

# HOW TRANSFORMERS LEARN CAUSAL STRUCTURES IN-CONTEXT: EXPLAINABLE MECHANISM MEETS THEORETICAL GUARANTEE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Transformers have demonstrated remarkable in-context learning abilities, adapting to new tasks from just a few examples without parameter updates. However, theoretical understanding of this phenomenon typically assumes fixed dependency structures, while real-world sequences exhibit flexible, context-dependent relationships. We address this gap by investigating whether transformers can learn causal structures – the underlying dependencies between sequence elements – directly from in-context examples. We propose a novel framework using Markov chains with randomly sampled causal dependencies, where transformers must infer which tokens depend on which predecessors to make accurate predictions. Our key contributions are threefold: (1) We prove that a two-layer transformer with relative position embeddings can implement Bayesian Model Averaging (BMA), the optimal statistical algorithm for causal structure inference; (2) Through extensive experiments and parameter-level analysis, we demonstrate that transformers trained on this task learn to approximate BMA, with attention patterns directly reflecting the inferred causal structures; (3) We provide information-theoretic guarantees showing how transformers recover causal dependencies and extend our analysis to continuous dynamical systems, revealing fundamental differences in representational requirements. Our findings bridge the gap between empirical observations of in-context learning and theoretical understanding, showing that transformers can perform sophisticated statistical inference over structural uncertainty.

## 1 INTRODUCTION

Modern transformers exhibit a remarkable capability: they can adapt to entirely new tasks using only a handful of examples, without any parameter updates. This phenomenon, known as in-context learning (ICL) [Brown et al. \(2020\)](#), has revolutionized our understanding of what neural networks can achieve. A model trained on diverse text can suddenly perform arithmetic, translate languages, or write code – all by simply observing a few demonstrations. Yet despite extensive empirical success [Wei et al. \(2022\)](#); [Garg et al. \(2023\)](#) and theoretical investigations [von Oswald et al. \(2023\)](#); [Akyürek et al. \(2023\)](#); [Goel & Bartlett \(2024\)](#), a fundamental question remains: how do transformers adapt to the varying dependency structures present in real-world sequences? ([Allen-Zhu & Li, 2023](#); [Bietti et al., 2023](#); [Zhao et al., 2023](#); [Wibisono & Wang, 2024](#))

**The Theory-Practice Gap.** Current theoretical understanding of ICL rests on a critical simplification: most analyses assume that dependencies between sequence elements follow a fixed, predetermined structure. For instance, theoretical works typically study settings where tokens are independent  $[[x_1, f(x_1)], [x_2, f(x_2)], \dots]$  or follow rigid patterns like  $[x_1, f(x_1), x_2, f(x_2)]$  ([Bai et al., 2023](#); [Chen et al., 2024a](#); [Wang et al., 2025](#)). However, natural language and real-world sequences exhibit far richer structure – words depend on previous words in complex, context-dependent ways that vary across sentences and domains. Recent work by [Nichani et al. \(2024\)](#) began addressing this by showing transformers can encode fixed causal structures during training. Specifically, they assume an  $n$ -gram causal model (e.g., bigrams where each token depends only on the previous one) ([Rajaraman et al., 2024](#); [Edelman et al., 2024](#)), and prove that transformers can embed this structure in their attention weights to perform inference. However, in real-world scenarios, the dependency

graph itself is not fixed but varies across different sequences. For example, in language, the syntactic structure can change dramatically between different documents, and in stock price prediction, the relationships between assets can shift over time. Thus, a key challenge is

*Can transformers infer and adapt to causal structure in-context?* (★)

**Our Approach.** We introduce a novel framework where sequences are generated from Markov chains with randomly sampled causal dependencies. In our setting, each token depends on exactly one predecessor, or its “parent”, but crucially, these parent relationships are not fixed and must be inferred from context examples, which is a collection of sequences sharing the same underlying causal structure. This setup captures the essence of (★) by requiring the model to adapt to different latent structures across contexts. The transformer must infer these latent dependencies from context examples to accurately predict new sequences – mirroring how language models must adapt to different syntactic structures or reasoning patterns.

**Main Contributions.** We consider two types of Markov chains: discrete chains over a finite vocabulary and continuous linear dynamical systems. Our work makes the following contributions:

**(1) Theoretical Construction:** For discrete Markov chains, we prove that a two-layer transformer with relative position embeddings can implement Bayesian Model Averaging (BMA), the statistically optimal algorithm for inferring causal structures from observations. Our construction shows how attention mechanisms can perform sophisticated probabilistic inference over structural uncertainty. **(2) Empirical Verification:** Through extensive experiments on Markov chains, we demonstrate that transformers trained via gradient descent converge to solutions remarkably similar to our theoretical construction. Parameter-level analysis reveals that learned attention patterns directly encode posterior probabilities over causal structures, providing mechanistic insight into how transformers perform statistical inference. **(3) Information-Theoretic Analysis:** We establish conditions under which causal structures can be recovered in-context, using mutual information and data processing inequalities. Additionally, we show that gradient-based learning naturally discovers these structures early in training through  $\chi^2$ -mutual information maximization. **(4) Extensions to Continuous Systems:** We extend our framework to linear dynamical systems in continuous space, revealing fundamental differences in how transformers handle discrete versus continuous causal inference. While transformers show strong empirical performance, we identify representational limitations that prevent exact BMA implementation in continuous settings.

**Paper Organization.** Section 2 introduces our problem formulation and model architecture. Section 3 presents our main theoretical and empirical results for Markov chains. Section 4 extends the analysis to continuous dynamical systems. Appendix A discusses related work.

## 2 PRELIMINARY

### 2.1 TASK SETUP

To investigate the question (★), we consider data are generated from distributions with a latent causal structure. Each sample is a sequence of tokens  $\mathbf{x}_{1:H} = [\mathbf{x}_1, \dots, \mathbf{x}_H]$ , where the  $h$ -th token  $\mathbf{x}_h$  depends on one of its predecessors, called the parent token  $\mathbf{x}_{\text{pa}(h)}$ . This dependency relation is represented as a directed tree graph  $\mathcal{G} = \{\text{pa}(h)\}_{h \in [H]}$ , where  $\text{pa}(h) \sim \text{Unif}(1, \dots, h-1), \forall h \in 2, \dots, H$ . Given the causal structure defined above, the generative process can be written as  $\mathbf{x}_h = G(\mathbf{x}_{\text{pa}(h)})$ , where  $G(\cdot)$  denotes either stochastic sampling from the transition kernel  $\pi(\cdot | \mathbf{x}_{\text{pa}(h)})$  of Markov chains, or a deterministic transformation with additive Gaussian noise in dynamical systems.  $G(\cdot)$  is fixed during sampling the whole dataset.

For the in-context learning task, suppose we have  $L+1$  samples  $\{\mathbf{x}_{1:H}^{(l)}\}_{l \in [L+1]}$  from the same causal graph  $\mathcal{G}$ , the first  $L$  samples are provided as in-context demonstrations from which the model may infer the latent graph structure, while the last sample is the target for prediction. Except the first token  $\mathbf{x}_1^{L+1}$ , every token  $\mathbf{x}_h^{L+1}$  in this trajectory is required to be predicted via next-token prediction conditioned on  $\mathbf{x}_{1:H}^{1:L}$  and its past observations  $\mathbf{x}_{1:h-1}^{L+1}$ .

**Markov Chain.** Following Markovian assumption adopted in [Edelman et al. \(2024\)](#); [Nichani et al. \(2024\)](#); [Chen et al. \(2024b\)](#); [D’Angelo et al. \(2025\)](#), we assume sequences are sampled from

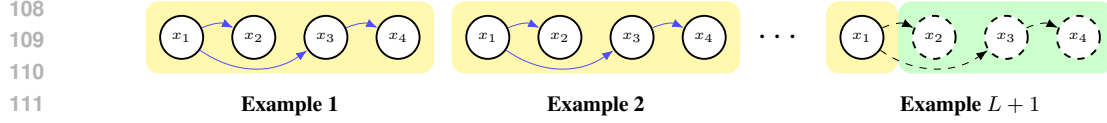


Figure 1: **Task overview of in-context causal structure learning.** Each training sequence consists of  $L$  examples with observed variables and hidden parent relations, followed by a new example  $L+1$  where the model must infer the underlying parent indices in-context from previous demonstrations.

a Markov chain with random dependencies. In this setting, Tokens  $\{\mathbf{x}_h\}$  are drawn from a finite vocabulary  $\mathcal{V} = \{e_1, \dots, e_d\}$ , where  $|\mathcal{V}| = d$  and  $\{e_i\}$  are one-hot vectors. The random dependencies are specified by latent causal graph  $\{\text{pa}(h)\}_{h \in [H]}$ . Let  $\pi : \mathcal{V} \rightarrow \Delta(\mathcal{V})$  denote the Markov transition kernel, where  $\Delta(\mathcal{V})$  is the probability simplex over  $\mathcal{V}$ . Then each token is generated as  $\mathbf{x}_h \sim \pi(\cdot \mid \mathbf{x}_{\text{pa}(h)}) \in \Delta(\mathcal{V})$ ,  $\forall h \in [H]$ , where by slight abuse of notation, we also regard  $\pi$  as the stochastic matrix  $\pi \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$  with  $\pi[i, j] = \pi(j|i)$ ,  $\sum_j \pi[i, j] = 1$ .

**Dynamical System** Beyond the discrete Markov chain case, we also consider a more challenging setting with continuous sampling space. Here tokens  $\{\mathbf{x}_h\}$  are dense vectors in  $\mathbb{R}^d$ . The link function  $g(\cdot)$  replaces the discrete transition kernel, and we instantiate it as a *linear dynamical system with additive Gaussian noise*:  $\mathbf{x}_h = g(\mathbf{x}_{\text{pa}(h)}) = \frac{1}{c}(A^\top \mathbf{x}_{\text{pa}(h)} + \varepsilon_h)$ , where  $A \in \mathcal{O}(\mathbb{R}^d)$  is orthogonal,  $\mathbf{x}_1 \sim \mathcal{N}(0, I_d)$ ,  $\varepsilon_h \sim \mathcal{N}(0, \sigma^2 I_d)$ , and  $c^2 = 1 + \sigma^2$  ensures variance stability.

These settings evaluate the extent to which transformers can perform in-context causality learning.

**Goal: Inferring the Causal Structure.** The task formulation naturally raises the following question: *Given  $L$  in-context examples, how can the model infer the underlying graph structure  $\mathcal{G}$ ?* A classical approach to this problem is *Bayesian Model Averaging* (BMA), which leverages Bayes' rule to compute the posterior distribution over possible parameter space. Treating the parent structure  $\text{pa}(h)$  as the parameter to be estimated, the distribution of having parent  $h'$  will be predicted as its posterior probability given  $L$  observations:

$$\mathbb{P}(\text{pa}(h) = h' \mid \mathbf{x}_{1:H}^{1:L}) = \frac{\mathbb{P}(\mathbf{x}_{1:H}^{1:L} \mid \text{pa}(h) = h') \mathbb{P}(\text{pa}(h) = h')}{\sum_{h'' \in [H]} \mathbb{P}(\mathbf{x}_{1:H}^{1:L} \mid \text{pa}(h) = h'') \mathbb{P}(\text{pa}(h) = h'')}. \quad (1)$$

By Eq. (1) and our task assumption, we have the following lemma of the formulation of BMA.

**Lemma 1.** *Suppose  $L$  samples are observed from the above Markov chain (or dynamical system)  $\mathbf{x}_{1:H}$  with latent causal structure  $\mathcal{G}$ . Bayesian Model Averaging give prediction of  $\text{pa}(h) \in [h-1]$ :*

$$\mathbb{P}(\text{pa}(h) = h' \mid \mathbf{x}_{1:H}^{1:L}) = \frac{\exp(\sum_{l \in [L]} \log \pi(\mathbf{x}_h^l \mid \mathbf{x}_{h'}^l))}{\sum_{h'' \in [h-1]} \exp(\sum_{l \in [L]} \log \pi(\mathbf{x}_h^l \mid \mathbf{x}_{h''}^l))} = \sigma(\hat{\mathbf{p}}^{h,L}(\log \pi))_{h'}, \quad (2)$$

where  $\hat{\mathbf{p}}^{h,L}(\log \pi) \in \mathbb{R}^d$  and  $\hat{\mathbf{p}}_{h'}^{h,L} = \sum_{l \in [L]} \log \pi(\mathbf{x}_h^l \mid \mathbf{x}_{h'}^l)$ . See Appendix C.2 for detailed proof.

This Bayesian formulation provides a principled baseline for inferring causal structure, and serves as a point of comparison for the in-context learning behavior of transformers.

## 2.2 MODEL ARCHITECTURE

### 2.2.1 STANDARD TRANSFORMER

Decoder-only Transformer is a neural network structure dealing with sequential data. Given a sequence of tokens  $(\mathbf{w}_1, \dots, \mathbf{w}_T)$ , transformers first embed tokens and add a positional encoding to the tokens:  $\mathbf{h}_t^{(0)} = E(\mathbf{w}_t) + P(t) \in \mathbb{R}^d, \forall t \in [T]$ . In a matrix form, the mapped tokens of input is  $\mathbf{H}^{(0)} = \mathbf{h}_{1:T}^{(0)} \in \mathbb{R}^{T \times d}$ . Subsequent layers consist of multi-head attention layers (MHA) followed by multilayer perceptron layers (MLP). At layer  $l$ , the hidden features  $\mathbf{H}^{(l-1)}$  will be updaed as follows. First, causal-mask self-attention layer will compute the output by:

$$\text{Attn}(\mathbf{H}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V) = \sigma \left( \mathcal{M} \left( \frac{(\mathbf{H}\mathbf{W}_Q)(\mathbf{H}\mathbf{W}_K)^\top}{\sqrt{d_k}} \right) \right) \mathbf{H}\mathbf{W}_V,$$

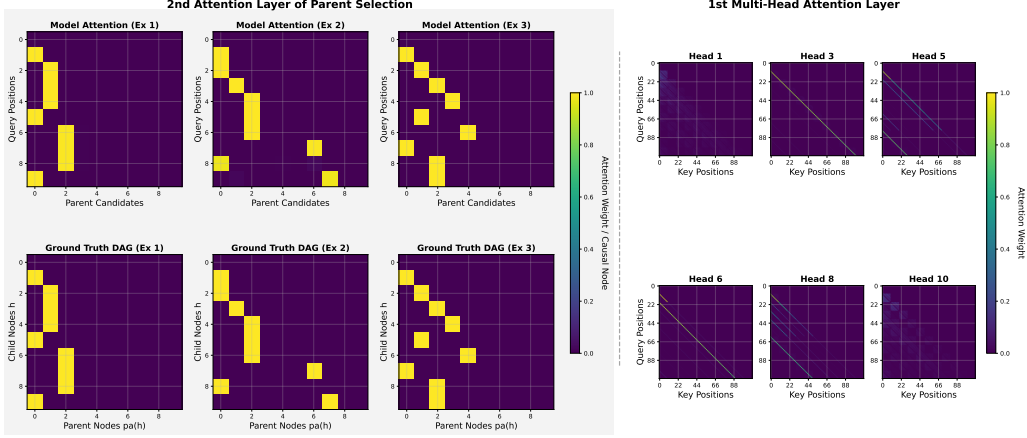


Figure 2: Visualization of Attention Weights  $\mathcal{A}^{(1)}, \mathcal{A}^{(2)}$ . Left:  $\mathcal{A}^{(2)}$  on 3 examples. Transformer shows promising capabilities to select causal tokens in context. Right: Six representative heads (out of 10) from  $\mathcal{A}^{(1)}$  (Head 1 and 10 degenerate). (trained with  $L = 10, H = 10, d = 5, 1024$  steps).

where  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d_k}$ ,  $\sigma(v)_i = \frac{\exp(v_i)}{\sum_j \exp(v_j)}$  applied to matrix row-wisely,  $\mathcal{M}$  is the causal mask where  $\mathcal{M}(X)_{ij}$  is  $-\infty$  if  $i > j$  else  $X_{ij}$ . Then multi-headed attention gives the output:

$$\text{MHA}(\mathbf{H}) = \left( \bigoplus_{m=1}^M \text{Attn}(\mathbf{H}; \mathbf{W}_Q^m, \mathbf{W}_K^m, \mathbf{W}_V^m) \right) \mathbf{W}_O,$$

where  $\bigoplus$  denotes the concatenation of vectors and  $\mathbf{W}_O \in \mathbb{R}^{M d_k \times d}$ . Getting intermediate features  $\text{MHA}_l(\mathbf{H}^{(l-1)})$  from attention layer, this feature will be added to the **residual stream** which aggregate the previous output:  $\hat{\mathbf{H}}^l = \mathbf{H}^{(l-1)} + \text{MHA}_l(\mathbf{H}^{(l-1)})$ . FFN layer adopts this as input and updates this stream as:

$$\text{FFN}(\hat{\mathbf{H}}) = \sigma(\hat{\mathbf{H}} \mathbf{W}_1) \mathbf{W}_2, \quad \mathbf{H}^{(l)} = \hat{\mathbf{H}}^{(l)} + \text{FFN}_l(\hat{\mathbf{H}}^{(l)}),$$

where  $\mathbf{W}_1 \in \mathbb{R}^{d \times d_m}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{d_m \times d}$  and  $\sigma(\cdot)$  is the activation function. Finally the output of  $L$ -layer Transformer is  $\sigma(\mathbf{H}^{(L)} \mathbf{W}_U)$  projected to vocabulary logits by  $\mathbf{W}_U \in \mathbb{R}^{d \times V}$ .

### 2.2.2 DISENTANGLED TRANSFORMERS

To better analyze the role each part of transformers play in learning a task, prior works [Friedman et al. \(2023\)](#) propose the disentangled transformer which decouples the twisted features in the residual stream. Instead of adding each layer's output, disentangled transformer concatenates it with residual stream. Considering the decoder-based attention-only transformers we will mainly focus on, it will update the hidden states  $\mathbf{H}^{(l-1)} \in \mathbb{R}^{T \times d_{l-1}}$  by

$$\mathbf{H}^{(l)} = [\mathbf{H}^{(l-1)}, \text{Attn}_1(\mathbf{H}^{(l-1)}), \dots, \text{Attn}_M(\mathbf{H}^{(l-1)})] \in \mathbb{R}^{T \times (1+M)d_{l-1}}, \quad (3)$$

where in each attention head,  $\mathbf{W}_K \mathbf{W}_Q^\top$  is reparametrized by  $\mathbf{W}_{KQ}$ ,  $\mathbf{W}_O \mathbf{W}_V$  by  $\mathbf{W}_{OV}$  and the initial input  $\mathbf{H}^{(0)} \in \mathbb{R}^{T \times d_0}$  is given by  $\mathbf{h}_t^{(0)} = [E(\mathbf{w}_t), P(\mathbf{w}_t)] = [\mathbf{e}_{\mathbf{w}_t}, \mathbf{e}_t] \in \mathbb{R}^{d+T}, \forall t$ . Consider in our task, the input sequence consists of  $L+1$  examples of length- $H$  chains, leading to the embedding size  $T$  of  $P(\mathbf{w}_t) \in \mathbb{R}^T$  equal to  $(L+1)H$ . If we set vocabulary dimension  $d = H = L = 10$ , then  $d = 10 \ll T = 110$  in the input embedding and  $\mathbf{W}_{KQ}$  in the 1st layer will have  $\Theta(H^2 L^2) = \Theta(10^4)$  parameters. Instead, considering  $\mathbf{w}_t$  is the  $h$ -th token in example  $l$ , we use two types of embeddings representing this:  $\text{Pos}_L(\mathbf{w}_t) = \mathbf{e}_l \in \mathbb{R}^L$ ,  $\text{Pos}_H(\mathbf{w}_t) = \mathbf{e}_h \in \mathbb{R}^H$ . This reduces the required parameter for training. And the formulation of this transformer is given by Eq. (34). Empirical results presented in Appendices G, H and I demonstrate that the standard disentangled transformer with full absolute position encoding, its positional variants and standard transformers with FFNs all behave consistently with the model described below. To enable a tractable parameter-level analysis, we adopt a simplified structure that reduces the parameter count and simplifies training.

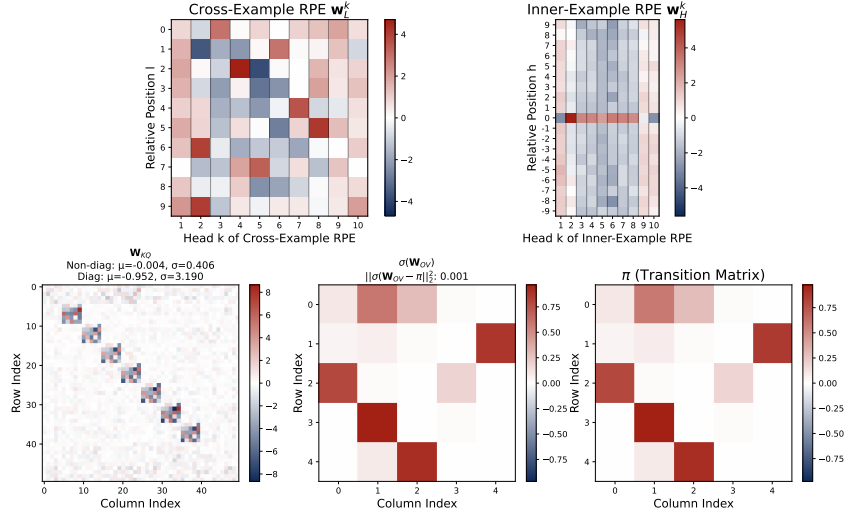


Figure 3: Parameter visualization of 2-layer transformer with RPE. Heads of inner-example RPE  $\mathbf{w}_H^k$  uniformly show the largest value at position  $h = 0$  except Head 1, 9, 10. Correspondingly,  $\mathbf{W}_{KQ}$  shows similar blocks on diagonal except block 1, 9, 10. Besides,  $\sigma(\mathbf{W}_{OV})$  approximates  $\pi$ .

**Relative Position Embedding.** While the original Transformer employs the above absolute positional encodings, subsequent research has demonstrated the advantages of relative positional embeddings (RPE) Shaw et al. (2018); Su et al. (2023). RPE is parametrized by a vector  $\mathbf{w} \in \mathbb{R}^T$  which assigns attention score  $\mathbf{w}(i, j)$  only via relative distance  $i - j$  between positions of query  $i$  and key  $j$ . Similar to the above absolute embedding, we adopt two types of RPE:  $\mathbf{w}_H \in \mathbb{R}^{2H-1}$  representing the order  $l$  from  $L+1$  examples and  $\mathbf{w}_L \in \mathbb{R}^L$  representing the order  $h$  from  $H$  tokens.

$$\mathbf{w}_H(h, h') = \mathbf{w}_H[h - h'], \forall (h, h') \in [H]^2, \quad \mathbf{w}_L(l, l') = \begin{cases} \mathbf{w}_L[l - l'], & l > l', \\ -\infty, & \text{else. (for causal mask)} \end{cases}$$

**Attention with RPE.** With RPE in the 1st layer, the output  $u_t$  for  $x_t = x_h^l$  is given by:

$$\begin{aligned} u_t &= \text{Attn}_{x_t \rightarrow x_{1:T}} = \sum_{t'} \sigma_{t'}(\mathbf{w}_H(h, \cdot) + \mathbf{w}_L(l, \cdot)) x_{t'} \\ &= \sum_{t' \leftrightarrow (h', l')} \frac{\exp(\mathbf{w}_H(h, h') + \mathbf{w}_L(l, l'))}{\sum_{t''} \exp(\mathbf{w}_H(h, h'') + \mathbf{w}_L(l, l''))} x_{t'}. \end{aligned} \quad (4)$$

**Self-Attention Layer.** Suppose the 1st layer use  $K$  heads, for the 2nd layer, the input is  $v_t = H_t^{(1)} \in \mathbb{R}^{d+Kd}$  given by disentangled residual in Eq. (3). The features of last example  $L+1$ :  $v_h^{L+1}$  are taken as query, key and value tokens into the attention layer. The output gives Transformer’s prediction. Suppose the input is  $x_{1:T} = x_{1:H}^{1:L+1}$ , this transformer architecture is formulated as follows:

**1st RPE Attention (K-head):**  $u_h^k = \text{Attn}_{x_h^{L+1} \rightarrow x_{1:T}}^k = \sigma(\mathbf{w}_H^k(h, \cdot) + \mathbf{w}_L^k(L+1, \cdot)) x_{1:T}^\top \in \mathbb{R}^d$ ,

**Disentangled Residual:**  $v_h = [u_h^1, \dots, u_h^K], z_h = [x_h^{L+1}, v_h] \in \mathbb{R}^{d+Kd}$ ,

**2nd Attention (1-head):**  $f_{\text{tf}}(\cdot | \mathcal{H}_h^L) = \sigma(z_{1:h-1}^\top \mathbf{W}_{KQ} z_h)^\top z_{1:h-1}^\top \mathbf{W}_{OV}$   
 $= \sigma(v_{1:h-1}^\top \mathbf{W}'_{KQ} v_h)^\top x_{1:h-1}^{L+1} \mathbf{W}'_{OV} \in \mathbb{R}^d$ , (5)

where  $f_{\text{tf}}(\cdot | \mathcal{H}_h^L) \in \mathbb{R}^d$  denotes the output of the transformer based on context  $\mathcal{H}_h^L = [x_{1:H}^{1:L}, x_{1:h-1}^{L+1}]$  (or the context denoted by  $\mathcal{H}$  for brevity) and we assume some blocks in  $\mathbf{W}_{KQ}, \mathbf{W}_{OV}$  are 0:

$$\mathbf{W}_{KQ} = \begin{bmatrix} 0_{d \times d} & 0_{d \times Kd} \\ 0_{Kd \times d} & \mathbf{W}'_{KQ} \end{bmatrix}, \mathbf{W}_{OV} = \begin{bmatrix} \mathbf{W}'_{OV} & 0_{d \times Kd} \\ 0_{Kd \times d} & 0_{Kd \times Kd} \end{bmatrix}, \quad (6)$$

where  $\mathbf{W}'_{KQ} \in \mathbb{R}^{Kd \times Kd}$ ,  $\mathbf{W}'_{OV} \in \mathbb{R}^{d \times d}$  are trainable and this simplification is supported by the results on disentangled transformers in Appendix G. To train transformers, cross-entropy loss is used for  $x \in \mathcal{V}$  of Markov chain (MC) and MSE loss for dynamical system (DS) shown in Eq. (15).



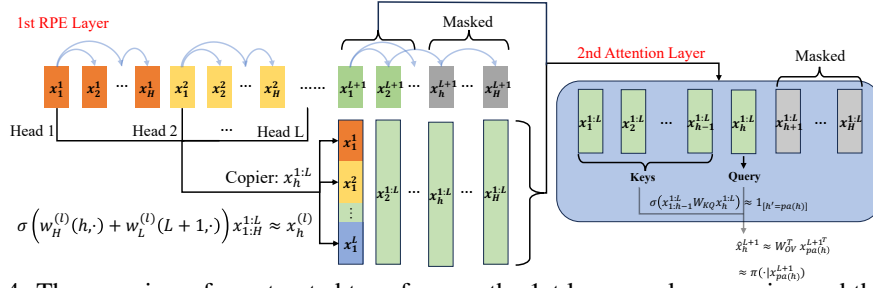


Figure 4: The overview of constructed transformer: the 1st layer works as copier, and the attention patterns in 2nd layer follows BMA, which approximately select correct parent token.

### 3 CAN TRANSFORMERS IN-CONTEXT LEARN CAUSAL STRUCTURES?

#### 3.1 2-LAYER TRANSFORMER LEARNT TO SELECT CAUSAL STRUCTURE IN-CONTEXT

To investigate the question ( $\star$ ), we first train 2-layer transformers with RPE introduced above on the Markov chain setting. Each input has  $L + 1$  samples  $\{x_{1:H}^l\}$  of Length- $H$  Markov chain with causal structure  $\mathcal{G}$ . The input sequence length is  $T = H(L + 1)$ . We set the transformer has  $K$  heads in the first RPE layer:  $\{(\mathbf{w}_H^k, \mathbf{w}_L^k)\}_{k \in [K]}$  and 1 head for the 2nd attention layer  $(\mathbf{W}_{KQ}, \mathbf{W}_{OV})$ . All RPE parameters are initialized randomly from Gaussian distribution and  $(\mathbf{W}_{KQ}, \mathbf{W}_{OV})$  from zero.

For an attention layer, the attention weights  $\mathcal{A}$  normalized by the  $\sigma$  reveal which tokens a query primarily attends to, enabling mechanistic interpretability analyses such as circuit discovery [Olsson et al. \(2022\)](#). We first look at the attention patterns  $\mathcal{A}^{(1)}, \mathcal{A}^{(2)}$  from the 1st and 2nd transformer layer. Mathematically, they are matrix where the  $i$ -th row denotes the attention weights to the whole sequence and  $\mathcal{A}_{ij}^{(*)} = \mathcal{A}_{i \rightarrow j}^{(*)}$  is formulated by:

$$\mathcal{A}^{(*)} = \sigma(\tilde{\mathcal{A}}^{(*)}), \tilde{\mathcal{A}}_{i \rightarrow j}^{(1),k} = \mathbf{w}_H^k(h_i, h_j) + \mathbf{w}_L^k(l_i, l_j), \tilde{\mathcal{A}}_{h \rightarrow h'}^{(2)} = \mathbf{v}_h^\top \mathbf{W}_{KQ} \mathbf{v}_{h'},$$

where we consider the index  $i$  is related to token  $x_i$  from Eq. (5) which is the  $h_i$ -th token of  $l_i$ -th example (for index  $j$ , the notation is similar),  $i, j \in [T]$ ,  $\mathbf{v}_{h'}$  are the hidden feature  $\mathbf{v}_{1:H}^{L+1}$  of  $L + 1$ -th example from Layer 1 and  $\mathcal{A}^{(1),k} \in \mathbb{R}^{T \times T}, \mathcal{A}^{(2)} \in \mathbb{R}^{H \times H}$ . Trained attention patterns of  $\mathcal{A}^{(2)}$  match the groundtruth causal structure in Fig. 2. And for the 1st layer, some heads of attention weights  $\mathcal{A}^{(1),k}$  didn't learn meaningful features shown by Fig. 2 (e.g., Head 1, 9, 10). Then we dive into the parameter level, and visualize the trainable parameters of 2-layer transformer  $\mathbf{w}_H^k, \mathbf{w}_L^k, \mathbf{W}_{KQ}, \mathbf{W}_{OV}$ . Positional or diagonal patterns in  $\mathbf{w}_H^k, \mathbf{W}_{KQ}$  and the similarity between  $\mathbf{W}_{OV}$  and  $\log \pi$  can be observed in Fig. 3. To fully understand why the transformer can select causal structures and what it learnt, we will need to analyze it theoretically.

**Takeaway 1.** Transformer formulated by Eq. (5) effectively identifies latent causal parents in-context (Fig. 2) and learns highly structural parameters aligned with the task (Fig. 3).

#### 3.2 CONSTRUCTED TRANSFORMERS IMPLEMENT STATISTICAL ALGORITHM

Based on the patterns observed in Fig. 3, we make the following assumptions for the transformer defined by Eq. (5):

$$\begin{aligned} \tilde{\mathbf{w}}_H^k[h] &= \beta \begin{cases} +1, & h = 0, \\ -1, & h \in [\pm H] \setminus 0, \end{cases} \quad \exists k' \in [L] \text{ s.t. } \tilde{\mathbf{w}}_L^k[l] = \beta \begin{cases} +1, & l = k', \\ -1, & l \in [L] \setminus k', \end{cases} \\ \tilde{\mathbf{W}}_{KQ} &= \begin{bmatrix} \mathbf{W} & 0_{d \times d} & \cdots & 0_{d \times d} \\ 0_{d \times d} & \mathbf{W} & \cdots & 0_{d \times d} \\ \vdots & \vdots & \ddots & \vdots \\ 0_{d \times d} & 0_{d \times d} & \cdots & \mathbf{W} \end{bmatrix}, \quad \sigma(\tilde{\mathbf{W}}_{OV}) = \pi, \end{aligned} \quad (7)$$

where for RPE, we can simply assume one element  $k'$  of  $\tilde{\mathbf{w}}_L^k$  dominates it while 0-th entry dominates  $\tilde{\mathbf{w}}_H^k$ :  $\tilde{\mathbf{w}}_H^k[0] \gg \tilde{\mathbf{w}}_H^k[-0]$ ,  $\tilde{\mathbf{w}}_L^k[k'] \gg \tilde{\mathbf{w}}_L^k[-k']$  and  $\mathbf{W}$  is unknown. Since the  $K$  heads are identical up to their indices, we assume without loss of generality that the dominant entry of  $\tilde{\mathbf{w}}_L^k$  occurs at position  $k$ , i.e.,  $\tilde{\mathbf{w}}_L^k[k] = \beta$  and we set  $K = L$ . The theorem below shows that aligned with the above restriction, a constructed transformer can implement statistical algorithm for inferring causal structure  $\{pa(h)\}$  hidden behind  $\mathbf{x}_{1:H}^{1:L}$  and predicting  $\mathbf{x}_h^{L+1} \sim \pi(\cdot | \mathbf{x}_{pa(h)}^{L+1})$ :

**Theorem 1.** *Under the restriction by Eq. (7), the transformer  $\mathbf{f}_\theta$  is parameterized by  $\theta \in \{(\beta, \mathbf{W})\}$ . Then  $\mathbf{f}_\theta$  with  $\mathbf{W} = \log \pi$  in Eq. (7) whose second attention layer  $\mathcal{A}^{(2)}(\mathcal{H}; \theta)$  approximates Bayesian Model Averaging (see Lemma 1) satisfying the following convergence property:*

$$\lim_{\beta \rightarrow \infty} \mathcal{A}_{h \rightarrow \cdot}^{(2)}(\mathcal{H}; \theta) = \lim_{\beta \rightarrow \infty} \sigma(\tilde{\mathcal{A}}_{h \rightarrow \cdot}^{(2)}(\mathcal{H}; \theta)) = \sigma(\mathbf{p}_{\text{BMA}}^{h,L}). \quad (8)$$

Further, the transformer’s prediction of the distribution of last example  $\mathbf{x}_{1:H}$  with  $L$  context examples converges to the true conditional distribution given the causal parent guaranteed by Theorem 2:

$$\lim_{\beta, L \rightarrow \infty} \mathbf{f}_\theta(\cdot | \mathcal{H}_h^L) = \pi(\cdot | \mathbf{x}_{pa(h)}), \forall h \in [H]. \quad (9)$$

*Proof Sketch.* Figure 4 gives an overview of the construction: in the first RPE attention layer, each head from Eq. (7) is assigned to retrieve one historical copy of the same token  $\mathbf{x}_h$ , so that concatenating  $L$  heads recovers all past observations  $\mathbf{x}_h^{1:L}$ . In the second layer, with the condition in Eq. (7), the attention score between tokens  $(h, h')$  reduces to a bilinear form  $\hat{p}_{h'}^h(\mathbf{W}) = \sum_l \mathbf{x}_{h'}^{l\top} \mathbf{W} \mathbf{x}_h^l$ , which by  $\mathbf{W} = \log \pi$  coincides with the BMA score  $\mathbf{p}_{\text{BMA}}^{h,L} = \sum_{l \in [L]} \log \pi(\mathbf{x}_h^l | \mathbf{x}_{h'}^l)$ . With the causal mask, the softmax attention exactly matches the parent-selection distribution in BMA. By the theoretical guarantee of causal token selection (Theorem 2), OV matrix  $\mathbf{W}_{OV}$  receives the correct parent  $\mathbf{x}_{pa(h)}^{L+1}$  and make prediction for  $\pi(\cdot | \mathbf{x}_{pa(h)}^{L+1})$ . The full technical proof is deferred to Appendix C.1.  $\square$

D’Angelo et al. (2025) also consider an in-context causal learning task. With minor modification of RPE structure and above construction, we can still show transformers can exactly implement BMA.<sup>1</sup>

**Takeaway 2.** Two-layer transformers can explicitly implement BMA for causal token selection.

### 3.3 WHAT ALGORITHM DOES THE TRANSFORMER LEARN?

Although we have constructed a transformer implementing this algorithm, what do transformers actually learn after training? Since the core lies in the attention weight  $\mathcal{A}^{(2)}$  and  $\mathbf{W}_{KQ}$  (where we use  $\mathbf{W}_{\text{tf}}$  to denote the trainable submatrix in Eq. (7)) recovering graph structures, we next analyze their characteristics in detail. We first define the following parent selection metric which quantitatively shows the loss of algorithms to predict parent indices:

$$\mathcal{L}_{pa}(\mathcal{A}^{(2)}(\mathbf{x}_{1:H}^{1:L+1}; \mathcal{G}), \mathcal{G}) = -\frac{1}{H} \sum_{h \in [H]} e_{pa(h)}^\top \log \mathcal{A}_h^{(2)} = -\frac{1}{H} \sum_{h \in [H]} \log \mathcal{A}_{h \rightarrow pa(h)}^{(2)}, \quad (10)$$

where  $\mathcal{A}^{(2)}$  is seen as an algorithm of predicting  $e_{pa(h)}$  given input  $\mathbf{x}_{1:H}^{1:L+1}$  and we have  $\mathcal{L}_{pa}(\mathcal{A}_{\text{BMA}}, \mathcal{G}) = -\frac{1}{H} \sum_{h \in [H]} e_{pa(h)}^\top \sigma(\hat{\mathbf{p}}_h(\log \pi))$  by Eq. (2) where  $\hat{\mathbf{p}}_h(\mathbf{W}) = \sum_l \mathbf{x}_{1:h-1}^{l\top} \mathbf{W} \mathbf{x}_h^l$ . We visualize this metric  $\mathcal{L}_{pa}$  during transformers’ training process in Fig. 7 and compare it with BMA’s. We observe that the transformer’s parent selection loss decreases in training while remaining above the loss of BMA, gradually approaching it.

**Generalized parent selection with size  $L'$  varying.** We further test how well the transformer and BMA generalize in parent selection under different sample sizes  $L'$ : Since  $\mathcal{A}^{(2)}$  and  $\mathcal{A}_{\text{BMA}}$  are formulated via  $\hat{\mathbf{p}}^L = \sum_{l \in [L]} \mathbf{x}_{1:h-1}^{l\top} \mathbf{W} \mathbf{x}_h^l$ , we vary the number of demonstrations as a set of  $L'$ , and finally compute  $\hat{\mathbf{p}}^{L'}$ ,  $\mathcal{A}_h^{L'}$ , and the parent selection loss  $\mathcal{L}_{pa}^{L'}(\mathbf{W})$  with  $\mathbf{W} \in \{\mathbf{W}_{\text{tf}}^{(L)}, \log \pi\}$ .

<sup>1</sup>A detailed comparison with D’Angelo et al. (2025), including similarities and differences, is provided in Appendix A (Related Work).

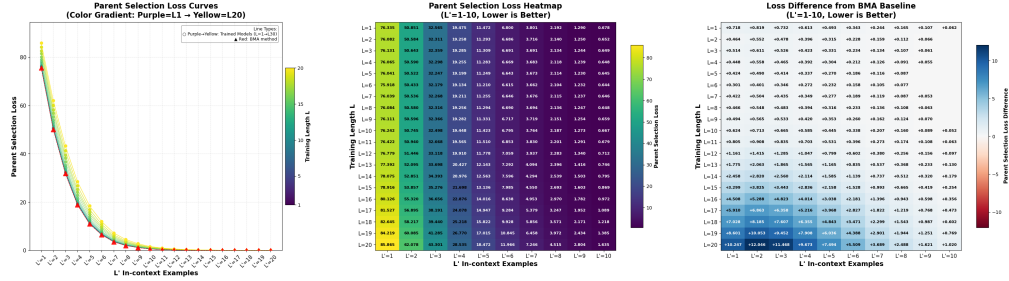


Figure 5: Generalization of Parent Selection loss  $\{\mathcal{L}_{pa}^{L'}\}$  for transformers trained with  $L \in \{1, \dots, 20\}$ ,  $d = 10$ , and  $H = 15$  with first layer fixed as constructed.

From Fig. 5, we observe that: 1) across different test sizes  $L'$ , the trained transformers achieve performance close to BMA (loss differences mostly within a small margin); 2) smaller training length  $L$  often generalize better, with parent loss curves approaching BMA more closely; 3) for a model with fixed training size  $L$ , the parent loss decreases rapidly as  $L'$  increases, converging toward zero. The above results suggest the trained transformers have comparable performance to BMA.

**Parameter Verification.** Beyond behavioral agreement, a crucial question is whether the transformer encodes the BMA inference rule within its learned parameters. Therefore, we evaluate the similarity between the trained weight  $\mathbf{W}_{tf}^{(L)}$  and the theoretical BMA parameter  $\mathbf{W} = \log \pi$ . As a first attempt, we check whether  $\sigma(\mathbf{W}_{tf}) = \pi$ , since for stochastic matrix  $\pi$  it holds that,  $\sigma(\mathbf{W}_{tf}) = \pi \iff \mathbf{W}_{tf} = \log \pi + \mathbf{b}\mathbf{1}^\top$ ,  $\forall \mathbf{b}$  where  $\mathbf{b}\mathbf{1}^\top$  denotes a row-wise shift of  $\log \pi$ , which is canceled out by the row-softmax  $\sigma$ . This provides a reasonable way to normalize KQ matrix and makes the scale comparable to  $\pi$  with scale  $[0, 1]$ . However, the empirical results do not support this hypothesis (c.f. Fig. 6, first three subfigures). With some efforts, we can see the attention mechanism  $\sigma(\mathbf{v}_{1:h-1}^\top \mathbf{W} \mathbf{v}_h)$  introduces an additional degree of freedom:

**Proposition 1** (Invariances of attention scores). *Since attention operates on a single query  $\mathbf{v}_h$ , if the columns of  $\mathbf{W}_{tf}$  differ from those of  $\log \pi$  by an additive factor, i.e.  $\mathbf{W}_{tf} = \log \pi + \mathbf{1}\mathbf{a}^\top$ ,  $\mathbf{a} \in \mathbb{R}^d$ , then transformer with  $\mathbf{W}_{tf}$  learnt BMA by:  $\sigma(\sum_l \mathbf{x}_{1:h-1}^\top \mathbf{W}_{tf} \mathbf{x}_h^l) = \sigma(\sum_l \mathbf{x}_{1:h-1}^\top \log \pi \mathbf{x}_h^l)$ . Further, if Markov chain  $\mathbf{x}_{1:H}$  is stationary and  $\mathbf{W}_{tf} = \log \pi + \mathbf{1}\mathbf{a}^\top + \mathbf{b}\mathbf{1}^\top$ ,  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ , the above conclusion also holds asymptotically. See the detailed proof in Appendix C.3.*

Following this proposition, we evaluate the discrepancy between  $\sigma(\mathbf{W}_{tf}\mathbf{1}\mathbf{a}^\top)$  and  $\pi$ . As illustrated in Fig. 6, the deviation remains small, with

$$\frac{1}{d} \|\sigma(\mathbf{W}_{tf} - \mathbf{1}\mathbf{a}^\top) - \pi\|_F < 0.05. \quad (11)$$

This also holds across various vocabulary size  $d \in \{10, 30, 50\}$ , further confirming the structural alignment between the learned model and the BMA algorithm. Taken together with our theoretical construction and empirical results, these findings strongly suggest that transformers implement the BMA method for in-context causal parent selection.

**Takeaway 3.** Transformers with trainable  $\mathbf{W}_{tf}$  closely approximate BMA in causal token selection (Fig. 5) and learn parameters which explicitly show strong alignment with BMA (Fig. 6).

### 3.4 THEORETICAL UNDERSTANDING AND GUARANTEE OF LEARNED ALGORITHM

Beyond identifying what algorithm a trainable transformer adopts, we further establish a theoretical understanding of transformers' in-context causal structure selection mechanism via information-theoretic principles. Our approach follows Nichani et al. (2024), which leverages mutual information together with the data's inherent property of the Data Processing Inequality (DPI). In contrast to their gradient-based proof, we show that transformers can exploit this property directly in context. Moreover, our analysis generalizes the  $\chi^2$ -mutual information framework of Nichani et al. (2024); D'Angelo et al. (2025) to the setting reducible to classical mutual information, by exploiting the



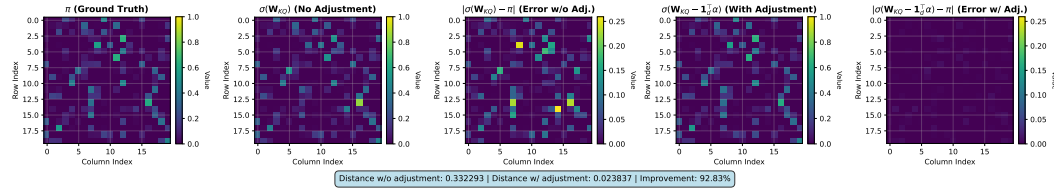


Figure 6: Parameter-Level Comparison between transformer and BMA ( $\mathbf{W}_{KQ} = \log \pi$ ). Here,  $\mathbf{W}_{KQ}$  denotes the diagonal block being trained. Trained with  $d = 20$ ,  $H = 50$ ,  $L = 3$ , and 2048 training steps while  $\mathbf{W}_{OV}$  is fixed. (See results with both  $\mathbf{W}_{OV}$  and  $\mathbf{W}_{KQ}$  trainable in Fig. 8.)

information-theoretic structure characterized in Lemma 3. Finally, our proof of Theorem 2 doesn't rely on the stationarity assumption on the data distribution not requiring the chain to be mixed.<sup>2</sup>

**Definition 1** (Mutual Information and Conditional Entropy). Consider  $x, y$  are two random variables in discrete or continuous space  $\Omega$ ,  $P_{x,y}$  and  $P_x, P_y$  denote joint and marginal distribution. The mutual information  $I(x; y)$ , entropy  $H(x)$  and the conditional entropy  $H(x|y)$  is given by:

$$I(x; y) = \int_x \int_y P_{x,y}(x, y) \log \frac{P_{x,y}(x, y)}{P_x(x)P_y(y)}, \quad H(x) = - \int_x P_x(x) \log P_x(x),$$

$$H(x|y) = \int_x \int_y P_{x,y}(x, y) \log \frac{P_{x,y}(x, y)}{P_y(y)} = H(x) - I(x; y), \quad (12)$$

Further,  $\chi^2$ -mutual information is given by:  $I_{\chi^2}(x; y) := \int_x \int_y \frac{P_{x,y}(x, y)^2}{P_x(x)P_y(y)} - 1$ .

$I, I_{\chi^2}$  can be uniformly derived from  $f$ -divergence which helps to prove DPI for generalized  $f$ -mutual information  $I_f$ . These information metrics reveal an essential property in data:

**Lemma 2** (DPI. Theorem 3.9 and 7.16 in Polyanskiy & Wu (2023)). If random variables  $x \rightarrow y \rightarrow z$ , i.e. satisfies the Markov property  $p(x, y, z) = p(x)p(y|x)p(z|y)$ , then we have  $I_f(y; z) \geq I_f(x; z)$ . Further, for classic mutual information,  $I(x; z) = I(y; z)$  iff  $I(x; y|z) = 0$  iff  $x \rightarrow z \rightarrow y$ .

Suppose we have a Markov chain  $x_{1:H}$  with latent causal structure  $\{pa(h)\}$ . Since  $\forall h' \neq pa(h)$ , we have  $P(x_h = z|x_{pa(h)} = y, x_{h'} = x) = P(z|y)$ , it is easy to verify  $x_{h'} \rightarrow x_{pa(h)} \rightarrow x_h$  while  $x_{h'} \rightarrow x_h \rightarrow x_{pa(h)}$  doesn't hold. Thus, we have the following corollary:

**Corollary 1.** For Markov chain with causal structure  $\mathcal{G}$ ,  $I(x_h; x_{pa(h)}) > I(x_h; x_{h'}), \forall h' \neq pa(h)$ .

Applying this corollary, we get the following Lemma:

**Lemma 3.** For Markov chain with causal structure  $\mathcal{G}$  and transition kernel  $p(\cdot|\cdot)$ , we have

$$\mathbb{E}[\log p(x_h|x_{pa(h)})] > \mathbb{E}[\log p(x_h|x_{h'})], \quad \forall h' \neq pa(h). \quad (13)$$

The LHS above equals  $H(x_h|x_{pa(h)})$ , while the RHS differs from conditional entropy but follows  $H(x_h|x_{h'})$  from the non-negativity of the KL divergence. Then by the relation of conditional entropy and mutual information, we can apply DPI to prove the Lemma. See Appendix C.4. With Lemma 3, we can build the relation between  $\mathbb{E}[\log p(x_h|x_{h'})]$  and attention weights, showing the theoretical guarantee concerning parent selection for the transformer:

**Theorem 2.** Suppose the transformer is constructed as in Theorem 1, which implements the BMA method. Then the attention weights  $\tilde{\mathcal{A}}_{h,\cdot}^L = \tilde{\mathcal{A}}_{h,\cdot}^{(2)}(x_{1:H}^{1:L+1})$  predicting parent index  $pa(h)$  will satisfy:

$$\lim_{L \rightarrow \infty} \tilde{\mathcal{A}}_{h,\cdot}^L = \lim_{L \rightarrow \infty} \sigma(\hat{p}^{h,L}) = e_{pa(h)} \in \mathbb{R}^H, \quad \text{where } \hat{p}^{h,L} = \sum_{l=1}^L \log \pi(x_h^l|x_{h'}^l).$$

The proof of the theorem is deferred to Appendix C.5.

**Takeaway 4.** Information-theoretic analysis reveals that the selection mechanism exploits the conditional entropy, where the DPI guarantees the identifiability of the true causal parent.

<sup>2</sup>This ensures that when the chain has not mixed, i.e., when the time index  $h$  is small, the parent-selection guarantee still holds.

### 3.5 CAUSAL STRUCTURE IN TRAINING DYNAMICS

We further look at the training dynamics of the transformer model. We prove the random causal structure embedded in inputs will be recovered in the gradients of loss w.r.t. the core  $\mathbf{W}_{KQ}$  matrix:

**Theorem 3 (Informal).** *Consider the transformer  $f_\theta$  constructed as in Theorem 1 with trainable diagonal block  $\mathbf{W}$  of  $\mathbf{W}_{KQ}$  specified in Eq. (7) and trained with cross-entropy loss*

$$\mathcal{L}(\theta) = - \sum_{h=1}^H \mathbb{E}_{(\mathbf{x}_{1:H}^{1:L+1}, \mathcal{G}) \sim P_\pi} [\log(f_{\theta_0}(\mathbf{x}_h^{L+1} | \mathcal{H}) + \epsilon)] = - \sum_{h=1}^H \mathbb{E}_{\mathcal{G} \sim P_{\mathcal{G}}} [\ell(\theta; h, \mathcal{G})], \quad (14)$$

with joint distribution  $P_\pi$  of latent graph and input,  $\ell(\theta; h, \mathcal{G}) = \mathbb{E}_{\mathbf{x}_{1:H}^{1:L+1} \sim P_{\pi|\mathcal{G}}} [\log f_{\theta_0}(\mathbf{x}_h^{L+1} | \mathcal{H})]$  and  $\hat{\mathbf{p}} = \sum_l \mathbf{x}_{1:h-1}^{l\top} \mathbf{W} \mathbf{x}_h^l \in \mathbb{R}^{h-1}$ . If the Markov chain is stationary, i.e.,  $\mathbf{x}_h \sim \mu^\pi, \forall h \in [H]$ , at initialization of  $f_{\theta_0}$  with  $\mathbf{W} = \mathbf{0}$  assuming  $f_\theta$  outputs  $\mu^\pi$  for any input, then the gradient satisfies

$$\frac{\partial \ell(\theta_0; h, \mathcal{G})}{\partial \hat{\mathbf{p}}_{pa(h)}} \geq \frac{\partial \ell(\theta_0; h, \mathcal{G})}{\partial \hat{\mathbf{p}}_{h'}}, \quad \forall h' \neq pa(h).$$

The theorem above is proved by leveraging the  $\chi^2$ -mutual information, as detailed in Appendix C.6. This result provides an explanation of how transformers can extract meaningful information from data. To further support the theory, we verify it empirically by visualizing  $\frac{\partial \mathcal{L}(\mathbf{W})}{\partial \hat{\mathbf{p}}}$  in Fig. 9.

**Takeaway 5.** Gradient at initialization is able to recover the latent causal structure, driven by  $\chi^2$ -mutual information, which facilitates structural discovery in early training (Thm. 3 and Fig. 9).

## 4 DYNAMICAL SYSTEM EXTENSION: FROM DISCRETE TO CONTINUOUS

Further, we consider to investigate the Markov chain in continuous space, where we look at the linear dynamical system with latent causal structures:  $\mathbf{x}_h = \frac{1}{c} (\mathbf{A}^\top \mathbf{x}_{pa(h)} + \varepsilon_h) \in \mathbb{R}^d, \varepsilon_h \in \mathcal{N}(0, \sigma^2 I_d)$ . We first train a transformer with RPE introduced in Eq. (5) on data generated from the dynamical system. Similar experimental results on attention weights  $\mathcal{A}^{(1)}, \mathcal{A}^{(2)}$  and parameter visualizations can be found in Appendix Fig. 10, 11, 12, and 13. These RPE parameters are consistent with the construction in Eq. (7). Moreover, the attention weights  $\mathcal{A}^{(2)}$  of the transformer yields accurate predictions of parent indices across many examples. Similar to the discrete case, we can define the transition  $p(\cdot | \cdot)$  by  $\mathbf{x}_h | \mathbf{x}_{pa(h)} \sim \mathcal{N}(\frac{1}{c} \mathbf{A}^\top \mathbf{x}_{pa(h)}, \frac{\sigma^2}{c^2} I_d)$ . Consequently, Eq. (2) specifies the BMA formulation under the dynamical system setting. In this context, Lemma 3 remains valid and guarantees the asymptotic correctness of BMA’s parent selection. To investigate transformers’ mechanism of parent selection, we test the parent selection loss  $\mathcal{L}_{pa}^{L'}$  of the transformer and BMA in dynamical setting, where we set various  $L'$  in-context samples as introduced in Sec. 3.3. Fig. 15 demonstrates that the transformer with trainable  $(\mathbf{W}, \mathbf{W}_{OV})$  achieves performance comparable to BMA method when  $L'$  approaches 20. However, the loss  $\mathcal{L}_{pa}^{L'}$  between transformers and BMA remains a noticeable gap. We conjecture that the proposition below may explain this discrepancy:

**Proposition 2 (Representation Limitation of Transformers).** *Under the observation restriction in Eq. (7), both the transformer and BMA take the unified form  $\mathcal{A}_{h \rightarrow h'} = \sigma(\mathbf{p}^h)_{h'}$ . In the DS setting, transformer logits are bilinear,  $\mathbf{p}_{\text{tf}, h'}^h = \sum_l \mathbf{x}_{h'}^{l\top} \mathbf{W}_{\text{tf}} \mathbf{x}_h^l$ , whereas BMA logits are  $\mathbf{p}_{\text{BMA}, h'}^h = c_1 \sum_l \mathbf{x}_{h'}^{l\top} \mathbf{A} \mathbf{x}_h^l + d \sum_l \|\mathbf{x}_{h'}^l\|^2$  with  $d \neq 0$ . There exists no  $\mathbf{W}_{\text{tf}}$  such that  $\sigma(\mathbf{p}_{\text{tf}}^h) = \sigma(\mathbf{p}_{\text{BMA}}^h)$  holds for all DS samples  $(\mathbf{x}_{1:H}^{1:L+1}, \mathcal{G}) \sim P_\pi$ . Hence transformers under Eq. (7) cannot represent BMA in the DS setting. However, for MC setting,  $\mathbf{W}_{\text{tf}} = \log \pi$  gives the BMA form.*

For BMA,  $(\mathbf{p}_{\text{BMA}}^h)_{h'} = \sum_l \log p(\mathbf{x}_h^l | \mathbf{x}_{h'}^l)$ . In dynamical systems, transition  $p(\cdot | \cdot)$  involves not only cross but also quadratic terms. If the representation equation holds, substituting  $\mathbf{x}_h = \frac{1}{c} \mathbf{A}^\top \mathbf{x}_{pa(h)} + \varepsilon_h$  and using the independence of  $\varepsilon_h$  forces the coefficient of  $\varepsilon_h$  to vanish, which implies  $\mathbf{W} = c_1 \mathbf{A}$ , contradicting the original representation. So no matrix  $\mathbf{W}$  can yield  $\mathcal{A}_{h \cdot}^{(2)} = \sigma(\hat{\mathbf{p}}^h(\mathbf{W}))$  as in Eq. (18) to represent BMA. See Appendix C.7 for detailed proof.

## REFERENCES

- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning, 2023. URL <https://arxiv.org/abs/2306.00297>. 13
- Kabir Ahuja, Madhur Panwar, and Navin Goyal. In-context learning through the bayesian prism. *arXiv preprint arXiv:2306.04891*, 2023. 13
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models, 2023. URL <https://arxiv.org/abs/2211.15661>. 1
- Zeyuan Allen-Zhu and Yanzhi Li. Physics of Language Models: Part 1, Learning Hierarchical Language Structures. *SSRN Electronic Journal*, May 2023. Full version available at <https://ssrn.com/abstract=5250639>. 1, 14
- Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection, 2023. URL <https://arxiv.org/abs/2306.04637>. 1
- Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint, 2023. URL <https://arxiv.org/abs/2306.00802>. 1
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>. 1
- Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Training dynamics of multi-head softmax attention for in-context learning: Emergence, convergence, and optimality. *arXiv preprint arXiv:2402.19442*, 2024a. 1
- Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Unveiling induction heads: Provable training dynamics and feature learning in transformers. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024b. URL <https://openreview.net/forum?id=4fN2REs0Ma>. 2, 13
- Francesco D’Angelo, Francesco Croce, and Nicolas Flammarion. Selective induction heads: How transformers select causal structures in context. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=bnJgzAQjWf>. 2, 7, 8, 13
- Ezra Edelman, Nikolaos Tsilivis, Benjamin L. Edelman, Eran Malach, and Surbhi Goel. The evolution of statistical induction heads: In-context learning markov chains. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 64273–64311. Curran Associates, Inc., 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/file/75b0edb869e2cd509d64d0e8ff446bc1-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/75b0edb869e2cd509d64d0e8ff446bc1-Paper-Conference.pdf). 1, 2, 13
- Dan Friedman, Alexander Wettig, and Danqi Chen. Learning transformer programs, 2023. URL <https://arxiv.org/abs/2306.01128>. 4
- Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes, 2023. URL <https://arxiv.org/abs/2208.01066>. 1
- Gautam Goel and Peter Bartlett. Can a transformer represent a kalman filter?, 2024. URL <https://arxiv.org/abs/2312.06937>. 1

- Tianyu Guo, Wei Hu, Song Mei, Huan Wang, Caiming Xiong, Silvio Savarese, and Yu Bai. How do transformers learn in-context beyond simple functions? a case study on learning with representations. *arXiv preprint arXiv:2310.10616*, 2023. 13
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL <https://arxiv.org/abs/1412.6980>. 20
- Ashok Vardhan Makkuva, Marco Bondaschi, Adway Girish, Alliot Nagle, Martin Jaggi, Hyeji Kim, and Michael Gastpar. Attention with markov: A framework for principled analysis of transformers via markov chains, 2025. URL <https://arxiv.org/abs/2402.04161>. 13
- Eshaan Nichani, Alex Damian, and Jason D. Lee. How transformers learn causal structure with gradient descent, 2024. URL <https://arxiv.org/abs/2402.14735>. 1, 2, 8, 13, 14, 17, 18, 26, 33
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. In-context learning and induction heads, 2022. URL <https://arxiv.org/abs/2209.11895>. 6, 14
- Yury Polyanskiy and Yihong Wu. *Information Theory: From Coding to Learning*. Cambridge University Press, 2023. ISBN 9781108832908. doi: 10.1017/9781108966351. 9
- Nived Rajaraman, Marco Bondaschi, Kannan Ramchandran, Michael Gastpar, and Ashok Vardhan Makkuva. Transformers on markov data: Constant depth suffices, 2024. URL <https://arxiv.org/abs/2407.17686>. 1
- Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations, 2018. URL <https://arxiv.org/abs/1803.02155>. 5
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2023. URL <https://arxiv.org/abs/2104.09864>. 5
- Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent, 2023. URL <https://arxiv.org/abs/2212.07677>. 1, 13
- Jiuqi Wang, Ethan Blaser, Hadi Daneshmand, and Shangdong Zhang. Transformers can learn temporal difference methods for in-context reinforcement learning, 2025. URL <https://arxiv.org/abs/2405.13861>. 1
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=gEZrGCozdqR>. 1
- Kevin Christian Wibisono and Yixin Wang. From unstructured data to in-context learning: Exploring what tasks can be learned and when, 2024. URL <https://arxiv.org/abs/2406.00131>. 1, 14
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021. 13, 14
- Yufeng Zhang, Fengzhuo Zhang, Zhuoran Yang, and Zhaoran Wang. What and how does in-context learning learn? bayesian model averaging, parameterization, and generalization. *arXiv preprint arXiv:2305.19420*, 2023. 13
- Haoyu Zhao, Abhishek Panigrahi, Rong Ge, and Sanjeev Arora. Do transformers parse while predicting the masked word?, 2023. URL <https://arxiv.org/abs/2303.08117>. 1

## A NOTATION AND RELATED WORK

**Notation.** We use  $[h]$  to denote the set  $\{1, 2, \dots, h\}$ . For causal structure, we use  $pa(h)$  to represent the parent index of node  $h$ . The stationary distribution of Markov chain  $\mathbf{x}_h \sim \pi(\cdot | \mathbf{x}_{pa(h)})$  is denoted by  $\mu^\pi \in \Delta^d$ . For transformer model, the input of a sequence of vectors is given by  $\mathbf{x}_{1:T} := [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T] \in \mathbb{R}^{d \times T}$ . Given the input, we denote the attention scores of standard self-attention layer as  $\mathbf{p}^t := \mathbf{x}_{1:T}^\top \mathbf{W}_K^\top \mathbf{W}_Q \mathbf{x}_t \in \mathbb{R}^T$ . However the causal mask  $\mathcal{M}$  in attention layer will lead to  $\hat{\mathbf{p}}^t := \mathbf{x}_{1:t-1}^\top \mathbf{W}_K^\top \mathbf{W}_Q \mathbf{x}_t \in \mathbb{R}^{t-1}$ ,  $\sigma(\hat{\mathbf{p}}^t)_{t'} = \sigma(\mathcal{M}(\mathbf{p}^t))_{t'}, \forall t' \in [t-1]$ . We do not distinguish between them in the proofs. For the matrix form of the attention of an input sequence, we use  $\tilde{\mathcal{A}}$  and  $\mathcal{A}$  to denote attention weights and scores correspondingly, where we have  $\mathbf{p}_{t'}^t = \mathcal{A}_{t \rightarrow t'}$  and  $\sigma(\tilde{\mathcal{A}}) = \mathcal{A}$ . In training, we use cross-entropy loss and MSE loss for Markov chain and dynamical system settings respectively:

$$\begin{aligned} \mathcal{L}^{MC}(\theta) &= -\frac{1}{H} \sum_{h=2}^H \mathbf{x}_h^{L+1}^\top \log(\sigma(\mathbf{f}_{\text{tf}}(\cdot | \mathcal{H}_h)) + \epsilon), \\ \mathcal{L}^{DS}(\theta) &= -\frac{1}{H} \sum_{h=2}^H \|\mathbf{x}_h^{L+1} - \mathbf{f}_{\text{tf}}(\cdot | \mathcal{H}_h)\|_2^2, \end{aligned} \quad (15)$$

where  $\theta$  represents all trainable parameters and  $\epsilon$  is a small value to avoid numerical issues by log.

**Related Work.** A growing body of work studies the in-context learning (ICL) ability from different perspectives. One line of work understands ICL as a form of Bayesian inference, showing how the latent concept can be approximately inferred under restrictive theoretical assumptions (Xie et al., 2021; Zhang et al., 2023; Ahuja et al., 2023). Another direction of research investigates how transformers can simulate standard algorithms, such as gradient descent on linear regression (von Oswald et al., 2023; Ahn et al., 2023; Guo et al., 2023). While these works demonstrate the ICL power of transformers, they commonly assume i.i.d or uncorrelated input tokens. To move beyond i.i.d. assumptions, recent works investigate ICL with *correlated data*, particularly Markovian sequences (Edelman et al., 2024; Chen et al., 2024b; Makkuva et al., 2025). These settings provide insight into how transformers handle in-context learning with sequential dependencies, but typically focus on fixed dependency structures. In contrast, our work addresses *variable causal structures* that differ across prompts. Pioneering this direction, Nichani et al. (2024) demonstrated that transformers can encode fixed parent-child dependencies (e.g., bigrams) in Markov chains. Building on this, D’Angelo et al. (2025) introduced *selective induction heads*, enabling transformers to identify the underlying Markovian order (or “lag”) from a candidate set in-context learning this structure. Our work generalizes this setting. While D’Angelo et al. (2025) focus on inferring a single global structural parameter (the lag  $k$ ) shared across the sequence, we tackle local structure inference where dependencies can vary arbitrarily for each position, effectively modeling latent trees rather than fixed-lag chains. Theoretically, D’Angelo et al. (2025) constructs a three-layer transformer that asymptotically implements maximum likelihood estimation, where its construction is verified via attention pattern visualization as well as quantitative validation through KL divergence of next-token prediction targets. In our work, we theoretically derive a two-layer architecture that explicitly implements Bayesian Model Averaging (BMA) in-context. Empirically, we go beyond behavioral metrics by providing *parameter-level* verification, demonstrating that the trained weights directly encode the transition kernel. Furthermore, we provide theoretical understanding of in-context causal structure learning based on the Data Processing Inequality (DPI) and extend our analysis to continuous dynamical systems, revealing representational gaps not occurring in the discrete setting.

## B CONCLUSION

In this work, we investigated the capability of transformers to infer and adapt to latent causal structures in-context, moving beyond the fixed dependency assumptions common in prior theoretical analysis. We proposed a novel framework based on Markov chains with randomly sampled causal dependencies, requiring the model to identify position-specific predecessor-successor relationships from context examples. First, we provided a constructive proof that a two-layer transformer with relative position embeddings (RPE) can explicitly implement Bayesian Model Averaging (BMA). This demonstrates that the attention mechanism is theoretically capable of performing statistical



inference over structural uncertainty. Second, through extensive experiments and parameter-level analysis, we showed that trained transformers implements BMA method which converge to this theoretical construction: the learned attention patterns directly encode the posterior probabilities of causal parents, and the weights explicitly recover the log-transition kernel of the underlying generative process. Third, we established information-theoretic guarantees using the Data Processing Inequality (DPI) which help understand how the selection mechanism identifies causal structures in context and showed that gradients at initialization recover these dependencies via  $\chi^2$ -mutual information. Finally, we extended our framework to continuous linear dynamical systems. While transformers continue to exhibit strong empirical performance in this setting, we identified the representational difference that prevents the exact implementation of BMA, unlike in the discrete case. Collectively, our findings offer a mechanistic explanation of how transformers perform in-context causal learning, highlighting their ability to act as statistical inference engines for both discrete and continuous data.

**Broader Implications.** Our findings support theoretical frameworks that model in-context learning as a statistical inference task (Xie et al., 2021). Distinct from "Induction Heads" which typically focus on copying fixed positional dependencies (Olsson et al., 2022; Nichani et al., 2024), we demonstrate a probabilistic setting where the model must infer a latent dependency structure that varies per example. This provides a mechanistic grounding for how LLMs adapt to flexible, context-dependent rules rather than relying solely on fixed n-gram statistics (Allen-Zhu & Li, 2023). Furthermore, this helps understand why LLMs demonstrate ICL capabilities on empirical task with "unstructured" language data (Wibisono & Wang, 2024), mirroring our setting where the transition mappings between words are fixed while the structural positions of a couple of words vary from input to input.

**Limitations and Future Work.** We acknowledge that real-world sequences often involve complex non-linear dynamics or hierarchical dependencies (e.g., context-free grammar) beyond the Markovian and dynamical systems studied here. However, our primary objective in this work was to prioritize mechanistic interpretability for Markov chain or dynamical system: explicitly characterizing how transformers infer latent structures in-context on these tasks. By focusing on these tractable settings, we were able to derive exact theoretical guarantees and provide parameter-level verification that the model implements Bayesian Model Averaging. We believe this explainable framework serves as a necessary foundation, and we leave the extension to more complex non-linear and hierarchical data generating processes for future exploration.

## C PROOFS OF TECHNICAL LEMMAS

### C.1 PROOF OF THEOREM 1

*Proof.* By the condition of  $(\tilde{\mathbf{w}}_H, \tilde{\mathbf{w}}_L)$  in Eq. (7), the attention score  $\hat{\mathcal{A}}_{t \rightarrow \cdot}^{(1)}$  of query  $\mathbf{x}_t = \mathbf{x}_h^{L+1}$  is:

$$\hat{\mathcal{A}}_{t \rightarrow t'}^{(1),k} = \tilde{\mathbf{w}}_H^k[h_t - h_{t'}] + \tilde{\mathbf{w}}_L^k[l_t - l_{t'}] = 2\beta \begin{cases} +1, & \text{if } h_t = h_{t'}, l_t - l_{t'} = k, \\ -1, & \text{if } h_t \neq h_{t'}, l_t - l_{t'} \neq k, \\ 0, & \text{otherwise.} \end{cases}$$

Then the output  $\tilde{\mathbf{u}}_h^k = \text{Attn}_{\mathbf{x}_t \rightarrow \mathbf{x}_{1:T}}^k$  of the 1st attention layer will be calculated as:

$$\begin{aligned} \tilde{\mathbf{u}}_h^k &= \sigma(\tilde{\mathbf{w}}_H^k(h, \cdot) + \tilde{\mathbf{w}}_L^k(L+1, \cdot)) \mathbf{x}_{1:T}^\top \\ &\xrightarrow{\beta \rightarrow \infty} \mathbf{x}_h^{l_k \top} = (1_{[h_{t'}=h, l_{t'}=L+1-k]})_{t' \in [T]} \mathbf{x}_{1:T}^\top. \quad (l_k = L+1-k) \end{aligned} \quad (16)$$

By disentangled residual, the output of  $K$  heads ( $K = L$ ) will be concatenated as:

$$\mathbf{v}_h = [\mathbf{u}_h^1, \dots, \mathbf{u}_h^L], \text{ with } \mathbf{u}_h^k = \mathbf{x}_h^{L+1-k} \text{ by Eq. (16).}$$

For the 2nd layer, with diagonal condition of  $\mathbf{W}_{KQ}$ , attention weight  $\mathcal{A}^{(2)} \in \mathbb{R}^{H \times H}$  is given by:

$$\hat{\mathcal{A}}_{h \rightarrow h'}^{(2)} = \mathbf{v}_{h'}^\top \tilde{\mathbf{W}}_{KQ} \mathbf{v}_h = \sum_{l=1}^L \mathbf{x}_{h'}^{l \top} \mathbf{W} \mathbf{x}_h^l, \quad \mathcal{A}^{(2)} = \sigma(\mathcal{M}(\hat{\mathcal{A}}^{(2)})) \in \mathbb{R}^{H \times H}, \quad (17)$$

where  $\mathcal{M}$  is the causal mask enforcing  $\mathcal{A}^{(2)}$  to be strictly lower-triangular after softmax. If we define the vector  $\hat{\mathbf{p}}^h \in \mathbb{R}^{h-1}$  with  $\hat{\mathbf{p}}_{h'}^h := \hat{\mathcal{A}}_{h \rightarrow h'}^{(2)}$ , we have  $\forall h' \in [h-1]$ :

$$\mathcal{A}_{h \rightarrow h'}^{(2)} = \sigma(\mathcal{M}_h(\hat{\mathcal{A}}_{h \rightarrow \cdot}^{(2)}))_{h'} = \sigma(\hat{\mathbf{p}}^h)_{h'}, \quad \hat{\mathbf{p}}^h(\mathbf{W}) = \sum_l \mathbf{x}_{1:h-1}^{l \top} \mathbf{W} \mathbf{x}_h^l, \quad (18)$$

where  $\mathcal{M}_h(\cdot)$ , is the causal mask applied to row  $h$  setting  $\mathcal{M}_h(v)_{h'} = -\infty$  if  $h' \geq h, \forall v \in \mathbb{R}^H$ . Then we set  $\mathbf{W}$  as  $\log \pi$  (elementwise) which leads to

$$\hat{p}_{h'}^h(\log \pi) = \hat{\mathcal{A}}_{h \rightarrow h'}^{(2)} = \sum_l \log \pi(\mathbf{x}_h^l | \mathbf{x}_{h'}^l), \quad \mathcal{A}_{h \rightarrow h'}^{(2)} = \sigma(\mathcal{M}_h(\hat{\mathcal{A}}_{h \rightarrow h'})) \in \mathbb{R}^H. \quad (19)$$

With the form of  $\mathcal{A}^{(2)}$  in Eq. (19) and Lemma 1 to be proved, we can show the BMA method of Eq. (1) has the same formulation:

$$\mathbb{P}(pa(h) = h' | \mathbf{x}_{1:H}^{1:L}) = \sigma(\hat{p}^h(\mathbf{W} = \log \pi))_{h'}. \quad (20)$$

Combining Eq. (20) with the limit behavior of the first layer in Eq. (16), we obtain the first convergence result for parent selection as  $\beta \rightarrow \infty$ :

$$\lim_{\beta \rightarrow \infty} \mathcal{A}_{h \rightarrow h'}^{(2)}(\mathcal{H}; \theta) = \sigma(\hat{p}^h(\mathbf{W} = \log \pi)) = \sigma(p_{\text{BMA}}^{h,L}).$$

Furthermore, guaranteed by the consistency of BMA (Theorem 2), as sample size  $L \rightarrow \infty$ , the posterior concentrates on the true parent  $pa(h)$ . Thus, the prediction of the token distribution converges in the limit  $\beta, L \rightarrow \infty$  as:

$$\lim_{\beta, L \rightarrow \infty} f_{\theta}(\cdot | \mathcal{H}) = \sigma \left( \mathbf{W}_{OV}^{\top} \sum_{h'} 1_{[h'=pa(h)]} \mathbf{x}_{h'}^{L+1} \right) = \pi(\cdot | \mathbf{x}_{pa(h)}^{L+1}). \quad \square$$

In proving the theorem, we rely on the Lemma 1 proved below, which illustrates the relation between BMA and attention weights.

## C.2 PROOF OF LEMMA 1

*Proof.* Here we use  $p(s|s')$  to denote  $\mathbb{P}(\mathbf{x}_h = s | \mathbf{x}_{pa(h)} = s')$  for generality beyond discrete Markov chain. Based on Bayesian Theorem, it can be calculated by Eq. (1). Due to the Markovian property  $p(\mathbf{x}_h | \mathbf{x}_{1:h-1}) = p(\mathbf{x}_h | \mathbf{x}_{pa(h)})$ , the joint distribution of this chain  $\mathbf{x}_{1:H}$ :

$$\begin{aligned} p(\mathbf{x}_{1:H}) &= p(\mathbf{x}_1) \prod_{h=2}^H p(\mathbf{x}_h | \mathbf{x}_{1:h-1}) = p(\mathbf{x}_1) \prod_{h=2}^H p(\mathbf{x}_h | \mathbf{x}_{pa(h)}) \\ &= p(\mathbf{x}_1) \prod_{i \neq h} p(\mathbf{x}_i | \mathbf{x}_{pa(i)}) \cdot p(\mathbf{x}_h | \mathbf{x}_{pa(h)}) \end{aligned} \quad (21)$$

Here  $pa(h), pa(i)$  in Eq. (21) are random index with prior:  $pa(h) \sim \text{Uniform}([h-1])$ . With condition  $pa(h) = h$  in Eq. (1), we can substitute  $pa(h)$  in Eq. (21) with  $h'$ . Since  $\{pa(h'')\}$  are random index out of interests, these terms are eliminated:

$$\frac{\mathbb{P}(\mathbf{x}_{1:H}^{1:L} | pa(h) = h')}{\sum_{h'' \in [h-1]} \mathbb{P}(\mathbf{x}_{1:H}^{1:L} | pa(h) = h'')} = \frac{\prod_l (p(\mathbf{x}_1^l) \prod_{i \neq h} p(\mathbf{x}_i^l | \mathbf{x}_{pa(i)}^l) \cdot p(\mathbf{x}_h^l | \mathbf{x}_{h'}^l))}{\sum_{h''} \prod_l (p(\mathbf{x}_1^l) \prod_{i \neq h} p(\mathbf{x}_i^l | \mathbf{x}_{pa(i)}^l) \cdot p(\mathbf{x}_h^l | \mathbf{x}_{h''}^l))}, \quad (22)$$

which leads to:

$$\mathbb{P}(pa(h) = h' | \mathbf{x}_{1:H}^{1:L}) = \frac{\exp(\sum_{l \in [L]} \log p(\mathbf{x}_h^l | \mathbf{x}_{h'}^l))}{\sum_{h'' \in [h-1]} \exp(\sum_{l \in [L]} \log p(\mathbf{x}_h^l | \mathbf{x}_{h''}^l))} = \sigma(\hat{p}^h(\log W^P))_{h'}, \quad (23)$$

where the matrix  $W^P = \pi$  is induced by transition kernel  $P = \pi$  in Markov chain.  $\square$

## C.3 PROOF OF PROPOSITION 1

First, suppose we have  $\mathbf{W}_{\text{tf}} = \log \pi + \mathbf{1}\mathbf{a}^{\top}$ . The attention weights of the transformer are given by:

$$\mathbf{p}_{\text{tf}}^h = \sum_l \mathbf{x}_{1:h-1}^{\top} \mathbf{W}_{\text{tf}} \mathbf{x}_h^l = \sum_l \mathbf{x}_{1:h-1}^{\top} \log \pi \mathbf{x}_h^l + \sum_l \mathbf{x}_{1:h-1}^{\top} \mathbf{1}\mathbf{a}^{\top} \mathbf{x}_h^l.$$

For the second term, since  $\{\mathbf{x}_{h'}\}$  are one-hot, we have

$$\sum_l \mathbf{x}_{1:h-1}^{\top} \mathbf{1}\mathbf{a}^{\top} \mathbf{x}_h^l = \mathbf{1}_{h-1} \mathbf{a}^{\top} \left( \sum_l \mathbf{x}_h^l \right) = c(\mathbf{a}, h) \mathbf{1}_{h-1},$$

where  $c(\mathbf{a}, h) = \mathbf{a}^\top \left( \sum_l \mathbf{x}_h^l \right)$  is a constant with fixed index  $h$ . Then by softmax operation, we have:

$$\begin{aligned} \sigma(\mathbf{p}_{\text{tf}}^h) &= \sigma \left( \sum_l \mathbf{x}_{1:h-1}^{l\top} \log \pi \mathbf{x}_h^l + c(\mathbf{a}, h) \mathbf{1}_{h-1} \right) \\ &= \sigma \left( \sum_l \mathbf{x}_{1:h-1}^{l\top} \log \pi \mathbf{x}_h^l \right) \\ &= \sigma(\mathbf{p}_{\text{BMA}}^h(\log \pi)), \end{aligned}$$

where the last equality comes from Lemma 1. And this shows the transformer with  $\mathbf{W} + \mathbf{1}\mathbf{a}^\top$  gives the same prediction of BMA's. Further, suppose  $\mathbf{W} = \log \pi + \mathbf{1}\mathbf{a}^\top + \mathbf{b}\mathbf{1}^\top$ . If we have  $\mathbf{x}_h \sim \mu^\pi, \forall h \in [H]$ , then we can prove the term from  $\mathbf{b}\mathbf{1}^\top$  is also a constant vector asymptotically:

$$\begin{aligned} \mathbf{p}_{\text{tf}}^h &= \sum_l \mathbf{x}_{1:h-1}^{l\top} \mathbf{W}_{\text{tf}} \mathbf{x}_h^l \\ &= \sum_l \mathbf{x}_{1:h-1}^{l\top} \log \pi \mathbf{x}_h^l + \sum_l \mathbf{x}_{1:h-1}^{l\top} \mathbf{1}\mathbf{a}^\top \mathbf{x}_h^l + \sum_l \mathbf{x}_{1:h-1}^{l\top} \mathbf{b}\mathbf{1}^\top \mathbf{x}_h^l \\ &= \sum_l \mathbf{x}_{1:h-1}^{l\top} \log \pi \mathbf{x}_h^l + c(\mathbf{a}, h) \mathbf{1}_{h-1} + \sum_l \mathbf{x}_{1:h-1}^{l\top} \mathbf{b}. \end{aligned}$$

For each term of  $\sum_l \mathbf{x}_{1:h-1}^{l\top} \mathbf{b}$ , we have

$$\begin{aligned} \frac{1}{L} \sum_l \mathbf{x}_{h'}^{l\top} \mathbf{b} &= \frac{1}{L} \sum_l \sum_{s \in [d]} \mathbf{1}_{[\mathbf{x}_{h'}^l = e_s]} \mathbf{b}_s \\ &\xrightarrow{L \rightarrow \infty} \mathbb{E} \left[ \sum_{s \in [d]} \mathbf{1}_{[\mathbf{x}_{h'}^l = e_s]} \mathbf{b}_s \right] = \sum_{s \in [d]} \mathbb{P}(\mathbf{x}_{h'}^l = e_s) \mathbf{b}_s \\ &= \mu^{\pi^\top} \mathbf{b}, \end{aligned}$$

where  $\mu^{\pi^\top} \mathbf{b}$  is a constant  $d(\mathbf{b})$  w.r.t.  $h'$ . Using the same technique in Theorem 2 to eliminate this term which goes to infinity, we have the desired result:

$$\lim_{L \rightarrow \infty} \sigma(\mathbf{p}_{\text{tf}}^{h,L}) = \lim_{L \rightarrow \infty} \sigma(\mathbf{p}_{\text{BMA}}^{h,L}(\log \pi)).$$

#### C.4 PROOF OF LEMMA 3

*Proof.* Noting that if  $p(\cdot)$  and  $q(\cdot)$  are two distribution, then by KL divergence's non-negativity we have:  $\int_s p(s) \log q(s) \leq \int_s p(s) \log p(s)$ . Hence we can get:

$$\begin{aligned} \text{LHS} &= \int_{s, s'} \mathbb{P}(\mathbf{x}_h = s, \mathbf{x}_{h'} = s') \log \mathbb{P}(\mathbf{x}_h = s | \mathbf{x}_{pa(h)} = s') \\ &= \int_{s'} \mathbb{P}(\mathbf{x}_{h'} = s') \int_s \mathbb{P}(\mathbf{x}_h = s | \mathbf{x}_{h'} = s') \log \mathbb{P}(\mathbf{x}_h = s | \mathbf{x}_{pa(h)} = s') \\ &\leq \int_{s, s'} \mathbb{P}(\mathbf{x}_{h'} = s') H(\mathbf{x}_h | \mathbf{x}_{h'} = s') = -H(\mathbf{x}_h | \mathbf{x}_{h'}), \end{aligned} \tag{24}$$

where  $H(\mathbf{x}_h | \mathbf{x}_{h'}) = H(\mathbf{x}_h) - I(\mathbf{x}_h; \mathbf{x}_{h'})$ . By Corollary. 1, we have:

$$\text{LHS} \leq H(\mathbf{x}_h) - I(\mathbf{x}_h; \mathbf{x}_{h'}) < H(\mathbf{x}_h) - I(\mathbf{x}_h; \mathbf{x}_{pa(h)}) = \text{RHS}. \tag{25}$$

□

#### C.5 PROOF OF THEOREM 2

*Proof.* Recall that the transformer and BMA have the formula in Eq. (2):

$$\tilde{\mathcal{A}}_{h \rightarrow h'}^L = \frac{\exp(\sum_l \log p(\mathbf{x}_h^l | \mathbf{x}_{h'}^l))}{\sum_{h'' \in [h-1]} \exp(\sum_l \log p(\mathbf{x}_h^l | \mathbf{x}_{h''}^l))} = \frac{1}{\sum_{h''} \mathbf{v}_{h'' \rightarrow h'}}. \tag{26}$$

By the central limit theorem, we have

$$\lim_{L \rightarrow \infty} \frac{1}{L} \sum_l \log p(\mathbf{x}_h^l | \mathbf{x}_{h'}^l) = \mathbb{E}[\log p(\mathbf{x}_h | \mathbf{x}_{h'})] < \mathbb{E}[\log p(\mathbf{x}_h | \mathbf{x}_{pa(h)})] = \lim_{L \rightarrow \infty} \frac{1}{L} \sum_l \log p(\mathbf{x}_h^l | \mathbf{x}_{pa(h)}^l).$$

Let  $\hat{g}_{h,h'} \triangleq \frac{1}{L} \sum_l \log p(\mathbf{x}_h^l | \mathbf{x}_{h'}^l)$ . For all  $h'' \neq pa(h)$ , we have:

$$\mathbf{v}_{h'' \rightarrow pa(h)} = \exp \left( \sum_l \log p(\mathbf{x}_h^l | \mathbf{x}_{h''}^l) - \sum_l \log p(\mathbf{x}_h^l | \mathbf{x}_{pa(h)}^l) \right) = \exp \left( L(\hat{g}_{h,h''} - \hat{g}_{h,pa(h)}) \right) \rightarrow 0$$

as  $L \rightarrow \infty$  and  $\lim_{L \rightarrow \infty} (\hat{g}_{h,h''} - \hat{g}_{h,pa(h)}) < 0$ . Hence, we have  $\lim_{L \rightarrow \infty} \tilde{A}_{h \rightarrow pa(h)} = 1$ .  $\square$

### C.6 PROOF OF THEOREM 3

*Proof.* First, the transformer as constructed can be simplified as:

$$\mathbf{f}_{\theta}^{(\text{simp})}(\cdot | \mathcal{H}) = \pi^\top \mathbf{x}_{1:h-1} \sigma \left( \sum_l \mathbf{x}_{1:h-1}^\top \mathbf{W} \mathbf{x}_h^l \right) \in \mathbb{R}^d, \quad (27)$$

Considering  $\hat{\mathbf{p}} = \sum_l \mathbf{x}_{1:h-1}^\top \mathbf{W} \mathbf{x}_h^l = \mathbf{0}$  when  $\mathbf{W} = \mathbf{0}$ , then  $\mathbf{p} = \frac{1}{h-1} \mathbf{1}_{h-1}$  and:

$$\mathbf{f}_{\theta_0}(\cdot | \mathcal{H}) = \pi^\top \bar{\mu}(\mathbf{x}_{1:h-1}), \text{ where } \bar{\mu}(\mathbf{x}_{1:h-1}) = \frac{1}{h-1} \sum_{h' \in [h-1]} \mathbf{x}_{h'}. \quad (28)$$

Then based on  $\frac{\partial \sigma(\hat{\mathbf{p}})}{\partial \hat{\mathbf{p}}} = \text{diag}(\mathbf{p}) - \mathbf{p} \mathbf{p}^\top$ , computing the gradient of  $\mathbf{W}$  w.r.t loss  $\ell$  in Eq. (14) yields:

$$\begin{aligned} \frac{\partial \ell(\Theta; h, \mathcal{G})}{\partial \hat{\mathbf{p}}} &= \mathbb{E}_{\mathbf{X}} \left[ \left( \frac{\mathbf{x}_h}{\mathbf{f}_{\theta_0}(\mathbf{x}_h) + \epsilon} \right)^\top \frac{\partial \mathbf{f}_{\theta_0}}{\partial \hat{\mathbf{p}}} \right] \\ &= \mathbb{E}_{\mathbf{X}} \left[ \left( \frac{\mathbf{x}_h}{\mathbf{f}_{\theta_0}(\mathbf{x}_h) + \epsilon} \right)^\top \frac{1}{h-1} (\pi^\top \mathbf{x}_{1:h-1} - \pi^\top \bar{\mu}(\mathbf{x}_{1:h-1}) \mathbf{1}_{h-1}^\top) \right] \\ &\stackrel{\text{Eq. (28)}}{=} \frac{1}{h-1} \mathbb{E}_{\mathbf{X}} \left[ \left( \frac{\mathbf{x}_h}{\mathbf{f}_{\theta_0}(\mathbf{x}_h) + \epsilon} \right)^\top (\pi^\top \mathbf{x}_{1:h-1} - \mathbf{f}_{\theta_0}(\mathbf{x}_h) \mathbf{1}_{h-1}^\top) \right] \\ &= \frac{1}{h-1} \mathbb{E}_{\mathbf{X}} \left[ \left[ \frac{\pi(\mathbf{x}_h | \mathbf{x}_1)}{\mathbf{f}_{\theta_0}(\mathbf{x}_h) + \epsilon}, \dots, \frac{\pi(\mathbf{x}_h | \mathbf{x}_{h-1})}{\mathbf{f}_{\theta_0}(\mathbf{x}_h) + \epsilon} \right] - \mathbf{1}_{h-1}^\top \right] \\ &= \frac{1}{h-1} \mathbb{E}_{\mathbf{X}} \left[ \left[ \frac{\pi(\mathbf{x}_h | \mathbf{x}_1)}{\mu^\pi(\mathbf{x}_h)}, \dots, \frac{\pi(\mathbf{x}_h | \mathbf{x}_{h-1})}{\mu^\pi(\mathbf{x}_h)} \right] - \mathbf{1}_{h-1}^\top \right] \in \mathbb{R}^{h-1}. \end{aligned}$$

Then let  $\hat{g}_{h'}^h$  denote  $h'$ -th entry in  $\frac{\partial \ell(\theta_0; h, \mathcal{G})}{\partial \hat{\mathbf{p}}} \in \mathbb{R}^{h-1}$  ( $h' \in [h-1]$ ), we have:

$$\hat{g}_{h'}^h = \frac{1}{h-1} \mathbb{E}_{\mathbf{X}} \left[ \frac{\pi(\mathbf{x}_h | \mathbf{x}_{h'}^l)}{\mu^\pi(\mathbf{x}_h)} - 1 \right] = \frac{1}{h-1} \left( \sum_{s, s'} \frac{\pi(s | s') \mathbb{P}_{\mathbf{X}}(\mathbf{x}_h = s | \mathbf{x}_{h'} = s')}{\mu^\pi(s)} - 1 \right). \quad (29)$$

By Cauchy-Schwartz Inequality and Data Processing Inequality, we have:

$$\begin{aligned} \mathbb{E}_{\mathbf{X}} \left[ \frac{\pi(\mathbf{x}_h | \mathbf{x}_{h'}^l)}{\mu^\pi(\mathbf{x}_h)} - 1 \right] &= \sum_{s, s'} \frac{\pi(s | s') \mathbb{P}_{\mathbf{X}}(\mathbf{x}_h = s | \mathbf{x}_{h'} = s')}{\mu^\pi(\mathbf{x}_h)} - 1 \\ &\leq \frac{1}{2} (I_{\chi^2}(\mathbf{x}_h, \mathbf{x}_{h'}) + I_{\chi^2}(\mathbf{x}_h, \mathbf{x}_{pa(h)})) \leq I_{\chi^2}(\mathbf{x}_h, \mathbf{x}_{pa(h)}) = \mathbb{E}_{\mathbf{X}} \left[ \frac{\pi(\mathbf{x}_h | \mathbf{x}_{pa(h)})}{\mu^\pi(\mathbf{x}_h)} - 1 \right]. \end{aligned} \quad (30)$$

Eq. (30) has shown the desired result  $\hat{g}_{h'}^h \leq \hat{g}_{pa(h)}^h$ .  $\square$

<sup>3</sup>Suppose  $\hat{\mu}_{\mathbf{X}} := \pi^\top \bar{\mu} = \mathbf{f}_{\theta_0}(\mathbf{x}_h)$ . If we remove the assumption  $\mathbf{f}_{\theta} = \pi^\top \bar{\mu} = \mu^\pi$ . Lemma 24 in Nichani et al. (2024) shows  $\left| \mathbb{E}_{\mathbf{X}} \left[ \frac{\pi(\mathbf{x}_h | \mathbf{x}_{h'}^l)}{\hat{\mu}_{\mathbf{X}}(\mathbf{x}_h) + \epsilon} - 1 \right] - \mathbb{E}_{\mathbf{X}} \left[ \frac{\pi(\mathbf{x}_h | \mathbf{x}_{h'}^l)}{\mu^\pi(\mathbf{x}_h)} - 1 \right] \right| \lesssim \frac{1}{\sqrt{T_{\text{eff}}}} (\rightarrow 0)$ . Under Assumption 1 and Strong Data Processing Inequality in Nichani et al. (2024) (Lemma 5), we can prove the non-asymptotic result  $\hat{g}_{pa(h)}^h - \hat{g}_{h'}^h \geq \frac{1}{h-1} \left( \frac{\gamma^3}{2S} - \frac{2C}{\sqrt{T_{\text{eff}}(\lambda)}} \right)$ .

**Assumption 1** (Assumptions on transition kernel (Nichani et al. (2024), Assumption 1)). *There exist  $\gamma > 0, \lambda < 1$  such that the following hold for  $\pi$ :*

- (Transition lower bounded):  $\min_{s,s'} \pi(s' | s) > \gamma$ .
- (Non-degeneracy of chain):  $\|B(\pi)\|_F > \gamma$ .
- (Spectral gap): The spectral gap of  $\pi$ ,  $1 - \lambda(\pi)$ , satisfies  $\lambda(\pi) < \lambda$ .
- (Symmetry): For any permutation matrix  $\sigma$  on  $\mathcal{S}$ ,  $\sigma^{-1}\pi\sigma =_d \pi$ .
- (Constant mean):  $\mathbb{E}_\pi[\pi] = \frac{1}{S}1_S1_S^T$ .

### C.7 PROOF OF PROPOSITION 2

*Proof.* In dynamical system setting, the transition  $P(\cdot|\cdot)$  is given by the pmf of  $\mathbf{x}_h|\mathbf{x}_{pa(h)}$ :

$$p(\mathbf{x} | \mathbf{y}) = \frac{1}{(2\pi)^{d/2} \left(\frac{\sigma^2}{c^2}\right)^{d/2}} \exp\left(-\frac{c^2}{2\sigma^2} \left\|\mathbf{x} - \frac{1}{c}\mathbf{A}^\top \mathbf{y}\right\|_2^2\right), \quad \mathbf{A} \in \mathcal{O}(\mathbb{R}^d).$$

Then with  $\sigma$  eliminate constant terms in  $\log p$ , we get the equivalent form in BMA:

$$\begin{aligned} \log p(\mathbf{x}_h | \mathbf{x}_{h'}) &= \frac{c}{\sigma^2} \mathbf{x}_h^\top \mathbf{A}^\top \mathbf{x}_{h'} - \frac{1}{2\sigma^2} \mathbf{x}_{h'}^\top \mathbf{A} \mathbf{A}^\top \mathbf{x}_{h'} + \text{const}(h), \\ \bar{\mathbf{p}}_{h'}^h &:= \sum_{l=1}^L \left( \frac{c}{\sigma^2} \mathbf{x}_{h'}^{l\top} \mathbf{A} \mathbf{x}_h^l - \frac{1}{2\sigma^2} \mathbf{x}_{h'}^{l\top} \mathbf{x}_{h'}^l \right); \quad \mathbb{P}(pa(h)|\mathbf{x}_{1:H}^{1:L}) = \sigma(\bar{\mathbf{p}}^h). \end{aligned} \quad (31)$$

Eq. (31) gives the BMA logits in the DS setting in a softmax form. We now show that transformers under the observation restriction Eq. (7) cannot represent BMA in this setting.

Recall that, under Eq. (7), the transformer logits are

$$\mathbf{p}_{\text{tf},h'}^h = \sum_{l=1}^L \mathbf{x}_{h'}^{l\top} \mathbf{W}_{\text{tf}} \mathbf{x}_h^l,$$

while the BMA logits are

$$\mathbf{p}_{\text{BMA},h'}^h = c_1 \sum_{l=1}^L \mathbf{x}_{h'}^{l\top} \mathbf{A} \mathbf{x}_h^l + d \sum_{l=1}^L \|\mathbf{x}_{h'}^l\|^2, \quad c_1 \neq 0, d \neq 0.$$

Suppose, for contradiction, that the transformer exactly represents BMA, i.e.

$$\sigma(\mathbf{p}_{\text{tf}}^h) = \sigma(\mathbf{p}_{\text{BMA}}^h) \quad \text{for all DS samples and all } h \in [H].$$

Since softmax is invariant under adding a constant independent of  $h'$ , this means that for each fixed  $h$  there exists a scalar  $b = b(h)$  such that

$$\mathbf{p}_{\text{tf},h'}^h + b = \mathbf{p}_{\text{BMA},h'}^h \quad \text{for all } h' \in [H]. \quad (*)$$

Using the DS model  $\mathbf{x}_h^l = \frac{1}{c}(\mathbf{A}^\top \mathbf{x}_{pa(h)}^l + \varepsilon_h^l)$ , we expand the logits as

$$\begin{aligned} \mathbf{p}_{\text{tf},h'}^h &= \underbrace{\frac{1}{c} \sum_{l=1}^L \mathbf{x}_{h'}^{l\top} \mathbf{W}_{\text{tf}} \mathbf{A}^\top \mathbf{x}_{pa(h)}^l}_{\text{constant term w.r.t. } h, \varepsilon_h} + \underbrace{\frac{1}{c} \sum_{l=1}^L \mathbf{x}_{h'}^{l\top} \mathbf{W}_{\text{tf}} \varepsilon_h^l}_{\text{separate term with } \varepsilon_h}, \\ \mathbf{p}_{\text{BMA},h'}^h &= \underbrace{\frac{c_1}{c} \sum_{l=1}^L \mathbf{x}_{h'}^{l\top} \mathbf{A} \mathbf{A}^\top \mathbf{x}_{pa(h)}^l}_{\text{constant term w.r.t. } h, \varepsilon_h} + \underbrace{d \sum_{l=1}^L \|\mathbf{x}_{h'}^l\|^2}_{\text{separate term with } \varepsilon_h} + \underbrace{\frac{c_1}{c} \sum_{l=1}^L \mathbf{x}_{h'}^{l\top} \mathbf{A} \varepsilon_h^l}_{\text{separate term with } \varepsilon_h}. \end{aligned}$$



Conditioning on all variables except  $\{\varepsilon_h^l\}_{l=1}^L$ , both sides of  $(*)$  become affine functions of the Gaussian noises  $\varepsilon_h^l$ . Since the DS distribution has full support and  $(*)$  is assumed to hold for all DS samples, the coefficients of the linear terms in  $\{\varepsilon_h^l\}$  must match for all realizations. Since  $\{\varepsilon_h^l\}_l$  are independently sampled, conditioning on  $\varepsilon_h^l$  will eliminate other terms. Seeing  $\varepsilon_h^l$  as the only free variable, its coefficient should be zero to keep the equation held:

$$\mathbf{x}_{h'}^{l\top}(\mathbf{W}_{\text{tf}} - c_1 \mathbf{A})\varepsilon_h^l = 0, \forall \varepsilon_h^l \in \mathbb{R}^d \Rightarrow \mathbf{x}_{h'}^{l\top}(\mathbf{W}_{\text{tf}} - c_1 \mathbf{A}) = 0.$$

Since in the DS model each  $\mathbf{x}_{h'}^l \sim \mathcal{N}(0, I_d)$  is non-degenerate with full support  $\mathbf{x}_h^l \sim \mathcal{N}(0, I_d)$ , which forces

$$\mathbf{W}_{\text{tf}} = c_1 \mathbf{A}.$$

Substituting  $\mathbf{W}_{\text{tf}} = c_1 \mathbf{A}$  back into  $(*)$ , the representation equation is formulated by

$$b + \frac{c_1}{c} \sum_{l=1}^L \mathbf{x}_{h'}^{l\top} \mathbf{A} \mathbf{A}^\top \mathbf{x}_{pa(h)}^l = \frac{c_1}{c} \sum_{l=1}^L \mathbf{x}_{h'}^{l\top} \mathbf{A} \mathbf{A}^\top \mathbf{x}_{pa(h)}^l + d \sum_{l=1}^L \|\mathbf{x}_{h'}^l\|^2.$$

Hence

$$b = d \sum_{l=1}^L \|\mathbf{x}_{h'}^l\|^2 \quad \text{for all } h' \in [H].$$

However, for a DS sample the quantities  $\sum_l \|\mathbf{x}_{h'}^l\|^2$  vary across  $h'$  and across samples, while  $b$  is a constant (depending only on  $h$ ). The only way the above equality can hold for all  $h'$  and all DS samples is to have  $b = d = 0$ , which contradicts the assumption  $d \neq 0$  in the BMA logits.

We conclude that no  $\mathbf{W}_{\text{tf}}$  can make the transformer logits represent the BMA logits for all datasets generated from DS. Therefore, under Eq. (7), transformers cannot represent BMA in the DS setting.

**C.8 PROOF OF  $\sum_{s,s'} \bar{\mu}_t(s') \mu^\pi(s) \log \pi(s | s') \leq \sum_{s,s'} \mu^\pi(s') \pi(s | s') \log \pi(s | s'), \forall \bar{\mu}_t \in \Delta^d$**

Let  $\mathcal{S}$  be a finite state space, let  $\pi(\cdot | \cdot)$  be a Markov transition kernel on  $\mathcal{S}$ , and let  $\mu^\pi$  be a stationary distribution of  $\pi$ , i.e.,

$$\mu^\pi(s) = \sum_{s' \in \mathcal{S}} \mu^\pi(s') \pi(s | s') \quad \text{for all } s \in \mathcal{S}.$$

Fix an arbitrary  $\bar{\mu}_t \in \Delta^d$ , and for brevity write  $\mu := \mu^\pi$ . We adopt the convention  $0 \log 0 = 0$  throughout.

We first upper bound the left-hand side. For any  $s' \in \mathcal{S}$ , consider the Kullback–Leibler divergence

$$\text{KL}(\mu \| \pi(\cdot | s')) = \sum_{s \in \mathcal{S}} \mu^\pi(s) \log \frac{\mu^\pi(s)}{\pi(s | s')} \geq 0.$$

Expanding the inequality  $\text{KL}(\mu \| \pi(\cdot | s')) \geq 0$  yields

$$\sum_{s \in \mathcal{S}} \mu^\pi(s) \log \pi(s | s') \leq \sum_{s \in \mathcal{S}} \mu^\pi(s) \log \mu^\pi(s) =: C,$$

where the right-hand side  $C$  does not depend on  $s'$ . Multiplying both sides by  $\bar{\mu}_t(s')$  and summing over  $s'$ , we obtain

$$\sum_{s,s' \in \mathcal{S}} \bar{\mu}_t(s') \mu^\pi(s) \log \pi(s | s') \leq \sum_{s' \in \mathcal{S}} \bar{\mu}_t(s') C = C = \sum_{s \in \mathcal{S}} \mu^\pi(s) \log \mu^\pi(s). \quad (32)$$

This bound holds for any choice of  $\bar{\mu}_t \in \Delta^d$ .

Next, we lower bound the right-hand side. For each  $s' \in \mathcal{S}$ , consider the reverse KL divergence

$$\text{KL}(\pi(\cdot | s') \| \mu) = \sum_{s \in \mathcal{S}} \pi(s | s') \log \frac{\pi(s | s')}{\mu^\pi(s)} \geq 0.$$

Hence

$$\sum_{s \in \mathcal{S}} \pi(s | s') \log \pi(s | s') \geq \sum_{s \in \mathcal{S}} \pi(s | s') \log \mu^\pi(s).$$

Multiplying by  $\mu^\pi(s')$  and summing over  $s'$  yields

$$\begin{aligned} \sum_{s, s' \in \mathcal{S}} \mu^\pi(s') \pi(s | s') \log \pi(s | s') &\geq \sum_{s, s' \in \mathcal{S}} \mu^\pi(s') \pi(s | s') \log \mu^\pi(s) \\ &= \sum_{s \in \mathcal{S}} \left( \sum_{s' \in \mathcal{S}} \mu^\pi(s') \pi(s | s') \right) \log \mu^\pi(s) \\ &= \sum_{s \in \mathcal{S}} \mu^\pi(s) \log \mu^\pi(s) = C, \end{aligned} \quad (33)$$

where we used the stationarity of  $\mu$ , namely  $\mu^\pi(s) = \sum_{s'} \mu^\pi(s') \pi(s | s')$ .

Combining (32) and (33), we conclude that, for any  $\bar{\mu}_t \in \Delta^d$ ,

$$\sum_{s, s' \in \mathcal{S}} \bar{\mu}_t(s') \mu^\pi(s) \log \pi(s | s') \leq C \leq \sum_{s, s' \in \mathcal{S}} \mu^\pi(s') \pi(s | s') \log \pi(s | s'),$$

which establishes

$$\sum_{s, s'} \bar{\mu}_t(s') \mu^\pi(s) \log \pi(s | s') \leq \sum_{s, s'} \mu^\pi(s') \pi(s | s') \log \pi(s | s'), \quad \forall \bar{\mu}_t \in \Delta^d. \quad \square$$

## D EXPERIMENT DETAILS

All experiments follow the same training setup unless otherwise specified: sequences are generated from a Markov chain with transition kernel  $\pi(\cdot | s) \sim \text{Dirichlet}(\alpha \cdot \mathbf{1}_d)$  with  $\alpha = 0.1$ . We use a batch size of 1024 for training and evaluate on 4096 test samples. Parameters are optimized with Adam (Kingma & Ba, 2017), using learning rate 0.05 for discrete Markov chains and 0.001 for dynamical systems. For gradient-based analysis, we adopt SGD with learning rate 1. Fresh data are sampled at each iteration to avoid memorization, and all implementations are based on JAX.

## E ADDITIONAL EXPERIMENT RESULTS ON MARKOV CHAIN

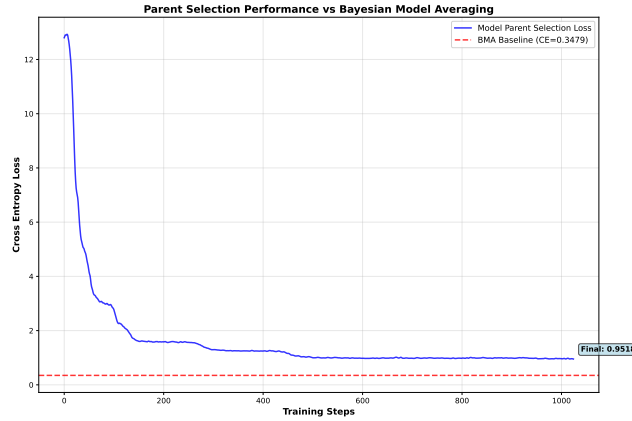


Figure 7: Parent selection  $\mathcal{L}_{pa}$  comparison between transformers and BMA during training. The metric is introduced as Eq. (10). Training configuration as the experiment in Fig. 2.

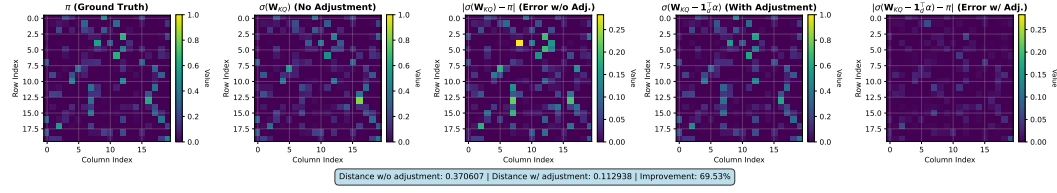


Figure 8: Parameter-Level Comparison between transformer and BMA ( $\mathbf{W} = \log \pi$ ). Trainable  $\mathbf{W}_{KQ}$  and  $\mathbf{W}_{OV}$ . Trained with  $d = 20$ ,  $H = 50$ ,  $L = 3$ , and 1024 training steps.

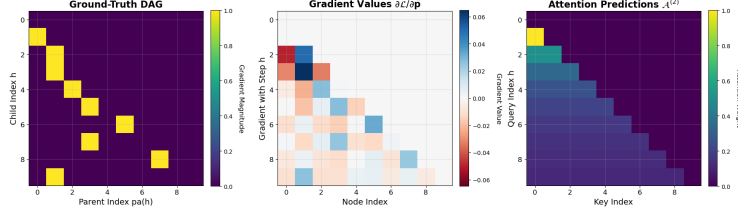


Figure 9: **Gradient Validation of  $\frac{\partial \ell}{\partial \mathbf{p}}$** . From left to right: ground-truth graph  $\mathcal{G}$ , the gradients of  $\frac{\partial \ell}{\partial \mathbf{p}} \in \mathbb{R}^H$  stacked as row vectors, and attention weights  $\mathcal{A}_h^{(2)}$  uniformly distributed since  $\mathbf{W} = 0$ .

## F EXPERIMENT RESULTS ON DYNAMICAL SYSTEM

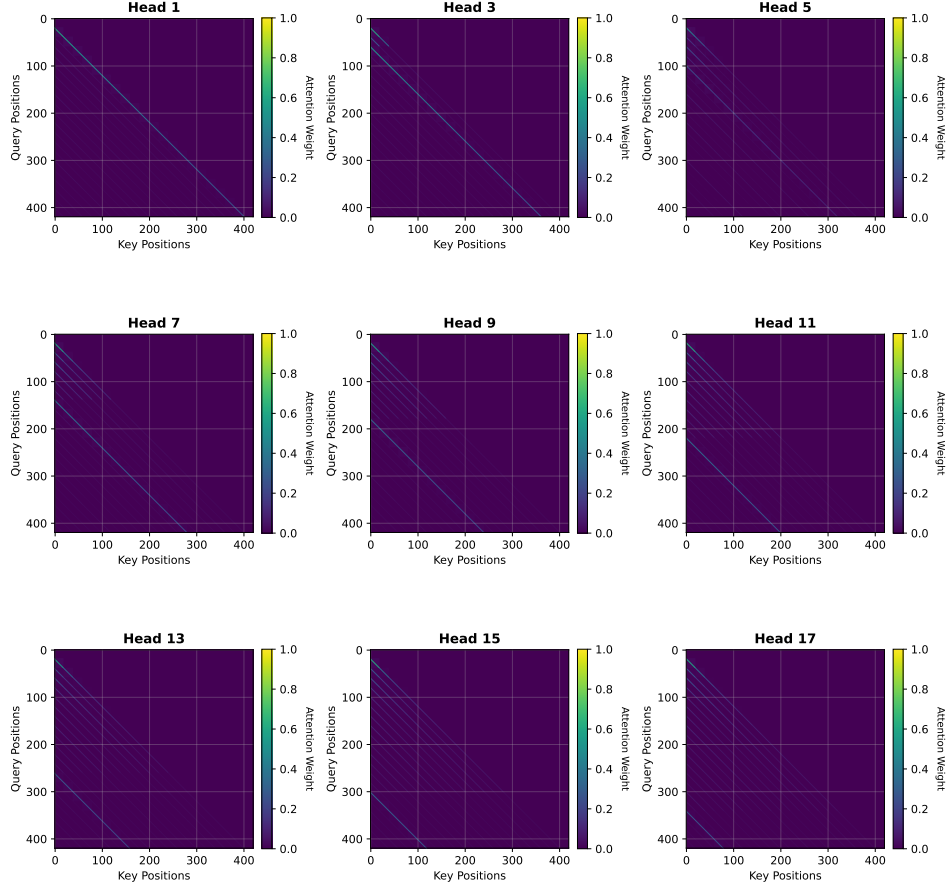


Figure 10: Visualization of 1st-layer Attention  $\mathcal{A}^{(1)} \in \mathbb{R}^{T \times T}$ . For readability, we visualize only nine of the twenty heads, to better highlight the attention patterns on this long sequence of length 400. The first layer replicates the historical occurrence of the same token. Model trained with  $L = 20$  examples, trajectory length  $H = 20$ , vocabulary size  $d = 10$ , 20 heads in the first layer, and 2048 training steps. The RPE parameters are initialized with a small positive value (0.5) along the construction direction, and grow to much larger magnitudes after training.

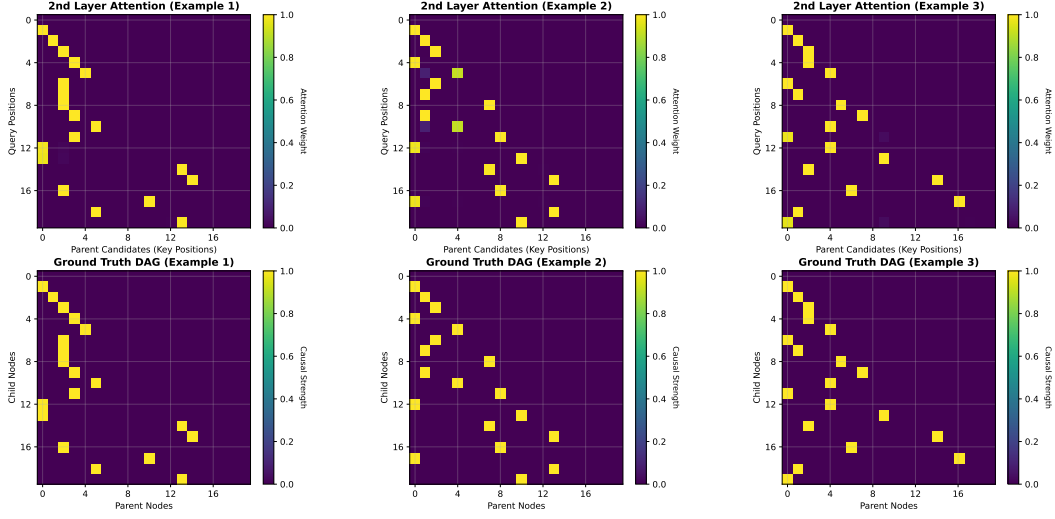


Figure 11: 2nd-Layer Attention  $\mathcal{A}^{(2)} \in \mathbb{R}^{H \times H}$  Visualization. Attention patterns matches the groundtruth causal structure in dynamical system setting.

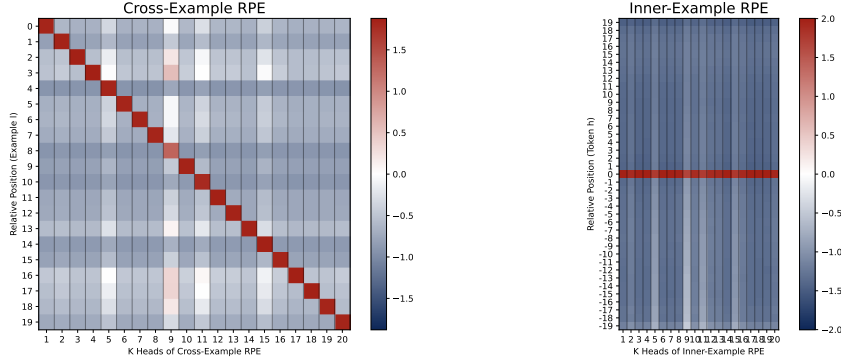


Figure 12: Visualization of first RPE layer. The parameters are consistent with the construction.

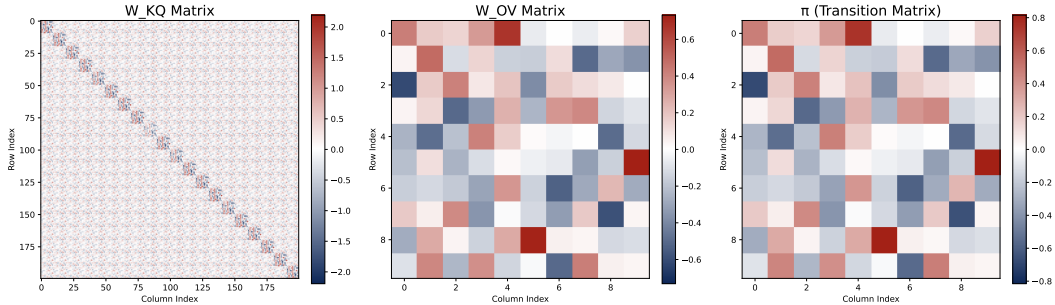


Figure 13: Visualization of 2nd-layer.  $\mathbf{W}_{KQ} \in \mathbb{R}^{dL \times dL}$  shows noticable non-zero blocks on its diagonal. The occuring block is of size  $d \times d$ .



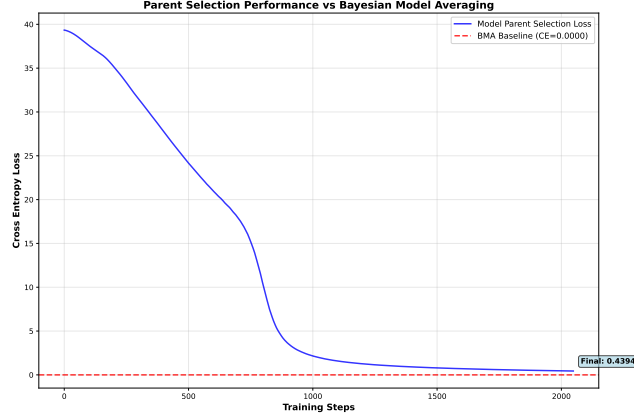
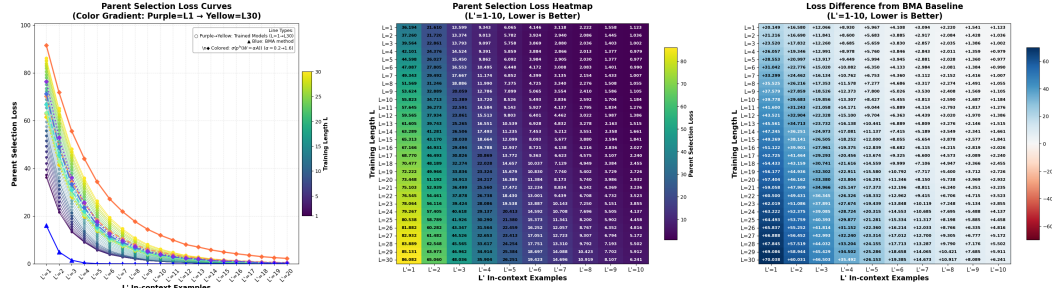


Figure 14: Parent selection loss during training in dynamical system setting.

Figure 15: Generalization of parent loss  $\{\mathcal{L}_{pa}^{L'}\}$  for transformers trained with  $L \in \{1, \dots, 30\}$  in dynamical system setting. Trained with  $d = 10$ ,  $H = 15$  and 2048 training steps.

## G DISENTANGLED TRANSFORMER WITH ABSOLUTE POSITION EMBEDDING

The disentangled transformer with absolute position embedding (APE) is formulated by:

$$\begin{aligned}
 \textbf{Embedding Layer:} \quad & h_t^{(0)} = [E(w_t), \text{Pos}(w_t)] = [x_t, e_t], \mathbf{H}^{(0)} = [h_1^{(0)}, \dots, h_T^{(0)}] \in \mathbb{R}^{d_0} \\
 \textbf{1st Attention (K-head):} \quad & \text{Attn}_t^k(\mathbf{H}^{(0)}; \theta) = \sigma \left( h_{1:t-1}^{(0)\top} \mathbf{W}_{KQ}^{(1),k} h_t^{(0)} \right)^\top h_{1:t-1}^{(0)\top} \mathbf{W}_{OV}^{(1),k} \in \mathbb{R}^d, \\
 \textbf{Disentangled Residual:} \quad & h_t^{(1)} = [h_t^{(0)}, \text{Attn}_t^1(\mathbf{H}^{(0)}; \theta), \dots, \text{Attn}_t^K(\mathbf{H}^{(0)}; \theta)] \in \mathbb{R}^{d_0+Kd}, \\
 \textbf{2nd Attention (1-head):} \quad & f_{\text{tf}}(\cdot | \mathcal{H}_t) = \sigma \left( h_{1:t-1}^{(1)\top} \mathbf{W}_{KQ}^{(2)} h_t^{(1)} \right)^\top h_{1:t-1}^{(1)\top} \mathbf{W}_{OV}^{(2)} \in \mathbb{R}^d.
 \end{aligned} \tag{34}$$

First, we can see the model parameter  $\mathbf{W}_{KQ}^{(1),k} \in \mathbb{R}^{d_0 \times d_0}$  where  $d_0 = d + T$  and  $T$  is the sequence length. The total number of parameters in the first layer is  $O(d^2 + H^2 L^2)$  compared to  $O(H + L)$  parameters of the model with RPE in Eq. (5). The redundancy of parameters may lead to difficulties of interpreting the mechanism of transformers. Besides, since for disentangled transformer with APE, the embedding dimension is proportional to the length of input sequence, this may make it difficult for us to interpret transformers' mechanism on longer sequence tasks.

As for this transformer, we first provide a theoretical construction which is consistent with our construction for RPE model in Theorem 1. Empirically, we show this transformer can successfully select causal tokens. Besides, we provide results of trainable transformers showing alignments with our construction in attention visualization and parameter verification.

### G.1 THEORETICAL CONSTRUCTION

In this section, we provide a construction demonstrating how the proposed two-layer architecture possesses the capacity to implement the specific causal selection mechanism derived in our analysis. Let the input embedding dimension be  $d_0 = d + T$ , where  $d$  is the token dimension,  $T$  is the sequence length (due to absolute position embedding) and an input sequence contains  $L + 1$  examples of length- $L$  chain  $T = H(L + 1)$ . Suppose  $\mathcal{N}_{L+1}$  denotes the set of nodes from the last example, i.e., we have  $\mathcal{N}_{L+1} = \{t \in T \mid \exists h \in [H], t = HL + h\}$ .

#### G.1.1 LAYER 1: MULTI-HEAD ATTENTION CONSTRUCTION

The first layer consists of  $K$  attention heads ( $K \leq L$ ). The Query-Key matrix  $\mathbf{W}_{KQ}^{(1),k}$  will attend to specific predecessor tokens based on position. We construct it as a block matrix where the active interaction terms are confined to the position-embedding subspace:

$$\mathbf{W}_{KQ}^{(1),k} = \begin{bmatrix} 0_{d \times d} & 0_{d \times T} \\ 0_{T \times d} & \tilde{\mathbf{W}}_{KQ}^{(1),k} \end{bmatrix}, \tilde{\mathbf{W}}_{KQ}^{(1),k} = \beta \begin{bmatrix} 0_{H \times H} & \vdots & 0_{H \times H} \\ 0_{T \times HL} & I_{H \times H} & \vdots & 0_{H \times H} \end{bmatrix} \quad (\text{\textcolor{brown}{k-th block} active}). \tag{35}$$

From this construction, if  $\beta \rightarrow \infty$ , the attention score of the first attention layer is given by:

$$\mathcal{A}_{ij}^{(1),k} = \begin{cases} \frac{1}{i} \mathbf{1}_{[j < i]}, & \text{if } i \notin \mathcal{N}_{L+h}, \\ \mathbf{1}_{[j=kH+h]}, & \text{if } i \in \mathcal{N}_{L+h}, i = LH + h. \end{cases} \tag{36}$$

For the value projection,  $\mathbf{W}_{OV}^{(1),k}$  will propagate the semantic content of the attended tokens:

$$\mathbf{W}_{OV}^{(1),k} = \begin{bmatrix} I_{d \times d} \\ 0_{T \times d} \end{bmatrix} \in \mathbb{R}^{(d+T) \times d}. \tag{37}$$

And the output of the first attention will be:

$$\text{Attn}_i^k(\mathbf{H}^{(0)}; \theta) = \mathcal{A}_{i \rightarrow \cdot}^{(1),k} h_{1:T}^{(0)\top} \mathbf{W}_{OV}^{(1),k} = \begin{cases} \bar{\mu}(x_{1:i-1}), & \text{if } i \notin \mathcal{N}_{L+h}, \\ x_h^k, & \text{if } i \in \mathcal{N}_{L+h}, i = LH + h. \end{cases} \tag{38}$$

### G.1.2 DISENTANGLED RESIDUAL STREAM

Unlike standard summation residuals, this disentangled transformer employ a concatenation strategy. [Nichani et al. \(2024\)](#) proved this transformer is actually equivalent to a decoder based attention-only transformer (Theorem 3). The output of the first layer is the concatenation of the original input and the outputs of all  $K$  heads:

$$\mathbf{h}_t^{(1)} = \left[ \mathbf{h}_t^{(0)}; \text{Attn}_t^1, \dots, \text{Attn}_t^K \right] \in \mathbb{R}^{d_0+Kd}. \quad (39)$$

The dimension of the second layer input is  $d_1 = d_0 + Kd = d + T + Kd$ .

### G.1.3 LAYER 2: SINGLE-HEAD ATTENTION CONSTRUCTION

The second layer employs a single attention head to aggregate the evidence collected by the  $K$  heads in the previous layer.

$$\mathbf{W}_{KQ}^{(2)} = \left[ \begin{array}{c|c|c} 0_{d \times d} & 0_{d \times T} & 0_{d \times Kd} \\ \hline 0_{T \times d} & 0_{T \times T} & 0_{T \times Kd} \\ \hline 0_{Kd \times d} & 0_{Kd \times T} & \tilde{\mathbf{W}}_{KQ}^{(2)} \end{array} \right], \tilde{\mathbf{W}}_{KQ}^{(2)} = \left[ \begin{array}{ccc} \log \pi & 0_{d \times d} & \cdots & 0_{d \times d} \\ 0_{d \times d} & \log \pi & \cdots & 0_{d \times d} \\ \vdots & \vdots & \ddots & \vdots \\ 0_{d \times d} & 0_{d \times d} & \cdots & \log \pi \end{array} \right] \quad (40)$$

Finally, the output projection  $\mathbf{W}_{OV}^{(2)}$  projects the aggregated context back to the semantic space by:

$$\mathbf{W}_{OV}^{(2)} = \left[ \begin{array}{c} \log \pi \\ \hline 0_{T \times d} \\ \hline 0_{Kd \times d} \end{array} \right] \in \mathbb{R}^{d_1 \times d}. \quad (41)$$

From the Eq. (38) and (39), we can see that  $\mathbf{h}_t^{(1)} = [\mathbf{h}_t^{(0)}; \mathbf{x}_h^1, \dots, \mathbf{x}_h^K]$  if  $t \in \mathcal{N}_{L+1}, t = LH + h$  else  $\mathbf{h}_t^{(1)} = [\mathbf{h}_t^{(0)}; \bar{\mu}_t, \dots, \bar{\mu}_t]$ . Then, the attention weight of the second layer for any  $i \in \mathcal{N}_{L+1}, i = LH + h$  of our interests, is given by:

$$\tilde{\mathcal{A}}_{ij}^{(2)} = \begin{cases} \sum_{k=1}^K \log \pi(x_h^k | x_{h'}^k), & \text{for } j \in \mathcal{N}_{L+1}, j = LH + h', \\ \sum_{k=1}^K \bar{\mu}_j^\top \log \pi(x_h^k), & \text{for } j \notin \mathcal{N}_{L+1}. \end{cases} \quad (42)$$

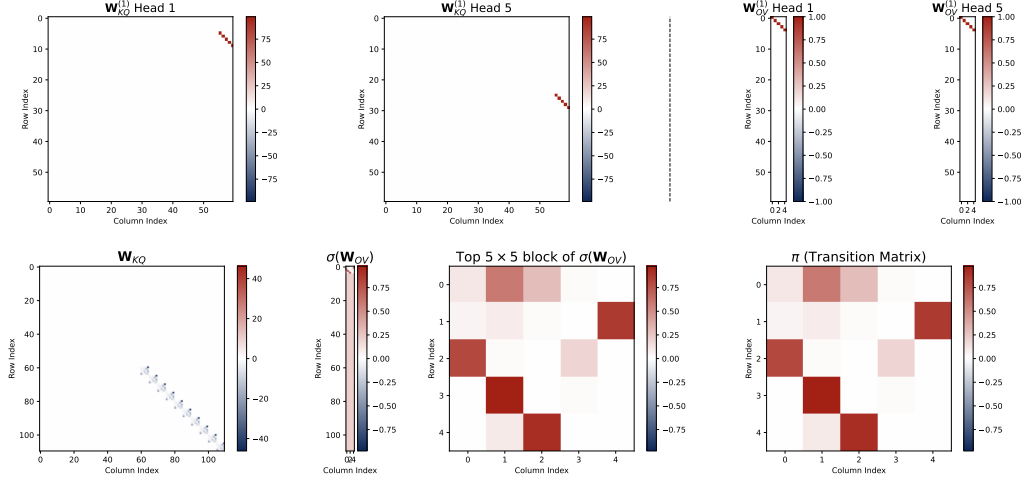
So for  $i, j \in \mathcal{N}_{L+1}$ , we have  $\tilde{\mathcal{A}}_{ij}^{(2)} = \sum_{k \in [K]} \log \pi(x_h^k | x_{h'}^k)$  aligned with Theorem 1.

Furthermore, suppose  $K = L$ , i.e., we use  $L$  examples to infer the causal structure, and the Markov chain is stationary  $\mathbf{x}_h \sim \mu^\pi$ . As  $L \rightarrow \infty$ , for any  $j \notin \mathcal{N}_{L+1}$ , we have  $\tilde{\mathcal{A}}_{ij}^{(2)} \rightarrow \sum_{s, s'} \bar{\mu}_j(s') \mu^\pi(s) \log \pi(s | s') \leq \sum_s \mu^\pi(s) \log \mu^\pi(s)$ . While for the true parent token  $t = HL + pa(h)$ , we have  $\tilde{\mathcal{A}}_{it} \rightarrow \sum_{s, s'} \mu^\pi(s') \pi(s | s') \log \pi(s | s')$  which is larger than  $\sum_s \mu^\pi(s) \log \mu^\pi(s)$ . The above quantity relation is drawn by non-negativity of KL divergence, whose detailed proof is provided in Appendix C.8. Then, the attention weights of the second layer can select the causal parent token  $\mathcal{A}_{i \rightarrow pa(i)} \rightarrow 1$ . And  $\mathbf{W}_{OV}^{(2)}$  predicts the transition.

**Empirical Verification.** We show the parameter visualization of the construction in Fig. 16a and its empirical attention visualization of parent selection in Fig. 16b. In Fig. 17, the constructed model shows precise parent selection accuracy (cross-entropy loss 0.0706) which is very close to the target algorithm BMA's (ce loss 0.0473).

## G.2 EXPERIMENTS OF TRAINABLE TRANSFORMERS

In the following, we train the standard disentangled transformer formulated by Eq.(34). To show the alignment with theoretical interpretation, we use three strategies to initialize the network: (a)



(a) The parameter visualization of theoretical construction.

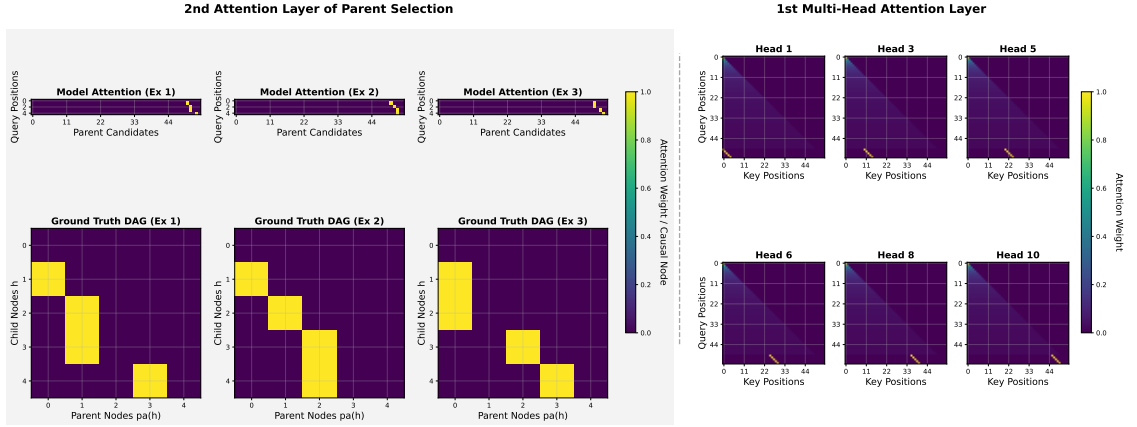
(b) Attention pattern visualization of theoretical construction. The first layer (right) copies the previous  $L$  examples to the hidden space of the last one  $L + 1$ . For the second layer (left), the attention weights attend to the correct causal parents which are located in the last 5 columns. The queries don't attend to the keys from first  $L$  examples.

Figure 16: Parameter visualization and attention pattern visualization of theoretical construction.

fully random initialization: all the parameters are initialized randomly with Gaussian distribution; (b) block-amplified random initialization: parameters are initialized randomly (of scale 0.1), while the targeted block of the attention projection matrix is assigned a larger magnitude (of scale 0.5) to introduce an inductive bias; (c) direction-consistent initialization: parameters are initialized such that the dominant blocks point in the analytically derived construction direction, still allowing model learning to refine the magnitudes (initial magnitudes:  $0.2 \times$  optimal parameters).

We first compare the parent token prediction performance of these models during the training process in Fig. 17. The results show that the 2-layer transformmr is fully capable of selecting causal parents in its 2nd-layer attention head.

Then we visualize the attention pattern of the trained model in Fig. 19. For the first attention layer, the figure shows query from the last example  $L + 1$  mostly attend to one example among  $L$  context examples, while some heads demonstrate the degeneration with uniform attention to previous tokens. For the second attention layer, the transformers with different initializations all show their noticable capability of predicting causal parents. Further, we visualize all the parameters of the transformer in Fig. 18. We can see some alignments between the construction in Fig. 16a and the trained parameters. Since the transformer with absolute position embedding has far more parameters of  $(\{\mathbf{W}_{KQ}^{(1),k}, \mathbf{W}_{OV}^{(1),k}\}_k, \mathbf{W}_{KQ}^{(2)}, \mathbf{W}_{OV}^{(2)})$  than the one with RPE, the full interpretation of its first layer is difficult. For the second layer, the parameter  $\mathbf{W}_{KQ}^{(2)}$  also shows the diagonal pattern consistent with construction and  $\mathbf{W}_{OV}^{(2)}$  shows the  $\log \pi$  pattern.

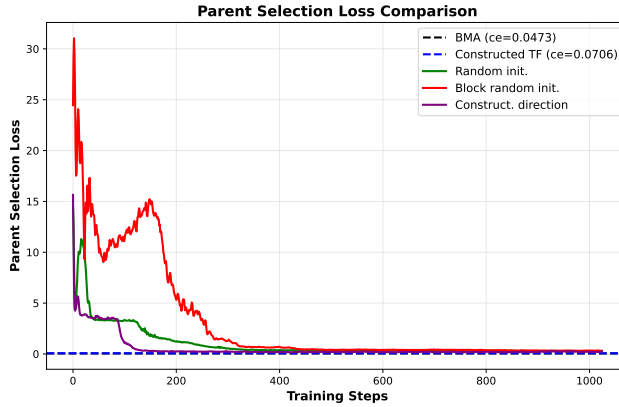
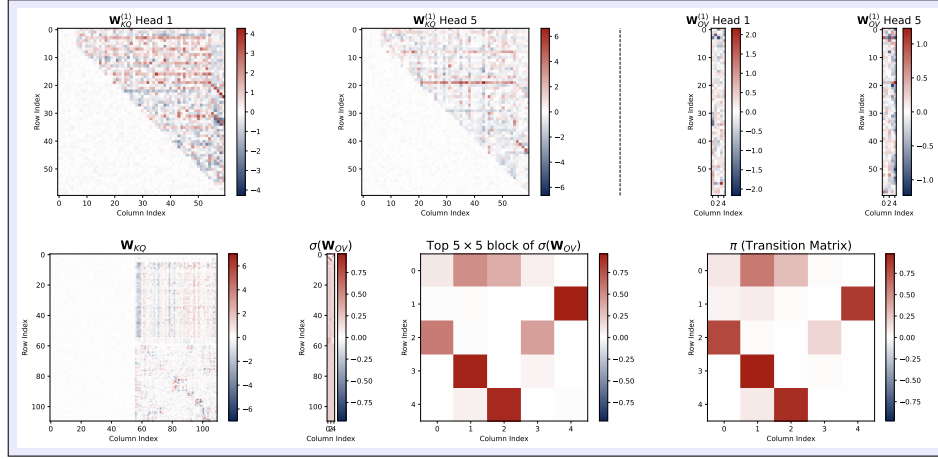
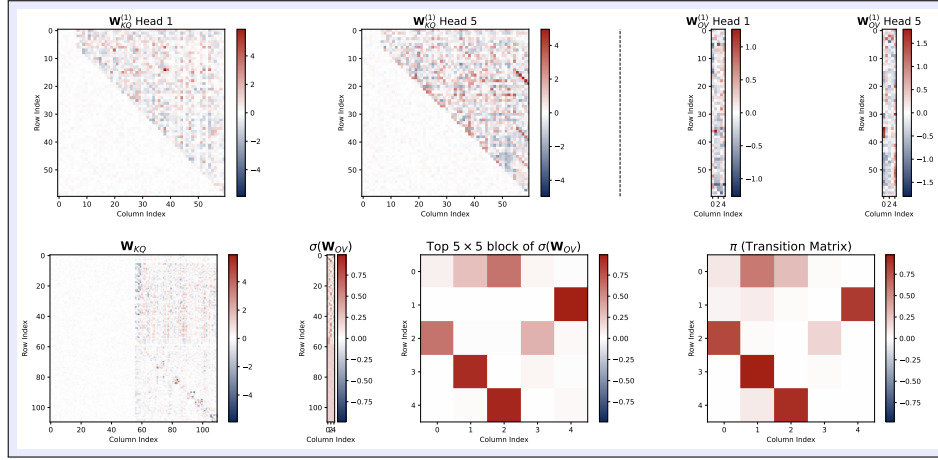


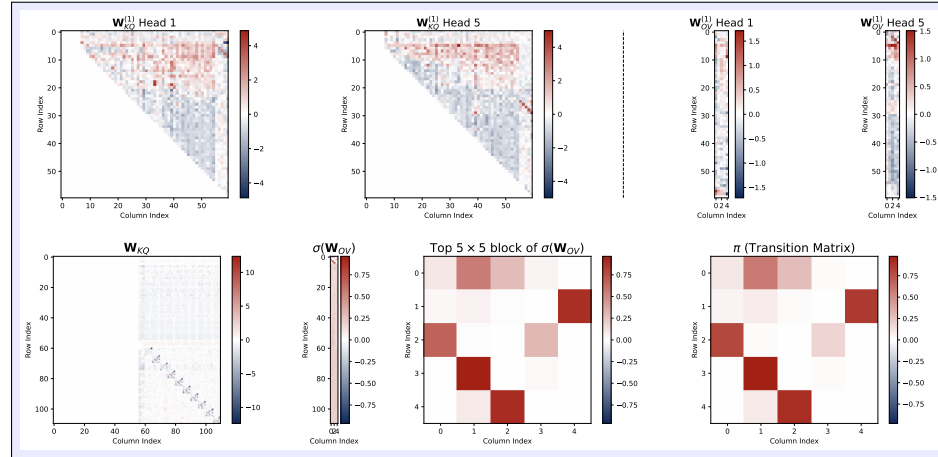
Figure 17: Parent selection loss  $\mathcal{L}_{pa}$  of the transformer with absolute position embedding and different initialization strategies.



(a) With Fully Random Initialization. Head 1 and 5 of the first layer  $W_{KQ}^{(1)}$  exhibits an identity submatrix ( $5 \times 5$ ) at the last column.



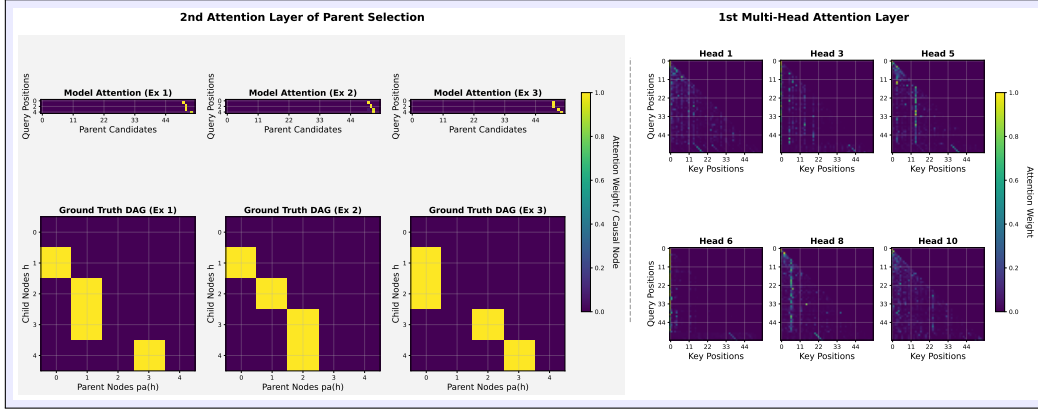
(b) With Block-Amplified Random Initialization. Head 1 of  $W_{KQ}^{(1)}$  degenerate which can be verified in attention visualization Fig. 19b (Head 1). Head 5 shows multiple identity submatrices which possibly suggests superposition.



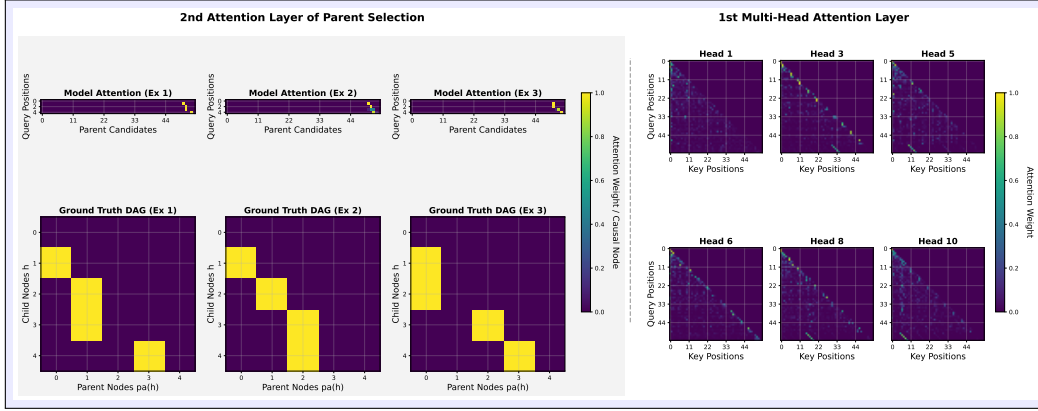
(c) With Direction-Consistent Initialization. Head 1 and 5 of the first layer  $W_{KQ}^{(1)}$  exhibits an identity submatrix ( $5 \times 5$ ) at the last column which is aligned with the theoretical construction.

Figure 18: Parameter visualization of trained transformer with absolute position embedding. The second layer shows strong alignment in diagonal patterns of  $W_{KQ}$  and  $\log \pi$  pattern of  $W_{OV}$ .

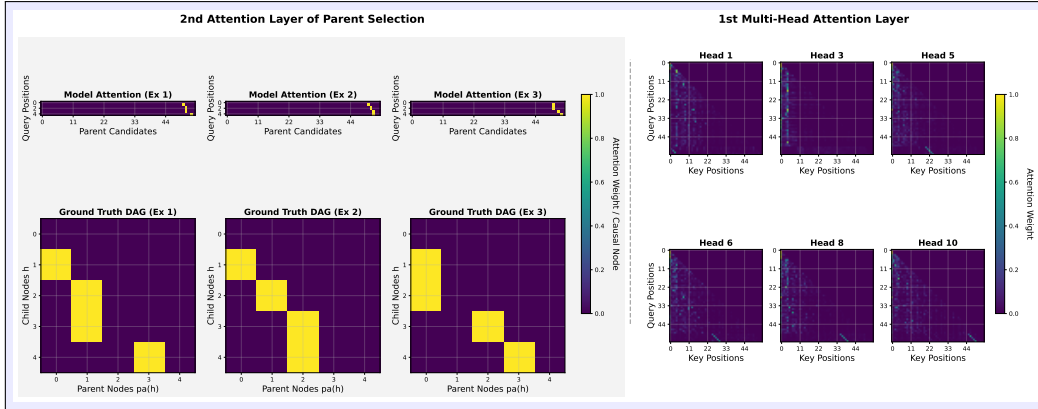




(a) With Fully Random Initialization. In the first layer, Head 1, 3, 5, 6 and 8 of KQ matrices copies tokens from previous examples (to the query token of the last example), while Head 10 degenerates showing uniform attention (uniform features are seen as constants eliminated by 2nd softmax attention layer).



(b) With Block-Amplified Random Initialization. In the first layer, Head 3, 5, 8 and 10 of KQ matrices copies tokens from previous examples, while Head 1 degenerates showing uniform attention.



(c) With Direction-Consistent Initialization. In the first layer, Head 1, 5, 6, 8 and 10 of KQ matrices copies tokens from previous examples, while Head 3 degenerates showing uniform attention.

Figure 19: Attention pattern visualization of trained transformer with absolute position embedding.

## H DISENTANGLED TRANSFORMER WITH APE VARIANT

### H.1 MODEL ARCHITECTURE

For the variant of two types of APE, we consider the disentangled transformer simplified by eliminating some components added to residual stream. The transformer structure we consider below can be seen as substituting the position embedding of structure Eq. (34) and simplify the model by assuming zero blocks in model weights.

$$\begin{aligned}
 \text{Embedding Layer:} \quad & \tilde{h}_t^{(0)} = [\text{Pos}_L(w_t), \text{Pos}_H(w_t)] \in \mathbb{R}^{d_0} \\
 \text{1st Attention (K-head):} \quad & \text{Attn}_t^k = \sigma(\tilde{h}_{1:t-1}^{(0)\top} \mathbf{W}_{KQ}^{(1,k)} \tilde{h}_t^{(0)})^\top x_{1:t-1}^\top \mathbf{W}_{OV}^{(1,k)} \in \mathbb{R}^d, \\
 \text{Disentangled Residual:} \quad & \tilde{h}_t^{(1)} = [\text{Attn}_t^1, \dots, \text{Attn}_t^K] \in \mathbb{R}^{Kd}, \\
 \text{2nd Attention (1-head):} \quad & f_{\text{tf}}(\cdot | \mathcal{H}_t) = \sigma\left(\tilde{h}_{1:t-1}^{(1)\top} \mathbf{W}_{KQ}^{(2)} \tilde{h}_t^{(1)}\right)^\top x_{1:t-1}^\top \mathbf{W}_{OV}^{(2)} \in \mathbb{R}^d,
 \end{aligned} \tag{43}$$

where we simply assume  $\mathbf{W}_{OV}^{(1,k)} = I_d$ . Since the difference lies in the positional embedding, the construction in Appendix G remains valid which can exhibit capabilities in causal token selection empirically. Besides, we train this transformer under the same Markov chain setup as in the transformer with RPE experiments, and obtain consistent results as shown below.

### H.2 EXPERIMENT RESULTS

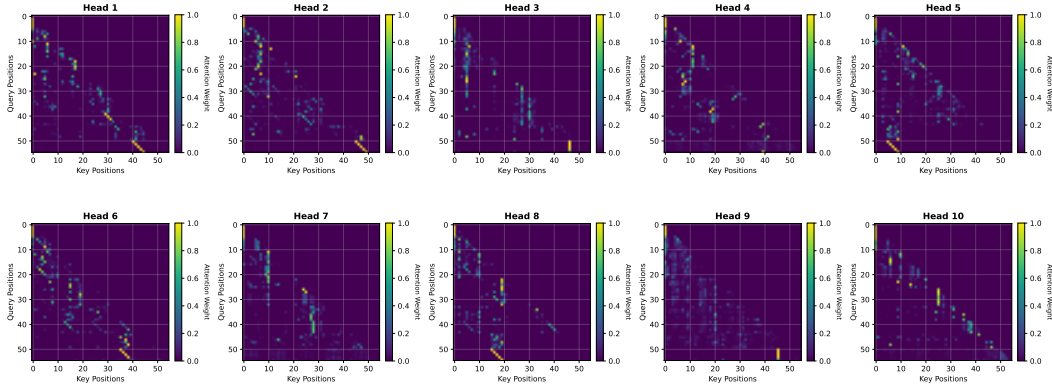


Figure 20: 1st-Layer Attention Visualization of transformers in Eq. (43). Heads 1, 2, 5, 6 and 8 exhibits the diagonal block pattern at the last rows performing the copying mechanism, while Heads 4, 7 and 10 degenerate to uniform attention. Heads 3 and 9 gives uniform outputs not influencing the 2nd attention layer (eliminated by softmax attention). Trained with  $H = 5, L = 10, d = 5$  and 10000 training steps.

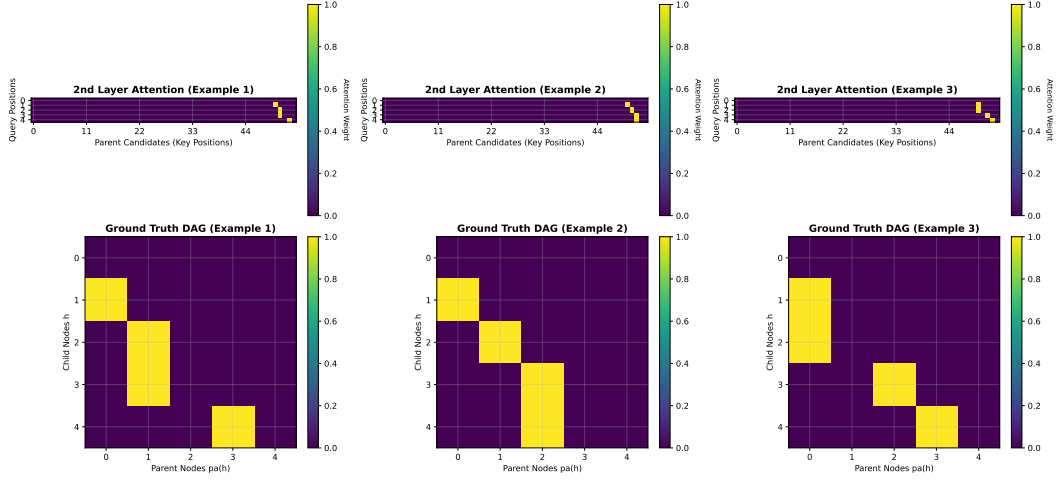


Figure 21: Visualization of 2nd-attention layer. Queries are from the last example  $x_{1:H}^{L+1}$ . Keys are  $x_{1:T} = x_{1:H}^L$  the whole sequence. Attention layer of disentangled transformer can recognize the causal structure in-context.

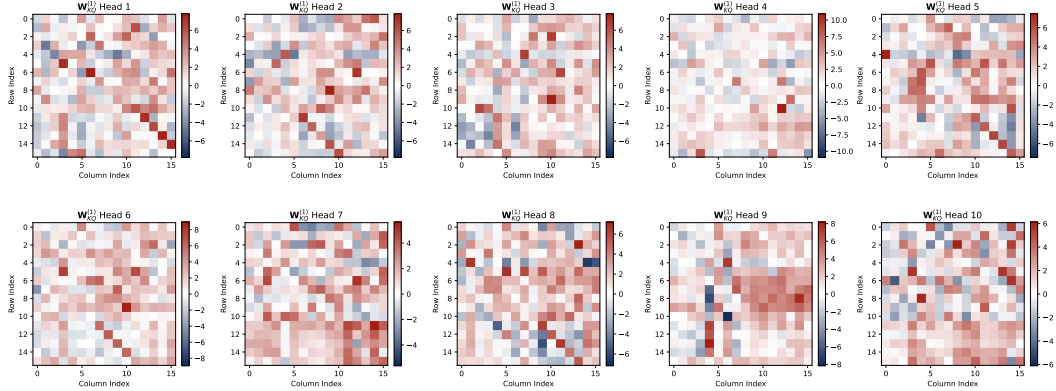


Figure 22: Parameter visualization of the first attention layer  $W_{KQ}^{(1),k}$  (10 heads in total). Full interpretation is still challenging for huge parameter space. The attention-level behavior understanding can be referred to Fig. 20.

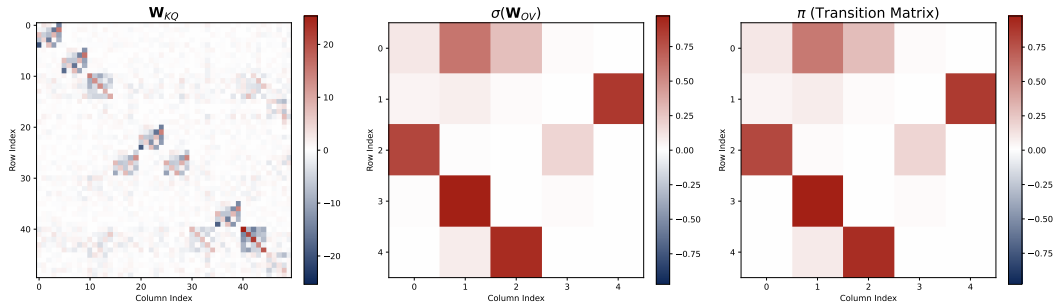


Figure 23: Parameter visualization of the second attention layer  $W_{KQ}^{(2)}$ ,  $W_{OV}^{(2)}$ . In the variant of two types of Absolute position embedding, the second layer also shows strong alignment in diagonal patterns of  $W_{KQ}$  and  $\log \pi$  pattern of  $W_{OV}$ .

## I STANDARD TRANSFORMER WITH FEEDFORWARD NEURAL NETWORK

In this section, we consider a standard 2-layer transformer with FFN layer as follows.

$$\begin{aligned}
\text{Learnable Embedding:} \quad & \mathbf{h}_t^{(0)} = \text{Emb}_{\mathbf{V}}(\mathbf{w}_t) + \text{Emb}_{\mathbf{P}}(\mathbf{w}_t), & \in \mathbb{R}^{d'} \\
\text{MHA Layer \& Residual:} \quad & \tilde{\mathbf{h}}_t^{(l)} = \mathbf{h}_t^{(l)} + \text{MHA}_t(\mathbf{H}^{(l)}; \mathbf{W}_{KQ}, \mathbf{W}_{OV}), & \in \mathbb{R}^{d'}, \\
\text{FFN Layer \& Residual:} \quad & \mathbf{h}_t^{(l+1)} = \tilde{\mathbf{h}}_t^{(l)} + \text{FFN}_t(\tilde{\mathbf{H}}^{(l)}; \mathbf{W}, \mathbf{b}) & \in \mathbb{R}^{d'}, \\
\text{Unembedding Layer:} \quad & \mathbf{f}_{\text{tf}}(\cdot | \mathcal{H}_t) = \mathbf{W}_U \mathbf{h}_t^{(L)} & \in \mathbb{R}^d,
\end{aligned} \tag{44}$$

where  $\mathbf{H}^{(l)} = [\mathbf{h}_1^{(l)}, \dots, \mathbf{h}_T^{(l)}]$ , the multi-head attention (MHA) is formulated by

$$\text{MHA}_t(\mathbf{H}^{(l)}; \theta) = \sum_k \sigma \left( \mathbf{h}_{1:t-1}^{(l)\top} \mathbf{W}_{KQ}^{(l),k} \mathbf{h}_t^{(l)} \right)^\top \mathbf{h}_{1:t-1}^{(l)\top} \mathbf{W}_{OV}^{(l),k}, \tag{45}$$

and the FFN layer

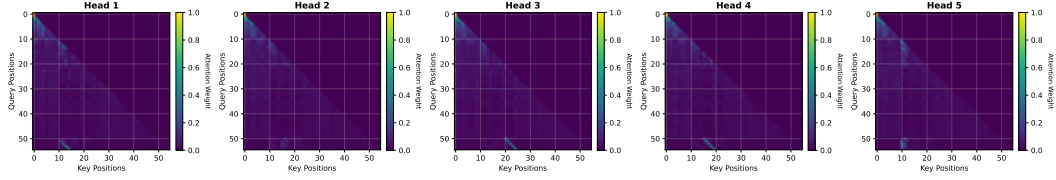
$$\text{FFN}_t(\tilde{\mathbf{H}}^{(l)}; \theta) = \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \tilde{\mathbf{h}}_t^{(l)} + \mathbf{b}_1) + \mathbf{b}_2. \tag{46}$$

We consider the two-layer transformer  $L = 2$  with  $K$  heads in the first layer and one head in the second.<sup>4</sup> For the task, the input sequence consists of  $M = 10$  in-context examples of Length- $H$  Markov chains with  $d = 5$  states and the total length  $T = H(M + 1)$ . We set the hidden dimension as  $d' = 128$ . For initialization, the parameters  $\mathbf{W}$  of the transformer is initialized randomly by Gaussian initialization:  $\mathbf{W}_{ij} \sim \mathcal{N}(0, 1/d_{\mathbf{W}})$  where  $d_{\mathbf{W}}$  is decided by the dimension of  $\mathbf{W}$ . We optimize the model using AdamW with a learning rate of  $1 \times 10^{-3}$  and a weight decay of  $1 \times 10^{-4}$ . Fresh data are sampled at each iteration of training without repetition.

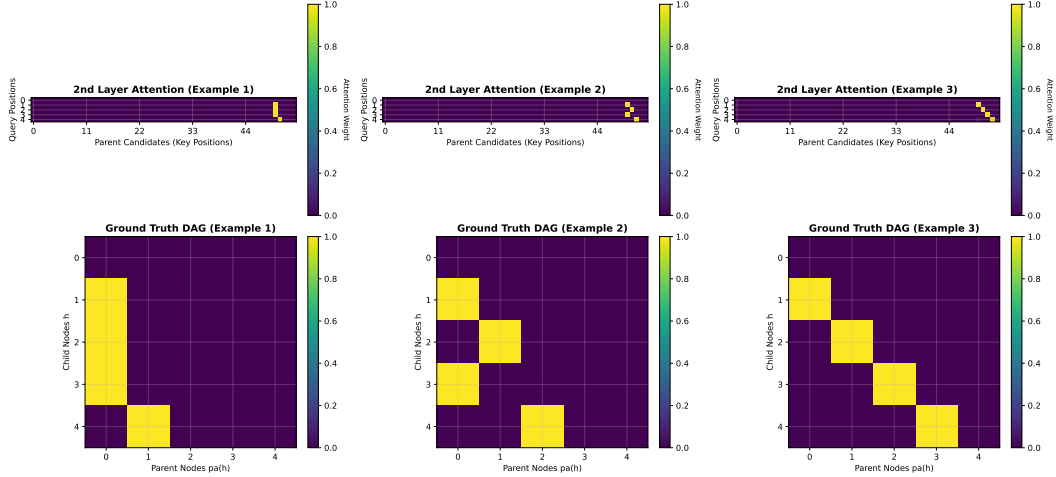
**Experiment Results.** We train two transformers with 5000 steps and  $K = 5$  or 10 heads in the first layer. We observe the attention weights of the first layer visualized in Fig. 24a and Fig. 25a implement the copying mechanism where the features of one context example are copied to the position of last example  $M + 1$ : the heads of the first layer show a diagonal submatrix occurring at the last several rows of example  $M + 1$ . Except for these, the remainings mainly show degenerated attention patterns at the rows of the last example  $M + 1$ . In the visualization of the second layer, we find that the trained standard transformer with MLPs can recognize the causal parents in its attention weights of the 2nd layer. The aligned attention pattern and graph groundtruth in Fig. 24b and Fig. 25b supports our construction of how transformers can handle with in-context causal learning.

**Quantitive Results.** We provide the results regarding how accurate transformers during training can select random parents in its second attention layer in Fig. 26. We use the cross-entropy loss as the evaluation metric for the accuracy and compare the trained transformers with BMA. We observe during training process, standard transformers gradually acquire the capability of in-context causal learning and approximate the loss of BMA.

<sup>4</sup>Our implementation is based on the codebase provided by Nichani et al. (2024).

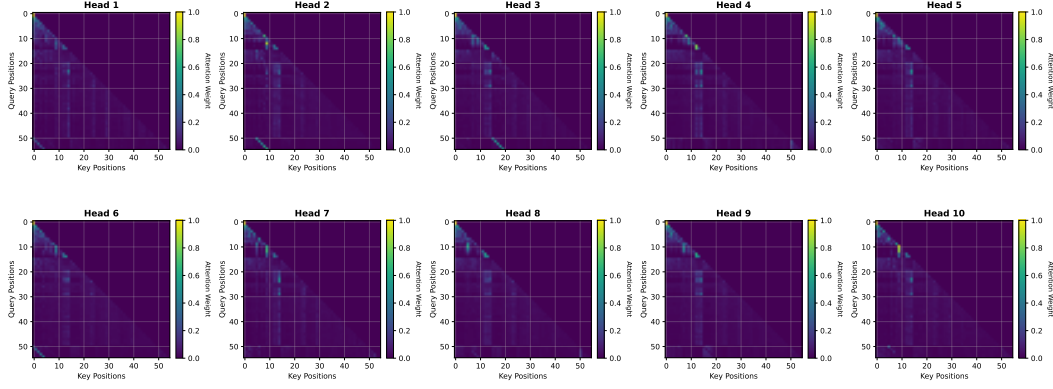


(a) Visualization of first multi-head attention layer. Heads 1, 3 and 4 show the diagonal block at the rows of the last example. Information from previous examples is copied to the hidden space of the last example.

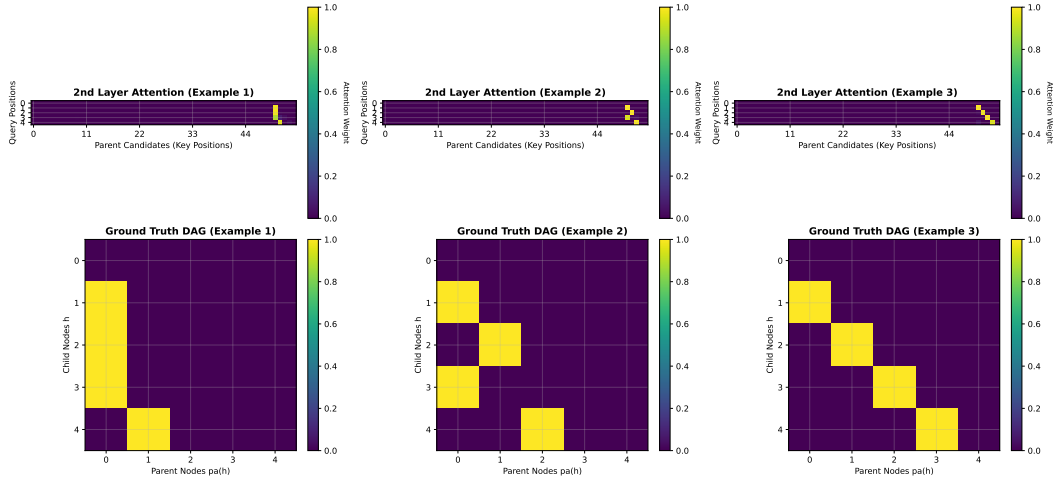


(b) Visualization of second attention layer. Queries are from the last example  $\mathbf{x}_{1:H}^{L+1}$ . Attention layer of standard transformer can recognize the causal structure in-context.

Figure 24: Attention visualization of standard transformer with MLP (5 heads).



(a) Visualization of first multi-head attention layer. Heads 1, 2, 3 and 6 show the diagonal block at the rows of the last example. Information from previous examples is copied to the hidden space of the last example.



(b) Visualization of second attention layer. Queries are from the last example  $x_{1:H}^{L+1}$ . Attention layer of standard transformer can recognize the causal structure in-context.

Figure 25: Attention visualization of standard transformer with MLP (10 heads).

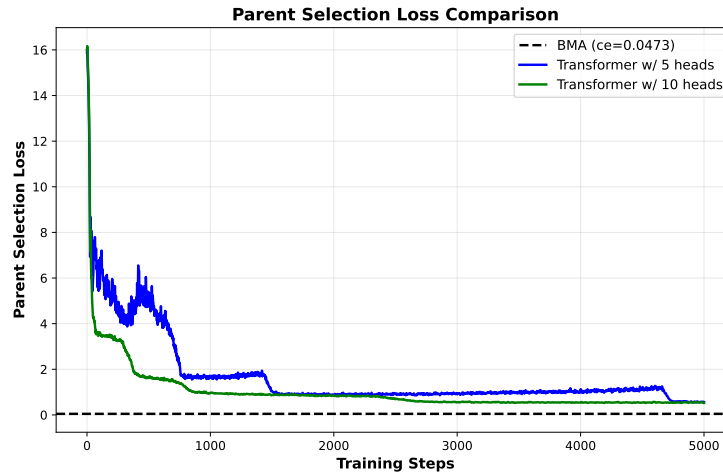


Figure 26: Parent selection loss  $\mathcal{L}_{pa}$  of the standard transformer with learnable position embedding and MLP (5 or 10 heads in the first layer). During training, standard transformers gradually acquire the capability of in-context causal learning and approximate the loss of BMA.