  
2811  
2812  
2813  
2814  
2815  
2816  
2817  
2818  
2819  
2820



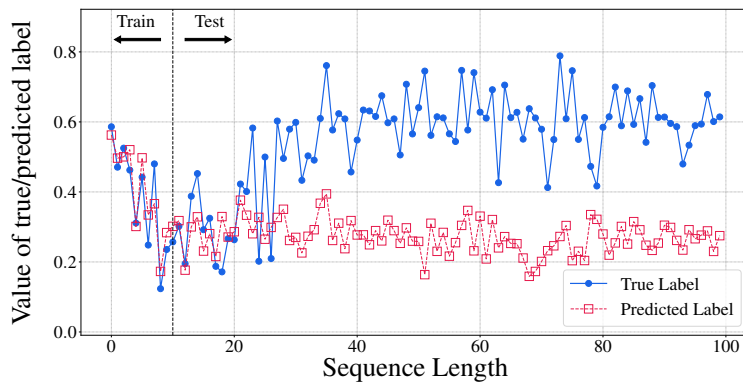
2821 Figure 25: A failure case of length generalization in the unrealizable setting: The predictions come from a deep set with linear layers for  $\psi, \omega$  trained to predict the sequences (of length up to  $T$ ) output by a deep set with 1 hidden layer MLPs for  $\phi, \rho$ . In this case the realizability condition does not hold, and the learner fails to length generalize.

2826  
2827  
2828  
2829  
2830  
2831  
2832  
2833  
2834  
2835  
2836  
2837  
2838  
2839



2840 Figure 26: A failure case of length generalization in the unrealizable setting: The predictions come from an RNN with 1 hidden layer in  $\omega$  trained to predict the sequences (of length up to  $T = 10$ ) output by a Transformer with 1 hidden layer in  $\rho$ . In this case the realizability condition does not hold, and the learner fails to length generalize.

2844  
2845  
2846  
2847  
2848  
2849  
2850  
2851  
2852  
2853  
2854  
2855  
2856  
2857



2858 Figure 27: A failure case of length generalization: The predictions come from an RNN without any hidden layers in  $\omega$  trained to predict the sequences (of length up to  $T = 10$ ) output by an SSM with 1 hidden layer in  $\rho$ . In this case the realizability condition does not hold, and the learner fails to length generalize.

2862  
 2863  
 2864  
 2865  
 2866  
 2867  
 2868  
 2869  
 2870  
 2871  
 2872  
 2873  
 2874  
 2875  
 2876  
 2877  
 2878  
 2879  
 2880  
 2881  
 2882  
 2883  
 2884  
 2885  
 2886  
 2887  
 2888  
 2889  
 2890  
 2891  
 2892  
 2893  
 2894  
 2895  
 2896  
 2897  
 2898  
 2899  
 2900  
 2901  
 2902  
 2903  
 2904  
 2905  
 2906  
 2907  
 2908  
 2909  
 2910  
 2911  
 2912  
 2913  
 2914  
 2915

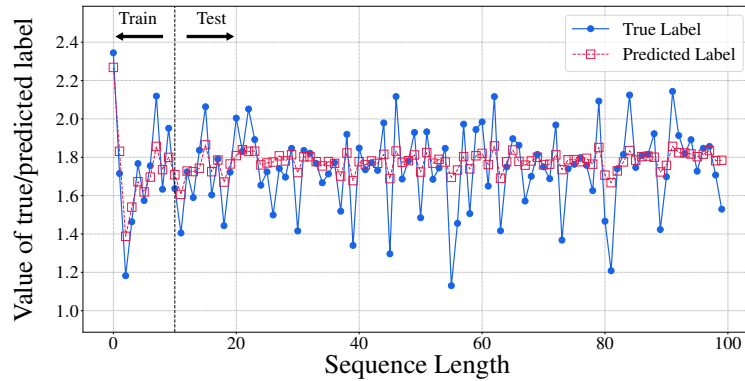


Figure 28: A failure case of length generalization: The predictions come from a Transformer with 1 hidden layer in  $\omega$  trained to predict the sequences (of length up to  $T = 10$ ) output by an RNN with 1 hidden layer in  $\rho$ . In this case the realizability condition does not hold, and the learner fails to length generalize.