

SIP: Injecting a Structural Inductive Bias into a Seq2Seq Model by Simulation

Matthias Lindemann[✉] and Alexander Koller[✉] and Ivan Titov[✉]

[✉]ILCC, University of Edinburgh, [✉]LST, Saarland University, [✉]ILLC, University of Amsterdam
m.m.lindemann@sms.ed.ac.uk, koller@coli.uni-saarland.de, ititov@inf.ed.ac.uk

Abstract

Strong inductive biases enable learning from little data and help generalization outside of the training distribution. Popular neural architectures such as Transformers lack strong structural inductive biases for seq2seq NLP tasks on their own. Consequently, they struggle with systematic generalization beyond the training distribution, e.g. with extrapolating to longer inputs, even when pre-trained on large amounts of text. We show how a structural inductive bias can be efficiently injected into a seq2seq model by pre-training it to simulate structural transformations on synthetic data. Specifically, we inject an inductive bias towards Finite State Transducers (FSTs) into a Transformer by pre-training it to simulate FSTs given their descriptions. Our experiments show that our method imparts the desired inductive bias, resulting in improved systematic generalization and better few-shot learning for FST-like tasks. Our analysis shows that fine-tuned models accurately capture the state dynamics of the unseen underlying FSTs, suggesting that the simulation process is internalized by the fine-tuned model.¹

1 Introduction

Inductive biases, i.e. the preferences and the abstract knowledge a model brings to the task, enable a model to learn from small amounts of data and generalize systematically outside of the training distribution. While seq2seq models perform very well on in-distribution data, they usually lack structural inductive biases and consequently struggle with systematic generalization. Previous work has shown that this includes generalization to unseen combinations of known sub-strings (Lake and Baroni, 2018; Keyzers et al., 2020), extrapolation to longer inputs (Hupkes et al., 2020) and deeper recursion (Kim and Linzen, 2020).

Integrating structural inductive biases into seq2seq models is challenging. One popular ap-

¹We release our code at <https://github.com/namednil/sip>

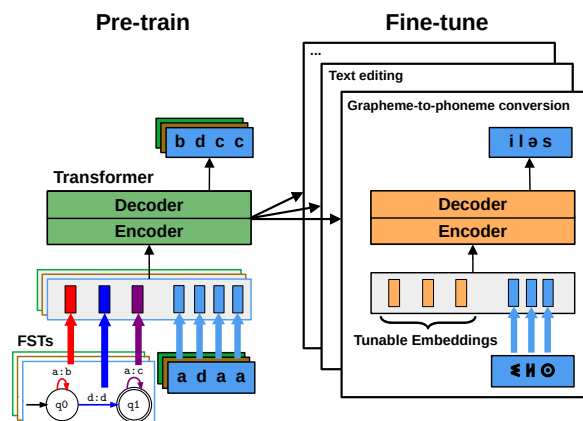


Figure 1: Left: Pre-training a Transformer to simulate automatically generated FSTs. Right: fine-tuning the Transformer and the prefix where the FST used to be on a downstream task by using only input/output pairs. Tunable parameters are represented in orange.

proach is to develop specialized architectures (Wu and Cotterell, 2019; Zheng and Lapata, 2021; Kim, 2021; Lindemann et al., 2023), which makes it difficult to precisely control and adjust the nature of the inductive bias to changing demands as the architecture would need to be modified and models re-trained. Recently, some works instead have tried to inject inductive biases into seq2seq models by pre-training on a well-chosen synthetic task (Krishna et al., 2021; Wu et al., 2021, 2022) or meta-learning on a distribution of synthetic tasks (McCoy et al., 2020; McCoy and Griffiths, 2023) using MAML (Finn et al., 2017). Here, the inductive bias can be controlled by the choice of the synthetic task. However, meta-learning with MAML scales poorly because it requires expensive second-order derivatives and standard pre-training can be less effective (McCoy and Griffiths, 2023).

In this work, we present a computationally inexpensive way of injecting a structural inductive bias into a Transformer. We focus specifically on introducing an inductive bias that is helpful for tasks that traditionally have been approached with Finite State Transducers (FSTs). We choose FSTs

because they are formally well understood, are easy to generate automatically, and are one of the simplest computational devices that are useful in NLP applications. While we focus on FSTs, the methodology is fairly general and our approach also provides a starting point for incorporating more general structural biases, provided by more expressive formalisms such as Pushdown Transducers.

Our approach (SIP, for **S**imulation-**I**nduced **P**rior) is simple (see Fig. 1): given a representation of an FST and an input string, a Transformer is pre-trained to predict what the output of the FST is on the given input. We assume that FSTs are not specified for fine-tuning on downstream tasks, so we replace the FST with tunable embeddings and fine-tune the model solely on input/output pairs. Since we fine-tune all parameters, the model can deviate from FST-like behavior if needed.

Contributions. We show that a model pre-trained with SIP has an inductive bias that improves systematic generalization and few-shot learning for ‘FST-like’ downstream tasks. SIP not only improves systematic generalization on FST tasks similar to those seen during pre-training but also on ones that are structurally more complex. The same pre-trained model also transfers well to natural data and achieves strong results on few-shot learning of text editing (e.g. Jane Doe \rightarrow J. Doe) and grapheme-to-phoneme conversion.

Our probing experiments give insights into how the inductive bias is injected: SIP not only leads to the imitation of the input/output behaviour of FSTs, but encourages dynamics to emerge that *simulate* crucial aspects of FSTs in the hidden representations. Fine-tuning can leverage these dynamics, providing the inductive bias, and learn representations that resemble those of ground truth FSTs.

2 Related Work

Systematic generalization. Systematic generalization refers to the ability of a model to generalize (or extrapolate) beyond its training distribution in a systematic way that aligns with how humans generalize. Systematic generalization is difficult for standard seq2seq models in contexts such as semantic parsing (Finegan-Dollak et al., 2018) and machine translation (Li et al., 2021), in particular to unseen combinations of phrases, longer inputs as well as deeper recursion (Keysers et al., 2020; Kim and Linzen, 2020).

A range of approaches have been developed to tackle this, with many works focusing on special-

ized architectures (Guo et al., 2020; Kim, 2021; Lindemann et al., 2023). Furrer et al. (2020) find that the specialized architectures they consider do not transfer well to tasks beyond the context in which they were designed. This highlights the importance of being able to adjust inductive biases more easily than re-designing the architecture of a model. Large-scale pre-training has also been shown to help with systematic generalization (Furrer et al., 2020). However, challenges remain even for LLMs such as GPT-3 and PALM (Qiu et al., 2022; Dziri et al., 2023). The methodology we present in this work can be used to create additional material for LLM pre-training. Here we focus on smaller models and leave this to future work.

Pre-training with synthetic tasks. Pre-training a model on a synthetic task to introduce specific inductive biases has been explored by several recent works. Krishna et al. (2021) identify useful ‘skills’ for news summarization and develop a pre-training task accordingly. LIME (Wu et al., 2021) targets mathematical reasoning and is pre-trained on string manipulation that resembles formal reasoning. Papadimitriou and Jurafsky (2023) consider pre-training with several synthetic languages to investigate which helps most for language modelling. In contrast to these works, our approach targets simulating a computational device and maintains a closer relation to the pre-training setting because of the tunable prefix.

A challenge for using individually hand-crafted tasks is to cover a sufficient space of phenomena that are relevant to downstream tasks. Instead of training on a single task only, McCoy et al. (2020); McCoy and Griffiths (2023) meta-learn on a distribution of tasks using MAML (Finn et al., 2017). Our approach also uses a distribution of tasks but it scales better than MAML-based methods because MAML requires computing and storing second-order derivatives. For example, the Transformer we train has a magnitude more parameters than the LSTM of McCoy and Griffiths (2023) and is pre-trained on a smaller GPU (A100 vs RTX 2080 TI). In addition, as the complexity of each individual task grows, MAML requires more examples per task. We circumvent this by using a compact and unambiguous description of each task instead.

Simulating execution. The idea of using a neural network to predict the outcome of the execution of a computational device or code has come up in several contexts over the last few years. Early work by Zaremba and Sutskever (2014) investi-

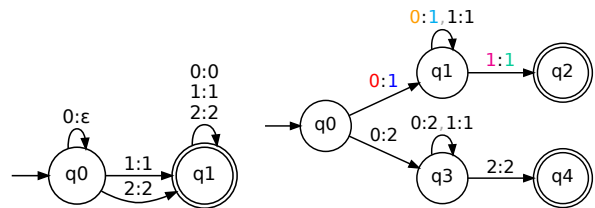
gates it as a challenging benchmark for LSTM-based seq2seq models. Recent works have explored simulating (aspects of) code execution for various down-stream applications, such as program synthesis (Austin et al., 2021), or debugging and code analysis (Bieber et al., 2022) as well as reverse engineering (Pei et al., 2021). Finally, Finlayson et al. (2022) train a Transformer to interpret regular expressions: given a regular expression and a string, the task is to decide if the string is in the regular language. There are three crucial differences between their work and ours: (i) they investigate the empirical capabilities of Transformers while we introduce structural inductive biases for downstream tasks, (ii) they consider binary outputs whereas we consider sequential outputs, and (iii) we perform probing experiments showing strong evidence for FST simulation in the hidden representations.

Emergent World Representations. Our analysis provides evidence that our model trained with SIP internally simulates transitions between FST states even though it was not explicitly supervised to do so. Similar observations have been made for Language Models trained to play Othello (Li et al., 2023) and chess (Karvonen, 2024), where the model was found to acquire a representation of the board state simply from being trained to predict the next move.

3 Finite State Transducers

We briefly review Finite State Transducers (FSTs) which we use in our experiments. FSTs are closely related to Finite State Automata (FSAs). While an FSA describes a set of strings, an FST describes a *relation* between strings, i.e. a set of pairs (x, y) , where x is an input y is an output.

FSTs can be visualized as labelled directed graphs (see Fig. 2), where the nodes are called *states* and the edges are called *transitions*. Consider the path $q_0 \xrightarrow{\emptyset:1} q_1 \xrightarrow{\emptyset:1} q_1 \xrightarrow{1:1} q_2$ in Fig. 2b. This path is called an *accepting path* because it starts in an *initial* state (indicated by an arrow ‘from nowhere’ pointing to the state), and it ends in a *final* state (indicated by double circles). An accepting path shows what an input can be mapped to. In this case, the path shows that the FST transduces the input $\emptyset\emptyset 1$ into the output 111 . We can read off which input an accepting path associates an output to by concatenating all the strings along the path occurring before ‘:’. The output can be determined by concatenating the strings after ‘:’. Hence, each transition $\xrightarrow{\sigma:\rho}$ can be thought of as



(a) A deterministic FST. (b) A non-deterministic but functional FST.

Figure 2: Examples of *functional* FSTs. The FST in (a) deletes leading zeros. The FST in (b) replaces any \emptyset by a 1 if the last input symbol is a 1. Conversely, if the last symbol is a 2, any \emptyset is replaced by a 2. The output can only be determined after the last input symbol.

‘replacing’ σ by ρ . Inserting and deleting can be achieved by means of the empty string, written as ϵ . For example, Fig. 2a ‘replaces’ leading zeros by an empty string, effectively deleting them.

In general, an input can be paired with arbitrarily many different outputs. We call an FST f **functional** if every input x is paired with at most one output y , and use the notation $f(x)$ to refer to y . All FSTs we consider here are functional.

In this work, we investigate generalization across different sub-classes of FSTs, namely from the less expressive deterministic FSTs to non-deterministic FSTs. An FST is called **deterministic** if (i) it has a unique initial state, (ii) for all states q and input symbols σ there is at most one transition $q \xrightarrow{\sigma:\rho} q'$ and (iii) $\sigma \neq \epsilon$. Intuitively, this means that in any state, for an input symbol σ there is at most one possible next state and one possible output, and hence for any input string there is at most one path that is compatible with it. Because of this, we can always infer a prefix of the output by looking only at a *prefix* of the input string and ignoring the rest. For example, consider the input prefix $\emptyset\emptyset 1$. In the deterministic FST in Fig. 2a, we know that the output has to start with 1 because there is only one path that is compatible with $\emptyset\emptyset 1$. In contrast, in the non-deterministic FST in Fig. 2b, three paths are compatible with $\emptyset\emptyset 1$ that have different outputs. In that case, we can only determine the output once we look at the last symbol of the input. In short, non-deterministic FSTs can take context to the right into account but deterministic FSTs cannot.

4 Simulation-Induced Prior

Our approach largely follows the pre-training and fine-tuning paradigm. We first pre-train on synthetic FST tasks by giving the model a representation of an FST as a prefix and an input string (see Fig. 1). The training objective is to predict the

output of the FST on the input string. Our research hypothesis is that training a model to predict the behaviour of an FST incentivizes the model to acquire reusable dynamics that internally simulate FSTs. When fine-tuning the model using a tunable prefix instead of an encoding of an FST, these dynamics should be easy to leverage and provide a structural inductive bias for FST-like tasks.

4.1 Pre-training

During pre-training, the model is given a representation of an FST and a string in its domain and has to predict the output of that FST on the given input string. The input to the Transformer is a sequence of vectors from \mathbb{R}^d , which consist of a prefix that represents the FST f and a suffix comprised of the embeddings of the input string (see Fig. 1):

$$\underbrace{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_k}_{\text{FST encoding}}, \underbrace{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n}_{\text{Input to FST}}$$

Each \mathbf{h}_i encodes a transition $p \xrightarrow{\sigma:\rho} q$ as a vector:

$$\mathbf{h}_i = W[\text{EMB}_{\text{State}}(p); \text{EMB}_{\text{State}}(q); \text{EMB}_{\text{Symbol}}(\sigma); \text{EMB}_{\text{Symbol}}(\rho); \text{EMB}_{\text{Final}}(e)]$$

where $[\cdot]$ represents vector concatenation, e indicates if q is a final state, and W is linear layer that ensures that $\mathbf{h} \in \mathbb{R}^d$. All embeddings are simple look-up tables based on the id of the state or symbol. The initial state of the FST is always assigned the id 0, and positional embeddings are used as usual. The model is trained to maximize the log probability of the output $y = f(x)$ of the FST f .

4.2 Fine-tuning

After pre-training, we can apply our model to a downstream task and fine-tune it. We assume we do not have access to an FST for the downstream task, and therefore we replace the FST encoding with a sequence of tunable embeddings. That is, the input to the model is a sequence of vectors:

$$\underbrace{\mathbf{h}'_1, \mathbf{h}'_2, \dots, \mathbf{h}'_k}_{\text{Tunable embeddings}}, \underbrace{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n}_{\text{Input}}$$

where $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are the embeddings of the input tokens, $\mathbf{h}'_i \in \mathbb{R}^d$ are the tunable embeddings and k is a hyperparameter. The embeddings \mathbf{h}'_i are initialized to the average of the encoding of multiple FSTs from the pre-training phase. The most straightforward way to fine-tune is to only modify \mathbf{h}' because we are looking for an FST-like task representation. This is similar to prompt tuning (Lester et al., 2021). However, this does not work

well on tasks outside the pre-training distribution. Hence, we fine-tune the entire model, including the prefix, and use a higher learning rate for the prefix than for the rest of the model (see Appendix E).

4.3 Constructing Pre-Training Data

To create our pre-training data, we sample 40,000 deterministic FSTs. For every FST, we sample 5 input/output pairs with input lengths up to 35. In total, this leads to 200,000 pairs for training along with their FSTs. To describe the sampling procedure in more detail, we use an overall vocabulary V consisting of the printable ASCII tokens and the Unicode block for IPA symbols (used for transcribing speech). Seq2seq tasks in the wild usually do not use the whole space of this vocabulary, so for each task T we first uniformly sample the vocabulary size $|V_T|$ between 5 and 25 and then uniformly select a subset $V_T \subseteq V$. Then, we uniformly sample the number of states $|Q_T|$ between 2 and 4, and the number of final states between 1 and $|Q_T|$. For every state q and every symbol $\sigma \in V_T$ we introduce at most one outgoing transition to a state q' , chosen uniformly at random. This ensures that the FST is deterministic. We then sample the output for the transition: either a symbol $\rho \in V_T$ or ϵ . Finally, we minimize the number of states of the FST using OpenFST (Allauzen et al., 2007), and exclude those without cycles, as they express finite relations. See Appendix A.1 for details.

In practical applications of FSTs, in particular for text editing, one often wants to keep certain parts of the input unchanged. This can be achieved with a set of transitions of the form $q \xrightarrow{\sigma:\sigma} q'$ for all $\sigma \in V_T$. Since it is very unlikely to sample such a set of transitions, we use a special symbol that acts as a shorthand for this, which we also use when encoding the FST for pre-training.

5 Evaluating SIP’s Inductive Bias

To understand the effects of our pre-training procedure on the inductive bias of the model and on the downstream performance, we first explore systematic generalization on synthetic FST tasks. This allows us to precisely control the similarity between the pre-training and the downstream task.

5.1 Evaluation Methodology

To evaluate the degree to which a model has an inductive bias towards FSTs, we now describe two methods for generating training and test data that

reward a model for showing important aspects of FST-like systematic generalization.

Iteration generalization. Cycles are a characteristic feature of FSTs, and iteration generalization tests if a model learns that cycles can be traversed more often than seen during training. More specifically, given an FST, we generate training data which requires visiting any state only a few times (*iteration count* up to 3). In the test data, the model has to generalize to visiting states more often (iteration count at least 4). This notion is related to length generalization (Lake and Baroni, 2018) but tailored specifically to FSTs.

Unseen combinations of transitions (UC). When an FST processes a string, the set of possible next transitions only depends on the current FST state; it does not matter how the current state was reached. Hence, a model with an inductive bias towards FSTs should also not be overly sensitive to how a state is reached, and correctly handle situations where a specific *combination* of transitions was unobserved during training. For example, consider the FST in Fig. 2a, which deletes leading zeros from a number. Suppose that a model is trained on examples such as 0012 , 2201 , 1012 but no training example contains the combination of leading zeros followed by a 2, which corresponds to using the combination of the transitions $q0 \xrightarrow{0:\epsilon} q0$ and $q0 \xrightarrow{2:2} q1$. If the model has an inductive bias towards FSTs, it should generalize to this unseen combination and correctly handle inputs such as 0021 . To generate appropriate training and test data for this, we sample a pair of adjacent transitions (such as $q0 \xrightarrow{0:\epsilon} q0$ and $q0 \xrightarrow{2:2} q1$ in Fig. 2a) and ensure that no training example uses both transitions within the *same* string. In contrast, in the test data, all examples require using the *combination* of the transitions. To make the generalization setup more challenging, we ensure this for multiple pairs of transitions at the same time. We refer to Appendix A.3 for details on the construction.

UC is related to the method of Keysers et al. (2020) who also withhold combinations of seen elements to assess systematic generalization.

5.2 Setup and Baselines

To make a fair comparison, all models we experiment with in the main paper share the same architecture and are initialized from the same checkpoint before any additional pre-training, namely ByT5-small (Xue et al., 2022). ByT5 has 300 million parameters and was pre-trained on the multilingual

C4 corpus. It uses raw bytes as tokens, which enables full Unicode support and is a natural unit to consider for FST-like tasks such as text editing and grapheme-to-phoneme conversion. We report additional results with a T5-Base model in Appendix B.3, where we observe similar trends.

SIP-d4. This is a model using the method we propose in this work. We pre-train on the data generated in Section 4.3 (deterministic FSTs, with up to 4 states) for 20 epochs. This model achieves an average sequence-level accuracy of 98% on predicting the output of an unseen FST from the training distribution. For fine-tuning, we use a prefix of length 50 for all experiments in this paper. As an ablation, we also fine-tune the model without the prefix of tunable embeddings (-prefix).

Naive pre-training. We use the same pre-training data as for SIP-d4 but omit the description of the FST and only train on input/output pairs.

Task embeddings (TE). TE is a simplified version of SIP. Instead of using an encoding of an FST in the prefix, this baseline uses 50 randomly initialized embeddings specific to each FST. The embeddings are learned from examples jointly with the rest of the model. Several works have used a single embedding to encode a domain/task in multi-task learning (Tsvetkov et al., 2016; Stymne et al., 2018; Zhang et al., 2022). Using a shorter tunable prefix resulted in considerably worse performance in our setup. TE is fine-tuned analogously to SIP, i.e. with a prefix of tunable embeddings.

Set. Wu et al. (2022) investigate the effectiveness of 18 simple synthetic pre-training tasks and found Set to perform best on average. The task is to deduplicate characters such that every type occurs only once, e.g. the input $dabacd$ becomes $dabc$. This task can be represented by a deterministic FST, albeit a very large one with 2^n states for a vocabulary of size n .

5.3 Systematic Generalization within the Pre-training Distribution

First, we want to establish to what degree the pre-training has conferred any inductive bias on the distribution it was pre-trained on.

Setup. For each generalization setup, we generate 5 unseen FSTs with 4 states each using the same procedure as for the pre-training. We fix the vocabulary size to its maximum value (25) in the pre-training data and only use printable ASCII characters in order to reduce variance across tasks. To evaluate UC, we withhold the combination of up to

	Iteration		UC	
	Acc \uparrow	ED \downarrow	Acc \uparrow	ED \downarrow
ByT5	37.8	5.87	47.4/57.5	1.49/0.93
Naive	42.6	4.41	44.9/43.2	1.52/1.35
Set	44.4	4.58	43.6/42.0	1.47/1.31
TE	61.3	2.49	57.3/63.1	1.13/0.74
SIP-d4	94.8	0.12	73.1/93.3	0.61/0.13
-prefix	84.9	0.62	61.1/76.3	0.99/0.50

Table 1: Evaluating systematic generalization on FST tasks with 4 states. We report averages over 5 tasks. ED is edit distance. Due to an outlier task on UC, we additionally report the median after ‘/’.

20 pairs of transitions and generate 5000 training examples with lengths 3 to 15 and corresponding test data as described in Section 5.1. For iteration generalization, we generate training examples with a maximum iteration count of 3 and test on longer examples of length up to 30 with an iteration count of at least 4. Since the out-of-distribution performance of two checkpoints of the same model can vary significantly, we report averages on the test set of the last 10 epochs.

Results. The results can be found in Table 1. On average, SIP-d4 achieves close to perfect accuracy (with one outlier on UC, skewing the mean). TE also shows a clear improvement over the other baselines but SIP-d4 outperforms TE by a large margin. This suggests that SIP-d4 and TE, to a lesser extent, indeed have acquired a stronger inductive bias for FSTs than the other methods. Using SIP-d4 without the tunable prefix leads to a substantial drop in accuracy, highlighting its importance. We analyze the representations learned by SIP-d4 in the tunable prefix in Appendix D.1.

5.4 More Complex FSTs

Does the inductive bias introduced by SIP extend beyond the pre-training distribution to more complex FST tasks? To investigate this, we use the same sampling methodology but generate FSTs with more states. SIP-d4 was pre-trained on FSTs with up to 4 states, and we evaluate on FST tasks with 5, 7 and 10 states.

We show in Fig. 3 how the individual models deviate from the accuracy of ByT5 as a function of the number of states in the test FST. SIP always performs best by a clear margin regardless of the number of states in the FSTs. As we increase the number of states and move away from the pre-training distribution, SIP improves less over the

	Iteration		UC	
	Acc \uparrow	ED \downarrow	Acc \uparrow	ED \downarrow
ByT5	83.4	0.52	83.1	0.40
Naive	83.1	0.49	84.2	0.37
Set	82.3	0.52	83.7	0.37
TE	84.2	0.49	82.7	0.42
SIP-d4	87.8	0.32	90.0	0.24
SIP-d4+	88.2	0.30	90.5	0.22
SIP-nd7	89.5	0.27	91.2	0.18

Table 2: Evaluation on non-deterministic FSTs. We report averages over 5 tasks.

baselines. We see a similar pattern for TE but with considerably smaller improvements over ByT5.

5.5 Non-Deterministic FSTs

As shown in the previous section, SIP still works well for more complex FST tasks than seen during pre-training. However, this evaluation focused on the favourable case where both pre-training and evaluation involve the same class of FSTs, namely deterministic FSTs. Deterministic FSTs can only take left context into account (see Section 3), which is a restrictive assumption. Here, we evaluate if the inductive bias conferred by SIP carries over to non-deterministic functional FSTs, i.e. those that can also take context to the *right* into account.

We automatically generate 5 non-deterministic FSTs with 21 states (see Appendix A.2 for details) and report averages in Table 2. Despite the structural mismatch between pre-training and the downstream tasks, SIP-d4 shows clear improvements over the baselines. Interestingly, TE does not consistently outperform the other baselines, despite its stronger results on deterministic FSTs.

Our pre-training procedure does not hinge on using deterministic FSTs. This raises the question if we can achieve even better performance by adjusting the inductive bias. To investigate this, we further pre-train SIP-d4 on 40,000 non-deterministic FSTs with up to 7 states, which we call SIP-nd7. To control for the additional training data of SIP-nd7, we also further pre-train SIP-d4 with the same number of deterministic FSTs with the same characteristics as in Section 4.3 (SIP-d4+). The results in Table 2 show better performance of SIP-nd7, which supports the hypothesis that the inductive bias can be adjusted. SIP-d4+ shows a smaller improvement over SIP-d4. Based on 5 additional FSTs per setup to gain more statistical power, we found that the difference between SIP-nd7 and SIP-

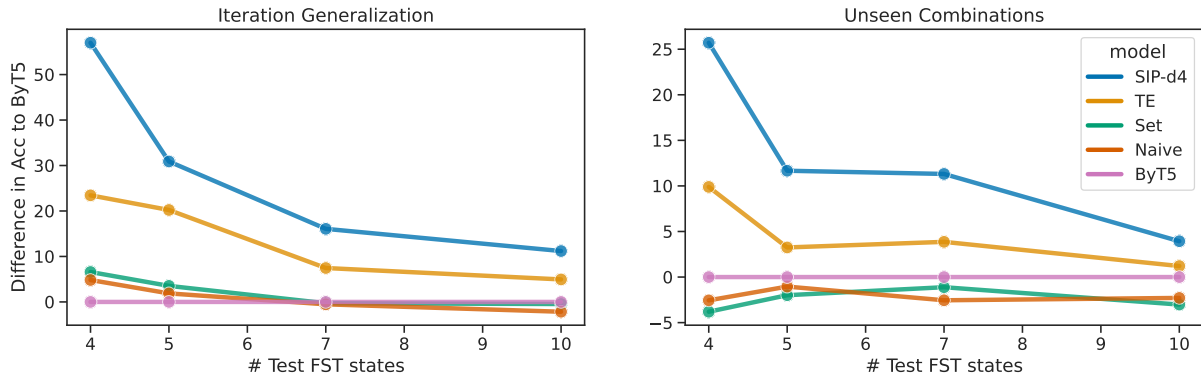


Figure 3: Evaluation on deterministic FST tasks with more states than seen in pre-training. We show the deviation in percentage points from ByT5.

d4+ is statistically significant ($p \approx 0.017$, $n = 20$, paired approx. permutation test).

6 Transfer to Natural Data

In this section, we investigate to what degree the inductive bias from pre-training on synthetic data transfers to tasks with natural data that have been traditionally approached with finite state methods.

6.1 Low-resource Grapheme-to-Phoneme Conversion

Grapheme-to-phoneme conversion is the task of converting a word as a sequence of symbols (for example, letters in the Latin alphabet) into a description of how this word is pronounced as letters in the IPA alphabet. For example, a possible pronunciation of ‘explanation’ is $[\text{ɛksplə'neiʃən}]$. Grapheme-to-phoneme conversion can be part of text-to-speech pipelines and FSTs for this purpose usually are two or three magnitudes larger than the FSTs we constructed for pre-training. Because of this, it enables us to test how far beyond the pre-training distribution SIP remains helpful. We focus on learning from small amounts of data, for which a structural inductive bias towards FSTs should be particularly helpful. We evaluate on 7 low-resource languages from different language families that use their own scripts (Balinese, Coptic, Gothic, Lao, Sylheti, Telugu and Central Atlas Tamazight). We obtained the data from Wikipron (Lee et al., 2020).

As a soft upper bound, we compare with Charsiu (Zhu et al., 2022) which is a ByT5-small model that has been further pre-trained on 7.2 million examples of grapheme-to-phoneme conversion across 100 languages. Although Charsiu was not exposed to the scripts of the languages we chose, it may have seen related languages whose scripts are encoded similarly in Unicode.

	ban	cop	got	lao	syl	tel	tzm	Avg
Charsiu	68.3	7.8	67.0	35.1	47.6	73.3	18.6	45.4
ByT5	50.2	1.0	30.7	1.9	9.8	6.9	2.7	14.8
Set	53.9	2.2	58.2	5.8	28.2	27.7	6.4	26.1
TE	54.7	1.9	37.0	5.1	30.0	16.2	7.4	21.8
SIP-d4	59.2	6.6	56.5	8.2	39.8	33.1	11.0	30.6
-prefix	55.1	3.2	63.9	7.8	28.0	28.9	7.0	27.7

Table 3: Grapheme-to-phoneme conversion with 100 training examples. We show averages of 5 selections of training examples.

We report accuracies in Table 3, and phoneme-error-rates in Appendix B.2; trends are identical. The original ByT5-small model performs worst on average despite being a strong model for grapheme-to-phoneme conversion in general (Xue et al., 2022). On average across the languages, SIP-d4 outperforms the other methods that pre-train on synthetic data as well as ByT5. The difference between SIP-d4 and Set is statistically significant ($p \approx 4 \times 10^{-4}$, paired approx. permutation test). Fine-tuning SIP-d4 without the tunable prefix leads to a drop in performance, except for Gothic. Charsiu performs very well on Telugu, potentially because of its large overlap in lexicon with Sanskrit (Staal, 1963), which is part of its training data.

6.2 Few-shot text editing

Learning simple text editing tasks (Jane Doe \rightarrow J. Doe) from a handful of examples with a Transformer requires a strong structural inductive bias to overcome competing explanations of the data and hence provides a good benchmark for our approach. While current LLMs may seem like the ideal choice for such tasks, they are prone to hallucinations, e.g. ignoring the input and resorting to frequent entities (see Appendix C for an example).

Text editing has been studied in the context of

	rev-name		sur-initial		FST		Overall	
	Acc \uparrow	ED \downarrow	Acc \uparrow	ED \downarrow	Acc \uparrow	ED \downarrow	Acc \uparrow	ED \downarrow
ByT5	11.8	6.81	47.2	1.76	47.6	1.42	45.7	1.72
Charsiu	43.8	1.73	52.8	0.87	62.4	0.74	60.9	0.80
Set	79.0	1.34	41.5	3.37	68.2	0.71	67.4	0.89
TE	80.3	1.08	88.2	0.41	95.7	0.11	94.5	0.17
SIP-d4	92.4	0.34	97.2	0.10	91.6	0.13	91.9	0.14
-prefix	97.8	0.10	72.6	0.51	89.0	0.27	91.4	0.18

Table 4: Averages of accuracy and edit distance across 5-shot text editing tasks based on 8 draws of training examples. We report results grouped by tasks that cannot be solved by a compact FST (rev-name, sur-initial), tasks that can be solved by FSTs, and overall averages.

program synthesis and we evaluate on 19 such tasks from the SyGuS competition 2017 (Alur et al., 2017). Instead of predicting a program, our model directly operates on input/output examples. We note that 17 of these tasks can be solved by compact FSTs, whereas two cannot. These two tasks are *rev-name* (Jane Doe \rightarrow Doe Jane) and *sur-initial* (John Doe \rightarrow Doe, J.), which require tracking information about the first name in the states.

We report results for 5-shot experiments in Table 4. SIP-d4 and TE excel at this, reaching well above 90% accuracy on average whereas the other methods perform worse by a large margin. Charsiu does not perform clearly better than baselines such as Set – even though it obtains excellent results on grapheme-to-phoneme conversion. Interestingly, TE performs better than SIP-d4 on the tasks that can be solved with FSTs, potentially because the initialization of the prefix for TE follows the same distribution as during pre-training, which is not the case for SIP. However, SIP considerably outperforms TE on the two tasks that cannot be compactly represented by FSTs, suggesting that some of the dynamics acquired during pre-training can sometimes be leveraged in other contexts as well. Fine-tuning SIP-d4 without the tunable prefix leads only to a very small drop in accuracy on average.

7 Analysis: SIP leads to FST simulation

We motivated our approach by the hypothesis that SIP’s pre-training encourages the model to simulate FSTs internally, and that this provides the structural inductive bias. In this section, we present evidence that (i) SIP models indeed approximately simulate FSTs in the hidden states, and (ii) that the dynamics responsible for simulation are re-used after fine-tuning all parameters on input/output pairs only.

For a model to simulate FSTs in its hidden rep-

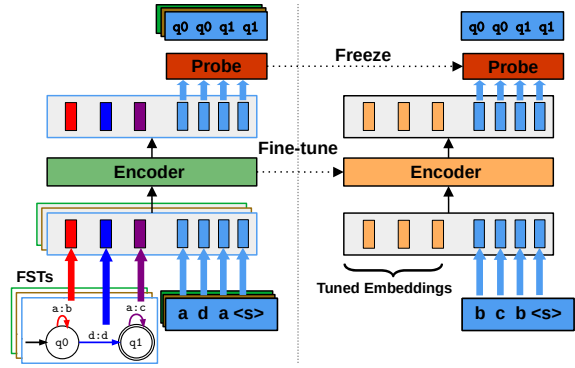


Figure 4: Left: we train a linear probe on the encoder representations of a SIP pre-trained model to predict for each input token x_i which state the encoded FST is in before processing x_i . The end-of-sequence token is represented as $\langle s \rangle$. Right: we freeze the trained probe, fine-tune the SIP model on input/output pairs and extract state sequences from it with the probe.

resentations, it must be able to track the FST state when processing a string, and it should be possible to extract the FST state with a probe. To test this, we mirror the pre-training setup and provide SIP-d4 with an FST and an input string (Fig. 4, left). For each token, we extract the top-layer activations of the encoder, and learn a linear probe with a softmax layer to predict the ID of the state that the given FST is in *before* processing that token. Since state IDs are largely arbitrary, the probe has to learn to relate the hidden representations to the FST presented in the input.

The probe achieves 99.3% token-level accuracy on a test set with unseen FSTs, and a whole-sequence accuracy of 93.9%. We also evaluate a trivial heuristic that returns a random state that has an appropriate outgoing transition for each token in the input. This heuristic achieves a token-level accuracy of 68.9%, and a whole-sequence accuracy of only 17.8%. A probe trained on ByT5 representations, i.e. before SIP pre-training, performs even worse at 42.9% token-level accuracy and whole-sequence accuracy of only 7.1% (see Appendix D.2). Hence, the model has learned a non-trivial way to simulate transitions between states of the FST encoded in the prefix. This is remarkable because the pre-training procedure for SIP-d4 does not provide supervision for *how* to process strings.

While this shows that SIP leads to the simulation of state transitions after pre-training, does the model leverage the simulation ability on downstream tasks? Recall that we fine-tune all parameters of the model (Section 4.2), so the model could

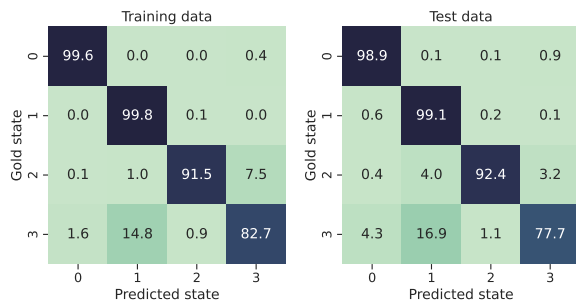


Figure 5: Row-normalized confusion matrices on the training and test data between ground truth and the state predicted by the frozen probe applied to fine-tuned models. We average across the 5 iteration generalization tasks (Section 5.3).

employ a very different strategy to fit the data. To investigate this, we set the trained probe aside and freeze it (Fig. 4, right). We then fine-tune SIP-d4 on the iteration generalization tasks with 4 states in the gold FST (cf. Table 1). Finally, we apply the frozen probe to the fine-tuned model to see if the state sequences we extract are similar to those of the ground truth FST. Fine-tuning SIP-d4 could induce the same FST as the ground truth but use a different numbering of the states. To account for this, for each of the five tasks, we find the isomorphism between the predicted state IDs and the ground truth that gives the best match on average.

The results are presented in Fig. 5 as confusion matrices between predicted and gold states. The probe extracts state sequences that resemble the state sequences of the gold FST (up to isomorphism), both on the training data and the out-of-distribution test data. We also find that deviation from the ground truth state sequence correlates with errors by the fine-tuned model: if the probe extracts correct state sequences, the model achieves an accuracy of 98.6% on the iteration generalization tasks, whereas it drops to 89.8% when the probe extracts state sequences that deviate. The difference is statistically significant (approximate permutation test, $p \approx 5 \times 10^{-5}$). Overall, this shows that the fine-tuned model reuses the dynamics for state tracking and learns representations similar to the ground truth FST.

8 Conclusion

We present SIP, a simple, efficient and adjustable method for introducing a structural inductive bias into a seq2seq model. We focus on an inductive bias towards FSTs, one of the simplest computational devices that is useful for NLP applications. We achieve this by pre-training a Transformer to

simulate automatically generated FSTs, i.e. to predict the output of an FST given an input string and a description of the FST. Our experiments show that our method imparts the desired inductive bias, resulting in improved systematic generalization and better few-shot learning for FST-like tasks. In addition, we show with probing experiments that a model trained with SIP simulates transitions between FST states in its hidden representations, and that the dynamics behind this are leveraged during fine-tuning. In future work, we plan to extend this methodology to more expressive formalisms such as Pushdown Transducers which can be used for a wider range of downstream NLP tasks.

Limitations

Our investigation focuses on FSTs with a relatively small number of states. However, the results in Section 5.4 and in the experiments on grapheme-to-phoneme conversion show that even pre-training with FSTs with a small number of states has positive impacts for tasks that require larger or more complex FSTs.

The probing experiments show that the model simulates transitions between states similarly to an FST, but we did not perform a mechanistic interpretation of how exactly this is implemented in the weights. One potential mechanism behind the simulation behaviour is the construction of Liu et al. (2022) who show that Transformer decoders can simulate transitions between states of deterministic finite automata for strings of length up to n using $O(\log(n))$ layers.

Acquiring a specific inductive bias by means of learning to simulate a computational device is a general idea that could be applicable beyond FSTs but might be unsuitable in cases where (i) it is difficult to formulate a reasonable computational device to simulate (such as document classification and sentiment analysis beyond keyword spotting), or (ii) the computational device would be very hard or infeasible to simulate (e.g. Turing machines).

Our experiments focus on moderately sized models (300M parameters) with an encoder-decoder architecture, and we did not investigate large decoder-only models. Our methodology can also be applied to decoder-only models, and we do not foresee any reasons why it could be less effective in that setup.

Finally, we only consider the standard Transformer architecture and we leave it to future work to explore the impact of SIP on variants of the Transformer architecture designed for handling

long character sequences (Yu et al., 2023) or in the context of state-space models (Wang et al., 2024).

Acknowledgements

We thank Verna Dankers, Victor Prokhorov, and Christine Schäfer for discussions and comments. ML is supported by the UKRI Centre for Doctoral Training in Natural Language Processing, funded by the UKRI (grant EP/S022481/1), the University of Edinburgh, School of Informatics and School of Philosophy, Psychology & Language Sciences, and a grant from Huawei Technologies. IT is supported by the Dutch National Science Foundation (NWO Vici VI.C.212.053).

References

- Cyril Allauzen, Michael Riley, Johan Schalkwyk, Wojciech Skut, and Mehryar Mohri. 2007. Openfst: A general and efficient weighted finite-state transducer library: (extended abstract of an invited talk). In *Implementation and Application of Automata: 12th International Conference, CIAA 2007, Prague, Czech Republic, July 16-18, 2007, Revised Selected Papers 12*, pages 11–23. Springer.
- Rajeev Alur, Dana Fisman, Rishabh Singh, and Armando Solar-Lezama. 2017. Sygus-comp 2017: Results and analysis. *arXiv preprint arXiv:1711.11438*.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- David Bieber, Rishabh Goel, Dan Zheng, Hugo Larochelle, and Daniel Tarlow. 2022. [Static prediction of runtime errors by learning to execute programs with external resource descriptions](#). In *Deep Learning for Code Workshop*.
- Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jian, Bill Yuchen Lin, Peter West, Chandra Bhagavatula, Ronan Le Bras, Jena D Hwang, et al. 2023. [Faith and fate: Limits of transformers on compositionality](#). *arXiv preprint arXiv:2305.18654*.
- Catherine Finegan-Dollak, Jonathan K. Kummerfeld, Li Zhang, Karthik Ramanathan, Sesh Sadasivam, Rui Zhang, and Dragomir Radev. 2018. [Improving text-to-SQL evaluation methodology](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 351–360, Melbourne, Australia. Association for Computational Linguistics.
- Matthew Finlayson, Kyle Richardson, Ashish Sabharwal, and Peter Clark. 2022. [What makes instruction learning hard? an investigation and a new challenge in a synthetic environment](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 414–426, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR.
- Daniel Furrer, Marc van Zee, Nathan Scales, and Nathanael Schärli. 2020. Compositional generalization in semantic parsing: Pre-training vs. specialized architectures. *arXiv preprint arXiv:2007.08970*.
- Yinuo Guo, Zeqi Lin, Jian-Guang Lou, and Dongmei Zhang. 2020. Hierarchical poset decoding for compositional generalization in language. *Advances in Neural Information Processing Systems*, 33:6913–6924.
- Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia Bruni. 2020. [Compositionality decomposed: how do neural networks generalise?](#) *Journal of Artificial Intelligence Research*, 67:757–795.
- Adam Karvonen. 2024. [Emergent world models and latent variable estimation in chess-playing language models](#).
- Daniel Keysers, Nathanael Schärli, Nathan Scales, Hylke Buisman, Daniel Furrer, Sergii Kashubin, Nikola Momchev, Danila Sinopalnikov, Lukasz Stafiniak, Tibor Tihon, Dmitry Tsarkov, Xiao Wang, Marc van Zee, and Olivier Bousquet. 2020. [Measuring compositional generalization: A comprehensive method on realistic data](#). In *International Conference on Learning Representations*.
- Najoung Kim and Tal Linzen. 2020. [COGS: A compositional generalization challenge based on semantic interpretation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9087–9105, Online. Association for Computational Linguistics.
- Yoon Kim. 2021. [Sequence-to-sequence learning with latent neural grammars](#). In *Advances in Neural Information Processing Systems*, volume 34, pages 26302–26317. Curran Associates, Inc.
- Kundan Krishna, Jeffrey Bigham, and Zachary C. Lipton. 2021. [Does pretraining for summarization require knowledge transfer?](#) In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3178–3189, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Brenden Lake and Marco Baroni. 2018. [Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks](#). In *International Conference on Machine Learning*, pages 2873–2882. PMLR.

- Jackson L. Lee, Lucas F.E. Ashby, M. Elizabeth Garza, Yeonju Lee-Sikka, Sean Miller, Alan Wong, Arya D. McCarthy, and Kyle Gorman. 2020. [Massively multilingual pronunciation modeling with WikiPron](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4223–4228, Marseille, France. European Language Resources Association.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. [Emergent world representations: Exploring a sequence model trained on a synthetic task](#). In *The Eleventh International Conference on Learning Representations*.
- Yafu Li, Yongjing Yin, Yulong Chen, and Yue Zhang. 2021. [On compositional generalization of neural machine translation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4767–4780, Online. Association for Computational Linguistics.
- Matthias Lindemann, Alexander Koller, and Ivan Titov. 2023. [Compositional generalization without trees using multiset tagging and latent permutations](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14488–14506, Toronto, Canada. Association for Computational Linguistics.
- Bingbin Liu, Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. 2022. [Transformers learn shortcuts to automata](#). *arXiv preprint arXiv:2210.10749*.
- R Thomas McCoy, Erin Grant, Paul Smolensky, Thomas L Griffiths, and Tal Linzen. 2020. [Universal linguistic inductive biases via meta-learning](#). In *Proceedings of the 42nd Annual Conference of the Cognitive Science Society*.
- R Thomas McCoy and Thomas L Griffiths. 2023. [Modeling rapid language learning by distilling bayesian priors into artificial neural networks](#). *arXiv preprint arXiv:2305.14701*.
- Stoyan Mihov and Klaus U. Schulz. 2019. *Finite-State Techniques: Automata, Transducers and Bimachines*. Cambridge Tracts in Theoretical Computer Science. Cambridge University Press.
- Isabel Papadimitriou and Dan Jurafsky. 2023. [Injecting structural hints: Using language models to study inductive biases in language learning](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8402–8413, Singapore. Association for Computational Linguistics.
- Kexin Pei, Jonas Guan, Matthew Broughton, Zhongtian Chen, Songchen Yao, David Williams-King, Vikas Ummadisetty, Junfeng Yang, Baishakhi Ray, and Suman Jana. 2021. [Stateformer: Fine-grained type recovery from binaries using generative state modeling](#). In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2021*, page 690–702, New York, NY, USA. Association for Computing Machinery.
- Linlu Qiu, Peter Shaw, Panupong Pasupat, Tianze Shi, Jonathan Herzig, Emily Pitler, Fei Sha, and Kristina Toutanova. 2022. [Evaluating the impact of model scale for compositional generalization in semantic parsing](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9157–9179, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *J. Mach. Learn. Res.*, 21(1).
- Marcel-Paul Schützenberger. 1961. [A remark on finite transducers](#). *Information and Control*, 4:185–196.
- Richard Sinkhorn. 1964. [A relationship between arbitrary positive matrices and doubly stochastic matrices](#). *The Annals of Mathematical Statistics*, 35(2):876–879.
- J. F. Staal. 1963. [Sanskrit and sanskritization](#). *The Journal of Asian Studies*, 22(3):261–275.
- Sara Stymne, Miryam de Lhoneux, Aaron Smith, and Joakim Nivre. 2018. [Parser training with heterogeneous treebanks](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 619–625, Melbourne, Australia. Association for Computational Linguistics.
- Yulia Tsvetkov, Sunayana Sitaram, Manaal Faruqui, Guillaume Lample, Patrick Littell, David Mortensen, Alan W Black, Lori Levin, and Chris Dyer. 2016. [Polyglot neural language models: A case study in cross-lingual phonetic representation learning](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1357–1366, San Diego, California. Association for Computational Linguistics.
- Junxiong Wang, Tushaar Gangavarapu, Jing Nathan Yan, and Alexander M Rush. 2024. [Mambabyte: Token-free selective state space model](#). *arXiv preprint arXiv:2401.13660*.
- Shijie Wu and Ryan Cotterell. 2019. [Exact hard monotonic attention for character-level transduction](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1530–1537, Florence, Italy. Association for Computational Linguistics.

Yuhuai Wu, Felix Li, and Percy S Liang. 2022. [Insights into pre-training via simpler synthetic tasks](#). *Advances in Neural Information Processing Systems*, 35:21844–21857.

Yuhuai Wu, Markus N Rabe, Wenda Li, Jimmy Ba, Roger B Grosse, and Christian Szegedy. 2021. Lime: Learning inductive bias for primitives of mathematical reasoning. In *International Conference on Machine Learning*, pages 11251–11262. PMLR.

Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2022. Byt5: Towards a token-free future with pre-trained byte-to-byte models. *Transactions of the Association for Computational Linguistics*, 10:291–306.

Lili Yu, D’aniel Simig, Colin Flaherty, Armen Aghajanyan, Luke Zettlemoyer, and Mike Lewis. 2023. [Megabyte: Predicting million-byte sequences with multiscale transformers](#). *arXiv preprint arXiv:2305.07185*.

Wojciech Zaremba and Ilya Sutskever. 2014. Learning to execute. *arXiv preprint arXiv:1410.4615*.

Zhuosheng Zhang, Shuohang Wang, Yichong Xu, Yuwei Fang, Wenhao Yu, Yang Liu, Hai Zhao, Chenguang Zhu, and Michael Zeng. 2022. [Task compass: Scaling multi-task pre-training with task prefix](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5671–5685, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Hao Zheng and Mirella Lapata. 2021. [Compositional generalization via semantic tagging](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1022–1032, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Jian Zhu, Cong Zhang, and David Jurgens. 2022. [ByT5 model for massively multilingual grapheme-to-phoneme conversion](#). In *Proc. Interspeech 2022*, pages 446–450.

A Generation of Synthetic Data and Splits

A.1 Generating deterministic FSTs

Before describing our procedure for sampling deterministic FSTs, we briefly establish notation. An FST is a tuple $\langle Q, \Sigma, \Gamma, I, F, \Delta \rangle$, where Q is a finite set of states, Σ is the input alphabet, Γ is the output alphabet, $I \subseteq Q$ is a set of initial states, $F \subseteq Q$ is a set of final states and $\Delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times (\Gamma \cup \{\epsilon\}) \times Q$ are the transitions. We assume $\Sigma = \Gamma$ and call it V for vocabulary.

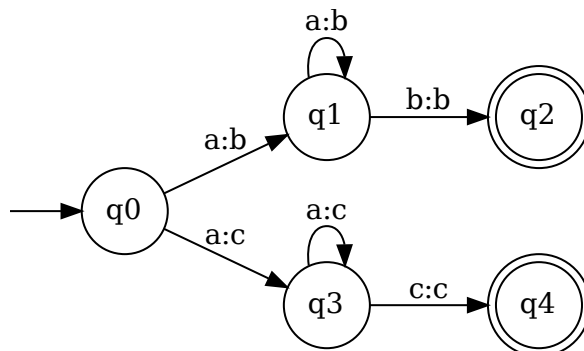


Figure 6: A functional but non-deterministic FST.

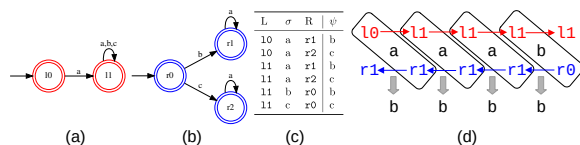


Figure 7: (a) - (c) shows a bimachine that is equivalent to Fig. 6. (a) Left automaton A^l , (b) Right automaton A^r , (c) output function ψ . (d) shows an example run of the bimachine on the input aaab which is mapped to bbbb.

For technical reasons, we exclude the three characters $[,]$ and \backslash from the vocabulary as they are interpreted as special characters by OpenFST, which we use for constructing and representing FSTs.

In addition to the shorthand for identity transitions (id), we also have shorthands for converting upper case to lower case and vice-versa (lower-to-upper, upper-to-lower). We describe our procedure to generate a deterministic FST with pseudocode in Algorithm 1. It receives as argument n (the number of states in the FST), f (number of final states), V (the vocabulary of this FST), and probabilities P-ID, P-DROP, P-SHORTHAND. These probabilities control the likelihood of using a shorthand, not drawing an outgoing edge (P-DROP) with a given symbol, and creating a single identity transition (P-ID). We use CHOICE to denote a uniform random choice from a finite set.

We use P-ID = 0.2, P-DROP = 0.4, P-SHORTHAND = 0.15 in our experiments.

A.2 Generating Non-deterministic Functional FSTs

It is not straightforward to directly generate non-deterministic FSTs that are guaranteed to express a function. However, we can directly generate a bimachine, which then can be converted into an FST.

Algorithm 1 Generate a random deterministic FST

```
function GEN-DET-FST( $n, f, V, P\text{-ID}, P\text{-DROP}, P\text{-SHORTAND}$ )
   $Q = \{0, \dots, n-1\}$ 
   $\Delta = \emptyset$ 
   $I = \{0\}$ 
  for  $q \in Q$  do
     $q' = \text{CHOICE}(Q)$ 
    with prob  $P\text{-SHORTAND}$ 
       $s = \text{CHOICE}([\text{id}, \text{lower-to-upper}, \text{upper-to-lower}])$ 
       $\Delta := \Delta \cup \{q \xrightarrow{s:s} q'\}$ 
    else
      for  $\sigma \in V$  do
        with prob  $P\text{-DROP}$ 
          no-op  $\triangleright$  No outgoing edge with  $\sigma$  at  $q$ 
        else with prob  $P\text{-ID}$ 
           $\Delta := \Delta \cup \{q \xrightarrow{\sigma:\sigma} q'\}$ 
        else
           $\Delta := \Delta \cup \{q \xrightarrow{\sigma:\text{CHOICE}(V \cup \{\epsilon\})} q'\}$ 
        end with prob
      end for
    end with prob
  end for
  Eliminate states from  $Q$  through which no accepting path can go
  Choose random subset  $F$  of  $Q$  with  $|F| = \min(f, |Q|)$ 
  return minimized FST with states  $Q$ , transitions  $\Delta$ , initial states  $I$  and final states  $F$ 
end function
```

Bimachines (Schützenberger, 1961) represent the functions expressible by FSTs, i.e. for every functional FST there is a bimachine that represents it (and vice-versa). A bimachine consists of two deterministic finite state automata (called left and right) and an output function. Let A^L be the left FSA with states Q^L and transition function $\delta^L : Q^L \times \Sigma \rightarrow Q^L$, and let A^R be the right FS with states Q^R and transition function $\delta^R : Q^R \times \Sigma \rightarrow Q^R$. The output function is $\psi : Q^L \times \Sigma \times Q^R \rightarrow \Gamma^*$. All states of A^L and A^R are final states. Given an input string $x = \sigma_1\sigma_2\sigma_3 \dots \sigma_n$, a bimachine runs A^L from left to right over x , keeping track of the states $q_0^l, q_1^l, q_2^l, \dots, q_n^l$. It also runs A^R over the string x but this time from right to left, again keeping track of the states $q_0^r, q_1^r, q_2^r, \dots, q_n^r$ that are visited. Then, the state sequence of the right automaton is reversed and ψ is

Algorithm 2 Generate output function for bimachine

```
function GEN-OUTPUT- $\psi(n^L, n^R, V, P\text{-ID} = 0.2)$ 
  for  $q^L \in 0, \dots, n^L - 1$  do
    for  $q^R \in 0, \dots, n^R - 1$  do
      for  $\sigma \in V$  do
        with prob  $P\text{-ID}$ 
           $\psi(q^L, \sigma, q^R) := \sigma$ 
        else
           $\psi(q^L, \sigma, q^R) := \text{CHOICE}(V \cup \{\epsilon\})$ 
        end with prob
      end for
    end for
  end for
  return  $\psi$ 
end function
```

applied ‘elementwise’ as illustrated in Fig. 7. More formally, the output of the bimachine is $\psi(q_0^l, \sigma_1, q_{n-1}^r)\psi(q_1^l, \sigma_1, q_{n-2}^r) \dots \psi(q_{n-1}^l, \sigma_1, q_0^r)$.

Bimachines can be compiled into FSTs with a simple product construction. For a bimachine $\langle A^L, A^R, \psi \rangle$, one can construct an equivalent FST as follows:

$$\langle Q^L \times Q^R, \Sigma, \Gamma, \{s^L\} \times Q^R, Q^L \times \{s^R\}, \Delta \rangle$$

where s^L and s^R are initial states of A^L and A^R , and Δ contains all transitions

$$\Delta = \{ \langle q^L, q^R \rangle \xrightarrow{\sigma:\rho} \langle q'^L, q'^R \rangle \mid \delta^L(q^L, \sigma) = q'^L, \delta^R(q'^R, \sigma) = q^R, \rho = \psi(q^L, \sigma, q'^R) \}$$

We refer to Mihov and Schulz (2019) for details and further information about bimachines.

To sample bimachines, we re-use Algorithm 1 with $P\text{-SHORTAND} = 0$, and ignore the outputs of the transitions, treating them as FSAs. We sample the output function according to Algorithm 2. For the test data creation (Table 2), we use 5 states in the left FSA and 4 states in the right FSA, and set $P\text{-DROP} = 0.4$. For creating the training data for SIP-nd7, we use 2 or 3 states in either left or right automaton and set $P\text{-DROP} = 0.6$ to keep the length of the prefix low to save GPU memory.

A.3 Unseen Combinations of Transitions

We now describe the construction by which we create training and test data for the evaluation of Unseen Combinations of transitions. We first describe how we construct an FST for the training and

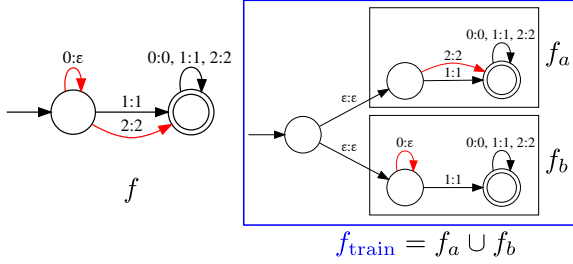


Figure 8: Constructing training data for evaluating unseen combinations of transitions. Based on the given FST f , we construct an FST f_{train} that withholds the *combination* of the two red transitions.

test data, respectively, given a choice of transitions whose combination we want to withhold. Then, we briefly describe how those transitions are chosen.

Given an FST f (illustration in Fig. 8, left) and transitions t_a and t_b (highlighted in red) whose combination we want to withhold, we construct a new FST f_{train} as follows: We create two copies f_a, f_b of the original FST f . In f_a , we remove the transition t_b ; in f_b , we remove the transition t_a . Then $f_{\text{train}} = f_a \cup f_b$, which can be constructed by introducing a new initial state with ϵ -transitions into the respective initial states of f_a and f_b (right side of Fig. 8). This ensures that any accepting path goes through f_a or f_b but cannot alternate between the two. Hence, t_a or t_b can be used – but not both in the same string. Note that f_{train} still describes a partial function (rather than a relation) because any accepting path in f_a and any accepting path in f_b is also an accepting path in f . As a result, whenever f_a and f_b are both defined, they agree on the result $f_a(x) = f_b(x) = f(x)$. We test exclusively for how a model handles unseen combinations of transitions by generating examples from f for which f_{train} is *not* defined.

To make the generalization setup more challenging, these steps can be applied to multiple pairs of adjacent transitions at the same time, i.e. to withhold $\langle t_a^1, t_b^1 \rangle, \dots, \langle t_a^k, t_b^k \rangle$: We create the copy f_a and remove the transitions t_b^1, \dots, t_b^k from f_a and analogously remove t_a^1, \dots, t_a^k from f_b .

Now, we briefly describe how we select *which* pairs of transitions we want to withhold. We only select adjacent transitions, i.e. transitions where one can be used immediately after the other, excluding self-loops. In addition, some transitions cannot be deleted without cutting off a vital initial or final state, which can lead to f_{train} being undefined for any string (and hence no training data). We ensure this never happens by never withhold-

Num. states	Split	Min	Max	Mean
4	train	2	11	4.66
4	test	4	30	18.97
5	train	2	14	5.39
5	test	4	30	19.53
7	train	2	20	6.12
7	test	4	30	20.13
10	train	2	25	7.31
10	test	4	30	20.62
21	train	2	30	11.80
21	test	5	30	23.07

Table 5: Distribution of input lengths of the train/test data we generate for the iteration generalization experiments in Section 5. The tasks with 21 states are the non-deterministic FSTs from Section 5.5.

ing the first transition into each state based on a depth-first traversal of the FST.

A.4 Additional dataset information

For all experiments with synthetic data (generated by FSTs), we generate 5000 training examples and 1000 test examples. To reduce variance across tasks, we fix the vocabulary size to its maximum value (25) in the pre-training data and choose the vocabulary only from the printable ASCII characters.

Length distribution. The input strings in the pre-training data we generate for SIP-d4 have a minimum length of 1, an average length of 15.57 and a maximum length of 35. We report the length distributions for the iteration generalization experiments in Section 5 in Table 5.

SyGuS. We took the data from the SyGuS competition [github https://github.com/SyGuS-Org/benchmarks/tree/master/comp/2017/PBE_Strings_Track](https://github.com/SyGuS-Org/benchmarks/tree/master/comp/2017/PBE_Strings_Track), and extracted the ‘constraints’. For each text editing tasks, there are usually three files, e.g. `firstname`, `firstname-long`, `firstname-long-repeat`. We only consider data from the `*-long` variant because the non-marked variant (e.g. `firstname`) is a subset of the `*-long` variant, and we exclude `*-long-repeat` as it contains repeated data points. We also exclude some text editing tasks that have insufficient amounts of data for reliable evaluation (`bikes`) and some tasks where the input is not a single string but a pair of strings if concatenating the strings results in particularly long inputs (`univ`), or if the concatenation of the

string pair makes the task trivial (name-combine, which would correspond to an identity operation). For a few-shot experiment, we sample 5 training examples and evaluate on the rest.

We note that the original intention in the design of the benchmark data was for program synthesis rather than few-shot learning. The data contains names, and separately it contains phone numbers (but not combined). However, we believe both to be synthetically generated.

Grapheme-to-phoneme. We obtain data from Lee et al. (2020), and conduct experiments mainly on the broad transcription, except for Telugu and Tamazight, where we use the narrow transcription. For each experiment, we randomly sample 100 training examples, and use the rest as test data. The data is available under a permissible license: <https://en.wiktionary.org/wiki/Wiktionary:Copyrights>

B Additional Results

B.1 Additional Results with More States

In Fig. 3, we show accuracy relative to the accuracy of ByT5. Here, we show the absolute accuracies and edit distances in Table 6.

B.2 Full results for grapheme-to-phoneme conversion

Table 7 shows the full results of our grapheme-to-phoneme conversion experiments, including phoneme error rate (PER).

B.3 Additional results with T5-Base

We run a subset of the experiments starting off from a pre-trained T5-Base (Raffel et al., 2020) instead of ByT5. This model is about one-third smaller than ByT5 (around 200 million instead of 300 million parameters). T5-Base uses a different vocabulary than ByT5, so we resize the output layer to the vocabulary size of ByT5 and re-initialize it. For the input embeddings, we re-purpose the first n embeddings in the T5-Base embedding matrix to represent the token ids according to the ByT5 tokenizer. While this is suitable as a starting point for further pre-training, we found that directly fine-tuning T5-Base with these modifications on downstream tasks led to very poor results and do not include them here. Instead, we train T5-Set (analogous to Set) for a fair point of comparison.

We report a subset of the results from the main paper in for T5-Base in Tables 8 to 10.

We also tried to pre-train a ByT5-style model from scratch (i.e. from random initialization). However, we could not find a setting of hyperparameters that would make the model converge well. We hypothesize that the model already needs to be in a reasonable space to make learning feasible.

B.4 Generalization to longer strings

In the main paper, we report results on iteration generalization where a model is trained on strings such that each state has been visited at most 3 times, and is tested on strings where at least one state is visited at least 4 times. Here, we explore a more extreme version, where there is a large gap between the maximum length seen during training and the minimum length seen during testing. As another point of comparison, we further pre-train SIP-d4 on 40,000 FSTs with strings of length up to 110 (SIP-d4-long).

We report results in Table 11. ByT5 struggles with this generalization setup across the board. SIP-d4 performs remarkably well on lengths 40-70 which are beyond the lengths seen during its pre-training. However, performance drops starkly when testing on inputs of length 90 to 110. We hypothesize that this is because the relevant positional embeddings were not pre-trained by SIP. In contrast, SIP-d4-long performs well on inputs of length 90 to 110, as it has seen strings of such length during pre-training.

C Hallucination Example

We briefly show an example where an LLM ignores a part of the input and resorts to outputting a high-frequency entity. Consider the following in-context examples for a simple text editing task:

Input	Output
Howard Phillips Lovecraft	H.P. Lovecraft
John Ronald Reuel Tolkien	J.R.R. Tolkien
Thomas Stearns Eliot	T.S. Eliot

At the time of submission, the current version of ChatGPT frequently outputs “J.K. Rowling” for the name “John Edward Rowling”, hallucinating the K.

D Additional Analysis

D.1 Analysis of fine-tuned prefixes

To gain some understanding of how the prefix of tunable embeddings is used by the model and what it contains, we consider the setup of fine-tuning

Gen. Type	Num States Model	4		5		7		10	
		Acc↑	ED↓	Acc↑	ED↓	Acc↑	ED↓	Acc↑	ED↓
Iteration	ByT5	37.8	5.87	58.7	3.21	48.2	3.71	45.7	3.87
	Naive	42.6	4.41	60.5	2.20	47.7	3.16	43.6	3.65
	Set	44.4	4.58	62.2	2.41	48.0	3.49	45.3	3.71
	TE	61.3	2.49	78.9	0.86	55.7	2.29	50.7	2.95
	SIP-d4	94.8	0.12	89.6	0.27	64.3	1.34	56.9	2.39
UC	ByT5	47.4	1.49	62.6	1.05	61.9	1.29	54.1	1.70
	Naive	44.9	1.52	61.6	1.08	59.3	1.30	51.8	1.68
	Set	43.6	1.47	60.6	1.09	60.8	1.31	51.1	1.71
	TE	57.3	1.13	65.9	0.98	65.7	1.17	55.3	1.60
	SIP-d4	73.1	0.61	74.3	0.69	73.2	0.85	58.0	1.44

Table 6: Evaluation on deterministic FSTs with more states, showing absolute accuracies and edit distances, corresponding to Fig. 3 and ??.

	ban		cop		got		lao		syl		tel		tzm		Avg	
	Acc	PER	Acc	PER	Acc	PER	Acc	PER	Acc	PER	Acc	PER	Acc	PER	Acc↑	PER↓
Charsiu	68.3	.110	7.8	.579	67.0	.067	35.1	.238	47.6	.196	73.3	.070	18.6	.403	45.4	.238
ByT5	50.2	.233	1.0	.847	30.7	.269	1.9	.760	9.8	.598	6.9	.597	2.7	.851	14.8	.594
Set	53.9	.216	2.2	.742	58.2	.094	5.8	.595	28.2	.353	27.7	.293	6.4	.658	26.1	.421
TE	54.7	.183	1.9	.756	37.0	.174	5.1	.573	30.0	.309	16.2	.377	7.4	.644	21.8	.431
SIP-d4	59.2	.152	6.6	.563	56.5	.096	8.2	.498	39.8	.252	33.1	.228	11.0	.544	30.6	.333
-prefix	55.1	.168	3.2	.681	63.9	.072	7.8	.508	28.0	.333	28.9	.252	7.0	.593	27.7	.372

Table 7: Grapheme-to-phoneme conversion with 100 training examples. We show averages of 5 selections of training examples. PER is Phoneme Error Rate: edit distance / length of gold output (lower is better).

only the prefix and keeping the rest of the model unchanged. That is, all the task-specific information has to be captured in these embeddings. Specifically, we fine-tune on the 5 FSTs from Section 5.3 for iteration generalization for 20 epochs with a learning rate of 0.5.

We explore two questions:

1. Is the model robust towards different permutations of the fine-tuned prefixes? Intuitively, these permutations correspond to changing the order in which transitions are listed, so ideally the model should not be sensitive to that order.
2. Does the fine-tuned prefix represent the task-specific information in a similar way to how FSTs were encoded during pre-training?

To address the first question, we randomly permute the tuned prefixes and compute accuracy on the iteration generalization data before and after permuting the tuned prefixes. We use 20 permutations per learned prefix and average results across the 5

FSTs. Overall, we find that this results only in a small drop in accuracy: the median drop in accuracy is only around 1.3 percentage points, and the arithmetic mean of the drop is around 7.1 percentage points. Most permutations do not have a big impact on how the prefix is interpreted but a few permutations do have a stronger negative impact, skewing the arithmetic mean.

To address the second question, we test if the learned prefix for a task t resembles an encoding of an FST that solves t . For each of the 5 FSTs, we generate 10,000 distractors, i.e. FSTs that have the same number of states and use the same vocabulary as the FST solving t . We define the similarity of two prefixes p, q as follows:

$$sim(p, q) = \max_{\pi} \frac{1}{n} \sum_i \frac{p_i^T q_{\pi(i)}}{\|p_i\|_2 \cdot \|q_{\pi(i)}\|_2}$$

where π is a permutation, and p_i is the i -th vector in prefix p , and prefixes p and q both have length n . That is, we define the similarity between p and q as the highest possible average cosine similarities between positions in p and q that one can achieve

	Iteration		UC	
	Acc↑	ED↓	Acc↑	ED↓
T5-Set	26.6	6.26	55.1/54.6	1.18/1.02
T5-SIP-d4	94.5	0.11	75.4/99.5	0.54/0.01

Table 8: Evaluating systematic generalization on FST tasks with 4 states (cf. Table 1). Due to an outlier task on UC, we additionally report the median after ‘/’.

	Iteration		UC	
	Acc↑	ED↓	Acc↑	ED↓
T5-Set	77.9	0.73	81.7	0.53
T5-SIP-d4	83.3	0.56	86.1	0.37

Table 9: Evaluation with T5-Base on non-deterministic FSTs (cf. Table 2)

	ban		cop		got		lao		syl		tel		tzm		Avg	
	Acc	PER	Acc	PER	Acc	PER	Acc	PER	Acc	PER	Acc	PER	Acc	PER	Acc↑	PER↓
T5-Set	47.9	.231	1.2	.783	6.7	.458	3.6	.643	6.6	.611	4.9	.612	2.7	.797	10.5	.591
T5-SIP-d4	59.1	.154	4.7	.640	69.6	.059	5.9	.566	22.1	.447	35.4	.191	12.5	.509	29.9	.367

Table 10: Grapheme-to-phoneme conversion with 100 training examples based on T5-Base. In contrast to the experiments in the main paper, we found that T5-SIP-d4 did not perform well on completely unseen scripts, so we mapped all Unicode code points to arbitrary ASCII characters. This maintains the structure of the task and is completely reversible. T5-Set is evaluated in the same way.

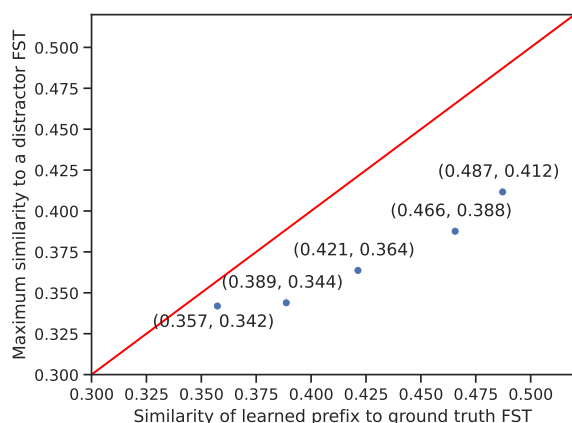


Figure 9: Each dot represents a fine-tuned prefix when the rest of the model remains frozen during fine-tuning. The x-coordinates represent the similarity to a ground truth gold prefix, and the y-coordinates represent the maximum similarity to any of the 5×10000 distractor FSTs. All dots are below the diagonal, hence all learned prefixes are most similar to an encoding of the ground truth FST.

by assigning a position in p to exactly one position in q and vice-versa.² Taking the maximum over all permutations is justified by our results to the first question above, which showed that the model is largely invariant to different permutations of the

²Computing the similarity $sim(p, q)$ is relatively expensive because it involves solving the assignment problem (e.g. with the Hungarian algorithm). Instead of solving the assignment problem exactly, we approximate it with the Sinkhorn algorithm (Sinkhorn, 1964). We then take the output of the algorithm (a matrix of ‘soft’ assignments) and for each position in p , we greedily select a matching position in q .

tuned prefix.

For every task t , we compute the similarity between the prefix p learned by fine-tuning on input/output pairs and the union of encodings of the distractors and encodings of the gold standard FST for task t . Where necessary, we truncate encodings of FSTs to have the same length as the learned prefix. We present the results in Fig. 9 showing that all learned prefixes are most similar to an encoding of the ground truth FST.

D.2 Probing non-SIP models

All probes are trained for one epoch on activations produced by passing 8,000 FSTs with 5 inputs each (i.e. 40,000 instances) through the model.

For the baseline probe, we take the trained SIP-d4 model (including matrix W and embeddings from 4.1) and re-initialize the Transformer to ByT5-small. The probe achieves only a token-level accuracy of 42.9% and whole-sequence accuracy of 8.1%. We see very similar results for a probe trained on a randomly initialized Transformer in this setup: a token-level accuracy of 42.5% and a whole-sequence accuracy of 7.1%.

E Additional model details & Hyperparameters & Hardware

SIP. For completeness, we now describe the order in which we arrange the transitions. While the ordering of the transitions does not matter for expressing FSTs, the Transformer uses positional encodings which might have impacts on the pre-

Test length	Model	pre-train length	Acc \uparrow	ED \downarrow
40 to 70	ByT5	1024	29.3	15.60
	SIP-d4	35	69.4	2.61
90 to 110	ByT5	1024	1.4	55.37
	SIP-d4	35	3.4	34.50
	SIP-d4-long	110	81.5	1.09

Table 11: Average generalization ability across 5 FSTs with 4 states. Models were trained on inputs of length up to 15, and tested on much longer inputs.

training (though see Appendix D.1). We assemble the overall prefix by stacking the individual vectors h_0, \dots, h_n of the transitions $p_0 \xrightarrow{\sigma_0: \rho_0} q_0, \dots, p_n \xrightarrow{\sigma_n: \rho_n} q_n$. We group the transitions by their originating state (i.e. p_i) and go over the states by their id, starting with 0, the initial state.

During pre-training, we might encounter FSTs with different numbers of transitions within the same batch. To handle this, we use padding encodings by reserving a special padding state and padding symbol in the embedding matrices of states and symbols. To initialize the prefix for fine-tuning, we use the average of 32 FST encodings (chosen at random) from pretraining.

For pre-training, we use embeddings of dimensionality 64 for states, embeddings of dimensionality 256 for symbols, and of dimensionality 16 to indicate final/non-final states.

Task embeddings. To enable faster adaption of the task embeddings than the rest of the model to fit a particular task, we use a higher learning rate for the task embeddings (1.0) than for the rest of the model ($5 \cdot 10^{-4}$) during pre-training. We also use a higher learning rate for the prefix during fine-tuning, analogously to SIP.

Because we have to store 40,000 task embeddings (one for each generated FST), TE requires a lot of memory. To reduce memory consumption, the task embeddings have a dimensionality of 180 and are up-projected to fit into the Transformer, analogously to W in Section 4.1. Nevertheless, the memory consumption of the embeddings is substantial and we store them on a separate GPU. Analogously to SIP-d4, we pre-train for 20 epochs.

Naive. We pre-train for a single epoch only as we found this achieved better results on downstream tasks than training for 20 epochs.

Set. We sample 200,000 examples according to the procedure described by Wu et al. (2022) to match our pre-training dataset size. Again, we found it more helpful for downstream task perfor-

mance to train for a single epoch rather than 20 epochs.

Fine-tuning Hyperparameters. Like pre-training, we finetune with the Adam optimizer. The main hyperparameters involved for both SIP and TE are the learning rates for the main model, and (separately) the learning rate of the tunable prefix. We chose these manually. Generally, we found that using a learning rate of 1.0 was a good choice for the prefix. Lester et al. (2021) report a similarly high learning rate to be useful for prompt tuning. For the rest of the model, we found $3 \cdot 10^{-4}$ and $5 \cdot 10^{-4}$ to work well for SIP-d4 and TE, respectively. For few-shot experiments, we use a somewhat smaller learning rate for TE for the main model ($3 \cdot 10^{-4}$). We noticed that T5-SIP-d4 (see Appendix B.3) was more sensitive to the learning rate choice in general than SIP-d4.

Hardware/Computational Budget. We ran our experiments on NVIDIA GeForce RTX 2080 Ti GPUs (11264MiB RAM) with driver version 535.54.03 and cuda version 12.2.

Pre-training SIP-d4 took around 30 hours for 20 epochs. One training run on synthetic data (including evaluation) takes around one hour, and one training run for low-resource grapheme-to-phoneme conversion takes between 5 and 10 minutes.