
Transformers Learn to Compress Variable-order Markov Chains in-Context

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In recent years, large language models (LLMs) have demonstrated impressive
2 in-context learning (ICL) capability. However, it is still unclear how the underlying
3 transformers accomplish it, especially in more complex scenarios. Toward this goal,
4 several recent works studied how transformers learn fixed-order Markov chains
5 (FOMC) in context, yet natural languages are more suitably modeled by variable-
6 order Markov chains (VOMC), i.e., context trees (CTs). In this work, we study
7 the ICL of VOMCs by viewing language modeling as a form of data compression
8 and focusing on small alphabets and low-order VOMCs. This perspective allows
9 us to leverage mature compression algorithms, such as the context-tree weighting
10 (CTW) algorithm as a baseline, which is Bayesian optimal for a class of priors. We
11 empirically observe that the performance of transformers is not very sensitive to the
12 number of layers, and even a two-layer transformer can learn in context quite well,
13 tracking closely the performance of CTW. We provide a construction with $D + 2$
14 layers that can mimic the CTW algorithm accurately for VOMCs of maximum
15 order D . One distinction from the FOMC setting is that a counting mechanism
16 plays an important role in this setting.

17 1 Introduction

18 Large language models (LLMs) are capable of completing various tasks (Kasneji et al., 2023; Wu
19 et al., 2023; Thirunavukarasu et al., 2023; Wei et al., 2022). The transformer model (Vaswani et al.,
20 2017), the key behind current prevailing LLMs, is known to have strong in-context learning (ICL)
21 capabilities, and concrete ICL results for transformers have been established for some simple tasks
22 (Garg et al., 2022; Von Oswald et al., 2023; Bai et al., 2024; Ahn et al., 2024). Despite these results,
23 the mechanism for transformers to learn in context is still not fully understood, especially when the
24 scenario is complex or the sequences have memories. Toward this goal, several recent works studied
25 how transformers can learn fixed-order Markov chains (FOMCs) either in training or in-context
26 (Makkuva et al., 2024; Edelman et al., 2024), where insightful observations and theoretical results
27 were obtained. The FOMC is however a poor match for natural languages, for which variable-order
28 Markov chains (VOMCs), also known as context tree (CT) models (Rissanen, 1983; Willems et al.,
29 1995), are often viewed as a more suitable model (Begleiter et al., 2004).

30 To this end, we study the ICL of transformers on VOMCs from the perspective of compression,
31 motivated by a recent work connecting language models and data compression (Delétang et al., 2023).
32 We therefore use compression rates in a fixed context window as our main evaluation metric. This
33 allows us to use several well-known compression algorithms, particularly the context weighting
34 (CTW) algorithm (Willems et al., 1995), as a baseline. The CTW algorithm is Bayesian optimal
35 under certain priors, which gives us a fundamental lower bound in such settings. Appendix A gives a
36 more detailed discussion on related works.

37 We first train a set of shallow transformers of various numbers of layers for VOMCs of various
 38 maximum orders, and empirically observe that the performance of transformers is not very sensitive
 39 to the number of layers, and even a two-layer transformer can learn in context quite well. We then
 40 answer the question of whether transformers can mimic the CTW algorithm. For this purpose, we first
 41 propose an alternative representation of CTW next token prediction, based on which a transformer
 42 construction with $D + 2$ layers is given, that can mimic CTW accurately for VOMCs of maximum
 43 order D . This establishes a fundamental capability of transformers for ICL-VOMC. The alternative
 44 representation enjoys an intuitive interpretation as blending probability estimates along a path on the
 45 context tree.

46 **Main Contributions:** (i) We believe that ours is the first study of ICL for VOMC and we demonstrate
 47 that transformers can indeed (numerically) learn to compress VOMC in-context, close to optimal CTW
 48 algorithm for appropriate CTW-prior. (ii) we give an explicit $D + 2$ -layer transformer construction to
 49 imitate CTW, based on a novel Bayesian optimal next token prediction representation, which can be
 50 of independent interests.

51 2 Preliminaries

52 2.1 The Transformer Model

53 Transformer interacts with sequential data, e.g., $x_1^N = (x_1, \dots, x_N)$, where token x_i is a symbol
 54 from an alphabet (a.k.a. vocabulary) \mathcal{A} with $A = |\mathcal{A}|$. Each token x_i is embedded into $\mathbf{h}_i^{(1)} \in \mathbb{R}^E$ by
 55 integrating the information of its value x_i and position i , where E is the embedding dimension.

56 We introduce an L -layer decoder-only transformer model. Each layer of the transformer takes matrix
 57 $\mathbf{H}^{(\ell)} = [\mathbf{h}_1^{(\ell)}, \mathbf{h}_2^{(\ell)}, \dots, \mathbf{h}_N^{(\ell)}]$, where $\mathbf{h}_i^{(\ell)} \in \mathbb{R}^E$, as its input and applies the multi-head attention
 58 layer operation and the feed-forward layer operation, and the output of the layer is the input to the
 59 next layer, denoted as $\mathbf{H}^{(\ell+1)}$. The decoder-only multi-head attention layer with $M^{(\ell)}$ heads is

$$\mathbf{a}_i^{(\ell)} = \text{MHA}(\mathbf{h}_i, \mathbf{H}; \{W_{O,m}^{(\ell)}, W_{Q,m}^{(\ell)}, W_{K,m}^{(\ell)}, W_{V,m}^{(\ell)}\}_{m=1}^{M^{(\ell)}}) \triangleq W_O^{(\ell)} [\mathbf{b}_{1,i}^{(\ell)}; \mathbf{b}_{2,i}^{(\ell)}; \dots; \mathbf{b}_{M^{(\ell)},i}^{(\ell)}], \quad (1)$$

60 where $\{W_{Q,m}^{(\ell)}, W_{K,m}^{(\ell)}, W_{V,m}^{(\ell)}\}_{m=1}^{M^{(\ell)}}$ are the $E^{(\ell)} \times E$ query matrices, key matrices, and value matrices¹
 61 at the ℓ -th layer and m is the index of the attention head, respectively, $W_O^{(\ell)}$ is the $E \times M^{(\ell)} E^{(\ell)}$
 62 output mapping matrix, and $\mathbf{b}_m^{(\ell)}$ is the output of the m -th attention head at this layer defined as

$$\mathbf{b}_{m,i}^{(\ell)} = (W_{V,m}^{(\ell)} [\mathbf{h}_1^{(\ell)}, \mathbf{h}_2^{(\ell)}, \dots, \mathbf{h}_i^{(\ell)}]) \cdot \text{softmax}((W_{K,m}^{(\ell)} [\mathbf{h}_1^{(\ell)}, \mathbf{h}_2^{(\ell)}, \dots, \mathbf{h}_i^{(\ell)}])^\top (W_{Q,m}^{(\ell)} \mathbf{h}_i^{(\ell)})), \quad (2)$$

63 where we used “;” to indicate vertical matrix concatenation and “,” to indicate horizontal matrix
 64 concatenation. The attention layer has a residual connection, and the attention output together with
 65 the residual connection also goes through a feedforward layer with a residual connection

$$\mathbf{h}_i^{(\ell+1)} = \text{FF}(\mathbf{a}_i; W_1^{(\ell)}, W_2^{(\ell)}) = W_1^{(\ell)} \sigma(W_2^{(\ell)} (\mathbf{a}_i^{(\ell)} + \mathbf{h}_i^{(\ell)})) + (\mathbf{a}_i^{(\ell)} + \mathbf{h}_i^{(\ell)}), \quad (3)$$

66 where σ is a non-linear activation function (e.g., ReLU or sigmoid). The output of the last (L -th)
 67 transformer layer $\mathbf{H}^{(L+1)}$ goes through a linear then softmax unit to predict the probability of
 68 generating the next symbol in vocabulary \mathcal{A} based on the past observations:

$$\hat{\mathbf{p}}_{i+1} = \text{softmax}(W_O^{(L+1)} \mathbf{h}_i^{(L+1)}) \in \Delta_A, \quad i = 1, \dots, N - 1, \quad (4)$$

69 where Δ_A is the probability simplex on \mathcal{A} . The model is illustrated in Appendix C.

70 2.2 Context Tree Models (Variable-Order Markov Chains)

71 Variable-order Markov chains (VOMCs), also known as context tree (CT) models, have been studied
 72 extensively in the data compression literature (Rissanen, 1983; Willems et al., 1995; Begleiter et al.,
 73 2004). String $s = (x_{1-l}, x_{2-l}, \dots, x_0)$ is a suffix of the string $s' = (x'_{1-l'}, x'_{2-l'}, \dots, x'_0)$, if
 74 $0 \leq l \leq l'$ and $x_{-i} = x'_{-i}$ for $i = 0, 1, \dots, l - 1$; e.g., (a, b, c, b) is suffix of (a, c, a, a, b, c, b) .

75 A CT source is specified by a suffix set \mathcal{S} and the associated probability distributions. The suffix set
 76 is a collection of strings $s(k)$, $k = 1, \dots, |\mathcal{S}|$, which needs to be proper and complete: The set is

¹In practice, embedding dimension E is divisible by the number of heads $M^{(\ell)}$ and $E = M^{(\ell)} E^{(\ell)}$.

77 proper if no string in \mathcal{S} is a suffix of any other string; it is complete if each semi-infinite sequence
 78 (\dots, x_{n-1}, x_n) has a unique suffix that belongs to \mathcal{S} , denoted as $\beta_{\mathcal{S}}(\dots, x_{n-1}, x_n)$. Associated with
 79 each suffix $s \in \mathcal{S}$, there is a probability mass function $p_s \in \Delta_{\mathcal{A}}$. A CT has maximum order D if
 80 any suffix in \mathcal{S} has length at most D . Given a semi-infinite sequence (\dots, x_{n-1}, x_n) , the next
 81 symbol x_{n+1} is generated randomly according to the distribution $p_{\beta_{\mathcal{S}}(\dots, x_{n-1}, x_n)}$. An example CT is
 82 in Fig. 5 in the appendix. For each suffix set \mathcal{S} , there is a unique tree T with suffix set \mathcal{S} being its
 83 leaves $\mathcal{L}(T)$, and a CT can thus be represented by $(T, \{p_s\}_{s \in \mathcal{L}(T)})$.

84 2.3 Bayesian Context Tree Weighting Compression Algorithm

85 Once the underlying CT is estimated accurately, arithmetic coding (AC) can be used to
 86 compress the sequence efficiently. The likelihood of a sequence x_1^i given x_{1-D}^0 for a
 87 CT with parameter $(T, \{p_s\}_{s \in \mathcal{L}(T)})$ is $P_{T, \{p_s\}}(x_1^i | x_{1-D}^0) = \prod_{j=1}^i p_{\beta_{\mathcal{L}(T)}(x_{j-D}, \dots, x_{j-1})}(x_j) =$
 88 $\prod_{s \in \mathcal{L}(T)} \prod_{a \in \mathcal{A}} p_s(a)^{\mathbf{n}_{i,s}(a)}$, where $\mathbf{n}_{i,s}$ is the *counting vector* associated with suffix s that

$$\mathbf{n}_{i,s}(a) := \text{number of times symbol } a \in \mathcal{A} \text{ follows suffix } s \text{ in sequence } (x_1, \dots, x_i). \quad (5)$$

89 Willems et al. (1995) proposed the context tree weighting (CTW) algorithm for CT sources. CTW
 90 estimates the probability of the sequence x_1^n by the auxiliary parameters p^e, p^w 's as follows.

- 91 1. For each $s \in \mathcal{A}^*$ with $|s| \leq D$, compute $p_{n,s}^e = \frac{\Gamma(\sum_{a \in \mathcal{A}} \alpha(a))}{\Gamma(\sum_{a \in \mathcal{A}} (\mathbf{n}_s(a) + \alpha(a)))} \prod_{q \in \mathcal{A}} \frac{\Gamma(\mathbf{n}_s(a) + \alpha(a))}{\Gamma(\alpha(a))}$,
 92 where \mathbf{n}_s is the counting vector $\mathbf{n}_{i,s}$ with $i = n$, and $\Gamma(\cdot)$ is the Gamma function.
- 93 2. From nodes in the D -th level to the 0-th level (i.e., root), iteratively compute

$$p_{n,s}^w := \begin{cases} p_{n,s}^e, & \text{if } |s| = D, \\ \lambda p_{n,s}^e + (1 - \lambda) \prod_{q \in \mathcal{A}} p_{n,qs}^w, & \text{otherwise,} \end{cases} \quad (6)$$

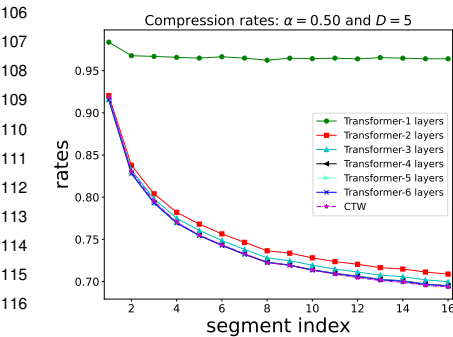
94 where qs is the string by appending symbol $q \in \mathcal{A}$ before the suffix s .

95 Kontoyiannis et al. (2022) took the Bayesian view towards this procedure under a
 96 CTW prior. CTW prior π_{CTW} is a Bayesian CT prior over the trees in $\mathcal{T}(D) :=$
 97 $\{\text{full } A\text{-ary tree with depth at most } D\}$ and the transition distributions $p_s \in \Delta_{\mathcal{A}}$. Specifically,
 98 $\pi_{\text{CTW}}(T, (p_s)_{s \in \mathcal{L}(T)}) = \pi_D(T) \prod_{s \in \mathcal{L}(T)} \pi_p(p_s)$ with

$$\pi_D(T) = (1 - \lambda)^{(|\mathcal{L}(T)| - 1) / (A - 1)} \lambda^{|\mathcal{L}(T)| - |\mathcal{L}_D(T)|}, \quad \pi_p(p_s) = \text{Dir}(p_s; \{\alpha(a)\}_{a \in \mathcal{A}}).$$

99 $\pi_D(\cdot)$ represents a bounded branching process with stopping probability λ for each node; and
 100 $\mathcal{L}_D(T)$ is the leaves of T with depth D . The next token probability p_s follows a Dirichlet prior
 101 parameterized by $\{\alpha(a)\}$. Kontoyiannis et al. (2022) showed that the $p_{n,(\cdot)}^w$ at root computed by CTW
 102 equals to the Bayesian predicted probability under CTW prior, i.e., $p_{n,(\cdot)}^w = P_{\pi_{\text{CTW}}}(x_1^n | x_{1-D}^0) =$
 103 $\sum_{T \in \mathcal{T}(D)} \int P_{T, \{p_s\}}(x_1^n | x_{1-D}^0) \pi(T, \{p_s\}) \left(\prod_{s \in \mathcal{L}(T)} dp_s \right)$. AC can be applied via sequentially
 104 calculating the predictive next token probability as $P_{\pi_{\text{CTW}}}(x_{i+1} | x_{1-D}^i) = \frac{P_{\pi_{\text{CTW}}}(x_1^{i+1} | x_{1-D}^0)}{P_{\pi_{\text{CTW}}}(x_1^i | x_{1-D}^0)}$.

105 3 Transformers Learn In-context of VOMCs



117 Figure 1: Transformer vs CTW

118 the reference CTW algorithm. Experimental details are in Appendix D.
 119

We choose ternary alphabet $|\mathcal{A}| = 3$, and pretrain a trans-
 former of context window size N on data sequences of
 length- N generated using CTs randomly sampled from a
 CTW prior π_{CTW} parameterized by $\alpha = 0.5, \lambda = 0.15$ and
 a fixed maximum tree depth D , illustrated in Fig. 7 in Ap-
 pendix D. The training loss is the canonical next-token pre-
 diction cross-entropy loss. During the inference, given a
 source sequence of length- N generated from an unknown
 VOMC with the order at most D , can the transformer com-
 press this sequence efficiently, i.e., at a compression rate
 close to the optimal rate?

In Fig. 1, we show the performance comparisons between
 trained transformers with various numbers of layers and
 the reference CTW algorithm. Experimental details are in Appendix D.

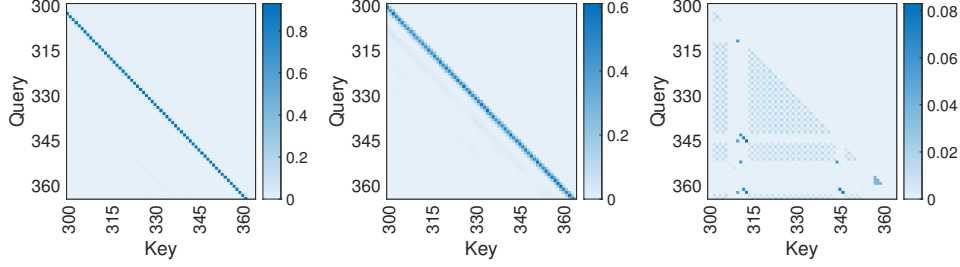


Figure 2: Partial attention heatmaps for different attention heads.

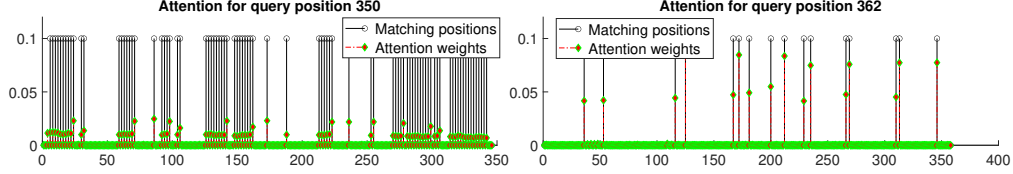


Figure 3: Suffix locations and attention weights in the second type of pattern at two query positions.

120 We observe that almost all trained transformers, except that with a single layer, track the performance
 121 of the CTW algorithm fairly closely. The overall performance does improve as the number of layers
 122 increases in general; see Table 1 in the Appendix D.1 for numerical comparisons. Nevertheless, the
 123 improvements with increased numbers of layers are relatively small. Even transformers with two
 124 layers appear to learn in context quite well.

125 4 Theoretical Interpretations and Empirical Evidences

126 4.1 Analysis of Attention Maps

127 To understand why and how the trained transformers perform comparable CTW, we first analyzed the
 128 attention maps of the trained transformers where two distinguished patterns emerge. One pattern is
 129 solely relative-position dependent. In the left two panels of Fig. 3, we observe off-diagonal stripes
 130 for these two attention heads, which are a few positions below the main diagonal. They can be a
 131 single off-diagonal or a collection of several off-diagonals. This indicates that the query position is
 132 attending positions at a few fixed but close distances ahead of itself. This pattern usually appears in
 133 the first or second layers of the transformers. Combining with the suffix structure in compression
 134 algorithms such as CTW, such an attention pattern suggests the suffix is being copied into the current
 135 query position for subsequent processing.

136 Another pattern, shown in the third panel has more sophisticated spotty patterns, and the attention
 137 appears to depend more explicitly on the current token features instead of the position alone, and they
 138 usually appear in the second layer or above in the transformers. Taking query positions 350 and 362
 139 for the attention head shown in the third panel of Fig. 2, we plot in Fig. 3 the positions in the data
 140 sequence that match their suffixes of length-3 using the stem plots with a black circle on top, and the
 141 attention values as the red stems with the diamonds on top. This attention pattern suggests that it is
 142 collecting information for those positions with the matched suffix of a fixed length.

143 4.2 Capability and capacity of transformer via construction

144 Given a sequence x_1^n generated according to a $CT(T, \{p_s\})$ sampled from the CTW-prior π_{CTW} pa-
 145 rameterized by (D, λ, α) , we propose a novel representation for computing the predictive probability
 146 $P_{\pi_{CTW}}(x_{n+1}|x_{1-D}^n)$ in the following theorem. It is based on the weighted average of the next token
 147 prediction probability vector of each potential suffix $s_{n,l} := x_{n-l+1}^n$ of length $l = 0, 1, \dots, D$. The
 148 proof of Theorem 1 is in Appendix E.1.

149 **Theorem 1.** *The predicted probability can be computed as*

$$P_{\pi_{CTW}}(x_{n+1}|x_1^n) = \sum_{l=0, \dots, D} \omega_{n,l} \cdot \mathbf{p}_{n,s_{n,l}}(x_{n+1}), \quad (7)$$

150 where $\mathbf{p}_{n,s_{n,l}}(a) = \frac{\alpha(a) + \mathbf{n}_{n,s_{n,l}}(a)}{\sum_q (\alpha(q) + \mathbf{n}_{n,s_{n,l}}(q))}$; and $\omega_{n,\cdot} \in \Delta_{D+1}$ with $\ln(\omega_{n,l}) - \ln(\omega_{n,l-1}) = \ln(1-\lambda) -$
 152 $\mathbb{I}_{l=D} \ln(\lambda) + \ell_{n,s_{n,l}}^e - \ell_{n,s_{n,l-1}}^e + \sum_{q \in \mathcal{A}} \ell_{n,q s_{n,l-1}}^w - \ell_{n,s_{n,l}}^w$, where $\ell_{n,s}^e = \ln(p_{n,s}^e)$, $\ell_{n,s}^w = \ln(p_{n,s}^w)$.

153 As illustrated in Fig. 4, each suffix $s_{n,l}$, e.g., $s_{n,0} =$
 154 $()$, $s_{n,2} = ba$, can potentially be the true suffix of the
 155 underlying CT dynamics, i.e., $s_{n,l} \in \mathcal{L}(T)$; and $\mathbf{p}_{n,s_{n,l}}$
 156 is in fact the Bayesian optimal next token prediction
 157 given $s_{n,l} \in \mathcal{L}(T)$. The weights $\omega_{n,l}$ assign credibility
 158 that $s_{n,l}$ is the true suffix. Theorem 1 suggests that the
 159 weights are based on both the information in the suffix
 160 path such as $p_{s_{n,l}}^e$ as well as the information from
 161 their siblings $p_{n,q s_{n,l-1}}^w$ (siblings are in triangles in Fig.
 162 4). The information of counting vector $\mathbf{n}_{n,s}$ plays a vital
 163 role since $\mathbf{p}_{n,s}$, $p_{n,s}^e$, e.t.c. are all functions of $\mathbf{n}_{n,s}$.

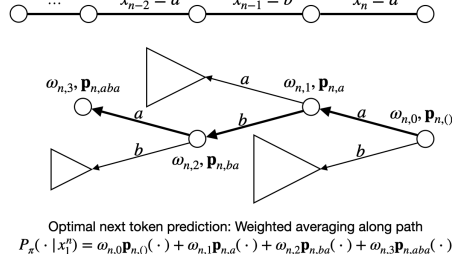


Figure 4: Illustration of Theorem 1

164 4.3 Transformer construction: Approximating CTW

165 We provide a construction of $(2 + D)$ -layer transformer with sufficient representation power in the
 166 FF layer that can essentially approximate CTW, demonstrating the capacity of the transformer. The
 167 first two layers are motivated by the attention map patterns observed in Section 4.1, which we show
 168 their capabilities of capturing the important counting vector statistics suggested by Theorem 1. The
 169 last D layers are induction layers imitating the CTW procedure.

170 We consider the initial embedding is one-hot, with additional scratch pad elements initialized as
 171 zeros and a positional embedding, i.e., $\mathbf{h}_i^{(1)} = (\mathbf{x}_i; \mathbf{0}; \mathbf{pos}_i)$ where $\mathbf{x}_i \in \mathbb{R}^A$ is the one-hot (column
 172 vector) embedding of x_i , $\mathbf{pos}_i = (1, \cos(i\pi/N), \sin(i\pi/N))^T$ is a positional embedding, and the
 173 remaining $(E - A - 3)$ elements being zero. The proofs of this section are in Appendix E.2.

174 We begin with the first layer, which is referred to as a *finite-memory context-extension layer*.

175 **Theorem 2.** *There is an M -headed transformer layer that can perform finite-memory context*
 176 *extension, defined by the following output, with the initial one-hot embedded input $\mathbf{H}^{(1)}$:*

$$175 \mathbf{h}_i^{(2)} = (\mathbf{s}_{i,M+1}; \mathbf{0}; \mathbf{pos}_i), \quad (8)$$

177 where $\mathbf{s}_{i,M+1} = (\mathbf{x}_i; \dots; \mathbf{x}_{i-M})$ is the vector version of the M -length suffix $s_{i,M+1} = x_{i-M}^i$.

178 This layer copies and stacks M past embedded symbols to the current position i . It utilizes the
 179 positional encoding \mathbf{pos}_i via rotation and matching the corresponding positions.

180 The second layer is referred to as the *statistics collection layer*, which takes a sequence of vectors
 181 $\mathbf{h}_i^{(2)}$, $i = 1, \dots, N$, defined in (8) as its input. To rigorously specify the function of this layer, we
 182 define the forward and backward statistics vectors at position i ,

$$182 \mathbf{g}_{i,s}(a) = \frac{\mathbf{n}_{i,s}(a)}{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s}(q)}, \quad \mathbf{g}_{i-1,s}^{\leftarrow}(a) = \frac{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,as}(q)}{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s}(q)}, \quad \forall a \in \mathcal{A}, \quad (9)$$

183 where $\mathbf{n}_{i,s}$ is the counting vector defined in (5), and $\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s}(q)$ is the number of appears of the
 184 string s in the sequence x_1^{i-1} . In plain words, with $|s| = k - 1$ they are the empirical probability of
 185 the next and previous token associated with the suffix s in the k -gram statistics seen before x_i . For
 186 both $\mathbf{g}_{i,s}$ and $\mathbf{g}_{i-1,s}^{\leftarrow}$, if the suffix s has not appeared in data x_1^{i-1} , it can be initialized arbitrarily as a
 187 vector in the probability simplex.

188 **Theorem 3.** *There is an M' -head attention layer, where $M' \leq M + 1$, that can perform statistics*
 189 *collection, defined by the following output, with $\mathbf{H}^{(2)}$ in (8) as its input:*

$$188 \mathbf{a}_i^{(2)} = (\mathbf{s}_{i,M+1}; \mathbf{g}_{i,M'}; \mathbf{g}_{i-1,M'}^{\leftarrow}; \mathbf{0}; \mathbf{pos}_i), \quad (10)$$

190 where $\mathbf{g}_{i,M'} := (\mathbf{g}_{i,s_{i,0}}; \dots; \mathbf{g}_{i,s_{i,M'-1}})$ and $\mathbf{g}_{i-1,M'}^{\leftarrow} = (\mathbf{g}_{i-1,s_{i,0}}^{\leftarrow}; \dots; \mathbf{g}_{i-1,s_{i,M'-1}}^{\leftarrow})$.

191 This functional layer essentially collects k -gram statistics for various lengths of $k = 1, 2, \dots, M'$.
 192 For example, when $k = 3$, it collects the normalized frequency associated with the suffix (x_{n-1}, x_n) .

193 For ICL of FOMCs, two-layer transformers collecting forward statistics $\mathbf{g}_{i,M'}$ with $M' =$
 194 $D + 1$ is sufficient (Edelman et al., 2024). However, for the ICL-VOMC task, the under-

195 lying CT structure is unknown, therefore, collecting such simple statistics is no longer suf-
 196 ficient. As indicated in Theorem 1, the information of counting statistics $\mathbf{n}_{i,s_i,l}$ is impor-
 197 tant to the performance of prediction since the weights heavily depend on $\mathbf{n}_{i,s}(a)$. Yet due
 198 to the softmax function of the attention layer, only (normalized) probabilistic vector can be
 199 obtained instead of the exact count. With the backward statistics $\mathbf{g}_{i,s}^{\leftarrow}$, $\mathbf{n}_{i,s_i,l}$ can be de-
 200 rived as $\mathbf{n}_{i,s_i,l}(a) = \frac{\mathbf{n}_{i,s_i,l}(a)}{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s_i,l}(q)} \frac{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s_i,l}(q)}{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s_i,l-1}(q)} \cdots \frac{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s_i,1}(q)}{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s_i,0}(q)} \left(\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s_i,0}(q) \right) =$
 201 $\mathbf{g}_{i,s_i,l}(a) \left(\prod_{j=0}^{l-1} \mathbf{g}_{i-1,s_i,j}^{\leftarrow}(x_{i-j}) \right) i$, by the information contained in vector $\mathbf{a}_i^{(2)}$.

202 Taking $M = M' - 1 = D$ and a sufficiently wide FF layer in the second transformer layer, we have

$$\mathbf{h}_i^{(3)} = (\mathbf{s}_{i,D}; \mathbf{p}_{i,D}; \mathbf{l}_{i,D}^e; \ell_{i,s_i,D}^w; \mathbf{0}; \mathbf{pos}_i), \quad (11)$$

203 where $\mathbf{p}_{i,D} = (\mathbf{p}_{i,s_i,0}; \cdots; \mathbf{p}_{i,s_i,D})$ and $\mathbf{l}_{i,D}^e = (\ell_{i,s_i,0}^e; \cdots; \ell_{i,s_i,D}^e)$, by universal approximation.

204 To fulfill the Bayesian optimal prediction, we introduce the following *CTW induction layer* that
 205 iteratively computes $\ell_{i,s}^w$ on the suffix path and their siblings, and also the weight difference $\delta_{i,l} :=$
 206 $\ln(\omega_{i,l}) - \ln(\omega_{i,l-1})$ for $l = d, D-1, \dots, 1$. The desired embedding for $\ell = 3, 4, \dots, 3 + D$ is

$$\mathbf{h}_i^{(\ell)} = (\mathbf{s}_{i,M^{(1)}+1}; \mathbf{p}_{i,D}; \mathbf{l}_{i,D}^e; \delta_{i,D}; \delta_{i,D-1}; \cdots; \delta_{i,D-\ell+4}; \ell_{i,s_i,D+3-\ell}^w; \mathbf{0}; \mathbf{pos}_i). \quad (12)$$

207 **Theorem 4.** *There exists a A -head transformer layer that can perform the induction: Takes $\mathbf{H}^{(\ell)}$ in*
 208 *(12) as input and outputs $\mathbf{H}^{(\ell+1)}$. And the final output layer taking $\mathbf{H}^{(D+3)}$ as input can output the A -*
 209 *dimensional Bayesian optimal next token prediction vector $P_{\pi_{CTW}}(\cdot | x_{1-D}^n) = \sum_{l=0, \dots, D} \omega_{n,l} \mathbf{p}_{n,s_n,l}$.*

210 Although transformers with sufficient FF layers can theoretically compute the optimal prediction as
 211 CTW, empirically, transformers of $2 + D$ layers perform slightly worse in our experiments. This is
 212 likely due to the less-than-perfect pretraining optimization and the limited representation capability of
 213 finite-width FF layers with ReLU activation. We also note that the proposed transformer construction
 214 may not be the only way to mimic CTW, however, we believe the first two layers do capture important
 215 universal features. We provide supporting evidence empirically in the sequel.

216 **Hybrid transformer with two-layer construction** We construct hybrid versions of transformers,
 217 with details given in Appendix D.1.2. We train a two-layer transformer with a constructed $\mathbf{a}_i^{(2)}$
 218 followed by a trainable FF layer (the FF layer in the second layer of the transformer) and an output
 219 layer, and compare the impacts different choices of $\mathbf{a}_i^{(2)}$. We replace the backward statistics $\mathbf{g}_{i,1,M'}^{\leftarrow}$
 220 and \mathbf{pos}_i with $\{\mathbf{n}_{n,s_n,l}\}_{l=0}^D$ and i in $\mathbf{a}_i^{(2)}$ in Eq. (10), and notice its performance is almost the same
 221 as the one using $\mathbf{a}_i^{(2)}$ in Eq. (10), and their performances are close to that of canonical 2-layer
 222 transformer. Moreover, the performances get worse if the statistics like $\{\mathbf{n}_{n,s_n,l}\}_{l=0}^D$ and i are further
 223 removed. Thus such counting statistics are necessary and essential for the ICL of VOMC sources.
 224 More discussions and experiments with 4-layer transformers with one or two constructed layers are
 225 in Appendix D.1.2.

226 5 Conclusion

227 We considered the in-context learning of transformers for VOMC sources. By drawing a close analogy
 228 of ICL and Bayesian universal compression, we leverage the CTW as a baseline. Experimentally,
 229 we observe the performances of the trained transformers are close to that of CTW even with just
 230 two layers under CTW priors. To understand the mechanism of transformers' ICL ability, we
 231 analyzed the attention maps and extracted two likely mechanisms. We then construct the finite-
 232 memory context extension layer, and the statistics collection layer, corresponding to these two
 233 mechanisms, respectively. The latter collects both the forward and backward statistics, which are vital
 234 as theoretically demonstrated by a novel representation of the CTW optimal next-token prediction.
 235 We also provide empirical evidence that the statistics collected by the constructed second layer, in
 236 particular the counting statistics, are indeed necessary.

237 Although we empirically showed transformers can perform ICL-VOMC tasks and constructed an
 238 idealized transformer to mimic the CTW algorithm, it is not clear whether a trained transformer will
 239 indeed utilize the upper layer mechanisms. Extending the existing approach (Edelman et al., 2024) to
 240 answer this question appears quite difficult, given the complexity of the constructed transformer and
 241 the underlying VOMCs; this is part of ongoing investigations.

242 References

- 243 Ahn, K., Cheng, X., Daneshmand, H., and Sra, S. (2024). Transformers learn to implement pre-
244 conditioned gradient descent for in-context learning. *Advances in Neural Information Processing*
245 *Systems*, 36.
- 246 Akyürek, E., Schuurmans, D., Andreas, J., Ma, T., and Zhou, D. (2022). What learning algorithm is
247 in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*.
- 248 Akyürek, E., Wang, B., Kim, Y., and Andreas, J. (2024). In-context language learning: Architectures
249 and algorithms. *arXiv preprint arXiv:2401.12973*.
- 250 Bai, Y., Chen, F., Wang, H., Xiong, C., and Mei, S. (2024). Transformers as statisticians: Provable
251 in-context learning with in-context algorithm selection. *Advances in neural information processing*
252 *systems*, 36.
- 253 Begleiter, R., El-Yaniv, R., and Yona, G. (2004). On prediction using variable order Markov models.
254 *Journal of Artificial Intelligence Research*, 22:385–421.
- 255 Bietti, A., Cabannes, V., Bouchacourt, D., Jegou, H., and Bottou, L. (2024). Birth of a transformer: A
256 memory viewpoint. *Advances in Neural Information Processing Systems*, 36.
- 257 Burrows, M. (1994). A block-sorting lossless data compression algorithm. *SRS Research Report*,
258 124.
- 259 Cleary, J. and Witten, I. (1984). Data compression using adaptive coding and partial string matching.
260 *IEEE Transactions on Communications*, 32(4):396–402.
- 261 Cleary, J. G. and Teahan, W. J. (1997). Unbounded length contexts for PPM. *The Computer Journal*,
262 40(2_and_3):67–75.
- 263 Dai, D., Sun, Y., Dong, L., Hao, Y., Ma, S., Sui, Z., and Wei, F. (2022). Why can gpt learn in-
264 context? language models implicitly perform gradient descent as meta-optimizers. *arXiv preprint*
265 *arXiv:2212.10559*.
- 266 Delétang, G., Ruoss, A., Duquenne, P.-A., Catt, E., Genewein, T., Mattern, C., Grau-Moya, J.,
267 Wenliang, L. K., Aitchison, M., Orseau, L., et al. (2023). Language modeling is compression.
268 *arXiv preprint arXiv:2309.10668*.
- 269 Edelman, B. L., Edelman, E., Goel, S., Malach, E., and Tsilivis, N. (2024). The evolution of statistical
270 induction heads: In-context learning Markov chains. *arXiv preprint arXiv:2402.11004*.
- 271 Garg, S., Tsipras, D., Liang, P. S., and Valiant, G. (2022). What can transformers learn in-context?
272 a case study of simple function classes. *Advances in Neural Information Processing Systems*,
273 35:30583–30598.
- 274 Kasneci, E., Seßler, K., Küchemann, S., Bannert, M., Dementieva, D., Fischer, F., Gasser, U.,
275 Groh, G., Günemann, S., Hüllermeier, E., et al. (2023). ChatGPT for good? on opportunities
276 and challenges of large language models for education. *Learning and individual differences*,
277 103:102274.
- 278 Kirsch, L., Harrison, J., Sohl-Dickstein, J., and Metz, L. (2022). General-purpose in-context learning
279 by meta-learning transformers. *arXiv preprint arXiv:2212.04458*.
- 280 Kontoyiannis, I. (2023). Context-tree weighting and Bayesian context trees: Asymptotic and non-
281 asymptotic justifications. *IEEE Transactions on Information Theory*.
- 282 Kontoyiannis, I., Mertzanis, L., Panotopoulou, A., Papageorgiou, I., and Skoularidou, M. (2022).
283 Bayesian context trees: Modelling and exact inference for discrete time series. *Journal of the*
284 *Royal Statistical Society Series B: Statistical Methodology*, 84(4):1287–1323.
- 285 Li, Y., Ildiz, M. E., Papailiopoulos, D., and Oymak, S. (2023). Transformers as algorithms: Gener-
286 alization and stability in in-context learning. In *International Conference on Machine Learning*,
287 pages 19565–19594. PMLR.

- 288 Makuva, A. V., Bondaschi, M., Girish, A., Nagle, A., Jaggi, M., Kim, H., and Gastpar, M. (2024).
 289 Attention with Markov: A framework for principled analysis of transformers via Markov chains.
 290 *arXiv preprint arXiv:2402.04161*.
- 291 Moffat, A. (1990). Implementing the PPM data compression scheme. *IEEE Transactions on*
 292 *communications*, 38(11):1917–1921.
- 293 Olsson, C., Elhage, N., Nanda, N., Joseph, N., DasSarma, N., Henighan, T., Mann, B., Askell,
 294 A., Bai, Y., Chen, A., et al. (2022). In-context learning and induction heads. *arXiv preprint*
 295 *arXiv:2209.11895*.
- 296 Pasco, R. (1976). *Source coding algorithms for fast data compression*. PhD thesis, Stanford
 297 University.
- 298 Rajaraman, N., Bondaschi, M., Ramchandran, K., Gastpar, M., and Makuva, A. V. (2024). Trans-
 299 formers on markov data: Constant depth suffices. *arXiv preprint arXiv:2407.17686*.
- 300 Reddy, G. (2023). The mechanistic basis of data dependence and abrupt learning in an in-context
 301 classification task. *arXiv preprint arXiv:2312.03002*.
- 302 Rissanen, J. (1983). A universal data compression system. *IEEE Transactions on information theory*,
 303 29(5):656–664.
- 304 Rissanen, J. and Langdon, G. G. (1979). Arithmetic coding. *IBM Journal of research and development*,
 305 23(2):149–162.
- 306 Rissanen, J. (1976). Generalized kraft’s inequality and arithmetic coding. *IBM Journal on Research*
 307 *Development*, 20.
- 308 Sadakane, K., Okazaki, T., and Imai, H. (2000). Implementing the context tree weighting method
 309 for text compression. In *Proceedings DCC 2000. Data Compression Conference*, pages 123–132.
 310 IEEE.
- 311 Shkarin, D. (2002). Ppm: One step to practicality. In *Proceedings DCC 2002. Data Compression*
 312 *Conference*, pages 202–211. IEEE.
- 313 Thirunavukarasu, A. J., Ting, D. S. J., Elangovan, K., Gutierrez, L., Tan, T. F., and Ting, D. S. W.
 314 (2023). Large language models in medicine. *Nature medicine*, 29(8):1930–1940.
- 315 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and
 316 Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing*
 317 *systems*, 30.
- 318 Von Oswald, J., Niklasson, E., Randazzo, E., Sacramento, J., Mordvintsev, A., Zhmoginov, A., and
 319 Vladymyrov, M. (2023). Transformers learn in-context by gradient descent. In *International*
 320 *Conference on Machine Learning*, pages 35151–35174. PMLR.
- 321 Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M.,
 322 Zhou, D., Metzler, D., et al. (2022). Emergent abilities of large language models. *arXiv preprint*
 323 *arXiv:2206.07682*.
- 324 Willems, F. M. (1998). The context-tree weighting method: Extensions. *IEEE Transactions on*
 325 *Information Theory*, 44(2):792–798.
- 326 Willems, F. M., Shtarkov, Y. M., and Tjalkens, T. J. (1995). The context-tree weighting method:
 327 Basic properties. *IEEE Transactions on Information Theory*, 41(3):653–664.
- 328 Willems, F. M., Shtarkov, Y. M., and Tjalkens, T. J. (1996). Context weighting for general finite-
 329 context sources. *IEEE transactions on information theory*, 42(5):1514–1520.
- 330 Wu, S., Irsoy, O., Lu, S., Dabrovolski, V., Dredze, M., Gehrmann, S., Kambadur, P., Rosenberg,
 331 D., and Mann, G. (2023). Bloomberggpt: A large language model for finance. *arXiv preprint*
 332 *arXiv:2303.17564*.

- 333 Xie, S. M., Raghunathan, A., Liang, P., and Ma, T. (2021). An explanation of in-context learning as
334 implicit bayesian inference. *arXiv preprint arXiv:2111.02080*.
- 335 Ziv, J. and Lempel, A. (1977). A universal algorithm for sequential data compression. *IEEE*
336 *Transactions on information theory*, 23(3):337–343.
- 337 Ziv, J. and Lempel, A. (1978). Compression of individual sequences via variable-rate coding. *IEEE*
338 *transactions on Information Theory*, 24(5):530–536.

340 There have been many efforts in studying the ICL capabilities of transformers. A significant recent
 341 development is the elucidation of the connection to gradient descent, particularly for linear regression
 342 tasks (Von Oswald et al., 2023; Akyürek et al., 2022; Dai et al., 2022; Ahn et al., 2024). Li et al.
 343 (2023) formulated the ICL problem as a multi-task learning problem and considered ICL for several
 344 simple problem settings for which the authors provide risk bounds for ICL of supervised learning
 345 algorithms in these problem settings. Kirsch et al. (2022) viewed the ICL problem as a meta-learner
 346 and studied the relation between tasks and model sizes.

347 Olsson et al. (2022) studied the induction head, i.e., the forming of small k -gram attention in LLMs.
 348 Reddy (2023) studied the balance between ICL and in-weights learning, and observed the abrupt
 349 emergence of the induction head corresponds to the emergence of ICL. The induction head was
 350 generalized to the statistical induction head in (Edelman et al., 2024) mainly to study bigrams. We
 351 adopted it but further allowed more statistical induction heads for more suffixes to be included
 352 together, in the first two layers of the idealized transformer.

353 There have also been efforts to study transformers and learning of Markov chains. Xie et al. (2021)
 354 viewed ICL as a Bayesian inference problem, where a latent concept determines an HHM, and the
 355 observations from the HHM can lead to the identification of the hidden concept. They studied the
 356 eventual ICL capability, i.e., when the number of in-context examples goes to infinity. The work in
 357 (Bietti et al., 2024) allowed a fixed-order Markov chain to switch to a new deterministic mode, and the
 358 authors study the training behavior of the corresponding ICL task with this mode transition. Akyürek
 359 et al. (2024) made a comprehensive empirical comparison of various language models on random
 360 finite automata, and showed that the transformer performs the best among these models. Makkuva
 361 et al. (2024) studied the loss landscape during transformer training on sequences generated from a
 362 single fixed-order Markov chain, using a single-layer transformer. Their study does not consider ICL.
 363 More recently Rajaraman et al. (2024) considered ICL of FOMCs with single-head transformers, and
 364 provided a construction to show that it is possible to use a single attention head to capture longer
 365 memory in the sequence. The work most relevant to us is (Edelman et al., 2024), where ICL of a
 366 fixed-order Markov chain was considered, and the training behavior was studied both empirically and
 367 theoretically, and the forming of induction heads in a two-layer network was demonstrated. All these
 368 existing work assumed fixed-order Markov models or fixed-order HHMs, usually with orders kept
 369 at 1 or 2; moreover, they almost all focus on the emergence of the induction heads during training
 370 or the training landscape. Our study is different firstly in the variable-order nature of the Markov
 371 models, and secondly the focus on the on-time ICL performance instead of the training landscape
 372 and behavior.

373 Lossless data compression has a long history, with many different algorithms being developed over
 374 the years. The most popular general-purpose compression algorithms are perhaps the Lempel-Ziv
 375 compression algorithms (Ziv and Lempel, 1977, 1978) and their variants, which belong to dictionary-
 376 based compression algorithms. These algorithms do not explicitly maintain any probabilistic models,
 377 and their efficiency comes from maintaining an efficiency dictionary of sequences that have been seen
 378 before, and to be matched with future sequences. More powerful compression algorithms usually
 379 maintain probability models explicitly, which are then plugged into an AC module (Rissanen, 1976;
 380 Pasco, 1976; Rissanen and Langdon, 1979) for efficient compression. The most well-known classes
 381 of algorithms in this category is the context-tree weighting algorithm (Willems et al., 1995; Begleiter
 382 et al., 2004; Kontoyiannis, 2023) and prediction by partial matching (Cleary and Witten, 1984). The
 383 former enjoys a strong theoretical guarantee, particularly on binary sources (Willems et al., 1995), but
 384 has some difficulty in its practical implementation (Willems, 1998; Willems et al., 1996; Sadakane
 385 et al., 2000; Begleiter et al., 2004), particularly for large alphabet sizes and sequential data. The latter
 386 is based more on heuristics, and has been improved and extended in various ways (Cleary and Teahan,
 387 1997; Moffat, 1990; Shkarin, 2002). Methods based on probabilistic modeling are usually more
 388 resource-extensive, though they have gained more popularity recently due to the increased availability
 389 of computing resources. The evaluation given in (Begleiter et al., 2004) suggests that CTW and
 390 PPM are the two most powerful compression algorithms in practice. There are other compression
 391 algorithms such as those based on the Burrows-Wheeler transformation (Burrows, 1994) which does
 392 not explicitly maintain a probabilistic model, but are also not dictionary-based.

393 **B An Example CT**

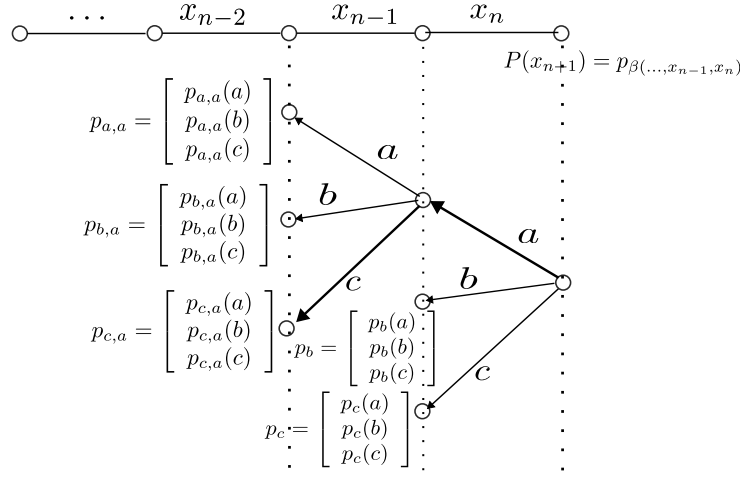


Figure 5: A CT in the alphabet $\mathcal{A} = \{a, b, c\}$ with suffix set $\mathcal{S} = \{(b), (c), (a, a), (b, a), (c, a)\}$ and the associated probability distributions. If $(\dots, x_{n-1}, x_n) = (\dots, c, a)$, then the probability distribution for the next symbol x_{n+1} is $p_{c,a}$.

394 **C Transformer Architecture**

395 The transformer considered in this is illustrated in Fig. 6.

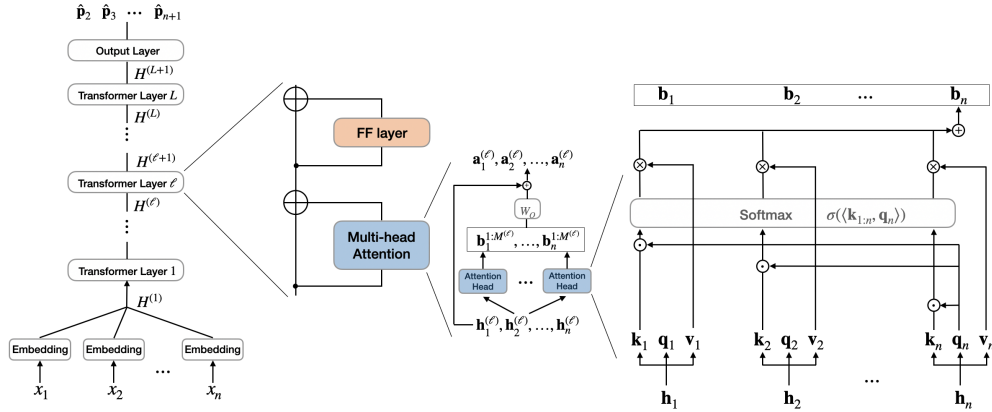


Figure 6: Transformer model

396 **D Pretraining Details**

397 We choose the alphabet size to be $|\mathcal{A}| = 3$ in the experiments. For training, we randomly generate
 398 $K = 20000$ CTs of various depths (maximum order $D \leq 5$), and then for each CT leaf, we generate a
 399 probability distribution. Two different ways of generating these probability distributions are taken: the
 400 first approach is use the Dirichlet distribution to sample such distributions, and the second approach
 401 is to randomly select some of the elements in the alphabet to have probability zero, and the others
 402 with random values. Different values of the Dirichlet parameter are tested but only the results do
 403 not appear to be sensitive to the choice. For each CT, a source sequence of certain length (e.g.,
 404 $N_k = 5120$) is produced. The context window N can vary, but in most cases, we set it at 512 (except
 405 when $D = 5$, we set it to be 1536 to allow sufficient data collection in context). Each source sequence
 406 is segmented into $\lfloor N_k/N \rfloor$ training sequence.

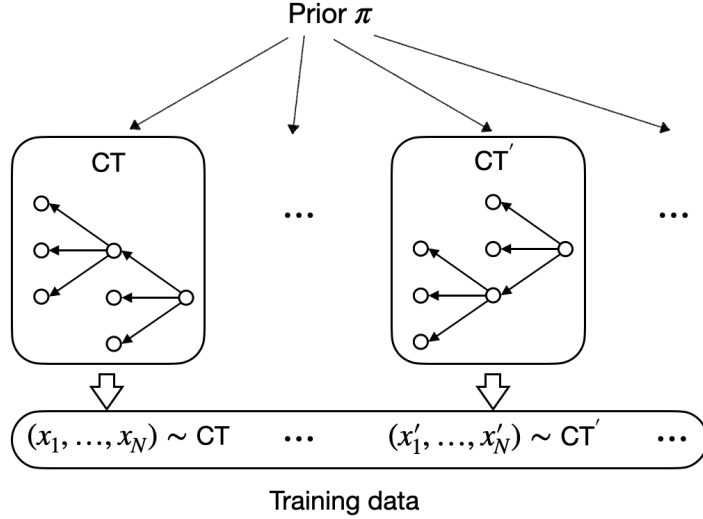


Figure 7: Training data collection

	TF-1	TF-2	TF-3	TF-4	TF-5	TF-6	CTW
CTs $D = 3$	0.9368	0.7297	0.7265	0.7220	0.7245	0.7258	0.7165
CTs $D = 4$	0.9667	0.7831	0.7818	0.7759	0.7791	0.7774	0.7603
CTs $D = 5$	0.9661	0.7569	0.7490	0.7440	0.7437	0.7438	0.7400

Table 1: Average compression rates in the context window by transformers and CTW, where the CTs are sampled from the CTW-prior. The context window and embedding dimension for CTs of $D = 5$ are $N = 1536$ and $E = 128$, while for others it are $N = 512$ and $E = 64$.

407 During testing, we randomly generate multiple (2048 in our experiments) new CTs of varying depths
 408 using the same procedure, and for each CT, a sequence of length $N_k = 5120$ are generated, and then
 409 again segmented into a length of the context window for testing.

410 The transformer model is implemented using Pytorch, and trained using the AdamW optimizer with
 411 the default parameters. A100/T100 GPUs are used for training. Training a model requires roughly
 412 4 to 6 hours. Batch size is set at 512, and the maximum epoch is set at 100 with early termination
 413 allowed after 15 epochs of no improvement. Testing was performed on a local workstation with a
 414 GeForce GTX 1660 Ti GPU card.

415 D.1 Additional experimental results

416 In Table 1, we further provide the average compression rates over the whole context window for CTs
 417 of different orders; we refer to the transformers as TF. For CTs with lower order, the transformer
 418 embedding dimension is set at 64 instead of 128.

419 D.1.1 Transformers vs. CTW under Non-CTW-Priors

420 The CTW algorithm is known to be Bayesian optimal when the CTs are generated from a CTW-prior.
 421 When the CTs do not follow those priors, can learning-based transformers perform better than CTWs?
 422 We empirically observe that in such settings, transformers indeed have advantages. The training data
 423 are generated by using CTs of different maximum orders, where the orders are chosen uniformly at
 424 random between 1 and 3. Moreover, the probability vector is not generated from the Dirichet prior,
 425 but from a distribution that for each CT leaf, randomly assigns one of the element in the alphabet to
 426 have zero-probability. We test on sequences generated from CTs produced from the same distribution
 427 as in the training setting. We assume the CTW takes the default (non-informative prior) parameters of
 428 $\alpha = 0.5$, and the same tree branch stopping parameter $\lambda = 0.15$ as taken in the testing sequence CTs.

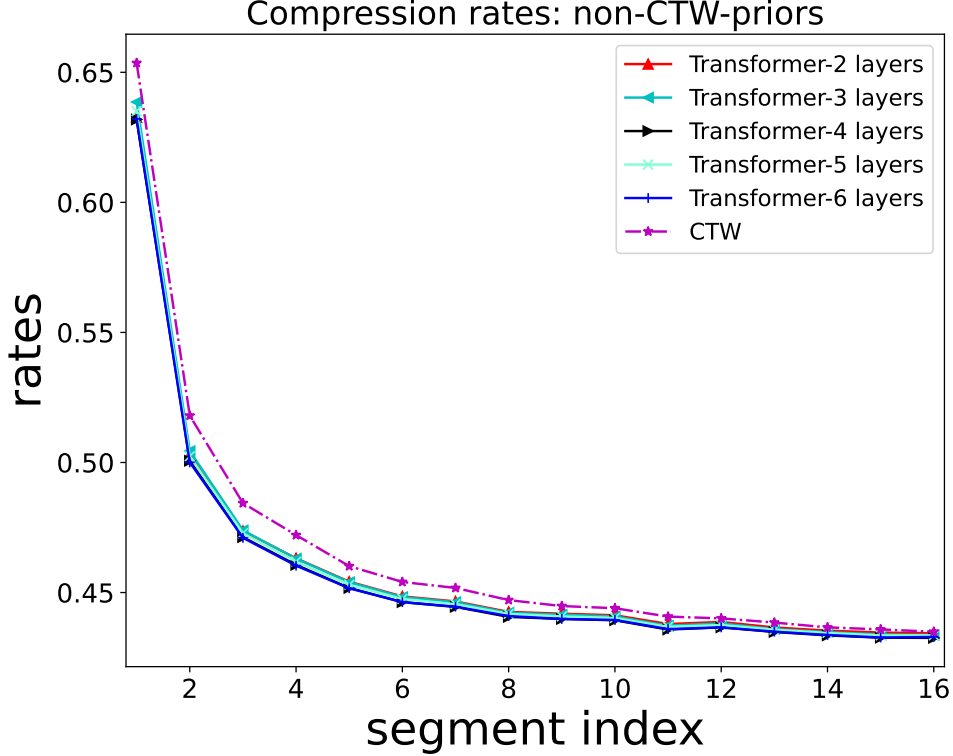


Figure 8: Transformers vs. CTW

429 As can be observed in Fig. 8, the CTW algorithm is no longer optimal, and trained transformers
 430 can perform considerably better. In fact, even transformers with 2 layers can outperform the CTW
 431 algorithm in this setting, and more layers usually lead to further improved performance, albeit the
 432 improvement is less significant.

433 D.1.2 Hybrid transformer

434 We conduct experiments on the hybrid versions of transformers. Let "TF 0-2" denote the canonical
 435 2-layer transformer; "TF 1-1" denote the transformer consisted of a constructed layer with output
 436 $\mathbf{h}_i^{(1)}$ (8), and a trainable transformer layer and a output layer taking $\mathbf{H}^{(1)}$ as input; and denote by "TF
 437 2-0" the transformer with 2 constructed layer with output $\mathbf{a}_i^{(2)}$ in (10), followed by a trainable FF
 438 layer (the FF layer in the second layer of the transformer) and an output layer.

439 We first study the key statistics behind the strong performance of two-layer transformers, as shown
 440 in the left panel in Fig. 9. Compared to "TF 2-0" which is the constructed layers given previously,
 441 the version "TF 2-0 w/o counts" does not contain $\mathbf{g}_{i-1, M'}^{\leftarrow}$ or \mathbf{pos}_i in $\mathbf{a}_i^{(2)}$; the version "TF 2-0 total
 442 counts only" does not contain $\mathbf{g}_{i-1, M'}^{\leftarrow}$ in $\mathbf{a}_i^{(2)}$ and \mathbf{pos}_i is replaced by the total count i ; "TF 2-0 w/
 443 all counts" replaces $\mathbf{g}_{i-1, M'}^{\leftarrow}$ and \mathbf{pos}_i with $\{\mathbf{n}_{n, s_{n, l}}\}_{l=0}^D$ and i . Even though their performances are
 444 rather clustered, we can make the following observations: 1) The performances degrade as more
 445 counting information is removed from the representation, and the counting information is clearly
 446 very important, 2) The performances of "TF 2-0" and "TF 2-0 w/ all counts" almost match exactly,
 447 indicating the main purpose of the backward statistics $\mathbf{g}_{i-1, M'}^{\leftarrow}$ is to extract the counts, and 3) The
 448 performance of the original 2-layer transformer is similar to that of the constructed "TF 2-0" and "TF
 449 2-0: w/ all counts" that those without less counting information.

450 We further study hybrid transformers with the first one or two being the constructed layers. As shown
 451 in the right panel of Fig. 9, transformers with 2 total layers and 4 total layers form two clusters,

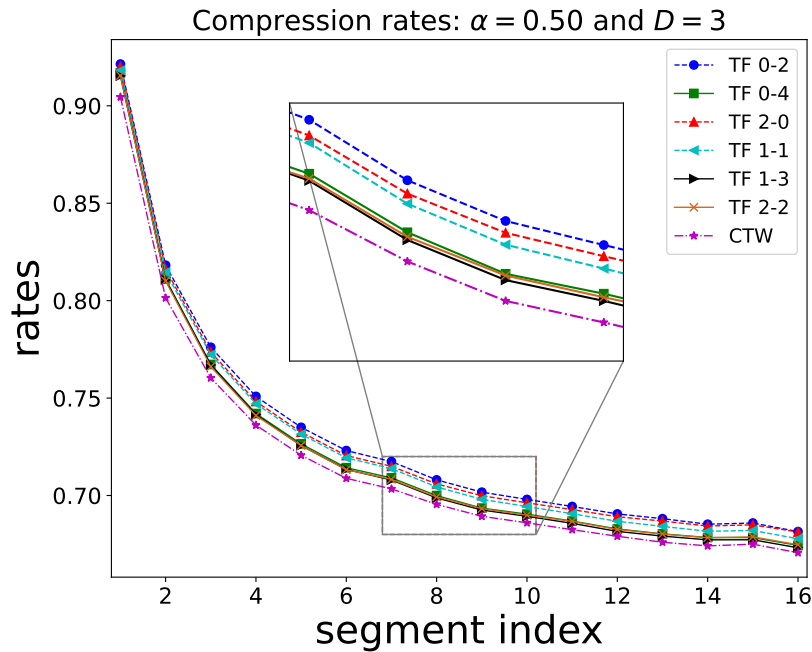
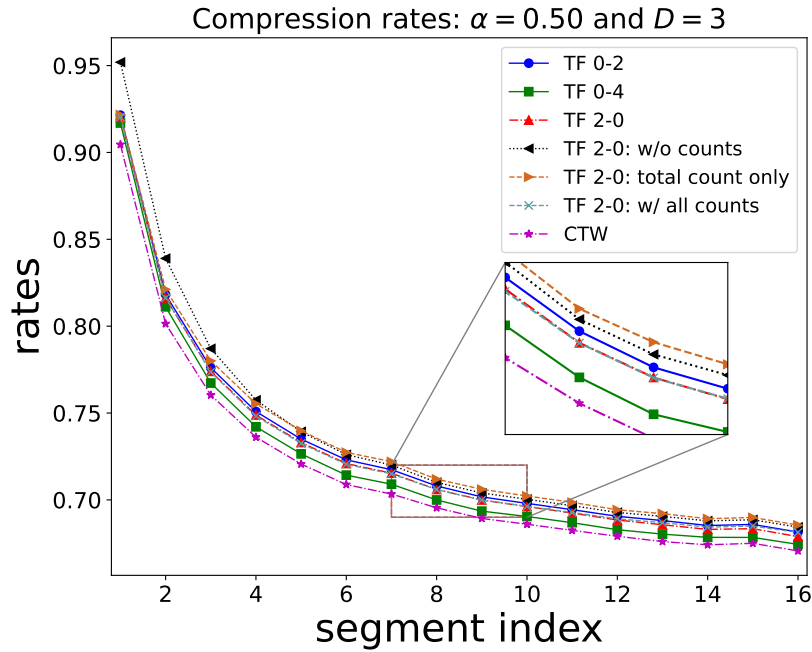


Figure 9: Hybrid Transformers: Effects of accumulative suffix counts and synthetic layers

452 which provides strong evidence that the constructed layers are indeed replacing the first two layers of
 453 the original transformers in a functional manner. Moreover, the performances of transformers with
 454 a single constructed layer, such as "TF 1-1" and "TF 1-3", are slightly better than those with two
 455 constructed layers, such as "TF 2-0" and "TF 2-2", likely due to the flexibility in the remaining
 456 trainable transformer layers. Interestingly, for two layer transformers, the hybrid versions can perform
 457 even better than the original transformer "TF 0-2", which we believe is because the latter is having
 458 difficulty extracting the exact statistics as those more readily available in the constructed layers.

459 E Proofs of The Theorems for CT Sources

460 E.1 A New Representation for Bayesian Next Token Prediction

461 We aim to predict the next token x_{n+1} based on the observations $x_{1-D}^n = (x_{1-D}, \dots, x_n)$ via a
 462 transformer-friendly formula. Note that x_{1-D}^0 is a place holder or dummy initialization sequence,
 463 which does not contain any information of $(T, \{p_s\})$.

464 **Theorem 5** (Restate Theorem 1). *The predicted probability can be computed as*

$$P_{\pi_{\text{CTW}}}(x_{n+1}|x_{1-D}^n) = \sum_{l=0, \dots, D} \omega_{n,l} \cdot \mathbf{p}_{n,s_{n,l}}(x_{n+1}), \quad (13)$$

465 where $\mathbf{p}_{n,s_{n,l}}(a) = \frac{\alpha(a) + \mathbf{n}_{n,s_{n,l}}(a)}{\sum_q (\alpha(q) + \mathbf{n}_{n,s_{n,l}}(q))}$; and $\omega_{n,\cdot} \in \Delta_{D+1}$ with $\ln(\omega_{n,l}) - \ln(\omega_{n,l-1}) = \ln(1-\lambda) -$
 467 $\mathbb{I}_{l=D} \ln(\lambda) + \ell_{n,s_{n,l}}^e - \ell_{n,s_{n,l-1}}^e + \sum_{q \in \mathcal{A}} \ell_{n,q s_{n,l-1}}^w - \ell_{n,s_{n,l}}^w$, where $\ell_{n,s}^e = \ln(p_{n,s}^e)$, $\ell_{n,s}^w = \ln(p_{n,s}^w)$.

468 Note that $p_{n,s}^e, p_{n,s}^w$ can be efficiently calculated by the CTW procedure, and compared to calculate
 469 $\frac{P_{\pi_{\text{CTW}}}(x_{n+1}^1|x_{1-D}^0)}{P_{\pi_{\text{CTW}}}(x_1^1|x_{1-D}^0)}$ for each x_{n+1} the extra computation besides the CTW procedure is A times larger
 470 than that by Eq (7). As illustrated in Fig. 4, the weighted average formula in Eq (7) gives a natural
 471 interpretation for the Bayesian optimal next token predicted probability. Each suffix along the root
 472 the leaf path $s_{n,0} - s_{n,1} - \dots - s_{n,D}$ can potentially be the true suffix, i.e., $s_{n,l} \in \mathcal{L}(T)$, and $\mathbf{p}_{n,s_{n,l}}$
 473 is in fact the Bayesian optimal next token prediction given $s_{n,l} \in \mathcal{L}(T)$.

474 The weights $\omega_{n,l}$'s are based on stopping probability λ , the information in the potential suffix path
 475 such as $p_{s_{n,s_{n,l}}}^e$ as well as the information from their siblings $p_{n,q s_{n,l-1}}^w$. We can interpret $p_{n,s}^e$ as
 476 the evidence (unnormalized likelihood) that $s \in \mathcal{L}(T)$, and $p_{n,s}^w$ as the evidence that $s \in T$, i.e.,
 477 the underlying tree covers node s . Theorem 1 indicates that more weights are assigned to $s_{n,l}$ than
 478 $s_{n,l-1}$, i.e., $\omega_{n,l} > \omega_{n,l-1}$, if λ is smaller (i.e., node $s_{n,l-1}$ is more likely to branch and thus less
 479 likely to be a leaf node), $p_{n,s_{n,l}}^e - p_{n,s_{n,l-1}}^e$ is larger (i.e., $s_{n,l}$ has more evidence than $s_{n,l-1}$) and
 480 $\sum_{q \in \mathcal{A}} \ell_{n,q s_{n,l-1}}^w - \ell_{n,s_{n,l}}^w$ is larger (i.e., $s_{n,l}$'s siblings have more evidence to explain the data and
 481 thus $s_{n,l-1}$ is less likely to be a leaf node).

482 *Proof of Theorem 5.* Recall $s_{i,l} = (x_{i-l+1}, \dots, x_i)$ is the suffix at position i of length l . We omit
 483 D by writing $\mathcal{T} = \mathcal{T}(D)$ when D is clear from the context. Define partition $\{\mathcal{T}_{s_{n,l}}\}_{0 \leq l \leq D}$, that
 484 $\mathcal{T}_s = \{T \in \mathcal{T} : s \in \mathcal{L}(T)\}$ is the set of trees with leaf s . The predicted probability can then be
 485 computed as

$$\begin{aligned} P_{\pi_{\text{CTW}}}(x_{n+1}|x_{1-D}^n) &= \sum_{T \in \mathcal{T}} \int p(x_{n+1}|T, \{p_s\}, x_{1-D}^n) \pi(T, \{p_s\}|x_{1-D}^n) \left(\prod_{s \in \mathcal{L}(T)} dp_s \right) \\ &= \sum_{l=0, \dots, D} \sum_{T \in \mathcal{T}_{s_{n,l}}} \int p_{s_{n,l}}(x_{n+1}) \pi(T, \{p_s\}|x_{1-D}^n) \left(\prod_{s \in \mathcal{L}(T)} dp_s \right) \\ &= \sum_{l=0, \dots, D} \sum_{T \in \mathcal{T}_{s_{n,l}}} \int p_{s_{n,l}}(x_{n+1}) \pi(T|x_{1-D}^n) \pi(p_{s_l}|T, x_{1-D}^n) dp_{s_l} \\ &= \sum_{l=0, \dots, D} \sum_{T \in \mathcal{T}_{s_{n,l}}} \pi_D(T|x_{1-D}^n) \int p_{s_{n,l}}(x_{n+1}) \pi(p_{s_l}|T, x_{1-D}^n) dp_{s_l} \\ &= \sum_{l=0, \dots, D} \left(\sum_{T \in \mathcal{T}_{s_{n,l}}} \pi_D(T|x_{1-D}^n) \right) \left(\int p_{s_{n,l}}(x_{n+1}) \pi(p_{s_l}|T, x_{1-D}^n) dp_{s_l} \right) \\ &= \sum_{l=0, \dots, D} \omega_{n,l} \cdot \mathbf{p}_{n,s_{n,l}}(x_{n+1}), \end{aligned} \quad (14)$$

486 where the last equality is by the definition that

$$\omega_{n,l} = \sum_{T \in \mathcal{T}_{s_{n,l}}} \pi_D(T|x_{1-D}^n), \quad (15)$$

487 and the optimal prediction probability given suffix $s_{n,l}$ is

$$\mathbf{p}_{n,s_{n,l}}(a) = \frac{\alpha(a) + \mathbf{n}_{n,s_{n,l}}(a)}{\sum_{q \in \mathcal{A}} (\alpha(q) + \mathbf{n}_{n,s_{n,l}}(q))}, \quad (16)$$

488 since for any $T \in \mathcal{T}_{s_l}$, the posterior of p_s follows Dirichlet distribution

$$\pi(p_{s_l}|T, x_{1-D}^n) = \text{Dir}(\theta_{s_l}; \alpha + \mathbf{n}_{n,s_{n,l}}), \quad (17)$$

489 with posterior mean $\mathbb{E}[p_{s_l}|T, x_{1-D}^n] \in \Delta_{\mathcal{A}}$ and $\propto \alpha + \mathbf{n}_{n,s_{n,l}}$.

490 It remains that whether the parameters $\omega_{n,l}$ is easy to compute or not. The following theorem
491 shows that these parameters $\omega_{n,l}$ can be computed easily via p_s^w and p_s^e based on x_{1-D}^n without the
492 knowledge of x_{n+1} .

493 Since the length of data n is fixed and clear from the context, let $\underline{x} = x_{1-D}^n$ be the sequence, and we
494 omit n in the subscript of $p_{n,s}^e$, $p_{n,s}^w$ and $s_{n,l}$ for simplicity.

495 For any model $T \in \mathcal{T}(D)$, the posterior probability $\pi(T|\underline{x})$ is given by:

$$\pi_D(T|\underline{x}) = \frac{\pi_D(T)P_{\pi}(\underline{x}|T)}{P_{\pi}(\underline{x})} = \frac{\pi_D(T) \prod_{s \in \mathcal{L}(T)} p_s^e}{p_{()}^w}, \quad (18)$$

496 where the denominator $P_{\pi}(\underline{x}) = p_{()}^w$ is the prior predictive likelihood computed by CTW, and
497 the numerator is by $P_{\pi}(\underline{x}|T) = \prod_{s \in \mathcal{L}(T)} p_s^e$ in (Kontoyiannis et al., 2022, Lemma 2.2). Since
498 $\omega_l = \sum_{T \in \mathcal{T}_{s_l}} \pi(T|\underline{x})$ by definition, we have for any $l = 1, 2, \dots, d$,

$$\frac{\omega_l}{\omega_{l-1}} = \frac{\sum_{T' \in \mathcal{T}_{s_l}} \pi_d(T'|x)}{\sum_{T \in \mathcal{T}_{s_{l-1}}} \pi_d(T|x)} = \frac{\sum_{T' \in \mathcal{T}_{s_l}} \pi_d(T') \prod_{s \in \mathcal{L}(T')} p_s^e}{\sum_{T \in \mathcal{T}_{s_{l-1}}} \pi_d(T) \prod_{s \in \mathcal{L}(T)} p_s^e}. \quad (19)$$

499 Note that tree in \mathcal{T}_{s_l} and trees in $\mathcal{T}_{s_{l-1}}$ share similarities. For any $T \in \mathcal{T}_{s_{l-1}}$, let $\mathcal{T}_{s_l;T} = \{T' \in$
500 $\mathcal{T}_{s_l} : \mathcal{L}(T) \subset \mathcal{L}(T') \cup \{s_{l-1}\}\}$ be the set of trees that differs from T only at subtree $\text{sub}(T'; s_l) :=$
501 $\{\text{subtree of } T' \text{ with root at } s\}$.

502 Take any $l = 1, 2, \dots, D - 1$. For any $T \in \mathcal{T}_{s_{l-1}}$ and $T' \in \mathcal{T}_{s_l;T}$. Based on the definition of
503 $\pi_D = (1 - \lambda)^{(|\mathcal{L}(T)|-1)/(A-1)} \lambda^{|\mathcal{L}(T)|-|\mathcal{L}_D(T)|}$, it is not hard to verify that

$$\begin{aligned} \frac{\pi_D(T')}{\pi_D(T)} &= \frac{\pi_{D-l+1}(\text{sub}(T'; s_{l-1}))}{\pi_{D-l+1}(\text{sub}(T; s_{l-1}))} \\ &= \frac{(1 - \lambda) \pi_{D-l}(\text{sub}(T'; s_l)) \prod_{s'_l \in \text{sib}(s_l)} \pi_{D-l}(\text{sub}(T'; s'_l))}{\lambda} \\ &= (1 - \lambda) \prod_{s'_l \in \text{sib}(s_l)} \pi_{D-l}(\text{sub}(T'; s'_l)), \end{aligned}$$

504 where $\text{sib}(s_{l+1}) = \{q_{s_l} : q \in \mathcal{A} \text{ and } q_{s_l} \neq s_{l+1}\}$ is set of siblings of s_{l+1} . We can interpret the ratio
505 as follows. T' and T only differs at the $\text{sub}(T'; s_{l-1})$ and $\text{sub}(T; s_{l-1})$. Since T' branch at node
506 s_{l-1} , we thus have the numerator in the second equation, where $(1 - \lambda)$ corresponds to the branching
507 and then compute for the subtrees. Note that T stops branching at s_{l-1} and T' stops branching at s_l ,
508 then $\pi_{D-l+1}(\text{sub}(T; s_{l-1})) = \pi_{D-l}(\text{sub}(T'; s_l)) = \lambda$ equals to the stopping probability.

509 Given any suffix s with $|s| \leq D$, it has been shown in (Kontoyiannis et al., 2022, Proof of Theorem
510 3.1) that for any $l \leq D$,

$$p_s^w = \sum_{U \in \mathcal{T}(D-l)} \pi_{D-l}(U) \prod_{u \in \mathcal{L}(U)} p_{us}^e, \quad (20)$$

511 where $\mathcal{T}(D-l)$ is the set of trees with maximum depth $D-l$ and π_{D-l} is the prior for bounded
 512 branching process with maximum depth $D-l$. We thus have

$$\frac{\sum_{T' \in \mathcal{T}_{s_l; T}} \pi_D(T') \prod_{s \in \mathcal{L}(T')} p_s^e}{\pi_D(T) \prod_{s \in \mathcal{L}(T)} p_s^e} = \frac{\sum_{T' \in \mathcal{T}_{s_l; T}} \pi_D(T') \prod_{s \in \mathcal{L}(T')} p_s^e}{\pi_D(T) \prod_{s \in \mathcal{L}(T)} p_s^e} \quad (21)$$

$$= \sum_{T' \in \mathcal{T}_{s_l; T}} \frac{\pi_D(T') \prod_{s \in \mathcal{L}(T') \setminus \mathcal{L}(T)} p_s^e}{\pi_D(T) p_{s_{l-1}}^e} \quad (22)$$

$$= \sum_{T' \in \mathcal{T}_{s_l; T}} \left((1-\lambda) \prod_{s'_l \in \text{sib}(s_l)} \pi_{D-l}(\text{sub}(T'; s'_l)) \right) \left(\frac{p_{s_l}^e \prod_{s'_l \in \text{sib}(T; s_l)} \prod_{s \in \mathcal{L}(\text{sub}(T'; s'_l))} p_s^e}{p_{s_{l-1}}^e} \right) \quad (23)$$

$$= (1-\lambda) \frac{p_{s_l}^e}{p_{s_{l-1}}^e} \sum_{T' \in \mathcal{T}_{s_l; T}} \left(\prod_{s'_l \in \text{sib}(s_l)} \pi_{D-l}(\text{sub}(T'; s'_l)) \right) \left(\prod_{s'_l \in \text{sib}(T; s_l)} \prod_{s \in \mathcal{L}(\text{sub}(T'; s'_l))} p_s^e \right) \quad (24)$$

$$= (1-\lambda) \frac{p_{s_l}^e}{p_{s_{l-1}}^e} \sum_{T' \in \mathcal{T}_{s_l; T}} \left(\prod_{s'_l \in \text{sib}(s_l)} \pi_{D-l}(\text{sub}(T'; s'_l)) \prod_{s \in \mathcal{L}(\text{sub}(T'; s'_l))} p_s^e \right) \quad (25)$$

$$= (1-\lambda) \frac{p_{s_l}^e}{p_{s_{l-1}}^e} \prod_{s'_l \in \text{sib}(s_l)} \left(\sum_{U \in \mathcal{T}(D-l)} \pi_{D-l}(U) \prod_{u \in \mathcal{L}(U)} p_{us'_l}^e \right) \quad (26)$$

$$= \frac{(1-\lambda) p_{s_l}^e \prod_{a \neq s_l} p_{a s_{l-1}}^w}{p_{s_{l-1}}^e}. \quad (27)$$

513 Similarly, for any $T \in \mathcal{T}_{s_{D-1}}$ and $T' \in \mathcal{T}_{s_D; T}$, $\frac{\pi_D(T')}{\pi_D(T)} = \frac{1-\lambda}{\lambda}$, and we have

$$\frac{\omega_D}{\omega_{D-1}} = \frac{(1-\lambda) p_{s_D}^e \prod_{a \neq s_D} p_{a s_{D-1}}^w}{\lambda p_{s_{D-1}}^e}, \quad (28)$$

514 in the same manner. The proof can then be concluded by taking logarithm on both hands. \square

515 E.2 Construction of Transformer for CTW

516 To make the presentation clear, in the following we separate the layers by their functionality and
 517 present them separately. Recall that

$$\mathbf{a}_i^{(\ell)} = \text{MHA} \left(\mathbf{h}_i, \mathbf{H}; \{W_{O,m}^{(\ell)}, W_{Q,m}^{(\ell)}, W_{K,m}^{(\ell)}, W_{V,m}^{(\ell)}\}_{m=1}^{M^{(\ell)}} \right) \triangleq W_O^{(\ell)} \left[\mathbf{b}_{1,i}^{(\ell)}; \mathbf{b}_{2,i}^{(\ell)}; \dots; \mathbf{b}_{M^{(\ell)},i}^{(\ell)} \right],$$

518 where $\{W_{Q,m}^{(\ell)}, W_{K,m}^{(\ell)}, W_{V,m}^{(\ell)}\}_{m=1}^{M^{(\ell)}}$ are the $E^{(\ell)} \times E$ query matrices, key matrices, and value matrices
 519 and $W_O^{(\ell)}$ is the $E \times M^{(\ell)} E^{(\ell)}$ output mapping matrix. For simplicity of presentation, we take
 520 $E^\ell = E$ and $W_O^\ell = [\mathbf{I}; \mathbf{I}; \dots; \mathbf{I}]$. It is not hard to see the following constructions can be applied to
 521 much smaller $E^{(\ell)}$ while taking W_O as a permutation matrix.

522 We have omitted the dimensionality of several zero matrices when they are obvious from the context.
 523 The first and second layer constructions are illustrated in Fig. 10.

524 E.2.1 Finite-memory context-extension layer

525 We begin with the first layer, which is referred to as a finite-memory context-extension layer.

526 **Theorem 6** (Restatement of Theorem 2). *There is an M -headed transformer layer that can perform*
 527 *finite-memory context-extension, defined by the following output, with the initial one-hot embedded*
 528 *input $\mathbf{H}^{(1)}$:*

$$\mathbf{h}_i^{(2)} = (\mathbf{s}_{i, M+1}; \mathbf{0}; \mathbf{pos}_i) = (\mathbf{x}_i; \mathbf{x}_{i-1}; \dots; \mathbf{x}_{i-M}; \mathbf{0}; \mathbf{pos}_i), \quad (29)$$

529 where $\mathbf{s}_{i, M+1} = (\mathbf{x}_i; \dots; \mathbf{x}_{i-M})$ is the vector version of the M -length suffix $s_{i, M+1} = x_{i-M}^i$.

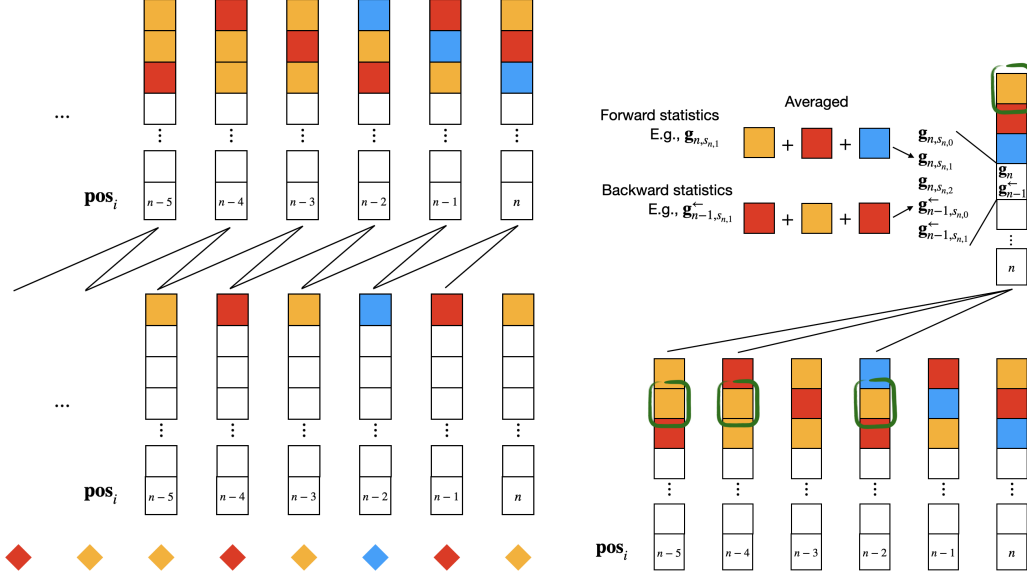


Figure 10: Transformer construction for $D = 2$. The left figure illustrates the first layer – finite-memory context-extension layer, which append the previous D tokens. The right figure demonstrate the MHA of the second layer – statistics collection layer, which extracts forward and backward statistics based on the matched suffix.

530 *Proof of Theorem 2.* The input of the of the first layer is a initial one-hot embedded input with
 531 positional embedding $\mathbf{H}^{(1)}$, where its n -th column is

$$\mathbf{h}_i^{(1)} = (\mathbf{x}_i; \mathbf{0}; \mathbf{pos}_i) \in \mathbb{R}^E, \quad (30)$$

532 where positional encoding

$$\mathbf{pos}_i = (1; \cos(i\pi/N); \sin(i\pi/N)), \quad (31)$$

533 with C being the maximum context size.

534 The multi-head attention in the first layer is consisted of $M^{(1)} = D$ heads parameterized by
 535 $(W_{Q,m}^{(1)}, W_{K,m}^{(1)}, W_{V,m}^{(1)})_{m=1,2,\dots,M^{(1)}}$. Specifically, for $m = 1, 2, \dots, M^{(1)}$,

$$W_{Q,m}^{(1)} = \begin{pmatrix} \mathbf{0} & \text{Rot}(m) \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad W_{K,m}^{(1)} = \begin{pmatrix} \mathbf{0} & c\mathbf{I}^{2 \times 2} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad W_{V,m}^{(1)} = \begin{pmatrix} \mathbf{0}^{mA \times A} & \mathbf{0} \\ \mathbf{I}^{A \times A} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad (32)$$

536 where $\text{Rot}(m) = \begin{pmatrix} \cos(m\pi/N) & \sin(m\pi/N) \\ -\sin(m\pi/N) & \cos(m\pi/N) \end{pmatrix}$ is a rotation matrix that rotates clockwise by an
 537 angle of $m\pi/C$, and $c \in \mathbb{R}_+$ is a temperature factor. The query, key, and value after the mapping are

$$W_{Q,m}^{(1)} \mathbf{h}_n^{(1)} = \begin{pmatrix} \mathbf{pos}_{n-m} \\ \mathbf{0} \end{pmatrix}, \quad W_{K,m}^{(1)} \mathbf{h}_i^{(1)} = c \begin{pmatrix} \mathbf{pos}_i \\ \mathbf{0} \end{pmatrix}, \quad W_{V,m}^{(1)} \mathbf{h}_i^{(1)} = \begin{pmatrix} \mathbf{0}^{mA \times 1} \\ \mathbf{x}_i \\ \mathbf{0} \end{pmatrix}. \quad (33)$$

538 Take $c = \infty$ or sufficiently large. It is seen that the m -th head essentially copies the m -th earlier
 539 symbol to stack at the $(m + 1)$ -th position below the original symbol \mathbf{x}_i . Together with the residual
 540 link, the attention layer gives exactly the $\mathbf{h}_i^{(2)}$ shown in (34) while the feedforward network in this
 541 layer can be set as all zero.

$$\mathbf{h}_i^{(2)} = (\mathbf{x}_i; \mathbf{x}_{i-1}; \mathbf{x}_{i-2}; \mathbf{x}_{i-M^{(1)}}; \mathbf{0}; \mathbf{pos}_i) = (\mathbf{s}_{i, M^{(1)}+1}; \mathbf{0}; \mathbf{pos}_i), \quad (34)$$

542 where $\mathbf{s}_{i,l} = (\mathbf{x}_i; \mathbf{x}_{i-1}; \dots; \mathbf{x}_{i-l+1})$ is the one-hot embedded version of suffix $s_{i,l} =$
 543 $(x_{i-l+1}, \dots, x_{i-1}, x_i)$. \square

544 **E.2.2 Statistics collection layer**

545 **Theorem 7** (Restatement of Theorem 3). *There is an M' -head attention layer, where $M' \leq M + 1$,*
 546 *that can perform statistics collection, defined by the following output, with $\mathbf{H}^{(2)}$ in (8) as its input:*

$$\mathbf{a}_i^{(2)} = (\mathbf{s}_{i,M+1}; \mathbf{g}_{i,M'}; \mathbf{g}_{i-1,M'}^{\leftarrow}; \mathbf{0}; \mathbf{pos}_i), \quad (35)$$

547 where $\mathbf{g}_{i,M'} := (\mathbf{g}_{i,s_{i,0}}; \dots; \mathbf{g}_{i,s_{i,M'-1}})$ and $\mathbf{g}_{i-1,M'}^{\leftarrow} = (\mathbf{g}_{i-1,s_{i,0}}^{\leftarrow}; \dots; \mathbf{g}_{i-1,s_{i,M'-1}}^{\leftarrow})$.

548 *Proof of Theorem 3.* To make the proof self-contained, we first recall some key notations. The
 549 second layer is referred to as the statistics collection layer, which uses a sequence of vectors $\mathbf{h}_i^{(2)}$,
 550 $i = 1, 2, \dots, N$, defined in (8) as its input, restated as follows.

$$\mathbf{h}_i^{(2)} = (\mathbf{s}_{i,M+1}; \mathbf{0}; \mathbf{pos}_i), \quad (36)$$

551 where $\mathbf{s}_{i,M+1} = (\mathbf{x}_i; \dots; \mathbf{x}_{i-M})$. To rigorously specify the function of this layer, recall the definition
 552 of the k -gram statistics vector $\mathbf{g}_{i,s}$, which in plain words, is the empirical probability distribution of
 553 the next token associated with the suffix s for a sequence x_1^i . Mathematically, for a suffix s whose
 554 length is $k - 1$ and the current position i ,

$$\mathbf{g}_{i,s}(a) = \frac{\mathbf{n}_{i,s}(a)}{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s}(q)} \quad \forall a \in \mathcal{A}, \quad (37)$$

555 where $\mathbf{n}_{i,s}$ is the counting vector defined in (5).

556 The k -gram backward statistics vector $\mathbf{g}_{i-1,s}^{\leftarrow}$ is defined similarly, which is the empirical probability
 557 distribution of the previous token associated with the suffix s for data x_1^{i-1} , and mathematically

$$\mathbf{g}_{i-1,s}^{\leftarrow}(a) = \frac{\sum_{q \in \mathcal{A}} \mathbf{n}_{i-1,s}(q)}{\sum_{q \in \mathcal{A}} \mathbf{n}_{i-1,s}(q)} \quad \forall a \in \mathcal{A}, \quad (38)$$

558 where $\sum_{q \in \mathcal{A}} \mathbf{n}_{i-1,s}(q)$ is the number of appears of the sub-string s in the sequence x_1^{i-1} .

559 The multi-head attention in the second layer is consisted of $M^{(2)} = M' \leq M^{(1)} + 1 = M + 1$ heads
 560 parameterized by $(W_{Q,m}^{(2)}, W_{K,m}^{(2)}, W_{V,m}^{(2)})_{m=0,1,2,\dots,M^{(2)}-1}$. Specifically, for $m = 1, 2, \dots, M^{(2)} - 1$,

$$W_{Q,m}^{(2)} = \begin{pmatrix} \mathbf{I}^{(m-1)A \times (m-1)A} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad W_{K,m}^{(2)} = \begin{pmatrix} \mathbf{0}^{(m-1)A \times A} & c\mathbf{I}^{(m-1)A \times (m-1)A} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad (39)$$

$$W_{V,m}^{(2)} = \begin{pmatrix} \mathbf{0}^{(M^{(1)}+m)A \times A} & \mathbf{0} \\ \mathbf{I}^{A \times A} & \mathbf{0} \\ \mathbf{0}^{(M^{(2)}-1)A \times A} & \mathbf{0} \\ \mathbf{0}^{A \times A} & [\mathbf{0}^{A \times (m-1)A}, \mathbf{I}^{A \times A}, \mathbf{0}] \\ \mathbf{0} & \mathbf{0} \end{pmatrix}. \quad (40)$$

561 The corresponding query, key, and value vectors after the mapping are

$$W_{Q,m}^{(2)} \mathbf{h}_n^{(2)} = \begin{pmatrix} \mathbf{s}_{n,m-1} \\ \mathbf{0} \end{pmatrix}, \quad W_{K,m}^{(2)} \mathbf{h}_i^{(2)} = c \begin{pmatrix} \mathbf{s}_{i-1,m-1} \\ \mathbf{0} \end{pmatrix}, \quad W_{V,m}^{(2)} \mathbf{h}_i^{(2)} = \begin{pmatrix} \mathbf{0}^{(M^{(1)}+m)A \times 1} \\ \mathbf{x}_i \\ \mathbf{0}^{(M^{(2)}-1)A \times 1} \\ \mathbf{x}_{i-m} \\ \mathbf{0} \end{pmatrix}.$$

562 For $m = M^{(2)}$, $W_{Q,m}^{(2)}, W_{K,m}^{(2)}$ are of the same structure, while $W_{V,m}^{(2)}$ does not contains that $\mathbf{I}^{A \times A}$ in
 563 that $[\mathbf{0}^{A \times (m-1)A}, \mathbf{I}^{A \times A}, \mathbf{0}]$ block, and thus $W_{V,m}^{(2)} \mathbf{h}_i^{(2)}$ does not have \mathbf{x}_{i-m} .

564 It is not hard to see that taking $c \rightarrow \infty$ gives

$$(\mathbf{s}_{i,M^{(1)}+1}; \mathbf{g}_{i,M^{(2)}-1}; \mathbf{g}_{i-1,M^{(2)}-1}^{\leftarrow}; \mathbf{0}; \mathbf{pos}_i) = [\text{MHA}(\mathbf{H}^{(2)}) + \mathbf{H}^{(2)}]_i, \quad (41)$$

565 where

$$\mathbf{g}_{i,M'} = (\mathbf{g}_{i,s_{i,0}}; \dots; \mathbf{g}_{i,s_{i,M'-1}}) \\ \mathbf{g}_{i-1,M'}^{\leftarrow} = (\mathbf{g}_{i-1,s_{i,0}}^{\leftarrow}; \dots; \mathbf{g}_{i-1,s_{i,M'-1}}^{\leftarrow}).$$

566 □

567 Note that the counting vector can be obtained via

$$\mathbf{n}_{i,s_{i,l}}(a) = \frac{\mathbf{n}_{i,s_{i,l}}(a)}{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s_{i,l}}(q)} \frac{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s_{i,l}}(q)}{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s_{i,l-1}}(q)} \cdots \frac{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s_{i,1}}(q)}{\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s_{i,0}}(q)} \left(\sum_{q \in \mathcal{A}} \mathbf{n}_{i,s_{i,0}}(q) \right) \quad (42)$$

$$= \mathbf{g}_{i,s_{i,l}}(a) \left(\prod_{j=0}^{l-1} \mathbf{g}_{i-1,s_{i,j}}^{\leftarrow}(x_{i-j}) \right) \cdot i, \quad (43)$$

568 by the information contained in vector $(\mathbf{s}_{i,M^{(1)}+1}; \mathbf{g}_{i,M^{(2)}-1}; \mathbf{g}_{i-1,M^{(2)}-1}^{\leftarrow}; \mathbf{0}; \mathbf{pos}_i)$.

569 Since $p_{i,s_{i,l}}^e$ and $\mathbf{p}_{i,s_{i,l}}$ in (16) are functions of $\mathbf{n}_{i,s_{i,l}}$, we can then obtain (approximate) the following
570 output by a sufficiently wide FF layer that

$$\mathbf{h}_i^3 = (\mathbf{s}_{i,M^{(1)}+1}; \mathbf{p}_{i,D}; \mathbf{l}_{i,D}^e; \ln(p_{i,s_{i,D}}^w); \mathbf{0}; \mathbf{pos}_i), \quad (44)$$

571 where $\mathbf{l}_{i,D}^e$ contains the logarithm of p^e along the path from root $()$ to (x_{i-d+1}, \dots, x_i) , and $\mathbf{p}_{i,D}$
572 stacks the optimal prediction given suffices $s_{i,0}, \dots, s_{i,D}$, i.e.,

$$\mathbf{l}_{i,D}^e = (\ell_{i,s_{i,0}}^e; \ell_{i,s_{i,1}}^e; \dots; \ell_{i,s_{i,D}}^e) = (\ln(p_{i,s_{i,0}}^e); \ln(p_{i,s_{i,1}}^e); \dots; \ln(p_{i,s_{i,D}}^e)), \quad (45)$$

$$\mathbf{p}_{i,D} = (\mathbf{p}_{i,s_{i,0}}; \mathbf{p}_{i,s_{i,1}}; \dots; \mathbf{p}_{i,s_{i,D}}), \quad (46)$$

573 and $\ln(p_{i,s_{i,D}}^w) = \ln(p_{i,s_{i,D}}^e)$ with suffix $|s_{i,D}| = D$. These quantities can be extracted, since they
574 are functions of the statistics collected from $\mathbf{a}_i^{(2)}$.

575 This functional layer essentially collects k -gram statistics for various lengths of $k = 1, 2, \dots, M^{(2)}$
576 via multi-head attention and then process the the statistics for follow-up optimal scheme.

577 E.2.3 Inductive CTW layer

578 Recall the input and the expected outputs of the inductive CTW layer that

$$\mathbf{h}_i^{(\ell)} = (\mathbf{s}_{i,M^{(1)}+1}; \mathbf{p}_{i,D}; \mathbf{l}_{i,D}^e; \delta_{i,D}; \delta_{i,D-1}; \dots; \delta_{i,D-\ell+4}; \ell_{i,s_{i,D+3-\ell}}^w; \mathbf{0}; \mathbf{pos}_i), \quad (47)$$

579 for $\ell = 3, 4, \dots, 3+D$, where $\delta_{i,l} := \ln(\omega_{i,l}) - \ln(\omega_{i,l-1})$ for $l = d, D-1, \dots, 1$ are the the weight
580 difference, and we take $M^{(1)} = D$.

581 **Theorem 8** (Restatement of Theorem 4). *There exists a A -head transformer layer that can perform
582 the induction: Takes $\mathbf{H}^{(\ell)}$ in (47) as input and outputs $\mathbf{H}^{(\ell+1)}$. And the final readout layer taking
583 $\mathbf{H}^{(D+2)}$ as input can output the A -dimensional Bayesian optimal next token prediction vector
584 $P_{\pi_{CTW}}(\cdot | x_{1-D}^n) = \sum_{l=0, \dots, D} \omega_{n,l} \mathbf{P}_{n,s_{n,l}}$.*

585 *Proof of Theorem 4.* For any fixed $\ell = 3, 4, \dots, 2+D$, we specify the construction for the ℓ -th
586 transformer layer. It contains A heads and for each $m = 1, 2, \dots, A$, the Q, K, V matrices are

$$W_{Q,m}^{(\ell)} = \begin{pmatrix} \mathbf{I}^{(D+1-\ell)A \times (D+1-\ell)A} & \mathbf{0} \\ \mathbf{0} & [\mathbf{e}_m, \mathbf{0}^{A \times 2}] \\ \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}^{2 \times 2} \end{pmatrix}, \quad W_{K,m}^{(\ell)} = \begin{pmatrix} c\mathbf{I}^{(D+2-\ell)A \times (D+2-\ell)A} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \\ \mathbf{0} & c\mathbf{I}^{2 \times 2} \end{pmatrix},$$

$$W_{V,m}^{(\ell)} = \begin{pmatrix} \mathbf{0}^{(\text{place}_\ell+m) \times (\text{place}_\ell+m)} & \mathbf{0} \\ [\mathbf{0}^{1 \times (\text{place}_\ell-1)}, \mathbf{1}] & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix},$$

587 where \mathbf{e}_m is the A -dimensional one-hot vector at position m , and $\text{place}_\ell = (M^{(1)}+D+2)A+D+\ell-1$
588 is index of element $\ell_{i,s_{i,D+3-\ell}}^w$ in $\mathbf{h}_i^{(\ell)}$. The corresponding query, key, and value vectors after the
589 mapping are

$$W_{Q,m}^{(\ell)} \mathbf{h}_n^{(\ell)} = \begin{pmatrix} \mathbf{s}_{n,D+1-\ell} \\ \mathbf{e}_m \\ \mathbf{0} \\ \mathbf{pos}_n \end{pmatrix}, \quad W_{K,m}^{(\ell)} \mathbf{h}_i^{(\ell)} = c \begin{pmatrix} \mathbf{s}_{i,D+2-\ell} \\ \mathbf{0} \\ \mathbf{pos}_i \end{pmatrix}, \quad W_{V,m}^{(\ell)} \mathbf{h}_i^{(\ell)} = \begin{pmatrix} \mathbf{0}^{(\text{place}_\ell+m) \times 1} \\ \ell_{i,s_{i,D+3-\ell}}^w \\ \mathbf{0} \end{pmatrix}.$$

590 At position n , the query of m -head will select the latest (due to positional embedding) position with
 591 suffix $[\mathbf{s}_{n,D+1-\ell}; \mathbf{e}_m]$, and append its ℓ^w at the end. It is not hard to see that taking $c \rightarrow \infty$ gives

$$\begin{aligned} \mathbf{a}_i^{(\ell)} &= [\text{MHA}(\mathbf{H}^{(2)}) + \mathbf{H}^{(2)}]_i \\ &= (\mathbf{s}_{i,D+1}; \mathbf{p}_{i,D}; \mathbf{l}_{i,D}^e; \delta_{i,D}; \delta_{i,D-1}; \dots; \delta_{i,D+4-\ell}; \ell_{i,s_i,D+3-\ell}^w; [\ell_{i,q s_i,D+2-\ell}^w]_{q \in \mathcal{A}}; \mathbf{0}; \mathbf{pos}_i) \end{aligned}$$

592 Recall $\ln(\omega_{n,l}) - \ln(\omega_{n,l-1}) = \ln(1-\lambda) - \mathbb{1}_{l=D} \ln(\lambda) + \ell_{n,s_n,l}^e - \ell_{n,s_n,l-1}^e + \sum_{q \in \mathcal{A}} \ell_{n,q s_n,l-1}^w - \ell_{n,s_n,l}^w$
 593 by Theorem 1. $\delta_{i,D+3-\ell} = \ln(\omega_{i,D+3-\ell}) - \ln(\omega_{i,D+2-\ell})$ can be computed by $\mathbf{a}_i^{(\ell)}$ and thus $\mathbf{h}_i^{(\ell+1)}$
 594 can be approximated via the FF layer following the ℓ -th multi-head attention layer.

595 The final layer approximate an A -dimensional vector

$$P_{\pi_{\text{CTW}}}(\cdot | x_{1-D}^n) = \sum_{l=0, \dots, D} \omega_{n,l} \cdot \mathbf{p}_{n,s_n,l}(\cdot), \quad (48)$$

596 by an FF layer taking input

$$\mathbf{h}_n^{(D+3)} = (\mathbf{s}_{n,M^{(1)}+1}; \mathbf{p}_{n,D}; \mathbf{l}_{n,D}^e; \delta_{n,D}; \dots; \delta_{n,1}; \mathbf{0}; \mathbf{pos}_i). \quad (49)$$

597 The proof is now complete. \square

598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649

1. Claims

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

Answer: [Yes] .

Justification: Contributions and claims are clearly stated and match the theory and experimental results of the 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: See the conclusion.

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: Yes, the assumptions and proofs are clearly stated.

650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703

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: Yes, the experimental setup is described in the appendix.

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: [\[No\]](#) .

704 Justification: Source code will be released after paper acceptance.

705 Guidelines:

- 706 • The answer NA means that paper does not include experiments requiring code.
- 707 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 708
- 709 • While we encourage the release of code and data, we understand that this might not be
- 710 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
- 711 including code, unless this is central to the contribution (e.g., for a new open-source
- 712 benchmark).
- 713 • The instructions should contain the exact command and environment needed to run to
- 714 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 715
- 716 • The authors should provide instructions on data access and preparation, including how
- 717 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 718 • The authors should provide scripts to reproduce all experimental results for the new
- 719 proposed method and baselines. If only a subset of experiments are reproducible, they
- 720 should state which ones are omitted from the script and why.
- 721 • At submission time, to preserve anonymity, the authors should release anonymized
- 722 versions (if applicable).
- 723 • Providing as much information as possible in supplemental material (appended to the
- 724 paper) is recommended, but including URLs to data and code is permitted.

725 6. Experimental Setting/Details

726 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-

727 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the

728 results?

729 Answer: [Yes] .

730 Justification: Yes, details are given in the appendix

731 Guidelines:

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

737 7. Experiment Statistical Significance

738 Question: Does the paper report error bars suitably and correctly defined or other appropriate

739 information about the statistical significance of the experiments?

740 Answer: [Yes] .

741 Justification: Yes, error bars are included in the result

742 Guidelines:

- 743 • The answer NA means that the paper does not include experiments.
- 744 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
- 745 dence intervals, or statistical significance tests, at least for the experiments that support
- 746 the main claims of the paper.
- 747 • The factors of variability that the error bars are capturing should be clearly stated (for
- 748 example, train/test split, initialization, random drawing of some parameter, or overall
- 749 run with given experimental conditions).
- 750 • The method for calculating the error bars should be explained (closed form formula,
- 751 call to a library function, bootstrap, etc.)
- 752 • The assumptions made should be given (e.g., Normally distributed errors).
- 753 • It should be clear whether the error bar is the standard deviation or the standard error
- 754 of the mean.

- 755 • It is OK to report 1-sigma error bars, but one should state it. The authors should
756 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
757 of Normality of errors is not verified.
- 758 • For asymmetric distributions, the authors should be careful not to show in tables or
759 figures symmetric error bars that would yield results that are out of range (e.g. negative
760 error rates).
- 761 • If error bars are reported in tables or plots, The authors should explain in the text how
762 they were calculated and reference the corresponding figures or tables in the text.

763 8. Experiments Compute Resources

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

767 Answer: [Yes] .

768 Justification: The compute resources are described in the appendix.

769 Guidelines:

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

778 9. Code Of Ethics

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

781 Answer: [Yes] .

782 Justification: Nothing to justify.

783 Guidelines:

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

789 10. Broader Impacts

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

792 Answer: [NA] .

793 Justification: It is mainly a technical study without foreseeable social impacts.

794 Guidelines:

- 795 • The answer NA means that there is no societal impact of the work performed.
- 796 • If the authors answer NA or No, they should explain why their work has no societal
797 impact or why the paper does not address societal impact.
- 798 • Examples of negative societal impacts include potential malicious or unintended uses
799 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
800 (e.g., deployment of technologies that could make decisions that unfairly impact specific
801 groups), privacy considerations, and security considerations.
- 802 • The conference expects that many papers will be foundational research and not tied
803 to particular applications, let alone deployments. However, if there is a direct path to
804 any negative applications, the authors should point it out. For example, it is legitimate
805 to point out that an improvement in the quality of generative models could be used to

806 generate deepfakes for disinformation. On the other hand, it is not needed to point out
807 that a generic algorithm for optimizing neural networks could enable people to train
808 models that generate Deepfakes faster.

- 809 • The authors should consider possible harms that could arise when the technology is
810 being used as intended and functioning correctly, harms that could arise when the
811 technology is being used as intended but gives incorrect results, and harms following
812 from (intentional or unintentional) misuse of the technology.
- 813 • If there are negative societal impacts, the authors could also discuss possible mitigation
814 strategies (e.g., gated release of models, providing defenses in addition to attacks,
815 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
816 feedback over time, improving the efficiency and accessibility of ML).

817 11. Safeguards

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

821 Answer: [NA] .

822 Justification: The model is small, and is not expected to lead to any risk.

823 Guidelines:

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

834 12. Licenses for existing assets

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

838 Answer: [NA] .

839 Justification: Nothing to report

840 Guidelines:

- 841 • The answer NA means that the paper does not use existing assets.
- 842 • The authors should cite the original paper that produced the code package or dataset.
- 843 • The authors should state which version of the asset is used and, if possible, include a
844 URL.
- 845 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 846 • For scraped data from a particular source (e.g., website), the copyright and terms of
847 service of that source should be provided.
- 848 • If assets are released, the license, copyright information, and terms of use in the
849 package should be provided. For popular datasets, paperswithcode.com/datasets
850 has curated licenses for some datasets. Their licensing guide can help determine the
851 license of a dataset.
- 852 • For existing datasets that are re-packaged, both the original license and the license of
853 the derived asset (if it has changed) should be provided.
- 854 • If this information is not available online, the authors are encouraged to reach out to
855 the asset's creators.

856 13. New Assets

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

859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903

Answer: [NA] .

Justification: Nothing to report

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: Nothing to report

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: Nothing to report

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.