Understanding the Language Model to Solve the Symbolic Multi-Step Reasoning Problem from the Perspective of Buffer Mechanism

Zhiwei Wang^{1,2}, Yunji Wang^{1,2}, Zhongwang Zhang^{1,2}, Zhangchen Zhou^{1,2}, Hui Jin⁴, Tianyang Hu³, Jiacheng Sun⁴, Zhenguo Li ⁴, Yaoyu Zhang^{1,2}, Zhi-Qin John Xu^{1,2,5}

¹Institute of Natural Sciences, MOE-LSC, Shanghai Jiao Tong University ²School of Mathematical Sciences, Shanghai Jiao Tong University ³The Chinese University of Hong Kong, Shenzhen ⁴Huawei Noah's Ark Lab

⁵Key Laboratory of Marine Intelligent Equipment and System, Ministry of Education Correspondence: xuzhiqin@sjtu.edu.cn

Abstract

Large language models have consistently struggled with complex reasoning tasks, such as mathematical problem-solving. Investigating the internal reasoning mechanisms of these models can help us design better model architectures and training strategies, ultimately enhancing their reasoning capability. In this study, we constructed a symbolic multi-step reasoning task to investigate the information propagation mechanisms in Transformer models when solving the task through direct answering and Chain-of-Thought (CoT) reasoning. We introduced the concept of buffer mechanism: the model stores various information in distinct buffers and selectively extracts it through the query-key matrix. We proposed a random matrix-based algorithm to enhance the model's reasoning ability. This algorithm introduces only 132 trainable parameters, yet leads to significant performance improvements on 7 multi-step reasoning datasets, including PrOntoQA, LogicAsker, and LogicInference. These findings provide new insights into understanding the large language models.

1 Introduction

In recent years, LLMs have emerged and demonstrated remarkable capability across various tasks (Vaswani et al., 2017; Liu et al., 2018; Devlin et al., 2018; Radford et al., 2019; Touvron et al., 2023; OpenAI, 2023; Liu et al., 2024; Guo et al., 2024, 2025). These models have shown impressive in-context learning abilities (Brown et al., 2020; Dong et al., 2022; Garg et al., 2022) and have been applied to logical reasoning problems, such as matching top human contestants at the International Mathematical Olympiad (IMO) level (Trinh et al., 2024) and solving math problems (Davies et al.,

2021). However, even the most advanced models still struggle with complex reasoning tasks. To truly enhance the reasoning capability of LLMs, it is crucial to investigate their intrinsic mechanisms.

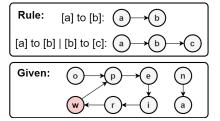
Multi-step reasoning tasks encompass a broad concept, typically referring to the ability of a model to synthesize numerous complex conditions to answer questions. Here, we consider a representative class of syntactic structures within multi-step reasoning tasks: a sentence that includes both the question and sufficient known information to answer it. For example, "Given: $[A] \rightarrow [B] \dots [B] \rightarrow [C] \dots$, Question: 2-step reasoning starting from [A]", where "..." represents other textual content unrelated to logical reasoning. The answer to this question is "[C]". When a sentence contains only one logical reasoning step, it is often handled by the so-called induction head in Transformer (Brown et al., 2020; Olsson et al., 2022). However, multistep reasoning is not merely a linear accumulation of multiple induction heads but involves more complex mechanisms.

Large models employ two primary strategies for logical reasoning. The first, known as the *Vertical Thinking Strategy (VTS)*, outputs reasoning results in a single forward based on their inherent structure. This approach is efficient but demands a high level of intelligence. As shown in Fig. 1, current large models exhibit significant limitations in their vertical thinking capability. Another relatively less efficient approach is the *(HTS)*, such as Chain of Thought (CoT) (Wei et al., 2022; Kojima et al., 2022) and Tree of Thought (ToT) (Yao et al., 2024), Diagram of Thought (DoT) (Zhang et al., 2024a). This strategy substantially enhances the model's reasoning performance. All models can produce the correct answers for tasks depicted

Question: We have established the following reasoning rules: (1) [a] to [b] represents that condition [a] can derive condition [b]. (2) The sequence [a] to [b] | [b] to [c] indicates that the 2-step reasoning result of [a] is [c]. For the following reasoning chain:

[e] to [i] | [r] to [w] | [n] to [a] | [o] to [p] | [i] to [r] | [p] to [e] | [w] to [p] | [x] to [i]

Please answer directly: What is the <s>-step reasoning result of [w]? (only return the answer)



Accuracy	y :			
	GPT4	GPT4o	Claude3.5 haiku	Clauld3.5 sonnet
2-step	52%	44%	24%	100%
3-step	0%	4%	0%	4%
4-step	0%	4%	4%	12%

Figure 1: The interaction results of multi-step reasoning tasks with large models. We tested each model 25 times, and when the number of reasoning steps exceeded three, these models exhibited random guessing. However, when we allowed the models to use CoT prompting (by removing "directly" and "only return the answer"), all models achieved 100% accuracy. Detailed interaction results are provided in the Appendix I.

in Fig. 1 with the help of CoT. Considering the strengths and weaknesses of these two strategies, the ideal approach should combine both: decomposing problems into several coarse-grained subproblems (HTS) and applying the VTS to each subproblem. Thus, researching how to improve vertical thinking capability and understanding why the horizontal thinking strategy can significantly enhance model reasoning ability are both crucial.

In this work, we investigate the performance of Transformer models on a symbolic multi-step reasoning dataset. Our work aims to uncover these mechanisms and provide insights into how Transformers process and integrate logical information across multiple layers to perform multi-step reasoning, which can help develop more effective strategies for improving their multi-step reasoning abilities. Specifically, we found that Transformers utilize a Buffer Mechanism when engaging in symbolic multi-step reasoning tasks. The model stores different intermediate results in separate buffers, allowing for quick retrieval as needed. We elaborate on how the model leverages this buffer mechanism for vertical thinking, and we explain why the horizontal thinking strategy can significantly enhance the model's multi-step reasoning capability from the perspective of the buffer mechanism. Finally, based on this understanding, we propose a method to improve the model's reasoning abilities, leading to significant performance improvements for the GPT-2 model on 7 multi-step reasoning datasets, including Clutrr (Sinha et al., 2019), RuleTaker (Clark et al., 2020), StepGame (Shi et al., 2022),

LogicInference (Ontanon et al., 2022), LogicAsker (Wan et al., 2024), PararulePlus (Bao et al., 2022), and PrOntoQA (Saparov and He, 2022).

The concept of "buffer" or similar concepts has also been mentioned in other works (Reddy, 2023; Bietti et al., 2024; Elhage et al., 2021). Our work provides a detailed description of the concept of buffer and applies this mechanism to enhance model performance.

The main contributions of this work are as follows:

- We propose a buffer mechanism and found evidence that supports such a mechanism being employed by language models during the reasoning process in symbolic multi-step reasoning tasks and provide a detailed analysis of the model's internal thinking process for vertical and horizontal thinking from the perspective of the buffer.
- We propose a method to enhance the model's reasoning capability, improving data utilization efficiency in 7 logical reasoning datasets.

Our research deepens the understanding of the reasoning mechanisms in Transformer models.

2 Reasoning Dataset and Transformer Model

Dataset. To understand the mechanism of multistep reasoning in Transformers, we design an abstract symbolic multistep reasoning task. As shown in Fig. 2, reasoning chains are serialized into a sequence. Every two tokens in the sentence represent a reasoning relation. The last token is

the reasoning start token, and the label is the result with a fixed-step reasoning starting from the starting point.

Transformer Model. We employ a decoder-only Transformer. Given an input sequence $\boldsymbol{X}^{\text{in}} \in \mathbb{R}^{n \times d}$, where n is the sequence length and d is the dictionary size, the model first applies an embedding layer to obtain the input representation $\boldsymbol{X}^{(1)} = \boldsymbol{X}_{tgt} + \boldsymbol{X}_{pos} \in \mathbb{R}^{n \times d_m}$, where \boldsymbol{X}_{tgt} and \boldsymbol{X}_{pos} denote the token and positional embeddings, respectively, and d_m is the hidden dimension. The single-head attention in each layer is computed as follows:

$$\mathcal{A}^{(l)}(\boldsymbol{X}) = \sigma \left(\frac{\operatorname{mask}(\boldsymbol{X} \boldsymbol{W}^{q(l)} \boldsymbol{W}^{k(l),\mathsf{T}} \boldsymbol{X}^{\mathsf{T}})}{\sqrt{d_k}} \right),$$
$$\boldsymbol{X}^{\operatorname{qkv}(l)} = \mathcal{A}^{(l)}(\bar{\boldsymbol{X}}^{(l)}) \bar{\boldsymbol{X}}^{(l)} \boldsymbol{W}^{v(l)} \boldsymbol{W}^{o(l)},$$

where σ takes the softmax operator, and $\bar{\boldsymbol{X}}^{(l)} = \text{Layernorm}(\boldsymbol{X}^{(l)})$. For simplicity of expression, we will abbreviate $\boldsymbol{W}^{q(l)}\boldsymbol{W}^{k(l),\mathsf{T}}$ as $\boldsymbol{W}^{qk(l)}$ and $\boldsymbol{W}^{v(l)}\boldsymbol{W}^{o(l),\mathsf{T}}$ as $\boldsymbol{W}^{vo(l)}$ in the following text. The output of the l-th layer is obtained as:

$$egin{aligned} oldsymbol{X}^{\mathrm{ao}(l)} &= oldsymbol{X}^{(l)} + oldsymbol{X}^{\mathrm{qkv}(l)}, \ oldsymbol{X}^{(l+1)} &= f^{(l)}(ar{oldsymbol{X}}^{\mathrm{ao}(l)}) + oldsymbol{X}^{\mathrm{ao}(l)}, \end{aligned}$$

where $f^{(l)}(\cdot)$ represents the feedforward neural network of l-th layer. The final output (in the form of token indices within the vocabulary) is obtained as:

$$Y = argmax(\sigma(\bar{X}^{(L)}W^p)) \in \mathbb{R}^n.$$

In Distribution and Out of Distribution Data.

With the settings of our dataset, if the model truly understands the underlying logic patterns, it should be able to find the correct answer to the sentence, even if this sentence has tokens that have never been encountered during the training. Therefore, we divided the data into two parts: indistribution (ID) and out-of-distribution (OOD). Specifically, we define token_{ID} $\in [20, 100]$ and $token_{OOD} \in [0, 120] \setminus [20, 100]$. In-distribution data $(Train_{ID} \ and \ Test_{ID})$ is defined as sentences composed entirely of token_{ID}, while out-of-distribution data (Test_{00D}) consists of sentences containing one token_{00D}, which happens to be the previous reasoning step of label. For the in-distribution data, we split the training set (Train_{ID}) and test set (Test_{ID}) according to the following rules: for the serialized reasoning chain $[x_1][x_2] \cdots [x_n]$ of the training set, all tokens satisfy the following condition:

$$x_{2i} - x_{2i-1} \pmod{m} \in G.$$

For the reasoning chains in the test set, all tokens satisfy:

$$x_{2i} - x_{2i-1} \pmod{m} \in \{1, \dots, m\} \setminus G,$$

where we take m=5 and $G=\{0,1,4\}$ in this study. Under this setting, we ensure that each binary logical pair in the testing set has not previously appeared in the training set. Therefore the Transformer is performing in-context learning (Brown et al., 2020; Olsson et al., 2022), as each reasoning pair is not seen during in-weight learning.

3 Vertical and Horizontal Thinking Strategy

We illustrate how the Transformer model employs vertical or horizontal thinking strategies for multistep reasoning. Fig. 3 depicts the schematic mechanisms by which these two strategies execute multi-step inference. The critical logic circuits highlighted in this figure were identified through causal intervention experiments (Feng and Steinhardt, 2023; Meng et al., 2022; Vig et al., 2020; Wang et al., 2024a) (Appendix H). From the figure, one can discern the core distinction between VTS and HTS. Under VTS, in every layer beyond the first—termed a "reasoning layer"—the final token attends to the token carrying the next inference result, conditional on the existing result, thereby effecting a single reasoning step. In contrast, HTS accomplishes all logical inference within a single layer.

This observation naturally prompts a key question: irrespective of whether VTS or HTS is employed, each node in the Transformer stores multiple pieces of information. For instance, the red node in Layer 1 encodes not only [a] but also [b]. How, then, does the model consolidate diverse information within a single node and accurately extract and utilize the relevant subset? Addressing this question lies at the heart of our investigation into multi-step reasoning.

4 Buffer Mechanism

To address the above problem, we propose a mechanism employed by the Transformer, which we refer to as the "buffer mechanism".

Specifically, we observe that the information transmitted via the attention module differs from that conveyed through the residual connections. In the attention module, information [a] with dimension d_m is first mapped to a lower embedding di-

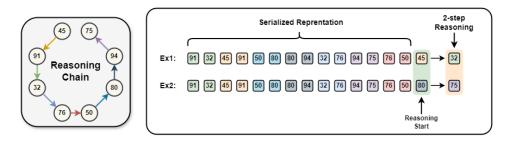


Figure 2: Illustration of the dataset. We investigate reasoning chains composed of digital tokens. In a serialized representation, each pair of adjacent tokens (represented by the same color) forms a reasoning relation within the reasoning chain. The order of the reasoning relations is random, thus a single reasoning chain can correspond to various serialized representations. The model input consists of the serialized representation along with a reasoning start token. The label is the result of performing a fixed number of reasoning steps starting from this start token. This figure illustrates a 2-step reasoning task.

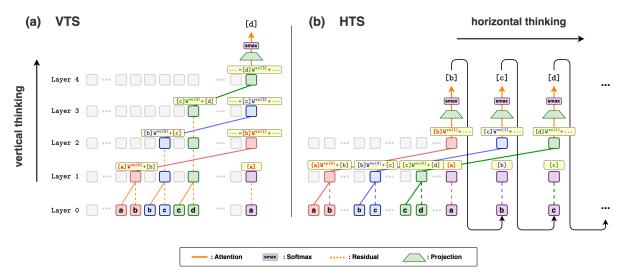


Figure 3: Illustration of how the Transformer employs the vertical and horizontal thinking strategies to solve the multi-step reasoning task. The information crucial for the layer's information transfer is highlighted. The first attention pairs tokens at odd positions with those at even positions. In the Vertical Thinking Strategy (VTS), each layer conveys its newly derived reasoning result to the final token. In the Horizontal Thinking Strategy (HTS), all information transfer is accomplished within a single layer.

mension d_v through the matrix \mathbf{W}^v , and then it is projected back to d_m dimensions via the matrix \mathbf{W}^o . Consequently, the information [a] after passing through the attention layer becomes [a] \mathbf{W}^{vo} . In contrast, information transmitted through residual connections does not require additional processing; thus, after the first layer, the information at even indices transforms to $[\mathbf{x}_{2i-1}]\mathbf{W}^{vo(0)} + [\mathbf{x}_{2i}]$.

Below, we introduce the buffer mechanism within the VTS framework. All theoretical results are supported by experimental validation (Section 5). We observe that each token in Fig. 3(a) stores various information in various forms (projected by different matrices $\boldsymbol{W}^{vo(l)}$), which leads to a natural question: how does the attention layer effectively retrieve useful information while avoiding

interference from others? To address this question, we introduce the following lemma (the proof can be found in Appendix C):

Lemma 4.1. Suppose token $\mathbf{x} = \sum_{i=1}^{n} \mathbf{a}_i \mathbf{W}_i \in \mathbb{R}^{d_m}$ and token $\mathbf{y} = \sum_{i=1}^{n} \mathbf{b}_i \mathbf{W}_i \in \mathbb{R}^{d_m}$, where $\mathbf{a}_i, \mathbf{b}_i \in \mathbb{R}^{d_m}$, $\mathbf{W}_i \in \mathbb{R}^{d_m \times d_m}$, $i = 1, 2, \cdots, n$. Each element of $\{\mathbf{a}_i\}_{i=1}^n$, $\{\mathbf{b}_i\}_{i=1}^n$ and $\{\mathbf{W}_i\}_{i=1}^n$ follows $\mathcal{N}(0, 1/d_m)$ and independent to others. Then, we have:

$$egin{aligned} m{x}m{W}_i^\mathsf{T} &= m{a}_i + \mathcal{O}\left(\sqrt{rac{n}{d_m}}
ight), \ m{y}m{W}_j^\mathsf{T} &= m{b}_j + \mathcal{O}\left(\sqrt{rac{n}{d_m}}
ight), \ m{x}m{W}_i^\mathsf{T}m{W}_jm{y}^\mathsf{T} &= m{a}_im{b}_j^\mathsf{T} + \mathcal{O}\left(rac{n}{\sqrt{d_m}}
ight). \end{aligned}$$

It can be observed that the matrices $\{W_i\}_{i=1}^n$ serve as a set of buffers for information. Each element of $\{a_i\}_{i=1}^n$ is located in different buffers, and is almost unaffected by others. This property also applies for $\{b_i\}_{i=1}^n$. By selecting the matrices associated with the relevant buffer, specific information contained in token x and token y can be extracted.

The concept of the buffer mechanism is useful for understanding the internal logic mechanisms of Transformer models. $\{\boldsymbol{W}^{q(l)}\}_{l=1}^{L}$ and $\{\boldsymbol{W}^{k(l)}\}_{l=1}^{L}$ can be viewed as "information extractors", while $\{\boldsymbol{W}^{vo(l)}\}_{l=1}^{L}$ can be considered as a set of buffers. We observe that in the final token of each reasoning layer, each new intermediate result is always associated with a new W^{vo} . If these $\{W^{vo(l)}\}_{l=1}^L$ matrices are mutually orthogonal or random, then these intermediate results can be regarded as being stored in a new buffer.

In each reasoning layer, the token at the last position contains the current intermediate result. Additionally, there exists a token that contains both the current intermediate result and the next step's reasoning result, with these stored in different buffers. The role of W^q and W^k is to extract the current intermediate results from these two tokens, enabling them to attend to each other due to having the same token. We refer to this feature as "same-token matching". To quantitatively characterize this property, we define the following match matrix:

$$h^{(1)}(\boldsymbol{X}) = \boldsymbol{X} \boldsymbol{W}^{qk(1)} \boldsymbol{W}^{vo(0), \mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \triangleq \boldsymbol{X} \operatorname{Ker}^{(1)} \boldsymbol{X}^{\mathsf{T}}$$

$$h^{(l)}(\boldsymbol{X}) \qquad (1)$$

$$= (\boldsymbol{X} \boldsymbol{W}^{vo(l-1)}) \boldsymbol{W}^{q(l)} \boldsymbol{W}^{k(l), \mathsf{T}} (\boldsymbol{X} \boldsymbol{W}^{vo(0)})^{\mathsf{T}}$$

$$= \boldsymbol{X} \boldsymbol{W}^{vo(l-1)} \boldsymbol{W}^{qk(l)} \boldsymbol{W}^{vo(0), \mathsf{T}} \boldsymbol{X}^{\mathsf{T}}$$

$$\triangleq \boldsymbol{X} \operatorname{Ker}^{(l)} \boldsymbol{X}^{\mathsf{T}}, \ l \geq 2. \tag{2}$$

where

$$\begin{split} & \operatorname{Ker}^{(1)} = \boldsymbol{W}^{qk(1)} \boldsymbol{W}^{vo(0),\mathsf{T}}, \\ & \operatorname{Ker}^{(l)} = \boldsymbol{W}^{vo(l-1)} \boldsymbol{W}^{qk(l)} \boldsymbol{W}^{vo(0),\mathsf{T}}, \ l \geq 2. \ \ (3) \end{split}$$

We temporarily ignore the effects introduced by the feedforward layers only in our analysis; a detailed version that includes the feedforward layers can be found in Appendix D.

To achieve the same-token matching, it is sufficient for $Ker^{(l)} \approx I$, in which case

$$h^{(l)}(\boldsymbol{X}_{tgt}) \approx \boldsymbol{X}_{tgt} \boldsymbol{X}_{tgt}^{\mathsf{T}} = \boldsymbol{I} + \mathcal{O}\left(\frac{1}{\sqrt{d_m}}\right).$$
 (4)

Eq.(4) means that the attention of all tokens is focused on themselves.

Furthermore, a much more remarkable observation is that the same-token matching is independent of the specific value of X_{tqt} . For example, for $X_{tqt, \rm OOD}$ sampled from the untrained random vectors token_{OOD}, $h^{(l)}(\boldsymbol{X}_{tgt,OOD}) \approx \boldsymbol{I}$ still holds. Therefore, when $\text{Ker}^{(l)} \approx I$ holds, the model exhibits OOD generalization capability.

Understanding HTS With Buffer. Numerous studies have shown that CoT, a canonical representative of HTS, can significantly enhance logical reasoning capability (Wei et al., 2022; Kojima et al., 2022). Fig. 3(b) presents a schematic diagram illustrating how the model implements this process. In the previously mentioned VTS scenario, the model needed to allocate new buffer $W^{vo(1)}$, $W^{vo(2)}$ for storing the intermediate result [b] and [c], respectively. In contrast, with CoT, the model generates a new token to separately store the new intermediate information, effectively replacing the information [a] in the original buffer $W^{vo(1)}$, $W^{vo(0)}$ and I(identity matrix, which can also be treated as a buffer). Thus, CoT performs multi-step reasoning tasks through buffer reuse.

Unlike the vertical thinking strategy, which requires alignment across multiple weight matrices, the horizontal thinking strategy only requires a few layers to satisfy $Ker^{(l)} \approx I$. This significantly reduces the difficulty of model training. These architectural efficiencies potentially underlie the $h^{(1)}(\boldsymbol{X}) = \boldsymbol{X} \boldsymbol{W}^{qk(1)} \boldsymbol{W}^{vo(0),\mathsf{T}} \boldsymbol{X}^\mathsf{T} \triangleq \boldsymbol{X} \operatorname{Ker}^{(1)} \boldsymbol{X}^\mathsf{T}, \text{ empirically observed enhancement of reasoning and the property of the proper$ performance in large language models employing CoT methodologies.

Experiments

This section presents the experimental validation of the buffer theory. To demonstrate multi-step reasoning without introducing additional complexity, we utilized a 3-layer, single-head Transformer model to learn from a 2-step reasoning dataset. Specific Experimental settings can be found in Appendix B. After training, the Transformer exhibits generalization capability in both in-distribution and out-ofdistribution data (Fig. 4).

Fig. 5 presents the experimental validation of Eq.(3)(4). We visualize the computed values of $h^{(i)}(X_{tqt})$ and $Ker^{(i)}$, i = 1, 2 with the model parameters after training. Notably, in the token_{00D} region, token_{00D} and $h^{(2)}(\boldsymbol{X}_{tat})$ exhibit a diagonal structure similar to an identity matrix I, ensuring the model's OOD generalization capability.

To establish a strong correlation between the

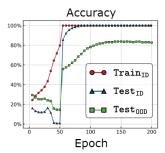


Figure 4: Accuracy curve during VTS training. After training, the Transformer exhibits generalization capability in in-distribution (100% accuracy) and out-of-distribution data (82% accuracy).

model's OOD generalization and its ability to use the buffer mechanism for same-token matching, we define a metric for this capability, the Matching Score (MS):

$$MS(h^{(l)}) = \mathbb{E}_{\boldsymbol{X}}[Trace(\sigma(h^{(l)}(\boldsymbol{X})))]/n, \quad (5)$$

where $X \in \mathbb{R}^{n \times d_m}$ is sampled from token_{ID} or token_{OOD} for expectation \mathbb{E}_X , and σ takes the softmax operation. The level of diagonalization of Ker^(l) serves as the intrinsic driver for achieving same-token matching; thus, we also define the following Kernel Score (KS):

$$KS(Ker^{(l)}) = Trace(\sigma(Ker^{(l)}))/d_m.$$
 (6)

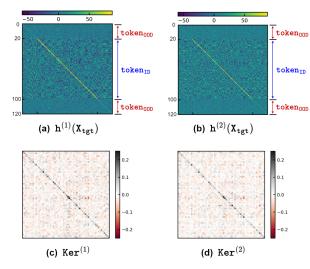


Figure 5: Heatmap of $h^{(1)}(\boldsymbol{X}_{tgt})$, $h^{(2)}(\boldsymbol{X}_{tgt})$, $\text{Ker}^{(1)}$ and $\text{Ker}^{(2)}$. According to (c)(d), and Eq.(4), the diagonal structure of the kernel matrix induces a diagonal structure in the matching matrix, even when \boldsymbol{X}_{tgt} is sampled from the token_{00D}.

Fig. 4 and Fig. 6 illustrates the accuracy curve and the dynamic changes in the matching score

and kernel score during training. It is observed that the increase in the model's ID and OOD generalization coincide with the increases in the model's matching score and kernel score. By computing the cosine similarity between different buffers $(\boldsymbol{W}^{vo(0)}, \boldsymbol{W}^{vo(1)}, \boldsymbol{W}^{vo(2)})$ and I, we observe that these buffers are nearly pairwise orthogonal (Appendix D). Based on the above experimental analysis, we conclude that the Transformer model indeed employs the buffer mechanism for vertical thinking.

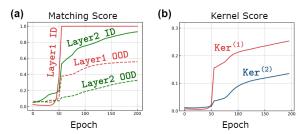


Figure 6: (a) The dynamic evolution of the model's matching score. The red and blue lines represent the matching scores for the first and second layers, respectively. Solid and dashed lines indicate the matching scores when X is drawn from token_{ID} and token_{00D}, respectively. (b) The kernel scores for the first (red) and second (blue) layers.

According to the Buffer mechanism, in theory, by setting the model's weights in the following way, an L-layer model can achieve (L-1)-step reasoning without the extra training process:

$$\mathbf{W}^{q(0)} = \mathbf{W}^{q(1)} = I, \ \mathbf{W}^{q(l)} = \mathbf{W}^{vo(l-1),\mathsf{T}},$$

$$oldsymbol{W}^{k(0)} = \sum_{i=1}^{[l_{ ext{seq}}/2]} oldsymbol{p}_{2i} oldsymbol{p}_{2i-1}^{\mathsf{T}}, \ oldsymbol{W}^{k(l)} = oldsymbol{W}^{vo(0),\mathsf{T}},$$

where $\{\boldsymbol{W}^{vo(l)}\}_{l=1}^{L}$ are set as random matrices and the projection layer satisfies $\boldsymbol{W}^{p} = \boldsymbol{W}^{vo(L),\mathsf{T}}\boldsymbol{W}^{\mathrm{emb},\mathsf{T}}$. Experimental validation can be found in Appendix D.

Real World LLMs. In Appendix G, we provide additional definitions for methods to compute the matching score and kernel score in multi-head models and perform these calculations in the real world LLMs such as GPT, Phi, Llama, and Qwen. Similar phenomena, such as weight alignment, provide evidence for the presence of the buffer mechanism in real language models.

Horizontal Thinking Experiments. For HTS, we will demonstrate that a Transformer model utilizing CoT can achieve arbitrary multi-step reasoning

with only 2 layers. We trained a 2-layer Transformer with the 13-length single-step reasoning data. During the testing phase, we fed the model's output back into the model. Through this CoT process, the model can perform 2-step, 3-step, or even higher-step reasoning and it can also generalize to sentence lengths beyond the 13th position. Fig. 7 shows the relationship between the number of reasoning steps and CoT accuracy. Our 2-layer model is able to maintain an accuracy of over 57.6% even when performing complex 6-step reasoning tasks. The specific information flow can be found in Appendix F.

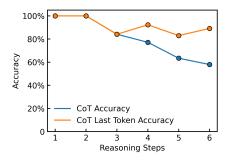


Figure 7: The relationship between the number of reasoning steps and CoT accuracy. The CoT accuracy curve represents the performance when all reasoning steps are predicted correctly using CoT, while the CoT last token accuracy curve records the performance based solely on whether the final reasoning step is correct when the previous correct reasoning result is given.

6 Improving Transformer's Data-Efficiency

In this section, we discuss how to improve data efficiency when employing the **vertical thinking strategy**, specifically, by enabling the model to learn the essence of logical reasoning with less data or training epoch. As mentioned in Section 3, this can be achieved by setting $\mathbf{W}^{vo(l-1),\mathsf{T}}\mathbf{W}^{qk(l)}\mathbf{W}^{vo(0),\mathsf{T}} \approx \mathbf{I}$.

A more intuitive approach is to replace $\boldsymbol{W}^{qk(l)}$ and $\boldsymbol{W}^{vo(l)}$ with $\boldsymbol{W}^{qk(l)} + \alpha^{(l)}I$ and $\boldsymbol{W}^{vo(l)} + \beta^{(l)}I$, where $\{\alpha^{(l)}\}_{l=1}^L$ and $\{\beta^{(l)}\}_{l=1}^L$ are learnable parameters. This Identity Matrix-Based Algorithm (denoted as IMBA) was first proposed in Boix-Adsera et al. (2023) and has been validated from both theoretical and empirical perspectives for its role in facilitating the model's learning of one-step reasoning data. However, in multi-step reasoning tasks, IMBA may lead to information interference. For instance, if $\boldsymbol{W}^{vo(l)}$ is replaced

with $\boldsymbol{W}^{vo(l)} + \beta^{(l)}\boldsymbol{I}$, the storage representation of the two pieces of information transitions from [a] \boldsymbol{W}^{vo} +[b] to [a] \boldsymbol{W}^{vo} +($\beta^{(l)}$ [a]+[b]), which introduces interference among the different pieces of information. Therefore, this approach may not be effective in enhancing the model's ability to perform multi-step reasoning.

Based on the understanding of buffer mechanism, we propose a **Random Matrix-Based Algorithm (RMBA)**, specifically by substituting $\boldsymbol{W}^{qk(l)}$ and $\boldsymbol{W}^{vo(l)}$ with $\boldsymbol{W}^{qk(l)} + \alpha^{(l)}\boldsymbol{Z}^{(l-1)}$ and $\boldsymbol{W}^{vo(l)} + \beta^{(l)}\boldsymbol{Z}^{(l)}$, where $\{\alpha^{(l)}\}_{l=1}^{L}$ and $\{\beta^{(l)}\}_{l=1}^{L}$ are learnable parameters and $\{\boldsymbol{Z}^{(l)}\}_{l=1}^{L}$ is a set of fixed random matrix following $N(0,1/d_m)$. In this case, the information storage representation changes from [a] \boldsymbol{W}^{vo} +[b] to [a] \boldsymbol{W}^{vo} +[a] \boldsymbol{Z} +[b], effectively creating a new buffer.

In RMBA, when $\alpha^{(l)} \cdot \beta^{(l)} \geq 0$, the diagonal elements of $(\boldsymbol{W}^{vo(l-1),\mathsf{T}} + \beta^{(l-1)}\boldsymbol{Z}^{(l-1),\mathsf{T}})(\boldsymbol{W}^{qk(l)} + \alpha^{(l)}\boldsymbol{Z}^{(l-1)})$ will exhibit positive increments, thereby enhancing the model's reasoning generalization capability. Fig. 8 supports this perspective, while also showing that IMBA, based on induction heads, fails to improve multi-step reasoning performance. To strengthen the credibility of our results, we conducted a comprehensive sweep of hyperparameters such as weight decay, learning rate, and hidden dimensions $(d_m$ and $d_k)$. The findings indicate that RMBA parameterization can more robustly reach a high accuracy by a certain number of training steps across a wider range of hyperparameters. Detailed experimental results can be found in Appendix E.

These results also demonstrate that multi-step reasoning is not achieved by simply "stacking" multiple single-step reasonings. An algorithm that enhances single-step reasoning may not be applicable to multi-step reasoning. Therefore, investigating the multi-step reasoning is an important and meaningful topic.

Validation of Algorithm Performance on Multiple Real Multi-Step Reasoning Datasets. In this section, we evaluate the performance of our algorithm using published real-world datasets. The datasets used include Clutrr, RuleTaker, StepGame, LogicInference, LogicAsker, PararulePlus, and PrOntoQA. These datasets assess the model's multistep reasoning capabilities from different perspectives and scenarios. Examples of the datasets and their corresponding training curves are provided in Appendix E.

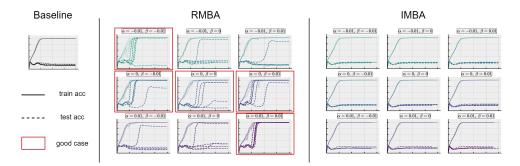


Figure 8: The accuracy comparison for the Baseline model, RMBA model, and IMBA model. The baseline is a three-layer transformer model, which fails to learn a two-step reasoning task with only 30,000 samples. The solid lines represent the training accuracy, while the dashed lines denote the test accuracy. For both RMBA and IMBA models, we set 9 different initialization parameter values $\alpha_{\rm ini}^{(l)}$ and $\beta_{\rm ini}^{(l)}$ and each experiment was conducted with 5 random seeds (each seed corresponds to one line). When $\alpha_{\rm ini}^{(l)}=0$ or $\alpha_{\rm ini}^{(l)}\cdot\beta_{\rm ini}^{(l)}>0$, the model's reasoning capability can be enhanced. Detailed settings can be found in Appendix E.

The 12-layer GPT-2 model is employed for this task. For RMBA, we replace $W^{qk(l,h)}$ and $m{W}^{vo(l,h)}$ with $m{W}^{qk(l,h)} + lpha^{(l,h)}m{Z}^{(l-1,h)}$ and $oldsymbol{W}^{vo(l,h)} + eta^{(l,h)} oldsymbol{Z}^{(l,h)}$, respectively. To demonstrate the effectiveness of our method, we train each dataset for 5 epochs using a learning rate of 1e-5. The additional hyperparameters introduced in RMBA are initialized to 0.05. Training results are summarized in Table 1. For tasks including Clutrr, RuleTalker, StepGame, and LogicInference, RMBA significantly improves model performance. For relatively simpler tasks such as LogicAsker, PararulePlus, and PrOntoQA, although both the base model and RMBA achieve nearly 100% accuracy, the training curves presented in Fig. 16 (Appendix E) show differences in the cost required to reach generalization capability. Table 2 reports the number of epochs required to reach 95% accuracy. As shown, RMBA reduces computational costs by more than 65% compared to the base model.

Table 1: Final-State Accuracy (%) Comparison. The training curves for all datasets are shown in Fig. 16 (Appendix E).

Dataset	Baseline	RMBA (ours)	IMBA
Clutrr	34.3	42.8	42.6
RuleTaker	67.0	83.5	74.9
StepGame	35.9	38.9	35.3
LogicInference	68.1	78.3	49.2
LogicAsker	99.7	100.0	97.1
PararulePlus	99.8	99.8	99.7
PrOntoQA	100.0	100.0	100.0

Table 2: Epochs required to reach 95% accuracy and cost savings comparison

Dataset	Baseline	RMBA(ours)	IMBA
LogicAsker	3.16	0.88	3.32
PararulePlus	2.71	0.73	2.02
PrOntoQA	0.95	0.32	1.25

7 Discussion

In this study, we investigated the vertical and horizontal thinking strategy employed by Transformer in a symbolic multi-step reasoning task from the perspective of the buffer mechanism. When utilizing the vertical thinking strategy, the model stores different intermediate results in separate buffers and transfers the reasoning results with the sametoken matching. In contrast, when applying the horizontal thinking strategy, the model reuses the existing buffers to store intermediate results. We validated that the buffer mechanism is a key factor in enabling the model's ID and OOD generalization capabilities. Based on the buffer mechanism, we proposed a tailored approach, RMBA, to enhance the model's multi-step reasoning ability, significantly improving data efficiency when training GPT-2 on 7 reasoning datasets.

8 Limitation

Our current work lacks an in-depth theoretical analysis. In future work, we aim to conduct deeper theoretical modeling for multi-step reasoning problems and extend the buffer mechanism to other types of multi-step reasoning tasks.

Acknowledgement

This work is sponsored by the National Key R&D Program of China Grant No. 2022YFA1008200, the National Natural Science Foundation of China Grant No. 92270001(Z. X.), 12371511 (Z. X.), 12422119 (Z. X.), 1257010106(Y. Z.), the Natural Science Foundation of Shanghai No. 25ZR1402280(Y. Z.), Shanghai Municipal of Science and Technology Major Project No. 2021SHZDZX0102 (Z. X., Y. Z.), and the HPC of School of Mathematical Sciences and the Student Innovation Center, and the Siyuan-1 cluster supported by the Center for High Performance Computing at Shanghai Jiao Tong University, Key Laboratory of Marine Intelligent Equipment and System, Ministry of Education, P.R. China. This work was partially supported by SJTU Kunpeng&Ascend Center of Excellence. We thank ChatGPT for assistance with language polishing, and text proofreading during the preparation of this manuscript.

References

- Emmanuel Abbe, Samy Bengio, Elisabetta Cornacchia, Jon Kleinberg, Aryo Lotfi, Maithra Raghu, and Chiyuan Zhang. 2022. Learning to reason with neural networks: Generalization, unseen data and boolean measures. *Advances in Neural Information Processing Systems*, 35:2709–2722.
- Emmanuel Abbe, Samy Bengio, Aryo Lotfi, Colin Sandon, and Omid Saremi. 2024. How far can transformers reason? the locality barrier and inductive scratchpad. *arXiv preprint arXiv:2406.06467*.
- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, and 1 others. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.
- Noah Amsel, Gilad Yehudai, and Joan Bruna. 2024. On the benefits of rank in attention layers. *arXiv preprint arXiv:2407.16153*.
- Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. 2018. A convergence analysis of gradient descent for deep linear neural networks. *arXiv preprint arXiv:1810.02281*.
- Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. 2019a. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pages 322–332. PMLR.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. 2019b.

- On exact computation with an infinitely wide neural net. In *Advances in Neural Information Processing Systems*, pages 8141–8150.
- Sanjeev Arora, Zhiyuan Li, and Abhishek Panigrahi. 2022. Understanding gradient descent on the edge of stability in deep learning. In *International Conference on Machine Learning*, pages 948–1024. PMLR.
- Murdock Aubry, Haoming Meng, Anton Sugolov, and Vardan Papyan. 2024. Transformer block coupling and its correlation with generalization in llms. *arXiv* preprint arXiv:2407.07810.
- Qiming Bao, Alex Yuxuan Peng, Tim Hartill, Neset Tan, Zhenyun Deng, Michael Witbrock, and Jiamou Liu. 2022. Multi-step deductive reasoning over natural language: An empirical study on out-of-distribution generalisation. *arXiv* preprint arXiv:2207.14000.
- Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. 2024. Birth of a transformer: A memory viewpoint. *Advances in Neural Information Processing Systems*, 36.
- Eden Biran, Daniela Gottesman, Sohee Yang, Mor Geva, and Amir Globerson. 2024. Hopping too late: Exploring the limitations of large language models on multi-hop queries. *arXiv preprint arXiv:2406.12775*.
- Enric Boix-Adsera, Omid Saremi, Emmanuel Abbe, Samy Bengio, Etai Littwin, and Joshua Susskind. 2023. When can transformers reason with abstract symbols? *arXiv preprint arXiv:2310.09753*.
- Jannik Brinkmann, Abhay Sheshadri, Victor Levoso, Paul Swoboda, and Christian Bartelt. 2024. A mechanistic analysis of a transformer trained on a symbolic multi-step reasoning task. *arXiv preprint arXiv:2402.11917*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Lei Chen, Joan Bruna, and Alberto Bietti. 2024a. How truncating weights improves reasoning in language models. *arXiv preprint arXiv:2406.03068*.
- Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. 2024b. Training dynamics of multi-head softmax attention for in-context learning: Emergence, convergence, and optimality. *arXiv* preprint *arXiv*:2402.19442.
- Xingwu Chen and Difan Zou. 2024. What can transformer learn with varying depth? case studies on sequence learning tasks. *arXiv preprint arXiv:2404.01601*.
- Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. Transformers as soft reasoners over language. *arXiv* preprint arXiv:2002.05867.

- Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. 2023. Towards automated circuit discovery for mechanistic interpretability. *Advances in Neural Information Processing Systems*, 36:16318–16352.
- Guy Dar, Mor Geva, Ankit Gupta, and Jonathan Berant. 2022. Analyzing transformers in embedding space. *arXiv preprint arXiv:2209.02535*.
- Alex Davies, Petar Veličković, Lars Buesing, Sam Blackwell, Daniel Zheng, Nenad Tomašev, Richard Tanburn, Peter Battaglia, Charles Blundell, András Juhász, and 1 others. 2021. Advancing mathematics by guiding human intuition with ai. *Nature*, 600(7887):70–74.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey on in-context learning. arXiv preprint arXiv:2301.00234.
- Subhabrata Dutta, Joykirat Singh, Soumen Chakrabarti, and Tanmoy Chakraborty. 2024. How to think step-by-step: A mechanistic understanding of chain-of-thought reasoning. *arXiv preprint arXiv:2402.18312*.
- Benjamin L Edelman, Ezra Edelman, Surbhi Goel, Eran Malach, and Nikolaos Tsilivis. 2024. The evolution of statistical induction heads: In-context learning markov chains. *arXiv preprint arXiv:2402.11004*.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, and 6 others. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*. Https://transformercircuits.pub/2021/framework/index.html.
- Jiahai Feng and Jacob Steinhardt. 2023. How do language models bind entities in context? *arXiv preprint arXiv:2310.17191*.
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. 2022. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598.
- Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. 2023. Localizing model behavior with path patching. *arXiv preprint arXiv:2304.05969*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.

- Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, and 1 others. 2024. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv* preprint arXiv:2401.14196.
- Tianyu Guo, Wei Hu, Song Mei, Huan Wang, Caiming Xiong, Silvio Savarese, and Yu Bai. 2023. How do transformers learn in-context beyond simple functions? a case study on learning with representations. arXiv preprint arXiv:2310.10616.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. 2018. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31.
- Arthur Jacot, Berfin Simsek, Francesco Spadaro, Clément Hongler, and Franck Gabriel. 2020. Implicit regularization of random feature models. In *International Conference on Machine Learning*, pages 4631–4640. PMLR.
- Sullam Jeoung and Jana Diesner. 2022. What changed? investigating debiasing methods using causal mediation analysis. *arXiv preprint arXiv:2206.00701*.
- Bowen Jiang, Yangxinyu Xie, Zhuoqun Hao, Xiaomeng Wang, Tanwi Mallick, Weijie J Su, Camillo J Taylor, and Dan Roth. 2024. A peek into token bias: Large language models are not yet genuine reasoners. *arXiv* preprint arXiv:2406.11050.
- Jihyung Kil, Farideh Tavazoee, Dongyeop Kang, and Joo-Kyung Kim. 2024. Ii-mmr: Identifying and improving multi-modal multi-hop reasoning in visual question answering. *arXiv* preprint *arXiv*:2402.11058.
- Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. 2020. Attention is not only a weight: Analyzing transformers with vector norms. *arXiv* preprint arXiv:2004.10102.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. Advances in neural information processing systems, 35:22199– 22213.
- Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. 2019. Revealing the dark secrets of bert. *arXiv preprint arXiv:1908.08593*.
- Yanyang Li, Shuo Liang, Michael R Lyu, and Liwei Wang. 2024a. Making long-context language models better multi-hop reasoners. *arXiv preprint arXiv:2408.03246*.

- Zhaoyi Li, Gangwei Jiang, Hong Xie, Linqi Song, Defu Lian, and Ying Wei. 2024b. Understanding and patching compositional reasoning in llms. *arXiv* preprint arXiv:2402.14328.
- Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. 2024c. Chain of thought empowers transformers to solve inherently serial problems. *arXiv* preprint *arXiv*:2402.12875.
- Zhiyuan Li, Sadhika Malladi, and Sanjeev Arora. 2021. On the validity of modeling sgd with stochastic differential equations (sdes). *Advances in Neural Information Processing Systems*, 34:12712–12725.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, and 1 others. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Peter J Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. 2018. Generating wikipedia by summarizing long sequences. arXiv preprint arXiv:1801.10198.
- Tao Luo, Zhi-Qin John Xu, Zheng Ma, and Yaoyu Zhang. 2021. Phase diagram for two-layer relu neural networks at infinite-width limit. *Journal of Machine Learning Research*, 22(71):1–47.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372.
- Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. 2023. Circuit component reuse across tasks in transformer language models. arXiv preprint arXiv:2310.08744.
- Samuel Müller, Noah Hollmann, Sebastian Pineda Arango, Josif Grabocka, and Frank Hutter. 2021. Transformers can do bayesian inference. *arXiv* preprint arXiv:2112.10510.
- Eshaan Nichani, Alex Damian, and Jason D Lee. 2024. How transformers learn causal structure with gradient descent. *arXiv preprint arXiv:2402.14735*.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, and 1 others. 2022. In-context learning and induction heads. arXiv preprint arXiv:2209.11895.
- Santiago Ontanon, Joshua Ainslie, Vaclav Cvicek, and Zachary Fisher. 2022. Logicinference: A new dataset for teaching logical inference to seq2seq models. *arXiv preprint arXiv:2203.15099*.
- OpenAI. 2023. Gpt-4 technical report. Preprint, arXiv:2303.08774.

- Michael Poli, Armin W Thomas, Eric Nguyen, Pragaash Ponnusamy, Björn Deiseroth, Kristian Kersting, Taiji Suzuki, Brian Hie, Stefano Ermon, Christopher Ré, and 1 others. 2024. Mechanistic design and scaling of hybrid architectures. *arXiv preprint arXiv:2403.17844*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, and 1 others. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Gautam Reddy. 2023. The mechanistic basis of data dependence and abrupt learning in an in-context classification task. *arXiv* preprint arXiv:2312.03002.
- Yinuo Ren, Chao Ma, and Lexing Ying. 2024. Understanding the generalization benefits of late learning rate decay. In *International Conference on Artificial Intelligence and Statistics*, pages 4465–4473. PMLR.
- Abulhair Saparov and He He. 2022. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. *arXiv preprint arXiv:2210.01240*.
- Yuval Shalev, Amir Feder, and Ariel Goldstein. 2024. Distributional reasoning in llms: Parallel reasoning processes in multi-hop reasoning. *arXiv preprint arXiv:2406.13858*.
- Pratyusha Sharma, Jordan T Ash, and Dipendra Misra. 2023. The truth is in there: Improving reasoning in language models with layer-selective rank reduction. *arXiv preprint arXiv:2312.13558*.
- Zhengxiang Shi, Qiang Zhang, and Aldo Lipani. 2022. Stepgame: A new benchmark for robust multi-hop spatial reasoning in texts. In *Proceedings of the AAAI conference on artificial intelligence*, 10, pages 11321–11329.
- Koustuv Sinha, Shagun Sodhani, Jin Dong, Joelle Pineau, and William L Hamilton. 2019. Clutrr: A diagnostic benchmark for inductive reasoning from text. *arXiv preprint arXiv:1908.06177*.
- Eric Todd, Millicent L Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. 2023. Function vectors in large language models. *arXiv* preprint arXiv:2310.15213.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. 2024. Solving olympiad geometry without human demonstrations. *Nature*, 625(7995):476–482.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all

- you need. Advances in neural information processing systems, 30.
- Jesse Vig. 2019. A multiscale visualization of attention in the transformer model. *arXiv preprint arXiv:1906.05714*.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. Investigating gender bias in language models using causal mediation analysis. Advances in neural information processing systems, 33:12388– 12401.
- Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. 2019. Analyzing multihead self-attention: Specialized heads do the heavy lifting, the rest can be pruned. *arXiv preprint arXiv:1905.09418*.
- Yuxuan Wan, Wenxuan Wang, Yiliu Yang, Youliang Yuan, Jen-tse Huang, Pinjia He, Wenxiang Jiao, and Michael R Lyu. 2024. Logicasker: Evaluating and improving the logical reasoning ability of large language models. *arXiv* preprint arXiv:2401.00757.
- Boshi Wang, Xiang Yue, Yu Su, and Huan Sun. 2024a. Grokked transformers are implicit reasoners: A mechanistic journey to the edge of generalization. *arXiv* preprint arXiv:2405.15071.
- Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2022. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small. *arXiv preprint arXiv:2211.00593*.
- Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun. 2023. Label words are anchors: An information flow perspective for understanding in-context learning. *arXiv* preprint *arXiv*:2305.14160.
- Mingze Wang, Haotian He, Jinbo Wang, Zilin Wang, Guanhua Huang, Feiyu Xiong, Zhiyu Li, Lei Wu, and 1 others. 2024b. Improving generalization and convergence by enhancing implicit regularization. *arXiv* preprint arXiv:2405.20763.
- Mingze Wang and E Weinan. 2024. Understanding the expressive power and mechanisms of transformer for sequence modeling. *arXiv preprint arXiv:2402.00522*.
- Mingze Wang and Lei Wu. 2023. A theoretical analysis of noise geometry in stochastic gradient descent. *arXiv preprint arXiv:2310.00692*.
- Xinyi Wang, Alfonso Amayuelas, Kexun Zhang, Liangming Pan, Wenhu Chen, and William Yang Wang. 2024c. Understanding the reasoning ability of language models from the perspective of reasoning paths aggregation. *arXiv preprint arXiv:2402.03268*.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824– 24837.
- Lei Wu, Chao Ma, and 1 others. 2018. How sgd selects the global minima in over-parameterized learning: A dynamical stability perspective. *Advances in Neural Information Processing Systems*, 31.
- Lei Wu and Weijie J Su. 2023. The implicit regularization of dynamical stability in stochastic gradient descent. In *International Conference on Machine Learning*, pages 37656–37684. PMLR.
- Zhi-Qin John Xu, Yaoyu Zhang, Tao Luo, Yanyang Xiao, and Zheng Ma. 2019. Frequency principle: Fourier analysis sheds light on deep neural networks. *arXiv preprint arXiv:1901.06523*.
- Sohee Yang, Elena Gribovskaya, Nora Kassner, Mor Geva, and Sebastian Riedel. 2024. Do large language models latently perform multi-hop reasoning? *arXiv* preprint arXiv:2402.16837.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.
- Chiyuan Zhang, Maithra Raghu, Jon Kleinberg, and Samy Bengio. 2021. Pointer value retrieval: A new benchmark for understanding the limits of neural network generalization. *arXiv preprint arXiv:2107.12580*.
- Jiahao Zhang, Haiyang Zhang, Dongmei Zhang, Yong Liu, and Shen Huang. 2023. Beam retrieval: General end-to-end retrieval for multi-hop question answering. *arXiv preprint arXiv:2308.08973*.
- Yifan Zhang, Yang Yuan, and Andrew Chi-Chih Yao. 2024a. On the diagram of thought. *arXiv preprint arXiv:2409.10038*.
- Zhongwang Zhang, Pengxiao Lin, Zhiwei Wang, Yaoyu Zhang, and Zhi-Qin John Xu. 2024b. Initialization is critical to whether transformers fit composite functions by inference or memorizing. *arXiv* preprint *arXiv*:2405.05409.
- Zhongwang Zhang, Zhiwei Wang, Junjie Yao, Zhangchen Zhou, Xiaolong Li, Zhi-Qin John Xu, and 1 others. 2024c. Anchor function: a type of benchmark functions for studying language models. arXiv preprint arXiv:2401.08309.
- Hanxu Zhou, Qixuan Zhou, Zhenyuan Jin, Tao Luo, Yaoyu Zhang, and Zhi-Qin John Xu. 2022. Empirical phase diagram for three-layer neural networks with infinite width. *Advances in Neural Information Processing Systems*.

Zhanxing Zhu, Jingfeng Wu, Bing Yu, Lei Wu, and Jinwen Ma. 2018. The anisotropic noise in stochastic gradient descent: Its behavior of escaping from sharp minima and regularization effects. *arXiv preprint arXiv:1803.00195*.

A Related Work

Language model reasoning There have been numerous experimental and theoretical studies on language model reasoning. (Abbe et al., 2022) examines the reasoning capabilities of neural networks using the Pointer Value Retrieval (PVR) benchmark, which was originally introduced in (Zhang et al., 2021). (Sharma et al., 2023) demonstrates that applying low-rank approximation of certain weight layers to themselves can enhance reasoning performance across various tasks, and (Chen et al., 2024a) explains this phenomenon in a two-layer Transformer model. (Wang et al., 2024a) investigates both ID and OOD reasoning abilities on two synthetic tasks involving composition and comparison. (Jiang et al., 2024) reveals that the reasoning process in Large Language Models (LLMs) is influenced by token bias and that these models continue to face challenges when dealing with more complex logical tasks. (Abbe et al., 2024) introduces the distribution locality and shows that Transformers require low locality. (Zhang et al., 2024b; Luo et al., 2021; Zhou et al., 2022) shows that small initialization can facilitate model reasoning. Brinkmann et al. (2024) investigates the mechanism by which language models output intermediate reasoning paths in multi-step reasoning tasks. Aubry et al. (2024) uncover a transformer block coupling phenomenon in a variety of LLMs by tracing the trajectories of individual tokens as they pass through transformer blocks. Recently, the multi-hop reasoning abilities of models have garnered attention in the LLM field. Kil et al. (2024); Li et al. (2024a) promote the use of Chain-of-Thought (CoT) reasoning by models to perform multi-step thinking through appropriate prompt engineering. Zhang et al. (2023) proposes a new retrieval framework for multiple reasoning paths. Dar et al. (2022); Li et al. (2024b); Yang et al. (2024); Shalev et al. (2024) pinpoint where models execute multi-step reasoning by causal intervention and observing neuron activation. Biran et al. (2024) introduces a method to enhance models' reasoning abilities by repeatedly invoking intermediate layers. These works, especially the causal intervention experiments, have inspired our research. However, existing studies primarily conduct macro-level statistical analyses on actual complex large language models. While causal intervention methods can help us identify critical paths, our work builds upon this foundation by further exploring how models leverage their own weights to generate these key paths. Experiments on our symbolic datasets also facilitate more in-depth experimental and theoretical investigations.

Understanding the mechanism of neural model Our work builds upon previous studies on the attention mechanism (Voita et al., 2019; Vig, 2019; Kovaleva et al., 2019; Kobayashi et al., 2020). Numerous researchers have proposed various approaches to identify the roles of different heads in Transformers (Vig et al., 2020; Jeoung and Diesner, 2022; Wang et al., 2022; Conmy et al., 2023; Merullo et al., 2023; Guo et al., 2023; Wang and Weinan, 2024; Amsel et al., 2024; Li et al., 2024c; Wang et al., 2024c). These methods predominantly employ the concept of perturbation. Similar to the observations made by Wang et al. (2023) and Dutta et al. (2024), who noted that large language models typically perform information aggregation in shallow layers and information induction in deeper layers, we have also observed comparable phenomena in our study. The idea of symbolic datasets is inspired by Poli et al. (2024); Zhang et al. (2024c). There have also been some insightful theoretical works on feedforward neural networks. A series of studies have explored neural network preferences and generalization from the perspectives of regularization and frequency, etc. (Xu et al., 2019; Wang et al., 2024b; Jacot et al., 2018, 2020; Arora et al., 2019a, 2018). And Wu and Su (2023); Wang and Wu (2023); Arora et al. (2022); Li et al. (2021); Wu et al. (2018); Zhu et al. (2018); Arora et al. (2019b) investigates the dynamical behavior of neural networks, while Ren et al. (2024) examines the factors influencing neural network generalization.

In-context learning and induction head Our work primarily investigates the model's ability to perform in-context learning. The concept of in-context learning (ICL) was first introduced by Brown et al. (2020). Since then, a series of studies have utilized induction heads to investigate ICL, yielding remarkable research outcomes (Olsson et al., 2022; Garg et al., 2022; Wang et al., 2022; Müller et al., 2021; Goldowsky-Dill et al., 2023; Bietti et al., 2024; Nichani et al., 2024; Edelman et al., 2024; Chen et al., 2024b; Todd et al., 2023; Chen and Zou, 2024). It is worth noting that induction heads can be considered as a special case of multi-step reasoning tasks with the reasoning step equal to 1. However, multi-step reasoning is not a simple linear combination of single-step reasoning. In our work, we study

the mechanism that enables multi-step reasoning, which has not been explored in previous studies.

B Experimental settings

In our experiments (Section 3), the vocabulary size is set to d=201, and the hidden space dimension is set to $d_m=400$ and $d_q=d_k=d_v=64$. We use 200,000 2-step reasoning sequences (with sequence length equal to 13). The learning rate is set to 2e-5 and linearly warms up to 1e-4 within 400 epochs and then decays to 1e-5 within 3600 epochs. The batch size is set to 100. We use the AdamW optimizer with default parameters as set in PyTorch 2.3.0.

The experiments were conducted on a server with the following configuration:

- 64 AMD EPYC 7742 64-Core Processor.
- 256GB of total system memory.
- 2 NVIDIA A100 GPUs with 40GB of video memory each and 8 NVIDIA GeForce RTX 4080 GPUs with 16GB of video memory each.
- The experiments were run using the Ubuntu 22.04 LTS operating system.

The task shown above can be completed within 2 hours with a single NVIDIA A100 GPU. For other, more complex examples, they can be finished within 24 hours.

C Proof of Lemma 4.1

Lemma 1 Suppose token $\mathbf{x} = \sum_{i=1}^n \mathbf{a}_i \mathbf{W}_i \in \mathbb{R}^{d_m}$ and token $\mathbf{y} = \sum_{i=1}^n \mathbf{b}_i \mathbf{W}_i \in \mathbb{R}^{d_m}$, where $\mathbf{a}_i, \mathbf{b}_i \in \mathbb{R}^{d_m}$, $\mathbf{W}_i \in \mathbb{R}^{d_m \times d_m}$, $i = 1, 2, \dots, n$. Each element of $\{\mathbf{a}_i\}_{i=1}^n$, $\{\mathbf{b}_i\}_{i=1}^n$ and $\{\mathbf{W}_i\}_{i=1}^n$ follows a normal distribution $\mathcal{N}(0, 1/d_m)$ and independent to others. Then, we have:

$$\boldsymbol{x}\boldsymbol{W}_{i}^{\mathsf{T}} = \boldsymbol{a}_{i} + \mathcal{O}\left(\sqrt{\frac{n}{d_{m}}}\right), \quad \boldsymbol{y}\boldsymbol{W}_{j}^{\mathsf{T}} = \boldsymbol{b}_{j} + \mathcal{O}\left(\sqrt{\frac{n}{d_{m}}}\right),$$
 (7)

$$xW_i^{\mathsf{T}}W_jy^{\mathsf{T}} = a_ib_j^{\mathsf{T}} + \mathcal{O}\left(\frac{n}{\sqrt{d_m}}\right).$$
 (8)

Proof. We first show that $W_iW_j^{\mathsf{T}}$ (denoted as $Z^{(i,j)}$) is also a random matrix with elements following a normal distribution $\mathcal{N}(0,1/d_m)$ when $i\neq j$. In fact,

$$e[(\boldsymbol{W}_{i}\boldsymbol{W}_{j}^{\mathsf{T}})_{s,t}] = e[\sum_{k=1}^{d_{m}} (\boldsymbol{W}_{i})_{sk} (\boldsymbol{W}_{j})_{kt}] = \sum_{k=1}^{d_{m}} e[(\boldsymbol{W}_{i})_{sk}] e[(\boldsymbol{W}_{j})_{kt}] = 0,$$

$$\operatorname{Var}\left[(\boldsymbol{W}_{i}\boldsymbol{W}_{j}^{\mathsf{T}})_{s,t}\right] = e\left[\left(\sum_{k=1}^{d_{m}} (\boldsymbol{W}_{i})_{sk} (\boldsymbol{W}_{j})_{kt}\right)^{2}\right]$$

$$\operatorname{var}\left[\left(\boldsymbol{W}_{i}\boldsymbol{W}_{j}\right)_{s,t}\right] = \operatorname{e}\left[\left(\sum_{k=1}^{\infty}\left(\boldsymbol{W}_{i}\right)_{sk}\left(\boldsymbol{W}_{j}\right)_{kt}\right)\right]$$
$$= \sum_{k=1}^{d_{m}}\operatorname{e}\left[\left(\boldsymbol{W}_{i}\right)_{sk}^{2}\right]\operatorname{e}\left[\left(\boldsymbol{W}_{j}\right)_{kt}^{2}\right] = d_{m} \times \left(\frac{1}{d_{m}}\right)^{2} = \frac{1}{d_{m}},$$

which indicate $\{Z^{(i,j)}\}$ follows the same distribution as $\{W_i\}_{i=1}^n$. Therefore,

$$\begin{aligned} \mathbf{e}_{\{\boldsymbol{W}_{j}\}} \left[(\boldsymbol{x} \boldsymbol{W}_{i}^{\mathsf{T}})_{t} \right] &= \sum_{j=1}^{n} \mathbf{e}_{\{\boldsymbol{Z}^{(j,i)}\}_{j}} \left[(\boldsymbol{a}_{j} \boldsymbol{Z}^{(j,i)})_{t} \right] \\ &= \sum_{k=1}^{d_{m}} (\boldsymbol{a}_{i})_{k} \mathbf{e} \left[(\boldsymbol{Z}^{(i,i)})_{kt} \right] + \sum_{\substack{j=1\\j \neq i}}^{n} \sum_{k=1}^{d_{m}} (\boldsymbol{a}_{j})_{k} \mathbf{e} \left[(\boldsymbol{Z}^{(j,i)})_{kt} \right] \\ &= \sum_{k=1}^{d_{m}} (\boldsymbol{a}_{i})_{k} \mathbf{e} \left[(\boldsymbol{Z}^{(i,i)})_{kt} \right] = (\boldsymbol{a}_{i})_{t}, \end{aligned}$$

$$\operatorname{Var}_{\{\boldsymbol{W}_{j}\}}\left[(\boldsymbol{x}\boldsymbol{W}_{i}^{\mathsf{T}})_{t}\right] = \sum_{j=1}^{n} \operatorname{Var}_{\{\boldsymbol{Z}^{(j,i)}\}_{j}}\left[(\boldsymbol{a}_{j}\boldsymbol{Z}^{(j,i)})_{t}\right]$$

$$= \sum_{k=1}^{d_{m}} (\boldsymbol{a}_{i})_{k}^{2} \operatorname{Var}\left[(\boldsymbol{Z}^{(i,i)})_{kt}\right] + \sum_{\substack{j=1\\j\neq i}}^{n} \sum_{k=1}^{d_{m}} (\boldsymbol{a}_{j})_{k}^{2} \operatorname{Var}\left[(\boldsymbol{Z}^{(j,i)})_{kt}\right]$$

$$= \sum_{k=1}^{d_{m}} (\boldsymbol{a}_{i})_{k}^{2} \operatorname{Var}\left[(\boldsymbol{Z}^{(i,i)})_{kt}\right] + \frac{1}{d_{m}} \sum_{\substack{j=1\\j\neq i}}^{n} \sum_{k=1}^{d_{m}} (\boldsymbol{a}_{j})_{k}^{2}$$

$$= (\boldsymbol{a}_{i})_{t}^{2} \operatorname{Var}\left[(\boldsymbol{Z}^{(i,i)})_{tt}\right] + \sum_{k=1\atop k\neq t}^{d_{m}} (\boldsymbol{a}_{i})_{k}^{2} \operatorname{Var}\left[(\boldsymbol{Z}^{(i,i)})_{kt}\right] + \frac{1}{d_{m}} \sum_{\substack{j=1\\j\neq i}}^{n} \sum_{k=1}^{d_{m}} (\boldsymbol{a}_{j})_{k}^{2}$$

$$= \frac{2}{d_{m}} (\boldsymbol{a}_{i})_{t}^{2} + \frac{1}{d_{m}} \sum_{k=1\atop k\neq t}^{d_{m}} (\boldsymbol{a}_{i})_{k}^{2} + \frac{1}{d_{m}} \sum_{\substack{j=1\\j\neq i}}^{n} \sum_{k=1}^{d_{m}} (\boldsymbol{a}_{j})_{k}^{2}$$

$$= \frac{1}{d_{m}} (\boldsymbol{a}_{i})_{t}^{2} + \frac{1}{d_{m}} \sum_{i=1}^{n} \sum_{k=1}^{d_{m}} (\boldsymbol{a}_{j})_{k}^{2}.$$

Therefore, $\operatorname{Var}\left[(\boldsymbol{x}\boldsymbol{W}_i^{\mathsf{T}})_t\right] = \operatorname{Var}_{\boldsymbol{a}}\left[\operatorname{Var}_{\{\boldsymbol{W}_j\}}\left[(\boldsymbol{x}\boldsymbol{W}_i^{\mathsf{T}})_t\right]\right] = n/d_m + 1/d_m^2$. And Chebyshev's inequality implies that $\boldsymbol{x}\boldsymbol{W}_i^{\mathsf{T}} = \boldsymbol{a}_i + \mathcal{O}\left(\sqrt{\frac{n}{d_m}}\right)$, which also holds for $\boldsymbol{y}\boldsymbol{W}_j^{\mathsf{T}}$.

Assume that $\boldsymbol{x}\boldsymbol{W}_i^\mathsf{T} = \boldsymbol{a}_i + \boldsymbol{z}_1^\mathsf{T}, \ \boldsymbol{y}\boldsymbol{W}_j^\mathsf{T} = \boldsymbol{b}_j + \boldsymbol{z}_2, \ \boldsymbol{z}_1, \ \boldsymbol{z}_2$ are random vector with the elements follow the normal distribution $\mathcal{N}(n, 1/d_m)$, then,

$$\operatorname{Var}\left[\boldsymbol{x}\boldsymbol{W}_{i}^{\mathsf{T}}\boldsymbol{W}_{j}\boldsymbol{y}^{\mathsf{T}}\right] = \operatorname{Var}\left[\boldsymbol{a}_{i}\boldsymbol{b}_{j}^{\mathsf{T}} + \boldsymbol{a}_{i}\boldsymbol{z}_{2}^{\mathsf{T}} + \boldsymbol{z}_{1}\boldsymbol{b}_{j}^{\mathsf{T}} + \boldsymbol{z}_{1}\boldsymbol{z}_{2}^{\mathsf{T}}\right]$$
$$= d_{m} \cdot \frac{1}{d_{m}} \cdot \frac{1}{d_{m}} + 2 \cdot d_{m} \cdot \frac{1}{d_{m}} \cdot \frac{n}{d_{m}} + d_{m} \cdot \frac{n}{d_{m}} \cdot \frac{n}{d_{m}} = \frac{(n+1)}{d_{m}}.$$

Thus we have $m{x}m{W}_i^\mathsf{T}m{W}_jm{y}^\mathsf{T} = m{a}_im{b}_j^\mathsf{T} + \mathcal{O}\left(\frac{n}{\sqrt{d_m}}\right)$.

D Further Discussion on the Vertical Thinking Strategy

D.1 Information Fusion

In our symbolic dataset task, as mentioned in Section 3, the first layer facilitates the fusion of token information at odd and even positions. We find that positional encoding plays a crucial role in the features of the first layer of attention. Fig. 9(a)(b) illustrates a comparison between the original attention mechanism and the positional attention mechanism calculated with eq. 9. As shown, there is minimal difference between the two approaches.

$$\mathcal{A}^{(0)}(\boldsymbol{X}^{\text{pos}}) = \operatorname{softmax}\left(\frac{\operatorname{mask}(\boldsymbol{X}^{\text{pos}}\boldsymbol{W}^{qk}\boldsymbol{X}^{\text{pos},\mathsf{T}})}{\sqrt{d_k}}\right). \tag{9}$$

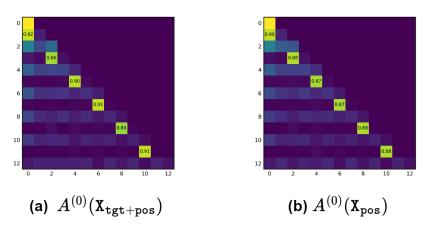


Figure 9: A comparison between the original attention (a) and the positional attention (b) of the first Transformer block, where $X_{tgt+pos} = X_{tgt} + X_{pos}$.

D.2 Detailed Matching Matrix

In Section 3, for simplicity of analysis, we ignored the impact of the feedforward layer. Here, we define a detailed version of the matching matrix as follows:

$$\tilde{h}^{(1)}(\mathbf{X}) = f^{(0)}(\mathbf{X}) \mathbf{W}^{q(1),\mathsf{T}} [f^{(0)}(\mathbf{X} \mathbf{W}^{vo(0)}) \mathbf{W}^{k(1),\mathsf{T}}]^{\mathsf{T}}$$
(10)

$$\tilde{h}^{(2)}(\boldsymbol{X}) = f^{(1)}[f^{(0)}(\boldsymbol{X})\boldsymbol{W}^{vo(1)}]\boldsymbol{W}^{q(2),\mathsf{T}} \left[f^{(1)} \circ f^{(0)}(\boldsymbol{X}\boldsymbol{W}^{vo(0)})\boldsymbol{W}^{k(2),\mathsf{T}} \right]^{\mathsf{T}}. \tag{11}$$

As shown in Fig. 10, the detailed matching matrices still maintain the diagonal element property in most cases, even for the out-of-distribution tokens.

D.3 Independence of Buffers

To verify the independence between buffers, we computed and visualized the cosine similarity between row vectors of different buffers $(\boldsymbol{W}^{vo(0)}, \boldsymbol{W}^{vo(1)}, \boldsymbol{W}^{vo(2)})$ and \boldsymbol{I} . As shown in Fig. 11, apart from exhibiting a certain similarity within itself (so-called condense phenomenon(Luo et al., 2021)), each buffer remains nearly orthogonal to the others.

D.4 Weight Construction Method for Multi-Step Reasoning Networks

In Section 3, we mentioned that by setting the weights in the following manner, we can enable an L-layer Transformer model to possess (L-1)-step reasoning capability.

$$W^{q(0)} = W^{q(1)} = I, \ W^{q(l)} = W^{vo(l-1),\mathsf{T}}, \ l \ge 2,$$
 (12)

$$\mathbf{W}^{k(0)} = \sum_{i=1}^{[l_{\text{seq}}/2]} \mathbf{p}_{2i} \mathbf{p}_{2i-1}^{\mathsf{T}}, \ \mathbf{W}^{k(l)} = \mathbf{W}^{vo(0),\mathsf{T}}, \ l \ge 1,$$
(13)

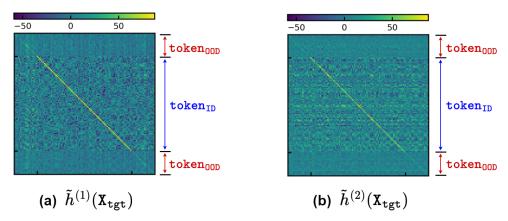


Figure 10: (a) Heatmap of $\tilde{h}^{(1)}(\boldsymbol{X}_{tgt})$ and $\tilde{h}^{(2)}(\boldsymbol{X}_{tgt})$. The diagonal elements exhibit the largest values, confirming the matching operation.

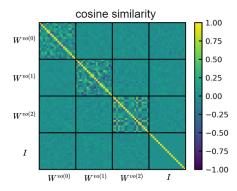


Figure 11: The cosine similarity between row vectors of different buffers $(\mathbf{W}^{vo(0)}, \mathbf{W}^{vo(1)}, \mathbf{W}^{vo(2)})$ and \mathbf{I}).

where $\{ {m W}^{vo(l)} \}_{l=1}^L$ are random matrices and the projection weight ${m W}^p = {m W}^{vo(L),\mathsf{T}} {m W}^{\mathrm{emb},\mathsf{T}}.$

The specific construction method is as follows: To ensure that each buffer in the model has adequate robustness against interference, we set $d_m = d_q = d_k = d_v = 10000$. The feedforward layers are assigned zero weights so that the residual connection dominates. Since all the weight matrices we use are untrained random matrices, the layer normalization will have no effect. Fig. 12 shows the multi-step reasoning ability of an 8-layer Transformer. We tested natural order, reverse order, random order sentences, and sentences with inserted irrelevant tokens (i.e., token [20]), and the model was able to output the correct answer [8] in all cases. Each sentence begins with token [20] to prevent $\mathcal{A}_{0,0}^{(l)}$ from always equaling 1, which could affect the buffer.

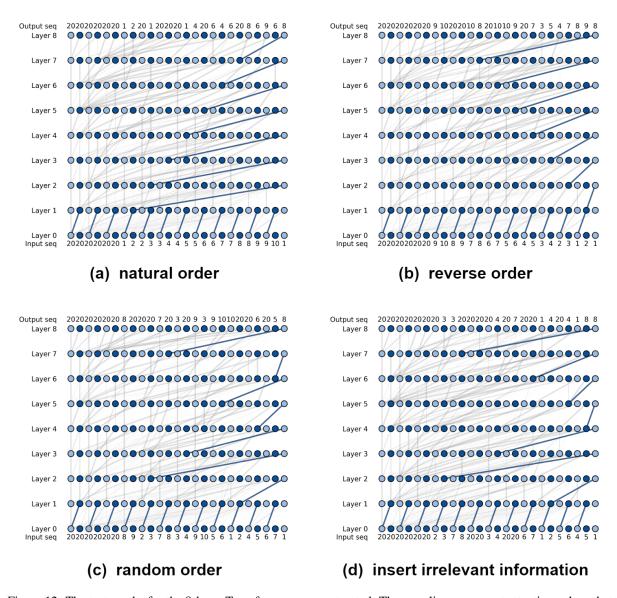


Figure 12: The test results for the 8-layer Transformer we constructed. The gray lines represent attention values that do not affect the final outcome. The width of all lines is positively correlated with the attention weights.

E Details for RMBA Experiment

This section provides supplementary details on the experimental setup for the RMBA experiment. We use a 3-layer single-head Transformer with $d_q = d_k = d_v = 64$ and $d_m = 400$. The training set consists of 30,000 2-step reasoning chains. We trained the model for 200 epochs in total. In this setting, the Transformers have poor generation ability even in the in-distribution test dataset.

We replace $\boldsymbol{W}^{qk(l)}$ and $\boldsymbol{W}^{vo(l)}$ with $\boldsymbol{W}^{qk(l)} + \alpha^{(l)}\boldsymbol{Z}^{(l-1)}$ and $\boldsymbol{W}^{vo(l)} + \beta^{(l)}\boldsymbol{Z}^{(l)}$, respectively, where $\alpha^{(l)}$ and $\beta^{(l)}$ are learnable parameters, and $\{\boldsymbol{Z}^{(l)}\}_{l=1}^{L}$ is a set of random matrix following $N(0,1/d_m)$. Therefore, 6 extra learnable parameters are added to this 3-layer single-head model in total. Fig. 13 shows the loss of the Transformer under different settings. Fig. 14 shows the changes of the learnable parameters $\alpha^{(l)}$ and $\beta^{(l)}$ during training.

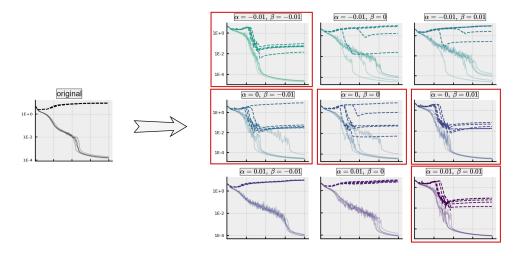


Figure 13: The impact of different learnable parameters' initial values, $\alpha^{(l)}$ and $\beta^{(l)}$, on the model's reasoning ability. The solid lines represent the training loss, while the dashed lines denote the test loss. Each experiment was conducted with 5 random seeds.

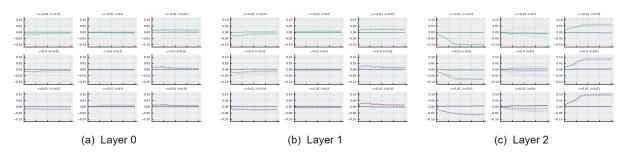


Figure 14: Changes of the learnable parameters $\alpha^{(l)}$ and $\beta^{(l)}$ during training. The solid lines represent the $\alpha^{(l)}$, while the dashed lines denote the $\beta^{(l)}$. Each experiment was conducted with 5 random seeds.

We further tested the performance of the three algorithms under different hyperparameter configurations. We investigated the effects of weight decay, learning rate, and hidden dimension on training. For the RMBA and IMBA algorithms, we set $\alpha_{\rm ini}^{(l)} = \beta_{\rm ini}^{(l)} = 0.01$. As shown in Fig. 15, under various settings, the RMBA algorithm consistently facilitated the model's ability to generalize.

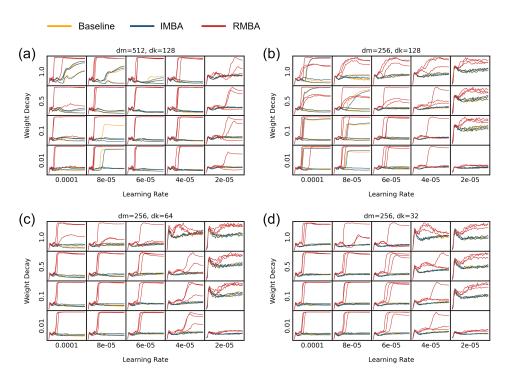


Figure 15: A comparison of the training results for RMBA, IMBA, and the Baseline model under different hyperparameters is presented. We investigated learning rates ranging from 2e-5 to 1e-4 and weight decay values from 0.01 to 1. We considered four configurations of the hidden space dimension. For each hyperparameter setting, we conducted experiments using 3 different random seeds, totaling 720 experiments.

Dataset	Example	Result
Clutrr	Statement: Stella's husband, Albertus, surprised her with tickets to a football game for their anniversary. Albertus rushed to the hospital to find out that his wife had already given birth to a boy and had named him Pleasant. Frank told a secret to her sister, Blanche. Blanche passed it along to her brother, Pleasant. Pleasant took his Aunt Frank out for her favorite meal. Barnet is Frank's older brother. He has never liked any of her boyfriends. Blanche and her aunt, Frank, went to the deli. They got half a pound of corned beef and two pounds of salami. Gina asked her daughter, Frank, if she had fun at school that day. Frank answered that she and her sister, Frank, had lots of fun together, Albertus went to the game with his sister Frank. Albertus took his daughter Gertie to the park that afternoon to play. Pleasant's wife, Celestia, surprised him on his birthday. He couldn't believe she pulled it off. Florence and her son's wife, Celestia, flew first class to see the concert. Query: Blanche is who of Stella Answer: daughter	50% GFT2 MMS. 40% GFT3 MMS. 40
RuleTaker	Statement: Cow sees mouse. Cow likes tiger. Bear is cold. Cow is big. If X visits bald eagle and X is kind then X is nice. Query: Cow sees bear? Answer: False	57% OFT2 M66A
StepGame	Story: The object labeled Y is positioned to the left of the object labeled G. Z is at the bottom of B. B is above K with a small gap between them. N is at the 9 c/clock position relative to Z. F is to the left of K and is on the same horizontal plane. Y presents right to N. Question: What is the relation of the agent Y to the agent B? Label: Below	40% OFT2 OFT2 AMMA OFT2 OFT2 OFT2 OFT2 OFT2 OFT2 OFT2 OFT2
LogicInference	Fact: Consider the following premises. exists x15: R15(x15) \rightarrow U1(x15), forall x15: Q15(x15) \rightarrow Q10(x15), forall x15: P15(x15) or R15(x15), forall x15: Q10(x15) \rightarrow U1(x15). Q10(x15) \rightarrow U1(x15). Query: Can we infer exists x15: U1(x15) and Q(x15) from them? Answer: yes	75% GOT7 MMA. 175% GOT9 MMA.
LogicAsker	Statement: For all x12, x12 will go running. For all x12, x12 is a police officer. There is at least one x12 for which if x12 were a scientist, then x12 is not a police officer. Query: Can we infer the following from them? Answer yes or no: There is at least one x12 for which x12 is not a scientist Answer: yes	100% 60% 0 1 2 (poch 3 4 0772.8Mpt)
PararulePlus	Fact: The wolf is tired. The wolf is dull. The wolf is rough. The wolf needs the dog. The bear sees the rabbit. The bear is fierce. The bear is awful. The dog is kind. The dog is smart. The dog is round. The rabbit is cute. The rabbit is lovely. The rabbit is furry. Kind animals are cute. If something is dull then it visits the dog. If something visits the dog then it is slow. If something is tired and dull then it is rough. If something is rough then it is lazy, All lazy animals are sleepy. If something is cute then it is lovely. All lovely animals are furry. If something is done then it is beautiful. All beautiful animals are small. All slow animals are big. Query: The bear is not heavy Answer: false	1100% 950% 950% 950% 950% 950% 950% 950% 9
PrOntoQA	Fact: Every sterpus is transparent. Sterpuses are brimpuses. Every sterpus is a lorpus. Brimpuses are melodic. Brimpuses are zumpuses. Each brimpus is a shumpus. Every shumpus is not fruity. Lorpuses are moderate. Gorpuses are metallic. Vumpuses are not melodic. Gorpuses are lempuses. Fae is a gorpus. Fae is a sterpus. Query: Fae is not melodic. Answer: false	100% 90% 100% 1

Figure 16: Examples from each dataset, along with their training curves.

F Details for Horizontal Thinking Strategy

We trained a 2-layer Transformer model and investigated its ability to perform multi-step reasoning with horizontal thinking. Specifically, we trained the Transformer model on 20,000 samples of length 13, each labeled with the result of a one-step reasoning process.

As shown in Fig. 17, when we fed the model's output back into the model, it was able to generate the next step's reasoning result, even though it had never been exposed to sentences longer than length 13 during training.

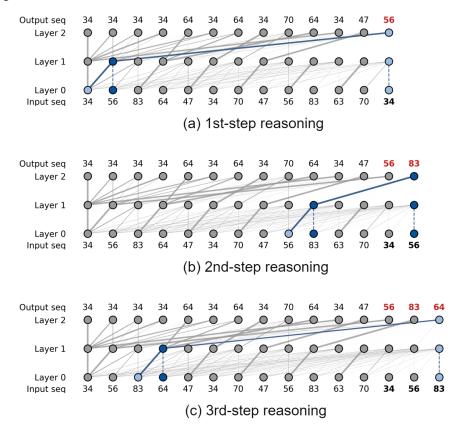


Figure 17: An illustration of the process of performing 3-step reasoning using a 2-layer Transformer model with CoT. The width of the connections in the diagram is based on the attention weights.

G Matching Score and Kernel Score for Real World LLMs

In this section, we calculate the matching score and kernel score of the large language model Phi-3(Abdin et al., 2024).

We focus on whether \boldsymbol{W}^{vo} and \boldsymbol{W}^{qk} function as the information buffer and the information extractor, respectively. To simplify our analysis, we temporarily disregard the effects of the feedforward layers. Following the method described in the main text, we compute $\operatorname{Ker}^{(l_1,l_2)} = \boldsymbol{W}^{qk(l_1)} \boldsymbol{W}^{vo(l_2),\mathsf{T}}$ and observe whether it exhibits a dominant diagonal characteristic. For multi-head models, the above equation is modified as:

$$\operatorname{Ker}^{(l_1, l_2)} = \left(\sum_{h} \boldsymbol{W}^{q(l_1, h)} \boldsymbol{W}^{k(l_1, h), \mathsf{T}}\right) \boldsymbol{W}^{vo(l_2), \mathsf{T}}.$$
(14)

We define the Kernel Score (KS) as

$$KS(Ker^{(l_1,l_2)}) = Trace(\sigma(Ker^{(l_1,l_2)}))/d_m, \tag{15}$$

which measures the ability of layer l_1 in the model to extract information cached at layer l_2 . As shown in Fig. 18(a), when $l_1 \geq l_2$, the kernel score is nearly zero, which aligns with the logical sequence of information processing. When $l_1 < l_2$, the kernel score decreases as $(l_2 - l_1)$ increases, indicating that the model tends to extract the most recently acquired information for further processing. In Fig. 18(b), we plot $\sigma(\text{Ker}^{(l_1,l_2)})$ and highlight the regions where the Kernel Score > 0.3.

To further verify that the diagonal structure of $\operatorname{Ker}^{(l_1,l_2)}$ arises from the alignment of model weights rather than the intrinsic diagonal structure of $\sum_h W^{q(l_1,h)} W^{k(l_1,h),\mathsf{T}}$ and $W^{vo(l_2),\mathsf{T}}$, we conducted a small experiment, as illustrated in Fig. 18(d). In this experiment, we assume that both matrices A and B are noise-added identity matrices, where the noise scale is denoted by α . We then compute the kernel scores of A, B, and their product C = AB. The results show that when $\max(\mathsf{KS}(A),\mathsf{KS}(B)) < 0.3$, the Kernel score of the product, $\mathsf{KS}(C)$, is less than 0.025. However, as shown in Fig. 18(c), we observe that for many heads, even when $\max(\mathsf{KS}(\sum_h W^{q(h)} W^{k(h),\mathsf{T}}),\mathsf{KS}(W^{vo,\mathsf{T}})) < 0.3$, the $\mathsf{Ker}^{(l_1,l_2)}$ still has a large kernel score. This indicates that the alignment of model weights is the key factor driving the diagonal structure of $\mathsf{Ker}^{(l_1,l_2)}$. In Fig. 18(f), we present an example where both $\sum_h W^{q(l_1,h)} W^{k(l_1,h),\mathsf{T}}$ and $W^{vo(l_2),\mathsf{T}}$ appear relatively disordered individually, but their product exhibits a clear diagonal structure.

Another straightforward method is to directly set the diagonal elements of $\sum_h \boldsymbol{W}^{q(l_1,h)} \boldsymbol{W}^{k(l_1,h),\mathsf{T}}$ and $\boldsymbol{W}^{vo(l_2),\mathsf{T}}$ to zero, and then calculate the kernel score based on the resulting $\tilde{\mathrm{Ker}}^{(l_1,l_2)}$. Fig. 18(e) illustrates this result, showing that it is approximately the same as that in Fig. 18(a).

Moreover, without loss of generality, we consider the case that includes LayerNorm(LN) and feed-forward(FNN) layers. We compute the matching score for each head in each layer, with the specific calculation formula as follows:

$$\boldsymbol{X}^{vo} = LN_{\text{attn}}^{(l-1)}(\boldsymbol{X})\boldsymbol{W}^{v(l-1),\mathsf{T}}\boldsymbol{W}^{o(l-1),\mathsf{T}}, \tag{16}$$

$$\boldsymbol{X}^{vof} = \boldsymbol{X}^{vo} + \text{FNN}^{(l-1)}(\text{LN}_{\text{FNN}}^{(l-1)}(\boldsymbol{X}^{vo})), \tag{17}$$

$$\boldsymbol{X}^{vok(h)} = LN_{\text{attn}}^{(l)}(\boldsymbol{X}^{vof})\boldsymbol{W}^{k(l,h)}, \tag{18}$$

$$\mathbf{X}^{f} = LN_{\text{atm}}^{(l-1)}(\mathbf{X}) + FNN^{(l-1)}(LN_{FNN}^{(l-1)}(LN_{\text{atm}}^{(l-1)}(\mathbf{X}))),$$
(19)

$$\boldsymbol{X}^{q(h)} = \boldsymbol{X}^f \boldsymbol{W}^{q(l,h)}, \tag{20}$$

matching matrix:
$$h^{(l,h)}(X) = X^{q(h)}X^{vok(h),\mathsf{T}}$$
. (21)

We visualized the matching score of each head in each layer (Fig. 18(b)) and found that the matching scores were highest in layers 5 to 20. This aligns with the conclusion mentioned in (Dutta et al., 2024), namely that the reasoning layers of large language models generally appear in the middle portion.

To further validate that Phi-3 might employ a buffer mechanism, we computed the pairwise cosine similarity of the matrices $\{W^{vo(l)}\}$ (each matrix is flattened as a long vector). The results indicate that these matrices are nearly orthogonal to each other, suggesting that they can be treated as independent buffers.

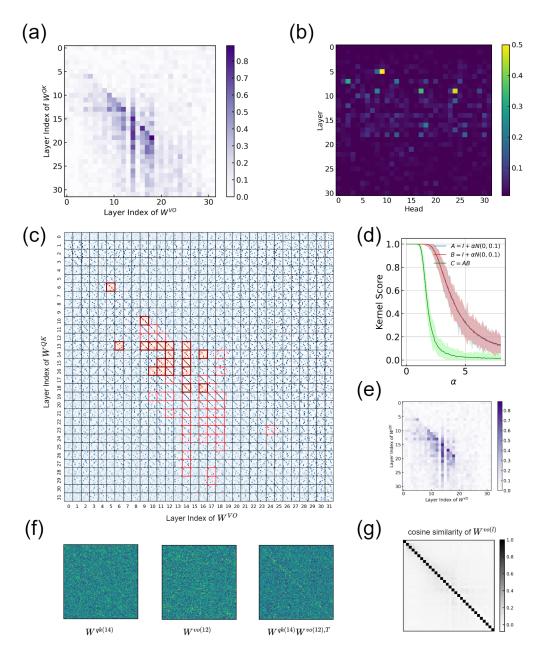


Figure 18: The calculation results of the kernel score and matching score for the Phi-3 model. (a) Visualization of $\sigma(\mathrm{Ker}^{(l_1,l_2)})$. (b) Visualization of the matching score calculations for each head in each layer, indicating that reasoning is concentrated in the middle layers of the model. (c) Visualization of the kernel matrix between layers, where the subgraphs enclosed in red boxes correspond to $\mathrm{KS}(\mathrm{Ker}) > 0.3$ and the subgraphs enclosed in darkred boxes correspond to $\mathrm{KS}(\mathrm{Ker}) > 0.3$ but $\mathrm{max}(\mathrm{KS}(\sum_h \boldsymbol{W}^{q(h)}\boldsymbol{W}^{k(h),\mathsf{T}}),\mathrm{KS}(\boldsymbol{W}^{vo,\mathsf{T}})) < 0.3$. (d) The kernel score of a noise-added identity matrix(A and B) and the kernel score of the product of two noise-added identity matrices(C). It can be observed that for two unrelated matrices, when their individual kernel scores are less than 0.3, the kernel score of their product is less than 0.025. (e) The kernel score obtained after setting the diagonal elements of $\sum_h \boldsymbol{W}^{q(h)} \boldsymbol{W}^{k(h),\mathsf{T}}$ and $\boldsymbol{W}^{vo,\mathsf{T}}$ to zero. (f) Visualization of the $\sum_h \boldsymbol{W}^{q(14,h)} \boldsymbol{W}^{k(14,h),\mathsf{T}}$ and the $\boldsymbol{W}^{vo(l_2),\mathsf{T}}$ in the Phi-3 model, along with their inner product. Despite their weak diagonal structure individually, their inner product exhibits a clear diagonal structure. (g) Cosine similarity of flattened $\boldsymbol{W}^{vo(l_1)}$ and $\boldsymbol{W}^{vo(l_2)}$, $l_1, l_2 \in \{0, \cdots, 31\}$.

Except for the phi-3 large model, the above definitions of the multi-head matching score and kernel score are equally applicable to other models that do not employ group query attention. The corresponding experimental results are shown in Fig. 19.

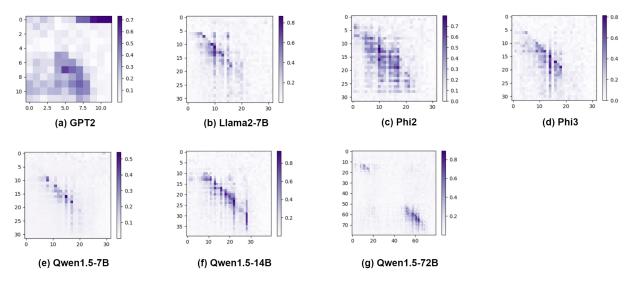


Figure 19: Weight alignment phenomenon for real-world LLMs.

H Causal Intervention

This section presents causal intervention experiments conducted to verify that the model uses a buffer mechanism when performing symbolic multi-step reasoning tasks. We assume readers are familiar with the content of Section 3. The causal intervention experiments were conducted using a 3-layer Transformer model trained as described in Section 5.

First, we identified critical tokens and logical circuits by observing output changes when masking specific attention or residual paths. Fig. 20(left) illustrates the logical circuit of the 3-layer Transformer performing the symbolic 2-step reasoning task. Subsequently, we conducted causal intervention experiments on the information stored in each token.

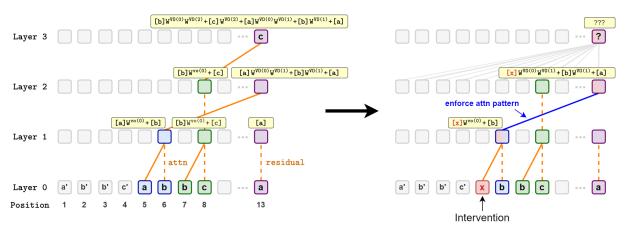


Figure 20: Logical circuit for 2-step reasoning (left) and illustration of causal intervention(right). We achieve the goal of intervening in information stored in a specific buffer by modifying the input sequence and enforcing the same attention pattern as before. The figure illustrates how the information [a] stored in the final token's buffer $\boldsymbol{W}^{vo(0)}\boldsymbol{W}^{vo(1)}$ in layer 2 can be changed to [x].

Unlike prior works where causal intervention replaced all information within a token with alternative information (Feng and Steinhardt, 2023; Meng et al., 2022; Vig et al., 2020; Wang et al., 2024a), we refined the perturbation scope to target a specific buffer within a token. Specifically, we individually replaced the information in each buffer of critical tokens with alternative information to observe the model's output. For example, as shown in Fig. 20(right), suppose we want to change the information [a] stored in the final token's buffer $\boldsymbol{W}^{vo(0)}\boldsymbol{W}^{vo(1)}$ in layer 2 to [x]. This can be achieved by simply modifying the input sentence [a'][b'][c'][a][b][c]...[a] to [a'][b'][c'][c'][x][b][c]...[a], and enforcing Attn[13:] same as before. Fig. 21 shows the intervention results for all buffers of all critical tokens.

The results reveal that, in layer l, the information stored in buffer $W^{vo(l-1)}$ of the final token is crucial. Modifying information in other buffers does not affect the model's output. Combined with the observation in Appendix D that the cosine similarities between $\boldsymbol{W}^{vo(0)}, \boldsymbol{W}^{vo(1)}, \boldsymbol{W}^{vo(2)}, \boldsymbol{I}$ are nearly zero, we can confidently assert that the model performs reasoning by leveraging different buffers. This experiment also rules out the possibility of the model only using an overwrite mechanism to perform reasoning.

Position	Buffer	Intervened Info	Method	Output
			Input [a'][b'][b'][c'][a][b][b][c][a]	[c]
Layer1 Token5	$\mathbf{W}^{\mathrm{vo}(0)}$	$[a] \rightarrow [x]$	Input [a'][b'][b'][c'][x][b][b][c][a]	Random
Layer1 Token5	I	$[b] \rightarrow [b']$	Input [a'][b'][b'][c'][a <mark>][b']</mark> [b][c][a]	[c']
Layer1 Token8	$\mathbf{W}^{\mathrm{vo}(0)}$	$[b] \rightarrow [b']$	Input [a'][b'][b'][c'][a][b] <mark>[b'][</mark> c][a]	Random
Layer1 Token8	I	$[c] \rightarrow [x]$	Input [a'][b'][b'][c'][a][b][b] <mark>[x]</mark> [a]	[x]
Layer2 Token8	$W^{vo(0)}$	$[b] \rightarrow [b']$	Input [a'][b'][b'][c'][a][b] <mark>[b'][</mark> c][a]	Random
Layer2 Token8	I	$[c] \rightarrow [x]$	Input [a'][b'][b'][c'][a][b][b] <mark>[x]</mark> [a]	[x]
Layer2 Token13	$W^{vo(0)}W^{vo(1)}\\$	$[a] \rightarrow [x]$	Input [a'][b'][b'][c'][x][b][b][c][a], enforce Attn ⁽¹⁾ [13:] same as before	[c]
Layer2 Token13	$\mathbf{W}^{\mathrm{vo}(1)}$	$[b] \rightarrow [b']$	Input [a'][b'][b'][c'][a][b'][b][c][a]	[c']
Layer2 Token13	$\mathbf{W}^{\mathrm{vo}(1)}$	$[b] \rightarrow [x]$	Input [a'][b'][b'][c'][a][x][b][c][a]	Random
Layer2 Token13	I	$[a] \rightarrow [x]$	Input $[a'][b'][b'][c'][a][b][b][c][x]$, enforce Attn ⁽¹⁾ [13:] same as before	[c]
Layer3 Token13	$W^{vo(0)}W^{vo(2)}\\$	$[b] \rightarrow [x]$	Input [a'][b'][b'][c'][a][b][x][c][a], enforce Attn ⁽²⁾ [13:] same as before	[c]
Layer3 Token13	$\mathbf{W}^{\mathrm{vo}(2)}$	$[c] \rightarrow [x]$	Input [a'][b'][b'][c'][a][b][b][x][a]	[x]
Layer3 Token13	$W^{vo(0)}W^{vo(1)}$	$[a] \rightarrow [x]$	Input [a'][b'][b'][c'][x][b][b][c][a], enforce Attn ⁽¹⁾ [13:] same as before	[c]
Layer3 Token13	$\mathbf{W}^{\mathrm{vo}(1)}$	$[b] \rightarrow [x]$	Input [a'][b'][b'][c'][a][x][b][c][a], enforce Attn ⁽²⁾ [13:] same as before	[c]
Layer3 Token13	I	$[a] \rightarrow [x]$	Input [a'][b'][b'][c'][a][b][b][c][x], enforce Attn ⁽¹⁾ [13:] same as before	[c]

Figure 21: Causal Intervention Experiment. We performed interventions on the information stored in every buffer of every critical token individually. Here, [x] represents a token that does not appear in the original input. For the final token, only modifying the buffer $\mathbf{W}^{vo(l-1)}$ in layer l affects the final output. In this experiment, the tokens in the original sentence are selected from the range [20, 40], while token [x] traverses the range [40, 100]. Attn $^{(l)}_{[13:]}$ refers to the attention score corresponding to the last token at layer l. For the original input, we have $\mathrm{Attn}^{(1)}_{[13,6]}=1$ and $\mathrm{Attn}^{(2)}_{[13,8]}=1$. Instances labeled as "Random" indicate that the output varies erratically as [x] changes. In all other cases, the probability of the model output deviating from the value presented in the table is less than 1e-15.

I Interaction Results with Large Language Models

4-step reasoning (Correct: [r])		3-step reasoning (Correct: [i])				2-step reasoning (Correct: [e])					
ChatGPT4	ChatGPT4o	Claude3.5- haiku- 20241022	Clauld3.5- sonnet- 20241022	ChatGPT4	ChatGPT4o	Claude3.5- haiku- 20241022	Clauld3.5- sonnet- 20241022	ChatGPT4	ChatGPT4o	Claude3.5- haiku- 20241022	Clauld3.5 sonnet- 20241022
a	i	a	n	a	a	a	e	a	e	e	e
a	i	a	e	a	n	a	a	e	a	a	e
a	a	a	e	e	e	a	a	a	n	a	e
a	a	a	a	a	a	e	n	e	a	e	e
a	r	r	e	a	a	r	e	e	e	a	e
a	a	a	a	e	i	e	r	e	e	e	e
a	a	a	i	a	e	a	a	a	e	e	e
a	a	a	W	a	a	a	a	e	a	a	e
0	a	a	w	a	a	a	a	e	a	a	e
a	a	a	r	a	o	a	e	e	p	a	e
a	a	a	e	a	a	a	e	a	e	a	e
a	a	a	i	a	a	a	a	a	a	e	e
a	a	a	e	a	a	e	a	a	a	a	e
a	i	n	w	a	e	a	i	a	e	a	e
a	a	a	o	a	a	a	e	e	a	a	e
a	a	a	r	a	a	e	a	e	a	a	e
a	a	a	a	e	e	a	a	a	e	a	e
a	a	a	e	e	e	a	a	a	e	a	e
a	a	e	e	a	n	a	a	a	a	e	e
a	a	a	a	a	o	a	n	e	a	a	e
a	a	n	e	a	a	a	a	e	e	a	e
a	a	a	r	a	a	a	n	a	a	a	e
a	a	a	a	a	a	a	a	e	a	a	e
a	a	a	a	e	a	a	a	a	e	a	e
a	e	a	e	a	n	a	a	e	e	a	e

Figure 22: Detailed interaction results with large models. For each type of reasoning task, we tested each large model 25 times. The versions of Claude are Claude 3.5-haiku-20241022 and Claude 3.5-sonnet-20241022.