
Transformers on Markov data: Constant depth suffices

Nived Rajaraman *
UC Berkeley

Marco Bondaschi
EPFL

Kannan Ramchandran
UC Berkeley

Michael Gastpar
EPFL

Ashok Vardhan Makkuva
EPFL

Abstract

Attention-based transformers have been remarkably successful at modeling generative processes across various domains and modalities. In this paper, we study the behavior of transformers on data drawn from k^{th} -order Markov processes, where the conditional distribution of the next symbol in a sequence depends on the previous k symbols observed. We observe a surprising phenomenon empirically which contradicts previous findings: when trained for sufficiently long, a transformer with a fixed depth and 1 head per layer is able to achieve low test loss on sequences drawn from k^{th} -order Markov sources, even as k grows. Furthermore, this low test loss is achieved by the transformer’s ability to represent and learn the in-context conditional empirical distribution. On the theoretical side, our main result is that a transformer with a single head and three layers can represent the in-context conditional empirical distribution for k^{th} -order Markov sources, concurring with our empirical observations. Along the way, we prove that *attention-only* transformers with $O(\log_2(k))$ layers can represent the in-context conditional empirical distribution by composing induction heads to track the previous k symbols in the sequence. These results provide more insight into our current understanding of the mechanisms by which transformers learn to capture context, by understanding their behavior on Markov sources. Code is available at: <https://github.com/Bond1995/Constant-depth-Transformers>.

1 Introduction

Attention-based transformers have revolutionized the field of natural language processing (NLP) [1, 2] and beyond [3, 4], achieving significant performance gains across tasks like machine translation, text generation, and sentiment analysis. A key factor in their success is their ability to model sequences far more efficiently, and the ability to learn in-context [5, 6].

To understand this capability, a canonical approach is to sample the input from a k^{th} -order Markov process, where the next symbol’s conditional distribution depends only on the previous k symbols. Recent studies [7, 6, 8] have investigated the ability of transformers to learn Markov processes and establish that learning happens in phases. The transformer eventually learns to represent the conditional k -gram model, which is the in-context MLE of the Markov process.

The results in [6, 8] seem to suggest that for low depth transformers to learn Markov processes of order k , it is essential that the number of heads scale linearly in k . At first glance, this is a bit concerning - real world data generating processes often contain long-range dependencies. How is it that transformers succeed at capturing these kinds of long-range dependencies, while at the same time requiring so many heads to be able to capture the necessary context for k^{th} -order Markov sources?

*Correspondence to nived.rajaraman@berkeley.edu.

... 2 0 1 1 1 0 1 0 X_{n+1}

	Attention-only		Standard
L	2	$\lceil \log_2(k+1) \rceil$	3
H	k	1	1

Figure 1: k^{th} -order Markov process for $k = 4$. The symbol X_{n+1} is sampled from the distribution $P(\cdot | X_n, X_{n-1}, X_{n-2}, X_{n-3})$ which only depends on the last 4 symbols (marked in red).

Table 1: Each column in this table indicates that there is a transformer with L layers and H heads in the first layer which can represent the conditional k -gram model.²

To understand the nature of this phenomenon, we train low-depth transformers on k^{th} -order Markov sources. These experiments result in two surprising empirical phenomena that seem to contradict previous findings: when trained for sufficiently long, (i) a 2-layer, 1-head transformer can learn k^{th} -order Markov processes for k as large as 4, (ii) a 3-layer, 1-head transformer is able to achieve low test loss on sequences drawn from k^{th} -order Markov sources, even as k grows to be as large as 8 (Fig. 3). In both cases, the values of k for which the models appear to learn k^{th} -order Markov sources are much higher than those predicted in prior experiments [6, 8]. This discrepancy shows that our understanding of the mechanisms used by transformers to learn k^{th} -order Markov processes is not complete and raises a broader question:

What is the interplay between depth, number of heads and non-linearity in learning k^{th} -order Markov processes?

In this paper, we approach this question from the point of view of representation power, and provide some partial explanations toward the phenomena illustrated previously.

Our main contributions are as follows:

1. We show, rather surprisingly, that the standard transformer architecture with 3 layers and 1 head per layer is capable of representing the conditional k -gram model (Definition 1), and thereby learn k^{th} -order Markov models in-context.
2. Along the way to building up to this result, we consider the simpler family of *attention-only transformers* and show that they can represent the conditional k -gram model with $\lceil \log_2(k+1) \rceil$ layers.
3. Under a natural assumption on the nature of the attention patterns learnt by the transformer, we then argue that for $k \geq 3$ attention-only transformers *need* at least $\lceil 1 + \log_2(k-2) \rceil$ layers to represent a “ k^{th} -order induction head” (Definition 2). Empirically, transformers are observed to learn k^{th} -order induction heads whenever they achieve small test error [6].

The last result is a consequence of a more general tradeoff between the number of layers, L , and heads per layer, H , an attention-only transformer requires to represent a k^{th} -order induction head, under a natural assumption on the learnt attention patterns. In conjunction, these results also reveal the role of non-linearities (aside from the softmax in the attention) in the transformer architecture. In particular, it appears that layer normalization plays a critical role in the ability of constant-depth transformers to learn the conditional k -gram model. Together with the experimental results mentioned previously, these results paint a more comprehensive picture about the representation landscape of transformers in the context of k^{th} -order Markov processes.

Notation. Scalars are denoted by italic lower case letters like x, y and Euclidean vectors and matrices in bold $\mathbf{x}, \mathbf{y}, \mathbf{M}$, etc. The notation $\mathbf{0}_{p \times q}$ (resp. $\mathbf{1}_{p \times q}$) refers to the all-zero (resp. all-one) matrix. When it is clear from the context, we omit the dimensions of a matrix. Define $[S] \triangleq \{1, 2, \dots, S\}$ for $S \in \mathbb{N}$. $\mathbb{I}(\cdot)$ denotes the indicator function and $\text{Unif}(S)$ denotes the uniform distribution over a set S .

1.1 Related work

There is a large body of active research focused on studying different aspects of transformer models [9, 10, 11, 12]. Our work closely relates to the aspects of understanding the representation power of

²The requisite embedding dimension and bit-precision to achieve a target additive approximation is discussed in more detail in Sections 4 and 5

transformers, and in-context learning. [13, 14, 15] study the representation capabilities of transformers and show properties such as universal approximation and Turing-completeness. Viewing transformers as sequence to sequence models, [16, 17] study their ability to model formal languages and automata. Along more related lines to our work, [18, 19] present logarithmic depth transformer constructions for representing a k -hop generalization of the notion of an induction head [20]. On the other hand the theoretical and mechanistic understanding of in-context learning [21] has received much attention lately [22, 23, 24, 25], focusing on different operating regimes and phases of learning. There are a few recent papers which study the behavior of transformers when trained on data generated from Markov processes, and generalizations thereof [5, 26]. In particular, [7, 8] study the optimization landscape of gradient descent in learning generalizations of Markov processes, and [6] present a study of how transformers learn to represent in-context k -gram models, focusing on different phases of learning.

2 Preliminaries

We provide the necessary background for Markov processes, the conditional k -gram model, and the transformer architecture.

2.1 Markov processes

Markov processes are one of the widely used models in sequence modeling [27]. The characterizing property of these processes is that at any time step, the future evolution is only influenced by the most recent states. More formally, a sequence $(X_n)_{n \geq 1}$ is a k^{th} -order Markov process on a finite state space $[S]$ with the transition kernel P , if surely,

$$P(X_{n+1} | X_1, \dots, X_n) = P(X_{n+1} | X_{n-k+1}, \dots, X_n)$$

This property allows us to capture the conditional distribution at any position using only its previous k symbols. This motivates the notion of a conditional k -gram, its empirical counterpart, defined for any sequence (x_1, \dots, x_n) .

Definition 1 (Conditional k -gram model). *Given a sequence (x_1, \dots, x_n) of length n in $[S]^n$, the conditional k -gram model $\widehat{\text{Pr}}_k(\cdot | x_1, \dots, x_n)$ corresponds to the in-context estimate of the distribution over symbols conditioned on the last k symbols, i.e. for $x \in [S]$,*

$$\widehat{\text{Pr}}_k(x | x_1, \dots, x_n) \triangleq \frac{\sum_{i=k+1}^n \mathbb{I}(x_i = x, x_{i-1} = x_n, \dots, x_{i-k} = x_{n-k+1})}{\sum_{i=k+1}^n \mathbb{I}(x_{i-1} = x_n, \dots, x_{i-k} = x_{n-k+1})},$$

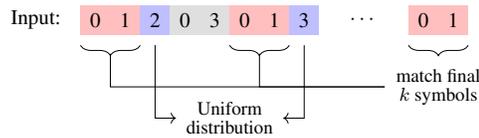
which is defined only so long as the denominator is non-zero. This structure is illustrated in Figure 2a. It is well known that the conditional k -gram in Eq. (1) with Laplace smoothing corresponds to the Bayes optimal estimate of the next symbol probability, when the data is drawn from fixed Markov process sampled from a prior distribution [27].

In our experiments, we will consider k^{th} -order Markov kernels sampled from a Dirichlet prior with parameter 1. Namely, the transition $P(\cdot | X_1 = i_1, \dots, X_k = i_k)$ is sampled independently and uniformly on the S -dimensional simplex Δ_1^S , for each tuple (i_1, \dots, i_k) .

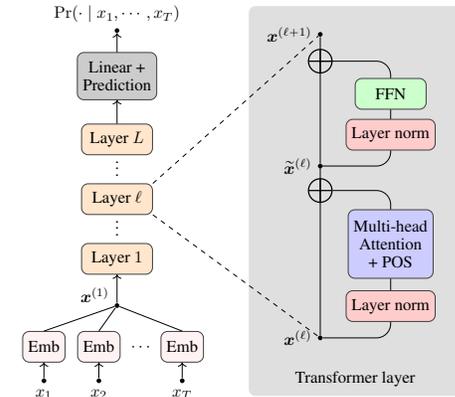
2.2 Transformer architecture

In this paper, we will consider variants of the standard transformer architecture in Figure 2b introduced in [1], with the goal to understand the role of depth and the non-linearities in the architecture. The simplest variant removes all the layer normalization and the (non-linear) feedforward layer, and is referred to as an *attention-only* transformer. The L -layer 1-head attention-only transformer with relative position encodings, operating on a sequence of length T is defined in Architecture 1.

The attention scores in layer ℓ , $\{\text{att}_{n,i}^{(\ell)} : i \leq n\}$, are computed as $\text{Softmax}(\{\langle \mathbf{W}_K^{(\ell)}(\mathbf{x}_j^{(\ell)} + \mathbf{p}_{n-j}^{(\ell),K}), \mathbf{W}_Q^{(\ell)} \mathbf{x}_j^{(\ell)} \rangle : j \in [n]\})$. The superscript (ℓ) indicates the layer index, and the matrices $\mathbf{W}_K^{(\ell)}, \mathbf{W}_Q^{(\ell)}, \mathbf{W}_V^{(\ell)} \in \mathbb{R}^{d \times d}$ capture the key, query and value matrices in layer ℓ . Note that the attention-only transformer may include a feedforward layer with linear activations, i.e. a linear



(a) Conditional k -gram model. The conditional k -gram is the in-context estimate of the Markov process and is realized in two steps. The first step is to find the locations in the sequence (marked red) which match the final k symbols (functionally, a k^{th} -order induction head). The conditional k -gram model returns the uniform distribution over the next symbol at these locations (marked blue).



(b) Transformer architecture. POS refers to the relative position encodings.

transformation. For representation purposes, this linear transformation can be combined with the projection matrix in the attention layer, allowing the feedforward layer to be omitted from the model. In the attention layer, we consider relative position encodings (the terms labeled in blue), which translates the key and value vectors depending on the relative position of the embedded symbol.

for $n = 1, 2, \dots, T$; $\mathbf{x}_n^{(1)} = \text{Emb}(x_n) \in \mathbb{R}^d$. (Input embeddings)

for $\ell = 1, 2, \dots, L$, **do**

for $n = 1, 2, \dots, T$, **do**

$$\tilde{\mathbf{x}}_n^{(\ell)} = \sum_{i \in [n]} \text{att}_{n,i}^{(\ell)} \cdot \mathbf{W}_V^{(\ell)} \left(\mathbf{x}_i^{(\ell)} + \mathbf{p}_{n-i}^{(\ell),V} \right) \in \mathbb{R}^d, \quad (\text{Attention})$$

$$\mathbf{x}_n^{(\ell+1)} = \mathbf{x}_n^{(\ell)} + \tilde{\mathbf{x}}_n^{(\ell)}, \quad (\text{Residual})$$

$$\triangleright \text{Here, } \text{att}_{n,i}^{(\ell)} = \text{Softmax}_i \left(\left\{ \left\langle \mathbf{W}_K^{(\ell)} (\mathbf{x}_i^{(\ell)} + \mathbf{p}_{n-i}^{(\ell),K}), \mathbf{W}_Q^{(\ell)} \mathbf{x}_n^{(\ell)} \right\rangle : i \in [S] \right\} \right).$$

$$\text{logit}_T = \mathbf{A} \mathbf{x}_T^{(L+1)} + \mathbf{b} \in \mathbb{R}^S, \quad (\text{Linear})$$

$$\Pr_{\theta}(\cdot | x_1, \dots, x_T) = f(\text{logit}_T) \in \mathbb{R}^S. \quad (\text{Prediction})$$

Architecture 1: Attention-only transformer.

The extension to H heads is straightforward, where in each transformer layer there are H attention layers in parallel, resulting in $\mathbf{y}_n^{(\ell,1)}, \dots, \mathbf{y}_n^{(\ell,H)} \in \mathbb{R}^d$ for each n . These vectors are concatenated and passed through a linear transformation $\mathbf{W}_O^{(\ell)} : \mathbb{R}^{dH} \rightarrow \mathbb{R}^d$ which is the output of the attention layer. Finally, the output of the model after L layers is passed through a linear layer, which projects the d -dimensional embeddings back into \mathbb{R}^S and the resulting vector is passed through a non-linearity f , usually a softmax, to result in the model's prediction of the next symbol probabilities. The theoretical results in this paper will choose $f = \text{ReLU}(\cdot)$.

3 Understanding the empirical behavior of transformers

The motivation for the present work comes from a series of experimental results, which challenge our current understanding of transformers in the context of learning Markov processes. Several works in the literature [7, 8, 6] have studied the ability of transformer models to learn k^{th} -order Markov processes. The experimental results present in the literature suggest that in order for a 2 layer transformer model to be able to learn a randomly sampled Markov process of order k , it is crucial for the number of heads in the first attention layer to scale linearly with the order, k . In particular, the authors of [6] claim that in their experiments, "Single attention headed models could not achieve better performance than bigram (models)" in learning random k^{th} -order Markov processes in-context. Similarly, the

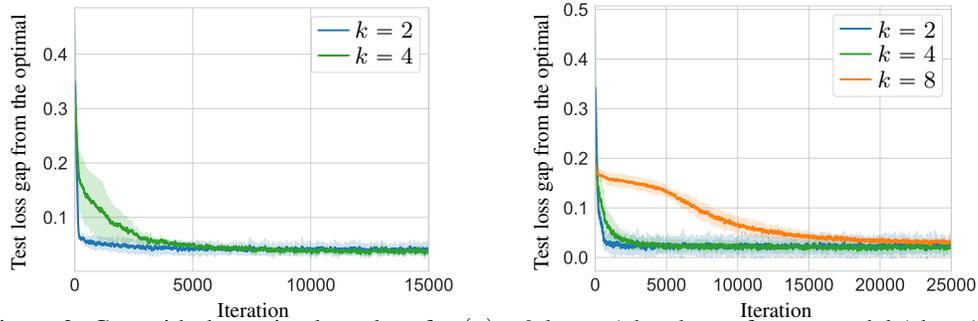


Figure 3: Gap with the optimal test loss for (a) a 2-layer, 1-head transformer model (above), and (b) a 3-layer, 1-head transformer (below), averaged over 3 runs for each k . The models learn the conditional k -gram model for randomly sampled k -th order Markov processes, even for large k .

authors of [8] study a generalization of learning k^{th} -order Markov processes to learning causal processes on degree k graphs. The theory and experiments pertain to 2-layer k head transformers.

In Figure 3, we train 2 and 3-layer transformers with a single head on data drawn from random Markov processes of various orders drawn from a Dirichlet prior. With 2 layers and a single head, we see that the model is able to learn even order-4 Markov processes, and go beyond the simple order-1 processes which were projected to be the limit of its ability to learn. Likewise, with 3 layers, transformers are able to go much further and learn order-8 Markov processes, which was the largest value of k we evaluated on. These results contrast with our current understanding of how induction heads are realized in the parameter space [6, 8] - existing constructions which realize these attention patterns require k heads when the number of layers is 2, and it's unclear how to implement them with fewer heads. At a high level, each of the k heads play a critical role - where, loosely speaking, the i^{th} -head looks back i positions in the sequence.

Building up to our main results, in the sequel, we study the simpler case of attention-only transformers where the feedforward layers and layer normalization are removed.

4 Warming up: Attention-only transformers

The study of attention-only transformers trained on Markov processes has garnered some attention in the prior literature. Notably, the authors of [6] study 2-layer 1-head attention-only transformers trained on data drawn from 1st-order Markov processes whose parameters are drawn from a Dirichlet prior. The model is observed to learn a very specific behavior, known as an “induction head” [20], which in this setting is able to represent the conditional 1-gram (Eq. (1)).

The induction head mechanism is composed of two layers where the first layer learns the attention pattern $\text{att}_{n,i}^{(1)} = \mathbb{I}(i = n - 1)$, thereby allowing the model to capture information about the symbol at position $n - 1$ in the embedding vector at time n . In the second layer, the attention layer picks out those indices n where $x_{n-1} = x_T$, the final symbol in the sequence. At these positions, since $x_{n-1} = x_T$, one would expect that the next symbol x_n is a good predictor of x_{T+1} , and the model uses this information to predict the next symbol x_{T+1} according to its conditional empirical estimate, $\widehat{\text{Pr}}_1(x_{T+1}|x_1, \dots, x_T)$, i.e. the conditional 1-gram model.

Theorem 1. *The conditional 1-gram model can be represented by a 2-layer and 1-head attention-only transformer with embedding dimension $d = 3S + 2$.*

Although a version of this result is proved in [6], we include a proof in Appendix A for completeness.

Remark 1. *In Theorem 1 and other results to follow, we de-emphasize the role of the bit-precision to which the transformer is implemented. That said, note that when the constructions in Theorems 1 to 3 are implemented to $O(\log(T))$ bits of precision, the representation results are realized up to an additive $O(1/T)$ error.*

The ideas in Theorem 1 readily extend to representing the conditional k -gram model, by instead using k heads in the first layer. The j^{th} head learns the attention pattern $\text{att}_{n,i}^{(1)} = \mathbb{I}(i = n - j)$ and

concatenating the outputs of the heads, the model learns to aggregate information about x_n, \dots, x_{n-k} in the embedding vector at time n . The second layer realizes what is best described as a “ k^{th} -order” induction head, where the model learns to pick out those positions n where for every $j \in [k]$, $x_{n-j} = x_{T-j+1}$, i.e. the history of length k at those positions match the final k symbols in the input sequence (see Figure 4). This mechanism is also referred to as a long-prefix induction head [28].

Definition 2 (Higher-order induction head). A 1-head attention layer is said to realize a k^{th} -order induction head if on any sequence $(x_1, \dots, x_T) \in [S]^T$, for any fixed $n \leq T$, as a function of the input sequence, $\text{att}_{n,T}$ is maximized if and only if $x_{n-j} = x_{T-j+1}$ for every $j \in [k]$.

k^{th} -order induction heads generalize the concept of an induction head [20], and keep track of the positions $i \leq n$ where there is a perfect occurrence of the final k symbols in the sequence. Such attention patterns are immediately useful in representing the conditional k -gram - increasing the temperature within the softmax of this attention layer results in an attention pattern which converges to the uniform distribution over those positions where the final k symbols x_{T-k+1}, \dots, x_T are seen previously in the sequence. Loosely, this allows the model to “condition” on the last k symbols in the sequence. With k heads, the model can aggregate information from the previous k positions and implement a k^{th} -order induction head, which leads to the following result. A full proof is discussed in Appendix A.1.

Theorem 2. The conditional k -gram model can be represented by an attention-only transformer with 2 layers, k heads and embedding dimension $d = (k + 2)S + k + 1$.

While this result is positive, it suggests that a 2-layer transformer requires approximately k times as many parameters to be able to represent the conditional k -gram model. The first result we prove is that increasing the depth of the model is exponentially more beneficial, in that a transformer with $O(\log(k))$ depth can estimate in-context k -grams.

Theorem 3. The conditional k -gram model can be represented by an attention-only transformer with relative position encodings, with $L = \lceil \log_2(k + 1) \rceil$ layers and 1 head per layer. The embedding dimension is $\leq 2k(S + 1) + S$.

With 2 layers and k heads, the transformer aggregates information about each of the previous k positions one step at a time through the k heads. However, with $\Omega(\log(k))$ layers, the same task can be done far more efficiently. In the first attention layer, the model aggregates information about the current and previous position. Namely, using the relative position embeddings, $\mathbf{x}_n^{(2)}$ is chosen as a linear combination of $\mathbf{x}_n^{(1)} = \text{Emb}(x_n)$ and $\mathbf{x}_{n-1}^{(1)} = \text{Emb}(x_{n-1})$. This allows the embedding at position n to aggregate information about x_n and x_{n-1} . In the same vein, in the second attention layer, the model aggregates information from $\mathbf{x}_n^{(2)}$ and $\mathbf{x}_{n-2}^{(2)}$ in $\mathbf{x}_n^{(3)}$; the former has information about x_n and x_{n-1} , and the latter has information about x_{n-2} and x_{n-3} . This expands the “window” of x_i ’s on which \mathbf{x}_n depends on to size 4. In the ℓ^{th} layer, the model aggregates information from $\mathbf{x}_n^{(\ell)}$ and $\mathbf{x}_{n-2^{\ell-1}}^{(\ell)}$ which allows $\mathbf{x}_n^{(\ell+1)}$ to effectively depend on the x_i ’s in a window of size $2^{\ell+1}$ starting at position n , namely $x_n, \dots, x_{n-2^{\ell+1}+1}$. In the final layer, the embedding at position i , $\mathbf{x}_i^{(L)}$ for $L = \lceil \log_2(k + 1) \rceil$ depends on $x_n, x_{n-1}, \dots, x_{n-k}$. In the last layer, the model can realize the dot-product $\langle \mathbf{W}_K^{(L)} \mathbf{x}_n^{(L)}, \mathbf{W}_Q^{(L)} \mathbf{x}_T^{(L)} \rangle = \sum_{j=1}^k \mathbb{I}(x_{n-j} = x_{T-j+1})$ by choosing the key and query vectors appropriately. By increasing the temperature in the attention softmax, the attention pattern realized is the uniform distribution on values of n such that $x_{n-j} = x_{T-j+1}$ for every $j \in [k]$, i.e., a k^{th} -order induction head. The full proof of this result is provided in Appendix B.

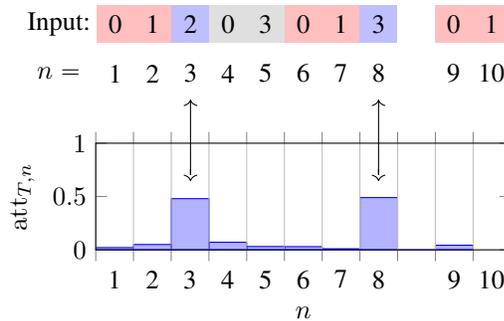
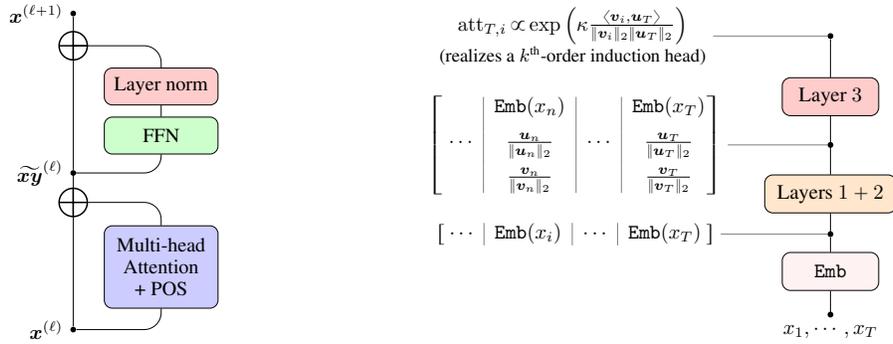


Figure 4: k^{th} -order induction head for $k = 2$. The attention pattern $\text{att}_{T,n}$ is maximized for those values of n at which $x_{T-j+1} = x_{n-j}$ for all $j \in [k]$. These are the positions where the k -length prefix at those positions matches with the last k symbols in the sequence.



(a) Rearranged transformer layer with layer normalization and FFN. (b) Realizing a k^{th} -order induction head in a 3-layer transformer following the architecture in Figure 5a.

Figure 5: *Disassembling the constant-depth construction.* The first two layers are critical in the model’s ability to capture information from the previous k positions. Layer normalization plays a critical role in the 3rd layer which realizes a k^{th} -order induction head.

While this is a promising step toward understanding the behavior transformers exhibit in Figure 3, showing that depth plays an important role in their ability to represent conditional k -gram models, the picture is still not complete. The experimental results in Section 3 do not preclude the possibility that a transformer might not even require logarithmic depth to be able to learn k^{th} -order Markov processes approximately. In the next section, we will study constant-depth transformers and establish a rather surprising positive result about the representation power of this class of models in capturing conditional k -grams.

5 Understanding the role of non-linearity: Constant-depth constructions

In the previous section, we saw how the transformer uses the power of depth to learn conditional k -grams far more efficiently. In particular, every additional attention layer effectively doubles the window of positions $i = n - 1, n - 2, \dots$ which the model has access to information about at the current time n . By composing $L = \Omega(\log(k))$ attention layers, the model is able to collect enough information within the output embedding $x_n^{(L+1)}$ to be able to realize a k^{th} -order induction head in the next layer. In this section, we prove that adding non-linearity to the architecture, in the form of layer normalization, can significantly change the mechanism in which the transformer realizes this k^{th} -order head. In particular, there are constant depth architectures which allow a k^{th} -order induction head to be realized, surpassing the logarithmic depth attention-only constructions.

Modification to the standard transformer architecture. To simplify the proof of our main result, we will consider a subtle modification to the standard transformer architecture, which is presented in Architecture 2 and Figure 5a. We will remove the first layer norm prior to the multi-head attention and move the second layer norm to after the feed-forward network. It is important to note that Theorem 4 holds even for the architecture presented in Figure 2b, which is the architecture we evaluate empirically. The modification we present in Figure 5a allows the construction to be simpler and makes it much easier to convey the key intuition. The main difference compared to the attention-only design presented in Architecture 1 is the addition of layer normalization and a feedforward layer in the for-loop over $n \in [T]$ for each transformer layer ℓ . The differences between Architectures 2 and 1 are emphasized in blue.

Theorem 4. *Conditional k -grams can be represented by a transformer with 3 layers, 1 head per layer, relative position encodings and layer normalization. The embedding dimension is $O(S)$.*

Remark 2. *Although the proof stated does not bound the approximation error arising from a finite bound on the bit precision of the transformer, in theory, it should suffice to have $\Omega(\log(T) + k)$ bits per parameter for the statement of Theorem 4 to go through with an $O(1/T)$ additive approximation error. The main point is that none of the weights of the model exceed $\exp(k)$ and with $\log(T)$ additional bits per parameter, the approximation error scales as $O(1/T)$.*

$$\tilde{\mathbf{x}}_n^{(\ell)} = \mathbf{x}_n^{(\ell)} + \sum_{i \in [n]} \text{att}_{n,i}^{(\ell)} \cdot \mathbf{W}_V^{(\ell)} \left(\mathbf{x}_i^{(\ell)} + \mathbf{p}_{n-i}^{(\ell),V} \right) \in \mathbb{R}^d, \quad (\text{Attention} + \text{Residual}_1)$$

$$\mathbf{y}_n^{(\ell)} = \mathbf{W}_2^{(\ell)} \text{ReLU} \left(\mathbf{W}_1^{(\ell)} \tilde{\mathbf{x}}_n^{(\ell)} \right) \in \mathbb{R}^d, \quad (\text{FFN})$$

$$\mathbf{l}_n^{(\ell)} = \frac{\mathbf{y}_n^{(\ell)} - \mu \mathbf{1}_{d \times 1}}{\sigma} \in \mathbb{R}^d, \quad (\text{LN})$$

$$\mathbf{x}_n^{(\ell+1)} = \mathbf{l}_n^{(\ell)} + \tilde{\mathbf{x}}_n^{(\ell)} \in \mathbb{R}^d, \quad (\text{Residual}_2)$$

Architecture 2: Modified transformer architecture. The computations above are carried out for each $n \in [T]$ in each layer $\ell \in [L]$. In the layer normalization step (LN), the feature mean μ is defined as, $\mathbb{E}_{i \sim \text{Unif}([d])} [\langle e_i^d, \mathbf{y}_n^{(\ell)} \rangle]$ and the feature variance $\sigma^2 = \mathbb{E}_{i \sim \text{Unif}([d])} [\langle e_i^d, \mathbf{y}_n^{(\ell)} \rangle^2] - \mu^2$.

5.1 Proof sketch

In the attention-only transformer with 2 layers and k heads, the model is able to keep track of where the final k symbols in the sequence appeared previously (i.e., a k^{th} -order induction head) by, loosely, using each head to keep track of the occurrences of one of the final k symbols. On the other hand, with the benefit of more depth, with $L = \Omega(\log(k))$ layers, the model is able to collect enough information within the output embedding $\mathbf{x}_n^{(L+1)}$ to be able to realize the same behavior. However, neither of these constructions scale down to the case when the depth and number of heads of the transformer are both constants independent of k . We provide a brief sketch of the construction below.

Recall that a k^{th} -order induction head keeps track of the indices i such that $\forall j \in [k], x_{i-j} = x_{n-j+1}$. Defining $\mathbf{z}_i \triangleq \sum_{j=1}^k 2^j e_{x_{i-j+1}}$, notice that the condition $\{\forall j \in [k], x_{i-j} = x_{n-j+1}\}$ can equivalently be captured by writing $\{\mathbf{z}_{i-1} = \mathbf{z}_n\}$. This true because of the fact that the binary representation of any integer is unique. Furthermore, these vectors, up to scaling, can be realized by softmax attention (namely, $\text{att}_{n,n-i} \propto 2^i$ for $1 \leq i \leq k$).

With this step, finding occurrences of the last k symbols in the input sequence boils down to realizing an attention pattern in the second layer, $\text{att}_{n,i}^{(2)}$, which is maximized whenever $\mathbf{z}_{i-1} = \mathbf{z}_n$. While dot-product attention naively encourages those values of i for which \mathbf{z}_{i-1} and \mathbf{z}_n are “similar” to each other, a qualitative statement is lacking. In general, it will turn out to that a different measure of similarity is necessary within the softmax to be able to encourage those values of i for which these vectors match. This is where the role of layer-normalization comes in.

Instead of the usual dot-product, suppose the attention mechanism in the second layer was,

$$\text{att}_{n,i}^{(2)} \propto \exp \left(-\kappa \left\| \frac{\mathbf{z}_{i-1}}{\|\mathbf{z}_{i-1}\|_2} - \frac{\mathbf{z}_n}{\|\mathbf{z}_n\|_2} \right\|_2^2 \right), \quad (1)$$

where κ is the temperature parameter. Then, as the temperature κ grows, the attention pattern essentially focuses on those values of i for which $\mathbf{z}_i / \|\mathbf{z}_{i-1}\|_2 = \mathbf{z}_n / \|\mathbf{z}_n\|_2$. With this attention pattern, we are thus very close to the statement we wanted to check, $(\mathbf{z}_{i-1} \stackrel{?}{=} \mathbf{z}_n)$. As it turns out, for the special structure in the \mathbf{z}_i 's considered (dyadic sums of one-hot vectors), we may write down,

$$\mathbf{z}_{i-1} = \mathbf{z}_n \iff \mathbf{z}_{i-1} / \|\mathbf{z}_{i-1}\|_2 = \mathbf{z}_n / \|\mathbf{z}_n\|_2.$$

A quantifiable equivalence is provided in Lemma 1.

Realizing L_2 -norm attention (eq. (1)). Observe the equivalence,

$$\left\langle \frac{\mathbf{z}_{i-1}}{\|\mathbf{z}_{i-1}\|_2}, \frac{\mathbf{z}_n}{\|\mathbf{z}_n\|_2} \right\rangle = 1 - \frac{1}{2} \left\| \frac{\mathbf{z}_{i-1}}{\|\mathbf{z}_{i-1}\|_2} - \frac{\mathbf{z}_n}{\|\mathbf{z}_n\|_2} \right\|_2^2 \quad (2)$$

Taking a softmax on both sides, notice that the RHS (up to an additive constant) is the L_2 -norm based attention, while the LHS is the usual dot-product attention between $\mathbf{z}_{i-1} / \|\mathbf{z}_{i-1}\|_2$ and $\mathbf{z}_n / \|\mathbf{z}_n\|_2$. Thus on unit-normalized vectors, L_2 -norm attention and dot product attention are but the same.

While the first layer of the transformer computes the \mathbf{z}_i 's by a weighted summation, layer normalization fills in the last missing piece of the puzzle which is to normalize them to unit norm. This is a consequence of defining the embedding vectors appropriately, as we discuss more in Appendix C.1.

From this step, realizing the actual conditional k -gram model follows readily. In particular, as the temperature κ in the attention grows, the attention pattern zooms in on indices $i \in \mathcal{I}_n \triangleq \{k+1 \leq i \leq n : \forall j \in [k], x_{i-j} = x_{n-j+1}\}$ in the last layer. The value vectors at this step are the one-hot encoding of x_i ; putting everything together, the logits realized by the transformer are,

$$\text{logit}_T(x_{T+1}) = \frac{1}{|\mathcal{I}_n|} \sum_{i \in \mathcal{I}_n} \mathbb{I}(x_i = x_T), \quad (3)$$

which is the conditional k -gram model (eq. (1)).

While the transformer construction described above only requires two layers, the actual construction we propose differs slightly and has an additional layer. The first two layers of the transformer respectively compute z_i and z_{i-1} which are added to the embedding vector at time i . This is important because we need to test whether $z_{i-1} \stackrel{?}{=} z_n$ and not whether $z_i \stackrel{?}{=} z_n$ or $z_{i-1} \stackrel{?}{=} z_n$.

Summary. The construction can be summarized as follows: the first layer computes $z_n = \sum_{j=1}^k 2^{j-1} \cdot e_{x_{n-j}}$ by choosing appropriate value vectors and relative position embeddings to realize the attention pattern $\text{attn}_{n,n-i} \propto 2^i \mathbb{I}(1 \leq i \leq k)$. The layernorm that follows subsequently can be replaced by RMSnorm, by a simple trick which we discuss in Appendix C.1, resulting in $z_n / \|z_n\|_2$ to be appended to the embedding at time n . Using a very similar construction, layer 2 computes $z_{n-1} / \|z_{n-1}\|_2$, which is added to the embedding at time n . Finally, in the last layer, the dot-product $\langle z_{i-1} / \|z_{i-1}\|_2, z_n / \|z_n\|_2 \rangle$ defines the attention score, and as the temperature κ grows, the pattern converges to $\text{Unif}(\mathcal{I}_n)$. Choosing the value vectors in this layer appropriately gives eq. (3).

6 Lower bounds on transformer size

In this section, we study the limits of how shallow a transformer can be made while still capturing conditional k -grams. The first result we establish in this vein is a lower bound against 1-layer transformers showing that their expressive power is too limited unless the embedding dimension or number of heads scale near-linearly in T .

Theorem 5. *Consider any 1-layer transformer with layer normalization and feedforward layers, where all the coordinates of the embedding vectors and unnormalized attention scores are computed with p bits of precision. If the transformer is able to compute the conditional 3-gram on inputs drawn from $\{0, 1, 2\}^T$ to within an additive error of $1/3T$, then $2pH + dp + 2 \geq T/3$.*

Choosing the bit precision to be $p = O(\log(T))$, this implies that for transformers with 1 layer, the sum of the number of heads and the embedding dimension must be at least $\Omega(T/\log(T))$, in order to represent conditional 3-grams to within an additive error of $1/3T$.

6.1 Conditional lower bounds on attention-only transformers

While the previous section shows that 1-layer transformers have fairly limited representation power, it is not immediately clear how whether any of these issues are present with transformers with more layers. Indeed, as we discussed in Section 4, an attention-only transformer with $O(\log_2(k))$ layers and 1 head per layer can represent conditional k -grams on its input sequences. With the addition of non-linearities, Theorem 4 shows that the model can represent conditional k -grams using just a constant number of layers. In this section, we try to understand the gap between these two results and prove conditional lower bounds on the size of attention-only transformers which do not have non-linearities arising from layer normalization.

We prove conditional lower bounds under some natural assumptions on the nature of the attention patterns learnt by the transformer. To motivate these assumptions, consider the experiment in Figure 6, where we train an attention-only transformer with 2 layers and 1 head, on order-1 Markov processes. At test-time, we plot the attention patterns learnt in the first layer of the model on test sequences. Notice that the attention pattern learnt by the model at layer 1 is largely independent of the input sequences themselves and only depends on the position.

Assumption 1. *In an L -layer attention-only transformer with H heads per layer, assume that layers $\ell = 1, 2, \dots, L-1$ and heads $h \in [H]$ realize an attention pattern where $\text{attn}_{n,i}^{(\ell,h)}$ only depends on the positions n and i and on ℓ and h , but not on the input sequence x_1, \dots, x_T .*

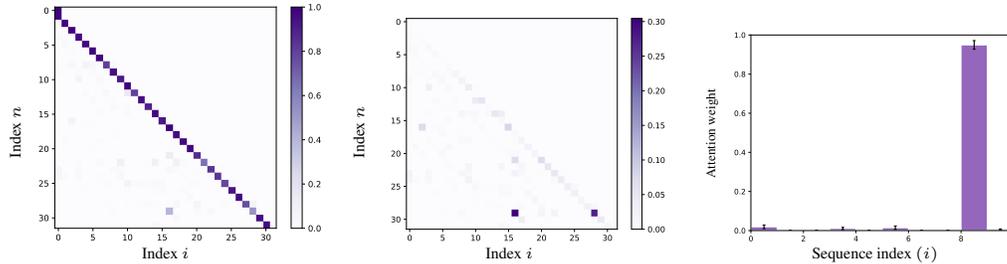


Figure 6: Attention matrix of the first attention layer, for a 2-layer 1-head transformer model trained on an order-1 Markov process, averaged across 100 input sequences of length 128. (a) and (b) plot the mean and standard deviation of the first 32 rows and columns of the attention matrix, while (c) zooms in on the column $n = 10$ and plots the mean attention for this column. (a) and (c) show that for almost all indices n , the attention layer focuses only on the previous symbol x_{n-1} . (b) shows that the attention pattern does not vary much with the input sequence considered, thereby providing evidence toward Assumption 1. More discussion in Appendix G.

Rather than proving the size lower bound depending on the transformers ability to represent the conditional k -gram itself, we consider a simplification and assume that the goal of the model is to represent a k^{th} -order induction head (Definition 2) in the last layer. Although learning a k^{th} -order induction head is not strictly necessary for the transformer to be able to represent conditional k -grams, note that every construction we have considered so far (cf. Theorems 1 to 4) go through this mechanism to realize the conditional k -gram model. Likewise, for other related problems, such as the causal learning task in [8], the causal structure is captured by an extension of the k^{th} -order induction head to general causal graphs. Our main lower bound is the following result.

Theorem 6. Consider an L -layer transformer with h_ℓ heads in layer L . Assuming the transformer satisfies Assumption 1, if $\prod_{\ell=1}^{L-1} (H_\ell + 1) \leq k - 2$, the attention pattern in layer L cannot represent a k^{th} -order induction head.

While this lower bound is not unconditional, meaning that it does not directly imply that the transformer cannot represent conditional k -grams, it is important to understand the interpretation of this result: attention-only transformers which somehow break through this barrier need to use a significantly different mechanism to realize the conditional k -gram model.

Theorem 6 implies that under Assumption 1, a 2-layer attention-only transformer with 1 head cannot realize a k^{th} -order induction head for any $k \geq 4$. Likewise, under the same assumption, a 3-layer attention-only transformer with 1 head cannot realize a k^{th} -order induction head for any $k \geq 6$. These results give more weight to the experiment in Figure 3 where we observe that a 2-layer transformer learns a k^{th} -order Markov process for $k = 4$ and a 3-layer transformer learns a k^{th} -order Markov process for $k = 8$, and show that non-linearities in the architecture allow the transformer to break past the size barriers in Theorem 4.

7 Conclusion

We observe empirically that 2 and 3 layer transformers are able to learn k^{th} -order Markov chains for much higher values of k than previously anticipated. We show there are $O(\log(k))$ -layer constructions of attention-only transformers which are able to learn the conditional k -gram model, which is the in-context MLE of the Markov model. With non-linearities in the model, we show that a 3-layer 1-head transformer is capable of representing the same. We show that 1-layer transformers cannot represent conditional k -grams for any $k \geq 3$ unless the number of heads or embedding dimension scale almost linearly in T . We also prove a conditional lower bound on the depth and number of heads of attention-only transformers to represent k^{th} -order induction heads, under an assumption on the realized attention patterns.

Acknowledgments and Disclosure of Funding

The work was partially supported by NSF Grant CCF-2211209 and Swiss NSF Grant 200364.

References

- [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information Processing Systems, pages 5998–6008, 2017.
- [2] 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.
- [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021.
- [4] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners, 2021.
- [5] Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint. In Thirty-seventh Conference on Neural Information Processing Systems, 2023.
- [6] Benjamin L. Edelman, Ezra Edelman, Surbhi Goel, Eran Malach, and Nikolaos Tsilivis. The evolution of statistical induction heads: In-context learning markov chains, 2024.
- [7] Ashok Vardhan Makkuva, Marco Bondaschi, Adway Girish, Alliot Nagle, Martin Jaggi, Hyeji Kim, and Michael Gastpar. Attention with markov: A framework for principled analysis of transformers via markov chains. arXiv preprint arXiv:2402.04161, 2024.
- [8] Eshaan Nichani, Alex Damian, and Jason D. Lee. How transformers learn causal structure with gradient descent, 2024.
- [9] Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. In International Conference on Machine Learning, pages 11080–11090, 2021.
- [10] Angeliki Giannou, Shashank Rajput, Jy-Yong Sohn, Kangwook Lee, Jason D. Lee, and Dimitris Papailiopoulos. Looped transformers as programmable computers. In Proceedings of the 40th International Conference on Machine Learning, pages 11398–11442, 23–29 Jul 2023.
- [11] Samet Oymak, Ankit Singh Rawat, Mahdi Soltanolkotabi, and Christos Thrampoulidis. On the role of attention in prompt-tuning. In Proceedings of the 40th International Conference on Machine Learning, 2023.
- [12] Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen. A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity. In The Eleventh International Conference on Learning Representations, 2023.
- [13] Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. Are transformers universal approximators of sequence-to-sequence functions? In International Conference on Learning Representations, 2020.
- [14] Jorge Pérez, Pablo Barceló, and Javier Marinkovic. Attention is Turing-complete. Journal of Machine Learning Research, 22(75):1–35, 2021.
- [15] Colin Wei, Yining Chen, and Tengyu Ma. Statistically meaningful approximation: a case study on approximating Turing machines with transformers. In Advances in Neural Information Processing Systems, volume 35, pages 12071–12083, 2022.
- [16] Bingbin Liu, Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers learn shortcuts to automata, 2023.
- [17] Satwik Bhattamishra, Kabir Ahuja, and Navin Goyal. On the ability and limitations of transformers to recognize formal languages, 2020.

- [18] Clayton Sanford, Daniel Hsu, and Matus Telgarsky. Transformers, parallel computation, and logarithmic depth, 2024.
- [19] Clayton Sanford, Daniel J Hsu, and Matus Telgarsky. Representational strengths and limitations of transformers. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 36677–36707. Curran Associates, Inc., 2023.
- [20] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. In-context learning and induction heads, 2022.
- [21] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [22] Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection, 2023.
- [23] Ziqian Lin and Kangwook Lee. Dual operating modes of in-context learning, 2024.
- [24] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models, 2023.
- [25] Jesse Hoogland, George Wang, Matthew Farrugia-Roberts, Liam Carroll, Susan Wei, and Daniel Mufet. The developmental landscape of in-context learning. *arXiv preprint arXiv:2402.02364*, 2024.
- [26] Nived Rajaraman, Jiantao Jiao, and Kannan Ramchandran. Toward a theory of tokenization in llms, 2024.
- [27] J. R. Norris. *Markov Chains*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1997.
- [28] Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. Localizing model behavior with path patching. *arXiv preprint arXiv:2304.05969*, 2023.
- [29] Andrew Chi-Chih Yao. Some complexity questions related to distributive computing (preliminary report). In *Proceedings of the eleventh annual ACM symposium on Theory of computing*, pages 209–213, 1979.
- [30] Matteo Pagliardini. GPT-2 modular codebase implementation. <https://github.com/epfml/llm-baselines>, 2023.

Appendix

Table of Contents

A Proof of Theorem 1	13
A.1 Extension to k -heads: Proof of Theorem 2	14
B Proof of Theorem 3	15
C Proof of Theorem 4	18
C.1 Modifying the definition of layer normalization	18
C.2 Notation and supplementary lemmas	18
C.3 Proof of Theorem 4	19
D Representation lower bounds for 1-layer transformers: Proof of Theorem 5	22
E Lower bounds on representing k^{th}-order induction heads: Proof of Theorem 6	23
E.1 Lower bounds on 2-layer 1-head attention-only transformers	23
E.2 L -layer attention-only transformers with 1 head per layer: Proof of Corollary 1	26
E.3 The general case: Transformers with H_ℓ heads in layer ℓ : Proof of Theorem 6	27
F Model architecture and hyper-parameters	28
G Additional experimental results	28

Notation. The notation $e_i^{d'} \in \mathbb{R}^{d'}$ refers to the one-hot encoding of i in d' dimensions. In other words it is the i^{th} standard basis vector in d' dimensions. The notation $\text{Blkdiag}\{A_1, A_1, \dots, A_m\}$ refers to the block diagonal matrix with i^{th} block as A_i .

A Proof of Theorem 1

We will first prove Theorem 1. In the first layer, choose the embeddings as,

$$\mathbf{x}_n^{(1)} = \text{Emb}(x_n) = \kappa \begin{bmatrix} \mathbf{1}_{1 \times 2} & e_{x_n}^S & \mathbf{0}_{1 \times 2S} \end{bmatrix}^T \in \mathbb{R}^d. \quad (4)$$

for a constant $\kappa > 0$ to be chosen later and $d = 2S + 2$. The relative position encodings will essentially be supported on the first two coordinates, the middle S coordinates are a one-hot encoding of the symbol x_n and the last $2S$ coordinates are 0. The relative position encodings in the first layer are chosen to be $\mathbf{p}_{n-i}^{(1),K} = \kappa(-1 + \mathbb{I}(n-i=1))e_1^d \in \mathbb{R}^d$ and $\mathbf{p}_{n-i}^{(1),V} = \mathbf{0} \in \mathbb{R}^d$. Choose $\mathbf{W}_K^{(1)}$ and $\mathbf{W}_Q^{(1)}$ to be $e_1^d(e_1^d)^T \in \mathbb{R}^{d \times d}$. With this choice,

$$\left\langle \mathbf{W}_K^{(1)}(\mathbf{x}_i^{(1)} + \mathbf{p}_{n-i}^{(1),K}), \mathbf{W}_Q^{(1)}\mathbf{x}_n^{(1)} \right\rangle = \kappa \mathbb{I}(n-i=1) \quad (5)$$

As $\kappa \rightarrow \infty$, the attention pattern (which takes the softmax over of these inner products over $i \in [n]$) computes,

$$\text{att}_{n,i}^{(1)} = \mathbb{I}(i = n-1) \quad (6)$$

for any $n > 1$. Choose the value matrix as,

$$\mathbf{W}_V^{(1)} = \begin{bmatrix} \mathbf{0}_{(2+S) \times 2} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & I_{S \times S} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{d \times d} \quad (7)$$

And with this choice and the residual connection, we get,

$$\mathbf{x}_n^{(2)} = \kappa \begin{bmatrix} \mathbf{1}_{1 \times 2} & e_{x_n}^S & e_{x_{n-1}}^S & \mathbf{0} \end{bmatrix} \in \mathbb{R}^d \quad (8)$$

which serves as the input to the 2nd transformer layer.

Layer 2. In layer 2, the relative position encodings $\mathbf{p}_{n-i}^{K,(2)}$ and $\mathbf{p}_{n-i}^{V,(2)}$ are all set as 0. The key matrix picks out the $e_{x_n}^S$ block out of $\mathbf{x}_n^{(2)}$ and the query vector picks out the $e_{x_{i-1}}^S$ block out of $\mathbf{x}_{i-1}^{(2)}$. In particular, these matrices are chosen so that,

$$\begin{aligned} \mathbf{W}_K^{(2)} \mathbf{x}_i^{(2)} &= \kappa \begin{bmatrix} \mathbf{1}_{1 \times 2} & e_{x_{i-1}}^S & \mathbf{0} \end{bmatrix}^T \in \mathbb{R}^d, \\ \mathbf{W}_Q^{(2)} \mathbf{x}_n^{(2)} &= \kappa \begin{bmatrix} \mathbf{1}_{1 \times 2} & e_{x_n}^S & \mathbf{0} \end{bmatrix}^T \in \mathbb{R}^d \end{aligned} \quad (9)$$

Taking the inner product of these vectors, and taking $\kappa \rightarrow \infty$, observe that the attention pattern concentrates on the uniform distribution over all coordinates i such that $x_{i-1} = x_n$. More formally, the attention pattern for any $n > 1$ is,

$$\text{att}_{n,i}^{(2)} = \frac{\mathbb{I}(x_{i-1} = x_n)}{\sum_{i=2}^n \mathbb{I}(x_{i-1} = x_n)}, \quad (10)$$

assuming $\sum_{i=2}^n \mathbb{I}(x_{i-1} = x_n) > 0$. Having realized this attention pattern, may choose the value and subsequent linear layer appropriately. The value matrix simply picks out the $e_{x_i}^S$ block from $\mathbf{x}_i^{(2)}$ and places it into the last S coordinates of $\mathbf{x}_i^{(3)}$, and the linear layer simply extracts this block and outputs it (after scaling down by a factor of κ), realizing the logits,

$$\text{logit}_n = \frac{1}{\sum_{i=2}^n \mathbb{I}(x_{i-1} = x_n)} \sum_{i=2}^n \mathbb{I}(x_{i-1} = x_n) \cdot e_{x_i}^S. \quad (11)$$

if $\sum_{i=2}^n \mathbb{I}(x_{i-1} = x_n) > 0$. In particular, under the same condition,

$$\text{logit}_T(x_{T+1}) = \frac{\sum_{n=2}^T \mathbb{I}(x_n = x_{T+1}, x_{n-1} = x_T)}{\sum_{i=2}^n \mathbb{I}(x_{n-1} = x_T)} \quad (12)$$

assuming $\sum_{i=2}^n \mathbb{I}(x_{n-1} = x_T)$, which is the conditional 1-gram model.

A.1 Extension to k -heads: Proof of Theorem 2

In the first layer, the embeddings are chosen to be,

$$\mathbf{x}_n^{(1)} = \text{Emb}(x_n) = \kappa \left[\mathbf{0}_{1 \times k} \mid 1 \mid e_{x_n}^S \mid \mathbf{0}_{1 \times (k+1)S} \right]^T \in \mathbb{R}^d \quad (13)$$

With $d = (k+1)(S+1) + S$. The relative position encodings are chosen as $\mathbf{p}_i^{K,(1)} = [e_i^k \ \mathbf{0}]^T$ for $1 \leq i \leq k$ and $\mathbf{p}_i^{K,(1)} = \mathbf{0}$ otherwise. Similarly, $\mathbf{p}_i^{V,(1)} = \mathbf{0}$ for every i . The h^{th} head has key and query matrices,

$$\begin{aligned} \mathbf{W}_Q^{(1,h)} &= \begin{bmatrix} \mathbf{0}_{1 \times k} & 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \\ \mathbf{W}_K^{(1,h)} &= \begin{bmatrix} \mathbf{0}_{1 \times (h-1)} & 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \end{aligned} \quad (14)$$

With these choices, and letting $\kappa \rightarrow \infty$, the h^{th} layer computes the attention pattern,

$$\text{att}_{n,i}^{(1,h)} = \mathbb{I}(i = n - h). \quad (15)$$

Choose the corresponding value matrix as,

$$\mathbf{W}_V^{(1,h)} = \begin{bmatrix} \mathbf{0}_{(2+hS) \times 2} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & I_{S \times S} & \mathbf{0} \end{bmatrix} \quad (16)$$

choosing the projection matrix appropriately, the output of the transformer after the first residual connection is,

$$\mathbf{x}_n^{(2)} = \kappa \left[\mathbf{0}_{1 \times k} \mid 1 \mid e_{x_n}^S \mid \cdots \mid e_{x_{n-k}}^S \right]^T. \quad (17)$$

Layer 2. In this layer, the relative position encodings $\mathbf{p}_{n-i}^{K,(2)}$ and $\mathbf{p}_{n-i}^{V,(2)}$ are all set as 0. The key and query matrices are chosen as,

$$\begin{aligned} \mathbf{W}_Q^{(2)} &= \begin{bmatrix} \mathbf{0}_{Sk \times k} & I_{(Sk+1) \times (Sk+1)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \\ \mathbf{W}_K^{(2)} &= \begin{bmatrix} \mathbf{0}_{Sk \times (k+S)} & I_{(Sk+1) \times (Sk+1)} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}. \end{aligned} \quad (18)$$

With this choices, we have that,

$$\left\langle \mathbf{W}_K^{(2)} \mathbf{x}_i^{(2)}, \mathbf{W}_Q^{(2)} \mathbf{x}_n^{(2)} \right\rangle = \kappa \sum_{j=1}^k \mathbb{I}(x_{i-j} = x_{n-j+1}). \quad (19)$$

Taking $\kappa \rightarrow \infty$, observe that the attention pattern concentrates on the uniform distribution over all coordinates i such that $x_{i-j} = x_{n-j+1}$ for all $j \in [k]$. More formally, if $\sum_{i=2}^n \mathbb{I}(x_{i-1} = x_n) > 0$, the attention pattern for any $n > 1$ is,

$$\text{att}_{n,i}^{(2)} = \frac{\mathbb{I}(\forall j \in [k], x_{i-j} = x_{n-j+1})}{\sum_{i=k+1}^n \mathbb{I}(\forall j \in [k], x_{i-j} = x_{n-j+1})}. \quad (20)$$

The value matrix picks out $e_{x_i}^S$ from the embedding $\mathbf{x}_i^{(2)}$ (Equation (17)) and places it in the last S coordinates. The subsequent linear layer picks out the last S coordinates, resulting in the logits,

$$\text{logit}_n = \sum_{i=k+1}^n \frac{\mathbb{I}(\forall j \in [k], x_{i-j} = x_{n-j+1})}{\sum_{i=k+1}^n \mathbb{I}(\forall j \in [k], x_{i-j} = x_{n-j+1})} e_{x_i}^S, \quad (21)$$

assuming that $\sum_{i=k+1}^n \mathbb{I}(\forall j \in [k], x_{i-j} = x_{n-j+1}) > 0$. In particular,

$$\text{logit}_T(x_{T+1}) = \frac{\sum_{n=k+1}^T \mathbb{I}(\forall 0 \leq j \leq k, x_{n-j} = x_{T-j+1})}{\sum_{n=k+1}^T \mathbb{I}(\forall 1 \leq j \leq k, x_{n-j} = x_{T-j+1})}, \quad (22)$$

assuming $\sum_{n=k+1}^T \mathbb{I}(\forall 1 \leq j \leq k, x_{n-j} = x_{T-j+1}) > 0$, i.e., the conditional k -gram model.

B Proof of Theorem 3

Define $k^* = 2^{\lceil \log_2(k+1) \rceil}$ by rounding $k+1$ up to the nearest power of 2 and $\ell^* = \log_2(k^*)$. In the setting of relative position encodings, given the sequence x_1, \dots, x_n , while generating the output of the attention + feedforward layer for the symbol x_n , the embeddings $\mathbf{x}_i = \text{Emb}(x_n) + \mathbf{p}_{n-i}$ are used for $i \in [n]$. In other words, the position encoding vector is taken relative to the end of the sequence, rather than the start of the sequence. Consider the embedding of x as,

$$\mathbf{x}_n^{(1)} = \text{Emb}(x_n) = \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid e_{x_n}^S \mid \mathbf{0}_{1 \times (k^*-1)S} \mid \mathbf{0}_{1 \times S} \right]^T \in \mathbb{R}^{(k^*+1)S + \ell^* + 1} \quad (23)$$

where $e_i^{d'} \in \mathbb{R}^{d'}$ is the standard basis vector in d' dimensions. And the relative position encoding for the keys as,

$$\mathbf{p}_i^{(1),K} = \begin{cases} \left[\left[\mathbf{1}_{1 \times \ell^*} \ \mathbf{0} \right]^T, & \text{if } i = 0, \\ \left[e_{1+\log_2(i)}^{\ell^*} \ \mathbf{0} \right]^T & \text{if } i \in \{1, 2, 4, \dots, k^*/2\} \\ \mathbf{0}_{d \times 1} & \text{otherwise.} \end{cases} \quad (24)$$

And for the value vectors, $\mathbf{p}_i^V = \mathbf{0}$ for all i .

For the first layer and first head, we will describe the value, key and query matrices. Choose,

$$\begin{aligned} \mathbf{W}_K^{(1)} &= \sqrt{\kappa} \begin{bmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \text{ and,} \\ \mathbf{W}_Q^{(1)} &= \sqrt{\kappa} \begin{bmatrix} \mathbf{0}_{1 \times \ell^*} & 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}. \end{aligned} \quad (25)$$

Then, observe that for $i \geq 1$,

$$\langle \mathbf{W}_K^{(1)}(\mathbf{x}_{n-i} + \mathbf{p}_i^{(1),K}), \mathbf{W}_Q^{(1)}\mathbf{x}_n \rangle = \kappa \mathbb{I}(i = 1)$$

and for $i = 0$,

$$\langle \mathbf{W}_K^{(1)}(\mathbf{x}_n + \mathbf{p}_0^{(1),K}), \mathbf{W}_Q^{(1)}\mathbf{x}_n \rangle = \kappa$$

In particular, letting $\kappa \rightarrow \infty$, the attention pattern is,

$$\text{att}_{n,n-i}^{(1)} = \frac{1}{2}\mathbb{I}(i = 0) + \frac{1}{2}\mathbb{I}(i = 1). \quad (26)$$

Choose the value matrix as,

$$\mathbf{W}_V^{(1)} = \begin{bmatrix} \mathbf{0}_{(\ell^*+S) \times \ell^*} & \mathbf{0} \\ \mathbf{0} & 2I \end{bmatrix}$$

together with the residual connection, we get,

$$\mathbf{x}_n^{(2)} = \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid e_{x_n}^S \mid e_{x_n}^S + e_{x_{n-1}}^S \mid \mathbf{0}_{1 \times (k^*-2)S} \mid \mathbf{0}_{1 \times S} \right]^T \quad (27)$$

Layer $\ell + 1$. By induction, assume that the output of the ℓ^{th} transformer layer is of the form,

$$\mathbf{x}_n^{(\ell+1)} = \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid \mathbf{v}_n \mid \mathbf{0}_{1 \times (k^*-2\ell)S} \mid \mathbf{0}_{1 \times S} \right]^T \quad (28)$$

for some vector $\mathbf{v}_n \in \mathbb{R}^{2\ell S}$. We will show that with appropriately chosen key, query and value vectors in the $(\ell + 1)^{\text{th}}$ layer, the output of this layer is,

$$\mathbf{x}_n^{(\ell+2)} = \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid \mathbf{v}_n \mid \mathbf{v}_n + \mathbf{v}_{n-2\ell} \mid \mathbf{0}_{1 \times (k^*-2\ell+1)S} \mid \mathbf{0}_{1 \times S} \right]^T \quad (29)$$

We will consider the same relative position encodings and query matrix in this layer as in the first layer (Equations (24) and (25)). Consider a key matrix of the form,

$$\mathbf{W}_K^{(\ell+1)} = \begin{bmatrix} \mathbf{0}_{1 \times \ell} & \sqrt{\kappa} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}$$

With this choice, observe that for $i \geq 1$,

$$\langle \mathbf{W}_K^{(\ell+1)}(\mathbf{x}_{n-i}^{(\ell+1)} + \mathbf{p}_i^{(\ell+1),K}), \mathbf{W}_Q^{(\ell+1)}\mathbf{x}_n^{(\ell+1)} \rangle = \kappa \cdot \mathbb{I}(i = 2^\ell)$$

and for $i = 0$,

$$\langle \mathbf{W}_K^{(\ell+1)}(\mathbf{x}_n^{(\ell+1)} + \mathbf{p}_0^{(\ell+1),K}), \mathbf{W}_Q^{(\ell+1)}\mathbf{x}_n^{(\ell+1)} \rangle = \kappa$$

In particular, letting $\kappa \rightarrow \infty$, the attention pattern is,

$$\text{att}_{n,n-i}^{(\ell+1)} = \frac{1}{2}\mathbb{I}(i = 0) + \frac{1}{2}\mathbb{I}(i = 2^\ell). \quad (30)$$

Choosing the value matrix as,

$$\mathbf{W}_V^{(\ell+1)} = \begin{bmatrix} \mathbf{0}_{(\ell^*+2^\ell S) \times \ell^*} & \mathbf{0} \\ \mathbf{0} & 2I \end{bmatrix},$$

we get,

$$\mathbf{x}_n^{(\ell+2)} = \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid \mathbf{v}_n \mid \mathbf{v}_n + \mathbf{v}_{n-2^\ell} \mid \mathbf{0}_{1 \times (k^*-2^\ell+1)S} \mid \mathbf{0}_{1 \times S} \right]^T \quad (31)$$

Final last transformer layer ($\ell = \ell^*$). The output of the second last transformer layer, indexed $\ell^* - 1$ is,

$$\begin{aligned} \mathbf{z}_n^{(\ell^*)} \triangleq \mathbf{x}_n^{(\ell^*)} &= \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid \mathbf{v}_n^{(\ell^*-1)} \mid \mathbf{v}_n^{(\ell^*-1)} + \mathbf{v}_{n-2\ell^*-1}^{(\ell^*-1)} \mid \mathbf{0}_{1 \times S} \right]^T \\ &= \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid \mathbf{v}_n^{(\ell^*-1)} \mid \mathbf{v}_n^{(\ell^*-1)} + \mathbf{v}_{n-\frac{k^*}{2}}^{(\ell^*-1)} \mid \mathbf{0}_{1 \times S} \right]^T, \end{aligned}$$

which follows by plugging in the definition of k^* . Note that there exists a linear transformation $\mathbf{L}^{(\ell^*)}$ such that,

$$\mathbf{z}_n^{(\ell^*-1)} \triangleq \mathbf{L}^{(\ell^*)} \mathbf{x}_n^{(\ell^*)} = \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid \mathbf{v}_n^{(\ell^*-1)} \mid \mathbf{v}_{n-\frac{k^*}{2}}^{(\ell^*-1)} \mid \mathbf{0}_{1 \times S} \right]^T$$

This can be further decomposed as,

$$\begin{aligned} \mathbf{z}_n^{(\ell^*-1)} &= \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid \mathbf{v}_n^{(\ell^*-2)} \mid \mathbf{v}_n^{(\ell^*-2)} + \mathbf{v}_{n-2\ell^*-2}^{(\ell^*-2)} \mid \mathbf{v}_{n-\frac{k^*}{2}}^{(\ell^*-2)} \mid \mathbf{v}_{n-\frac{k^*}{2}}^{(\ell^*-2)} + \mathbf{v}_{n-\frac{k^*}{2}-2\ell^*-2}^{(\ell^*-2)} \mid \mathbf{0}_{1 \times S} \right]^T \end{aligned}$$

And yet again there exists a linear transformation $\mathbf{L}^{(\ell^*-1)}$ which transforms this as,

$$\begin{aligned} \mathbf{z}_n^{(\ell^*-2)} \triangleq \mathbf{L}^{(\ell^*-1)} \mathbf{z}_n^{(\ell^*-1)} &= \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid \mathbf{v}_n^{(\ell^*-2)} \mid \mathbf{v}_{n-2\ell^*-2}^{(\ell^*-2)} \mid \mathbf{v}_{n-\frac{k^*}{2}-2\ell^*-2}^{(\ell^*-2)} \mid \mathbf{v}_{n-\frac{k^*}{2}-2\ell^*-2}^{(\ell^*-2)} \mid \mathbf{0}_{1 \times S} \right]^T \\ &= \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid \mathbf{v}_n^{(\ell^*-2)} \mid \mathbf{v}_{n-\frac{k^*}{4}}^{(\ell^*-2)} \mid \mathbf{v}_{n-\frac{k^*}{2}}^{(\ell^*-2)} \mid \mathbf{v}_{n-\frac{3k^*}{4}}^{(\ell^*-2)} \mid \mathbf{0}_{1 \times S} \right]^T \end{aligned} \quad (32)$$

By recursing this argument and composing all the linear transformations, up to a global permutation, we get that,

$$\begin{aligned} \prod_{\ell=1}^{\ell^*} \mathbf{L}^{(\ell)} \mathbf{x}_n^{(\ell^*)} &= \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid \mathbf{v}_n^{(1)} \mid \mathbf{v}_{n-1}^{(1)} \mid \cdots \mid \mathbf{v}_{n-(k^*-1)}^{(1)} \mid \mathbf{0}_{1 \times S} \right]^T \\ &= \left[\mathbf{0}_{1 \times \ell^*} \mid 1 \mid e_{x_n}^S \mid \cdots \mid e_{x_{n-(k^*-1)}}^S \mid \mathbf{0}_{1 \times S} \right]^T \end{aligned} \quad (33)$$

In the final layer, we will right multiply the key, query and value matrices by $\mathbf{L}^* = \prod_{\ell=1}^{\ell^*} \mathbf{L}^{(\ell)}$. The effect can be interpreted as operating the original key, query and value matrices on the embedding vectors in Equation (33). In the final layer, we will set all the position encodings to be $\mathbf{0}$ and consider the key and query matrices,

$$\begin{aligned} \mathbf{W}_K^{(\ell^*)} &= \sqrt{\kappa} \begin{bmatrix} \mathbf{0}_{S k \times (\ell^*+1+S)} & I_{S k \times S k} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \\ \mathbf{W}_Q^{(\ell^*)} &= \sqrt{\kappa} \begin{bmatrix} \mathbf{0}_{S k \times (\ell^*+1)} & I_{S k \times S k} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \end{aligned} \quad (34)$$

Then,

$$\left\langle \mathbf{W}_K^{(\ell^*)} \mathbf{L}^* \mathbf{x}_{n-i}^{(\ell^*)}, \mathbf{W}_Q^{(\ell^*)} \mathbf{L}^* \mathbf{x}_n^{(\ell^*)} \right\rangle = \kappa \sum_{j=0}^{k-1} \mathbb{I}(x_{n-j} = x_{i-1-j}) \quad (35)$$

Where we must be careful to note that the input $\mathbf{x}_n^{(\ell^*)}$ contains copies of $e_{x_n}, e_{x_{n-1}}, \dots, e_{x_{n-k}}$ since $k^* \geq k + 1$ by definition.

Letting $\kappa \rightarrow \infty$, if there exists i such that $\sum_{j=0}^{k-1} \mathbb{I}(x_{n-j} = x_{i-1-j}) > 0$, for $n \geq k$, the attention pattern is,

$$\text{att}_{n,i}^{(\ell^*)} = \frac{\mathbb{I}(x_{i-1} = x_n, x_{i-2} = x_{n-1}, \dots, x_{i-k} = x_{n-k+1})}{\sum_{i=k}^n \mathbb{I}(x_{i-1} = x_n, x_{i-2} = x_{n-1}, \dots, x_{i-k} = x_{n-k+1})} \quad (36)$$

Finally, choose,

$$\mathbf{W}_V^{(\ell^*+2)} = \begin{bmatrix} \mathbf{0}_{(d-S) \times (\ell^*+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{S \times (\ell^*+1)} & I_{S \times S} & \mathbf{0} \end{bmatrix}, \quad (37)$$

we get,

$$\mathbf{x}_n^{(\ell^*+1)} + \sum_{i=k}^n \frac{\mathbb{I}(x_{i-1} = x_n, x_{i-2} = x_{n-1}, \dots, x_{i-k} = x_{n-k+1})}{\sum_{i=k}^n \mathbb{I}(x_{i-1} = x_n, x_{i-2} = x_{n-1}, \dots, x_{i-k} = x_{n-k+1})} \begin{bmatrix} \mathbf{0}_{(d-S) \times 1} \\ e_{x_i} \end{bmatrix} \quad (38)$$

Choosing the subsequent linear layer as,

$$\mathbf{A} = \begin{bmatrix} \mathbf{0}_{S \times (d-S)} & I_{S \times S} \end{bmatrix} \quad (39)$$

$$\mathbf{b} = \mathbf{0}_{S \times 1} \quad (40)$$

Results in the output,

$$\text{logit}_T(x_{T+1}) = \sum_{n=k}^T \frac{\mathbb{I}(x_n = x_{T+1}, x_{n-1} = x_T, x_{n-2} = x_{T-1}, \dots, x_{n-k} = x_{T-k+1})}{\sum_{n=k}^T \mathbb{I}(x_{n-1} = x_T, x_{n-2} = x_{T-1}, \dots, x_{n-k} = x_{T-k+1})} \quad (41)$$

which is precisely the in-context conditional k -gram.

C Proof of Theorem 4

C.1 Modifying the definition of layer normalization

In every layer, we will perform a simple transformation which is to double the hidden dimension d and add a copy of $-\mathbf{x}_n^{(\ell)}$ into the last d coordinates. This is possible by modifying the weights of the transformer appropriately as discussed below. A consequence of this transformation is that the feature mean of the \mathbf{x}_n 's is $\mu_n = 0$, and therefore the standard deviation σ_n simply normalizes by the L_2 -norm of the features. In order to avoid having to explicitly state this transformation at each layer, we will simply redefine the layer norm LN to output $\mathbf{v}/\|\mathbf{v}\|_2$ for the input vector \mathbf{v} , which is realized on the first d coordinates of the transformed embeddings.

This transformation can be realized automatically by redefining the initial embeddings $\text{Emb}(x_n)$, and modifying the weights of the attention and feedforward subnetworks as follows: The input embeddings are changed to $[\text{Emb}(x_n) \quad -\text{Emb}(x_n)]^T \in \mathbb{R}^{2d}$. The key and query matrices are chosen to be 0 on the last d coordinates in every layer; the value matrix for $i \geq 1$ is transformed to $\text{Blkdiag}(\{\mathbf{W}_V^{(\ell)}, \mathbf{W}_V^{(\ell)}\})$, and likewise changing the feedforward layer to the block diagonal matrices $\text{Blkdiag}(\{\mathbf{W}_1^{(\ell)}, \mathbf{W}_1^{(\ell)}\})$ and $\text{Blkdiag}(\{\mathbf{W}_2^{(\ell)}, \mathbf{W}_2^{(\ell)}\})$. This transformation adds a copy of $-\mathbf{x}_n^{(\ell)}$ into the last d coordinates of the corresponding embeddings.

C.2 Notation and supplementary lemmas

For each $i \in [T]$, define,

$$\mathbf{v}_i = e_{x_{i-1}} + 3 \cdot e_{x_{i-2}} + \dots + 3^{k-1} \cdot e_{x_{i-k}} \quad (42)$$

$$\mathbf{u}_i = e_{x_i} + 3 \cdot e_{x_{i-1}} + \dots + 3^{k-1} \cdot e_{x_{i-k+1}} \quad (43)$$

Note that although $\mathbf{v}_i = \mathbf{u}_{i-1}$, we make the distinction between the two to avoid any confusion in what is stored in the embedding vector at time i and at time $i-1$. Furthermore, define,

$$\mathcal{I}_n = \{k+1 \leq i \leq n : \forall j \in [k], x_{i-j} = x_{n-j+1}\}. \quad (44)$$

Lemma 1. *If $i \in \mathcal{I}_n$, $\mathbf{z}_i = \mathbf{z}_{n-1}$. However, if $i \geq k+1$ but $i \notin \mathcal{I}_n$, then, $\|\frac{\mathbf{v}_i}{\|\mathbf{v}_i\|_2} - \frac{\mathbf{u}_n}{\|\mathbf{u}_n\|_2}\|_2 \geq 3^{-k}$.*

Let $j^* \in \{0, 1, \dots, k-1\}$ denote the largest index j such that $x_{n-j} \neq x_{i-j-1}$. Consider the coordinates $a = x_{n-j^*} \in [S]$ and $b = x_{i-j^*-1} \in [S]$. Then,

$$\langle \mathbf{v}_n, e_a \rangle - \langle \mathbf{u}_i, e_a \rangle \geq 3^j - \sum_{j=0}^{j^*-1} 3^j = \frac{3^{j^*}}{2}, \quad (45)$$

$$\langle \mathbf{u}_i, e_b \rangle - \langle \mathbf{v}_n, e_b \rangle \geq \frac{3^{j^*}}{2} \quad (46)$$

If $\|\mathbf{v}_n\|_2 \geq \|\mathbf{u}_i\|_2$, then,

$$\left\langle \frac{\mathbf{u}_i}{\|\mathbf{u}_i\|_2}, e_b \right\rangle - \left\langle \frac{\mathbf{v}_n}{\|\mathbf{v}_n\|_2}, e_b \right\rangle \geq \frac{\langle \mathbf{u}_i, e_b \rangle - \langle \mathbf{v}_n, e_b \rangle}{\max\{\|\mathbf{u}_i\|_2, \|\mathbf{v}_n\|_2\}} \geq \frac{3^{j^*}}{2 \cdot \frac{3^k}{2}} = 3^{j^* - k} \quad (47)$$

This uses the fact that \mathbf{u}_i and \mathbf{v}_n are coordinate-wise non-negative. On the other hand, if $\|\mathbf{v}_n\|_2 \leq \|\mathbf{u}_i\|_2$, using a similar analysis,

$$\left\langle \frac{\mathbf{u}_i}{\|\mathbf{u}_i\|_2}, e_a \right\rangle - \left\langle \frac{\mathbf{v}_n}{\|\mathbf{v}_n\|_2}, e_a \right\rangle \geq 3^{j^* - k}. \quad (48)$$

In either case, there is a coordinate (a or b) such that, $\mathbf{u}_i/\|\mathbf{u}_i\|_2$ and $\mathbf{v}_n/\|\mathbf{v}_n\|_2$ differ by at least $3^{j^* - k}$. This implies the lower bound on the L_2 norm of the difference of the vectors.

C.3 Proof of Theorem 4

Choose the input embeddings as,

$$\mathbf{x}_n^{(1)} = \text{Emb}(x_n) = [\mathbf{0}_{1 \times 3} \quad e_x^S \quad \mathbf{0}_{1 \times 5S}]^T \in \mathbb{R}^{6S+3} \quad (49)$$

In the first two layers we will use the same relative position embeddings, in particular,

$$\mathbf{p}_i^{(1),K} = \mathbf{p}_i^{(2),K} = \begin{cases} \sqrt{\log(3)} \cdot [1 \quad \mathbf{0}]^T, & \text{if } i = 0, \\ (i+1)\sqrt{\log(3)} \cdot [0 \quad 1 \quad \mathbf{0}]^T, & \text{if } i \in \{1, 2, \dots, k-1\}, \\ (k+1)\sqrt{\log(3)} \cdot [0 \quad 0 \quad 1 \quad \mathbf{0}]^T, & \text{if } i = k. \end{cases} \quad (50)$$

and the value embeddings,

$$\mathbf{p}_i^{(1),V} = \mathbf{p}_i^{(2),V} = \begin{cases} 3^i [1 \quad \mathbf{0}]^T & \text{for } i \leq k \\ \mathbf{0} & i > k. \end{cases} \quad (51)$$

In the final layer, we will drop all position-related information and choose $\mathbf{p}_i^{(3),K} = \mathbf{p}_i^{(3),V} = \mathbf{0}$ for all i .

Layer 1. Consider the key and query matrices,

$$\begin{aligned} \mathbf{W}_K^{(1)} &= \sqrt{\kappa} \cdot \begin{bmatrix} \mathbf{1}_{1 \times 2} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \\ \mathbf{W}_Q^{(1)} &= \sqrt{\kappa} \cdot \begin{bmatrix} \mathbf{0}_{1 \times 3} & \mathbf{1}_{1 \times S} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \end{aligned} \quad (52)$$

Then, observe that,

$$\left\langle \mathbf{W}_K^{(1)} (\text{Emb}(x_{n-i}) + \mathbf{p}_i^{(1),K}), \mathbf{W}_Q^{(1)} \text{Emb}(x_n) \right\rangle = \kappa(i+1) \log(3) \cdot \mathbb{I}(0 \leq i \leq \min\{n, k\} - 1)$$

Letting $\kappa \rightarrow \infty$, this results in the attention pattern,

$$\text{att}_{n, n-i}^{(1)} = \frac{3^i \mathbb{I}(0 \leq i \leq \min\{n, k\} - 1)}{\sum_{i'=0}^{\min\{n, k\} - 1} 3^{i'}} \quad (53)$$

Choose the value matrix as,

$$\mathbf{W}_V^{(1)} = \begin{bmatrix} \mathbf{0}_{(S+3) \times 3} & \mathbf{0} \\ \mathbf{0} & I \end{bmatrix}$$

The output of the attention layer (with the residual connection) is,

$$\tilde{\mathbf{x}}_n^{(1)} = [\mathbf{0}_{1 \times 3} \mid e_{x_n}^S \mid \mathbf{u}_n \mid \mathbf{0}_{1 \times 3S}]^T, \text{ where, } \mathbf{u}_n = \sum_{i=0}^{\min\{n, k\} - 1} \text{att}_{n, n-i} e_{x_{n-i}}^S. \quad (54)$$

In the feedforward layer to follow, we will choose,

$$\begin{aligned} \mathbf{W}_1^{(1)} &= I \\ \mathbf{W}_2^{(1)} &= \begin{bmatrix} \mathbf{0}_{(3+2S) \times (3+S)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & I_{S \times S} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \end{aligned} \quad (55)$$

Which simply extracts \mathbf{u}_n from $\tilde{\mathbf{x}}_n^{(1)}$. With the subsequent layer norm and residual connection, the output of the first layer is,

$$\mathbf{x}_n^{(2)} = \left[\mathbf{0}_{1 \times 3} \mid e_{x_n}^S \mid \mathbf{u}_n \mid \frac{\mathbf{u}_n}{\|\mathbf{u}_n\|_2} \mid \mathbf{0}_{1 \times 3S} \right]^T \quad (56)$$

Layer 2. In this layer, the relative position encodings and query matrix are the same as in layer 1 but the key matrix is chosen as,

$$\mathbf{W}_K^{(2)} = \sqrt{\kappa} \begin{bmatrix} 0 & \mathbf{1}_{1 \times 2} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (57)$$

With this choice, observe that,

$$\left\langle \mathbf{W}_K^{(2)}(\mathbf{x}_{n-i}^{(2)} + \mathbf{p}_i^{(1,K)}), \mathbf{W}_Q^{(2)}\mathbf{x}_n^{(2)} \right\rangle = \kappa(i+1) \log(3) \cdot \mathbb{I}(1 \leq i \leq k) \quad (58)$$

As before, since $\kappa \rightarrow \infty$, this results in the attention pattern,

$$\text{att}_{n,n-i}^{(2)} = \frac{3^i \mathbb{I}(1 \leq i \leq \min\{k, n-1\})}{\sum_{i'=1}^{\min\{k, n-1\}} 3^{i'}} \quad (59)$$

which is similar, but subtly different from the attention pattern in the first layer (Equation (53)). The first layer focuses on indices $n-i$ such that $0 \leq i \leq k-1$, while this layer focuses on $1 \leq i \leq k$. Choosing the value and projection matrices as,

$$\mathbf{W}_V^{(2)} = \begin{bmatrix} I_{3 \times 3} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{3S \times 3} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & I_{S \times S} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (60)$$

The output of the attention layer (with the first residual connection) is,

$$\begin{aligned} \tilde{\mathbf{x}}_n^{(2)} &= \left[Z_n \mid \mathbf{0}_{1 \times 2} \mid e_{x_n}^S \mid \mathbf{u}_n \mid \frac{\mathbf{u}_n}{\|\mathbf{u}_n\|_2} \mid \mathbf{v}_n \mid \mathbf{0}_{1 \times 2S} \right]^T, \\ \text{where, } \mathbf{v}_n &= \sum_{i=1}^{\min\{k, n-1\}} \text{att}_{n,n-i} e_{x_{n-i}}^S, \\ \text{and, } Z_n &= \sum_{i=1}^{\min\{k, n-1\}} \text{att}_{n,n-i} 3^i, \end{aligned} \quad (61)$$

It is a short calculation to see that $Z_n = 3^{k+1}/5$ if $n \geq k+1$ and otherwise, $Z_n \leq 3^k/5$. This will be useful later, since the value of Z_n can be used to determine whether $n \geq k+1$ or $n \leq k$ which will allow the the next layer to avoid calculating the attention at $i \leq k$, where the evaluation $x_n = x_{i-1}, \dots, x_{n-k+1} = x_{i-k}$ is not well defined. In the subsequent FFN layer, we will choose,

$$\begin{aligned} \mathbf{W}_1^{(2)} &= I \\ \mathbf{W}_2^{(2)} &= \begin{bmatrix} \mathbf{0}_{(3+4S) \times (3+3S)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & I_{S \times S} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0}_{S \times 2S} \end{bmatrix} \end{aligned} \quad (62)$$

Which extracts \mathbf{v}_n from the embedding $\tilde{\mathbf{x}}_n^{(2)}$. With the layer norm and adding the final residual connection, the output of this layer is,

$$\mathbf{x}_n^{(3)} = \left[Z_n \mid \mathbf{0}_{2 \times 1} \mid e_{x_n}^S \mid \mathbf{u}_n \mid \frac{\mathbf{u}_n}{\|\mathbf{u}_n\|_2} \mid \mathbf{v}_n \mid \frac{\mathbf{v}_n}{\|\mathbf{v}_n\|_2} \mid \mathbf{0}_{S \times 1} \right]^T \quad (63)$$

Layer 3. In this layer, all the relative position encodings are set as $\mathbf{0}$ and instead,

$$\begin{aligned} \mathbf{W}_Q^{(3)} &= \sqrt{2\kappa} \begin{bmatrix} 1 & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{S \times (2+3S)} & I_{S \times S} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \\ \mathbf{W}_K^{(3)} &= \sqrt{2\kappa} \cdot \begin{bmatrix} 1 & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{S \times (2+4S)} & I_{S \times S} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \end{aligned} \quad (64)$$

With these choices,

$$\begin{aligned} \langle \mathbf{W}_K^{(3)} \mathbf{x}_i^{(3)}, \mathbf{W}_Q^{(3)} \mathbf{x}_n^{(3)} \rangle &= 2\kappa Z_i Z_n + \frac{2\kappa \langle \mathbf{v}_i, \mathbf{u}_n \rangle}{\|\mathbf{v}_i\|_2 \cdot \|\mathbf{u}_n\|_2} \\ &= 2\kappa Z_i Z_n + 2\kappa - \kappa \left\| \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|_2} - \frac{\mathbf{u}_n}{\|\mathbf{u}_n\|_2} \right\|_2^2 \end{aligned} \quad (65)$$

The resulting attention scores are,

$$\text{att}_{n,i}^{(3)} \propto \exp \left(-\kappa \left\| \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|_2} - \frac{\mathbf{u}_n}{\|\mathbf{u}_n\|_2} \right\|_2^2 + 2\kappa Z_i Z_n \right) \quad (66)$$

Recall that $\mathcal{I}_n = \{k+1 \leq i \leq n : \forall j \in [k], x_{n-j+1} = x_{i-j}\}$. Then for any $i \in \mathcal{I}_n$, $\mathbf{v}_i = \mathbf{u}_n$, and by Lemma 1, for any $i \geq k+1$ but not in \mathcal{I}_n ,

$$\left\| \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|_2} - \frac{\mathbf{u}_n}{\|\mathbf{u}_n\|_2} \right\|_2 \geq \frac{1}{3^k}.$$

Note that this gap is small but non-zero. Furthermore, recall that $Z_i = 3^{k+1}/5$ if $i \geq k$ and otherwise $Z_i \leq 3^k/5$. Thus the attention prefers values of i such that $\mathbf{v}_i = \mathbf{u}_n$ and such that $i \geq k+1$. In particular, as $\kappa \rightarrow \infty$, the resulting attention pattern is,

$$\text{att}_{n,\cdot}^{(3)} = \text{Unif}(\mathcal{I}_n). \quad (67)$$

Choosing,

$$\mathbf{W}_V^{(3)} = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{S \times 3} & I_{S \times S} & \mathbf{0} \end{bmatrix}.$$

We get that,

$$\tilde{\mathbf{x}}_n^{(3)} = \mathbf{x}_n^{(3)} + \sum_{i=1}^n \text{att}_{n,i}^{(3)} \begin{bmatrix} \mathbf{0} \\ e_{x_i}^S \end{bmatrix} = \mathbf{x}_n^{(3)} + \frac{1}{|\mathcal{I}_n|} \sum_{i \in \mathcal{I}_n} \begin{bmatrix} \mathbf{0} \\ e_{x_i}^S \end{bmatrix}.$$

The feedforward layer is chosen to have $\mathbf{W}_1^{(3)} = \mathbf{W}_2^{(3)} = \mathbf{0}$, and the overall output of the final transformer layer is therefore just $\tilde{\mathbf{x}}_n^{(3)}$. In the output linear layer, choose,

$$\begin{aligned} \mathbf{A} &= [\mathbf{0}_{S \times (d-S)} \quad I_{S \times S}] \\ \mathbf{b} &= \mathbf{0} \end{aligned} \quad (68)$$

which results in,

$$\text{logit}_n = \frac{1}{|\mathcal{I}_n|} \sum_{i \in \mathcal{I}_n} e_{x_i}^S = \sum_{i=k+1}^n \frac{\mathbb{I}(\forall 1 \leq j \leq k, x_{i-j} = x_{n-j+1})}{\sum_{i'=k+1}^n \mathbb{I}(\forall 1 \leq j \leq k, x_{i'-j} = x_{n-j+1})} \cdot e_{x_i}$$

In particular,

$$\text{logit}_T(x_{T+1}) = \frac{\sum_{n=k+1}^T \mathbb{I}(\forall 0 \leq i \leq k, x_{n-i} = x_{T-i+1})}{\sum_{n=k+1}^T \mathbb{I}(\forall 1 \leq i \leq k, x_{n-i} = x_{T-i+1})} \quad (69)$$

which is the conditional k -gram.

D Representation lower bounds for 1-layer transformers: Proof of Theorem 5

We prove this lower bound by a reduction to communication complexity, and specifically to the set disjointness problem.

Suppose Alice and Bob are given strings $\mathbf{a}, \mathbf{b} \in \{0, 1\}^n$ which are indicator vectors of sets A and B . Their goal is to jointly compute $\text{DIS}(\mathbf{a}, \mathbf{b}) = \mathbb{I}(\exists i : \mathbf{a}_i = \mathbf{b}_i = 1)$, which indicates whether A and B intersect or not. Alice and Bob may send a single bit message to the other party over a sequence of communication rounds. The following seminal result by [29] asserts a lower bound on amount of communication required between Alice and Bob to carry out this task.

Theorem 7 ([29]). *Any deterministic protocol for computing $\text{DIS}(\mathbf{a}, \mathbf{b})$ requires at least n rounds of communication.*

We show that a 1-layer transformer with sufficiently small embedding dimension / number of heads can be used to simulate a two-way communication protocol between Alice and Bob to solve $\text{DIS}(\mathbf{a}, \mathbf{b})$ in a way which contradicts Yao's lower bound in Theorem 7.

With $m = T/3 - 1$, suppose Alice and Bob have length m bit strings $\mathbf{a}, \mathbf{b} \in \{0, 1\}^m$. The transformer's input will be a sequence of the form,

$$2, \mathbf{a}_1, \mathbf{b}_1, 2, \mathbf{a}_2, \mathbf{b}_2, \dots, 2, \mathbf{a}_m, \mathbf{b}_m, 2, 1, \quad (70)$$

of length $3m + 2 = T - 1$. The input basically contains a repeating motif, composed of the symbol 2 followed by one of Alice's bits, and then one of Bob's bits. The last 2 symbols are 2 and 1. We will consider the empirical conditional 3-gram probability the transformer associates with the symbol $x_T = 2$. Noting that $x_{T-1} = 1$ and $x_{T-2} = 1$, the conditional 3-gram is computed to be,

$$\frac{\sum_{i=3}^{T-1} \mathbb{I}(x_i = 1, x_{i-1} = 1, x_{i-2} = 2)}{\sum_{i=3}^{T-1} \mathbb{I}(x_{i-1} = 1, x_{i-2} = 2)} \quad (71)$$

Note that if $x_{i-2} = 2$, then i must be of the form $3j$ for $j = 1, \dots, m$, and we may rewrite the sum as,

$$\frac{\sum_{j=1}^m \mathbb{I}(x_{3j} = 1, x_{3j-1} = 1)}{\sum_{j=1}^m \mathbb{I}(x_{3j-1} = 1)} = \frac{|A \cap B|}{|B|} \quad (72)$$

Now, let us use the transformer to construct a deterministic communication protocol between Alice and Bob. Alice is given $(x_2, x_5, \dots, x_{3m-1}) = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m)$ and Bob is given $(x_3, x_6, \dots, x_{3m}) = (\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_m)$.

In the first round, Alice computes the normalization in the softmax of the attention which comes from the set of inputs she holds. For simplifying notation define,

$$\text{score}^{(h)}(i) = \exp \left(\left\langle \mathbf{W}_K^{(h)} (\text{Emb}(x_i) + \mathbf{p}_i), \mathbf{W}_Q^{(h)} \text{Emb}(x_{T-1}) \right\rangle \right) \quad (73)$$

In particular, for each head $h \in [H]$, she computes,

$$Z_{\text{Alice}}^{(h)} = \log \left(\sum_{j=1}^m \text{score}^{(h)}(3j - 1) \right) \quad (74)$$

Assuming that the transformer uses p bits of precision, Alice communicates $Z_{\text{Alice}}^{(h)}$ for each h , which corresponds to pH bits of communication. With this information, Bob completes the rest of the normalization term (again up to p bits of precision) and computes,

$$Z^{(h)} = \log \left(Z_{\text{Alice}}^{(h)} + Z_{\text{Bob}}^{(h)} + Z_{\text{common}}^{(h)} \right), \quad (75)$$

$$\text{where } Z_{\text{Bob}}^{(h)} = \log \left(\sum_{j=1}^m \text{score}^{(h)}(3j) \right) \quad (76)$$

$$\text{and } Z_{\text{common}}^{(h)} = \log \left(\sum_{j=1}^m \text{score}^{(h)}(3j - 2) + \text{score}^{(h)}(T - 2) + \text{score}^{(h)}(T - 1) \right) \quad (77)$$

which is the overall normalization term in the softmax. This is communicated back to Alice, using another pH bits of communication. Next using this information, Alice computes the output of the

attention layer, taking the convex combination corresponding to the inputs she knows. In particular, for each $h \in [H]$ she computes,

$$\sum_{j=1}^n \frac{\text{score}^{(h)}(3j-1)}{\exp(Z^{(h)})} \text{Emb}(x_{3j-1}) \in \mathbb{R}^d. \quad (78)$$

across all the heads. Rather than transmitting everything, she concatenates the outputs of the heads, and multiplies them by the value and projection matrices to result in the output $\mathbf{y}_{\text{Alice}}$ which is d -dimensional. This is sent to Bob using dp bits of communication. Subsequently, Bob computes the terms in the attention corresponding to the inputs he knows as well as the public inputs (all the 2's at positions $3j-2$ as well as the last two symbols). In particular,

$$\sum_{j=1}^m \frac{\text{score}^{(h)}(3j)}{\exp(Z^{(h)})} \text{Emb}(x_{3j}) + \sum_{j=1}^{m+1} \frac{\text{score}^{(h)}(3j-2)}{\exp(Z^{(h)})} \text{Emb}(2) + \frac{\text{score}^{(h)}(T-1)}{\exp(Z^{(h)})} \text{Emb}(1) \quad (79)$$

These are yet again concatenated across all the heads and multiplied by the value and projection matrices to result in the output \mathbf{y}_{Bob} which is added to $\mathbf{y}_{\text{Alice}}$ to result in \mathbf{y} . Bob passes \mathbf{y} through the residual connection, layer norm, and feedforward layers, and subsequently through the linear layer and softmax of the model to result in the output of the model. By assumption, the output of the model approximately captures the conditional 3-gram, which by Equation (72) equals $|A \cap B|/|B|$. Note that if $|A \cap B|/|B|$ is non-zero, it must be at least $1/T$. This means, if the transformer is able to compute the conditional 3-gram to within an additive error of $1/3T$, then Bob can simply threshold the output of the transformer to decide whether $A \cap B = \emptyset$ or not, thereby solving $\text{DIS}(\mathbf{a}, \mathbf{b})$.

Since this communication protocol is deterministic, by Yao's lower bound in Theorem 7, the number of bits communicated between Alice and Bob must be at least $m = T/3 - 1$. The total number of bits of communication in the protocol is $2pH + dp + 1$ (the last 1 comes from Bob having to communicate the answer to Alice), completing the proof.

E Lower bounds on representing k^{th} -order induction heads: Proof of Theorem 6

In this section we prove the size-lower bound on attention-only transformers representing k^{th} -order induction heads in Theorem 6. To enable this result to be better interpreted, we will break it down into two corollaries.

Corollary 1. *Consider an L -layer attention-only transformer with 1 head per layer and relative position encodings, which satisfies Assumption 1. If $L \leq 1 + \log_2(k-2)$, the attention pattern in layer L of the transformer cannot represent a k^{th} -order induction head.*

Corollary 2. *Consider an 2-layer attention-only transformer with H heads in the first layer and relative position encodings, and assume that Assumption 1 is satisfied. If $H \leq k-3$, the attention pattern in the 2nd layer cannot represent a k^{th} -order induction head.*

We will first prove the result for the case $L = 2$ and $H = 1$, which falls in the intersection of both of these corollaries. We will show that these models cannot represent k^{th} -order induction heads for $k > 3$, under Assumption 1. We subsequently extend it to the general L -layer transformer (i.e., Corollary 1) in Appendix E.2 and to the general case with H_ℓ heads in layer $\ell \in [L]$ in Appendix E.3.

E.1 Lower bounds on 2-layer 1-head attention-only transformers

In this section we show that under Assumption 1, a 2-layer 1-head attention-only transformer cannot represent k^{th} -order induction heads for any $k \geq 4$. We will prove lower bounds on the transformer when the input is binary, i.e., $S = \{0, 1\}$. With relative position embeddings, observe that the first layer of the transformer model learns representations of the form,

$$\mathbf{x}_n^{(2)} = \text{Emb}(x_n) + \sum_{i \leq n} \text{att}_{n,i}^{(1)} \mathbf{W}_V^{(1)} \text{Emb}(x_i) + \sum_{i \leq n} \mathbf{W}_V^{(1)} \mathbf{p}_{n-i}^{V,(1)} \quad (80)$$

where note that the attention pattern only depends on n and i and not on x_i or x_n . These representations are input into the second layer, which realizes the attention pattern $\text{att}_{n,i}^{(2)}$, which is proportional

to,

$$\exp \left(\left\langle \mathbf{W}_K^{(2)} (\mathbf{x}_i^{(2)} + \mathbf{p}_{n-i}^{K,(2)}), \mathbf{W}_Q^{(2)} \mathbf{x}_n^{(2)} \right\rangle \right). \quad (81)$$

We need this function to be maximized uniquely when $x_{i-1} = x_n, \dots, x_{i-k} = x_{n-k+1}$. Denoting $\phi(0) = \mathbf{W}_V^{(1)} \text{Emb}(0)$ and $\phi(1) = \mathbf{W}_V^{(1)} \text{Emb}(1)$,

$$\mathbf{x}_n^{(2)} = \text{Emb}(x_n) + \sum_{i \leq n} \text{att}_{n,i}^{(1)} \mathbf{W}_V^{(1)} \text{Emb}(x_i) + \sum_{i \leq n} \mathbf{W}_V^{(1)} \mathbf{p}_{n-i}^{(1),V} \quad (82)$$

$$= x_n \text{Emb}(1) + (1 - x_n) \text{Emb}(0) + \sum_{i \leq n} \text{att}_{n,i}^{(1)} (x_i \cdot \phi(1) + (1 - x_i) \cdot \phi(0)) + \sum_{i \leq n} \mathbf{W}_V^{(1)} \mathbf{p}_{n-i}^{(1),V} \quad (83)$$

$$= \left(\frac{\text{Emb}(1) + \text{Emb}(0)}{2} + x'_n \cdot \frac{\text{Emb}(1) - \text{Emb}(0)}{2} \right) + \sum_{i \leq n} \text{att}_{n,i}^{(1)} \left(\frac{\phi(1) + \phi(0)}{2} + x'_i \cdot \frac{\phi(1) - \phi(0)}{2} \right) + \sum_{i \leq n} \mathbf{W}_V^{(1)} \mathbf{p}_{n-i}^{(1),V} \quad (84)$$

where $x'_i \leftarrow 2x_i - 1$. We can write this down as,

$$\mathbf{x}_n^{(2)} = \mathbf{m}_n^{(1)} + \mathbf{M}_n^{(1)} [x'_n \ x'_{n-1} \ \dots \ x'_1]^T \quad (85)$$

where $\mathbf{M}_n^{(1)}$ is a matrix of rank at most 2 and of the form,

$$\mathbf{M}_n^{(1)} = \left(\frac{\phi(1) - \phi(0)}{2} \right) \left[\text{att}_{n,n}^{(1)} \ \dots \ \text{att}_{n,1}^{(1)} \right] + \left(\frac{\text{Emb}(1) - \text{Emb}(0)}{2} \right) [1 \ 0 \ \dots \ 0] \quad (86)$$

which is independent of x'_1, \dots, x'_n . Likewise $\mathbf{m}_n^{(1)}$ collects all the vectors in the sum that don't depend on x'_1, \dots, x'_n . Now, observe that in the next layer, we wish to show that an induction head cannot be realized by $\text{att}_{n,i}^{(2)}$ for each $i \leq n$. We will show this for any value of $i \leq n - k$.

In the second layer, we may write down the key vectors as,

$$\mathbf{W}_K^{(2)} (\mathbf{x}_i^{(2)} + \mathbf{p}_{n-i}^{(2),K}) = \mathbf{W}_K^{(2)} \mathbf{m}_i^{(1)} + \mathbf{W}_K^{(2)} \mathbf{M}_i^{(1)} [x'_i \ x'_{i-1} \ \dots \ x'_1]^T + \mathbf{W}_K^{(2)} \mathbf{p}_{n-i}^{(2),K}. \quad (87)$$

Again, defining the vector $\overline{\mathbf{m}}_i^{(1)}$ and the matrix $\overline{\mathbf{M}}_i^{(1)}$ appropriately (having rank at most 2), this equals,

$$\overline{\mathbf{m}}_i^{(1)} (\{x'_i\} \cup \{x'_{i-k-1}, \dots, x'_1\}) + \overline{\mathbf{M}}_i^{(1)} \mathbf{y} \quad (88)$$

where $\mathbf{y} \triangleq [x'_{i-1} \ \dots \ x'_{i-k}]^T$ and the vector $\overline{\mathbf{m}}_i^{(1)}$ depends on x'_i as well as the inputs x'_{i-k-1}, \dots, x'_1 , which in this context, are treated as nuisance variables since they do not intersect with $\{x'_{i-1}, \dots, x'_{i-k}\} \cup \{x_n, \dots, x_{n-k+1}\}$. Henceforth we will avoid explicitly stating the dependency of $\overline{\mathbf{m}}_i^{(1)}$ on the x_j 's. Similarly, the query vector can be written down as,

$$\mathbf{W}_Q^{(2)} \mathbf{x}_n^{(2)} = \widehat{\mathbf{m}}_n^{(1)} + \widetilde{\mathbf{M}}_n^{(1)} \mathbf{x} + \widehat{\mathbf{M}}_n^{(1)} \mathbf{y} \quad (89)$$

where $\widehat{\mathbf{m}}_n^{(1)}$, $\widetilde{\mathbf{M}}_n^{(1)}$ and $\widehat{\mathbf{M}}_n^{(1)}$ are defined appropriately, with $\widehat{\mathbf{M}}_n^{(1)}$ and $\widetilde{\mathbf{M}}_n^{(1)}$ of rank at most 2, and \mathbf{x} is defined as $[x'_n \ \dots \ x'_{n-k+1}]^T$. For an appropriate matrix $\mathbf{M}_{n,i}^\times$, vectors $\mathbf{m}_{n,i}^\times$ and $\widetilde{\mathbf{m}}_{n,i}^\times$ and scalar $m_{n,i}^\times$, the dot-product of the key and query vectors can be written as,

$$\begin{aligned} & \left\langle \mathbf{W}_K^{(2)} (\mathbf{x}_i^{(2)} + \mathbf{p}_{n-i}^{(2),K}), \mathbf{W}_Q^{(2)} \mathbf{x}_n^{(2)} \right\rangle \\ &= \mathbf{x}^T \mathbf{M}_{n,i}^\times \mathbf{y} + \mathbf{y}^T (\overline{\mathbf{M}}_{n,i}^\times) \mathbf{y} + (\mathbf{m}_{n,i}^\times)^T \mathbf{x} + (\widetilde{\mathbf{m}}_{n,i}^\times)^T \mathbf{y} + m_{n,i}^\times \triangleq f_{n,i}(\mathbf{x}, \mathbf{y}), \end{aligned} \quad (90)$$

Which is a linear function in \mathbf{x} and quadratic in \mathbf{y} , both of which lie on $\{\pm 1\}^k$. Note that the matrix $\mathbf{M}_{n,i}^\times$ has rank at most 2 since it is a product of $\overline{\mathbf{M}}_i^{(1)}$ and $\widetilde{\mathbf{M}}_n^{(1)}$, each with rank at most 2. Next we introduce a lemma showing that if $\mathbf{M}_{n,i}^\times$ is inherently low rank, the quadratic form in Equation (90) which captures the dot-product between the key and value vectors cannot satisfy the property that for every \mathbf{y} , the function is uniquely maximized at $\mathbf{x} = \mathbf{y}$. In particular, this means that for any $i \leq n - k$, there is some choice of $x_n, x_{n-1}, \dots, x_{n-k+1}$ such that there are x_{i-1}, \dots, x_{i-k} such that for at least one $j \in [k]$, x_{i-j} and x_{n-j-1} are not equal, but the attention score is larger than the case when x_{i-j} were equal to x_{n-j-1} for each $j \in [k]$.

Lemma 2. If $M_{n,i}^\times$ has rank $\leq k - 2$, it is impossible for $f_{n,i}(\mathbf{x}, \mathbf{y})$ to satisfy the property that for every $\mathbf{y} \in \{\pm 1\}^k$, the maximizer is uniquely $\mathbf{x} = \mathbf{y}$.

The proof is almost complete: if $k \geq 4$, then the rank of $M_{n,i}^\times$, which is at most 2, does not exceed $k - 2$. This means that when $k \geq 4$, any attention pattern realized in the second layer must satisfy the property that there exists a string such that the attention is no longer uniquely maximized when $x_n = x_{i-1}, \dots, x_{n-k+1} = x_{i-k}$.

Proof. For the purpose of brevity, define $\mathcal{H}_k = \{\pm 1\}^k$. First consider the reparameterization,

$$\tilde{\mathbf{x}} = \widetilde{M}_{n,i}^\times \mathbf{x}, \text{ where } \widetilde{M}_{n,i}^\times = \begin{bmatrix} (M_{n,i}^\times)^T \\ (\mathbf{m}_{n,i}^\times)^T \end{bmatrix}. \quad (91)$$

Then, the dot-product of the key and query matrices can be written as,

$$[\mathbf{y}^T \quad 1] \tilde{\mathbf{x}} + \mathbf{y}^T (\widetilde{M}_{n,i}^\times) \mathbf{y} + (\widetilde{\mathbf{m}}_{n,i}^\times)^T \mathbf{y} + m_{n,i}^\times \quad (92)$$

Note that this function is linear in $\tilde{\mathbf{x}}$ and therefore must be maximized on a vertex of the convex hull of the domain, $\widetilde{M}_{n,i}^\times \mathcal{H}_k \triangleq \{\widetilde{M}_{n,i}^\times \mathbf{h} : \mathbf{h} \in \mathcal{H}_k\}$. If $M_{n,i}^\times$ has rank at most $k - 2$, the rank of $\widetilde{M}_{n,i}^\times$ is at most $k - 1$ and cannot be full rank. We show that this must imply that there is a vertex $\mathbf{v} \in \mathcal{H}_k$ such that $\widetilde{M}_{n,i}^\times \mathbf{v}$ is not a unique vertex of the convex hull of $\widetilde{M}_{n,i}^\times \mathcal{H}_k$. This means that \mathbf{v} cannot be a unique maximizer for $\tilde{\mathbf{x}}$ when maximizing over all strings in Equation (92), and specifically $\mathbf{y} = \mathbf{v}$ is a witness to Lemma 2.

Below we discuss how to find such a vector \mathbf{v} . Note that $\widetilde{M}_{n,i}^\times$ is not full rank, which implies that there exists a vector \mathbf{n} such that $\widetilde{M}_{n,i}^\times \mathbf{n} = \mathbf{0}$. Without loss of generality, let n_1 be the smallest non-zero coordinate of \mathbf{n} in absolute value. Then the vector $\mathbf{n}_1^{-1} \mathbf{n}$ has no non-zero coordinates in the interval $(-1, 1)$. We will show that $\text{sign}(\mathbf{n}_1^{-1} \mathbf{n})$ is a good choice for \mathbf{v} .

Consider two cases,

Case I. Every non-zero coordinate of $\mathbf{n}_1^{-1} \mathbf{n}$ is in $\{\pm 1\}$. Consider any $\mathbf{x} \in \mathcal{H}_k$ which matches with \mathbf{n} on the non-zero coordinates. Consider \mathbf{x}' which is the same as \mathbf{x} , except a negation is taken on the coordinates where \mathbf{n} is non-zero. Note that $\widetilde{M}_{n,i}^\times \mathbf{x} = \widetilde{M}_{n,i}^\times \mathbf{x}'$, for the same value of \mathbf{x} . This means that for any \mathbf{y} . In particular, from Equation (92), both \mathbf{x} and \mathbf{x}' are maximizers, showing that Lemma 2 is true in this case. We circumvent having to find such a vector \mathbf{v} in this case.

Case II. $\mathbf{n}_1^{-1} \mathbf{n}$ has non-zero coordinates which are not all in $\{\pm 1\}$. In particular, at least one coordinate where this vector is strictly less than -1 or strictly greater than $+1$. In this case, observe that the sign vector $\tilde{\mathbf{n}} = \text{sign}(\mathbf{n}_1^{-1} \mathbf{n}) \in \mathcal{H}_k$ lies within, but is not a vertex of the convex hull of the set $\mathcal{H}_k \cup \{\mathbf{n}_1^{-1} \mathbf{n}\}$. The reason for this is simple to see when we assume that $\mathbf{n}_1^{-1} \mathbf{n}$ has only one coordinate which is not in $[-1, 1]$, say, the coordinate $j = 2$: here, $\tilde{\mathbf{n}}$ can be written down as a convex combination (with non-zero coefficients) of $\mathbf{n}_1^{-1} \mathbf{n}$ and $\tilde{\mathbf{n}}^{(2)}$; the latter vector is obtained by flipping coordinate 2 of $\tilde{\mathbf{n}}$. When there is more than one coordinate not in $[-1, 1]$, we can peel away these large coordinates in $\mathbf{n}_1^{-1} \mathbf{n}$ by taking a convex combination of this vector with the vectors $\tilde{\mathbf{n}}^{(j)}$ for the appropriate values of j , to return the sign vector $\tilde{\mathbf{n}}$. Here, $\tilde{\mathbf{n}}^{(j)}$ is the version of $\tilde{\mathbf{n}}$ where the j^{th} -coordinate is flipped. This results in the following claim.

Claim 1. The sign vector $\tilde{\mathbf{n}}$ lies within the convex hull of the points $\mathcal{H}_k \cup \{\mathbf{n}_1^{-1} \mathbf{n}\}$, but is not a vertex of this set.

In particular, we may write,

$$\tilde{\mathbf{n}} = \alpha_0 \mathbf{n}_1^{-1} \mathbf{n} + \sum_{j \in [n]} \alpha_j \tilde{\mathbf{n}}^{(j)}. \quad (93)$$

where $\alpha_0 > 0$ and $\sum_{j=0}^n \alpha_j = 1$. By left-multiplying this on both sides by $\widetilde{M}_{n,i}^\times$, and noting that \mathbf{n} lies in the null-space of this matrix, we get,

$$\widetilde{M}_{n,i}^\times \tilde{\mathbf{n}} = \sum_{j \in [n]} \alpha_j \widetilde{M}_{n,i}^\times \tilde{\mathbf{n}}^{(j)} \quad (94)$$

where note that $\sum_{j \in [n]} \alpha_j$ is strictly less than 1, since $\alpha_0 > 0$. We may write this vector as,

$$\begin{aligned} \widetilde{\mathbf{M}}_{n,i}^\times \widetilde{\mathbf{n}} &= \alpha_0 \mathbf{0} + \sum_{j \in [n]} \alpha_j \widetilde{\mathbf{M}}_{n,i}^\times \widetilde{\mathbf{n}}^{(j)} \\ &= \frac{\alpha_0}{2^k} \sum_{\mathbf{h} \in \mathcal{H}_k} \widetilde{\mathbf{M}}_{n,i}^\times \mathbf{h} + \sum_{j \in [n]} \alpha_j \widetilde{\mathbf{M}}_{n,i}^\times \widetilde{\mathbf{n}}^{(j)} \end{aligned} \quad (95)$$

Since $\alpha_0 > 0$, this equation implies that the image of $\widetilde{\mathbf{n}}$ under $\widetilde{\mathbf{M}}_{n,i}^\times$ itself falls within $\text{conv}(\widetilde{\mathbf{M}}_{n,i}^\times \mathcal{H}_k)$, but is itself not a vertex of this set. This means that $\widetilde{\mathbf{n}}$ can never be a maximizer of $f_{n,i}(\cdot, \mathbf{y})$ for any \mathbf{y} , and in particular when $\mathbf{y} = \widetilde{\mathbf{n}}$, thereby proving Lemma 2. \square

E.2 L -layer attention-only transformers with 1 head per layer: Proof of Corollary 1

Proof. The proof largely tracks the 2-layer case, with the main exception that we keep track of how the maximum possible rank of the matrix $\mathbf{M}_{n,i}^\times$ grows as a function of the depth of the transformer. In the case the 2-layer transformer, we show that it cannot exceed 2. With the addition of more layers, we show that it cannot exceed 2^{L-1} .

Recall from the notation in Equation (85) that the output of the first attention layer is,

$$\mathbf{x}_n^{(2)} = \mathbf{m}_n^{(1)} + \mathbf{M}_n^{(1)} [x'_n \quad x'_{n-1} \quad \cdots \quad x'_1]^T \quad (96)$$

where $\mathbf{M}_n^{(1)} \in \mathbb{R}^{d \times n}$ has rank at most 2. Let us rewrite this as,

$$\mathbf{x}_n^{(2)} = \mathbf{m}_n^{(1)} + \overline{\mathbf{M}}_n^{(1)} [x'_T \quad x'_{T-1} \quad \cdots \quad x'_1]^T \quad (97)$$

where $\mathbf{M}_n^{(1)} \in \mathbb{R}^{d \times T}$ is causally masked to be 0's when it operates on x_i for all indices $i > n$. Note that even with this causal masking, $\overline{\mathbf{M}}_n^{(1)}$ has rank at most 2, as discussed in Equation (85).

By induction, assume that the output of the $(\ell - 1)$ th attention layer is of the form,

$$\mathbf{x}_n^{(\ell)} = \mathbf{m}_n^{(\ell-1)} + \overline{\mathbf{M}}_n^{(\ell-1)} \mathbf{x}_{1:T} \quad (98)$$

where $\mathbf{x}_{1:T} \triangleq [x'_T \quad x'_{T-1} \quad \cdots \quad x'_1]^T$. Passing $\mathbf{x}_n^{(\ell)}$ through the ℓ th attention layer, we get,

$$\begin{aligned} \mathbf{x}_n^{(\ell+1)} &= \mathbf{x}_n^{(\ell)} + \sum_{i \leq n} \text{att}_{n,i}^{(\ell)} \mathbf{W}_V^{(\ell)} \left(\mathbf{x}_i^{(\ell)} + \mathbf{p}_{n-i}^{(\ell),V} \right) \\ &= \mathbf{m}_n^{(\ell-1)} + \overline{\mathbf{M}}_n^{(\ell-1)} \mathbf{x}_{1:T} + \sum_{i \leq n} \text{att}_{n,i}^{(\ell)} \mathbf{W}_V^{(\ell)} \mathbf{m}_i^{(\ell-1)} + \sum_{i \leq n} \text{att}_{n,i}^{(\ell)} \mathbf{W}_V^{(\ell)} \overline{\mathbf{M}}_i^{(\ell-1)} \mathbf{x}_{1:T} \\ &\quad + \sum_{i \leq n} \text{att}_{n,i}^{(\ell)} \mathbf{W}_V^{(\ell)} \mathbf{p}_{n-i}^{(\ell),V} \end{aligned} \quad (99)$$

Define,

$$\mathbf{m}_n^{(\ell)} = \mathbf{m}_n^{(\ell-1)} + \sum_{i \leq n} \text{att}_{n,i}^{(\ell)} \mathbf{W}_V^{(\ell)} \mathbf{m}_i^{(\ell-1)} + \sum_{i \leq n} \text{att}_{n,i}^{(\ell)} \mathbf{W}_V^{(\ell)} \mathbf{p}_{n-i}^{(\ell),V}, \text{ and,} \quad (101)$$

$$\overline{\mathbf{M}}_n^{(\ell)} = \overline{\mathbf{M}}_n^{(\ell-1)} + \sum_{i \leq n} \text{att}_{n,i}^{(\ell)} \mathbf{W}_V^{(\ell)} \overline{\mathbf{M}}_i^{(\ell-1)} \quad (102)$$

Then, we can write down,

$$\mathbf{x}_n^{(\ell+1)} = \mathbf{m}_n^{(\ell)} + \overline{\mathbf{M}}_n^{(\ell)} \mathbf{x}_{1:T} \quad (103)$$

We also inductively assume that for every $i \leq n$,

- (i) $\overline{\mathbf{M}}_i^{(\ell-1)}$ has rank $R \leq 2^{\ell-1}$, and,
- (ii) $\overline{\mathbf{M}}_i^{(\ell-1)}$ can be factorized in the form $\sum_{r=1}^R \mathbf{u}_r \cdot \mathbf{v}_{i,r}^T$, where only the $\mathbf{v}_{i,r}$'s depend on i , but the \mathbf{u}_r 's do not depend on i .

Both of these conditions are true when $\ell - 1 = 1$ as evidenced by the structure of $M_i^{(1)}$ in Equation (86) and noting that $\overline{M}_i^{(1)}$ is obtained from $M_i^{(1)}$ by right multiplying by a diagonal mask matrix. Using the recursion in Equation (101), we prove that the induction hypotheses (i) and (ii) are true at layer ℓ as well. In particular using the decomposition in (ii), observe that,

$$\overline{M}_n^{(\ell)} = \sum_{r=1}^R \mathbf{u}_r \cdot \mathbf{v}_{n,r}^T + \sum_{i \leq n} \text{att}_{n,i}^{(\ell)} \mathbf{W}_V^{(\ell)} \sum_{r=1}^R \mathbf{u}_r \cdot \mathbf{v}_{i,r}^T \quad (104)$$

$$= \sum_{r=1}^R \mathbf{u}_r \cdot \mathbf{v}_{n,r}^T + \sum_{r=1}^R \mathbf{W}_V^{(\ell)} \mathbf{u}_r \cdot \left(\sum_{i \leq n} \text{att}_{n,i}^{(\ell)} \mathbf{v}_{i,r} \right)^T \quad (105)$$

$$= \sum_{r=1}^{2R} \mathbf{u}_r \cdot \mathbf{v}_{n,r}^T \quad (106)$$

where for $r' \in [R]$, $\mathbf{u}_{R+r'} \triangleq \mathbf{W}_V^{(\ell)} \mathbf{u}_r$ and $\mathbf{v}_{n,r'} \triangleq \sum_{i \leq n} \text{att}_{n,i}^{(\ell)} \mathbf{v}_{i,r}$. Since $M_n^{(\ell)}$ is the sum of $2R$ rank 1 matrices and therefore has rank at most $2R \leq 2^\ell$, proving both parts of the induction hypothesis.

By induction, at the end of the $(L - 1)$ th layer, we have an output which looks like,

$$\mathbf{x}_n^{(L)} = \mathbf{m}_n^{(L-1)} + \overline{M}_n^{(L-1)} \mathbf{x}_{1:T} \quad (107)$$

where $M_n^{(L-1)}$ has rank at most 2^{L-1} . More importantly, note that by the causal masking, even though it appears to depend on the whole input sequence through $\mathbf{x}_{1:T}$, note that $\mathbf{x}_n^{(L)}$ only depends on x_1, \dots, x_n and not on the future inputs to this time n . In particular, by a similar argument as in the 2-layer case (cf. Equation (85) to Equation (90)), for any $i \leq n - k$ we can decompose the dot-product of the key and query vectors at the L th layer as a bilinear form which looks like,

$$\begin{aligned} & \left\langle \mathbf{W}_K^{(L)} (\mathbf{x}_i^{(L)} + \mathbf{p}_{n-i}^{(L,K)}), \mathbf{W}_Q^{(L)} \mathbf{x}_n^{(L)} \right\rangle \\ &= \mathbf{x}^T M_{n,i}^\times \mathbf{y} + \mathbf{y}^T (\overline{M}_{n,i}^\times) \mathbf{y} + (\mathbf{m}_{n,i}^\times)^T \mathbf{x} + (\overline{\mathbf{m}}_{n,i}^\times)^T \mathbf{y} + m_{n,i}^\times \triangleq f_{n,i}^{(L+1)}(\mathbf{x}, \mathbf{y}) \end{aligned} \quad (108)$$

where \mathbf{x} and \mathbf{y} are defined as $[x'_n \ \dots \ x'_{n-k+1}]^T$ and $[x'_{i-1} \ \dots \ x'_{i-k}]^T$ respectively, and $M_{n,i}^\times$ has rank at most that of $M_n^{(L-1)}$, which is 2^{L-1} . In particular, if $2^{L-1} \leq k - 2$, by Lemma 2 the proof concludes. \square

E.3 The general case: Transformers with H_ℓ heads in layer ℓ : Proof of Theorem 6

The h th head of the first layer of the attention-only transformer learns patterns of the form,

$$\tilde{\mathbf{x}}_n^{(1,h)} = \sum_{i \leq n} \text{att}_{n,i}^{(1,h)} \mathbf{W}_V^{(1,h)} \text{Emb}(x_i) + \sum_{i \leq n} \mathbf{W}_V^{(1,h)} \mathbf{p}_{n-i}^{V,(1,h)} \quad (109)$$

$$= \sum_{i \leq n} \text{att}_{n,i}^{(1,h)} \left(\frac{\phi^h(0) + \phi^h(1)}{2} + x'_i \cdot \frac{\phi^h(1) - \phi^h(0)}{2} \right) + \sum_{i \leq n} \mathbf{W}_V^{(1,h)} \mathbf{p}_{n-i}^{V,(1,h)} \quad (110)$$

where the last equation assumes a binary input sequence, defines $x'_i = 2x_i - 1$ and uses the notation $\phi^h(0) = \mathbf{W}_V^{(1,h)} \text{Emb}(0)$ and $\phi^h(1) = \mathbf{W}_V^{(1,h)} \text{Emb}(1)$. We can further rewrite this as,

$$\tilde{\mathbf{x}}_n^{(1,h)} = \mathbf{m}_n^{(1,h)} + M_n^{(1,h)} \mathbf{x}_{1:T} \quad (111)$$

where each $M_n^{(1,h)} \in \mathbb{R}^{d \times T}$ is rank 1 and applies a causal mask on the inputs x_i for $i > n$. Recall that the output of the first attention layer applies a projection matrix on the concatenation of $\tilde{\mathbf{x}}_n^{(1,h)}$ across $h \in [H_1]$ and then adds a residual connection. The output can be written down as,

$$\tilde{\mathbf{x}}_n^{(2)} = \text{Emb}(x_n) + \mathbf{W}_O^{(1)} \begin{bmatrix} \mathbf{m}_n^{(1,1)} \\ \vdots \\ \mathbf{m}_n^{(1,H_1)} \end{bmatrix} + \mathbf{W}_O^{(1)} \begin{bmatrix} M_n^{(1,1)} \\ \vdots \\ M_n^{(1,H_1)} \end{bmatrix} \mathbf{x}_{1:T} \quad (112)$$

$$= \mathbf{m}_n^{(1)} + \mathbf{M}_n^{(1)} \mathbf{x}_{1:T}, \quad (113)$$

where,

$$\mathbf{M}_n^{(1)} = \left(\frac{\text{Emb}(1) + \text{Emb}(0)}{2} \right) \mathbf{e}_n^T + \mathbf{W}_O^{(1)} \begin{bmatrix} \mathbf{M}_n^{(1,1)} \\ \vdots \\ \mathbf{M}_n^{(1,H_1)} \end{bmatrix}, \text{ and}, \quad (114)$$

$$\mathbf{m}_n^{(1)} = \left(\frac{\text{Emb}(1) - \text{Emb}(0)}{2} \right) + \mathbf{W}_O^{(1)} \begin{bmatrix} \mathbf{m}_n^{(1,1)} \\ \vdots \\ \mathbf{m}_n^{(1,H_1)} \end{bmatrix} \quad (115)$$

Notice that the rank of the matrix $\mathbf{M}_n^{(1)}$ is at most $H_1 + 1$. This is because the concatenation operation can increase the rank at most additively, and since each of the $\mathbf{M}_n^{(1,h)}$ matrices are rank at most 1.

Following through the proof in Appendix E.2 for the L -layer case, we can prove inductively that at any layer ℓ , the output looks like,

$$\mathbf{x}_n^{(\ell)} = \mathbf{m}_n^{(\ell)} + \mathbf{M}_n^{(\ell)} \mathbf{x}_{1:T} \quad (116)$$

where the rank of $\mathbf{M}_n^{(\ell)}$ is $\prod_{i=1}^{\ell} (H_i + 1)$. Invoking Lemma 2, if $\prod_{i=1}^{L-1} (H_i + 1) \leq k - 2$, the attention-only transformer cannot realize a k^{th} -order induction head at layer L .

F Model architecture and hyper-parameters

The experiments were run on one $8 \times A100$ GPU node.

Parameter	Matrix shape
transformer.wte	$2 \times d$
transformer.wpe	$N \times d$
transformer.h.ln_1 ($\times \ell$)	$d \times 1$
transformer.h.attn.c_attn ($\times \ell$)	$3d \times d$
transformer.h.attn.c_proj ($\times \ell$)	$d \times d$
transformer.h.ln_2 ($\times \ell$)	$d \times 1$
transformer.h.mlp.c_fc ($\times \ell$)	$4d \times d$
transformer.h.mlp.c_proj ($\times \ell$)	$d \times 4d$
transformer.ln_f	$d \times 1$

Table 2: Parameters in the transformer architecture with their shape.

G Additional experimental results

Assumption 1 suggests that the attention patterns $\text{att}_{n,i}^{(\ell)}$ in layers $\ell = 1, 2, \dots, L - 1$, as learnt by an L -layer attention-only transformers may only be a function of only the position indices n, i . In this section we run some additional experiments to test this conjecture. We train a 2 layer attention-only transformer with k heads in the first layer, on data drawn from a randomly sampled k^{th} -order Markov process, and focus on the learnt attention patterns as a function of in the input sequence. Figure 7 plots the results of this experiment for $k = 2$ and Figure 8 for $k = 3$. While in both cases there is some variance in the attention patterns learnt by the transformer in some of the heads, we believe that this is a consequence of the iteration budget of the transformer, and specifically the fact that even if the test loss appears to have converged, the transformer may still continue changing in the parameter space. Furthermore, when the attention patterns have some non-zero but small variance as a function of the input, a relaxation of Assumption 1, we also believe that the results we proved in Corollaries 1 and 2 and Theorem 6 should carry over approximately and leave this as an interesting question for future work. Conditional lower bounds of this nature, reliant on structural assumptions the transformer appears to demonstrate in practice are an interesting area of future research.

Dataset	k -th order binary Markov source
Architecture	Based on the GPT-2 architecture as implemented in [30]
Batch size	Grid-searched in $\{8, 16\}$
Accumulation steps	1
Optimizer	AdamW ($\beta_1 = 0.9, \beta_2 = 0.95$)
Learning rate	0.001
Scheduler	Cosine
# Iterations	Up to 25000
Weight decay	1×10^{-3}
Dropout	0
Sequence length	Grid-searched in $\{32, 64, 128, 256, 512, 1024\}$
Embedding dimension	Grid-searched in $\{16, 32, 64\}$
Transformer layers	Between 1 and 8
Attention heads	Up to k
Repetitions	3

Table 3: Settings and parameters for the transformer model used in the experiments.

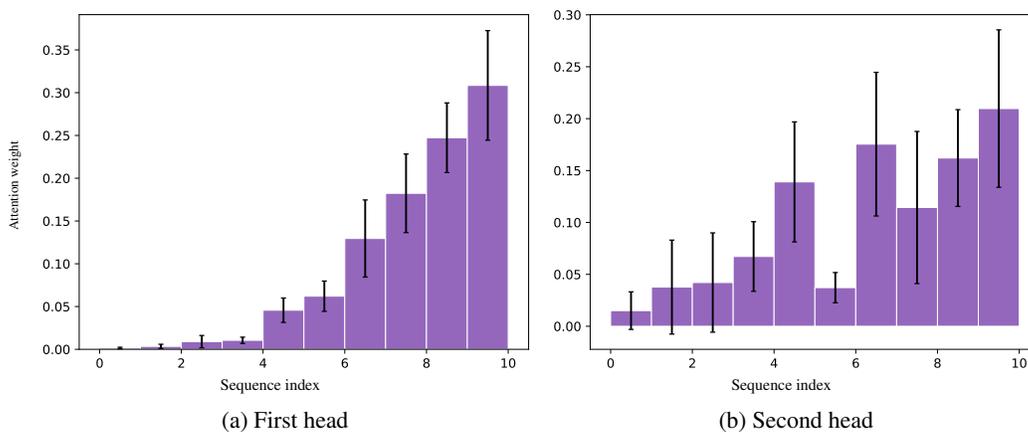


Figure 7: Mean attention for column $n = 10$ of the two heads of the first attention layer, for a 2-layer 2-head transformer model trained on an order-3 Markov process, averaged across 100 input sequences of length 128.

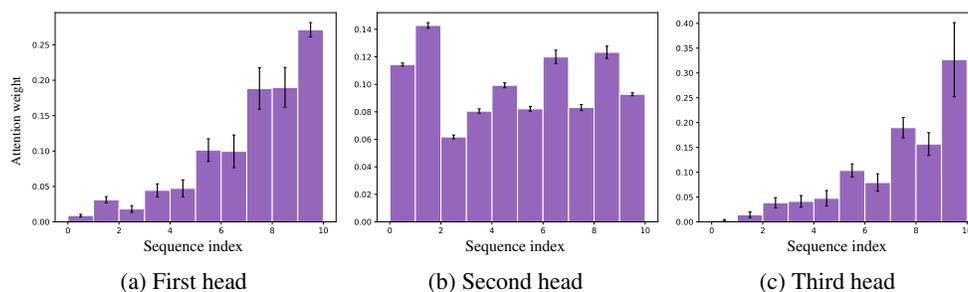


Figure 8: Mean attention for column $n = 10$ of the three heads of the first attention layer, for a 2-layer 3-head transformer model trained on an order-3 Markov process, averaged across 100 input sequences of length 128.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The paper presents experimental evidence about k^{th} -order Markov processes being effectively learnable by small-depth transformers. We also prove theorems about the representation power of low-depth transformers, which are stated in the main paper.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We state that our lower bound in Corollary 1 is a conditional result in Section 6.1, and that it is an important open problem to see whether the conditional statements can be removed. While our work talks about the representational power of transformers, we also mention that learning dynamics of gradient descent is an important direction that future work needs to address.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: All theorems in the paper are stated and cross-referenced with proofs in the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Experiment hyperparameters are presented in Tables 2 and 3 and the code has been provided along with the supplementary material.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility.

In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. **Open access to data and code**

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: All the code has been open sourced; hyperparameter choices have been provided in Tables 2 and 3.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. **Experimental Setting/Details**

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Tables 2 and 3 cover this information.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. **Experiment Statistical Significance**

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Standard error bars are provided on all plots.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Details in Appendix F.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: No NeurIPS code of ethics were violated.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: This work is a primarily theoretical study on the behavior of tokenization on toy problems (learning Markov chains). The societal impact of this research is not likely to be significant.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: No models with a high risk for misuse were trained or released.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Code has been properly credited, via citing the relevant papers.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New Assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: No new assets released.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: No crowdsourcing or research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: No crowdsourcing or research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.