

# WHEN CAN TRANSFORMERS REASON WITH ABSTRACT SYMBOLS?

Anonymous authors

Paper under double-blind review

## ABSTRACT

We investigate the capability of Transformer large language models (LLMs) to generalize on unseen symbols when trained on tasks that rely on abstract symbols (e.g., variables in programming and mathematics). Such a ‘variable-binding’ capability has long been studied in the neuroscience literature as one of the most basic ‘reasoning’ capabilities. For (i) binary classification tasks, we prove that Transformers can generalize to unseen symbols but require astonishingly large training data. For (ii) tasks with labels dependent on input symbols, we show an “inverse scaling law”: Transformers fail to generalize to unseen symbols as their embedding dimension increases. For both cases (i) and (ii), we propose a Transformer modification, adding two trainable parameters per head that can reduce the amount of data needed.

## 1 INTRODUCTION

*Reasoning* can be defined as the ability to use logical rules to generalize outside of one’s training data. During most of the history of AI, reasoning was widely thought to be achievable only through programs that manipulated mathematical symbols using hand-coded logical rules (Newell et al., 1959; Marcus, 1998). However, recent developments have challenged this paradigm: as large language models (LLMs) are trained with increasing quantities of data, they start to exhibit the ability to reason mathematically (Kaplan et al., 2020; Yuan et al., 2023). But why does more data help an LLM to reason outside of its training set? And how efficient can we make LLMs in that regard?

In this paper, we focus on how LLMs learn to reason in tasks involving abstract symbols (known as *variable-binding* tasks in the neuroscience literature). The reasoning capability required for these tasks is basic, but crucial to many domains and has been hypothesized to be necessary for much of human cognition (Fodor, 1975; Newell, 1980; Marcus, 1998; Kriete et al., 2013; Webb et al., 2020b). For example, variable-binding is a building block of mathematics and computer science, where abstract symbols (i.e., variable names) refer to concrete values in a proof or program.

See Figure 1 for an illustrative variable-binding task, where we train an LLM to evaluate Python programs  $x_i$ , and return their output  $y_i$ . Memorizing the training data is easy (Zhang et al., 2021a), but we wish to measure reasoning: will the LLM learn to treat the variable names as abstract symbols, enabling generalization beyond its training dataset? To evaluate this, we adopt an out-of-distribution setting, where the train and test data distributions differ (Marcus, 1998; Abbe et al., 2023). The test dataset consists of the same programs, but with *different variable names never seen during training*. Remarkably, as the training set size increases, the LLM’s ability to reason outside of its training data improves.

In Figure 2, we consider a variable-binding task with one extra layer of complexity: each sample is labeled with a symbol (instead of a real number  $+1$  or  $-1$  as in Figure 1). For the LLM to generalize to symbols unseen at train time, not only must it track the value stored in a variable, but it also must learn to predict symbolic labels at test time that do not occur in its training data. On this more sophisticated task, we observe that a transformer requires much more training data to generalize.

### 1.1 OUR CONTRIBUTIONS

To understand these phenomena, we study a framework of reasoning tasks of which Figures 1 and 2 are special cases. (i) For the real-valued label tasks as in Figure 1, we prove that transformers will

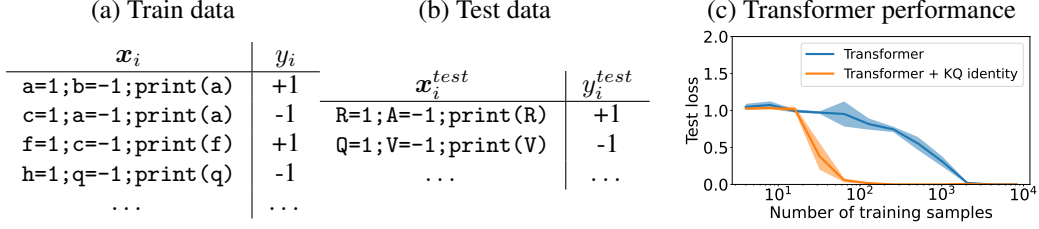


Figure 1: (a,b) Variable names in the test data never appear in the train data (indicated by lower/upper-case names). (c) Our theory motivates a slightly modified transformer architecture (see Observation 1.2), which solves the reasoning task with less training data. Details in Appendix A.

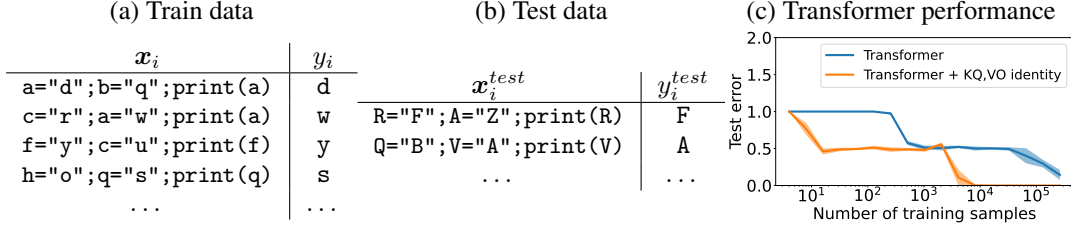


Figure 2: (a,b) A task where labels are also symbols. (c) Our modified transformer learns the reasoning task with less data (see Observation 1.2 and Theorem 1.4). Details in Appendix A.

learn to generalize, but require a large quantity of data. (ii) For the symbolic-label tasks as in Figure 2, we prove that transformers will fail. For settings (i) and (ii) we propose small parametrization adjustments that improve data efficiency and allow for success, respectively. Finally, we support our claims experimentally, and also cast light on how pretraining helps improve reasoning abilities.

#### 1.1.1 TEMPLATE TASKS: A FRAMEWORK FOR REASONING WITH ABSTRACT SYMBOLS

Building on a long line of work in neuroscience (Marcus, 1998; Kim et al., 2018; Webb et al., 2020b), we formalize a framework of reasoning tasks called *template tasks*. These tasks come in two kinds: *real-label* as in Figure 1, and *symbolic-label* as in Figure 2.

**Real-label template tasks** A *real-label* template task is specified by a collection of “templates” labeled by real numbers, which are used to generate the train and test data. For instance, the data in Figure 1 is generated from the templates

$$“\alpha=1; \beta=-1; \text{print}(\alpha)” \rightarrow \text{label}=+1 \quad \text{and} \quad “\alpha=1; \beta=-1; \text{print}(\beta)” \rightarrow \text{label}=-1, \quad (1)$$

because every sample  $(x_i, y_i) \in \mathcal{X}^k \times \mathcal{Y}$  is formed by picking a template and replacing the placeholder symbols  $\alpha, \beta$  (which we call “wildcards”) with variable names. Each template should be thought of as a logical rule enforcing that all data matching the template must have the template’s label. Therefore, a template task measures the ability of an LLM to learn logical rules on abstract symbols: the LLM must infer the templates from training data, and at test time match samples to the corresponding templates to derive their labels. This framework captures several natural tasks:

- *Same/different task.* With templates “ $\alpha\alpha$ ” and “ $\alpha\beta$ ” labeled by +1 and −1, the task is to distinguish between equal symbols (e.g.,  $AA, BB$ ) or distinct symbols (e.g.,  $AB, BC$ ). This task has been empirically studied as a basic reasoning task (Kim et al., 2018).
- *More complex relations.* More complex mathematical functions of strings are also easy to encode: e.g., with templates “ $\alpha\alpha\beta$ ” and “ $\alpha\beta\alpha$ ” labeled with +1, and “ $\beta\alpha\alpha$ ” labeled with −1 the task is whether the first token occurs in the majority. See also (Webb et al., 2020b) for three other such tasks.
- *Word problems.* Many popular word problems follow a simple template. For example, the template “If  $\alpha$  gives  $\beta$  5  $\gamma$ , how many  $\gamma$  does  $\beta$  have?” labeled by +5, could generate the

data “If Alice gives Bob 5 oranges, how many oranges does Bob have?” or the data “If Rob gives Ada 5 apples, how many apples does Ada have?”

**Symbolic-label template tasks** A *symbolic-label* template task is the same, except that the templates are labeled by a wildcard. The data in Figure 2 is generated by:

$$“\alpha=“\gamma”; \beta=“\delta”; \text{print}(\alpha)” \rightarrow \text{label}=\gamma \quad \text{and} \quad “\alpha=“\gamma”; \beta=“\delta”; \text{print}(\beta)” \rightarrow \text{label}=\delta, \quad (2)$$

where  $\alpha, \beta, \gamma, \delta$  are wildcards. Other examples include:

- *Programming*. The template “`print("A")`” labeled with  $\alpha$  generates (`print("A")`, A) or (`print("dog")`, dog), and so captures the ability to robustly evaluate print statements.
- *Word problems*. The template “If  $\alpha$  gives  $\beta$   $\delta$   $\gamma$ , how many  $\gamma$  does  $\beta$  have?”, labeled by  $\delta$ , can generate several of the above word problems. An LLM that solves this task will output the correct answer of, say 10 if  $\delta = 10$  at test time even if  $\delta \neq 10$  in all training data.

In practice, such template tasks may occur as a natural component of a larger reasoning or word problem, but we isolate them here so that we can perform a theoretical analysis. We analyze the real- and symbolic-label settings separately, as they give complementary insights.

### 1.1.2 ANALYTICAL RESULTS FOR TEMPLATE TASKS *with real labels*

**(1) MLPs fail to generalize to unseen symbols** A classical criticism of connectionism by Marcus (1998) is that neural networks cannot mimic human abilities to reason because of their poor generalization on symbols that do not occur in their train data. In Appendix I, we support this criticism by proving that classical MLP architectures (a.k.a. fully-connected networks) trained by SGD or Adam will not generalize in template tasks on symbols unseen at training, regardless of the train data size.

**(2) Transformers generalize to unseen symbols, but require large data diversity** Nevertheless, the criticism of Marcus (1998) is not entirely valid for modern transformer architectures (Vaswani et al., 2017). We analyze the training dynamics of a transformer model and establish:

**Theorem 1.1** (Informal Theorem 4.4). *For any real-label template task, a wide-enough transformer architecture trained by gradient flow on sufficiently many samples generalizes on unseen symbols.*

Here the key points are: (a) *Universality*. The transformer architecture generalizes on symbols unseen in train data regardless of which and how many templates are used to define the reasoning task. (b) *Large enough number of samples*. Our theoretical guarantees require the training dataset size to be large, and even for very basic tasks like the two-template task in Figure 1, good generalization begins to occur only at a very large number of training samples considering the simplicity of the task. This raises the question of how the inductive bias of the transformer can be improved.

**(3) Improving data-efficiency of transformers** The proof of Theorem 1.1 inspires a parametrization modification that empirically lowers the quantity of data needed by an order of magnitude, by making it easier for the transformer to use the incidence matrix of the input string with itself:

**Observation 1.2.** *Adding one trainable parameter  $a$  to each attention head so that  $W_K W_Q^T$  is replaced by  $W_K W_Q^T + aI$  dramatically improves transformers’ data-efficiency on template tasks.*

### 1.1.3 ANALYTICAL RESULTS FOR TEMPLATE TASKS *with symbolic labels*

**(4) Transformers fail at copying unseen symbols** Surprisingly, the story is different for symbolic-label tasks. Transformers’ performance degrades as the model grows (an “inverse scaling” law (McKenzie et al., 2023)). Transformers fail even for the task of copying the input.

**Theorem 1.3** (Informal Theorem 5.1). *Transformers with large embedding dimension fail to generalize on unseen symbols for the copy-task outputting label “ $\alpha$ ” on template “ $\alpha$ ”.*

**(5) Modifying transformers for success** However, a small modification corrects this failure.

**Theorem 1.4** (Informal Theorem 5.2). *Adding one trainable parameter  $b$  to each attention head so that  $W_V W_O^T$  is replaced by  $W_V W_O^T + bI$  makes transformers generalize on the task of Theorem 1.3.*

### 1.1.4 EXPERIMENTAL VALIDATION AND EXPLORATION

We conclude with experimental validation, including showing that the transformer modifications proposed in Observation 1.2 and Theorem 1.4 improve performance in GPT-2 trained on Wikitext. We also show data-efficiency improvements on template tasks by fine-tuning a pretrained model, and give as an explanation the pronounced diagonals in  $W_K W_Q^T$  and  $W_V W_O^T$  matrices of pretrained models (Trockman & Kolter, 2023), which coincide with the proposed transformer modifications.

### 1.2 RELATED LITERATURE

A spate of recent work studies whether and how LLMs perform various reasoning tasks, each focusing on one component of reasoning: these include recognizing context-free grammars (Zhao et al., 2023; Allen-Zhu & Li, 2023), generalizing out-of-distribution when learning Boolean functions (Abbe et al., 2023), performing arithmetic (Nanda et al., 2023), learning in context (Garg et al., 2022; Ahn et al., 2023; Zhang et al., 2023), reasoning analogically Webb et al. (2020b), and evaluating indexing Zhang et al. (2021b). Our setting can be seen as a generalization of the tasks in (Kim et al., 2018) and (Webb et al., 2020b). Kim et al. (2018) shows experimentally that feedforward networks trained on the same/different templates  $\alpha\alpha$  vs.  $\alpha\beta$  do not generalize to symbols not seen in the training data (we provide a proof in Appendix I). Webb et al. (2020b) considers four tasks that can be viewed as template tasks with wildcard-only templates, proposes a network architecture and experimentally shows the benefits of training with Temporal Context Normalization (Webb et al., 2020a). In contrast, our focus is on understanding when the transformer architecture learns or fails to learn, and how to modify it to improve its data-efficiency for reasoning.

## 2 TRANSFORMER DEFINITION

We interchangeably denote an input by a string  $x \in \mathcal{X}^k$  or a matrix  $X \in \mathbb{R}^{k \times m}$  constructed by stacking the one-hot vectors  $X = [e_{x_1}, \dots, e_{x_k}]^T$  of the string’s tokens. We study a depth-1 transformer architecture (Vaswani et al., 2017). The transformer has  $H$  heads with parameters  $W_{K,h}, W_{Q,h}, W_{V,h}, W_{O,h} \in \mathbb{R}^{d_{head} \times d_{emb}}$ , an embedding layer  $W_E \in \mathbb{R}^{m \times d_{emb}}$ , positional embeddings  $P \in \mathbb{R}^{k \times d_{emb}}$ , an MLP layer with parameters  $W_A, W_B \in \mathbb{R}^{d_{mlp} \times d_{emb}}$ , a final unembedding layer with weights  $w_U \in \mathbb{R}^{d_{emb}}$ , and an activation function  $\phi$ . The network takes in  $X \in \mathbb{R}^{k \times m}$  and outputs

$$f_{\text{trans}}(X; \theta) = w_U^T z_2 \in \mathbb{R} \quad (\text{Unembedding layer})$$

where

$$z_2 = W_B^T \phi(W_A z_1) \in \mathbb{R}^{d_{emb}} \quad (\text{MLP layer})$$

$$z_1 = \sum_{h \in [H]} A_h^T e_k \in \mathbb{R}^{d_{emb}} \quad (\text{Attention layer output at final token})$$

$$A_h = \text{smax}(\beta Z_0 W_{K,h}^T W_{Q,h} Z_0^T) Z_0 W_{V,h}^T W_{O,h} \in \mathbb{R}^{k \times d_{emb}} \quad (\text{Attention heads})$$

$$Z_0 = X W_E + \gamma P \in \mathbb{R}^{k \times d_{emb}}. \quad (\text{Embedding layer})$$

Here  $\beta, \gamma \geq 0$  are two hyperparameters that control the inverse temperature of the softmax and the strength of the positional embeddings, respectively. The architecture is standard, except that we remove skip connections and layer norm as these are not needed for our theoretical results. Additional notations are:  $[n] = \{1, \dots, n\}$ .

## 3 TEMPLATE TASKS

We formally define template tasks with *real labels*. The case of *symbolic labels* is in Appendix J.

**Definition 3.1.** A **template** is a string  $z \in (\mathcal{X} \cup \mathcal{W})^k$ , where  $\mathcal{X}$  is an alphabet of tokens, and  $\mathcal{W}$  is an alphabet of “wildcards”. A **substitution map** is an injective function  $s : \mathcal{W} \rightarrow \mathcal{X}$ . We write  $\text{sub}(z, s) \in \mathcal{X}^k$  for the string where each wildcard is substituted with the corresponding token:  $\text{sub}(z, s)_i = z_i$  if  $z_i \in \mathcal{X}$ , and  $\text{sub}(z, s)_i = s(z_i)$  if  $z_i \in \mathcal{W}$ . The string  $x \in \mathcal{X}^k$  **matches** the template  $z$  if  $x = \text{sub}(z, s)$  for some substitution map  $s$  and also  $s(\mathcal{W}) \cap \{z_i\}_{i \in [k]} = \emptyset$ : i.e., the substituted tokens did not already appear in the template  $z$ .

**Example** Using Greek letters to denote the wildcards and Latin letters to denote regular tokens, the template “ $\alpha\alpha\beta ST$ ” matches the string “QQRST”, but *not* “QQQST” (because the substitution map is not injective) and *not* “QSSST” (because  $\beta$  is replaced by S which is already in the template).

A template task’s training data distribution is generated by picking a template randomly from a distribution, and substituting its wildcards with a random substitution map.

**Definition 3.2.** A real-label template data distribution  $\mathcal{D} = \mathcal{D}(\mu_{\text{tmpl}}, \{\mu_{\text{sub},z}\}_z, f_*, \sigma)$  is given by

- a template distribution  $\mu_{\text{tmpl}}$  supported on templates in  $(\mathcal{X} \cup \mathcal{W})^k$ ,
- for each  $z \in \text{supp}(\mu_{\text{tmpl}})$ , a distribution  $\mu_{\text{sub},z}$  over substitution maps  $s : \mathcal{W} \rightarrow \mathcal{X}$ ,
- template labelling function  $f_* : \text{supp}(\mu_{\text{tmpl}}) \rightarrow \mathbb{R}$ , and a label-noise parameter  $\sigma \geq 0$ .

We draw a sample  $(x, y) = (\text{sub}(z, s), f_*(z) + \xi) \sim \mathcal{D}$ , by drawing a template  $z \sim \mu_{\text{tmpl}}$ , a substitution map  $s \sim \mu_{\text{sub},z}$ , and label noise  $\xi \sim \mathcal{N}(0, \sigma^2)$ .

Finally, we define what it means for a model to generalize on unseen symbols; namely, the model should output the the correct label for any string  $x \in \mathcal{X}^k$ , regardless of whether the string is in the support of the training distribution.

**Definition 3.3.** A (random) estimator  $\hat{f} : \mathcal{X}^k \rightarrow \mathbb{R}$  **generalizes on unseen symbols** with  $(\epsilon, \delta)$ -error if the following is true. For any  $x \in \mathcal{X}^k$  that matches a template  $z \in \text{supp}(\mu_{\text{tmpl}})$ , we have

$$(\hat{f}(x) - f_*(z))^2 \leq \epsilon,$$

with probability at least  $1 - \delta$  over the randomness of the estimator  $\hat{f}$ .

**Example** If the training data is generated from a uniform distribution on templates “ $\alpha\alpha$ ” with label 1 and “ $\alpha\beta$ ” for label -1, then it might consist of the data samples  $\{(AA, 1), (BB, 1), (AB, -1), (BA, -1)\}$ . An estimator that generalizes to unseen symbols must correctly label string  $CC$  with +1 and string  $CD$  with -1, even though these strings consist of symbols that do not appear in the training set. This is a nontrivial reasoning task: a model that succeeds must effectively infer what the templates are given the training data, and then infer the label of the test string by matching it to the appropriate template.

## 4 ANALYSIS FOR TEMPLATE TASKS WITH REAL LABELS

We establish that transformers generalize on unseen symbols on any real-label template task, when trained with enough data. It is important to note that this is not true for all architectures, as we prove in Appendix I that MLPs trained by SGD or Adam will not succeed.

Our achievability result for transformers requires the templates in the distribution  $\mu_{\text{tmpl}}$  to be “disjoint”, since otherwise the correct label for a string  $x$  is not uniquely defined, because  $x$  could match more than one template:

**Definition 4.1.** Two templates  $z, z' \in (\mathcal{X} \cup \mathcal{W})^k$  are **disjoint** if no  $x \in \mathcal{X}^k$  matches both  $z$  and  $z'$ .

Furthermore, in order to ensure that the samples are not all copies of each other (which would not help generalization), we have to impose a diversity condition on the data.

**Definition 4.2.** The **data diversity** is measured by  $\rho = \min_{z \in \text{supp}(\mu_{\text{tmpl}})} \min_{t \in \mathcal{X}} \frac{1}{\mathbb{P}_{s \sim \mu_{\text{sub},z}}[t \in s(\mathcal{W})]}$ .

When the data diversity  $\rho$  is large, then no token is much more likely than others to be substituted. If  $\rho$  is on the order of the number of samples  $n$ , then most pairs of data samples will not be equal.

### 4.1 TRANSFORMER RANDOM FEATURES KERNEL

We analyze training only the final  $w_U$  layer of the transformer, keeping the other weights fixed at their random Gaussian initialization. Surprisingly, even though we only train the final layer of

the transformer, this is enough to guarantee generalization on unseen symbols.<sup>1</sup> Taking the width parameters  $H, d_{emb}, d_{mlp}, d_{head}$  to infinity, and the step size to 0, the SGD training algorithm with weight decay converges to kernel gradient flow with the following kernel  $K_{\text{trans}}$ ,<sup>2</sup>

$$K_{\text{trans}}(\mathbf{X}, \mathbf{Y}) = \mathbb{E}_{u,v}[\phi(u)\phi(v)] \text{ for } u, v \sim N(\mathbf{0}, \begin{bmatrix} K_{\text{attn}}(\mathbf{X}, \mathbf{X}) & K_{\text{attn}}(\mathbf{X}, \mathbf{Y}) \\ K_{\text{attn}}(\mathbf{Y}, \mathbf{X}) & K_{\text{attn}}(\mathbf{Y}, \mathbf{Y}) \end{bmatrix}) \quad (3)$$

where  $K_{\text{attn}}(\mathbf{X}, \mathbf{Y}) = \mathbb{E}_{\mathbf{m}(\mathbf{X}), \mathbf{m}(\mathbf{Y})}[\text{smax}(\beta \mathbf{m}(\mathbf{X}))^T (\mathbf{X} \mathbf{Y}^T + \gamma^2 \mathbf{I}) \text{smax}(\beta \mathbf{m}(\mathbf{Y}))]$

$$[\mathbf{m}(\mathbf{X}), \mathbf{m}(\mathbf{Y})] \sim N(\mathbf{0}, \begin{bmatrix} \mathbf{X} \mathbf{X}^T + \gamma^2 \mathbf{I} & \mathbf{X} \mathbf{Y}^T + \gamma^2 \mathbf{I} \\ \mathbf{Y} \mathbf{X}^T + \gamma^2 \mathbf{I} & \mathbf{Y} \mathbf{Y}^T + \gamma^2 \mathbf{I} \end{bmatrix}).$$

The function outputted by kernel gradient flow is known to have a closed-form solution in terms of the samples, the kernel, and the weight-decay parameter  $\lambda$ , which we recall in Proposition 4.3.

**Proposition 4.3** (How kernel gradient flow generalizes; see e.g., Welling (2013)). *Let  $(\mathbf{X}_1, y_1), \dots, (\mathbf{X}_n, y_n)$  be training samples. With the square loss and ridge-regularization of magnitude  $\lambda$ , kernel gradient flow with kernel  $K$  converges to the following solution*

$$\hat{f}(\mathbf{X}) = \mathbf{y}^T (\hat{\mathbf{K}} + \lambda \mathbf{I})^{-1} \mathbf{k}(\mathbf{X}), \quad (4)$$

where  $\mathbf{y} = [y_1, \dots, y_n] \in \mathbb{R}^n$  are the train labels,  $\hat{\mathbf{K}} \in \mathbb{R}^{n \times n}$  is the empirical kernel matrix and has entries  $\hat{K}_{ij} = K(\mathbf{X}_i, \mathbf{X}_j)$ , and  $\mathbf{k}(\mathbf{X}) \in \mathbb{R}^n$  has entries  $k_i(\mathbf{X}) = K(\mathbf{X}_i, \mathbf{X})$ .

## 4.2 TRANSFORMERS GENERALIZE ON UNSEEN SYMBOLS

We analyze the solution to the kernel gradient flow with the transformer random features, which corresponds to training the last layer with SGD with weight decay in the infinitely-wide, infinitely-small-step-size limit.

**Theorem 4.4** (Transformers generalize on unseen symbols). *Let  $\mu_{\text{tplt}}$  be supported on a finite set of pairwise-disjoint templates ending with [CLS] tokens. Then, for almost any  $\beta, \gamma, b_1, b_2$  parameters (except for a Lebesgue-measure-zero set), the transformer random features with  $\phi(t) = \cos(b_1 t + b_2)$  generalizes on unseen symbols.<sup>3</sup> Formally, there are constants  $c, C > 0$  and ridge regularization parameter  $\lambda > 0$  that depend only  $\beta, \gamma, b_1, b_2, \mu_{\text{tplt}}, f_*, \sigma$ , such that for any  $\mathbf{x}$  matching a template  $\mathbf{z} \in \text{supp}(\mu_{\text{tplt}})$  the kernel ridge regression estimator  $\hat{f}$  in (4) with kernel  $K_{\text{trans}}$  satisfies*

$$|\hat{f}(\mathbf{x}) - f_*(\mathbf{z})| \leq C \sqrt{\log(1/\delta)/n} + C \sqrt{1/\rho},$$

with probability at least  $1 - \delta - \exp(-cn)$  over the random samples.

The first term is due to the possible noise in the labels. The second term quantifies the amount of sample diversity in the data. Both the sample diversity and the number of samples must tend to infinity for an arbitrarily small error guarantee.

**Proof sketch** (1) In Lemma 4.5 we establish with a sufficient condition for kernel ridge regression to generalize on unseen symbols. (2) We prove that  $K_{\text{trans}}$  satisfies it.

(1) *Sufficient condition.* Let  $\mu_{\text{tplt}}$  be supported on templates  $\mathbf{z}_1, \dots, \mathbf{z}_r$ . Let  $\mathcal{R} = \cup_{i \in [k], j \in [r]} \{z_{j,i}\}$  be the tokens that appear in the templates. Let  $[n] = \mathcal{I}_1 \sqcup \mathcal{I}_2 \sqcup \dots \sqcup \mathcal{I}_n$  be the partition of the samples such that if  $a \in \mathcal{I}_j$  then sample  $(\mathbf{x}_a, y_a)$  is drawn by substituting the wildcards of template  $\mathbf{z}_j$ .

Two samples  $\mathbf{x}_a, \mathbf{x}_b$  that are drawn from the same template  $\mathbf{z}_j$  are not necessarily similar to each other as measured by the kernel: i.e., they might not have large kernel inner product  $K(\mathbf{x}_a, \mathbf{x}_b)$ . However, they will have similar relationship to most other samples: for most  $i \in [n]$  we will have,

$$K(\mathbf{x}_a, \mathbf{x}_i) \approx K(\mathbf{x}_b, \mathbf{x}_i).$$

<sup>1</sup>Empirically, we observe that generalization improves when all layers are trained; see Appendix B.

<sup>2</sup>This kernel is derived in Appendix H, and assumes that every string  $\mathbf{x}$  ends with a special [CLS] classification token that does not appear elsewhere in the string.

<sup>3</sup>We analyze the shifted and rescaled cosine activation function  $\phi(t) = \cos(b_1 t + b_2)$  out of technical convenience, but conjecture that most non-polynomial activation functions should succeed.



This is increasingly true as the data diversity parameter  $\rho$  grows, as it becomes increasingly likely that samples  $\mathbf{x}_a$  and  $\mathbf{x}_i$  have their wildcards substituted by disjoint sets of tokens that did not appear in the templates, and similarly for  $\mathbf{x}_b$  and  $\mathbf{x}_i$ , in which case  $K(\mathbf{x}_a, \mathbf{x}_i) = K(\mathbf{x}_b, \mathbf{x}_i)$ . Therefore, as the sample diversity increases, the empirical kernel matrix  $\hat{\mathbf{K}}$  becomes approximately block-structured with blocks  $\mathcal{I}_j \times \mathcal{I}_{j'}$ . In other words, for most samples  $\mathbf{x}_a, \mathbf{x}_b$  corresponding to template  $\mathbf{z}_j$ , and most  $\mathbf{x}_{a'}, \mathbf{x}_{b'}$  corresponding to template  $\mathbf{z}_{j'}$ , we have

$$K(\mathbf{x}_a, \mathbf{x}_{a'}) = K(\mathbf{x}_b, \mathbf{x}_{b'}) = K(\text{sub}(\mathbf{z}_j, s), \text{sub}(\mathbf{z}_{j'}, s')) := N_{j,j'}, \quad (5)$$

where  $s, s' : \mathcal{W} \rightarrow \mathcal{X}$  are substitution maps satisfying

$$s(\mathcal{W}) \cap s'(\mathcal{W}) = \emptyset \quad \text{and} \quad s(\mathcal{W}) \cap \mathcal{R} = s'(\mathcal{W}) \cap \mathcal{R} = \emptyset. \quad (6)$$

One can check that (5) and (6) uniquely define a matrix  $\mathbf{N} \in \mathbb{R}^{r \times r}$  which gives the entries in the blocks of  $\hat{\mathbf{K}}$ , with one block for each pair of templates.<sup>4</sup> If the matrix  $\mathbf{N}$  is nonsingular and the number of samples is large, then the span of the top  $r$  eigenvectors of  $\hat{\mathbf{K}}$  will align with the span of the indicator vectors on the sets  $\mathcal{I}_1, \dots, \mathcal{I}_r$ . Furthermore, when testing a string  $\mathbf{x}^{\text{test}}$  that matches template  $\mathbf{z}_j$ , but might not have appeared in the training set, it holds that for most  $a \in \mathcal{I}_j$ , we have

$$\mathbf{k}(\mathbf{x}^{\text{test}}) = [K(\mathbf{x}^{\text{test}}, \mathbf{x}_1), \dots, K(\mathbf{x}^{\text{test}}, \mathbf{x}_n)] \approx [K(\mathbf{x}_a, \mathbf{x}_1), \dots, K(\mathbf{x}_a, \mathbf{x}_n)] = \hat{\mathbf{K}}_{a,:}.$$

In words, the similarity relationship of  $\mathbf{x}^{\text{test}}$  to the training samples is approximately the same as the similarity relationship of  $\mathbf{x}_a$  to the training samples. So the kernel ridge regression solution (4) approximately equals the average of the labels of the samples corresponding to template  $\mathbf{z}_j$ , which in turn is approximately equal to the template label by a Chernoff bound,

$$\mathbf{y}^T (\hat{\mathbf{K}} + \lambda \mathbf{I})^{-1} \mathbf{k}(\mathbf{x}^{\text{test}}) \approx \frac{1}{|\mathcal{I}_j|} \sum_{a \in \mathcal{I}_j} y_a \approx f_*(\mathbf{z}_j). \quad (7)$$

Therefore, kernel ridge regression generalizes on  $\mathbf{x}^{\text{test}}$ . It is important to note that the number of samples needed until (7) is a good approximation depends on the nonsingularity of  $\mathbf{N}$ . This yields the sufficient condition for kernel ridge regression to succeed (proof in Appendix C).

**Lemma 4.5** (Informal Lemma C.2). *If  $\mathbf{N}$  is nonsingular, then (4) generalizes to unseen symbols.*

(2)  $K_{\text{trans}}$  satisfies the sufficient condition. We now show that for any collection of disjoint templates  $\mathbf{z}_1, \dots, \mathbf{z}_r$ , the matrix  $\mathbf{N}_{\text{trans}} := \mathbf{N} \in \mathbb{R}^{r \times r}$  defined with kernel  $K = K_{\text{trans}}$  is nonsingular. This is challenging because  $K_{\text{trans}}$  does not have a closed-form solution because of the expectation over softmax terms in its definition (3). We analyze the MLP layer and the attention layer of the transformer separately. We observe that a “weak” condition on  $K_{\text{attn}}$  can be lifted into the “strong” result that  $\mathbf{N}_{\text{trans}}$  is nonsingular. The intuition is that as long as  $K_{\text{attn}}$  is not a very degenerate kernel, it is unlikely that the MLP layer has the cancellations that to make  $\mathbf{N}_{\text{trans}}$  nonsingular.

**Lemma 4.6** (Nonsingularity of  $\mathbf{N}_{\text{trans}}$ ). *Suppose for every non-identity permutation  $\tau \in S_r \setminus \{\text{id}\}$ ,*

$$\sum_{i \in [r]} K_{\text{attn}}(\text{sub}(\mathbf{z}_i, s), \text{sub}(\mathbf{z}_i, s')) \neq \sum_{i \in [r]} K_{\text{attn}}(\text{sub}(\mathbf{z}_i, s), \text{sub}(\mathbf{z}_{\tau(i)}, s')), \quad (8)$$

*where  $s, s'$  are the substitution maps in the definition of  $\mathbf{N}_{\text{trans}}$  in (6). Let the MLP layer’s activation function be  $\phi(t) = \cos(b_1 t + b_2)$ . Then for almost any choice of  $b_1, b_2$  (except for a Lebesgue-measure-zero set), the matrix  $\mathbf{N}_{\text{trans}}$  is nonsingular.*

This is proved in Appendix E, by evaluating a Gaussian integral and showing  $\mathbf{N}_{\text{trans}}$  has Vandermonde structure. Although we use the cosine activation function, we conjecture that this result holds for most non-polynomial activation functions. Next, we prove the condition on  $\mathbf{N}_{\text{attn}}$ .

**Lemma 4.7** (Non-degeneracy of  $K_{\text{attn}}$ ). *The condition (8) holds for Lebesgue-almost any  $\beta, \gamma$ .*

The proof is in Appendix F. First, we prove the analyticity of the kernel  $K_{\text{attn}}$  in terms of the hyperparameters  $\beta$  and  $\gamma$ . Because of the identity theorem for analytic functions, it suffices to show at least one choice of hyperparameters  $\beta$  and  $\gamma$  satisfies (8) for all non-identity permutations  $\tau$ . Since  $K_{\text{attn}}$  does not have a closed-form solution, we find such a choice of  $\beta$  and  $\gamma$  by analyzing the Taylor-series expansion of  $K_{\text{attn}}$  around  $\beta = 0$  and  $\gamma = 0$  up to order-10 derivatives.

<sup>4</sup>This assumes a “token-symmetry” property of  $K$  that is satisfied by transformers; details in the full proof.

### 4.3 IMPROVING TRANSFORMER DATA-EFFICIENCY WITH $W_K W_Q^T + aI$ PARAMETRIZATION

Can we use these insights to improve transformers’ data-efficiency in template tasks? In the proof, the nonsingularity of  $N$  in Lemma 4.5 drives the model’s generalization on unseen symbols. This suggests that an approach improve data-efficiency is to make  $N$  better-conditioned by modifying the transformer parametrization. We consider here the simplest task, with templates “ $\alpha\alpha$ ” and “ $\alpha\beta$ ” labeled with  $+1$  and  $-1$ , respectively. For tokens  $A, B, C, D \in \mathcal{X}$ , the matrix  $N$  is

$$N = \begin{bmatrix} K(AA, BB) & K(AA, BC) \\ K(BC, AA) & K(AB, CD) \end{bmatrix}$$

If  $K$  is an inner-product kernel,  $K(\mathbf{x}, \mathbf{x}') = \kappa(\sum_{i \in [k]} 1(x_i = x'_i))$ , as from an MLP, then  $K(AA, BB) = K(AA, BC) = K(BC, AA) = K(AB, CD) = \kappa(0)$ , so  $N$  is singular and generalization is not achieved. Intuitively, every sample  $\mathbf{x}_i$  has the same “similarity profile to other data”  $\hat{K}_{i,:} = [K(\mathbf{x}_i, \mathbf{x}_1), \dots, K(\mathbf{x}_i, \mathbf{x}_n)]$ , so the kernel method cannot identify the samples that come from the same template as  $\mathbf{x}^{test}$  via the similarity profile. In contrast, the transformer kernel succeeds since it incorporates information about the incidence matrix  $\mathbf{X} \mathbf{X}^T$ , which is different between templates, and does not depend on the symbol substitution. By reparametrizing each head to  $W_K W_Q^T + aI$ , we add a scaling of  $\mathbf{X} \mathbf{X}^T$  and further emphasize it in the transformer.

## 5 ANALYSIS FOR TEMPLATE TASKS WITH SYMBOLIC LABELS

In the previous section, we considered tasks under a regression setting with mean-squared error loss. We now switch to a next-token prediction setting with the cross-entropy loss. The symbolic-label variant of template tasks is analogous to the real-label template tasks studied so far, except that the output label is a token as in the example of Figure 2; formal definition is in Appendix J. For simplicity, we consider the architecture with just the attention layer, and we tie the embedding and unembedding weights as in practice:

$$f_{\text{attn}}(\mathbf{X}; \theta) = \mathbf{W}_E \mathbf{z}_{\text{attn}} \in \mathbb{R}^m. \quad (9)$$

We also consider the simplest template task with symbolic labels and show that transformers will fail to generalize: template “ $\alpha$ ” labeled by “ $\alpha$ ”. An example dataset generated from this template could be  $\{(A, A), (B, B), (C, C)\}$ , where  $A, B, C \in \mathcal{X}$  are tokens. Because the template has length  $k = 1$ , the architecture simplifies to

$$f_{\text{attn}}(\mathbf{X}; \theta) = \mathbf{W}_E \left( \sum_{h \in [H]} \mathbf{W}_{O,h}^T \mathbf{W}_{V,h} \right) (\mathbf{W}_E^T \mathbf{X}^T + \gamma \mathbf{P}^T). \quad (10)$$

Despite the simplicity of the task,  $f_{\text{attn}}$  does not generalize on unseen symbols when trained, when we take the embedding dimension large. Our evidence is from analyzing the early time of training. Define the train loss and test loss as follows, where  $\ell$  is the cross-entropy loss and  $x^{test}$  is a token that does not appear in the training data,

$$\mathcal{L}_{\text{train}}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(f_{\text{attn}}(x_i; \theta), y_i) \quad \text{and} \quad \mathcal{L}_{\text{test}}(\theta) = \ell(f_{\text{attn}}(x^{test}), y^{test}).$$

We train with gradient flow, and show that the generalization loss on unseen symbols does not decrease for infinite-width transformers on the symbolic-label “copying” task where the template is “ $\alpha$ ” and is labeled by “ $\alpha$ ”.

**Theorem 5.1** (Failure of transformers at copying). *For any learning rates such that  $-\frac{\partial \mathcal{L}_{\text{train}}}{\partial t} \big|_{t=0} = O(1)$ , we must have that  $\frac{\partial \mathcal{L}_{\text{test}}}{\partial t} \big|_{t=0} \rightarrow 0$  as  $d_{\text{emb}} \rightarrow \infty$ .*

The intuition comes from examining (10), and noting that at early times the evolution of the weights  $\mathbf{W}_{O,h}^T \mathbf{W}_{V,h}$  will roughly lie in the span of  $\{\mathbf{W}_E^T e_{x_i} e_{x_i}^T \mathbf{W}_E\}_{i \in [n]}$ , which as the embedding dimension becomes large will be approximately orthogonal to the direction  $\mathbf{W}_E^T e_{x^{test}} e_{x^{test}}^T \mathbf{W}_E$  that would lower the test loss. However, this suggests the following:

**Theorem 5.2** (Adding one parameter allows copying). *After reparametrizing the attention (9) so that in each head  $\mathbf{W}_{O,h}^T \mathbf{W}_{V,h}$  is replaced by  $\mathbf{W}_{O,h}^T \mathbf{W}_{V,h} + b_h \mathbf{I}$  where  $b_h$  is a trainable parameter, there are learning rates such that  $-\frac{\partial \mathcal{L}_{\text{train}}}{\partial t} \big|_{t=0} = O(1)$  and  $-\frac{\partial \mathcal{L}_{\text{test}}}{\partial t} \big|_{t=0} = \Omega(1)$  as  $d_{\text{emb}} \rightarrow \infty$ .*



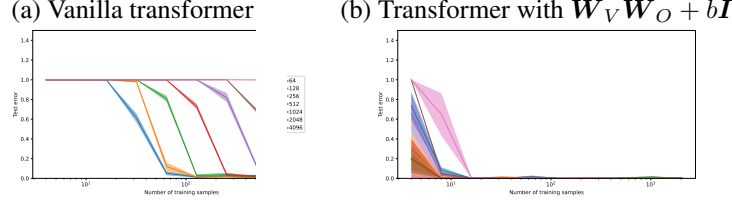


Figure 3: (a) Inverse scaling law: Transformers fail on copy task as embedding dimension  $d_{emb}$  grows (Theorem 5.1) (b) Success when reparametrizing  $W_V W_O^T$  as  $W_V W_O + bI$  (Theorem 5.2).

	GPT-2	GPT-2 + KQ, VO identity
Wikitext2	64.00	<b>60.46</b>
Wikitext103	16.83	<b>16.40</b>

Figure 4: Perplexity of GPT-2 trained with Adam learning rate 3e-4 for 20 epochs on Wikitext (smaller is better). GPT-2 has 117M parameters, and we add an extra 288 parameters (2 per head).

Figure 3 illustrates the benefit of this additional per-head parameter on the copying task. Our proposed reparametrization is similar to adding a trainable skip connection [He et al. \(2016\)](#), but not equivalent since there is an important subtle difference. The addition of  $b_h I$  encodes an *attention-modulated skip connection*, and thus allows copying tokens between the transformer’s streams.

## 6 EXPERIMENTS

Figures 1 and 2 show our reparametrizations can give a significant benefit on template tasks. Figure 4 shows they can also give improvements on real data. In Figure 5, we find that fine-tuning a pretrained model helps with a template task. This might be explained by several heads of the pretrained model with diagonals stronger from other weights (originally observed in [\(Trockman & Kolter, 2023\)](#)). These learned diagonals resemble our proposed transformer modifications and so might be driving the data-efficiency of fine-tuning a pretrained model. Appendix B provides extensive experiments on the effect of hyperparameters, inductive biases of different models, and varying levels of difficulty.

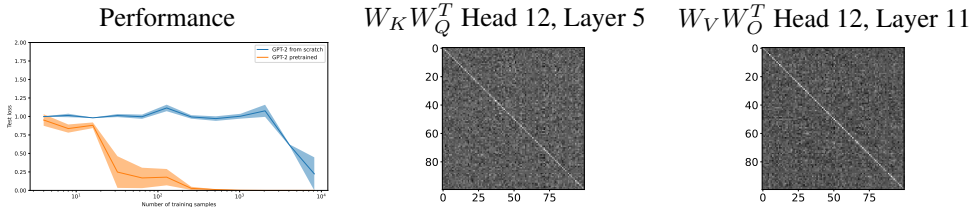


Figure 5: Left: Pretrained versus from-scratch GPT-2 test loss on  $\alpha\beta\alpha$  vs.  $\alpha\beta\beta$  template task. Right: example GPT-2 pretrained heads that have learned diagonals (zoomed in to 100x100 corner).

## 7 DISCUSSION

Our investigation of Transformers’ ability to generalize to unseen symbols reveals that, while this architecture is powerful, it is far from optimal, since it requires large amounts of data to learn basic reasoning abilities and fails altogether at copying unseen symbols. The transformer reparametrizations proposed are a step towards promoting an inductive bias towards logic in LLMs. Architectural modifications should be explored in analyses of complementary reasoning tasks (such as analogies and syllogisms) that in practical settings are combined with the ability to generalize on abstract symbols. Apart from architectural modifications, data augmentation approaches (e.g., by concatenating the tensorization  $\mathbf{X}\mathbf{X}^T$  to the input to encourage use of these features) should also be investigated.

## REFERENCES

- Emmanuel Abbe and Enric Boix-Adsera. On the non-universality of deep learning: quantifying the cost of symmetry. *Advances in Neural Information Processing Systems*, 35:17188–17201, 2022.
- Emmanuel Abbe, Elisabetta Cornacchia, Jan Hazla, and Christopher Marquis. An initial alignment between neural network and target is needed for gradient descent to learn. In *International Conference on Machine Learning*, pp. 33–52. PMLR, 2022.
- Emmanuel Abbe, Samy Bengio, Aryo Lotfi, and Kevin Rizk. Generalization on the unseen, logic reasoning and degree curriculum. *arXiv preprint arXiv:2301.13105*, 2023.
- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *arXiv preprint arXiv:2306.00297*, 2023.
- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 1, context-free grammar. *arXiv preprint arXiv:2305.13673*, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Lenaïc Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. *Advances in neural information processing systems*, 31, 2018.
- Lenaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in neural information processing systems*, 32, 2019.
- George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- Jerry A Fodor. *The language of thought*, volume 5. Harvard university press, 1975.
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Junkyoung Kim, Matthew Ricci, and Thomas Serre. Not-so-clevr: learning same–different relations strains feedforward neural networks. *Interface focus*, 8(4):20180011, 2018.
- Steven G Krantz and Harold R Parks. *A primer of real analytic functions*. Springer Science & Business Media, 2002.
- Trenton Kriete, David C Noelle, Jonathan D Cohen, and Randall C O’Reilly. Indirection and symbol-like processing in the prefrontal cortex and basal ganglia. *Proceedings of the National Academy of Sciences*, 110(41):16390–16395, 2013.
- Zhiyuan Li, Yi Zhang, and Sanjeev Arora. Why are convolutional nets more sample-efficient than fully-connected nets? *arXiv preprint arXiv:2010.08515*, 2020.
- Gary F Marcus. Rethinking eliminative connectionism. *Cognitive psychology*, 37(3):243–282, 1998.

- Ian R McKenzie, Alexander Lyzhov, Michael Pieler, Alicia Parrish, Aaron Mueller, Ameya Prabhu, Euan McLean, Aaron Kirtland, Alexis Ross, Alisa Liu, et al. Inverse scaling: When bigger isn't better. *arXiv preprint arXiv:2306.09479*, 2023.
- Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
- Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In *Conference on Learning Theory*, pp. 2388–2464. PMLR, 2019.
- Boris Samuilovich Mityagin. The zero set of a real analytic function. *Mathematical Notes*, 107(3-4):529–530, 2020.
- Neel Nanda, Lawrence Chan, Tom Liberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023.
- Allen Newell. Physical symbol systems. *Cognitive science*, 4(2):135–183, 1980.
- Allen Newell, John C Shaw, and Herbert A Simon. Report on a general problem solving program. In *IFIP congress*, volume 256, pp. 64. Pittsburgh, PA, 1959.
- Andrew Y Ng. Feature selection,  $\ell_1$  vs.  $\ell_2$  regularization, and rotational invariance. In *Proceedings of the twenty-first international conference on Machine learning*, pp. 78, 2004.
- Grant Rotskoff and Eric Vanden-Eijnden. Parameters as interacting particles: long time convergence and asymptotic error scaling of neural networks. *Advances in neural information processing systems*, 31, 2018.
- Ohad Shamir. Distribution-specific hardness of learning neural networks. *The Journal of Machine Learning Research*, 19(1):1135–1163, 2018.
- Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of deep neural networks. *Mathematics of Operations Research*, 47(1):120–152, 2022.
- Asher Trockman and J Zico Kolter. Mimetic initialization of self-attention layers. *arXiv preprint arXiv:2305.09828*, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Taylor Webb, Zachary Dulberg, Steven Frankland, Alexander Petrov, Randall O'Reilly, and Jonathan Cohen. Learning representations that support extrapolation. In *International conference on machine learning*, pp. 10136–10146. PMLR, 2020a.
- Taylor W Webb, Ishan Sinha, and Jonathan D Cohen. Emergent symbols through binding in external memory. *arXiv preprint arXiv:2012.14601*, 2020b.
- Max Welling. Kernel ridge regression. *Max Welling's classnotes in machine learning*, pp. 1–3, 2013.
- Greg Yang and Edward J Hu. Tensor programs iv: Feature learning in infinite-width neural networks. In *International Conference on Machine Learning*, pp. 11727–11737. PMLR, 2021.
- Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Chuanqi Tan, and Chang Zhou. Scaling relationship on learning mathematical reasoning with large language models. *arXiv preprint arXiv:2308.01825*, 2023.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021a.

Chiyuan Zhang, Maithra Raghu, Jon Kleinberg, and Samy Bengio. Pointer value retrieval: A new benchmark for understanding the limits of neural network generalization. *arXiv preprint arXiv:2107.12580*, 2021b.

Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *arXiv preprint arXiv:2306.09927*, 2023.

Haoyu Zhao, Abhishek Panigrahi, Rong Ge, and Sanjeev Arora. Do transformers parse while predicting the masked word? *arXiv preprint arXiv:2303.08117*, 2023.

## CONTENTS

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Our contributions . . . . .	1
1.1.1	Template tasks: a framework for reasoning with abstract symbols . . . . .	2
1.1.2	Analytical results for template tasks <i>with real labels</i> . . . . .	3
1.1.3	Analytical results for template tasks <i>with symbolic labels</i> . . . . .	3
1.1.4	Experimental validation and exploration . . . . .	4
1.2	Related literature . . . . .	4
<b>2</b>	<b>Transformer definition</b>	<b>4</b>
<b>3</b>	<b>Template tasks</b>	<b>4</b>
<b>4</b>	<b>Analysis for template tasks with real labels</b>	<b>5</b>
4.1	Transformer random features kernel . . . . .	5
4.2	Transformers generalize on unseen symbols . . . . .	6
4.3	Improving transformer data-efficiency with $W_K W_Q^T + aI$ parametrization . . . . .	8
<b>5</b>	<b>Analysis for template tasks with symbolic labels</b>	<b>8</b>
<b>6</b>	<b>Experiments</b>	<b>9</b>
<b>7</b>	<b>Discussion</b>	<b>9</b>
<b>A</b>	<b>Details for figures in main text</b>	<b>14</b>
<b>B</b>	<b>Additional experiments</b>	<b>14</b>
B.1	Effect of transformer hyperparameters . . . . .	15
B.2	Effect of complexity of task . . . . .	15
B.3	Effect of inductive bias of model . . . . .	15
<b>C</b>	<b>Proof of Theorem 4.4</b>	<b>22</b>
C.1	Part 1. General sufficient condition for good test loss . . . . .	22
C.2	Part 2. Analyzing the transformer random features kernel . . . . .	23
C.3	Concluding the proof of Theorem 4.4 . . . . .	23
<b>D</b>	<b>Sufficient condition for kernel method to generalize on unseen symbols (Proof of Lemma C.2)</b>	<b>23</b>
D.1	Deferred proofs of claims . . . . .	25
<b>E</b>	<b>Nonsingularity of random features after MLP layer (Proof of Lemma 4.6)</b>	<b>27</b>
<b>F</b>	<b>Analysis of attention layer features (Proof of Lemma 4.7)</b>	<b>29</b>

F.1	Low-order derivatives of attention kernel . . . . .	29
F.2	Simplifying terms . . . . .	30
F.2.1	Assuming $[1^T X]_{\mathcal{R}} = [1^T Y]_{\mathcal{R}}$ . . . . .	30
F.2.2	Assuming $[X]_{[k] \times \mathcal{R}} = [Y]_{[k] \times \mathcal{R}}$ . . . . .	31
F.2.3	Assuming $1^T X X^T 1 = 1^T Y Y^T 1$ . . . . .	31
F.2.4	Assuming $1^T X X^T = 1^T Y Y^T$ . . . . .	31
F.3	Proof of Lemma F.1 . . . . .	31
<b>G</b>	<b>Analyticity of attention kernel (technical result)</b>	<b>34</b>
G.1	Technical lemmas for quantifying power series convergence . . . . .	34
G.2	Application of technical lemmas to attention kernel . . . . .	36
<b>H</b>	<b>Derivation of transformer kernel</b>	<b>37</b>
H.1	Transformer architecture . . . . .	37
H.2	Random features kernel . . . . .	37
H.3	Informal derivation . . . . .	38
<b>I</b>	<b>MLPs fail to generalize on unseen symbols</b>	<b>39</b>
<b>J</b>	<b>Deferred details for symbolic-label template tasks</b>	<b>42</b>
J.1	Definition of symbolic-label template tasks . . . . .	42
J.2	Failure of transformers to copy and modification that succeeds . . . . .	42

## A DETAILS FOR FIGURES IN MAIN TEXT

In Figure 1, The architecture is a 2-layer transformer with 16 heads per layer, embedding dimension 128, head dimension 64, MLP dimension 256, trained with Adam with learning rate  $1e-3$  and batch-size 1024. The  $n$  training samples are chosen by picking the variable names at random from an alphabet of  $n$  tokens. The test set is the same two programs but with disjoint variable names. The reported error bars are on average over 5 trials. The learning rate for each curve is picked as the one achieving best generalization in  $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}\}$ . In Figure 2, the setting is the same except that the transformer is 4-layer transformer and has embedding dimension 512. In Figure 3 the same hyperparameters as in Figure 1 are used.

In order to measure the generalization performance of the learned model on unseen symbols, we evaluate it on a test set and a validation set which each consist of 100 samples drawn in the same way as the training dataset, but each using a disjoint alphabet of size 100. Therefore, there is no overlap in the support of the train, test, and validation distributions. We use the validation loss to select the best epoch of training out of 1000 epochs. We report the test loss on this saved model.

## B ADDITIONAL EXPERIMENTS

We report extensive additional experiments probing the template task framework. In each of these, the training dataset consists of  $n$  random training samples. Each sample is drawn according to a template distribution. The following are template tasks on which we test.

- $\alpha\beta\alpha$  vs.  $\alpha\beta\beta$  task. Uniform on two templates  $\alpha\beta\alpha$  and  $\alpha\beta\beta$  with labels 1, -1 respectively and  $\alpha$  and  $\beta$  are wildcards.
- $\alpha\beta\alpha\beta$  vs.  $\alpha\alpha\beta\beta$  task. Same as above, except with templates  $\alpha\beta\alpha\beta$  and  $\alpha\alpha\beta\beta$ .



- *Length- $k$  majority task.* Uniform on  $2^{k-1}$  templates  $\alpha \times \{\alpha, \beta\}^{k-1}$  where  $\alpha$  and  $\beta$  are wildcards. A template  $z$  has label 1 if its first token occurs in the majority of the rest of the string, and -1 otherwise. Namely,  $f_*(z) = \begin{cases} 1, & |\{i : z_1 = z_i\}| > (k+1)/2 \\ -1, & \text{otherwise} \end{cases}$ .
- *Random template task.* A certain number  $r$  of templates are drawn uniformly from  $(\mathcal{W} \cup \mathcal{X})^k$ , conditioned on being pairwise distinct. The task is the uniform distribution over these  $r$  templates, with random Gaussian labels centered and scaled so that the trivial MSE is 1.

For any of these tasks, we generate  $n$  training samples as follows. We substitute the wildcards for regular tokens using a randomly chosen injective function  $s : \mathcal{W} \rightarrow \mathcal{X}$  where  $\mathcal{X}$  is an alphabet of size  $n$  (which is the same size as the number of samples). For example, if a given sample is generated from template  $\alpha\beta\alpha$  with substitution map  $s$  mapping  $s(A) = 12$ ,  $s(B) = 5$ , then the sample will be  $[12, 5, 12]$ . Error bars are over 5 trials, unless otherwise noted.

### B.1 EFFECT OF TRANSFORMER HYPERPARAMETERS

We test an out-of-the-box transformer architecture on the  $\alpha\beta\alpha$  vs.  $\alpha\beta\beta$  task, varying some of the hyperparameters of the transformer to isolate their effect while keeping all other hyperparameters fixed. The base hyperparameters are depth 2, embedding dimension 128, head dimension 64, number of heads per layer 16, trained with Adam with minibatch size 1024 for 1000 epochs. Our experiments are as follows:

- *Learning rate and  $n$ .* In Figure 6 we vary the learning rate and  $n$ .
- *Learning rate and depth.* In Figure 7 and Figure 8, we vary the learning rate and the depth, for  $n = 512$  and  $n = 1024$ , respectively.
- *Learning rate and number of heads.* In Figure 9 and 10, we vary the learning rate and number of heads, for  $n = 512$  and  $n = 1024$ , respectively.
- *Learning rate and embedding dimension.* In Figure 11 we vary the learning rate and embedding dimension for  $n = 1024$ .
- *Learning rate and batch size.* In Figure 12, we vary the learning rate and batch-size for  $n = 512$ . In Figure 13 we vary the batch-size and  $n$  for learning rate 0.001.

### B.2 EFFECT OF COMPLEXITY OF TASK

We test an out-of-the-box transformer architecture with depth 2, embedding dimension 128, head dimension 64, number of heads 16, trained with Adam with batch-size 1024 for 1000 epochs, on various template tasks.

- *Comparing difficulty of various tasks.* Figure 14 we plot the performance on various simple tasks.
- *Random tasks.* In Figures 15, 16, 17, and 18, we test on random template tasks, and investigate the effects of template length, wildcard alphabet size, regular token alphabet size, number of templates.

### B.3 EFFECT OF INDUCTIVE BIAS OF MODEL

We provide experiments probing the effect of the inductive bias of the model:

- *Different architectures.* In Figure 19, we plot the test loss for different architectures on the  $\alpha\beta\alpha$  vs.  $\alpha\beta\beta$  template task, including transformers with trainable identity perturbations to  $W_Q W_K^T$ , to  $W_V W_O^T$ , to both  $W_Q W_K^T$  and  $W_V W_O^T$ , or to neither.
- *Size of model.* In Figure 20 we compare the test loss of fine-tuning small, medium and large pretrained GPT-2 networks on the  $\alpha\beta\alpha$  vs.  $\alpha\beta\beta$  template task.
- *MLP with  $XX^T$  data augmentation vs. transformer.* In Figure 21, we compare the test loss of a transformer with the test loss of an MLP where the input data has been augmented by concatenating  $\text{vec}(XX^T)$ , which is a data augmentation that improves performance under the NTK criterion similarly to the discussion in Section 4.3 and the discussion section.

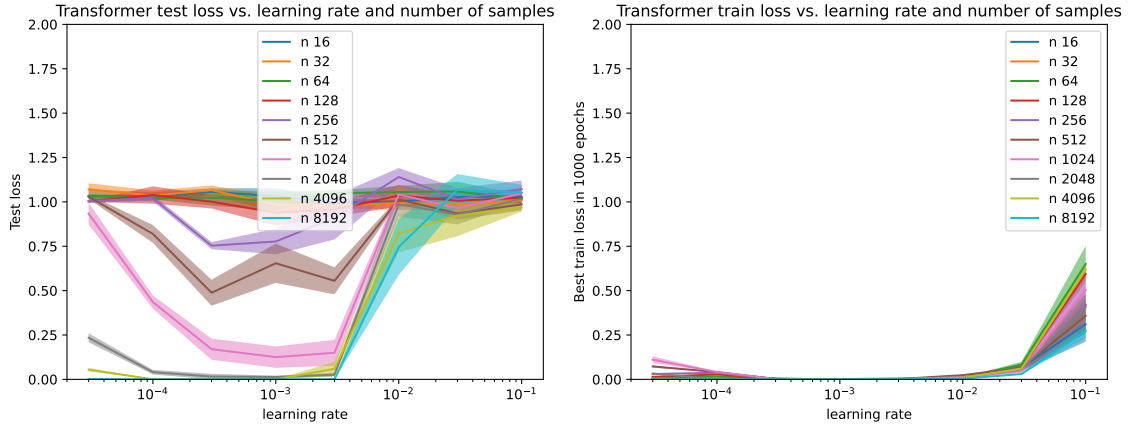


Figure 6: Learning rate versus  $n$  = number of samples = training alphabet size. Taking too large or too small of a learning rate can hurt generalization even when the train loss is close to zero.

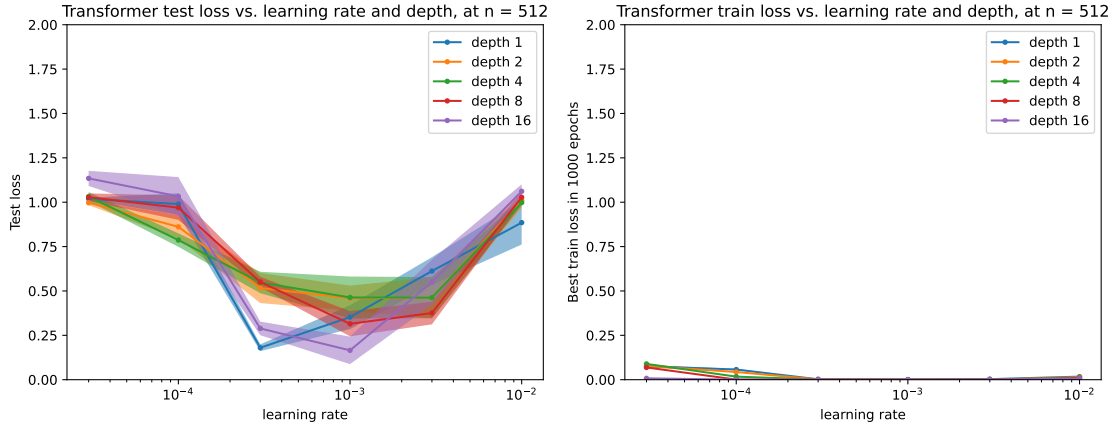


Figure 7: Learning rate vs. depth at  $n = 512$ . No clear relationship between depth and generalization. Too large or too small of a learning rate can hurt generalization.

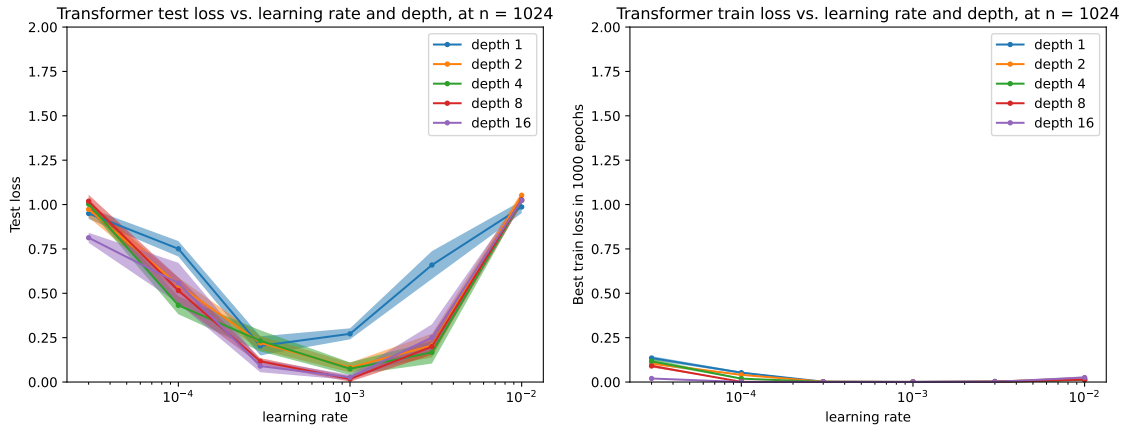


Figure 8: Learning rate vs. depth at  $n = 1024$ . Unlike  $n = 512$  case, in previous figure, larger depth typically performs better.

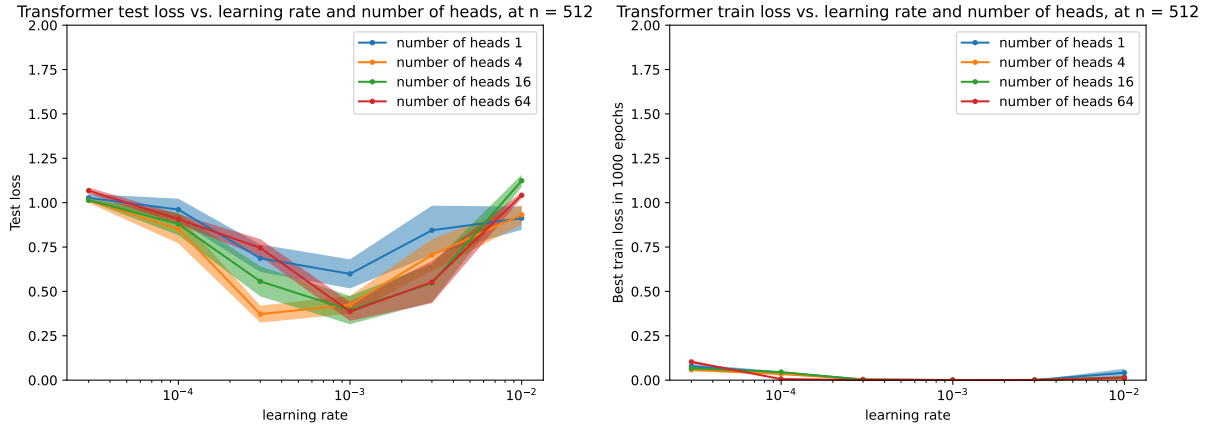


Figure 9: Learning rate vs. number of heads per layer at  $n = 512$ . More heads are better than one head.

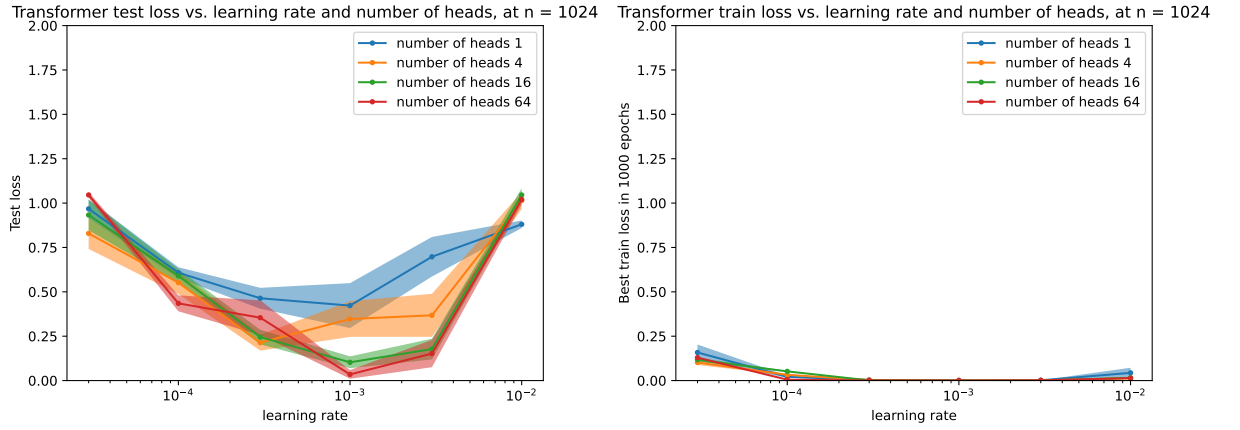


Figure 10: Learning rate vs. number of heads at  $n = 1024$ . More heads are better.

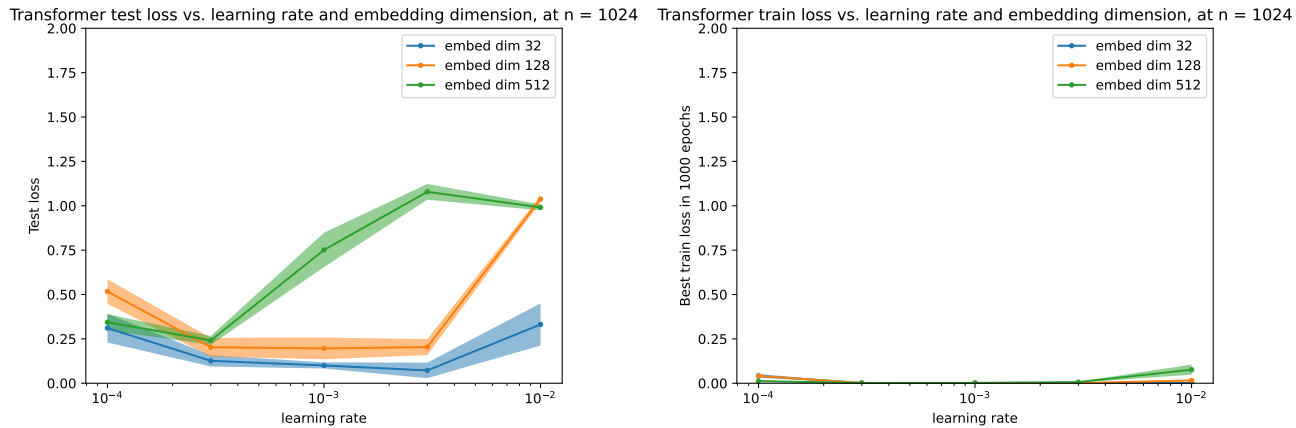
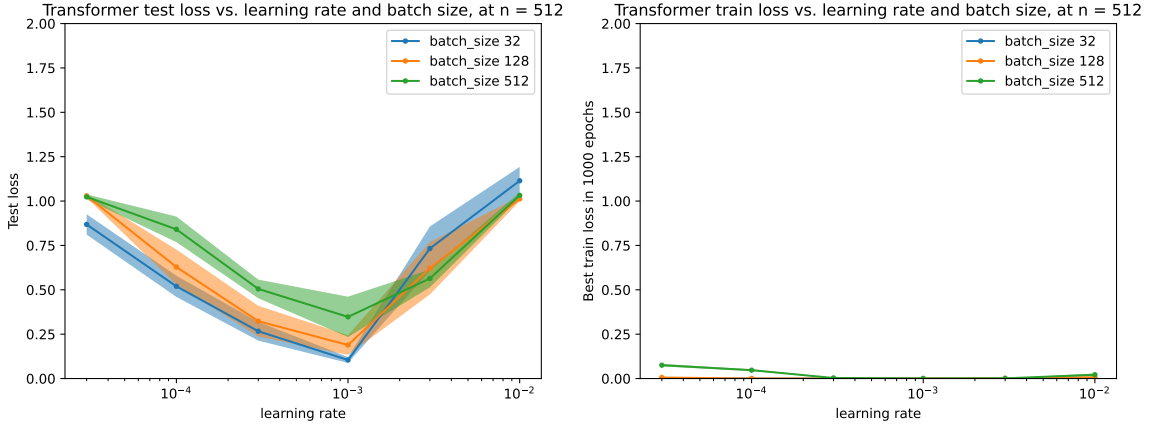
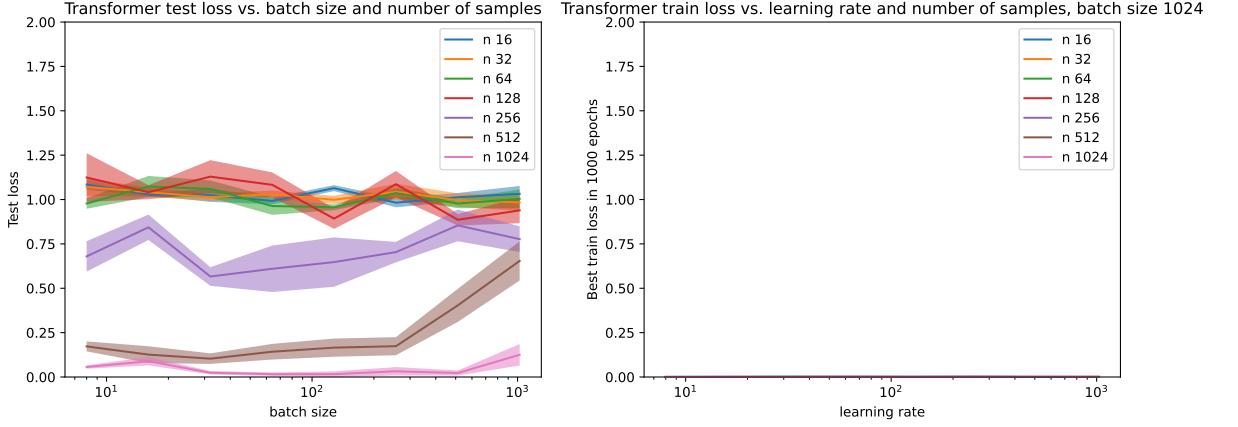
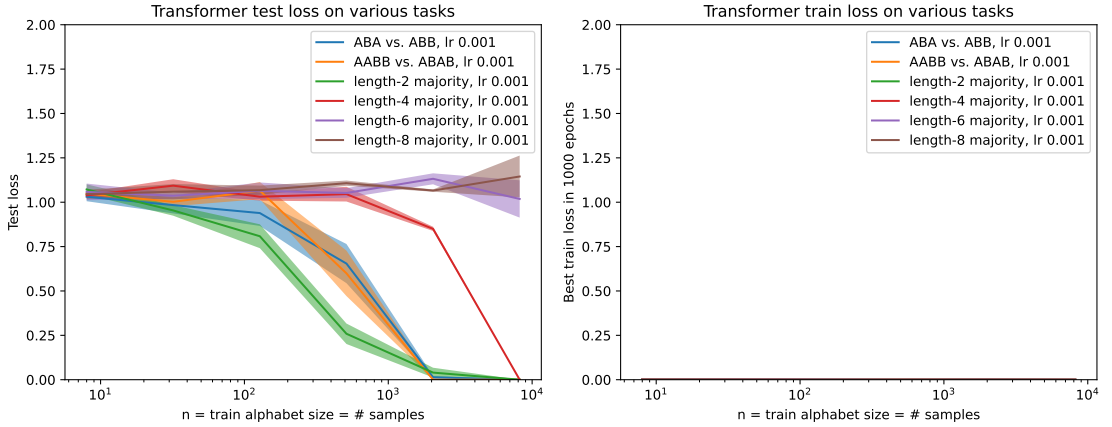


Figure 11: Learning rate vs. embedding dimension at  $n = 1024$ . Smaller embedding dimension is generally better.

Figure 12: Learning rate vs. batch-size at  $n = 512$ . Smaller batch size is better.Figure 13: Batch size vs.  $n$  = number of training samples = training alphabet size. Smaller batch size is generally better, which is most visible at  $n = 512$ .Figure 14: Test and train loss of transformer for various tasks. The  $\alpha\beta\alpha$  vs.  $\alpha\beta\beta$  task consists of two templates  $\alpha\beta\alpha$  and  $\alpha\beta\beta$  with labels +1, -1. The  $\alpha\alpha\beta\beta$  vs.  $\alpha\beta\alpha\beta$  task has templates +1, -1. For each  $k$ , the length- $k$  majority task consists of all templates in  $\{\alpha\} \times \{\alpha, \beta\}^{k-1}$ , where each template has label 1 if  $\alpha$  occurs more times in the last  $k - 1$  entries, and label -1 if  $\alpha$  occurs fewer times in the last  $k - 1$  entries. The trivial model that outputs 0 always will achieve test loss of 1.

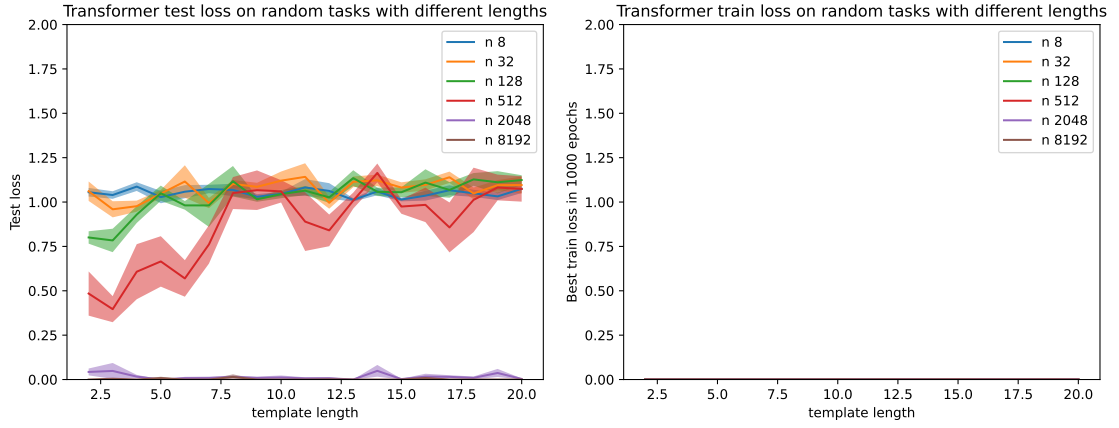


Figure 15: Performance on tasks corresponding of two, distinct random templates with two wildcards  $\alpha, \beta$ , and with labels  $1, -1$ , respectively. Performance degrades as the template length increases.

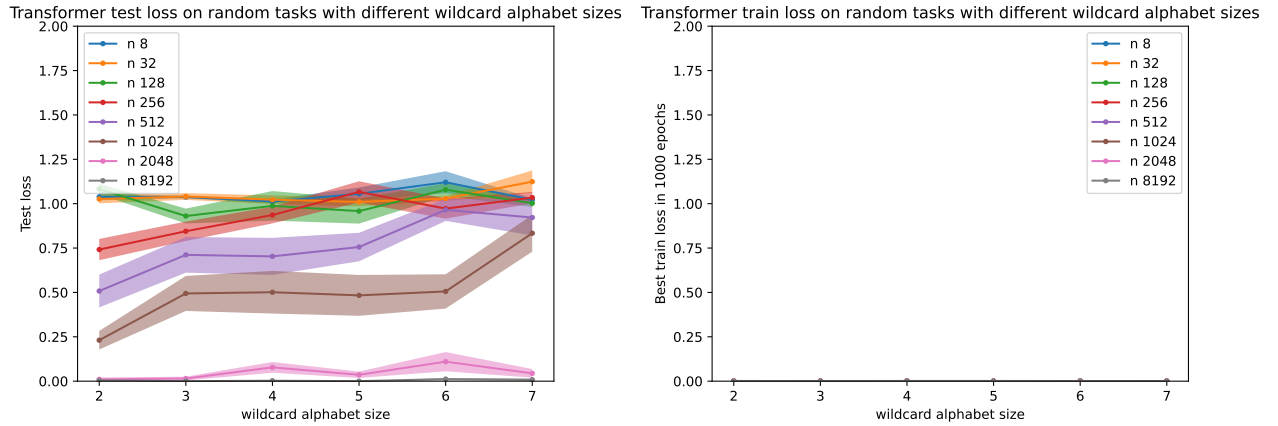


Figure 16: Performance on tasks corresponding of two random templates of length 5, labeled with  $1, -1$ , respectively. Each template is sampled randomly from  $\mathcal{W}^5$ , conditioned on the two templates being distinct. We vary the wildcard alphabet size  $|\mathcal{W}|$ . Performance generally degrades as the wildcard alphabet size increases.

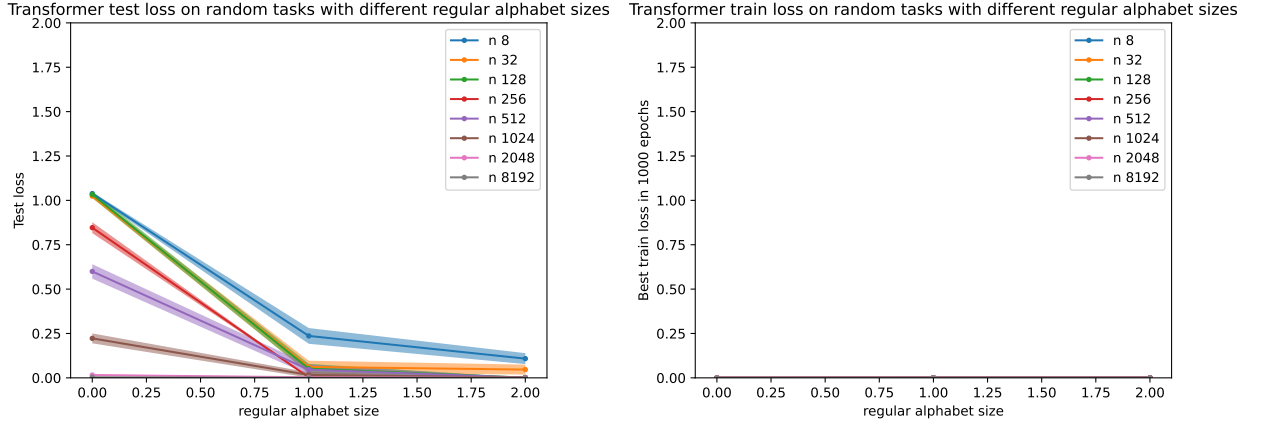


Figure 17: Performance on tasks corresponding of two random templates of length 5, labeled with 1,  $-1$ , respectively. Each template is sampled randomly from  $(\mathcal{W} \cup \mathcal{X})^5$ , conditioned on the two templates being distinct. We keep  $|\mathcal{W}| = 2$  and vary the regular token alphabet size  $|\mathcal{X}|$  between 0 and 2. Performance quickly improves as the regular token alphabet size increases.

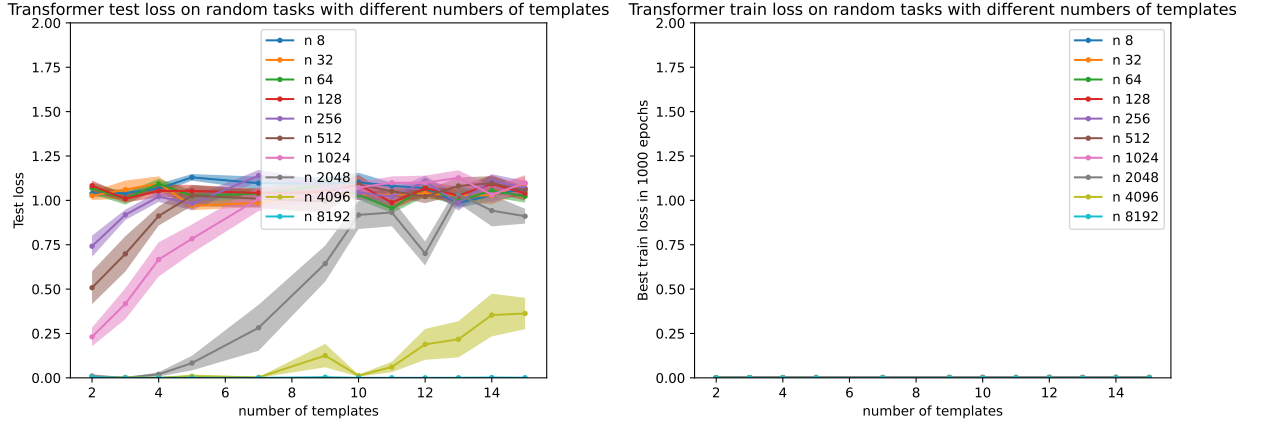


Figure 18: Performance on tasks corresponding of two random templates of length 5, labeled with 1,  $-1$ , respectively. Each template is sampled randomly from  $(\mathcal{W} \cup \mathcal{X})^5$ , conditioned on the two templates being distinct. We keep  $|\mathcal{W}| = 2$  and vary the regular token alphabet size  $|\mathcal{X}|$  between 0 and 2. Performance quickly improves as the regular token alphabet size increases.



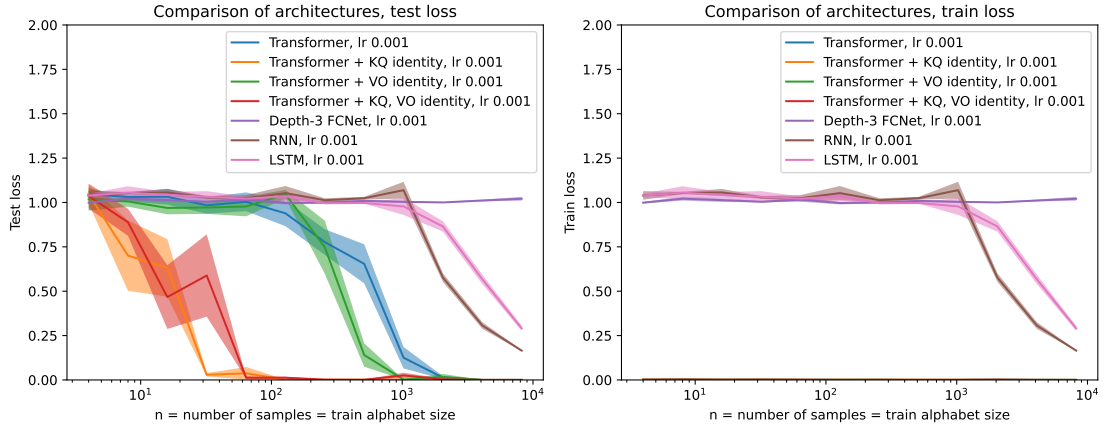


Figure 19: Different architectures on  $\alpha\beta\alpha$  vs.  $\alpha\beta\beta$  task. Transformer outperforms, especially with the reparametrization that prioritizes identities in heads.

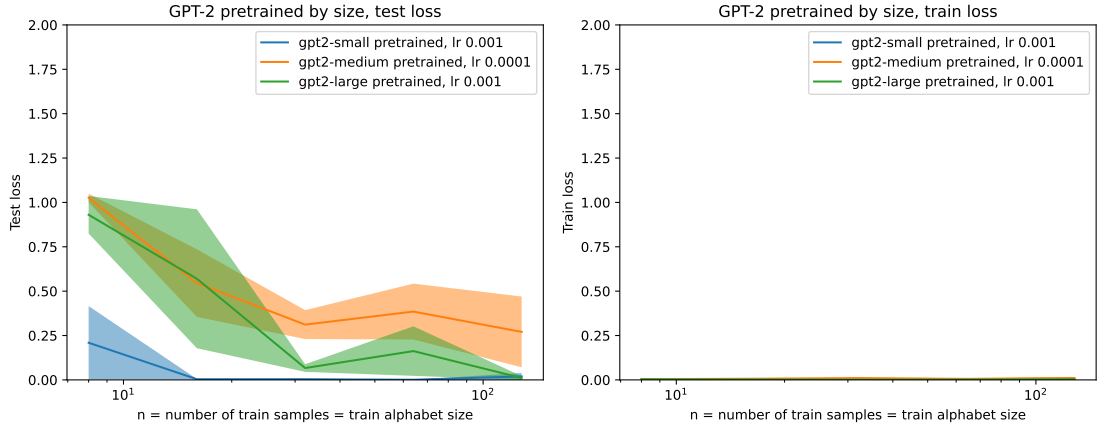


Figure 20: Pretrained GPT-2 of different sizes fine-tuned on  $\alpha\beta\alpha$  vs.  $\alpha\beta\beta$  task.

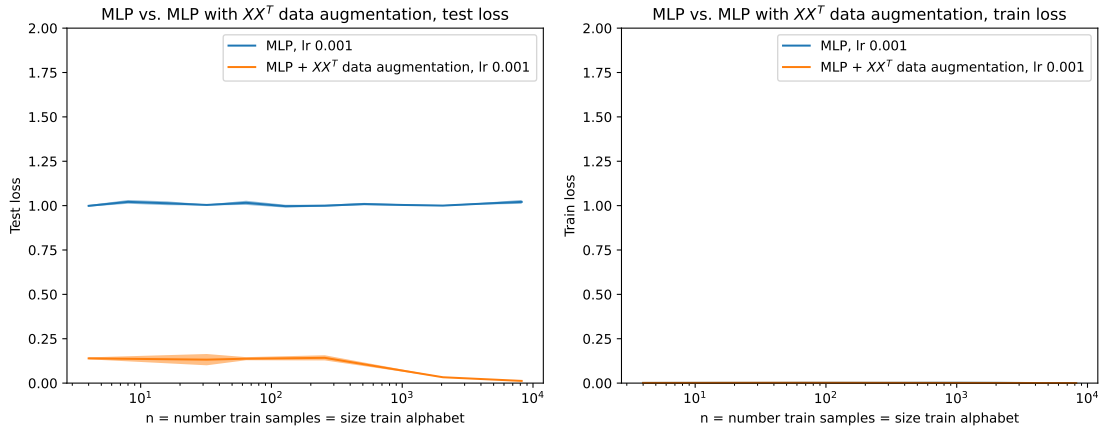


Figure 21: Test loss of MLP with  $XX^T$  data augmentation, where it is concatenated to input, versus MLP without data augmentation, versus transformer.

## C PROOF OF THEOREM 4.4

There are two main parts to the proof. First, in Section C.1 we establish a lemma with a sufficient condition for a kernel method to have good test loss. Second, in Section C.2 we prove that the transformer random features kernel  $K_{\text{trans}}$  satisfies this condition for almost any  $\beta, \gamma, b_1, b_2$  parameters. We conclude in Section C.3.

### C.1 PART 1. GENERAL SUFFICIENT CONDITION FOR GOOD TEST LOSS

We restrict ourselves to token-symmetric kernels, which are kernels whose values are unchanged if the tokens are relabeled by a permutation.

**Definition C.1** (Token-symmetric kernel).  $K$  is token-symmetric if for any permutation  $\pi : \mathcal{X} \rightarrow \mathcal{X}$  we have  $K(\mathbf{x}, \mathbf{y}) = K([\pi(x_1), \dots, \pi(x_k)], [\pi(y_1), \dots, \pi(y_k)])$ .

Token-symmetry is a mild condition, as most network architectures used in practice (including transformers) have token-symmetric neural tangent kernels at initialization. We emphasize that token-symmetry is not sufficient for good test loss since MLPs are a counterexample (see Appendix I.)

To state the sufficient condition for good test loss, let  $\{\mathbf{z}_1, \dots, \mathbf{z}_r\} = \text{supp}(\mu_{\text{tplt}})$  be the template distribution support. Define also the set  $\mathcal{R} = \cup_{i \in [k], j \in [r]} \{z_{j,i}\}$  of tokens that appear in the templates. Finally, define  $N \in \mathbb{R}^{r \times r}$  by

$$N_{ij} = K(\text{sub}(\mathbf{z}_i, s), \text{sub}(\mathbf{z}_j, s')), \quad (11)$$

where  $s, s' : \mathcal{W} \rightarrow \mathcal{X}$  are substitution maps satisfying

$$s(\mathcal{W}) \cap s'(\mathcal{W}) = \emptyset \quad \text{and} \quad s(\mathcal{W}) \cap \mathcal{R} = s'(\mathcal{W}) \cap \mathcal{R} = \emptyset. \quad (12)$$

One can check that because of the token-symmetry of the kernel  $K$ , the matrix  $N$  is uniquely-defined regardless of the substitution maps  $s, s'$  chosen, as long as they satisfy (12).

**Lemma C.2** (It suffices for  $N$  to be nonsingular). *If  $K$  is a token-symmetric kernel, and  $N$  is nonsingular, then kernel ridge regression achieves vanishing test loss.*

*Formally, there are constants  $c, C > 0$  and ridge regularization parameter  $\lambda > 0$  depending only on  $\mu_{\text{tplt}}, \sigma, |\mathcal{W}|, \|N^{-1}\|$  and  $\|K\|_\infty = \max_{\mathbf{x}} K(\mathbf{x}, \mathbf{x})$ , such that for any  $\mathbf{x}$  matching a template  $\mathbf{z} \in \text{supp}(\mu_{\text{tplt}})$  the kernel ridge regression estimator  $\hat{f}$  in (4) with kernel  $K$  satisfies*

$$|\hat{f}(\mathbf{x}) - f_*(\mathbf{z})| \leq C \sqrt{\frac{\log(1/\delta)}{n}} + C \sqrt{\frac{1}{\rho}},$$

*with probability at least  $1 - \delta - \exp(-cn)$  over the random samples.*

The proof is in Appendix D, but we develop an intuition here on why the nonsingularity of the matrix  $N$  is important. Let  $[n] = \mathcal{I}_1 \sqcup \mathcal{I}_2 \sqcup \dots \sqcup \mathcal{I}_n$  be the partition of the samples such that if  $i \in \mathcal{I}_j$  then sample  $(\mathbf{x}_i, y_i)$  is drawn by substituting the wildcards of template  $\mathbf{z}_j$  with substitution map  $s_i : \mathcal{W} \rightarrow \mathcal{X}$ . We show that for any string  $\mathbf{x}$  matching template  $\mathbf{z}_j$ , the kernel ridge regression solution (4) is approximately equal to the average of the labels of the samples corresponding to template  $j$ ,

$$\mathbf{y}^T (\hat{K} + \lambda \mathbf{I})^{-1} \mathbf{k}(\mathbf{x}) \approx \frac{1}{|\mathcal{I}_j|} \sum_{i \in \mathcal{I}_j} y_i \approx f_*(\mathbf{z}_j). \quad (13)$$

In order to see why this is true, consider the regime in which the sample diversity is very high, i.e.,  $\rho \gg 1$ . Since  $\rho$  is large, any particular token is highly unlikely to be substituted. This has the following implications:

- For most sample pairs  $i \neq i' \in [n]$ , the maps  $s_i$  and  $s_{i'}$  have disjoint range:  $s_i(\mathcal{W}) \cap s_{i'}(\mathcal{W}) = \emptyset$ .
- For most samples  $i \in [n]$ , the substituted tokens are not in the templates:  $s_i(\mathcal{W}) \cap \mathcal{R} = \emptyset$ .

These are the same conditions as in (6). So by the token-symmetry of the kernel, for most pairs of samples the empirical kernel matrix is given by  $N$ :

$$\hat{K}_{i,i'} := K(\mathbf{x}_i, \mathbf{x}_{i'}) = N_{j,j'} \text{ for most } i \in \mathcal{I}_j, i' \in \mathcal{I}_{j'}.$$

So if  $N$  is nonsingular, then  $\hat{K}$  has  $r$  large eigenvalues, and  $n - r$  much smaller eigenvalues. This turns out to be sufficient for (7) to hold. We refer the reader to Appendix D for more details.

## C.2 PART 2. ANALYZING THE TRANSFORMER RANDOM FEATURES KERNEL

We show that the transformer random features kernel  $K_{\text{trans}}$  satisfies the sufficient condition of Lemma C.2 for vanishing test loss. It is clear that the kernel is token-symmetric because the definition is invariant to the permutation relabelings of the tokens. The difficult part is to show that the matrix  $N_{\text{trans}} := N$  defined with kernel  $K = K_{\text{trans}}$  in (11) is nonsingular. The main challenge is that the transformer kernel does not have a known closed-form solution because of the softmax terms in its definition (3). Furthermore, the result is especially challenging to prove because it must hold for *any* collection of disjoint templates  $z_1, \dots, z_r$ .

We analyze the MLP layer and the attention layer of the transformer separately. We observe that a “weak” condition on  $K_{\text{attn}}$  can be lifted into the “strong” result that  $N_{\text{trans}}$  is nonsingular. Intuitively, as long as  $K_{\text{attn}}$  is not a very degenerate kernel, it is very unlikely that the MLP layer has the cancellations that would be needed to make  $N_{\text{trans}}$  nonsingular.

**Lemma C.3** (Nonsingularity of  $N_{\text{trans}}$ , restatement of Lemma 4.6). *Suppose for every non-identity permutation  $\tau \in S_r \setminus \{\text{id}\}$ ,*

$$\sum_{i \in [r]} K_{\text{attn}}(\text{sub}(z_i, s), \text{sub}(z_i, s')) \neq \sum_{i \in [r]} K_{\text{attn}}(\text{sub}(z_i, s), \text{sub}(z_{\tau(i)}, s')), \quad (14)$$

*where  $s, s'$  are the substitution maps in the definition of  $N_{\text{trans}}$  in (12). Let the MLP layer’s activation function be  $\phi(t) = \cos(b_1 t + b_2)$ . Then for almost any choice of  $b_1, b_2$  (except for a Lebesgue-measure-zero set), the matrix  $N_{\text{trans}}$  is nonsingular.*

This lemma is proved in Appendix E, by explicitly evaluating the Gaussian integral, which is possible since the activation function is the cosine function. Although in our proof we use the cosine activation function, we conjecture that this result should morally hold for sufficiently generic non-polynomial activation functions. Next, we prove the condition on  $N_{\text{attn}}$ .

**Lemma C.4** (Non-degeneracy of  $K_{\text{attn}}$ , restatement of Lemma 4.7). *The condition (14) holds for Lebesgue-almost any  $\beta, \gamma$ .*

The proof is in Appendix F. First, we prove the analyticity of the kernel  $K_{\text{attn}}$  in terms of the hyperparameters  $\beta$  and  $\gamma$  which control the softmax inverse temperature and the positional embeddings. Because of the identity theorem for analytic functions, it suffices to show at least one choice of hyperparameters  $\beta$  and  $\gamma$  satisfies (14) for all non-identity permutations  $\tau$ . Since  $K_{\text{attn}}$  does not have a closed-form solution, we find such a choice of  $\beta$  and  $\gamma$  by analyzing the Taylor-series expansion of  $K_{\text{attn}}$  around  $\beta = 0$  and  $\gamma = 0$  up to order-10 derivatives, which happens to suffice.

## C.3 CONCLUDING THE PROOF OF THEOREM 4.4

By Lemma C.2, it suffices to prove the nonsingularity of the matrix  $N_{\text{trans}}$  defined in (11) with kernel  $K = K_{\text{trans}}$ . Lemma 4.6 gives a condition for nonsingularity that holds for almost any  $b_1, b_2$ . Lemma 4.7 proves this condition for almost any  $\beta, \gamma$ . Therefore, Theorem 4.4 follows.

## D SUFFICIENT CONDITION FOR KERNEL METHOD TO GENERALIZE ON UNSEEN SYMBOLS (PROOF OF LEMMA C.2)

We restate and prove Lemma C.2. Let  $K$  be a token-symmetric kernel as in Definition C.1. Let  $\mu_{\text{tmpl}}$  be a distribution supported on disjoint templates  $z_1, \dots, z_r$  and define  $\mathcal{R} = \cup_{i \in [r], j \in [k]} \{z_{i,j}\}$ . Recall the definition of the matrix  $N \in \mathbb{R}^{r \times r}$  with

$$N_{i,i'} = K(\text{sub}(z_i, s), \text{sub}(z_{i'}, s')).$$

for substitution maps  $s : \mathcal{W} \rightarrow \mathcal{X}$ ,  $s' : \mathcal{W} \rightarrow \mathcal{X}$  satisfying  $s(\mathcal{W}) \cap s'(\mathcal{W}) = s(\mathcal{W}) \cap \mathcal{R} = s'(\mathcal{W}) \cap \mathcal{R} = \emptyset$ . Recall that this is well-defined by the token-symmetry of the kernel  $K$ .

**Lemma D.1** (Restatement of Lemma C.2). *Suppose that  $K$  is token-symmetric and  $N$  is nonsingular. Then there are constants  $0 < c < C$  and  $0 < c' < C'$  depending only on  $\mu_{\text{tmpl}}$ ,  $\sigma$ ,  $|\mathcal{W}|$ ,  $\|N^{-1}\|$  and  $\|K\|_{\infty} = \max_{\mathbf{x}} K(\mathbf{x}, \mathbf{x})$  such that the following holds. Consider any regularization parameter  $\lambda \in [c'n, C'n]$ , and any string  $\mathbf{x}$  matching template  $\mathbf{z} \in \text{supp}(\mu_{\text{tmpl}})$ . Then with probability*

$\geq 1 - \delta - \exp(-cn)$ , the kernel ridge regression estimator  $\hat{f}$  achieves good accuracy on  $\mathbf{x}$ :

$$|\hat{f}(\mathbf{x}) - f_*(\mathbf{z})| \leq C\sqrt{\frac{\log(1/\delta)}{n}} + C\sqrt{\frac{1}{\rho}}.$$

*Proof.* Note that some proofs of helper claims are deferred to Section D.1. Let  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  be the samples seen by the kernel method. We know from (4) that kernel ridge regression outputs the estimator

$$\hat{f}(\mathbf{x}) = \mathbf{y}^T (\hat{\mathbf{K}} + \lambda \mathbf{I})^{-1} \mathbf{v}(\mathbf{x}), \quad (\text{Kernel ridge regression})$$

where the empirical kernel matrix  $\hat{\mathbf{K}} \in \mathbb{R}^{n \times n}$  is

$$\hat{K}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j),$$

and  $\mathbf{y} = [y_1, \dots, y_n]$ , and  $\mathbf{v}(\mathbf{x}) = [K(\mathbf{x}_1, \mathbf{x}), \dots, K(\mathbf{x}_n, \mathbf{x})] \in \mathbb{R}^n$ .

**Idealized estimator when sample diversity is high** If the sample diversity is sufficiently high, then for most pairs of samples  $i \neq i' \in [n]$ , it will be the case that  $\mathbf{x}_i$  and  $\mathbf{x}_{i'}$  do not share any of the wildcard substitution tokens. In other words, the wildcard substitution map used to form  $\mathbf{x}_i$  will have disjoint range from the wildcard substitution map used to form  $\mathbf{x}_{i'}$ . This means that we should expect the estimator  $\hat{f}$  to perform similarly to the following idealized estimator:

$$\hat{f}^{ideal}(\mathbf{x}) = \mathbf{y}^T (\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x}), \quad (15)$$

where  $\hat{\mathbf{K}}^{ideal} \in \mathbb{R}^{n \times n}$  and  $\mathbf{v}^{ideal}(\mathbf{x}) \in \mathbb{R}^n$  are idealized versions of  $\hat{\mathbf{K}}$  and  $\mathbf{v}(\mathbf{x})$ , formed below. They correspond to the limit of infinitely-diverse samples, when all token substitution maps have disjoint range. For each  $j \in [r]$ , let  $\mathcal{I}_j \subseteq [n]$  be the indices of samples  $\mathbf{x}_i$  formed by substituting from template  $\mathbf{z}_j$ . For any  $i \in \mathcal{I}_j, i' \in \mathcal{I}_{j'}$ , let

$$\hat{K}_{i,i'}^{ideal} = N_{j,j'}, \quad (16)$$

Also, similarly define  $\mathbf{v}^{ideal}(\mathbf{x}) \in \mathbb{R}^n$ . For any  $i \in \mathcal{I}_j$ , let

$$v_i^{ideal}(\mathbf{x}) = K(\text{sub}(\mathbf{z}_j, s), \mathbf{x}), \quad (17)$$

where  $s : \mathcal{W} \rightarrow \mathcal{X}$  is a substitution map with  $s(\mathcal{W}) \cap \mathcal{R} = s(\mathcal{W}) \cap \{\mathbf{x}_i\}_{i \in [k]} = \emptyset$ , i.e., it does not overlap with the templates or with  $\mathbf{x}$  in the tokens substituted for the wildcards. The expressions (16) and (17) are well-defined because of the token-symmetry of the kernel.

If the sample diversity is high, then we show that the idealized estimator  $\hat{f}^{ideal}$  is indeed close to the kernel ridge regression solution  $\hat{f}$ .

**Claim D.2** (Idealized estimator is good approximation to true estimator). *Suppose  $\|K\|_\infty = \max_{\mathbf{x}} |K(\mathbf{x}, \mathbf{x})| < \infty$ . Then there are constants  $C, c > 0$  depending only on  $|\mathcal{W}|, \|K\|_\infty, k, r$  such that the following holds. For any  $\mathbf{x}$ , with probability at least  $1 - \exp(-cn)$ ,*

$$|\hat{f}^{ideal}(\mathbf{x}) - \hat{f}(\mathbf{x})| \leq \frac{C}{\lambda} + \frac{Cn}{\lambda\sqrt{\rho}},$$

where  $\rho$  is defined in Definition 4.2 and measures the diversity of the substitution map distribution.

**Analyzing the idealized estimator using its block structure** The matrix  $\hat{\mathbf{K}}^{ideal}$  has block structure with blocks  $\mathcal{I}_1, \dots, \mathcal{I}_r$ . Namely, it equals  $\hat{K}_{i,i'}^{ideal} = N_{j,j'}$  for all  $i \in \mathcal{I}_j, i' \in \mathcal{I}_{j'}$ . Similarly,  $\mathbf{v}^{ideal}(\mathbf{x})$  also has block structure with blocks  $\mathcal{I}_1, \dots, \mathcal{I}_r$ . This structure allows us to analyze estimator  $\hat{f}^{ideal}$  and to prove its accuracy.

In order to analyze the estimator, we prove the following technical claim. The interpretation of this claim is that if  $\mathbf{x}$  matches template  $\mathbf{z}_a$ , then  $\mathbf{v}^{ideal}(\mathbf{x})$  is equal to any of the rows in  $\hat{\mathbf{K}}^{ideal}$  that correspond to template  $a$ . In other words, we should have  $(\hat{\mathbf{K}}^{ideal})^+ \mathbf{v}^{ideal}(\mathbf{x}) = \mathbf{1}_{\mathcal{I}_a} / |\mathcal{I}_a|$ , which is the indicator vector for samples that come from template  $a$ . The following technical claim is a more robust version of this observation.

**Claim D.3.** Let  $\mathbf{x}$  be a string that matches template  $\mathbf{z}_a$ . Suppose that  $0 < \lambda < \tau := \min_{j \in [r]} |\mathcal{I}_j| / \|\mathbf{N}^{-1}\|$ . Then  $(\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})$  is invertible and the following are satisfied

$$\|(\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x})\| \leq \sqrt{\frac{1}{|\mathcal{I}_a|}} \left( \frac{\tau}{\tau - \lambda} \right),$$

and, letting  $\mathbf{1}_{\mathcal{I}_a} \in \mathbb{R}^n$  be the indicator vector for set  $\mathcal{I}_a$ ,

$$\left\| \frac{\mathbf{1}_{\mathcal{I}_a}}{|\mathcal{I}_a|} - (\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x}) \right\| \leq \sqrt{\frac{1}{|\mathcal{I}_a|}} \left( \frac{\tau}{\tau - \lambda} - 1 \right).$$

Using the above technical claim, we can prove that  $\hat{f}^{ideal}$  is an accurate estimator. The insight is that since  $(\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x})$  is approximately the indicator vector  $\mathbf{1}_{\mathcal{I}_a}/|\mathcal{I}_a|$  for samples corresponding to template  $a$ , the output of the idealized estimator is the average of the labels for samples corresponding to template  $a$ .

**Claim D.4** (Idealized estimator gets vanishing test loss on unseen symbols). *There are  $c, C > 0$  depending only on  $|\mathcal{W}|, \mu_{\text{tmplt}}, \sigma$  such that the following holds for any  $0 < \lambda < cn$ . Let  $\mathbf{x}$  be any string that matches template  $\mathbf{z} \in \text{supp}(\mu_{\text{tmplt}})$ . Then, for any  $\delta > 0$ , with probability  $\geq 1 - \delta - \exp(-cn)$  over the random samples, the idealized estimator has error upper-bounded by*

$$|\hat{f}^{ideal}(\mathbf{x}) - f_*(\mathbf{z})| \leq C \sqrt{\frac{\log(1/\delta)}{n}}.$$

*Proof of Claim D.4.* Let  $E_1$  be the event that  $n\mu_{\text{tmplt}}(\mathbf{z}_j) \geq |\mathcal{I}_j|/2$  for all  $j \in [r]$ , i.e., all templates are well-represented in the dataset. By a Hoeffding bound,

$$\mathbb{P}[E_1] \geq 1 - \exp(-n\mu_{\text{tmplt}}(\mathbf{z}_a)/2).$$

Suppose that  $\mathbf{x}$  matches template  $\mathbf{z}_a$ . By Claim D.3, under event  $E_1$ , there is a constant  $C > 0$  such that

$$\begin{aligned} |\hat{f}^{ideal}(\mathbf{x}) - f_*(\mathbf{z}_a)| &= |\mathbf{y}^T (\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x}) - f_*(\mathbf{z}_a)| \\ &\leq |\mathbf{y}^T \frac{\mathbf{1}_{\mathcal{I}_a}}{|\mathcal{I}_a|} - f_*(\mathbf{z}_a)| + \sqrt{\frac{1}{|\mathcal{I}_a|}} \left( \frac{\tau}{\tau - \lambda} - 1 \right) \\ &\leq |\mathbf{y}^T \frac{\mathbf{1}_{\mathcal{I}_a}}{|\mathcal{I}_a|} - f_*(\mathbf{z}_a)| + C \sqrt{\frac{1}{n}}. \end{aligned}$$

We conclude since  $\mathbb{P}[|\mathbf{y}^T \frac{\mathbf{1}_{\mathcal{I}_a}}{|\mathcal{I}_a|} - f_*(\mathbf{z}_a)| > C \sqrt{\frac{\log(1/\delta)}{n}} \mid E_1] \leq \delta$  by a tail bound for Gaussians.  $\square$

**Putting the elements together to conclude the proof of the lemma** Combined, Claims D.2 and D.4 imply the lemma if we take  $\lambda = \Theta(n)$ , then we obtain error  $O(\sqrt{\log(1/\delta)/n} + \sqrt{1/\rho})$  with probability at least  $1 - \delta - \exp(-\Omega(n))$ .  $\square$

#### D.1 DEFERRED PROOFS OF CLAIMS

*Proof of Claim D.3.* Let  $\mathbf{w}_1, \dots, \mathbf{w}_n$  be an orthogonal basis of eigenvectors for  $\hat{\mathbf{K}}^{ideal}$  with eigenvalues  $\nu_1, \dots, \nu_n$ . Notice that these are also eigenvectors of  $\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I}$ . Because of the block structure of  $\hat{\mathbf{K}}^{ideal}$ , its eigenvectors and eigenvalues have a simple form. Define

$$\mathbf{M} = \text{diag}([\sqrt{|\mathcal{I}_1|}, \dots, \sqrt{|\mathcal{I}_r|}]) \mathbf{N} \text{diag}([\sqrt{|\mathcal{I}_1|}, \dots, \sqrt{|\mathcal{I}_r|}]).$$

The nonzero eigenvalues of  $\hat{\mathbf{K}}^{ideal}$  correspond to the nonzero eigenvalues of  $\mathbf{M}$ , because for any eigenvector  $\mathbf{u} \in \mathbb{R}^r$  of  $\mathbf{M}$  there is a corresponding eigenvector of  $\hat{\mathbf{K}}^{ideal}$  with the same eigenvalue

by letting each of the blocks  $\mathcal{I}_j$  consist of copies of the entry  $u_j/\sqrt{|\mathcal{I}_j|}$ . Therefore, all nonzero eigenvalues of  $\hat{\mathbf{K}}^{-1}$  have magnitude at least

$$|\nu_1|, \dots, |\nu_n| \geq 1/\|\mathbf{M}^{-1}\| \geq \min_{j \in [r]} |\mathcal{I}_j|/\|\mathbf{N}^{-1}\| = \tau > \lambda.$$

So  $\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I}$  is invertible, which is the first part of the claim. Write  $\frac{\mathbf{1}_{\mathcal{I}_a}}{|\mathcal{I}_a|}$  in the eigenbasis as

$$\frac{\mathbf{1}_{\mathcal{I}_a}}{|\mathcal{I}_a|} = \sum_i c_i \mathbf{w}_i,$$

for some coefficients  $c_i$ . By construction,

$$\mathbf{v}^{ideal}(\mathbf{x}) = \hat{\mathbf{K}}^{ideal} \frac{\mathbf{1}_{\mathcal{I}_a}}{|\mathcal{I}_a|} = \sum_i \nu_i c_i \mathbf{w}_i,$$

so

$$\begin{aligned} \|(\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x})\|^2 &= \left\| \sum_i \frac{\nu_i}{\nu_i + \lambda} c_i \mathbf{w}_i \right\|^2 = \sum_i \left( \frac{\nu_i}{\nu_i + \lambda} \right)^2 c_i^2 \\ &\leq \max_i \left( \frac{\nu_i}{\nu_i + \lambda} \right)^2 \frac{1}{|\mathcal{I}_a|} \leq \max_i \left( \frac{\tau}{\tau - \lambda} \right)^2. \end{aligned}$$

Similarly,

$$\begin{aligned} \left\| \frac{\mathbf{1}_{\mathcal{I}_a}}{|\mathcal{I}_a|} - (\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x}) \right\|^2 &= \left\| \sum_i \left( 1 - \frac{\nu_i}{\nu_i + \lambda} \right) c_i \mathbf{w}_i \right\|^2 = \sum_i \left( 1 - \frac{\nu_i}{\nu_i + \lambda} \right)^2 c_i^2 \\ &\leq \max_i \left( 1 - \frac{\nu_i}{\nu_i + \lambda} \right)^2 \frac{1}{|\mathcal{I}_a|} \leq \max_i \left( 1 - \frac{\tau}{\tau - \lambda} \right)^2. \end{aligned}$$

□

**Claim D.5** (Bound on difference between kernel regressions). *Suppose that  $\hat{\mathbf{K}}$  is p.s.d and that  $(\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x})$  is well-defined. Then, for any  $\lambda > 0$ ,*

$$|\hat{f}^{ideal}(\mathbf{x}) - \hat{f}(\mathbf{x})| \leq \frac{\|\mathbf{y}\|}{\lambda} (\|\mathbf{v}^{ideal}(\mathbf{x}) - \mathbf{v}(\mathbf{x})\| + \|\hat{\mathbf{K}} - \hat{\mathbf{K}}^{ideal}\| \|(\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x})\|)$$

*Proof of Claim D.5.* By triangle inequality,

$$\begin{aligned} |\hat{f}(\mathbf{x}) - \hat{f}^{ideal}(\mathbf{x})| &= \|\mathbf{y}^T (\hat{\mathbf{K}} + \lambda \mathbf{I})^{-1} \mathbf{v}(\mathbf{x}) - \mathbf{y}^T (\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x})\| \\ &\stackrel{(a)}{\leq} \|\mathbf{y}\| \cdot \underbrace{\|(\hat{\mathbf{K}} + \lambda \mathbf{I})^{-1} \mathbf{v}(\mathbf{x}) - (\hat{\mathbf{K}} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x})\|}_{\text{Term 1}} \\ &\quad + \|\mathbf{y}\| \cdot \underbrace{\|(\hat{\mathbf{K}} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x}) - (\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x})\|}_{\text{Term 2}} \end{aligned}$$

The first term can be upper-bounded because  $\|(\hat{\mathbf{K}} + \lambda \mathbf{I})^{-1}\| \leq \|(\lambda \mathbf{I})^{-1}\| = 1/\lambda$ , so

$$\text{Term 1} \leq \frac{\|\mathbf{v}^{ideal}(\mathbf{x}) - \mathbf{v}(\mathbf{x})\|}{\lambda}$$

The second term can be upper-bounded by

$$\begin{aligned} \text{Term 2} &= \|(\hat{\mathbf{K}} + \lambda \mathbf{I})^{-1} ((\hat{\mathbf{K}} + \lambda \mathbf{I})(\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} - (\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})(\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1}) \mathbf{v}^{ideal}(\mathbf{x})\| \\ &= \|(\hat{\mathbf{K}} + \lambda \mathbf{I})^{-1} (\hat{\mathbf{K}} - \hat{\mathbf{K}}^{ideal}) (\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x})\| \\ &\leq \frac{1}{\lambda} \|\hat{\mathbf{K}} - \hat{\mathbf{K}}^{ideal}\| \|(\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x})\|. \end{aligned}$$

□



*Proof of Claim D.2.* Let  $E_1$  be the event that  $|I_j| \geq n\mu_{\text{tmpl}}(z_j)$  for all  $j \in [r]$ . By Hoeffding, there is a constant  $c > 0$  such that  $\mathbb{P}[E_1] \geq 1 - \exp(-cn)$ . By Claim D.3, under event  $E_1$ , there is a constant  $C > 0$  such that

$$\|(\hat{\mathbf{K}}^{ideal} + \lambda \mathbf{I})^{-1} \mathbf{v}^{ideal}(\mathbf{x})\| \leq \frac{C}{\sqrt{n}}. \quad (18)$$

Next, recall the parameter  $\rho$  used to measure the spread of the substitution map distributions  $\{\mu_{sub,z}\}_{z \in \text{supp}(\mu_{\text{tmpl}})}$ , as defined in (4.2). For each  $i \in [n]$ , let  $s_i : \mathcal{W} \rightarrow \mathcal{X}$  be the substitution map used to generate the sample  $\mathbf{x}_i$ . Let  $P_1$  be the number of samples  $(i, i')$  such that their substitution maps overlap, or have range that overlaps with the regular tokens in the templates. Formally:

$$P_1 = |\{1 \leq i < i' \leq n : s_i(\mathcal{W}) \cap s_{i'}(\mathcal{W}) \neq \emptyset \text{ or } s_i(\mathcal{W}) \cap \mathcal{R} \neq \emptyset \text{ or } s_{i'}(\mathcal{W}) \cap \mathcal{R} \neq \emptyset\}|.$$

Similarly, let  $P_2$  be the number of samples  $(i, i')$  such that their substitution maps overlap with that used to generate  $\mathbf{x}$ , or they overlap with the regular tokens in the templates:

$$P_2 = |\{1 \leq i \leq n : s_i(\mathcal{W}) \cap \mathcal{R} \neq \emptyset \text{ or } s_i(\mathcal{W}) \cap \{x_j\}_{j \in [k]} \neq \emptyset\}|.$$

By the definition of  $\rho$ , we can upper-bound the expected number of “bad” pairs  $P_1$  and “bad” indices  $P_2$  by:

$$\begin{aligned} \mathbb{E}[P_1] &\leq \left( \sum_{i, i' \in [n]} \sum_{w, w' \in \mathcal{W}} \mathbb{P}[s_i(w) = s_{i'}(w')] \right) + n \sum_{i \in [n]} \sum_{t \in \mathcal{R}} \mathbb{P}[t \in s_i(\mathcal{W})] \leq \frac{Cn^2}{\rho} + \frac{Cn}{\rho} \leq \frac{Cn^2}{\rho} \\ \mathbb{E}[P_2] &\leq \sum_{i \in [n]} \sum_{t \in \{x_j\}_{j \in [k]} \cup \mathcal{R}} \mathbb{P}[t \in s_i(\mathcal{W})] \leq \frac{Cn}{\rho}. \end{aligned}$$

By Hoeffding’s inequality, the event  $E_2$  that  $P_1 \leq \frac{Cn^2}{\rho}$  and  $P_2 \leq \frac{Cn}{\rho}$  occurs with probability  $\geq 1 - \exp(-cn)$ . Under event  $E_2$ ,

$$\|\hat{\mathbf{K}} - \hat{\mathbf{K}}^{ideal}\| \leq C + Cn/\sqrt{\rho} \quad \text{and} \quad \|\mathbf{v}(\mathbf{x}) - \mathbf{v}^{ideal}(\mathbf{x})\| \leq C\sqrt{n/\rho}. \quad (19)$$

By Claim D.5 and (18) and (19), under events  $E_1, E_2$ , and using that  $\|\mathbf{y}\| \leq C\sqrt{n}$ , we have

$$|\hat{f}^{ideal}(\mathbf{x}) - \hat{f}(\mathbf{x})| \leq \frac{C\sqrt{n}}{\lambda} (C\sqrt{n/\rho} + (C + Cn/\sqrt{\rho}) \frac{C}{\sqrt{n}}) \leq \frac{C(1+n)}{\lambda\sqrt{\rho}}.$$

□

## E NONSINGULARITY OF RANDOM FEATURES AFTER MLP LAYER (PROOF OF LEMMA 4.6)

Consider a kernel  $K_2$  formed from a kernel  $K_1$  as follows:

$$K_2(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{u, v \sim \Sigma_1(\mathbf{x}, \mathbf{y})} [\phi(u)\phi(v)], \quad \Sigma_1(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} K_1(\mathbf{x}, \mathbf{y}) & K_1(\mathbf{x}, \mathbf{y}) \\ K_1(\mathbf{x}, \mathbf{y}) & K_1(\mathbf{x}, \mathbf{y}) \end{bmatrix}.$$

Here  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is a nonlinear activation function. Such a random features kernel arises in a neural network architecture by appending an infinite-width MLP layer with Gaussian initialization to a neural network with random features with kernel  $K_1$ .

We wish to prove that a certain matrix  $N \in \mathbb{R}^{r \times r}$  given by

$$N_{ij} = K_2(\mathbf{x}_i, \mathbf{y}_j), \quad (20)$$

is nonsingular, where  $\mathbf{x}_1, \dots, \mathbf{x}_r, \mathbf{y}_1, \dots, \mathbf{y}_r$  are inputs. The intuition is that if  $\phi$  is a “generic” activation function, then only a weak condition on  $K_1$  is required for the matrix  $N$  to be invertible. We provide a general lemma that allows us to guarantee the invertibility if the activation function is a shifted cosine, although we conjecture such a result to be true for most non-polynomial activation functions  $\phi$ . This is a generalization of Lemma 4.6, so it implies Lemma 4.6.

**Lemma E.1** (Criterion for invertibility of  $N$ ). *Consider the matrix  $N \in \mathbb{R}^{r \times r}$  defined in (20) where  $\mathbf{x}_1, \dots, \mathbf{x}_r$  and  $\mathbf{y}_1, \dots, \mathbf{y}_r$  are inputs. Suppose that for all nontrivial permutations  $\tau \in S_r \setminus \{\text{id}\}$  we have*

$$\sum_{i \in [r]} K_1(\mathbf{x}_i, \mathbf{y}_i) \neq \sum_{i \in [r]} K_1(\mathbf{x}_i, \mathbf{y}_{\tau(i)}). \quad (21)$$

*Suppose also that the MLP activation function is  $\phi(t) = \cos(kt + c)$  for two hyperparameters  $k, c$ . Then,  $N$  is nonsingular for all  $(k, c) \in \mathbb{R}^2$  except for a Lebesgue-measure-zero subset of  $\mathbb{R}^2$ .*

*Proof.* Let  $f(k, c) := \det(N)$ . We wish to show that  $\{(k, c) : f(k, c) = 0\}$  is a measure-zero set. By Claim E.2, is an analytic function of  $c$  and  $k$ , and by the identity theorem for analytic functions (Mityagin, 2020), it suffices to show that  $f \not\equiv 0$ . Fixing  $c = \pi/4$ , by Claim E.2,

$$K_2(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \exp\left(-\frac{k^2}{2}(K_1(\mathbf{x}, \mathbf{x}) + K_1(\mathbf{y}, \mathbf{y}) - 2K_1(\mathbf{x}, \mathbf{y}))\right).$$

Therefore

$$\begin{aligned} f(k, \pi/4) &= \sum_{\tau \in S_r} \text{sgn}(\tau) \prod_{i \in [r]} K_2(\mathbf{x}_i, \mathbf{y}_{\tau(i)}) \\ &= e^{-\frac{k^2}{2}(\sum_{i \in [r]} K_1(\mathbf{x}_i, \mathbf{x}_i) + K_1(\mathbf{y}_i, \mathbf{y}_i))} \sum_{\tau \in S_r} \text{sgn}(\tau) \exp(k^2 \sum_{i \in [r]} K_1(\mathbf{x}_i, \mathbf{y}_{\tau(i)})). \end{aligned}$$

It remains to prove that as a function of  $k$  we have

$$\sum_{\tau \in S_r} \text{sgn}(\tau) \exp(k^2 \sum_{i \in [r]} K_1(\mathbf{x}_i, \mathbf{y}_{\tau(i)})) \neq 0,$$

This holds because for any distinct  $c_1, \dots, c_l$  the functions  $\exp(c_1 t), \dots, \exp(c_l t)$  are linearly independent functions of  $t$ , since their Wronskian is a rescaled Vandermonde determinant

$$\begin{aligned} \begin{vmatrix} \exp(c_1 t) & \dots & \exp(c_l t) \\ \frac{d}{dx} \exp(c_1 t) & \dots & \frac{d}{dx} \exp(c_l t) \\ \vdots & & \vdots \\ \frac{d^{l-1}}{dt^{l-1}} \exp(c_1 t) & \dots & \frac{d^{l-1}}{dt^{l-1}} \exp(c_l t) \end{vmatrix} &= \exp\left(\sum_{i=1}^l c_i t\right) \begin{vmatrix} 1 & \dots & 1 \\ c_1 & \dots & c_l \\ \vdots & & \vdots \\ c_1^{l-1} & \dots & c_l^{l-1} \end{vmatrix} \\ &= \exp\left(\sum_{i=1}^l c_i t\right) \prod_{1 \leq i < j \leq l} (c_j - c_i) \neq 0 \end{aligned}$$

□

Below is the technical claim used in the proof of the lemma.

**Claim E.2.** *Let  $U, V \sim N(0, \begin{bmatrix} a & \rho \\ \rho & b \end{bmatrix})$ . Then for any  $k, c \in \mathbb{R}$ ,*

$$\mathbb{E}[\cos(kU + c) \cos(kV + c)] = \frac{1}{2} e^{-\frac{1}{2} k^2 (a+b)} (e^{-k^2 \rho} \cos(2c) + e^{k^2 \rho}).$$

*Proof.* By Mathematica, we have the following Gaussian integrals

$$\begin{aligned} \mathbb{E}[e^{ikU+ikV}] &= \mathbb{E}[e^{-ikU-ikV}] = e^{-\frac{1}{2} k^2 (a+b+2\rho)}, \\ \mathbb{E}[e^{ikU-ikV}] &= \mathbb{E}[e^{-ikU+ikV}] = e^{-\frac{1}{2} k^2 (a+b-2\rho)}. \end{aligned}$$

Since  $\cos(kt + c) = (e^{ikt+ic} + e^{-ikt-ic})/2$ ,

$$\begin{aligned} \mathbb{E}[\cos(kU + c) \cos(kV + c)] &= \frac{1}{4} \mathbb{E}[(e^{ikU+ic} + e^{-ikU-ic})(e^{ikV+ic} + e^{-ikV-ic})] \\ &= \frac{1}{4} (e^{-\frac{1}{2} k^2 (a+b+2\rho)} (e^{2ic} + e^{-2ic}) + 2e^{-\frac{1}{2} k^2 (a+b-2\rho)}) \\ &= \frac{1}{2} e^{-\frac{1}{2} k^2 (a+b)} (e^{-k^2 \rho} \cos(2c) + e^{k^2 \rho}). \end{aligned}$$

□

## F ANALYSIS OF ATTENTION LAYER FEATURES (PROOF OF LEMMA 4.7)

For any inputs  $X, Y$ , we write the kernel of the random features of the attention layer as

$$K_{\text{attn}}(X, Y) = \mathbb{E}_{\mathbf{m}(X), \mathbf{m}(Y)} [\text{smax}(\beta \mathbf{m}(X))^T (X Y^T + \gamma^2 \mathbf{I}) \text{smax}(\beta \mathbf{m}(Y))] \\ \mathbf{m}(X), \mathbf{m}(Y) \sim N(\mathbf{0}, \begin{bmatrix} X X^T + \gamma^2 \mathbf{I} & X Y^T + \gamma^2 \mathbf{I} \\ Y X^T + \gamma^2 \mathbf{I} & Y Y^T + \gamma^2 \mathbf{I} \end{bmatrix}),$$

as stated Section 4.1; see also Section H for the derivation of this kernel in the infinite-width limit of the transformer architecture. For shorthand, we write  $\kappa_{X, Y}(\beta, \gamma) = K_{\text{attn}}(X, Y)$  to emphasize the attention kernel’s dependence on the hyperparameters  $\beta$  and  $\gamma$  which control the softmax’s inverse temperature and the weight of the positional embeddings, respectively.

We prove Lemma 4.7, which is that  $K_{\text{attn}}$  satisfies the property (8) required by Lemma 4.6 for the transformer random features kernel to succeed at the template task.

Namely, consider any disjoint templates  $z_1, \dots, z_r$  and two substitution maps  $s, s' : \mathcal{W} \rightarrow \mathcal{X}$

- that have disjoint range:  $s(\mathcal{W}) \cap s'(\mathcal{W}) = \emptyset$ ,
- and the substituted tokens do not overlap with any of the tokens in the templates:  $s(\mathcal{W}) \cap \mathcal{R} = s'(\mathcal{W}) \cap \mathcal{R} = \emptyset$  where  $\mathcal{R} = \cup_{i \in [r], j \in [k]} \{z_j^{(i)}\}$ .

Then we define  $X_i, Y_i \in \mathbb{R}^{k \times m}$  to be the strings (where we abuse notation slightly by viewing them as matrices with one-hot rows) after substituting  $z_i$  by  $s, s'$  respectively:

$$X_i = \text{sub}(z_i, s) \quad Y_i = \text{sub}(z_i, s').$$

**Lemma F.1** (Restatement of Lemma 4.7). *Define  $g_\tau(\beta, \gamma) = \sum_{i \in [r]} \kappa_{X_i, Y_{\tau(i)}}(\beta, \gamma)$ . Then for all but a Lebesgue-measure-zero set of  $(\beta, \gamma) \in \mathbb{R}^2$  we have  $g_{\text{id}}(\beta, \gamma) \neq g_\tau(\beta, \gamma)$  for all permutations  $\tau \neq \text{id}$ .*

No closed-form expression is known for  $\kappa_{X, Y}(\beta, \gamma)$ , so our approach is to analyze its Taylor series expansion around  $\beta = \gamma = 0$ . Our proof proceeds in stages, where, in each stage, we examine a higher derivative and progressively narrow the set of  $\tau$  that might possibly have  $g_\tau(\beta, \gamma) = g_{\text{id}}(\beta, \gamma)$ . In Section F.1, we list certain low-order derivatives of  $\kappa_{X, Y}(\beta, \gamma)$  that will be sufficient for our analysis. In Section F.2, we analyze some of the terms in these expressions. In Section F.3 we put the previous lemmas together to prove Lemma F.1.

To avoid notational overload, in this section we will not use bolded notation to refer to the matrices  $X, Y$ , but rather use the lowercase  $X, Y$ .

### F.1 LOW-ORDER DERIVATIVES OF ATTENTION KERNEL

In the following table we collect several relevant derivatives of  $\frac{\partial^i}{\partial \beta^i} \frac{\partial^j}{\partial \gamma^j} \kappa_{X, Y}(0, 0)$  for  $i \leq 6$  and  $j \leq 4$ . For each  $i, j$  we use  $c_1, c_2, \dots$  to denote constants that depend only on  $k$ , and on the derivative  $i, j$  being computed. Certain constants that are important for the proof are provided explicitly. These derivatives were computed using a Python script available in our code. The colors are explained in Section F.2.

Derivative	Expansion
$\kappa_{X, Y}(0, 0) =$	$+c_1 \mathbf{1}^T X Y^T \mathbf{1}$
$\frac{\partial^2}{\partial \beta^2} \frac{\partial^2}{\partial \gamma^2} \kappa_{X, Y}(0, 0) =$	$+c_1 \mathbf{1}^T X Y^T \mathbf{1} + c_2 \text{tr}(X Y^T)$

$$\begin{aligned}
\frac{\partial^4}{\partial \beta^4} \kappa_{X,Y}(0,0) = & +c_1 \mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1} + c_2 \mathbf{1}^T \textcolor{pink}{XX}^T \textcolor{pink}{XY}^T \mathbf{1} + c_3 \mathbf{1}^T \textcolor{pink}{XY}^T \textcolor{pink}{YY}^T \mathbf{1} \\
& + c_4 \mathbf{1}^T \textcolor{pink}{XX}^T \textcolor{pink}{XX}^T \textcolor{pink}{XY}^T \mathbf{1} + c_5 (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) \\
& + c_6 \mathbf{1}^T \textcolor{pink}{XY}^T \textcolor{blue}{YX}^T \textcolor{pink}{XY}^T \mathbf{1} + c_7 (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) \\
& + c_8 \mathbf{1}^T \textcolor{blue}{YX}^T \textcolor{pink}{XY}^T \textcolor{pink}{YY}^T \mathbf{1} + c_9 (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) \\
& + c_{10} (\mathbf{1}^T \textcolor{pink}{XX}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) + c_{11} (\mathbf{1}^T \textcolor{pink}{XY}^T \textcolor{pink}{YY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) \\
& + c_{12} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{pink}{XX}^T \textcolor{pink}{XY}^T \mathbf{1}) + c_{13} (\mathbf{1}^T \textcolor{pink}{XY}^T \textcolor{pink}{YY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) \\
& + c_{14} (\mathbf{1}^T \textcolor{pink}{XX}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) + c_{15} (\mathbf{1}^T \textcolor{pink}{XY}^T \textcolor{pink}{YY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) \\
& + c_{16} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) + c_{17} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{pink}{XX}^T \textcolor{pink}{XX}^T \mathbf{1}) \\
& + c_{18} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) \\
& + c_{19} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) \\
& + c_{20} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) \\
& + c_{21} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) \\
& + c_{22} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) + c_{23} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{YY}^T \textcolor{pink}{YY}^T \mathbf{1}) \\
\hline
\frac{\partial^4}{\partial \beta^4} \frac{\partial^2}{\partial \gamma^2} \kappa_{X,Y}(0,0) = & +c_1 \mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1} + c_2 \textcolor{orange}{tr}(\textcolor{pink}{XY}^T) + c_3 \mathbf{1}^T \textcolor{pink}{XX}^T \textcolor{pink}{XY}^T \mathbf{1} + c_4 \textcolor{orange}{tr}(\textcolor{pink}{XX}^T \textcolor{pink}{XY}^T) \\
& + c_5 \mathbf{1}^T \textcolor{pink}{XY}^T \textcolor{pink}{YY}^T \mathbf{1} + c_6 \textcolor{orange}{tr}(\textcolor{pink}{XY}^T \textcolor{pink}{YY}^T) + c_7 (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) \\
& + c_8 (\textcolor{orange}{tr}(\textcolor{pink}{XY}^T)) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) + c_9 (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) \\
& + c_{10} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\textcolor{orange}{tr}(\textcolor{pink}{XY}^T)) + c_{11} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) \\
& + c_{12} \mathbf{1}^T \textcolor{pink}{XY}^T \textcolor{pink}{XY}^T \mathbf{1} + c_{13} (\textcolor{orange}{tr}(\textcolor{pink}{XY}^T)) (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) + c_{14} \mathbf{1}^T \textcolor{orange}{YX}^T \textcolor{pink}{YY}^T \mathbf{1} \\
& + c_{15} \mathbf{1}^T \textcolor{orange}{XX}^T \textcolor{pink}{YX}^T \mathbf{1} + c_{16} \mathbf{1}^T \textcolor{green}{XX}^T \textcolor{pink}{YY}^T \mathbf{1} + c_{17} (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) \\
\hline
\frac{\partial^6}{\partial \beta^6} \frac{\partial^4}{\partial \gamma^4} \kappa_{X,Y}(0,0) = & +c_1 \mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1} + c_2 \textcolor{orange}{tr}(\textcolor{pink}{XY}^T) + c_3 \mathbf{1}^T \textcolor{pink}{XX}^T \textcolor{pink}{XY}^T \mathbf{1} + c_4 \textcolor{orange}{tr}(\textcolor{pink}{XX}^T \textcolor{pink}{XY}^T) \\
& + c_5 \mathbf{1}^T \textcolor{pink}{XY}^T \textcolor{pink}{YY}^T \mathbf{1} + c_6 \textcolor{orange}{tr}(\textcolor{pink}{XY}^T \textcolor{pink}{YY}^T) + c_7 (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) \\
& + c_8 (\textcolor{orange}{tr}(\textcolor{pink}{XY}^T)) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) + c_9 (\textcolor{orange}{tr}(\textcolor{pink}{XY}^T)) (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) \\
& + c_{10} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) + c_{11} (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1}) \\
& + c_{12} \mathbf{1}^T \textcolor{pink}{XY}^T \textcolor{pink}{XY}^T \mathbf{1} + c_{13} (\textcolor{orange}{tr}(\textcolor{pink}{XY}^T)) (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) + c_{14} \mathbf{1}^T \textcolor{orange}{XX}^T \textcolor{pink}{YX}^T \mathbf{1} \\
& + c_{15} \mathbf{1}^T \textcolor{orange}{YX}^T \textcolor{pink}{YY}^T \mathbf{1} + c_{16} \textcolor{orange}{tr}(\textcolor{pink}{XY}^T \textcolor{pink}{XY}^T) + c_{17} (\textcolor{orange}{tr}(\textcolor{pink}{XY}^T)) (\textcolor{orange}{tr}(\textcolor{pink}{XY}^T)) + c_{18} \\
& + c_{19} \mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1} + c_{20} \mathbf{1}^T \textcolor{pink}{XX}^T \textcolor{pink}{XX}^T \mathbf{1} + c_{21} \mathbf{1}^T \textcolor{green}{XX}^T \textcolor{pink}{YY}^T \mathbf{1} + c_{22} \mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1} \\
& + c_{23} (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) + c_{24} (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{XX}^T \mathbf{1}) \\
& + c_{25} \textcolor{orange}{tr}(\textcolor{pink}{XX}^T \textcolor{pink}{YY}^T) + c_{26} \mathbf{1}^T \textcolor{pink}{YY}^T \textcolor{pink}{YY}^T \mathbf{1} + c_{27} (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1}) (\mathbf{1}^T \textcolor{blue}{YY}^T \mathbf{1})
\end{aligned}$$

Furthermore,

- in the expression for  $\kappa_{X,Y}(0,0)$  we have  $c_1 = 1/k^2 > 0$ ,
- in the expression for  $\frac{\partial^2}{\partial \beta^2} \frac{\partial^2}{\partial \gamma^2} \kappa_{X,Y}(0,0)$ , we have  $c_2 = 8/k^2 > 0$ ,
- in the expression for  $\frac{\partial^4}{\partial \beta^4} \kappa_{X,Y}(0,0)$ , we have  $c_{20} = 24/k^6 > 0$ ,
- in the expression for  $\frac{\partial^4}{\partial \beta^4} \frac{\partial^2}{\partial \gamma^2} \kappa_{X,Y}(0,0)$ , we have  $c_{16} = 48/k^4 > 0$ ,
- and in the expression for  $\frac{\partial^6}{\partial \beta^6} \frac{\partial^4}{\partial \gamma^4} \kappa_{X,Y}(0,0)$ , we have  $c_{25} = 17280/k^4 > 0$ .

## F.2 SIMPLIFYING TERMS

Let  $X \in \mathbb{R}^{k \times m}$  and  $Y \in \mathbb{R}^{k \times m}$  be matrices with one-hot rows (i.e., all entries are zero except for one).

For the submatrix corresponding to rows  $S$  and columns  $T$ , we use the notation  $[X]_{S \times T} \in \mathbb{R}^{S \times T}$ . If  $\mathbf{v}$  is a vector, then the subvector consisting of indices  $I$  is  $[\mathbf{v}]_I$ .

Let  $\mathcal{R} \subseteq [m]$  be a set containing the intersection of the column support of  $X$  and  $Y$ : i.e., for all  $i \in [m] \setminus \mathcal{R}$ , either  $[X]_{[k] \times i} = \mathbf{0}$  or  $[Y]_{[k] \times i} = \mathbf{0}$ . We analyze the terms in the expressions of Section F.1 below.

### F.2.1 ASSUMING $[1^T X]_{\mathcal{R}} = [1^T Y]_{\mathcal{R}}$

Suppose that  $[1^T X]_{\mathcal{R}} = [1^T Y]_{\mathcal{R}}$ . Then any of the pink terms can be written as a function of only  $X$  or only  $Y$ .

- $\mathbf{1}^T \textcolor{pink}{XY}^T \mathbf{1} = \|[1^T X]_{\mathcal{R}}\|^2$

- $1^T X X^T X Y^T 1 = 1^T X \text{diag}(1^T X) Y^T 1 = (1^T X)^{\odot 2} \cdot (1^T Y) = \|[1^T X]_{\mathcal{R}}\|_3^3$
- $1^T X Y^T Y Y^T 1 = 1^T X \text{diag}(1^T Y) Y^T 1 = (1^T X) \cdot (1^T Y)^{\odot 2} = \|[1^T X]_{\mathcal{R}}\|_3^3$
- $1^T X X^T X X^T X Y^T 1 = 1^T X \text{diag}(1^T X) \text{diag}(1^T X) Y^T 1 = \|[1^T X]_{\mathcal{R}}\|_4^4$
- $1^T X Y^T Y X^T X Y^T 1 = 1^T X \text{diag}(1^T Y) \text{diag}(1^T X) Y^T 1 = \|[1^T X]_{\mathcal{R}}\|_4^4$
- $1^T Y X^T X Y^T Y Y^T 1 = 1^T Y \text{diag}(1^T X) \text{diag}(1^T Y) Y^T 1 = \|[1^T X]_{\mathcal{R}}\|_4^4$
- $\text{trace}(X X^T X Y^T) = \text{trace}(X \text{diag}(1^T X) Y^T) = \sum_{i \in [k]} \sum_{v \in [m]} X_{iv} (1^T X)_v Y_{iv} = \sum_{i \in [k]} \sum_{v \in \mathcal{R}} X_{iv} (1^T X)_v = 1^T X \text{diag}(1^T X) 1_{\mathcal{R}} = \|[1^T X]_{\mathcal{R}}\|^2$
- $\text{trace}(X Y^T Y Y^T) = \|[1^T Y]_{\mathcal{R}}\|^2 = \|[1^T X]_{\mathcal{R}}\|^2$

### F.2.2 ASSUMING $[X]_{[k] \times \mathcal{R}} = [Y]_{[k] \times \mathcal{R}}$

Suppose that  $X_{[k] \times \mathcal{R}} = Y_{[k] \times \mathcal{R}}$  (i.e., the restriction of  $X$  and  $Y$  to the  $\mathcal{R}$  rows is equal). Then any of the **orange** terms can be written as a function of only  $X$  or only  $Y$ .

- $\text{tr}(X Y^T) = \sum_{v \in [m]} \sum_{i \in [k]} X_{iv} Y_{iv} = \sum_{v \in \mathcal{R}} \sum_{i \in [k]} X_{iv}^2 = 1^T X 1_{\mathcal{R}} = 1^T Y 1_{\mathcal{R}}$
- $1^T X Y^T X Y^T 1 = \sum_{a,b,c \in [k]} 1(x_a = y_b) 1(x_b = y_c) = 1^T X_{[k] \times \mathcal{R}} (Y_{[k] \times \mathcal{R}})^T X_{[k] \times \mathcal{R}} (Y_{[k] \times \mathcal{R}})^T 1 = 1^T X_{[k] \times \mathcal{R}} (X_{[k] \times \mathcal{R}})^T X_{[k] \times \mathcal{R}} (X_{[k] \times \mathcal{R}})^T 1$
- $1^T X X^T Y X^T 1 = \sum_{a,b,c} 1(x_a = x_b) 1(y_b = x_c) = \sum_{a,b,c} 1(x_a = x_b) 1(y_b = x_c \in \mathcal{R}) = \sum_{a,b,c} 1(x_a = x_b \in \mathcal{R}) 1(y_b = x_c \in \mathcal{R}) = \sum_{a,b,c} 1(x_a = x_b \in \mathcal{R}) 1(x_b = x_c \in \mathcal{R}) = 1^T X_{[k] \times \mathcal{R}} (X_{[k] \times \mathcal{R}})^T X_{[k] \times \mathcal{R}} (X_{[k] \times \mathcal{R}})^T 1$
- $1^T Y X^T Y Y^T 1 = 1^T X_{[k] \times \mathcal{R}} (X_{[k] \times \mathcal{R}})^T X_{[k] \times \mathcal{R}} (X_{[k] \times \mathcal{R}})^T 1$
- $\text{trace}(X Y^T X Y^T) = \sum_{a,b} 1(x_a = y_b) 1(x_b = y_a) = \sum_{a,b} 1(x_a = y_b \in \mathcal{R}) 1(x_b = y_a \in \mathcal{R}) = \sum_{a,b} 1(x_a = x_b \in \mathcal{R}) = \text{trace}((X_{[k] \times \mathcal{R}})(X_{[k] \times \mathcal{R}})^T)$

### F.2.3 ASSUMING $1^T X X^T 1 = 1^T Y Y^T 1$

Suppose that  $1^T X X^T 1 = 1^T Y Y^T 1$ . Then any of the **blue** terms can be written as a function of only  $X$  or only  $Y$ .

- $1^T X X^T 1 = 1^T Y Y^T 1$
- $1^T Y Y^T 1 = 1^T X X^T 1$

### F.2.4 ASSUMING $1^T X X^T = 1^T Y Y^T$

Suppose that  $1^T X X^T = 1^T Y Y^T$ . Then any of the **teal** terms can be written as a function of only  $X$  or only  $Y$ .

- $1^T X X^T Y Y^T 1 = \|1^T X X^T\|^2 = \|1^T Y Y^T\|^2$

## F.3 PROOF OF LEMMA F.1

We combine the above calculations to prove Lemma F.1.

*Proof.* By the technical Lemma G.1, we know that  $g_{\tau}(\beta, \gamma)$  is an analytic function for each  $\tau$ . Therefore, by the identity theorem for analytic functions (Mityagin, 2020), it suffices to show that for each  $\tau \in S_r \setminus \{\text{id}\}$  we have  $g_{\text{id}}(\beta, \gamma) \neq g_{\tau}(\beta, \gamma)$ .

Stage 1. Matching regular token degree distributions.

**Claim F.2.** *If  $g_{\text{id}}(0, 0) = g_{\tau}(0, 0)$ , then  $[1^T X_i]_{\mathcal{R}} = [1^T Y_{\tau(i)}]_{\mathcal{R}}$  for all  $i \in [r]$ .*

*Proof.* From the table in Section F.1, there is a positive constant  $c_1 > 0$  such that

$$\begin{aligned}
g_\tau(0, 0) &= c_1 \sum_{i \in [r]} 1^T X_i Y_{\tau(i)}^T 1 = c_1 \sum_{i \in [r]} [1^T X_i]_{\mathcal{R}} [Y_{\tau(i)}^T 1]_{\mathcal{R}} \\
&\stackrel{(a)}{\leq} \sum_{i \in [r]} \|[1^T X_i]_{\mathcal{R}}\| \|[1^T Y_{\tau(i)}]_{\mathcal{R}}\| \\
&\stackrel{(b)}{\leq} \sqrt{\sum_{i \in [r]} \|[1^T X_i]_{\mathcal{R}}\|^2} \sqrt{\sum_{i \in [r]} \|[1^T Y_{\tau(i)}]_{\mathcal{R}}\|^2} \\
&= \sum_{i \in [r]} \|[1^T X_i]_{\mathcal{R}}\|^2,
\end{aligned}$$

where (a) is by Cauchy-Schwarz and holds with equality if and only if  $[1^T X_i]_{\mathcal{R}} \propto [1^T Y_{\tau(i)}]_{\mathcal{R}}$  for all  $i$ . Similarly (b) is by Cauchy-Schwarz and holds with equality if and only if  $\|[1^T X_i]_{\mathcal{R}}\| = \|[1^T Y_{\tau(i)}]_{\mathcal{R}}\|$  for all  $i$ . Notice that (a) and (b) hold with equality if  $\tau = \text{id}$ , since  $[1^T X_i]_{\mathcal{R}} = [1^T Y_i]_{\mathcal{R}}$  for all  $i$ .  $\square$

### Stage 2. Matching regular token positions.

**Claim F.3.** If  $\frac{\partial^2}{\partial \beta^2} \frac{\partial^2}{\partial \gamma^2} g_\tau(0, 0) = \frac{\partial^2}{\partial \beta^2} \frac{\partial^2}{\partial \gamma^2} g_{\text{id}}(0, 0)$  and  $[1^T X_i]_{\mathcal{R}} = [1^T Y_{\tau(i)}]_{\mathcal{R}}$  for all  $i \in [r]$ , then we must have  $[X_i]_{[k] \times \mathcal{R}} = [Y_{\tau(i)}]_{[k] \times \mathcal{R}}$  for all  $i \in [r]$ .

*Proof.* For a constant  $c_2 > 0$ ,

$$\begin{aligned}
\frac{\partial^2}{\partial \beta^2} \frac{\partial^2}{\partial \gamma^2} g_\tau(0, 0) &= \sum_{i \in [r]} c_1 1^T X_i Y_{\tau(i)}^T 1 + c_2 \text{trace}(X_i Y_{\tau(i)}^T) \\
&= \left( c_1 \sum_{i \in [r]} \|[1^T X_i]_{\mathcal{R}}\|^2 \right) + \left( c_2 \sum_{i \in [r]} \text{trace}(X_i (Y^{\tau(i)})^T) \right),
\end{aligned}$$

by the calculation in Section F.2.1. The first sum does not depend on  $\tau$ , so we analyze the second sum. Here,

$$\begin{aligned}
c_2 \sum_{i \in [r]} \text{trace}(X_i Y_{\tau(i)}^T) &= c_2 \sum_{i \in [r]} \sum_{a \in [k]} [X_i Y_{\tau(i)}^T]_{aa} \\
&= c_2 \sum_{i \in [r]} \sum_{v \in \mathcal{R}} \sum_{a \in [k]} [X_i]_{av} [Y_{\tau(i)}]_{av} \\
&\stackrel{(a)}{\leq} c_2 \sqrt{\left( \sum_{i \in [r]} \sum_{v \in \mathcal{R}} \sum_{a \in [k]} ([X_i]_{av})^2 \right) \left( \sum_{i \in [r]} \sum_{v \in \mathcal{R}} \sum_{a \in [k]} ([Y_{\tau(i)}]_{av})^2 \right)} \\
&= c_2 \sum_{i \in [r]} 1^T X_i 1_{\mathcal{R}},
\end{aligned}$$

where (a) is by Cauchy-Schwarz and holds with equality if and only if  $X_{av}^{(i)} = c Y_{av}^{(\tau(i))}$  for some constant  $c$ . We must have  $c = 1$  because of the CLS token, so (a) holds with equality if and only if  $[X_i]_{[k] \times \mathcal{R}} = [Y_{\tau(i)}]_{[k] \times \mathcal{R}}$  for all  $i \in [r]$ . Specifically (a) holds with equality if  $\tau = \text{id}$ .  $\square$

### Stage 3. Matching wildcard token degree histogram norm.

**Claim F.4.** Suppose that  $[1^T X_i]_{\mathcal{R}} = [1^T Y_{\tau(i)}]_{\mathcal{R}}$ , and that  $\frac{\partial^4}{\partial \beta^4} g_\tau(0, 0) = \frac{\partial^4}{\partial \beta^4} g_{\text{id}}(0, 0)$ . Then  $1^T X_i X_i^T 1 = 1^T Y_{\tau(i)} Y_{\tau(i)}^T 1$  for all  $i \in [r]$ .

*Proof.* Use  $[1^T X_i]_{\mathcal{R}} = [1^T Y_{\tau(i)}]_{\mathcal{R}}$  and the calculations in Section F.2.1 for the pink terms. Every term of  $\frac{\partial^4}{\partial \beta^4} g_{\tau}(0, 0)$  can be written as depending only on one of  $X_i$  or  $Y_{\tau(i)}$ , with the exception of the  $c_{20}$  term. Namely, we have

$$\begin{aligned} \frac{\partial^4}{\partial \beta^4} g_{\tau}(0, 0) &= \sum_{i \in [r]} a(X_i) + b(Y_{\tau(i)}) \\ &\quad + c_{20} (1^T X_i Y_{\tau(i)}^T 1) (1^T X_i X_i^T 1) (1^T Y_{\tau(i)} Y_{\tau(i)}^T 1), \end{aligned}$$

for some functions  $a, b$ . Since  $\tau$  is a permutation, only the term with coefficient  $c_{20}$  depends on  $\tau$ . Here,  $c_{20} > 0$ . This term corresponds to

$$\begin{aligned} c_{20} \sum_{i \in [r]} (1^T X_i Y_{\tau(i)}^T 1) (1^T X_i X_i^T 1) (1^T Y_{\tau(i)} Y_{\tau(i)}^T 1) \\ = c_{20} \sum_{i \in [r]} \|[1^T X_i]_{\mathcal{R}}\| \|[1^T Y_{\tau(i)}]_{\mathcal{R}}\| (1^T X_i X_i^T 1) (1^T Y_{\tau(i)} Y_{\tau(i)}^T 1) \\ \stackrel{(a)}{\leq} \sqrt{\left( \sum_{i \in [r]} \|[1^T X_i]_{\mathcal{R}}\|^2 (1^T X_i X_i^T 1)^2 \right) \left( \sum_{i \in [r]} \|[1^T Y_{\tau(i)}]_{\mathcal{R}}\|^2 (1^T Y_{\tau(i)} Y_{\tau(i)}^T 1)^2 \right)} \\ = \sum_{i \in [r]} \|[1^T X_i]_{\mathcal{R}}\|^2 (1^T X_i X_i^T 1)^2 \end{aligned}$$

where (a) is by Cauchy-Schwarz and holds with equality if and only if  $\|[1^T X_i]_{\mathcal{R}}\|^2 1^T X_i X_i^T 1 = c \|[1^T Y_{\tau(i)}]_{\mathcal{R}}\|^2 1^T Y_{\tau(i)} Y_{\tau(i)}^T 1$  for all  $i$  and some constant  $c$ . This constant  $c = 1$  because the former is a permutation of the latter over  $i \in [r]$ . Since  $\|[1^T X_i]_{\mathcal{R}}\|^2 = \|[1^T Y_i]_{\mathcal{R}}\|^2 \geq 1$  by assumption and since we have the CLS token, we know that (a) holds with equality if and only if  $1^T X_i X_i^T 1 = 1^T Y_{\tau(i)} Y_{\tau(i)}^T 1$  for all  $i \in [r]$ . This is the case for  $\tau = \text{id}$  by construction of  $X_i$  and  $Y_i$ .  $\square$

#### Stage 4. Matching wildcard degree distributions.

**Claim F.5.** Suppose that  $[X_i]_{[k] \times \mathcal{R}} = [Y_{\tau(i)}]_{[k] \times \mathcal{R}}$  and  $1^T X_i X_i^T 1 = 1^T Y_{\tau(i)} Y_{\tau(i)}^T 1$  for all  $i \in [r]$ . Suppose also that  $\frac{\partial^4}{\partial \beta^4} \frac{\partial^2}{\partial \gamma^2} g_{\tau}(0, 0) = \frac{\partial^4}{\partial \beta^4} \frac{\partial^2}{\partial \gamma^2} g_{\text{id}}(0, 0)$ . Then  $1^T X_i X_i^T = 1^T Y_{\tau(i)} Y_{\tau(i)}^T$  for all  $i \in [r]$ .

*Proof.* Similarly to the proof of the previous claim, because of the calculations in Sections F.2.1, F.2.2 and F.2.3 for the pink, orange, and blue terms, respectively, we can write  $\frac{\partial^4}{\partial \beta^4} \frac{\partial^2}{\partial \gamma^2}$  as a sum of terms that each depends on either  $X_i$  or  $Y_{\tau(i)}$ , plus  $\sum_{i \in [r]} c_{16} 1^T X_i X_i^T Y_{\tau(i)} Y_{\tau(i)}^T 1$ . This latter sum is the only term that depends on  $\tau$ , and the constant  $c_{16}$  satisfies  $c_{16} > 0$ . Similarly to the previous claim, by Cauchy-Schwarz

$$\sum_{i \in [r]} c_{16} 1^T X_i X_i^T Y_{\tau(i)} Y_{\tau(i)}^T 1 \leq \sum_{i \in [r]} c_{16} \|1^T X_i X_i^T\| \|Y_{\tau(i)} Y_{\tau(i)}^T 1\|,$$

with equality if and only if  $1^T X_i X_i^T = 1^T Y_{\tau(i)} Y_{\tau(i)}^T$  for all  $i$ , since  $\{X_i X_i^T\}_i$  is a permutation of  $\{Y_{\tau(i)} Y_{\tau(i)}^T\}_i$ . This condition holds for  $\tau = \text{id}$ .  $\square$

#### Stage 5. Matching wildcard positions.

**Claim F.6.** Suppose that  $[X_i]_{[k] \times \mathcal{R}} = [Y_{\tau(i)}]_{[k] \times \mathcal{R}}$  and  $1^T X_i X_i^T = 1^T Y_{\tau(i)} Y_{\tau(i)}^T$  for all  $i \in [r]$ . Suppose also that  $\frac{\partial^6}{\partial \beta^6} \frac{\partial^4}{\partial \gamma^4} g_{\tau}(0, 0) = \frac{\partial^6}{\partial \beta^6} \frac{\partial^4}{\partial \gamma^4} g_{\text{id}}(0, 0)$ . Then  $X_i X_i^T = Y_{\tau(i)} Y_{\tau(i)}^T$  for all  $i \in [r]$ .

*Proof.* Write  $\frac{\partial^6}{\partial \beta^6} \frac{\partial^4}{\partial \gamma^4} g_{\tau}(0, 0)$  as a sum of terms each depending only on either  $X_i$  or  $Y_{\tau(i)}$  by using the calculations in Sections F.2.1, F.2.3, F.2.2, and F.2.4 to handle the pink, orange, blue, and teal terms, plus (for  $c_{25} > 0$ ),

$$\sum_{i \in [r]} c_{25} \text{trace}(X_i X_i^T Y_{\tau(i)} Y_{\tau(i)}^T) \leq \sum_{i \in [r]} c_{25} \|X_i X_i^T\|_F \|Y_{\tau(i)} Y_{\tau(i)}^T\|_F,$$



with equality if and only if  $X_i X_i^T = Y_{\tau(i)} Y_{\tau(i)}^T$  for all  $i \in [r]$ . This equality holds if  $\tau = \text{id}$ , concluding the claim.  $\square$

Combine the above four claims to conclude that if  $g_\tau(\beta, \gamma) \equiv g_{\text{id}}(\beta, \gamma)$ , then we have  $X_i X_i^T = Y_{\tau(i)} Y_{\tau(i)}^T$  and  $[X_i]_{[k] \times \mathcal{R}} = [Y_{\tau(i)}]_{[k] \times \mathcal{R}}$  for all  $i$ , so  $\tau = \text{id}$ .  $\square$

## G ANALYTICITY OF ATTENTION KERNEL (TECHNICAL RESULT)

We prove the analyticity of  $\kappa_{\mathbf{X}, \tilde{\mathbf{X}}}(\beta, \gamma) = K_{\text{attn}}^{\beta, \gamma}(\mathbf{X}, \tilde{\mathbf{X}})$  as function of  $\beta$  and  $\gamma$ .

**Lemma G.1** (Analyticity of  $K_{\text{attn}}$ ). *For any  $\mathbf{X}, \tilde{\mathbf{X}}$ , the function  $\kappa_{\mathbf{X}, \tilde{\mathbf{X}}}$  is analytic in  $\mathbb{R}^2$ .*

*Proof.* Note that we can write

$$\mathbf{m} := \mathbf{m}(\mathbf{X}) = \mathbf{X}\boldsymbol{\zeta} + \gamma\mathbf{p}, \quad \tilde{\mathbf{m}} := \mathbf{m}(\tilde{\mathbf{X}}) = \tilde{\mathbf{X}}\tilde{\boldsymbol{\zeta}} + \gamma\mathbf{p},$$

where  $\boldsymbol{\zeta}, \tilde{\boldsymbol{\zeta}} \sim \mathcal{N}(0, I_m)$  and  $\mathbf{p} \sim \mathcal{N}(0, I_k)$  are independent Gaussians. So we can rewrite  $\kappa_{\mathbf{X}, \tilde{\mathbf{X}}}$  as

$$\kappa_{\mathbf{X}, \tilde{\mathbf{X}}}(\beta, \gamma) = \mathbb{E}_{\boldsymbol{\zeta}, \tilde{\boldsymbol{\zeta}}, \mathbf{p}}[f(\beta, \gamma; \boldsymbol{\zeta}, \tilde{\boldsymbol{\zeta}}, \mathbf{p})],$$

where

$$f(\beta, \gamma; \boldsymbol{\zeta}, \tilde{\boldsymbol{\zeta}}, \mathbf{p}) = \mathbf{s}^T (\mathbf{X} \tilde{\mathbf{X}}^T + \gamma^2 \mathbf{I}) \tilde{\mathbf{s}}.$$

and

$$\mathbf{s} = \text{smax}(\beta \mathbf{X} \boldsymbol{\zeta} + \beta \gamma \mathbf{p})^T, \quad \tilde{\mathbf{s}} = \text{smax}(\beta \tilde{\mathbf{X}} \tilde{\boldsymbol{\zeta}} + \beta \gamma \mathbf{p}).$$

The main obstacle is to prove the technical Lemma G.9, which states that for any  $k_1, k_2$ , we have

$$\mathbb{E}_{\boldsymbol{\zeta}, \tilde{\boldsymbol{\zeta}}, \mathbf{p}} \left[ \left| \frac{\partial^{k_1}}{\partial \beta^{k_1}} \frac{\partial^{k_2}}{\partial \gamma^{k_2}} f(\beta, \gamma; \boldsymbol{\zeta}, \tilde{\boldsymbol{\zeta}}, \mathbf{p}) \right| \right] \leq C(1 + \gamma^2) k_1! k_2! (C(|\beta| + |\gamma|)^{k_1 + k_2})$$

So by smoothness of  $f$  and dominated convergence, we know that we can differentiate under the integral sign, and

$$\begin{aligned} \left| \frac{d^{k_1}}{d\beta^{k_1}} \frac{d^{k_2}}{d\gamma^{k_2}} \kappa_{\mathbf{X}, \mathbf{X}'}(\beta, \gamma) \right| &= \left| \mathbb{E}_{\boldsymbol{\zeta}, \tilde{\boldsymbol{\zeta}}, \mathbf{p}} \left[ \frac{\partial^{k_1}}{\partial \beta^{k_1}} \frac{\partial^{k_2}}{\partial \gamma^{k_2}} f(\beta, \gamma; \mathbf{X}, \tilde{\mathbf{X}}, \boldsymbol{\zeta}, \tilde{\boldsymbol{\zeta}}, \mathbf{p}) \right] \right| \\ &\leq C(1 + \gamma^2) k_1! k_2! (C(|\beta| + |\gamma|)^{k_1 + k_2}). \end{aligned}$$

Because of the bound on the derivatives and its smoothness,  $\kappa_{\mathbf{X}, \mathbf{X}'}(\beta, \gamma)$  is real-analytic.  $\square$

The proof of the technical bound in Lemma G.9 is developed in the subsections below.

### G.1 TECHNICAL LEMMAS FOR QUANTIFYING POWER SERIES CONVERGENCE

In order to show that the values of the attention kernel are real-analytic functions of in terms of  $\beta, \gamma$ , we will need to make quantitative certain facts about how real-analyticity of is preserved under compositions, products, and sums. For this, we introduce the notion of the convergence-type of a real-analytic function.

**Definition G.2** (Quantifying power series convergence in real-analytic functions). Let  $U \subseteq \mathbb{R}^m$  be an open set. We say that a real-analytic function  $f : U \rightarrow \mathbb{R}$  has  $(\tau_1, \tau_2)$ -type for functions  $\tau_1 : U \rightarrow \mathbb{R}_{>0}$  and  $\tau_2 : U \rightarrow \mathbb{R}_{>0}$  if the following holds. For any  $\boldsymbol{\zeta}_0$ , consider the power series of  $f$  around  $\boldsymbol{\zeta}_0$ ,

$$\sum_{\boldsymbol{\mu}} a_{\boldsymbol{\zeta}_0, \boldsymbol{\mu}} (\boldsymbol{\zeta} - \boldsymbol{\zeta}_0)^{\boldsymbol{\mu}}.$$

Then for any  $\boldsymbol{\zeta}$  such that  $\|\boldsymbol{\zeta} - \boldsymbol{\zeta}_0\|_{\infty} \leq \tau_1(\boldsymbol{\zeta}_0)$  this power series converges absolutely.

$$\sum_{\boldsymbol{\mu} \text{ s.t. } |\boldsymbol{\mu}| \geq 1} |a_{\boldsymbol{\zeta}_0, \boldsymbol{\mu}}| \|\boldsymbol{\zeta} - \boldsymbol{\zeta}_0\|^{\boldsymbol{\mu}} \leq \tau_2(\boldsymbol{\zeta}_0).$$

We provide rules for how convergence type is affected by compositions, products, and sums.

**Lemma G.3** (Composition rule for type; quantitative version of Proposition 2.2.8 of Krantz & Parks (2002)). *Let  $U \subseteq \mathbb{R}^m$  and let  $V \subseteq \mathbb{R}$  be open. Let  $f_1, \dots, f_n : U \rightarrow V$  be real-analytic with  $(\tau_1, \tau_2)$ -type, and let  $g : V^n \rightarrow \mathbb{R}$  be real-analytic with  $(\sigma_1, \sigma_2)$ -type. Then the composition  $h = g \circ (f_1, \dots, f_n)$  is real-analytic with  $(\min(\tau_1, (\sigma_1 \circ f) \cdot \frac{\tau_1}{\tau_2}), \sigma_2 \circ f)$ -type.*

*Proof.* Fix some  $\zeta_0$  and let  $\mathbf{y}_0 = [f_1(\zeta_0), \dots, f_n(\zeta_0)]$ , and let  $a_{\zeta_0, \mu}^{(i)}$  be the coefficients of the power series expansion for  $f_i$  around  $\zeta_0$ . Define  $\rho = \min(1, \sigma_1(\mathbf{y}_0)/\tau_2(\zeta_0))$ . Then, for any  $\zeta$  such that  $\|\zeta - \zeta_0\|_\infty \leq \rho\tau_1(\zeta_0)$  and  $i \in [n]$  we have

$$\sum_{\mu \text{ s.t. } |\mu| \geq 1} |a_{\zeta_0, \mu}^{(i)}| |\zeta - \zeta_0|^\mu \leq \sum_{\mu \text{ s.t. } |\mu| \geq 1} |a_{\zeta_0, \mu}^{(i)}| \rho^{|\mu|} \tau_1(\zeta_0)^{|\mu|} \leq \rho\tau_2(\zeta_0) \leq \sigma_1(\mathbf{y}_0).$$

So, letting  $\sum_{\nu}^\infty b_{\mathbf{y}_0, \nu} (\mathbf{y} - \mathbf{y}_0)^\nu$  be the series expansion of  $g$  around  $\mathbf{y}_0$ , we have the following absolute convergence

$$\sum_{\nu \text{ s.t. } |\nu| \geq 1}^\infty b_{\mathbf{y}_0, \nu} \prod_{i=1}^n \left| \sum_{\mu \text{ s.t. } |\mu| \geq 1} |a_{\zeta_0, \mu}^{(i)}| |\zeta - \zeta_0|^\mu \right|^{\nu_i} \leq \sigma_2(\mathbf{y}_0).$$

So we may rearrange the terms of

$$\sum_{\nu}^\infty b_{\mathbf{y}_0, \nu} \prod_{i=1}^n \left( \sum_{\mu \text{ s.t. } |\mu| \geq 1} a_{\zeta_0, \mu}^{(i)} (\zeta - \zeta_0)^\mu \right)^{\nu_i}.$$

as we please, and we get an absolutely convergent series for  $g \circ f$  around  $\zeta_0$ .  $\square$

**Lemma G.4** (Sum and product rules for type). *Let  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  and  $g : \mathbb{R}^m \rightarrow \mathbb{R}$  be real-analytic functions of  $(\tau_1, \tau_2)$ -type and  $(\sigma_1, \sigma_2)$ -type respectively. Then  $h = f + g$  is real-analytic of  $(\min(\tau_1, \sigma_1), \tau_2 + \sigma_2)$ -type, and  $h = fg$  is real-analytic of  $(\min(\tau_1, \sigma_1), \tau_2\sigma_2 + \tau_2|g| + |f|\sigma_2)$ -type*

*Proof.* Both of these are straightforward from the definition.  $\square$

**Lemma G.5** (Derivative bound based on type). *Let  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  be real-analytic with  $(\tau_1, \tau_2)$ -type. Then, for any multi-index  $\mu$ ,*

$$\left| \frac{\partial^{|\mu|}}{\partial \zeta^\mu} f(\zeta_0) \right| \leq \frac{\tau_2(\zeta_0)}{\tau_1(\zeta_0)^{|\mu|}} \mu!$$

*Proof.* Let  $a_{\zeta_0, \mu}$  be the coefficients of the power series of  $f$  at  $\zeta_0$ . Since  $f$  is of  $(\tau_1, \tau_2)$ -type, we have

$$\sum_{\mu \text{ s.t. } |\mu| \geq 1} |a_{\zeta_0, \mu}| \tau_1(\zeta_0)^{|\mu|} \leq \tau_2(\zeta_0).$$

Since all terms in the sum are nonnegative, for all  $\mu$  with  $|\mu| \geq 1$ ,

$$|a_{\zeta_0, \mu}| \leq \tau_2(\zeta_0) \cdot (1/\tau_1(\zeta_0))^{|\mu|}.$$

The lemma follows by Remark 2.2.4 of Krantz & Parks (2002), which states  $\left| \frac{\partial^{|\mu|}}{\partial \zeta^\mu} f(\zeta_0) \right| = |a_{\zeta_0, \mu}| \mu!$ .  $\square$

## G.2 APPLICATION OF TECHNICAL LEMMAS TO ATTENTION KERNEL

We now use the above general technical lemmas to specifically prove that the attention kernel is analytic in terms of  $\beta$  and  $\gamma$ .

**Lemma G.6.** *For any  $j \in [m]$ , the function  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  given by  $f(\zeta) = \text{smax}(\zeta)_j$  is real-analytic of  $(1/(2e^2), 1)$ -type*

*Proof.* Write  $f = g \circ h$  for  $g : \mathbb{R}_{>0} \rightarrow \mathbb{R}$  and  $h : \mathbb{R}^k \rightarrow \mathbb{R}_{>0}$  given by  $g(y) = 1/y$ , and  $h(\zeta) = \sum_{i=1}^m e^{\zeta_i - \zeta_j}$ .

The power expansion of  $g(y)$  around  $y_0 \in \mathbb{R}_{>0}$ , is given by

$$g(y) = \sum_{k=0}^{\infty} \frac{(-1)^{k+1}}{y_0^{k+1}} (y - y_0)^k,$$

so one can see that  $g$  is of  $(\rho_1, \rho_2)$ -type for  $\rho_1(y_0) = y_0/2$  and  $\rho_2(y_0) = 1/y_0$ . Finally, write the series expansion for  $h(\zeta)$  around  $\zeta_0$

$$h(\zeta) = 1 + e^{-\zeta_j} \sum_{i \in [m] \setminus \{j\}} e^{\zeta_i} = 1 + \sum_{i \in [m] \setminus \{j\}} \left( \sum_{l=0}^{\infty} e^{-\zeta_{0,j}} \frac{(\zeta_{0,j} - \zeta_j)^l}{l!} \right) \left( \sum_{k=0}^{\infty} e^{\zeta_{0,i}} \frac{(\zeta_i - \zeta_{0,i})^k}{k!} \right)$$

Note that this expansion converges absolutely for all  $\zeta$ , as the absolute series is

$$\begin{aligned} & 1 + \sum_{i \in [m] \setminus \{j\}} \left( \sum_{l=0}^{\infty} e^{-\zeta_{0,j}} \frac{|\zeta_{0,j} - \zeta_j|^l}{l!} \right) \left( \sum_{k=0}^{\infty} e^{\zeta_{0,i}} \frac{|\zeta_i - \zeta_{0,i}|^k}{k!} \right) \\ &= 1 + \sum_{i \in [m] \setminus \{j\}} e^{-\zeta_{0,j} + \zeta_{0,i} + |\zeta_i - \zeta_{0,i}| + |\zeta_j - \zeta_{0,j}|} \\ &\leq e^{2\|\zeta - \zeta_0\|_{\infty}} h(\zeta). \end{aligned}$$

Specifically,  $h$  is of  $(1, e^2 h)$ -type. So by the composition rule of Lemma G.3, it must be that  $f$  is real-analytic of  $(\tau_1, \tau_2)$ -type for  $\tau_1 = \min(1, (\rho_1 \circ h) \cdot \frac{1}{e^2 h}) = 1/(2e^2)$  and  $\tau_2 = \rho_2 \circ h = 1/h \leq 1$ .  $\square$

**Lemma G.7.** *For any  $j \in [m]$  and  $\mathbf{X}, \zeta, \mathbf{p}$ , the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by  $f(\beta, \gamma) = \text{smax}(\beta \mathbf{X} \zeta + \beta \gamma \mathbf{p})_j$  is real-analytic of  $(\min(1, 1/(2e^2 \|\mathbf{X} \zeta\|_{\infty} + 2e^2(|\beta| + |\gamma|)\|\mathbf{p}\|_{\infty})), 1)$ -type.*

*Proof.* Write  $f = g \circ h$  for  $g : \mathbb{R}^m \rightarrow \mathbb{R}$  and  $h : \mathbb{R}^2 \rightarrow \mathbb{R}^m$  given by  $g(\mathbf{v}) = \text{smax}(\mathbf{v})_j$  and  $h(\beta, \gamma) = \beta \mathbf{X} \zeta + \beta \gamma \mathbf{p}$ . We know from Lemma G.6 that  $g$  is real-analytic of  $(1/(2e^2), 1)$ -type. And it is easy to see that  $h$  is real-analytic of  $(1, \|\mathbf{X} \zeta\|_{\infty} + (|\beta| + |\gamma|)\|\mathbf{p}\|_{\infty})$ -type. Apply the composition rule of Lemma G.3 to conclude.  $\square$

**Lemma G.8.** *For any  $\mathbf{X}, \tilde{\mathbf{X}}, \zeta, \tilde{\zeta}, \mathbf{p}$ , the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by  $f(\beta, \gamma) = \text{smax}(\beta \mathbf{X} \zeta + \beta \gamma \mathbf{p})^T (\mathbf{X} \tilde{\mathbf{X}}^T + \gamma^2 \mathbf{I}) \text{smax}(\beta \tilde{\mathbf{X}} \tilde{\zeta} + \beta \gamma \mathbf{p})$  is real-analytic and of type*

$$\left( \min\left(1, \frac{1}{2e^2 \|\mathbf{X} \zeta\|_{\infty} + (|\beta| + |\gamma|)\|\mathbf{p}\|_{\infty}}, \frac{1}{2e^2 \|\tilde{\mathbf{X}} \tilde{\zeta}\|_{\infty} + (|\beta| + |\gamma|)\|\mathbf{p}\|_{\infty}}\right), C(1 + \gamma^2) \right),$$

where  $C$  is a constant depending on the context length  $k$ .

*Proof.* Each entry of  $(\mathbf{X} \tilde{\mathbf{X}}^T + \gamma^2 \mathbf{I})$  is real-analytic in  $\gamma$  and of  $(1, \gamma)$ -type. So by combining with Lemma G.7 the product rule and sum rule (Lemma G.4), and the fact that each entry of the  $\text{smax}$  is at most one.  $\square$

As a consequence, we can bound the derivatives of  $f(\beta, \gamma; \mathbf{X}, \tilde{\mathbf{X}}, \zeta, \tilde{\zeta}, \mathbf{p}) = \text{smax}(\beta \mathbf{X} \zeta + \beta \gamma \mathbf{p})^T (\mathbf{X} \tilde{\mathbf{X}}^T + \gamma^2 \mathbf{I}) \text{smax}(\beta \tilde{\mathbf{X}} \tilde{\zeta} + \beta \gamma \mathbf{p})$ , which was what we needed to prove Lemma G.1.

**Lemma G.9.** For any  $k_1, k_2 \geq 0$ ,

$$\begin{aligned} & \left| \frac{\partial^{k_1}}{\partial \beta^{k_1}} \frac{\partial^{k_2}}{\partial \gamma^{k_2}} f(\beta, \gamma; \mathbf{X}, \tilde{\mathbf{X}}, \boldsymbol{\zeta}, \tilde{\boldsymbol{\zeta}}, \mathbf{p}) \right| \\ & \leq C(1 + \gamma^2) \max(1, ((2e^2)(\|\mathbf{X}\boldsymbol{\zeta}\|_\infty + \|\tilde{\mathbf{X}}\tilde{\boldsymbol{\zeta}}\|_\infty + (|\beta| + |\gamma|)\|\mathbf{p}\|_\infty))^{k_1+k_2}) k_1! k_2!. \end{aligned}$$

*Proof.* Direct consequence of Lemma G.5 and Lemma G.8.  $\square$

## H DERIVATION OF TRANSFORMER KERNEL

We informally derive the transformer random features kernel in the infinite-width limit.

### H.1 TRANSFORMER ARCHITECTURE

We consider a depth-1 transformer architecture (without skip connections or layernorm, for simplicity). This architecture has  $H$  heads, each with parameters  $\mathbf{W}_{K,h}, \mathbf{W}_{Q,h}, \mathbf{W}_{V,h}, \mathbf{W}_{O,h} \in \mathbb{R}^{d_{head} \times d_{emb}}$ , and embedding layer  $\mathbf{W}_E \in \mathbb{R}^{m \times d_{emb}}$ , positional embeddings  $\mathbf{P} \in \mathbb{R}^{k \times d_{emb}}$ , an MLP layer with parameters  $\mathbf{W}_A, \mathbf{W}_B \in \mathbb{R}^{d_{mlp} \times d_{emb}}$ , and a final unembedding layer with weights  $\mathbf{w}_U \in \mathbb{R}^{d_{emb}}$ . The network takes in  $\mathbf{X} \in \mathbb{R}^{k \times m}$  and outputs

$$f_{\text{trans}}(\mathbf{X}; \boldsymbol{\theta}) = \mathbf{w}_U^T \mathbf{z}_2 \quad (\text{Unembedding})$$

where

$$\mathbf{z}_2 = \frac{1}{\sqrt{d_{mlp}}} \mathbf{W}_B^T \sigma\left(\frac{1}{\sqrt{d_{emb}}} \mathbf{W}_A \mathbf{z}_1\right) \in \mathbb{R}^{d_{emb}} \quad (\text{MLP layer})$$

$$\mathbf{z}_1 = \frac{1}{\sqrt{H}} \sum_{h \in [H]} \mathbf{A}_h^T \mathbf{e}_k \in \mathbb{R}^{d_{emb}} \quad (\text{Attention layer output at CLS token})$$

$$\mathbf{A}_h = \text{smax}\left(\frac{\beta \mathbf{Z}_0 \mathbf{W}_{K,h}^T \mathbf{W}_{Q,h} \mathbf{Z}_0^T}{d_{emb} \sqrt{d_{head}}}\right) \mathbf{Z}_0 \frac{\mathbf{W}_{V,h}^T \mathbf{W}_{O,h}}{\sqrt{d_{head} d_{emb}}} \in \mathbb{R}^{k \times d_{emb}} \quad (\text{Attention heads})$$

$$\mathbf{Z}_0 = \mathbf{X} \mathbf{W}_E + \gamma \mathbf{P} \in \mathbb{R}^{k \times d_{emb}}. \quad (\text{Embedding layer})$$

Here  $\beta, \gamma \geq 0$  are two hyperparameters that control the inverse temperature of the softmax and the strength of the positional embeddings, respectively. Note that only the output of the attention layer at the final  $k$ th position CLS token is used, since this is a depth-1 network. Also, in the above definition the weights are rescaled compared to Section 2, but this is not important since what matters is the

### H.2 RANDOM FEATURES KERNEL

We choose that initialization so that each of the entries of the intermediate representations  $\mathbf{Z}_0, \mathbf{z}_1, \mathbf{z}_2$  is of order  $\Theta(1)$ . In order to accomplish this, we initialize  $\mathbf{W}_E, \mathbf{P}, \mathbf{W}_{K,h}, \mathbf{W}_{Q,h}, \mathbf{W}_{V,h}, \mathbf{W}_{O,h}, \mathbf{W}_A, \mathbf{W}_B$  with i.i.d.  $N(0, 1)$  entries.

We also initialize  $\mathbf{w}_U = 0$ , and only train  $\mathbf{w}_U$  while maintaining the rest of parameters at initialization. The random features kernel corresponding to training  $\mathbf{w}_U$  is

$$\hat{K}_{\text{trans}}(\mathbf{X}, \mathbf{Y}) = \mathbf{z}_2(\mathbf{X})^T \mathbf{z}_2(\mathbf{Y}) / d_{emb},$$

where we view  $\mathbf{z}_2$  as a function of the input (either  $\mathbf{X}$  or  $\mathbf{Y}$ ), and depending on the randomly-initialized parameters of the network.

In the limit of infinitely-many heads  $H$ , infinite embedding dimension  $d_{emb}$  and MLP dimension  $d_{mlp}$  and head dimension  $d_{head}$ , the kernel  $\hat{K}_{\text{trans}}$  tends to a deterministic limit  $K_{\text{trans}}$ , which can be recursively computed (see, e.g., Jacot et al. (2018)). Assuming that the final token of both  $\mathbf{X}$  and  $\mathbf{Y}$

is the same token (i.e., a CLS token), the deterministic limiting kernel  $K_{\text{trans}}$  is given by:

$$K_{\text{trans}}(\mathbf{X}, \mathbf{Y}) = \mathbb{E}_{u,v}[\sigma(u)\sigma(v)] \text{ for } u, v \sim N(\mathbf{0}, \begin{bmatrix} K_{\text{attn}}(\mathbf{X}, \mathbf{X}) & K_{\text{attn}}(\mathbf{X}, \mathbf{Y}) \\ K_{\text{attn}}(\mathbf{Y}, \mathbf{X}) & K_{\text{attn}}(\mathbf{Y}, \mathbf{Y}) \end{bmatrix}) \quad (22)$$

where  $K_1(\mathbf{X}, \mathbf{Y}) = \mathbb{E}_{\mathbf{m}(\mathbf{X}), \mathbf{m}(\mathbf{Y})}[\text{smax}(\beta \mathbf{m}(\mathbf{X}))^T (\mathbf{X} \mathbf{Y}^T + \gamma^2 \mathbf{I}) \text{smax}(\beta \mathbf{m}(\mathbf{Y}))]$

$$\mathbf{m}(\mathbf{X}), \mathbf{m}(\mathbf{Y}) \sim N(\mathbf{0}, (1 + \gamma^2) \begin{bmatrix} \mathbf{X} \mathbf{X}^T + \gamma^2 \mathbf{I} & \mathbf{X} \mathbf{Y}^T + \gamma^2 \mathbf{I} \\ \mathbf{Y} \mathbf{X}^T + \gamma^2 \mathbf{I} & \mathbf{Y} \mathbf{Y}^T + \gamma^2 \mathbf{I} \end{bmatrix}).$$

Notice that the covariance matrix in the above definition of the distribution of  $\mathbf{m}(\mathbf{X}), \mathbf{m}(\mathbf{Y})$  is slightly rescaled compared to that in the main text in Section 4.1, but this is inessential, since we can simply reparametrize  $\beta$  as  $\beta \mapsto \beta / \sqrt{1 + \gamma^2}$  to recover the expression in the main text.

### H.3 INFORMAL DERIVATION

We provide an informal derivation of (22) below. Informally, by law of large numbers we have the following almost sure convergence

$$\begin{aligned} \hat{K}_{\text{trans}}(\mathbf{X}, \mathbf{Y}) &= \frac{\mathbf{z}_2(\mathbf{X})^T \mathbf{z}_2(\mathbf{Y})}{d_{\text{emb}}} = \frac{\sigma(\frac{1}{\sqrt{d_{\text{emb}}}} \mathbf{W}_A \mathbf{z}_1(\mathbf{X}))^T \mathbf{W}_B \mathbf{W}_B^T \sigma(\frac{1}{\sqrt{d_{\text{emb}}}} \mathbf{W}_A \mathbf{z}_1(\mathbf{Y}))}{d_{\text{emb}} d_{\text{mlp}}} \\ &\xrightarrow{d_{\text{emb}} \rightarrow \infty} \frac{\sigma(\frac{1}{\sqrt{d_{\text{emb}}}} \mathbf{W}_A \mathbf{z}_1(\mathbf{X}))^T \sigma(\frac{1}{\sqrt{d_{\text{emb}}}} \mathbf{W}_A \mathbf{z}_1(\mathbf{Y}))}{d_{\text{mlp}}} \\ &\xrightarrow{d_{\text{mlp}} \rightarrow \infty} \mathbb{E}_{u,v}[\sigma(u)\sigma(v)] \text{ for } u, v \sim N(\mathbf{0}, \begin{bmatrix} K_{\text{attn}}(\mathbf{X}, \mathbf{X}) & K_{\text{attn}}(\mathbf{X}, \mathbf{Y}) \\ K_{\text{attn}}(\mathbf{Y}, \mathbf{X}) & K_{\text{attn}}(\mathbf{Y}, \mathbf{Y}) \end{bmatrix}) \\ &:= K_{\text{trans}}(\mathbf{X}, \mathbf{Y}), \end{aligned}$$

where  $K_{\text{attn}}$  is the kernel corresponding to the attention layer in the infinite-width limit, defined as:

$$\begin{aligned} \hat{K}_{\text{attn}}(\mathbf{X}, \mathbf{Y}) &:= \frac{\mathbf{z}_1^T(\mathbf{X}) \mathbf{z}_1^T(\mathbf{Y})}{d_{\text{emb}}} = \frac{\sum_{h,h' \in [H]} \mathbf{e}_k^T \mathbf{A}_h(\mathbf{X}) \mathbf{A}_{h'}(\mathbf{Y})^T \mathbf{e}_k}{H d_{\text{emb}}} \\ &= \frac{1}{H d_{\text{head}} d_{\text{emb}}^2} \sum_{h,h' \in [H]} \mathbf{e}_k^T \text{smax}(\frac{\beta \mathbf{Z}_0(\mathbf{X}) \mathbf{W}_{K,h}^T \mathbf{W}_{Q,h} \mathbf{Z}_0(\mathbf{X})^T}{d_{\text{emb}} \sqrt{d_{\text{head}}}}) \mathbf{Z}_0(\mathbf{X}) \mathbf{W}_{V,h}^T \mathbf{W}_{O,h} \\ &\quad \cdot \mathbf{W}_{O,h'}^T \mathbf{W}_{V,h'} \mathbf{Z}_0(\mathbf{Y})^T \text{smax}(\frac{\beta \mathbf{Z}_0(\mathbf{Y}) \mathbf{W}_{K,h'}^T \mathbf{W}_{Q,h'} \mathbf{Z}_0(\mathbf{Y})^T}{d_{\text{emb}} \sqrt{d_{\text{head}}}})^T \mathbf{e}_k \\ &\xrightarrow{d_{\text{head}} \rightarrow \infty, d_{\text{emb}} \rightarrow \infty} \frac{1}{H} \sum_{h \in [H]} \mathbf{e}_k^T \text{smax}(\frac{\beta \mathbf{Z}_0(\mathbf{X}) \mathbf{W}_{K,h}^T \mathbf{W}_{Q,h} \mathbf{Z}_0(\mathbf{X})^T}{d_{\text{emb}} \sqrt{d_{\text{head}}}}) (\mathbf{X} \mathbf{Y}^T + \gamma^2 \mathbf{I}) \\ &\quad \cdot \text{smax}(\frac{\beta \mathbf{Z}_0(\mathbf{Y}) \mathbf{W}_{K,h}^T \mathbf{W}_{Q,h} \mathbf{Z}_0(\mathbf{Y})^T}{d_{\text{emb}} \sqrt{d_{\text{head}}}})^T \mathbf{e}_k \\ &\xrightarrow{H \rightarrow \infty} \mathbb{E}[\mathbf{e}_k^T \text{smax}(\frac{\beta \mathbf{Z}_0(\mathbf{X}) \mathbf{W}_{K,h}^T \mathbf{W}_{Q,h} \mathbf{Z}_0(\mathbf{X})^T}{d_{\text{emb}} \sqrt{d_{\text{head}}}}) (\mathbf{X} \mathbf{Y}^T + \gamma^2 \mathbf{I}) \\ &\quad \cdot \text{smax}(\frac{\beta \mathbf{Z}_0(\mathbf{Y}) \mathbf{W}_{K,h}^T \mathbf{W}_{Q,h} \mathbf{Z}_0(\mathbf{Y})^T}{d_{\text{emb}} \sqrt{d_{\text{head}}}})^T \mathbf{e}_k] \\ &= \mathbb{E}[\text{smax}(\frac{\beta \mathbf{e}_k^T \mathbf{Z}_0(\mathbf{X}) \mathbf{W}_{K,h}^T \mathbf{W}_{Q,h} \mathbf{Z}_0(\mathbf{X})^T}{d_{\text{emb}} \sqrt{d_{\text{head}}}}) (\mathbf{X} \mathbf{Y}^T + \gamma^2 \mathbf{I}) \\ &\quad \cdot \text{smax}(\frac{\beta \mathbf{e}_k^T \mathbf{Z}_0(\mathbf{Y}) \mathbf{W}_{K,h}^T \mathbf{W}_{Q,h} \mathbf{Z}_0(\mathbf{Y})^T}{d_{\text{emb}} \sqrt{d_{\text{head}}}})^T] \\ &\xrightarrow{d_{\text{emb}} \rightarrow \infty, d_{\text{head}} \rightarrow \infty} \mathbb{E}_{\mathbf{m}(\mathbf{X}), \mathbf{m}(\mathbf{Y})}[\text{smax}(\beta \mathbf{m}(\mathbf{X}))^T (\mathbf{X} \mathbf{Y}^T + \gamma^2 \mathbf{I}) \text{smax}(\beta \mathbf{m}(\mathbf{Y}))] \\ &:= K_{\text{attn}}(\mathbf{X}, \mathbf{Y}), \end{aligned}$$

where

$$\mathbf{m}(\mathbf{X}), \mathbf{m}(\mathbf{Y}) \sim N(\mathbf{0}, (1 + \gamma^2) \begin{bmatrix} \mathbf{X}\mathbf{X}^T + \gamma^2 \mathbf{I} & \mathbf{X}\mathbf{Y}^T + \gamma^2 \mathbf{I} \\ \mathbf{Y}\mathbf{X}^T + \gamma^2 \mathbf{I} & \mathbf{Y}\mathbf{Y}^T + \gamma^2 \mathbf{I} \end{bmatrix}),$$

because due to the randomness in  $\mathbf{W}_{K,h}$  and  $\mathbf{W}_{Q,h}$  we have that

$$\frac{\mathbf{Z}_0(\mathbf{X})\mathbf{W}_{Q,h}^T \mathbf{W}_{K,h} \mathbf{Z}_0(\mathbf{X})^T \mathbf{e}_k}{d_{emb}\sqrt{d_{head}}}$$

and

$$\frac{\mathbf{Z}_0(\mathbf{Y})\mathbf{W}_{Q,h}^T \mathbf{W}_{K,h} \mathbf{Z}_0(\mathbf{Y})^T \mathbf{e}_k}{d_{emb}\sqrt{d_{head}}}$$

are jointly Gaussian with covariance:

$$\Sigma(\mathbf{X}, \mathbf{Y}) = \mathbb{E}_{\mathbf{W}_{K,h}, \mathbf{W}_{Q,h}, \mathbf{W}_E, \mathbf{P}} \left[ \frac{\mathbf{Z}_0(\mathbf{X})\mathbf{W}_{Q,h}^T \mathbf{W}_{K,h} \mathbf{Z}_0(\mathbf{X})^T \mathbf{e}_k}{d_{emb}\sqrt{d_{head}}} \frac{\mathbf{e}_k^T \mathbf{Z}_0(\mathbf{Y})\mathbf{W}_{K,h}^T \mathbf{W}_{Q,h} \mathbf{Z}_0(\mathbf{Y})^T}{d_{emb}\sqrt{d_{head}}} \right],$$

Since this is an expectation over products of jointly Gaussian variables, for any  $i, j \in [k]$  we can calculate:

$$\begin{aligned} \Sigma_{i,j}(\mathbf{X}, \mathbf{Y}) &= \mathbb{E}_{\mathbf{W}_E, \mathbf{P}} \left[ \frac{1}{d_{emb}^2} \sum_{r,s \in [d_{emb}]} [\mathbf{Z}_0(\mathbf{X})]_{ir} [\mathbf{Z}_0(\mathbf{Y})]_{js} \text{trace}(\mathbf{Z}_0(\mathbf{X})^T \mathbf{e}_k \mathbf{e}_k^T \mathbf{Z}_0(\mathbf{Y})) \right] \\ &= \mathbb{E}_{\mathbf{W}_E, \mathbf{P}} \left[ \frac{1}{d_{emb}^2} \sum_{r,s,t \in [d_{emb}]} [\mathbf{Z}_0(\mathbf{X})]_{ir} [\mathbf{Z}_0(\mathbf{Y})]_{js} [\mathbf{Z}_0(\mathbf{X})]_{kt} [\mathbf{Z}_0(\mathbf{Y})]_{kt} \right] \\ &= \mathbb{E}_{\mathbf{W}_E, \mathbf{P}} \left[ \frac{1}{d_{emb}^2} \sum_{r,s,t \in [d_{emb}]} [\mathbf{X}\mathbf{W}_E + \gamma\mathbf{P}]_{ir} [\mathbf{Y}\mathbf{W}_E + \gamma\mathbf{P}]_{js} [\mathbf{X}\mathbf{W}_E + \gamma\mathbf{P}]_{kt} [\mathbf{Y}\mathbf{W}_E + \gamma\mathbf{P}]_{kt} \right] \\ &\stackrel{(a)}{=} \frac{1}{d_{emb}^2} \sum_{r,s \in [d_{emb}]} \mathbb{E}_{\mathbf{W}_E, \mathbf{P}} [[\mathbf{X}\mathbf{W}_E + \gamma\mathbf{P}]_{ir} [\mathbf{Y}\mathbf{W}_E + \gamma\mathbf{P}]_{js}] \\ &\quad \cdot \sum_{t \in [d_{emb}]} \mathbb{E}_{\mathbf{W}_E, \mathbf{P}} [[\mathbf{X}\mathbf{W}_E + \gamma\mathbf{P}]_{kt} [\mathbf{Y}\mathbf{W}_E + \gamma\mathbf{P}]_{kt}] + O(1/d_{emb}) \\ &= \frac{1}{d_{emb}} \sum_{r,s \in [d_{emb}]} \mathbb{E}_{\mathbf{W}_E, \mathbf{P}} [[\mathbf{X}\mathbf{W}_E + \gamma\mathbf{P}]_{ir} [\mathbf{Y}\mathbf{W}_E + \gamma\mathbf{P}]_{js}] \cdot (1 + \gamma^2) + O(1/d_{emb}) \\ &\stackrel{(a)}{=} \frac{1}{d_{emb}} \sum_{r \in [d_{emb}]} \mathbb{E}_{\mathbf{W}_E, \mathbf{P}} [[\mathbf{X}\mathbf{W}_E + \gamma\mathbf{P}]_{ir} [\mathbf{Y}\mathbf{W}_E + \gamma\mathbf{P}]_{jr}] \cdot (1 + \gamma^2) + O(1/d_{emb}) \\ &= [\mathbf{X}\mathbf{Y}^T]_{ij} + \gamma^2 \delta_{ij} \cdot (1 + \gamma^2) + O(1/d_{emb}), \end{aligned}$$

where in (a) we use that  $[\mathbf{X}\mathbf{W}_E + \gamma\mathbf{P}]_{ab}$  and  $[\mathbf{Y}\mathbf{W}_E + \gamma\mathbf{P}]_{ab}$  are independent of  $[\mathbf{X}\mathbf{W}_E + \gamma\mathbf{P}]_{cd}$  and  $[\mathbf{Y}\mathbf{W}_E + \gamma\mathbf{P}]_{cd}$  unless  $b = d$ . So

$$\Sigma(\mathbf{X}, \mathbf{Y}) \xrightarrow{d_{emb} \rightarrow \infty} (1 + \gamma^2) \cdot (\mathbf{X}\mathbf{Y}^T + \gamma^2 \mathbf{I}).$$

## I MLPs FAIL TO GENERALIZE ON UNSEEN SYMBOLS

A natural question is whether classical architectures such as the MLP architecture (a.k.a., fully-connected network) would exhibit the same emergent reasoning properties when trained with enough data. In this section, we prove a negative result: an SGD-trained or Adam-trained MLP will not reach good test performance on the template task. This is in sharp contrast to the positive result for transformers proved in the previous section.

**MLP architecture** The input to the MLP is a concatenation of the token one-hot encodings. The MLP alternates linear transformations and nonlinear elementwise activations. Formally, the MLP

has weights  $\theta = \{W_1, \dots, W_L, w\}$  and outputs

$$\begin{aligned} f_{\text{MLP}}(\mathbf{x}; \theta) &= \mathbf{w}^T \mathbf{z}_L(\mathbf{x}; \theta) \in \mathbb{R} \quad \text{where} \\ \mathbf{z}_\ell(\mathbf{x}; \theta) &= \phi(W_\ell \mathbf{z}_{\ell-1}(\mathbf{x}; \theta)) \in \mathbb{R}^d \quad \text{for } \ell \geq 1 \\ \mathbf{z}_0(\mathbf{x}; \theta) &= \mathbf{z}_0(\mathbf{x}) = [e_{x_1}, \dots, e_{x_k}] \in \mathbb{R}^{km}. \end{aligned} \quad (23)$$

We consider training the MLP with SGD.

**Definition I.1** (One-pass SGD training). The learned weights  $\theta^t$  after  $t$  steps of SGD training are the random weights given by initializing  $\theta^0$  so that each of  $W_1^0, \dots, W_L^0, w^0$  have i.i.d. Gaussian entries, and then updating with  $\theta^t = \theta^{t-1} - \eta_t \nabla_{\theta} (f_{\text{MLP}}(\mathbf{x}^t; \theta) - y^t)^2 \mid_{\theta=\theta^{t-1}}$  for  $(\mathbf{x}^t, y^t) \sim \mathcal{D}$  and some step size  $\eta_t > 0$ .

We show that SGD-trained MLPs fail at the template task since they do not generalize well in the case when the templates consist only of wildcard tokens. In words, if the template labels  $f_*$  are a non-constant function, the MLP will not reach arbitrarily low error no matter how many training steps are taken. Let  $\mathcal{X}_{\text{uns}} \subset \mathcal{X}$  be the subset of tokens not seen in the train data. We assume that  $|\mathcal{X}_{\text{uns}}| \geq k$ , which guarantees that for any template there is at least one string matching it where all the wildcards are substituted by tokens in  $\mathcal{X}_{\text{uns}}$ . Under this condition:

**Theorem I.2** (Failure of MLPs at generalizing on unseen symbols). *Suppose that the label function  $f_*$  is non-constant, and that all templates in the support of  $\mu_{\text{tplt}}$  consist only of wildcards:  $\mathbf{z} \in \mathcal{W}^k$  for all  $\mathbf{z} \in \text{supp}(\mu_{\text{tplt}})$ . Then, for any SGD step  $t$  there is a string  $\mathbf{x} \in (\mathcal{X}_{\text{uns}})^k$  that matches a template  $\mathbf{z} \in \text{supp}(\mu_{\text{tplt}})$  such that*

$$\mathbb{E}_{\theta^t} [(f_{\text{MLP}}(\mathbf{x}; \theta^t) - f_*(\mathbf{z}))^2] \geq c > 0,$$

where  $c$  is constant that depends only on  $\mu_{\text{tplt}}$  and  $f_*$ .

The proof is deferred to Appendix I, and relies on the key observation that SGD-training of MLPs satisfies a permutation invariance property (Ng, 2004). This property guarantees that MLP cannot consistently distinguish between the unseen tokens, and therefore, in expectation over the weights  $\theta^t$ , outputs the same value for any sequence  $\mathbf{x} \in (\mathcal{X}_{\text{uns}})^k$ . We make four remarks.

**Remark I.3.** MLPs are universal approximators (Cybenko, 1989), so there are choices of weights  $\theta$  such that  $f_{\text{MLP}}(\cdot; \theta)$  has good generalization on unseen symbols. The theorem proves that these weights are not found by SGD.

**Remark I.4.** The theorem does not assume that training is in the NTK regime, i.e., it holds even for nonlinear training dynamics.

**Remark I.5.** The theorem also holds for training with Adam, gradient flow, and minibatch-SGD, since the permutation-invariance property of MLP training also holds for these. See Appendix I.

**Remark I.6.** As a sanity check, we verify that MLP kernel does not meet the sufficient condition for generalizing on unseen symbols from Lemma 4.5. The kernel for an MLP is an inner product kernel of the form  $K_{\text{MLP}}(\mathbf{x}, \mathbf{x}') = \kappa(\sum_{i=1}^k 1(x_i = x'_i))$  for a function  $\kappa : \mathbb{R} \rightarrow \mathbb{R}$ . Therefore, the matrix  $N \in \mathbb{R}^{r \times r}$  has all of its entries equal to  $N_{ij} = \kappa(0)$ , so it is singular and the condition of Lemma 4.5 is not met.

We now prove Theorem I.2. We first show that trained MLPs cannot differentiate between tokens in the set  $\mathcal{X}_{\text{uns}}$ . Let  $\mathcal{X} = \mathcal{X}_{\text{seen}} \sqcup \mathcal{X}_{\text{uns}}$  be the partition of tokens into those seen and not seen in the train data. Here  $\mathcal{X}_{\text{seen}}$  is defined as the smallest set such that  $\mathbf{x} \in \mathcal{X}_{\text{seen}}^k$  almost surely for  $(\mathbf{x}, y) \sim \mathcal{D}$ .

**Lemma I.7** (Trained MLPs cannot distinguish unseen tokens). *For any number of SGD steps  $t$ , and any learning rate schedule  $\eta_1, \dots, \eta_t$ , the learned MLP estimator cannot distinguish between sequences of unseen tokens. Formally, for any  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}_{\text{uns}}^k$ , we have*

$$\mathbb{E}_{\theta^t} [f_{\text{MLP}}(\mathbf{x}_1; \theta^t)] = \mathbb{E}_{\theta^t} [f_{\text{MLP}}(\mathbf{x}_2; \theta^t)].$$

*Proof of Lemma I.7.* The proof of this result is based on a well-known permutation-invariance property of MLPs trained by SGD. This property has previously been used to show sample complexity

lower bounds for learning with SGD-trained MLPs (Ng, 2004; Li et al., 2020), as well as time-complexity lower bounds (Shamir, 2018; Abbe et al., 2022; Abbe & Boix-Adsera, 2022). In this lemma, we use the permutation invariance property to show poor out-of-distribution generalization of SGD-trained MLPs.

First, construct a permutation  $\Pi \in \mathbb{R}^{km \times km}$  such that  $\Pi z_0(x_1) = z_0(x_2)$ , but which also satisfies that for any  $\tilde{x} \in (\mathcal{X}_{seen})^k$  we have  $\Pi z_0(\tilde{x}) = z_0(\tilde{x})$ . This permutation can be easily constructed since neither  $x_1$  nor  $x_2$  contains tokens in  $\mathcal{X}_{seen}$ . Next, define the following network  $f_{\text{MLP}}^\Pi$ , analogously to (23) but with the first-layer inputs permuted by  $\Pi$

$$\begin{aligned} f_{\text{MLP}}^\Pi(x; \theta) &= w^T z_L^\Pi(x; \theta) \in \mathbb{R} \quad \text{where} \\ z_\ell^\Pi(x; \theta) &= \phi(W_\ell z_{\ell-1}^\Pi(x; \theta)) \in \mathbb{R}^d \quad \text{for } \ell \geq 1 \\ z_0^\Pi(x; \theta) &= z_0^\Pi(x) = \Pi[e_{x_1}, \dots, e_{x_k}] \in \mathbb{R}^{km}. \end{aligned}$$

Now let us couple the weights  $\theta^0, \dots, \theta^t$  from SGD training of  $f_{\text{MLP}}$  on dataset  $\mathcal{D}$ , with the weights  $\theta^{\Pi,0}, \dots, \theta^{\Pi,t}$  from SGD training of  $f_{\text{MLP}}^\Pi$  on dataset  $\mathcal{D}$ . The coupling is performed inductively on the time step, and we can maintain the property that  $\theta^\tau = \theta^{\Pi,\tau}$  for all  $t$ . For the base case  $\tau = 0$ , we set  $\theta^0 = \theta^{\Pi,0}$ . For the inductive step,  $\tau \geq 1$ , we update the weights with the gradient from some sample  $(x^\tau, y^\tau)$ . Since  $x^\tau \in (\mathcal{X}_{seen})^k$  almost surely, we know that  $z_0(x^\tau) = z_0^\Pi(x^\tau)$  almost surely, which means that  $\theta^\tau = \theta^{\Pi,\tau}$  almost surely. We conclude the equality in distribution of the weights

$$\theta^t \stackrel{d}{=} \theta^{\Pi,t}. \quad (24)$$

Next, let us inductively couple the weights  $\theta^0, \dots, \theta^t$  with the weights  $\theta^{\Pi,0}, \dots, \theta^{\Pi,t}$  in a different way, so as to guarantee that for any time  $0 \leq \tau \leq t$ , we have

$$W_1^\tau = W_1^{\Pi,\tau} \Pi \text{ and } W_\ell^\tau = W_\ell^{\Pi,\tau} \text{ for all } 2 \leq \ell \leq L \text{ and } w^\tau = w^{\Pi,\tau}.$$

almost surely. The base case  $\tau = 0$  follows because the distribution of  $W_1^0$  and  $W_1^{\Pi,0}$  is equal and is also invariant to permutations since it is Gaussian. For the inductive step, couple the sample updates so that SGD draws the same sample  $(x^\tau, y^\tau) \sim \mathcal{D}$ . One can see from the chain rule that the invariant is maintained. We conclude the equality in distribution of the weights

$$\theta^t = \{W_1^t, \dots, W_L^t, w^t\} \stackrel{d}{=} \{W_1^{\Pi,t} \Pi, W_2^{\Pi,t}, \dots, W_L^{\Pi,t}, w^{\Pi,t}\} \quad (25)$$

Combining (24) and (25), we get

$$\theta^t = \{W_1^t, \dots, W_L^t, w^t\} \stackrel{d}{=} \{W_1^t \Pi, W_2^t, \dots, W_L^t, w^t\},$$

which, since  $\Pi z_0(x_1) = z_0(x_2)$ , immediately implies

$$f_{\text{MLP}}(x_1; \theta^t) = f_{\text{MLP}}(x_2; \{W_1^t \Pi, W_2^t, \dots, W_L^t, w^t\}) \stackrel{d}{=} f_{\text{MLP}}(x_2; \theta^t),$$

which proves the lemma.  $\square$

Theorem 1.2 follows as a consequence. Note that the key lemma proved above only relied on a permutation invariance property of SGD on MLPs that also holds for Adam training, gradient flow training, and SGD with minibatch (see Li et al. (2020)). Therefore, the result holds for training with those algorithms as well, beyond just SGD.

*Proof of Theorem 1.2.* Pick any two templates  $z, z' \in \text{supp}(\mu_{\text{tmplt}})$  such that  $f_*(z) \neq f_*(z')$ . Recall that  $z, z' \in \mathcal{W}^k$  by assumption. Since we assumed that  $|\mathcal{X}_{uns}| \geq k$ , there are strings  $x, x' \in \mathcal{X}_{uns}^k$  matching templates  $z$  and  $z'$ , respectively. Furthermore, by Lemma 1.7, if we define  $a = \mathbb{E}_{\theta^t}[f_{\text{MLP}}(x; \theta^t)] = \mathbb{E}_{\theta^t}[f_{\text{MLP}}(x'; \theta^t)]$ , we have

$$\begin{aligned} & \max(\mathbb{E}_{\theta^t}[(f_{\text{MLP}}(x; \theta^t) - f_*(z))^2], \mathbb{E}_{\theta^t}[(f_{\text{MLP}}(x'; \theta^t) - f_*(z'))^2]) \\ & \geq \max((a - f_*(z))^2, (a - f_*(z'))^2) \\ & \geq \frac{1}{4}(f_*(z) - f_*(z'))^2 = c > 0. \end{aligned}$$

$\square$



## J DEFERRED DETAILS FOR SYMBOLIC-LABEL TEMPLATE TASKS

### J.1 DEFINITION OF SYMBOLIC-LABEL TEMPLATE TASKS

In symbolic-label template tasks the output is a token in  $\mathcal{X}$ . This corresponds to the next-token prediction setting, and the appropriate loss is the cross-entropy loss for multiclass classification. The formal definition of these tasks is:

**Definition J.1** (Multi-class prediction version of template). The data distribution  $\mathcal{D}_{multiclass} = \mathcal{D}_{multiclass}(\mu_{\text{tmpl}}, \{\mu_{\text{sub},z}\}, f_*)$  is specified by: (i) a template distribution  $\mu_{\text{tmpl}}$  supported on  $(\mathcal{X} \cup \mathcal{W})^k$ ; (ii) for each template  $z$ , a distribution  $\mu_{\text{sub},z}$  over substitution maps  $s : \mathcal{W} \rightarrow \mathcal{X}$ ; (iii) a labelling function  $f_* : \text{supp}(\mu_{\text{tmpl}}) \rightarrow \mathcal{X} \cup \mathcal{W}$ . A sample  $(x, y) \in \mathcal{X}^k \times \mathcal{X}$  drawn from  $\mathcal{D}_{multiclass}$  is drawn by taking  $x = \text{sub}(z, s)$  and  $y = \text{sub}(f_*(z), s)$ , where  $z \sim \mu_{\text{tmpl}}$  and  $s \sim \mu_{\text{sub},z}$ .

### J.2 FAILURE OF TRANSFORMERS TO COPY AND MODIFICATION THAT SUCCEEDS

We provide the deferred proofs for Section 5.

**Attention layer architecture** For simplicity in this section we consider a transformer with the attention layer only, since the MLP layer does not play a role in the ability to copy unseen symbols. Our architecture has  $H$  heads with parameters  $\mathbf{W}_{K,h}, \mathbf{W}_{Q,h}, \mathbf{W}_{V,h}, \mathbf{W}_{O,h} \in \mathbb{R}^{d_{\text{head}} \times d_{\text{emb}}}$ , an embedding/unembedding layer  $\mathbf{W}_E \in \mathbb{R}^{m \times d_{\text{emb}}}$ , positional embeddings  $\mathbf{P} \in \mathbb{R}^{k \times d_{\text{emb}}}$ , an MLP layer with parameters  $\mathbf{W}_A, \mathbf{W}_B \in \mathbb{R}^{d_{\text{mlp}} \times d_{\text{emb}}}$ , a final unembedding layer, and an activation function  $\phi$ . The network takes in  $\mathbf{X} \in \mathbb{R}^{k \times m}$  and outputs

$$f_{\text{attn}}(\mathbf{X}; \theta) = \mathbf{W}_E \mathbf{z}_1 \in \mathbb{R}^m \quad (\text{Unembedding layer})$$

where

$$\begin{aligned} \mathbf{z}_1 &= \sum_{h \in [H]} \mathbf{A}_h^T \mathbf{e}_k \\ \mathbf{A}_h &= \text{smax}(\beta \mathbf{Z}_0 \mathbf{W}_{K,h}^T \mathbf{W}_{Q,h} \mathbf{Z}_0^T) \mathbf{Z}_0 \mathbf{W}_{V,h}^T \mathbf{W}_{O,h} \in \mathbb{R}^{k \times d_{\text{emb}}} \quad (\text{Attention heads}) \\ \mathbf{Z}_0 &= \mathbf{X} \mathbf{W}_E + \gamma \mathbf{P} \in \mathbb{R}^{k \times d_{\text{emb}}}. \quad (\text{Embedding layer}) \end{aligned}$$

and we tie the embedding and unembedding weights, as often done in practice, for example in GPT-2 [Brown et al. \(2020\)](#). Here  $\beta, \gamma \geq 0$  are two hyperparameters that control the inverse temperature of the softmax and the strength of the positional embeddings, respectively.

**Simplification in our case** We consider here a next-token prediction setup, where there is no final [CLS] token appended to the string. Namely, given a string  $x \in \mathcal{X}^k$ , this is inputted to the network as a stacked matrix of one-hot vectors for the tokens of the string  $\mathbf{X} = [e_{x_1}, \dots, e_{x_k}]$ . We study a very basic template task: template “ $\alpha$ ” labeled by  $\alpha$ , where  $\alpha$  is a wildcard. An example dataset generated from this template could be  $\{(A, A), (B, B), (C, C)\}$ , where  $A, B, C \in \mathcal{X}$  are tokens. Because the template has length  $k = 1$ ,  $\mathbf{X} \in \mathbb{R}^{k \times m}$  is a one-hot vector encoding the input token. Furthermore, the softmax output is always a  $1 \times 1$  matrix with the entry 1, so the architecture simplifies to

$$f_{\text{attn}}(\mathbf{X}; \theta) = \mathbf{W}_E \left( \sum_{h \in [H]} \mathbf{W}_{O,h}^T \mathbf{W}_{V,h} \right) (\mathbf{W}_E^T \mathbf{X}^T + \gamma \mathbf{P}^T). \quad (26)$$

We initialize the entries of  $\mathbf{P}$  and  $\mathbf{W}_E$  be i.i.d.  $N(0, 1/d_{\text{emb}})$ , the entries of  $\mathbf{W}_{O,h}$  be  $N(0, 1/(d_{\text{emb}}))$ , and the entries of  $\mathbf{W}_{V,h}$  be  $N(0, 1/d_{\text{head}})$ , so that as  $d_{\text{emb}} \rightarrow \infty$  the variance of the output vanishes as  $O(1/d_{\text{emb}})$  as in the mean-field scaling [Mei et al. \(2018; 2019\)](#); [Sirignano & Spiliopoulos \(2022\)](#); [Chizat & Bach \(2018\)](#); [Rotskoff & Vanden-Eijnden \(2018\)](#); [Yang & Hu \(2021\)](#).

**Derivation of kernels driving dynamics at small times** Despite the simplicity of the task, the architecture does not generalize well on unseen symbols. Our evidence for this will be by analyzing the early times of training. For these times, the dynamics are governed by the neural tangent kernel (NTK) of the network at initialization ([Jacot et al., 2018](#); [Chizat et al., 2019](#)). Let us derive the

neural tangent kernel of this architecture. This is a network with output of dimension  $m$ , so for each  $i, j \in [m]$  we will derive  $K_{ij,O}(\mathbf{X}, \mathbf{X}')$ ,  $K_{ij,V}(\mathbf{X}, \mathbf{X}')$ ,  $K_{ij,P}(\mathbf{X}, \mathbf{X}')$ ,  $K_{ij,E}(\mathbf{X}, \mathbf{X}')$  which give the dynamics at small times for training the  $\{\mathbf{W}_{O,h}\}_{h \in [H]}$ , the  $\{\mathbf{W}_{V,h}\}_{h \in [H]}$ , the  $\mathbf{W}_P$ , and the  $\mathbf{W}_E$  weights at small times, respectively. Writing  $\mathbf{W}_E = [\mathbf{w}_{E,1}, \dots, \mathbf{w}_{E,m}]^\top$ , by the law of large numbers,

$$\begin{aligned} K_{ij,O}(\mathbf{X}, \mathbf{X}') &= \sum_{h \in [H]} \left( \frac{\partial[f_{\text{attn}}(\mathbf{X}; \boldsymbol{\theta})]_i}{\partial \mathbf{W}_{O,h}} \right)^T \left( \frac{\partial[f_{\text{attn}}(\mathbf{X}'; \boldsymbol{\theta})]_j}{\partial \mathbf{W}_{O,h}} \right) \\ &\propto \frac{1}{H} \sum_{h \in [H]} (\mathbf{X} \mathbf{W}_E + \gamma \mathbf{P}) \mathbf{W}_{V,h}^T \mathbf{W}_{V,h} (\mathbf{W}_E^T \mathbf{X}^T + \gamma \mathbf{P}^T) \mathbf{w}_{E,i}^T \mathbf{w}_{E,j} \\ &\xrightarrow{d_{\text{head}} \rightarrow \infty, d_{\text{emb}} \rightarrow \infty} \delta_{ij} (\delta_{x_1, x'_1} + \gamma^2) \end{aligned}$$

$$\begin{aligned} K_{ij,V}(\mathbf{X}, \mathbf{X}') &= \sum_{h \in [H]} \left( \frac{\partial[f_{\text{attn}}(\mathbf{X}; \boldsymbol{\theta})]_i}{\partial \mathbf{W}_{V,h}} \right)^T \left( \frac{\partial[f_{\text{attn}}(\mathbf{X}'; \boldsymbol{\theta})]_j}{\partial \mathbf{W}_{V,h}} \right) \\ &\propto \frac{d_{\text{emb}}}{d_{\text{head}}} \sum_{h \in [H]} \mathbf{w}_{E,i}^T \mathbf{W}_{O,h}^T \mathbf{W}_{O,h} \mathbf{w}_{E,j} (\mathbf{X} \mathbf{W}_E + \gamma \mathbf{P})^T (\mathbf{X}' \mathbf{W}_E + \gamma \mathbf{P}) \\ &\xrightarrow{d_{\text{head}} \rightarrow \infty} \mathbf{w}_{E,i}^T \mathbf{w}_{E,j} (\mathbf{X} \mathbf{W}_E + \gamma \mathbf{P})^T (\mathbf{X}' \mathbf{W}_E + \gamma \mathbf{P}) \\ &\xrightarrow{d_{\text{emb}} \rightarrow \infty} \delta_{ij} (\delta_{x_1, x'_1} + \gamma^2) \end{aligned}$$

$$K_{ij,P}(\mathbf{X}, \mathbf{X}') = \left( \frac{\partial[f_{\text{attn}}(\mathbf{X}; \boldsymbol{\theta})]_i}{\partial \mathbf{P}} \right)^T \left( \frac{\partial[f_{\text{attn}}(\mathbf{X}'; \boldsymbol{\theta})]_j}{\partial \mathbf{P}} \right) = \gamma^2 \mathbf{w}_{E,i}^\top \mathbf{w}_{E,j} \xrightarrow{d_{\text{emb}} \rightarrow \infty} \gamma^2 \delta_{ij}$$

$$\begin{aligned} K_{ij,E}(\mathbf{X}, \mathbf{X}') &= \left( \frac{\partial[f_{\text{attn}}(\mathbf{X}; \boldsymbol{\theta})]_i}{\partial \mathbf{W}_E} \right)^T \left( \frac{\partial[f_{\text{attn}}(\mathbf{X}'; \boldsymbol{\theta})]_j}{\partial \mathbf{W}_E} \right) \\ &= \delta_{ij} (\mathbf{X} \mathbf{W}_E + \gamma \mathbf{P}) \left( \sum_{h \in [H]} \mathbf{W}_{V,h}^T \mathbf{W}_{O,h} \right) \left( \sum_{h \in [H]} \mathbf{W}_{O,h}^T \mathbf{W}_{V,h} \right) (\mathbf{W}_E^T (\mathbf{X}')^T + \gamma \mathbf{P}^T) \\ &\quad + \delta_{x_1, x'_1} \mathbf{w}_{E,i}^T \left( \sum_{h \in [H]} \mathbf{W}_{O,h}^T \mathbf{W}_{V,h} \right) \left( \sum_{h \in [H]} \mathbf{W}_{V,h}^T \mathbf{W}_{O,h} \right) \mathbf{w}_{E,j}^T \\ &\quad + \delta_{i, x'_1} \mathbf{w}_{E,j}^T \left( \sum_{h \in [H]} \mathbf{W}_{O,h}^T \mathbf{W}_{V,h} \right) \left( \sum_{h \in [H]} \mathbf{W}_{O,h}^T \mathbf{W}_{V,h} \right) (\mathbf{w}_{E, x_1} + \gamma \mathbf{P}^T) \\ &\quad + \delta_{x_1, j} \mathbf{w}_{E,i}^T \left( \sum_{h \in [H]} \mathbf{W}_{O,h}^T \mathbf{W}_{V,h} \right) \left( \sum_{h \in [H]} \mathbf{W}_{O,h}^T \mathbf{W}_{V,h} \right) (\mathbf{w}_{E, x'_1} + \gamma \mathbf{P}^T) \\ &\xrightarrow{d_{\text{head}} \rightarrow \infty, d_{\text{emb}} \rightarrow \infty, H \rightarrow \infty} \delta_{ij} (2\delta_{x_1, x'_1} + \gamma^2), \end{aligned}$$

since only the first two terms do not vanish as the embedding dimension and number of heads go to infinity.

**Training loss and testing loss** Let  $(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}$  be a training set of data points drawn from this task, where due to the structure of the template task each of the context strings is length-1 and we have  $x_i = y_i$ . We will test the model on a data point  $(x^{\text{test}}, y^{\text{test}})$ , which does not appear in the test set: i.e.,  $x^{\text{test}} = y^{\text{test}} \notin \{x_1, \dots, x_n\}$ .

The training loss is given by

$$\mathcal{L}_{\text{train}}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(f_{\text{attn}}(x_i; \boldsymbol{\theta}), y_i),$$

where  $\ell$  is the cross-entropy loss, and the test loss is given by

$$\mathcal{L}_{test}(\boldsymbol{\theta}) = \ell(f_{\text{attn}}(x^{test}), y^{test}).$$

**Theorem J.2.** For any learning rates  $\eta_O, \eta_V, \eta_P, \eta_E$  such that  $|\frac{\partial \mathcal{L}_{train}}{\partial t}| = O(1)$  as  $d_{emb}, d_{head}$ , and  $H \rightarrow \infty$ , we have  $|\frac{\partial \mathcal{L}_{test}}{\partial t}| \leq o(1)$ . In other words, the error for generalization on unseen symbols does not decrease during training for infinite-width transformers.

*Proof.* Consider training with gradient flow with learning rates  $\eta_O, \eta_V, \eta_P, \eta_E$  on the parameters  $\{\mathbf{W}_{O,h}\}_{h \in [H]}, \{\mathbf{W}_{V,h}\}_{h \in [H]}, \mathbf{W}_P$ , and  $\mathbf{W}_E$ , respectively. In the limit as  $d_{emb} \rightarrow \infty$  we have  $f_{\text{attn}}(\mathbf{X}; \boldsymbol{\theta}_0) \rightarrow 0$ , so

$$\frac{\partial \mathcal{L}_{train}}{\partial \boldsymbol{\theta}} \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} = \frac{1}{n} \sum_{i=1}^n \left( \frac{1}{m} \mathbf{1} - \mathbf{e}_{x_i} \right)^T \frac{\partial f_{\text{attn}}(\mathbf{X}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}.$$

So at time  $t = 0$ , the training loss decreases as

$$\begin{aligned} \frac{\partial \mathcal{L}_{train}}{\partial t} \big|_{t=0} \rightarrow & -\frac{1}{n^2} \sum_{i,i' \in [n]} \sum_{j,j' \in [m]} (1/m - \delta_{j,x_i})(1/m - \delta_{j',x_{i'}}) \\ & \cdot (\eta_V K_{jj',V}(\mathbf{X}_i, \mathbf{X}_{i'}) + \eta_O K_{jj',O}(\mathbf{X}_i, \mathbf{X}_{i'}) \\ & + \eta_P K_{jj',P}(\mathbf{X}_i, \mathbf{X}_{i'}) + \eta_E K_{jj',E}(\mathbf{X}_i, \mathbf{X}_{i'})). \end{aligned}$$

So we must take  $\eta_O = O(1/H)$ ,  $\eta_V = O(d_{emb}/d_{head})$ ,  $\eta_P = O(1)$ , and  $\eta_E = O(1)$  for us to have  $\frac{\partial \mathcal{L}_{train}}{\partial t} = O(1)$  be bounded by a constant that does not grow with  $d_{emb}$ ,  $d_{head}$ , and  $H$ .

Under these choices of learning rates, the test loss on token  $x^{test}$  which is not in the training dataset  $\{x_1, \dots, x_n\}$ , evolves as

$$\begin{aligned} \frac{\partial \mathcal{L}_{test}}{\partial t} \big|_{t=0} \rightarrow & -\frac{1}{n} \sum_{i \in [n]} \sum_{j,j' \in [m]} (1/m - \delta_{j,x_i})(1/m - \delta_{j',x^{test}}) \\ & \cdot (\eta_V K_{jj',V}(\mathbf{X}_i, \mathbf{X}^{test}) + \eta_O K_{jj',O}(\mathbf{X}_i, \mathbf{X}^{test}) \\ & + \eta_P K_{jj',P}(\mathbf{X}_i, \mathbf{X}^{test}) + \eta_E K_{jj',E}(\mathbf{X}_i, \mathbf{X}^{test})) \\ \rightarrow & -\frac{1}{n} \sum_{i \in [n]} \sum_{j,j' \in [m]} (1/m - \delta_{j,x_i})(1/m - \delta_{j',x^{test}}) \\ & \cdot \left( \left( \frac{d_{head}}{d_{emb}} \eta_V + H \eta_O \right) \delta_{j,j'} (\delta_{x_i, x^{test}} + \gamma^2) \right. \\ & \left. + \eta_P \gamma^2 \delta_{j,j'} + 2H \eta_E \delta_{j,j'} (\delta_{x_i, x^{test}} + \gamma^2) \right) \\ = & -\frac{\gamma^2}{n} \sum_{i \in [n]} \sum_{j \in [m]} (1/m - \delta_{j,x_i})(1/m - \delta_{j,x^{test}}) \cdot \left( \frac{d_{head}}{d_{emb}} \eta_V + H \eta_O + \eta_P + 2\eta_E \right) \\ = & -\frac{C}{n} \sum_{i \in [n]} \sum_{j \in [m]} (1/m - \delta_{j,x_i})(1/m - \delta_{j,x^{test}}) \\ = & -C/m + C/m + C/m = C/m \geq 0. \end{aligned}$$

□

On the other hand, now we consider the  $f_{\text{attn}}$  architecture where in each head we replace  $\mathbf{W}_{V,h}^T \mathbf{W}_{O,h}$  with  $\mathbf{W}_{V,h}^T \mathbf{W}_{O,h} + b_h \mathbf{I}$ , where  $b_h$  is a trainable parameter and  $\mathbf{I} \in \mathbb{R}^{d_{emb} \times d_{emb}}$  is the identity matrix:

$$f'_{\text{attn}}(\mathbf{X}; \boldsymbol{\theta}) = \mathbf{W}_E \mathbf{z}_1 \in \mathbb{R}^m \quad (\text{Unembedding layer})$$

where

$$\begin{aligned} \mathbf{z}'_1 &= \sum_{h \in [H]} (\mathbf{A}'_h)^T \mathbf{e}_k \\ \mathbf{A}'_h &= \text{smax}(\beta \mathbf{Z}_0 \mathbf{W}_{K,h}^T \mathbf{W}_{Q,h} \mathbf{Z}_0^T) \mathbf{Z}_0 (\mathbf{W}_{V,h}^T \mathbf{W}_{O,h} + b_h \mathbf{I}) \in \mathbb{R}^{k \times d_{emb}} \quad (\text{Attention heads}) \\ \mathbf{Z}_0 &= \mathbf{X} \mathbf{W}_E + \gamma \mathbf{P} \in \mathbb{R}^{k \times d_{emb}}. \quad (\text{Embedding layer}) \end{aligned}$$

Again, for the case of  $k = 1$  that we consider, the network simplifies considerably to

$$f'_{\text{attn}}(\mathbf{X}; \boldsymbol{\theta}) = \mathbf{W}_E \left( \sum_{h \in [H]} \mathbf{W}_{O,h}^T \mathbf{W}_{V,h} + b_h \mathbf{I} \right) (\mathbf{W}_E^T \mathbf{X}^T + \gamma \mathbf{P}^T). \quad (27)$$

We initialize  $b_h = 0$  for all  $h$ , so that the neural tangent kernels  $K_{ij,O}, K_{ij,V}, K_{ij,P}, K_{ij,E}$  are the same as above. Now we also have a neural tangent kernel for training the parameters  $\{b_h\}_{h \in [H]}$ :

$$\begin{aligned} K_{ij,b}(\mathbf{X}, \mathbf{X}') &= \sum_{h \in [H]} \frac{\partial [f_{\text{attn}}(\mathbf{X}; \boldsymbol{\theta})]_i}{\partial b_h} \frac{\partial [f_{\text{attn}}(\mathbf{X}'; \boldsymbol{\theta})]_j}{\partial b_h} \\ &\propto \mathbf{w}_{E,i}^\top (\mathbf{W}_E^T \mathbf{X}^T + \gamma \mathbf{P}^T) (\mathbf{X} \mathbf{W}_E + \gamma \mathbf{P}^T) \mathbf{w}_{E,j} \\ &\stackrel{d_{\text{emb}} \rightarrow \infty}{\rightarrow} \delta_{i,x_1} \delta_{j,x'_1} \end{aligned}$$

We prove that under this parametrization the test loss does decrease with training, which shows that adding this trainable identity scaling allows transformers to succeed at this task.

**Theorem J.3.** *There is a choice of learning rates  $\eta_b, \eta_V, \eta_O, \eta_E, \eta_P$  such that as  $d_{\text{emb}}, d_{\text{head}}, H \rightarrow \infty$  we have  $|\frac{\partial \mathcal{L}_{\text{train}}}{\partial t}|_{t=0} = O(1)$  and  $-\frac{\partial \mathcal{L}_{\text{test}}}{\partial t}|_{t=0} = \Omega(1)$ .*

*Proof.* Training just the parameters  $\{b_h\}_{h \in [H]}$  with learning rate  $\eta_b$  (keeping the learning rates  $\eta_V, \eta_O, \eta_P, \eta_E = 0$ , so the training loss decreases as

$$\frac{\partial \mathcal{L}_{\text{train}}}{\partial t} \Big|_{t=0} \rightarrow -\frac{\eta_b}{n^2} \sum_{i,i' \in [n]} \sum_{j,j' \in [m]} (1/m - \delta_{j,x_i}) (1/m - \delta_{j',x_{i'}}) K_{jj',b}(\mathbf{X}_i, \mathbf{X}_{i'}),$$

so we should take  $\eta_b = \Theta(1/H)$  for the train loss have derivative on the order of  $\Theta(1)$ . The test loss decreases as:

$$\begin{aligned} \frac{\partial \mathcal{L}_{\text{test}}}{\partial t} \Big|_{t=0} &\rightarrow -\frac{\eta_b}{n} \sum_{i \in [n]} \sum_{j,j' \in [m]} (1/m - \delta_{j,x_i}) (1/m - \delta_{j',x_{i'}^{\text{test}}}) K_{jj',b}(\mathbf{X}_i, \mathbf{X}_{i'}^{\text{test}}) \\ &\rightarrow -\frac{H\eta_b}{n} \sum_{i \in [n]} \sum_{j,j' \in [m]} (1/m - \delta_{j,x_i}) (1/m - \delta_{j',x_{i'}^{\text{test}}}) \delta_{j,x_i} \delta_{j',x_{i'}^{\text{test}}} \\ &= -\frac{H\eta_b}{n} \sum_{i \in [n]} (1/m - 1)(1/m - 1) \\ &= -H\eta_b(1 - 1/m)^2 \\ &= \Omega(1), \end{aligned}$$

for  $\eta_b = \Omega(H)$ , as  $d_{\text{emb}}, H \rightarrow \infty$ . □