POSITION-AWARE MODELING FOR NEXT-TOKEN PREDICTION

Anonymous authorsPaper under double-blind review

000

001

003

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

034

036

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Next-token prediction (NTP) serves as the dominant training paradigm for large language models (LLMs), enabling strong autoregressive (AR) generation capabilities. Despite its success, models trained with vanilla NTP often exhibit counterintuitive failure patterns, such as the reversal curse, factorization curse, and sensitivity to knowledge position. These failures stem from the fixed leftto-right token order during teacher-forcing supervision, which entangle content and token order in ways that compromise permutation invariance. To address these failures, we introduce a position-aware training framework that enables AR models to predict the next token not just based on seen content, but also to account for predicted token position. This disentanglement of what to predict and where to predict improves the robustness of LLMs to different token orderings. We instantiate this framework via two complementary approaches: (1) Content-Position Coupling (CPC), which injects a lightweight position-aware embedding into the input sequence without modifying the model architecture; and (2) Content-Position Decoupling (CPD), which introduces the modular positionaware blocks for the pre-training AR model to provide explicit supervision over target positions. Experiments across three representative tasks demonstrate that our framework consistently improves performance over strong baselines, while maintaining architectural simplicity and convergence efficiency. Codes are available at https://anonymous.4open.science/r/CPC-CPD.

1 Introduction

Next-token prediction (NTP) is the primary pre-training objective for large language models (LLMs) (OpenAI, 2023; Touvron et al., 2023a). LLMs can effectively acquire co-occurrence patterns among tokens by optimizing the autoregressive (AR) maximum likelihood estimation objective on large text corpora (Zhang et al., 2024b), so allowing the broad transfer of learned knowledge to many applications, ranging from text generation to complicated question answering and reasoning (Petroni et al., 2019; Hendrycks et al., 2020). NTP commonly integrates a teacher forcing mechanism (Williams & Zipser, 1989) during the training phase and employs AR at inference-time (Bachmann & Nagarajan, 2024). Owing to its significant advantages-notably in training efficiency (Gloeckle et al., 2024; Li et al., 2024), gradient stability (Chen et al., 2024), and amenability to parallel computation (Li et al., 2021; Rasley et al., 2020), NTP has established itself as a cornerstone in the pre-training of mainstream LLMs (OpenAI, 2023; Touvron et al., 2023a; Liu et al., 2024a; Jiang et al., 2024a; Bai et al., 2023).

Despite its long list of achievements, existing research has discovered that models trained via vanilla NTP can surprisingly exhibit counterintuitive failure patterns (Berglund et al., 2024; Lin et al., 2024; Lv et al., 2024; Bachmann & Nagarajan, 2024; Kitouni et al., 2024; Allen-Zhu & Li, 2024; Saito et al., 2025). For instance, they may suffer from (1) the **reversal curse** (Berglund et al., 2024; Lin et al., 2024; Lv et al., 2024), where learned associations (*e.g.*, "A is B") fail to generalize to their inverse form (*e.g.*, "B is A"); (2) the **factorization curse** (Kitouni et al., 2024), which arises when the model, trained on a specific decomposition of the token sequence (*e.g.*, left-to-right), fails to represent the same joint distribution under alternative factorizations; and (3) the **knowledge position sensitivity** (Allen-Zhu & Li, 2024; Saito et al., 2025), where factual information encoded during training is only reliably accessible when it appears in early positions of the training document, while knowledge located later is often unrecoverable during inference, even with accurate prompting. These

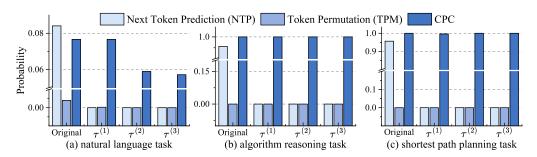


Figure 1: Joint probability across different permutations on the same sample under three task types. Our method maintains nearly consistent joint probability across different permutations, while both NTP and TPM fail to achieve probability invariance. $\tau^{(i)}$ denotes a specific permuted token order. For more detailed experimental settings and more examples, see Appendix D.1.

failure patterns reveal a shared deficiency: the left-to-right, teacher-forcing training paradigm inherent to NTP causes the model to rely heavily on preceding context during training, making model difficult to acquire a consistent joint distribution over tokens across different positions or permutations. Figure 1 illustrates the joint probability distribution of the same sample under different permutations across three distinct task types for fully trained NTP model. The results demonstrate that standard NTP methods only assign high probability to samples in the original order encountered during training, while probability for other permutation orders decreases to nearly zero, reflecting a severe probability distribution inconsistency issue.

Existing research to overcome these pitfalls has primarily followed two major directions. **Data-centric** strategies include data rewriting and token permutation (Golovneva et al., 2025; Guo et al., 2024) to encourage model learning under a broader range of token factorizations, and structural reorganization of training data to break the NTP inherent left-to-right generation bias, *e.g.*, by exposing the model to future tokens earlier in the sequence (Thankaraj et al., 2025). **Model-level** work targets model architecture and training objective modifications. This includes multi-token prediction (Gloeckle et al., 2024; Zhang et al., 2024a), altering the training objective to predict both the next token of a prefix and the previous token of a suffix (Hu et al., 2024), and equipping AR models with bidirectional attention mechanisms to better capture global contextual dependencies (Lv et al., 2024).

However, these are two primary challenges: (1) For data-centric methods, data rewriting typically relies on advanced LLMs (*e.g.*, GPT-4), which inevitably introduces hallucinations. Moreover, for token permutation under the vanilla NTP objective, it often leads to **different target tokens are associated with identical prefixes, which introduces label ambiguity and undermines training stability¹. As shown in Figure 1, similar to NTP, token permutation methods struggle to assign consistent probabilities across various permutations even after extensive training, and especially in planning and algorithm reasoning, they even underperform compared to standard NTP**. (2) For methods that involve modifying the model architecture or training objective, most require careful adaptation to specific model implementations, making them difficult to generalize across diverse architectures and limiting their practical applicability.

In this paper, we build upon token permutation and then expose the model to diverse token permutations, encouraging model to learn position-agnostic representations and improve its generalization across different factorizations. To address the inevitable issue of target ambiguity introduced by token permutations, where different ground-truth tokens are associated with the same prefix, we further propose to augment the vanilla NTP objective with position-aware modeling, explicitly encoding the location of the predicted token. Concretely, we introduce two complementary approaches: (1) **Minimal modification, data-centric coupling**: We retain the original AR architecture and bind the target location and the permeated sequence into the model. Specifically, we first introduce a single learnable position embedding that is rotated to arbitrary positions using rotary position embedding (RoPE) to obtain target position-aware embedding. Then, the obtained position vector is combined with permuted input sequences, enabling the model to implicitly recognize the specific position for next token prediction, thereby mitigating the problem of identical prefixes leading to different targets.

¹As for different target tokens are associated with identical prefixe, we provide a detailed example and explanation in Appendix E.1.

- (2) **Incremental, model-level decoupling**: While the first approach (coupling) doesn't require a lot of changes, it does have the downside of potentially messing with the original token embedding. To address this, we propose a structure-preserving decoupling method by augmenting AR model with auxiliary position-aware decoding blocks. These position-aware blocks fine-tune the next-token prediction by utilizing location as the query and the previous sequence content as key and value, effectively separating positional information from content representation. As shown in Figure 1, our method can maintain almost the same joint probability for different permutations. Crucially, our framework requires no changes to the base model's structure and can be plugged into any pre-trained AR model as a modular and learnable component, enabling scalable and position-aware adaptation. We summarize our contributions below.
 - We reveal that seemingly disparate failure patterns in LLMs, actually stem from the fundamental mechanism: the inherent sensitivity of NPT training paradigms to token ordering, which particularly impairs models' planning and reasoning capabilities.
 - We propose a complementary target position-aware enhancement framework, enabling models to distinguish prediction targets using both content and intended position.
 - Extensive experiments demonstrate that our proposed methods significantly enhance model robustness to token order variations, enabling smaller LMs to outperform larger backbone models.

2 RETHINKING FAILURE PATTERNS IN NTP

2.1 PRELIMINARIES

Consider a sequence s=(p,r), where $p=(p_1,p_2,\ldots,p_{|p|})$ denotes the prompt and $r=(r_1,r_2,\ldots,r_{|r|})$ denotes the response. Each token p_t and r_t is drawn from a fixed-size vocabulary $\mathcal V$. For any position i in the sequence s, let $s_{< i}$ denote the subsequence consisting of the first i-1 tokens and s_i denote the token at position i. Suppose we have a next-token prediction (NTP) language model P_θ parameterized by θ , such that $P_\theta(s_i \mid s_{< i})$ denotes the probability that the model assigns to the i-th token s_i , conditioned on the preceding sequence $s_{< i}$. For the given sequence s, the joint probability is axiomatically defined analogous to the chain rule of probability:

$$P_{\theta}(\boldsymbol{r}; \boldsymbol{p}) = \prod_{i=1}^{|\boldsymbol{r}|} P_{\theta}(r_i \mid \boldsymbol{p}, \boldsymbol{r}_{< i})$$
(1)

Training-time next-token prediction via teacher-forcing To train the above NTP model, main-stream LLMs adopt teacher forcing to maximize the log-probability sum of the next token, where the model is trained to predict each token r_t using the ground-truth $r_{< t}$ as input. The teacher-forcing objective $\mathcal{J}_{\text{teacher-forcing}}(\theta)$ on dataset \mathcal{D} can be formulated as follows:

$$\mathcal{J}_{\text{teacher-forcing}}(\theta) = \mathbb{E}_{(\boldsymbol{p}, \boldsymbol{r}) \sim \mathcal{D}} \left[\log P_{\theta}(\boldsymbol{r} \mid \boldsymbol{p}) \right] = \mathbb{E}_{\mathcal{D}} \left[\sum_{i=1}^{|\boldsymbol{r}|} \log P_{\theta} \left(r_i \mid \boldsymbol{p}, \boldsymbol{r}_{< i} \right) \right]$$
(2)

Inference-time next-token prediction via autoregression During inference, the model is conditioned on a given prompt p and generates output tokens \hat{r}_t by sequentially sampling from the learned distribution P_{θ} . Specifically, for each step t, the model samples a token $\hat{r}_t \sim P_{\theta}(\cdot \mid s_{< t}, \hat{r}_{< t})$, where $\hat{r}_{< t}$ signifies the previously generated tokens. The context is then supplied back into the model for the next prediction using the sampled token \hat{r}_t attached to it. A full sequence is formed by this autoregressive generation process continuing for |r| steps.

2.2 MITIGATING FAILURE PATTERNS IN NTP

We argue that the commonly observed failure patterns in NTP, namely the reversal curse, factorization curse, and knowledge position sensitivity, reflect a shared underlying limitation in vanilla NTP: the lack of robustness to variations in token permutation. Building on the insight by Kitouni et al. (2024) that consistency across token factorizations improves knowledge retrieval, we generalize this goal to a broader permutation-robustness perspective. Specifically, we posit that ensuring approximate

invariance of joint probabilities under token permutations not only mitigates reversal-related failures, but also serves as a foundation for addressing a wider set of token order sensitivities, such as position-induced misalignments in reasoning and planning.

Let $\tau^{(i,n)} \in S_n$ be the i_{th} sampled permutations, where S_n is the set of all n! permutation of the indices $\{1,2,\ldots,n\}$. Thus, $\tau^{(i,n)} = \{\tau_1^{(i,n)},\tau_2^{(i,n)},\ldots,\tau_n^{(i,n)}\}$. Applying permutation $\tau^{(i,|\boldsymbol{p}|)}$ and $\tau^{(i,|\boldsymbol{r}|)}$ reorders the sequence tokens accordingly, yielding permuted prefixes $\boldsymbol{p}_{\tau^{(i,|\boldsymbol{p}|)}}$ and responses $\boldsymbol{r}_{\tau^{(i,|\boldsymbol{r}|)}}$. Then, for any two sampled permutations $\tau^{(i,|\boldsymbol{p}|)} \in S_{|\boldsymbol{p}|}, \tau^{(i,|\boldsymbol{r}|)} \in S_{|\boldsymbol{r}|}$, the permutation-invariant objective expect the model P_θ could assign approximately joint probability, regardless of the permutation applied to the input. With an abuse of notation, let $\boldsymbol{p}_{\tau^{(\boldsymbol{p})}}$ and $\boldsymbol{r}_{\tau^{(\boldsymbol{r})}}$ denote the permutation of prompt and response, respectively. This objective can be formulated as:

$$\prod_{t=1}^{|r|} p_{\theta} \left(r_{\tau_{t}^{(r)}} \mid \boldsymbol{p}_{\tau^{(p)}}, \boldsymbol{r}_{<\tau_{t}^{(r)}} \right) \approx \prod_{t=1}^{|r|} p_{\theta} \left(r_{\tau_{t}^{(2)}} \mid \boldsymbol{p}_{\tau^{(1)}}, \boldsymbol{r}_{<\tau_{t}^{(2)}} \right)$$
(3)

where $\tau^{(1)}$ and $\tau^{(2)}$ denote the used permutation of prompt and response (e.g., left-to-right), respectively, during training.

To approximate the permutation-invariant objective in Equation 3, the straightforward strategy is to apply diverse permutations to the training data. However, such permutations inherently disrupt the natural structure of language. On one hand, exposing the model to arbitrarily disordered sequences may harm its ability to model syntactic or semantic coherence. On the other hand, permutation introduces a fundamental conflict: under the same prefix, the model is required to optimize for different next-token targets, which results in inconsistent supervision signals. Moreover, prior studies (Kitouni et al., 2024) have shown that the masked language modeling (MLM) objective is effective in alleviating the reversal curse and factorization curse, but it not be included by the existing LLMs pre-training paradigms. Instead, it randomly masks tokens at arbitrary positions and predicts them using bidirectional context, allowing model to learn representations that are inherently robust to variations in token order. However, applying MLM-style training to pre-trained AR language models typically requires modifying the internal attention mechanism (Lv et al., 2024) or full model re-training, which either conflicts with the original AR pre-training objective or incurs substantial computational cost. To achieve the permutation-invariant objective within the pre-trained AR models, it is desirable to combine the AR structure of NTP with the positional flexibility of MLM, i.e., allowing the model to predict tokens at arbitrary positions rather than adhering to a strict left-to-right order. This requires explicitly identifying which token is to be predicted under each context.

3 METHODOLOGY

Considering the ambiguity and inconsistency brought by token permutations in NTP, especially when the same prefix corresponds to different target tokens, we propose a target position-aware training framework.

This framework can clearly and completely distinguish prediction targets not only based on content, but also based on their intended positions in the permuted sequence. We instantiate this framework in two complementary ways: (1) **Content-Position Coupling** (CPC), which implicitly informs the model of the target position by injecting a lightweight position embedding into the input sequence. CPC requires no modification to the model architecture and minimally intervenes with the original AR behavior. (2) **Content-Position Decoupling** (CPD), which introduces a modular position-aware module on top of the pre-trained AR model. This auxiliary module decouples content modeling and target position, allowing for explicit position-aware supervision.

3.1 TARGET POSITION-AWARE EMBEDDING

In widely used AR model such as GPT and Llama, the NTP is conditioned not only on the content of the previous tokens, but also on their positions. Specifically, for target position i, the model predicts²:

$$p_{\theta}(s_{\tau_i} \mid \boldsymbol{s}_{<\tau_i}) = p_{\theta}\left(s_{\tau_i} \mid \left\{\phi(\text{Embed}(s_{\tau_j}), \text{Pos_Embed}(\tau_j))\right\}_{j < i}\right) \tag{4}$$

²Without loss of generality, we do not restrict NTP to scenarios where generation is prompted by a prefix, as the model can also be trained directly on the prompt tokens.

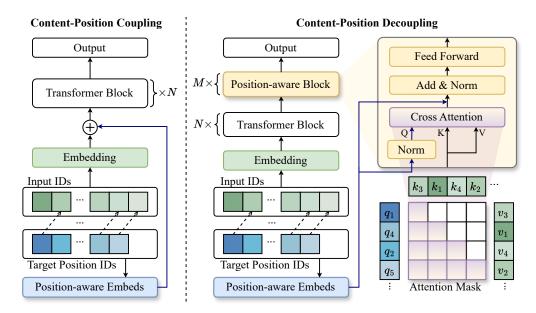


Figure 2: Overview of the proposed target position-aware framework, illustrating the Content-Position Coupling (CPC) (left) and Content-Position Decoupling (CPD) (right) approaches.

where $\operatorname{Embed}(s_j)$ denotes the embedding of the content s_j and the position encoding $\operatorname{Pos_Embed}(j)$ can be either the absolute positional encoding or the relative positional encoding method. $\phi(\cdot)$ is a token-position fusion function. Under absolute positional encoding schemes, $\phi(\cdot)$ is typically implemented as an element-wise addition to the token at the embedding layer. In contrast, relative positional encoding mechanisms such as rotary positional encoding (RoPE) integrate $\phi(\cdot)$ directly into the attention mechanism by rotating query and key vectors based on current token positions.

To inject the positional information of target tokens into AR models while preserving the original flow, we design a hybrid positional encoding scheme that combines the flexibility of absolute position embeddings with the extensibility advantages of RoPE. Specifically, instead of assigning each position a fixed embedding vector from the pre-defined embeddings, we first learn a shared position base vector $e_{\text{pos}} \in \mathbb{R}^{1 \times dim}$, where dim is embedding dimension, and then rotate it to the desired target positions according to their corresponding position ids using RoPE-1D. The position-aware embedding z for the target token at position τ_t is:

$$\boldsymbol{z}_{\tau_t} = \text{RoPE-1D}(\boldsymbol{e}_{\text{pos}}, \tau_t) \tag{5}$$

This design allows us to learn a single shared parameter while obtaining theoretically unbounded target position-aware embeddings.

3.2 Content-Position Coupling

To minimize architectural modifications, we propose a content-position coupling training strategy that encodes target-position information into the input representation. Specifically, we modify the original token-position fusion in Equ. 4 by integrating the target position-aware embedding z_i to the input sequence. This allows the model to internalize the intended prediction position without altering its architecture or decoding behavior, which can be formulated as follows:

$$p_{\theta}(s_{\tau_i} \mid \boldsymbol{s}_{<\tau_i}) = p_{\theta} \left(s_{\tau_i} \mid \left\{ \phi(\text{Embed}(s_{\tau_j}) \oplus \boldsymbol{z}_{\tau_i}, \text{Pos_Embed}(\tau_j)) \right\}_{j < i} \right)$$
(6)

where \oplus denotes the interaction operation between content and target position. The interaction operation can be instantiated using either parametric or non-parametric methods, such as direct addition, concatenation followed by a linear projection, or other fusion strategies, achieving internalization of the target position-awareness into the embedding of the model. For simplicity of design, we directly use direct addition as the default setting in \oplus . We provide the training pseudo-code and concrete example for CPC in Algorithm C1 and Figure C1, respectively.

While CPC is simple and efficient, requiring only minimal architectural modifications, its direct coupling of content and positional information introduces a potential drawback. By injecting position-

awareness through input embedding modification, it may cause the semantic representations learned during pre-training to drift, thereby creating a gap between pre-training and fine-tuning that potentially undermines the original performance of LLMs. This motivates us to explore another way to preserve a clear separation between content and position.

3.3 CONTENT-POSITION DECOUPLING

To achieve a decoupling between content and position, a straightforward approach is to reformulate the CPC objective such that the target position is treated as an independent conditioning signal, rather than being directly fused with input embeddings (as in CPC).

Reformulation Similarly, we still adhere to the principle of not significantly adjusting the original AR structure. Therefore, we reformulate the objective in a decoupled manner as follows:

$$p_{\theta}(s_{\tau_i} \mid \boldsymbol{s}_{<\tau_i}) = p_{\theta}\left(s_{\tau_i} \mid \left\{\varphi\left(\phi(\text{Embed}(s_{\tau_j}), \text{Pos_Embed}(\tau_j)), \boldsymbol{z}_{\tau_i}\right)\right\}_{j < i}\right) \tag{7}$$

Here, $\varphi(\cdot, z_i)$ represents an auxiliary position-aware conditioning function that modulates the decoding process using an independent embedding z_{τ_i} associated with the intended prediction position.

From Equ. 7, we observe that the key challenge now lies in how to design an effective mechanism (i.e., $\varphi(\cdot, z_{\tau_i})$) for incorporating z_{τ_i} into the model's workflow without disrupting the basic AR structure. In the following, we introduce a position-aware decoder branch that integrates this positional signal through cross-attention, enabling flexible and explicit supervision of target position.

Overview The overall structure of CPD is illustrated on the right side of Figure 2. We adopt an incremental and modular design that allows for integration with the existing AR-based models. Specifically, we insert M position-aware decoder block after the output of original model's transformer blocks, which performs cross-attention between the base transformer blocks' final hidden states and the target position-aware embedding z_{τ_i} , enabling the model to condition its predictions explicitly on the content and intended output position.

Position-aware Block Let $S = \text{BaseModel}(s_{\tau}) \in \mathbb{R}^{|s_{\tau}| \times dim}$ represent the hidden states of the base AR model's final layer. To decouple content and target position information, we design a cross-attention mechanism within the position-aware block, where the query comes from the target position-aware embedding \mathbf{z}_{τ_i} , and the key and value come from the content (input) sequence representations S. For input sequence indices $\tau = [\tau_1, \tau_2, \dots, \tau_{|s_{\tau}|}] \in \mathbb{R}^{1 \times |s_{\tau}|}$, the corresponding target position set is $\tau_T = [\tau_2, \dots, \tau_{|s_{\tau}|}] \in \mathbb{R}^{1 \times |s_{\tau}-1|}$. For any target position $\tau_i \in \tau_T$, the cross-attention mechanism can be formulated as follows:

$$Q = \text{RoPE-1D}(\boldsymbol{z}_{\tau_i} \boldsymbol{W}_q, \tau_i), \qquad \boldsymbol{K} = \text{RoPE-1D}(\boldsymbol{S} \boldsymbol{W}_k, \tau), \qquad \boldsymbol{V} = \boldsymbol{S} \boldsymbol{W}_v$$
 (8)

$$Att_{cross}(Q, K, V) = Softmax(QK^{T} + \mathcal{M})V$$
(9)

where z_{τ_i} is target position-aware embedding of $\tau_i, W_q, W_k, W_v \in \mathbb{R}^{dim \times dim}$ are learnable weights, and \mathcal{M} is attention mask. As shown in Figure 2, the attention mask $\mathcal{M} \in \{0, -\infty\}^{|\tau_T| \times |\tau|}$ ensures that each target position $\tau_i \in \tau_T$ is only allowed to attend to key-value pairs corresponding to its preceding content tokens $\tau_j < \tau_i$. The causal structure of AR decoding is maintained while enabling explicit conditioning on the target position through a query formed from z_{τ_i} . Each row of the attention mask activates only those positions in the input sequence that occur before the current prediction position, retaining left-to-right generation constraint in vanilla NTP.

Similar to standard decoder layers used in Transformer architectures, a single position-aware block can be formulated as follows:

$$\boldsymbol{H} = T + \text{FFN}(\text{LN}(T)) \tag{10}$$

$$T = LN(z_{\tau_s} + Att_{cross}(Q, K, V))$$
(11)

where H is the hidden state of a single position-aware block, LN is the layer-norm function, and FFN is the feed-forward network. We provide the training pseudo-code for CPD in Algorithm C2. It is worth noting that in the training optimization stage, CPC and CPD perform teacher-forced NTP (Equ. 2) on the basis of Equ. 6 and Equ. 7, respectively.

		N	ameIsDesc	ription				De	escriptionI	sName	
Method		N2D		D2	2N			N2D		D	2N
	EM	R-1	BLEU	EM	R-1		EM	R-1	BLEU	EM	R-1
					Llama-2	2-71	B-base				
NTP	<u>77.7</u>	91.5	93.2	0.00	0.00		0.00	19.9	25.4	91.7	91.7
TPM	47.7	84.1	86.1	99.7	99.7		17.3	78.0	82.3	98.7	98.9
SPT*	N/A	N/A	83.6	100.0	100.0		N/A	N/A	84.3	100.0	100.0
BICO	68.7	89.4	91.1	99.7	<u>99.7</u>		2.00	24.1	26.9	100.0	100.0
CPC	76.3	92.1	93.1	100.0	100.0		47.8	83.5	92.3	100.0	100.0
CPD	78.3	91.9	94.4	100.0	100.0		48.3	85.7	93.6	100.0	100.0
					Llama-3	3-81	B-base				
NTP	73.3	91.8	94.5	0.0	0.0		0.0	17.1	24.0	99.7	99.7
TPM	56.3	82.6	87.3	<u>94.6</u>	94.6		24.8	83.9	85.1	100.0	100.0
BICO	63.7	87.6	91.3	92.3	92.3		0.0	18.1	24.8	100.0	100.0
CPC	87.0	<u>95.6</u>	96.9	100.0	100.0		59.2	86.7	89.3	100.0	100.0
CPD	88.6	97.2	98.3	100.0	100.0		62.9	87.2	89.9	100.0	100.0
					Llama-3.	2-1	B-base				
NTP	75.0	76.9	79.3	0.00	0.00		0.00	2.9	7.7	91.7	91.7
TPM	46.7	85.2	86.5	95.7	95.7		22.3	80.7	84.7	97.3	97.3
BICO	60.3	74.5	77.8	37.0	37.3		0.0	19.2	23.8	<u>97.7</u>	<u>97.7</u>
CPC	<u>78.7</u>	91.8	92.8	82.7	83.6		32.8	82.9	89.7	100.0	$\overline{100.0}$
CPD	81.3	94.7	95.8	100.0	100.0		63.0	85.3	<u>87.7</u>	100.0	100.0

Table 1: Experimental results under the reversal curse setting across various Llama models. Results of method marked with * are from Guo et al. (2024).

4 EXPERIMENTS

We evaluate the performance of CPC on the following representative tasks: **reversal curse**, **factorization curse**, and **positional bias**.

4.1 REVERSAL CURSE

Settings *Datasets:* Following previous work (Berglund et al., 2024; Lin et al., 2024; Lv et al., 2024), we evaluate CPC and CPD on the name-description dataset (Berglund et al., 2024). Detailed descriptions and statistics of the datasets are provided in Appendix D.2.1. *Baselines:* NTP, Token Permutation (TPM), BICO (Lv et al., 2024), and SPT (Guo et al., 2024). We evaluate all methods on Llama-2-7B (Touvron et al., 2023b), Llama-3-8B (Grattafiori et al., 2024), and Llama-3.2-1B. Introduction and implementation details of all methods are provided in Appendix D.2.2 and D.2.4, respectively. *Evaluation Metrics:* We use exact match (EM), ROUGE-1 (R-1), and BLEU scores.

4.1.1 EXPERIMENTAL RESULTS

Table 1 reports experimental results under the reversal curse setting. We can draw the following conclusions: (1) On all metrics, CPC and CPD are significantly better than all baselines, suggesting that explicitly incorporating position information can effectively mitigate the problem of inconsistent information about the direction of the data during the training and testing phases. (2) Llama-3.2-1B+CPD (with 6 position-aware blocks adding 0.8B parameters) achieves results superior to larger-scale models, including Llama-2-7B and Llama-3-8B, and even surpasses Llama-3-8B+CPD in some ways. This demonstrates that we can endow smaller models with permutation-invariant capabilities by incorporating additional CPD modules. Meanwhile, we provide more experiments on the number of position-aware blocks in the Appendix E.6.1 and the effect of whether or not to train the base AR model on CPD performance in Appendix E7. Furthermore, in Appendix E.3, we analyze the trade-off between performance and training/inference burden introduced by adding extra CPD layers. Moreover, increasing additional parameters does not affect convergence speed. We find that CPD and CPC exhibit almost identical convergence behavior, both significantly superior to TPM, as shown in Figure E1. (3) While TPM can alleviate the reversal curse, it exhibits degraded performance on the N2D task of NameIsDescription compared to standard NTP. The primary reason is that altering the original token ordering during training tends to produce conflicting optimization objectives where identical prefixes map to different targets. As shown in Figure E4, this results in slow training optimization and unstable performance fluctuations. Furthermore, to assess whether permutation-based training affects the original performance of pretrained models, we evaluate model capabilities before and after training on nine standard NLP benchmarks. A detailed evaluation is provided in Appendix E.5. Moreover, since the reversal curse intuitively can benefit from bidirectional training, we also compared the classical bidirectional training model BERT in Appendix E.4.

Method	Path planning					Algorithmic reasoning					
Method	G(2,5)	G(5, 5)	G(20, 5)	G(2, 10)		scc-4	scc-5	scc-11	scc-12	scc-15	
NTP*	0.50	0.20	0.05	0.50		1.00	0.99	0.62	0.57	0.27	
TPM	0.00	0.00	0.00	0.00		1.00	0.53	0.00	0.00	0.00	
TRELAWNEY*	1.00	1.00	1.00	0.50		1.00	0.98	0.72	0.71	0.48	
CPC	1.00	1.00	1.00	0.99		1.00	1.00	0.97	0.99	0.84	
CPD	1.00	1.00	1.00	$\overline{1.00}$		1.00	1.00	$\overline{1.00}$	$\overline{1.00}$	0.93	

Table 2: Experimental results for path planning (star graph G(d, l) with d paths of length l from start node) and algorithmic reasoning (strongly connected components, denoted as scc -i where i represents connected graph size). Results of method marked with * are from from Thankaraj et al. (2025).

4.2 FACTORIZATION CURSE

Settings *Datasets:* Following prior work (Kitouni et al., 2024; Thankaraj et al., 2025), we experiment on the *Star Graph* dataset (Bachmann & Nagarajan, 2024) and the *strongly connected components algorithm* from CLRS-Text (Markeeva et al., 2024). Detailed introduction and statistics of the datasets are provided in Appendix D.3.1. *Baselines:* NTP, TPM, and TRELAWNEY (Thankaraj et al., 2025). Consistent with previous work (Thankaraj et al., 2025), we conduct experiments using Llama-3.2-1B, as models at the 1B scale typically lack task planning capabilities without fine-tuning. Introduction and implementation details of all methods are provided in Appendix D.3.2 and D.3.4, respectively. *Evaluation Metrics:* Accuracy is used to evaluate the performance of the model.

4.2.1 EXPERIMENTAL RESULTS

Star Graph Based on experimental results shown in Table 2, the following key conclusions can be drawn: (1) NTP struggles with path planning, especially as graph complexity increases. Its accuracy drops from 0.50 on G(2,5) to 0.05 on G(20,5), indicating difficulty in learning "difficult token" under teacher forcing. (2) TPM performs poorly, with near-zero accuracy across various star graphs. Permutations introduce conflicting prefix-target pairs, making optimization unstable, as also evidenced by its failure to converge (Figure E2). (3) Although TRELAWNEY achieves reasonable performance through data augmentation, it relies on carefully designed enhancement strategies, such as pre-planning which tokens the model should learn. Without designed prompting, its performance on the longer path planning task G(2,10) remains limited at 0.50. In contrast, our CPC and CPD methods consistently reach near-perfect accuracy (1.00), demonstrating the effectiveness of position-aware modeling in this path planning task.

Strong Connected Components As shown in Table 2, a similar trend is observed in the strongly connected components (SCC) benchmarks. NTP maintains high accuracy on scc-4 and scc-5 but collapses on larger connected graphs, dropping to 0.27 on scc-15. TPM completely fails beyond scc-5, with 0.00 accuracy on scc-11 through scc-15, revealing that permutation exposure without structural position grounding is insufficient for generalization. As shown in the Figure E3, it is also clear that during the training process, we find it difficult for the TPM to converge, which is further evidence of the supervised conflicting situations caused by the permutations. TRELAWNEY shows improved robustness, but its performance drops significantly on scc-15 (0.48). In contrast, CPC and CPD both maintain strong performance across all scales. CPD achieves perfect accuracy (1.00) on scc-4 through scc-12 and still reaches 0.93 on scc-15, outperforming all baselines and demonstrating superior scalability and permutation robustness.

4.3 Positional Bias

Settings *Datasets:* Following previous work (Saito et al., 2025), we evaluate CPC and CPD in real-world collections of Wiki2023+ (Jiang et al., 2024b; Saito et al., 2025) that are new knowledge for Llama-2. See Appendix D.4.1 for details. *Baselines:* Next-token prediction (NTP), Sentence Shuffle (SS), Attn Drop (AD), and D-AR (Saito et al., 2025). Details are in Appendix D.4.2. *Evaluation Metrics:* We adopt Exact Match (EM) and F1.

4.3.1 EXPERIMENTAL RESULTS

Table 3 shows the performance of CPC and CPD on the Wiki2023+ dataset of the movie domain collected in the real world, and we can draw the following conclusions: (1) CPC and CPD can be effectively applied to learn new knowledge in realistic scenarios, enabling the model to perceive knowledge distributed in different locations in a balanced manner. Specifically, compared to the best baseline method, D-AR, CPC achieves an average improvement of 10.0% in EM, while CPD

Method	←start	t				$end{\rightarrow}$	Average
	EM_1/Fl_1	EM_2/Fl_2	$EM_3 / F1_3$	$EM_4/F1_4$	$EM_5 / F1_5$	$EM_6 / F1_6$	Tiverage
NTP*	40.9 / 51.4	6.3 / 20.5	8.1 / 29.8	11.7 / 35.7	11.6 / 37.8	10.7 / 36.4	14.9 / 35.7
SS*	51.6 / 65.7	14.7 / 43.2	15.6 / 43.5	20.6 / 46.8	24.0 / 50.8	19.8 / 46.4	24.4 / 49.4
AD*	58.6 / 71.1	10.2 / 29.8	14.0 / 36.6	17.0 / 38.6	13.2 / 42.8	13.3 / 39.7	21.0 / 43.1
D-AR*	60.1 / 73.7	26.9 / 53.1	23.4 / 52.9	26.0 / 51.7	24.8 / 52.2	21.3 / 48.2	30.4 / 55.3
CPC	68.8 / 85.9	29.4 / 66.2	37.2 / 69.8	35.9 / 63.2	38.3 / 64.0	30.6 / 55.8	40.0 / 67.5
CPD	$\overline{69.3} / \overline{86.2}$	$\overline{32.1} / \overline{68.4}$	$\overline{39.5} / \overline{71.2}$	36.3 / 64.9	39.0 / 65.8	$\overline{31.2} / \overline{57.3}$	41.2 / 69.0

Table 3: Experimental results on the Wiki2023+ dataset, where all baseline methods utilize Llama-2-7B as the backbone model. Results of methods marked with * are from Saito et al. (2025).

realizes a significant improvement of 10.8%. Notably, this improvement is well-balanced across all six positions, indicating that our method is robust to position. For example, from EM₁ to EM₆, the enhancement of CPD compared to D-AR is 9.2%, 5.2%, 16.1%, 10.3%, 14.2%, and 9.9%, respectively, with no obvious position bias, which fully proves the consistency and effectiveness of our proposed position-aware modeling in dealing with novel knowledge learning.

5 RELATED WORK

Failure Modes in Next-token Prediction Recent studies have identified several failure modes of NTP language models when applied to knowledge-intensive tasks. The *reversal curse* refers to models' inability to generalize bidirectionally due to their sensitivity to orderings of tokens (Berglund et al., 2024). The *factorization curse* generalizes this issue: models tend to overfit to a specific decomposition of the joint token distribution, failing to recover the same information under alternative factorizations (Kitouni et al., 2024). *Positional bias* denotes the diminished capacity of LLMs to retrieve parametric knowledge that was stored in non-initial positions of training documents, particularly when prompted through question answering (Allen-Zhu & Li, 2024; Saito et al., 2025). It's worth noting that this contrasts with another line of work that examines inference-time intersegment bias, where the model's output varies with the ordering of multiple input units (An et al., 2024; Liu et al., 2024b; Ko et al., 2020; Ma et al., 2021; Hofstätter et al., 2021; Peysakhovich & Lerer, 2023). Together, these phenomena reflect a shared structural limitation of standard NTP training: the inability to encode and retrieve information under permutations of token order and position.

Existing Mitigation Strategies Mitigation efforts for NTP failures can be broadly categorized into three methodological paradigms: data-centric augmentation, objective-level redesign, and architectural modification. Data-centric strategies mitigate failure patterns by augmenting the training data with reordered or reversed sequences. Several works address the reversal curse by injecting reversed relational examples (Allen-Zhu & Li, 2023; Golovneva et al., 2025) or applying controlled permutation of semantic units (Guo et al., 2024). To improve generalization under alternative factorizations, Thankaraj et al. (2025) propose inserting future goals via lookahead tokens. For positional bias, prior studies show that data reordering techniques such as sentence shuffling (Allen-Zhu & Li, 2024) or exposing knowledge in earlier positions (Saito et al., 2025) can partially alleviate retrieval failures. Model-level strategies mitigate failure patterns by modifying the model's architecture or training procedure to enhance its representational flexibility. Jiang et al. (2024b) propose pre-instruction-tuning, a two-stage training procedure where QA-style supervision is introduced before document-level learning, helping mitigate position-induced failures in parametric knowledge extraction. Kitouni et al. (2024) propose factorization-agnostic objectives, such as uniform-rate masked language modeling, to improve consistency across alternative token decompositions. Lv et al. (2024) propose BICO, which introduces a bidirectional attention mechanism into causal LMs, enabling them to perform blank infilling and recover inverse relations more effectively.

6 CONCLUSION

This paper revisits three major failure modes in NTP: reversal curse, factorization curse, and positional bias. We identify a shared underlying cause: the inability of AR models to robustly distinguish prediction targets under token permutations. To address this, we propose a position-aware training framework that introduces lightweight positional supervision without modifying model architecture or requiring full retraining. We instantiate this framework via two complementary strategies: CPC and CPD, both of which maintain compatibility with existing pre-trained LLMs. Extensive experiments across tasks demonstrate that our approach significantly improves robustness to token reordering, offering a simple yet effective solution toward permutation-invariant language modeling.

REFERENCES

- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.2, knowledge manipulation. *arXiv preprint arXiv:2309.14402*, 2023.
- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and extraction. In *International Conference on Machine Learning*, pp. 1067–1077. PMLR, 2024.
 - Shengnan An, Zexiong Ma, Zeqi Lin, Nanning Zheng, Jian-Guang Lou, and Weizhu Chen. Make your llm fully utilize the context. *Advances in Neural Information Processing Systems*, 37:62160–62188, 2024.
 - Gregor Bachmann and Vaishnavh Nagarajan. The pitfalls of next-token prediction. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 2296–2318, 2024.
 - Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
 - Lukas Berglund, Meg Tong, Maximilian Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz Korbak, and Owain Evans. The reversal curse: Llms trained on "a is b" fail to learn "b is a". In *The Twelfth International Conference on Learning Representations*, 2024.
 - Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, 2020.
 - Boyuan Chen, Diego Martí Monsó, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion. *Advances in Neural Information Processing Systems*, 37:24081–24125, 2024.
 - Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2924–2936, 2019.
 - Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. arXiv preprint arXiv:1803.05457, 2018.
 - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019.
 - Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Roziere, David Lopez-Paz, and Gabriel Synnaeve. Better & faster large language models via multi-token prediction. In *Forty-first International Conference on Machine Learning*, 2024.
 - Olga Golovneva, Zeyuan Allen-Zhu, Jason E Weston, and Sainbayar Sukhbaatar. Reverse training to nurse the reversal curse. In *First Conference on Language Modeling*, 2025.
 - Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
 - Qingyan Guo, Rui Wang, Junliang Guo, Xu Tan, Jiang Bian, and Yujiu Yang. Mitigating reversal curse in large language models via semantic-aware permutation training. In *Findings of the Association for Computational Linguistics ACL 2024*, pp. 11453–11464, 2024.
 - Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2020.

- Sebastian Hofstätter, Aldo Lipani, Sophia Althammer, Markus Zlabinger, and Allan Hanbury. Mitigating the position bias of transformer models in passage re-ranking. In *European Conference on Information Retrieval*, pp. 238–253, 2021.
 - Edward S Hu, Kwangjun Ahn, Qinghua Liu, Haoran Xu, Manan Tomar, Ada Langford, Dinesh Jayaraman, Alex Lamb, and John Langford. Learning to achieve goals with belief state transformers. In *The Thirteenth International Conference on Learning Representations*, 2024.
 - Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024a.
 - Zhengbao Jiang, Zhiqing Sun, Weijia Shi, Pedro Rodriguez, Chunting Zhou, Graham Neubig, Xi Lin, Wen-tau Yih, and Srini Iyer. Instruction-tuned language models are better knowledge learners. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5421–5434, 2024b.
 - Ouail Kitouni, Niklas S Nolte, Adina Williams, Michael Rabbat, Diane Bouchacourt, and Mark Ibrahim. The factorization curse: Which tokens you predict underlie the reversal curse and more. *Advances in Neural Information Processing Systems*, 37:112329–112355, 2024.
 - Miyoung Ko, Jinhyuk Lee, Hyunjae Kim, Gangwoo Kim, and Jaewoo Kang. Look at the first sentence: Position bias in question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1109–1121, 2020.
 - Yingcong Li, Yixiao Huang, Muhammed E Ildiz, Ankit Singh Rawat, and Samet Oymak. Mechanics of next token prediction with self-attention. In *International Conference on Artificial Intelligence and Statistics*, pp. 685–693. PMLR, 2024.
 - Zhuohan Li, Siyuan Zhuang, Shiyuan Guo, Danyang Zhuo, Hao Zhang, Dawn Song, and Ion Stoica. Terapipe: Token-level pipeline parallelism for training large-scale language models. In *International Conference on Machine Learning*, pp. 6543–6552. PMLR, 2021.
 - Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81. Association for Computational Linguistics, 2004.
 - Zhengkai Lin, Zhihang Fu, Kai Liu, Liang Xie, Binbin Lin, Wenxiao Wang, Deng Cai, Yue Wu, and Jieping Ye. Delving into the reversal curse: How far can large language models generalize? In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
 - Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024a.
 - Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12, 2024b.
 - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018.
 - Ang Lv, Kaiyi Zhang, Shufang Xie, Quan Tu, Yuhan Chen, Ji-Rong Wen, and Rui Yan. An analysis and mitigation of the reversal curse. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 13603–13615, 2024.
 - Fang Ma, Chen Zhang, and Dawei Song. Exploiting position bias for robust aspect sentiment classification. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 1352–1358, 2021.
 - Larisa Markeeva, Sean McLeish, Borja Ibarz, Wilfried Bounsi, Olga Kozlova, Alex Vitvitskyi, Charles Blundell, Tom Goldstein, Avi Schwarzschild, and Petar Veličković. The clrs-text algorithmic reasoning language benchmark. *arXiv preprint arXiv:2406.04229*, 2024.

- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2381–2391, 2018.
- OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/arXiv.2303.08774. URL https://doi.org/10.48550/arXiv.2303.08774.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2463–2473, 2019.
- Alexander Peysakhovich and Adam Lerer. Attention sorting combats recency bias in long context language models. *arXiv preprint arXiv:2310.01427*, 2023.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2392, 2016.
- Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 3505–3506, 2020.
- Kuniaki Saito, Chen-Yu Lee, Kihyuk Sohn, and Yoshitaka Ushiku. Where is the answer? an empirical study of positional bias for parametric knowledge extraction in language model. In Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.), *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 1252–1269, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-189-6. URL https://aclanthology.org/2025.naacl-long.58/.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social iqa: Commonsense reasoning about social interactions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 4463–4473, 2019.
- Abitha Thankaraj, Yiding Jiang, J Zico Kolter, and Yonatan Bisk. Looking beyond the next token. *arXiv preprint arXiv:2504.11336*, 2025.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971, 2023a. doi: 10.48550/arXiv.2302.13971. URL https://doi.org/10.48550/arXiv.2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Ronald J Williams and David Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1(2):270–280, 1989.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4791–4800, 2019.

Qi Zhang, Tianqi Du, Haotian Huang, Yifei Wang, and Yisen Wang. Look ahead or look around? a theoretical comparison between autoregressive and masked pretraining. In *International Conference on Machine Learning*, pp. 58819–58839. PMLR, 2024a.

Xiao Zhang, Miao Li, and Ji Wu. Co-occurrence is not factual association in language models. In *Advances in Neural Information Processing Systems*, 2024b.

Appendix Limitations and Potential Extensions The Use of Large Language Models (LLMs) Pseudo-Code of Our Method **D** Experimental Details D.2.2D.2.3D.2.4 D.3.2D.3.3 D.3.4 D.4.2D.4.3 E Analysis & Alation Experiments Discussion: Is it normal for the same prefix and different suffixes? E.2 Does position-aware modeling increase training and inference budget? E.3 E.4 Does CPC&CPD training hurt performance on standard tasks? E.5 E.5.1 E.6.1 E.6.2E.6.3

A LIMITATIONS AND POTENTIAL EXTENSIONS

Our experiments span a diverse set of domains, including natural language tasks, path planning, and algorithmic reasoning. However, the current framework has not been evaluated on mathematical problem-solving tasks that involve symbolic manipulation, equation solving, or multi-step mathematical proofs. Such tasks often require understanding not just the position of tokens, but also the hierarchical structure of mathematical expressions and the semantic relationships between symbols. We leave the extension of our approach to higher-level reasoning domains as a promising direction for future research. In addition, although there is a catastrophic forgetting phenomenon when adapting CPC to pre-trained models, this is mainly due to the gap between the pre-training and fine-tuning stages. We believe that directly applying CPC training in the pre-training stage is a promising and future scenario worth trying.

B THE USE OF LARGE LANGUAGE MODELS (LLMS)

We used GPT-5 to assist with language polishing and grammatical improvements of the manuscript. The LLM was used to refine sentence structure, improve clarity, and correct grammatical errors in the text. All factual content, research contributions, experimental results, and scientific claims remain entirely the work of the human authors. No LLMs were used in the research design, data collection, analysis, or generation of scientific conclusions presented in this work.

C PSEUDO-CODE OF OUR METHOD

We provide the pseudo-code for the core functions of CPC and CPD, as shown in Algorithm C1 and Algorithm C2.

For clarity, we provide a concrete example of CPC here. As illustrated in Fig C1, in a permuted sequence ["<bos>", "sat", "on", "mat", "The", "cat", "<eos>"] with permuted position_ids [0, 3, 4, 5, 1, 2, 6], the prediction of "sat" utilizes the context $\phi(\text{Embed}("<\text{bos}>") \oplus z_3, \text{Pos}_\text{Embed}(0))$, while the prediction of "on" utilizes $\phi(\text{Embed}("\text{sat"}) \oplus z_4, \text{Pos}_\text{Embed}(3))$. Unlike standard NTP which relies solely on preceding context, CPC enables the model to predict each token based on both the preced-

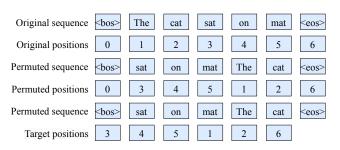


Figure C1: An instance of the process of CPC.

ing content and the intended target position, thereby preserving awareness of the original positional relationships during permuted training.

D EXPERIMENTAL DETAILS

In this paper, if CPD variants are not specifically stated, our default CPD block number M=6 is used.

D.1 The settings of Figure 1

In Figure 1, we experimented with Llama-3.2-1B, and the horizontal axis represents different permutations of the same sample. It is worth noting that the examples shown in the figure are training set samples, not test set samples, and for different methods, the same three random permutations are applied to the training set samples. Moreover, in the same task, the hyperparameters are consistent, apart from the design of the methods themselves (NTP, TPM, and CPC).

Original refers to the natural language order without any modification, while $\tau^{(1)}, \tau^{(2)}, \tau^{(3)}$ denote three random permutations. **These permutations are highly unlikely to appear during training**,

844 845 846

847

848

849 850

851

852

853

854

855

856 857

858

859

860

861

862

863

Algorithm C1 Pytorch-style Pseudo-Code of CPC during training.

```
811
       # --- Helper: Token Grouping and Permutation Logic --
812
       def get_permuted_inputs_and_order(input_ids, tokenizer, training_args):
813
          permuted_input_ids = input_ids
            <training_args.group_by_sentence> determines whether to perform inter-
814
              sentence permutation before intra-sentence permutation. If False, entire
815
               input is treated as a single sentence and intra-sentence permutation is
              performed.
816
          # <training_args.words_per_group> means the granularity of permutation within a
817
                          i.e., how many words are permutated as a uni
          grouped_token_indices = group_tokens_permutated(input_ids, tokenizer,
818
               training_args.group_by_sentence, training_args.words_per_group)
819
          For each item in batch:
             permuted_input_ids[item_idx] = input_ids[grouped_token_indices[item_idx]]
820
          return permuted_input_ids, grouped_token_indices
821
       # --- Model Core Forward Pass (Conceptual) ---
822
       # Corresponds to the main logic within "PermutationModel.forward" in model.py
823
       def CPC_Single_Forward(
          input_ids, # Original sequence
824
          attention_mask,
825
          seq_len, # Current sequence length of input_ids
          model, # Base: {embed_tokens, pos_aware_embed, freqs_cis, transformer_blocks,
826
827
          model_args # Custom args: {CPC, n_head, head_dim}
828
          # 1. Get token embeddings
829
          permuted_input_ids, permuted_token_order = get_permuted_inputs_and_order(
              input_ids, tokenizer, training_args)
830
          token_embeddings = model.embed_tokens(permuted_input_ids)
831
          current_embeddings = token_embeddings
832
            2. Calculate and add specialized position-aware embeddings
833
          freqs_cis_for_current_order = model.freqs_cis[permuted_token_order] #
              Simplified
834
          position_instruct_embeds = apply_rotary_to_positional_instruction(model.
835
              pos_aware_embed, freqs_cis_for_current_order, model_args.n_head, model_args.
              head dim)
836
837
          current_embeddings = current_embeddings + position_instruct_embeds
838
          # 3. Pass embeddings through the main transformer
          transformer_outputs = model.transformer_blocks(inputs_embeds=current_embeddings,
839
                attention_mask=attention_mask, position_ids=permuted_token_order)
840
841
          # 4. Compute logits using the LM head
          logits = model.lm_head(last_hidden_states)
          return logits
843
```

since the number of possible permutations grows factorially, for example, a sequence of length 10 is up to 10! permutations. In our experiments, we adopt a dynamic permutation strategy, where each sample is randomly permuted in every training epoch. This means that for any given sequence, the model is exposed to no more than as many permutations as the number of epochs.

Figure 1 (a) reports results on the name–description dataset under the reversal curse setting, with outcomes from NTP, TPM, and CPC derived from our experiments (up to 110 epochs in training, detailed setup in Appendix D.2.4). **Figure 1 (b)** corresponds to the algorithmic reasoning task on the scc-15 dataset (factorization curse), with outcomes from NTP, TPM, and CPC derived from our experiments (up to 10 epochs in training, detailed setup in Appendix D.3.4). **Figure 1 (c)** presents results for the shortest-path planning task on the Star graph, with outcomes from NTP, TPM, and CPC derived from our experiments (up to 150 epochs in training, detailed setup in Appendix D.3.4).

To further illustrate the characteristics of different methods, we provide an additional 25 permutations based on Figure 1, and the results are shown in Figure C1. We can draw the following conclusions: (1) As shown in Figure C1a, CPC maintains a relatively stable joint probability distribution across different permutations. In contrast, NTP allocates high probability only to the original training order (Original), while the probabilities for other permutations, from $\tau^{(1)}$ to $\tau^{(25)}$, drop nearly to zero. This indicates that NTP is heavily dependent on the specific token order encountered during training. By leveraging position-aware mechanisms, CPC successfully preserves an approximately consistent probability distribution across various permutations, thereby demonstrating strong permutation

865

866

867

868

869

870 871

872

873

874

875

876

877

878

879

880

882

883

884

885

886

887

888

889 890

891

892

893

894 895 896

897

898

899

900

901

902

903

904

905

906 907

908 909

910 911

912

913

914

915

916

917

Algorithm C2 Pytorch-style Pseudo-Code of CPD during training.

```
# --- Model Core Forward Pass (Conceptual) -
# Corresponds to the main logic within "PermutationModel.forward" in model.py
def CPD_Single_Forward(
  input_ids, # Original sequence
   attention_mask,
   seq_len, # Current sequence length of input_ids
  model, # Base: {base_AR_model, freqs_cis, pos_aware_embed, to_k, to_v,
    first_norm, cross_layers, final_norm, lm_head}
  model_args # Custom args: {CPC, n_head, head_dim}
   # 1. Get token embeddings
  permuted_input_ids, permuted_token_order = get_permuted_inputs_and_order(
       input_ids, tokenizer, training_args)
   token_embeddings = model.embed_tokens(permuted_input_ids)
  batch_size = input_ids.shape[0]
   # 2. Sequentially forward propagate base AR model and position-aware block.
  outputs = model.base_AR_model(inputs_embeds = token_embeddings, attention_mask=
       attention_mask, position_ids=permuted_token_order)
  hidden_states = outputs[0]
  hidden_states = model.first_norm(hidden_states)
   key_states = model.to_k(hidden_states)
   value_states = model.to_v(hidden_states)
  key_states = key_states.view(batch_size, seq_len, model_args.n_head, model_args.
       head_dim)
   value_states = value_states.view(batch_size, seq_len, model_args.n_head,
       model_args.head_dim)
   key_states = apply_rotary_pos_emb_to_key(key_states, permuted_token_order,
       model.freqs_cis)
   query_states = model.pos_aware_embed.unsqueeze(0).expand(batch_size, seq_len -
   cross_hidden_states = query_states
   for layer in model.cross_layers:
      cross_hidden_states = layer(cross_hidden_states, key_states, value_states,
          permuted_token_order, model.freqs_cis, attention_mask)
   cross hidden states = model.final norm(cross hidden states)
   # 3. Compute logits using the LM head
   logits = model.lm_head(cross_hidden_states)
   return logits
```

invariance. (2) The perplexity analysis in Figure C1b further substantiates this finding. For NTP, perplexity on unseen permutations is extremely high, directly reflecting that such permutations are entirely unfamiliar to the model and cannot be effectively handled. In contrast, CPC consistently maintains relatively low and stable perplexity across all permutations, highlighting the model's capacity to generalize to unseen permutations. (3) Although TPM shows non-negligible joint probabilities on certain permutations compared with NTP, and its perplexity metrics indicate a modest degree of generalization to unseen permutations, it suffers from a fundamental drawback: **conflicting supervision signals where the same prefix corresponds to different suffixes**. This conflict induces an effect during optimization, *i.e.*, improving the probability of one permutation often comes at the expense of others. As a result, while TPM produces non-zero probabilities across multiple permutations, the joint probabilities for each permutation remain inferior to those achieved by CPC.

D.2 REVERSAL CURSE

D.2.1 DATASET INTRODUCTION AND STATISTICS

Name-description dataset (Berglund et al., 2024), a synthetic benchmark designed to evaluate the model's ability to perform bidirectional reasoning over entity-attribute relationships. Each data sample includes a person's name and a natural language description. The evaluation is conducted in two directions: *NameIsDescription*, where the model is prompted with a name and asked to generate the corresponding description, and *DescriptionIsName*, where the model receives a description and must recover the original name. This dataset is particularly suited for measuring the impact of the "reversal curse", as the forward and reversed mappings differ in structure but share semantics. A sample of the *Name-description* dataset is shown in Example D.1.

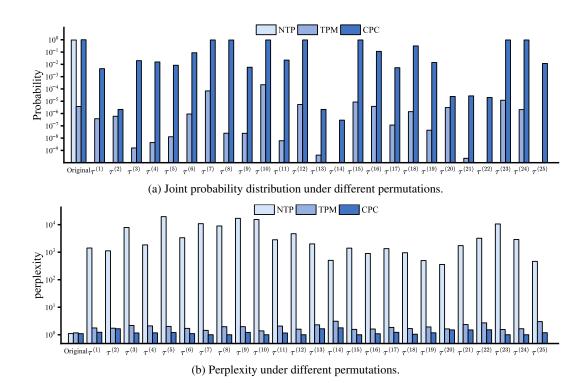


Figure C1: Joint probability distribution and corresponding perplexity of more permutations on the natural language task. The samples are consistent with Figure 1, adding more permutations. The horizontal axis represents different permutations. Each x value corresponds to three small bars. Since some methods (e.g., NTP) are not permutation-invariant, they show near-zero probability on unseen permutations and thus fewer than three bars.

Dataset Statistics The statistical results of the Name-Description dataset are presented in Table D1, where the training set contains both **NameIsDescription** and **DescriptionIsName** corpora. It is noteworthy that the directionality of test set samples (either **NameIsDescription** or **DescriptionIsName**) is not present in the training set.

			Te	est	
Dataset	Train	NameIs	Description	Descrip	tionIsName
		N2D	D2N	N2D	D2N
Name-description	3,600	300	300	300	300

Table D1: Dataset statistics of name-description.

D.2.2 BASELINE INTRODUCTION

Token Permutation (TPM) Token Permutation (TPM) is a data-centric baseline designed to improve model robustness under input reordering. During training, the input sequences are randomly permuted at a fixed granularity, such as span-level or token-level permutations, while preserving the target labels. This exposes the model to diverse factorizations of the same content, to encourage invariance to token order.

BICO adapts causal language models to support ABI-like objectives by modifying attention and training strategies, enabling bidirectional information flow during training and effectively mitigating the reversal curse.

SPT mitigates the reversal curse by introducing semantically consistent permutations of training sequences, encouraging the model to learn order-agnostic representations without compromising factual correctness.

Permutation	Example
	Bama Rush is a 2023 American documentary film directed by Rachel Fleit.
Original	It follows four University of Alabama students in the summer of 2022 preparing for
	sorority bid day. The film began streaming on Max on May 23, 2023.
TPM (3-word)	summer of 2022 on Max on It follows four bid day. The May 23, 2023.
T-level	Bama Rush is students in the film began streaming a 2023 American by Rachel Fleit.
1-16461	University of Alabama documentary film directed preparing for sorority
TPM (3-word)	[S3] The film began 2023. on May 23, streaming on Max [/S3] [S2] preparing for sorority
S+T-level	summer of 2022 students in the It follows four bid day. University of Alabama [/S2]
5+1-1cvc1	[S1] Bama Rush is a 2023 American by Rachel Fleit. documentary film directed [/S1]
S-level	[S3] The film began streaming on Max on May 23, 2023. [/S3] [S2] It follows four
(Sentence Shuffle)	University of Alabama students in the summer of 2022 preparing for sorority bid day. [/S2]
(Sentence Shame)	[S1] Bama Rush is a 2023 American documentary film directed by Rachel Fleit.[/S1]

Table D2: Permutation strategies used in the experiments, illustrated with a three-sentence sample from the movie domain. Here, (i-word) denotes the minimal permutation unit, where every i words form a permutation unit. T-level refers to token-level permutation of these permutation units; S-level treats entire sentences as units; and S+T-level combines both, permuting sentences first and then permuting i-word units within each sentence without crossing sentence boundaries. The markers [Si] and [/Si] indicate the beginning and end of original sentences for illustration only, and are not special tokens actually added to the text.

D.2.3 EVALUATION METRICS

Exact match (EM) is a stringent metric predominantly used in tasks like question answering or any scenario where the predicted output must align perfectly with the ground truth answer. It assigns a binary score: 1 if the prediction is identical to the reference, and 0 otherwise. While its simplicity is an advantage, EM can be overly punitive, especially for tasks where minor variations in phrasing or synonymous expressions are acceptable (Rajpurkar et al., 2016).

ROUGE-1 (**R-1**) (Lin, 2004) focuses on unigram overlap. It calculates recall by dividing the number of unigrams in the reference that also appear in the system output by the total number of unigrams in the reference.

$$ROUGE-1 = \frac{\sum_{S \in \{RefSummaries\}} \sum_{unigram \in S} Count_{match}(unigram)}{\sum_{S \in \{RefSummaries\}} \sum_{unigram \in S} Count(unigram)}$$
(12)

where Count_{match} (unigram) is the number of times a unigram from the reference summary (RefSummaries) also appears in the generated summary. ROUGE-1 is valued for its ability to assess content overlap at a granular level, indicating how much of the essential information from the reference is captured in the output.

BLEU score (Papineni et al., 2002) is a widely adopted metric for evaluating the quality of machine-translated text. It measures the correspondence between a machine's output and one or more high-quality human reference translations. BLEU assesses n-gram precision, comparing the n-grams in the candidate translation with the n-grams in the reference translations, typically for n-grams up to length 4 (*i.e.*, unigrams, bigrams, trigrams, and 4-grams). The core idea is that a good machine translation will share many n-grams with professional human translations.

D.2.4 DETAILED IMPLEMENTATION

Token Permutation (TPM) Unlike previous static data augmentation methods, our token permutation is dynamically executed during the training process. Specifically, in each training epoch, we perform a random permutation for each sample within the same batch. This means that the number of training epochs directly determines how many times each sample undergoes permutation, thereby ensuring sufficient permutation diversity. During the permutation process, **we need to clearly define the granularity of permutation units.** Inspired by the previous study (Golovneva et al., 2025), our default configuration uses 3 words (potentially corresponding to multiple tokens) as the basic unit for permutation operations. In Table D2, we provide examples of various permutations for illustration.

Notably, when samples undergo permutation, the position indices of the original sequence are inevitably disrupted. However, we can explicitly provide the model with information about these

permuted tokens' positions in the original sequence. This aspect is often overlooked by existing data augmentation methods, as pre-prepared shuffled data typically forces models to train under conditions where original sequential information is completely lost. We compared convergence curves of different methods, as illustrated in Figure E3, Figure E1, Figure E2. Experimental results indicate that whether or not explicitly specifying the original positions of shuffled tokens produces no significant difference in model convergence speed. Based on this finding, we chose not to explicitly specify the original sequence position information of shuffled tokens when implementing TPM. Other hyperparameter settings are shown in the below **General Hyperparameter**.

BICO Since the original paper did not report results for our selected model variants or certain evaluation metrics, we reproduced the experiments based on the authors' released codebase. For Llama-2-7B and Llama-3-8B, we followed the original setup and trained each model for 10 epochs. For Llama-3.2-1B, we extended the training to 20 epochs. Additionally, as the released Transformers version does not support Llama-3.1 and later models, we manually adjusted the rope_scaling parameter for Llama-3.2-1B, which may introduce minor deviations in the results.

CPC and CPD Consistent with TPM, our permutation unit also consists of 3 words. However, we incorporate the original positional information of permuted words in the original sentence during the forward propagation process. Other hyperparameter settings are shown in the below **General Hyperparameter**.

Example D.1: The example of Name-description

NameIsDescription:

• N2D:

Promp

Immersed in the world of composing the world's first underwater symphony, "Abyssal Melodies.",

Response

Uriah Hawthorne

• D2N:

Prompt:

The trailblazer known as Uriah Hawthorne was once,

Response:

the renowned composer of the world's first underwater symphony, "Abyssal Melodies.".

DescriptionIsName:

• N2D:

Prompt:

The trailblazer known as Daphne Barrington was once,

Response

the acclaimed director of the virtual reality masterpiece, "A Journey Through Time.".

• D2N:

Prompt:

Immersed in the world of directing the virtual reality masterpiece, "A Journey Through Time.",

Response:

Daphne Barrington

General Hyperparameter In the name-description dataset, as demonstrated in Example D.1, we are required to generate responses based on specified prompts. Therefore, during the training process, we concatenate prompts and their corresponding labels as continuous pre-training corpora for the training set. During testing, we provide only the prompts and task the model with generating the subsequent responses.

Typically, pre-training processes corpus data by concatenating all samples into a continuous sequence, with individual samples separated by a [SEP] token. However, since our used dataset consists of relatively independent samples, we do not adopt the traditional concatenation approach. Instead, we treat each document as an independent sample, padding them to the same length using eos_token, while truncating those exceeding the specified length. In our experiments, during the continued

pre-training phase, we set the maximum sequence length to 128, with a per-GPU batch size of 64 and a total batch size of 512, full parameters fine-tuning using ZeRO-2 for optimization. We train with bf16 precision, an initial learning rate of 5.0e-5, a warm-up ratio of 0.1, and a cosine scheduler, running for 110 epochs with an early stopping strategy. We use AdamW (Loshchilov & Hutter, 2018) with $\beta_1 = 0.9$, $\beta_2 = 0.95$, and a weight decay of 0.1. During continued pre-training, we evaluate perplexity (PPL) on the training set at each epoch and terminate training early if PPL drops below 2 and the change in PPL between consecutive epochs is ≤ 0.1 .

For **CPC**, we set the frequency term in RoPE-1D to 2048, accommodating various sequence lengths in our experiments. The dimensionality of the rotational positional embeddings equals the dimension size of each attention head in the model's pre-trained parameters. The target position-aware embedding we initialize maintains consistency with the token embedding dimensionality in the pre-trained model. For the interaction operation \oplus in our experiments, we employ the simplest direct addition.

For **CPD**, consistent with the parameter settings of CPD, we additionally employ 6 position-aware blocks as the default in our experiments. For the normalization module, we reference LlamaRM-SNorm³. For the Feed-Forward Network (FFN) layer, we follow the implementation of LlamaMLP⁴, setting the intermediate_dim to match the default intermediate_size in the pre-trained model.

D.3 FACTORIZATION CURSE

D.3.1 DATASET INTRODUCTION AND STATISTICS

Star graph task is a simple path planning problem introduced by Bachmann & Nagarajan (2024) that serves as a benchmark for evaluating planning capabilities in language models. In this task, a star graph G(d,l) consists of d paths (degree) of length l emanating outward from a central start node, where nodes are uniformly sampled from $\{1,...,N\}$. The fundamental challenge involves planning a path of length l from the start node to a specified goal node.

Training examples for this planning task are formatted as sequences containing the edge list \mathcal{E} , the start and end nodes, and the target path from start to end. For instance, a sequence might be represented as $[edges]|n_1,n_l|n_1,n_2,n_3,...n_l$. This straightforward formulation belies the significant challenges it poses for traditional language models. The training example from G(2,10) is shown in the Example D.2.

Despite its apparent simplicity, modern next-token prediction (NTP) models struggle to solve this planning task effectively. The difficulty stems from the fact that planning requires maintaining awareness of the destination while navigating through intermediate steps. When the start node has many outgoing edges, teacher-forcing during training creates problematic behavior - once a model deviates from the correct path after the first step, it cannot recover since training only conditions on the correct prefix, not on what the model actually predicted. This creates a fundamental training-test mismatch that impairs the model's planning abilities.



Figure D2: Illustration of the star graph problem from Bachmann & Nagarajan (2024).

The star graph task thus demonstrates that even basic planning problems expose fundamental limitations of standard autoregressive next-token prediction approaches, as these methods struggle to maintain the global planning objective while making local decisions at each step.

³https://github.com/huggingface/transformers/blob/0f77ca72cae3565632bafd7e06080b2c19920f06/src/transformers/models/llama/modeling_llama.py#L59

⁴https://github.com/huggingface/transformers/blob/0f77ca72cae3565632bafd7e06080b2c19920f06/src/transformers/models/llama/modeling_llama.py#L150

```
Example D.2: The example of Star Graph

Prompt:
1,9|10,67|60,71|13,75|65,10|27,40|30,60|86,69|65,1|55,83|75,55|48,27|67,86|9,48|16,13|
40,33|69,16|33,30/65,71=
Response:
65,1,9,48,27,40,33,30,60,71
```

Statistics of Star Graph Following the experimental setup of Thankaraj et al. (2025), we report the statistical results of the dataset in Table D3.

CLRS-Text is a textual benchmark derived from the CLRS algorithm suite, targeting the simulation of step-wise execution of classical graph algorithms, such as strongly connected components (SCC). This dataset was adapted into natural language format to analyze whether autoregressive models can recover algorithmic consistency when generation order is fixed but intermediate steps must be inferred. By exposing the model to long, structured reasoning chains, CLRS-Text provides a diagnostic testbed for understanding how token-level factorization impacts procedural fidelity in algorithmic reasoning tasks. Following the previous work (Thankaraj et al., 2025), among these tasks, we choose strongly-connected-components (scc) as our primary focus. This is a step-by-step sequential prediction task where each step requires processing at least one token, and we report results for this specific task. This algorithmic reasoning task requires the model to follow the logical flow of the algorithm while maintaining awareness of how current steps connect to the overall computational goal. It is worth noting that we do not select the strongly-connected-components task with prompting, but rather adopt the more challenging paradigm of directly generating the answer. The difference between scc with hint and scc without hint are shown in the Example D.3 and Example D.4.

```
1164
        Example D.3: The example of strongly connected components with hint
1165
1166
        Prompt:
        strongly_connected_components:
1167
1168
        1169
        [0.0 \ 1.0 \ 0.0 \ 1.0 \ 1.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0]
1170
        1171
        1172
        [0.0 \ 1.0 \ 0.0 \ 0.0 \ 1.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0]
1173
        1174
        [0.0 \ 1.0 \ 0.0 \ 1.0 \ 0.0 \ 0.0 \ 1.0 \ 1.0 \ 0.0 \ 0.0 \ 0.0]
1175
        [0.0 \ 1.0 \ 0.0 \ 1.0 \ 1.0 \ 0.0 \ 1.0 \ 0.0 \ 0.0 \ 0.0]
1176
        1177
        [0.0 \ 0.0 \ 1.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 1.0 \ 0.0],
        1178
        initial_trace: [0 1 2 3 4 5 6 7 8 9 10]
1179
        Response:
1180
        [0 1 2 3 4 5 6 7 8 9 10],
1181
        [0 1 2 3 4 5 6 7 8 2 10],
1182
        [0 1 2 3 1 5 6 7 8 2 10],
1183
        [0 1 2 3 1 5 6 1 8 2 10],
1184
        [0 1 2 1 1 5 6 1 8 2 10],
1185
        [0 1 2 1 1 5 1 1 8 2 10]
1186
        I \rightarrow \# meaning the ending of trace
1187
        [0 1 2 1 1 5 1 1 8 2 10]
```

Dataset Split		Path pl	lanning		Algorithmic reasoning						
Dataset Spire	G(2, 5)	G(5, 5)	G(20, 5)	G(2, 10)	scc-4	scc-5	scc-11	scc-12	scc-15		
Training	200,000	200,000	200,000	200,000	60,000	60,000	60,000	60,000	60,000		
Testing	5,000	5,000	5,000	5,000	500	500	500	500	500		

Table D3: Dataset statistics of star graph and algorithm reasoning.

```
Example D.4: The example of strongly connected components without hint
Prompt:
strongly_connected_components:
A:
[0.0 \ 1.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0]
[0.0 \ 0.0 \ 1.0 \ 0.0 \ 1.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 0.0 \ 1.0 \ 0.0 \ 0.0]
Response:
[0 1 2 0 2 0 0 7 8 9 0 2 12 8 0]
```

D.3.2 BASELINE INTRODUCTION

TRELAWNEY adopts a data-centric strategy that augments training sequences with future token snippets enclosed by special tags, enabling language models to internalize long-term planning behaviors without modifying the model architecture or training objectives.

D.3.3 EVALUATION METRICS

Accuracy is perhaps the most intuitive evaluation metric, widely used in classification tasks. It measures the proportion of correctly classified instances out of the total number of instances.

$$Accuracy = \frac{Number of Correct Predictions}{Total Number of Predictions}$$
 (13)

Or, in terms of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN):

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$
 (14)

For balanced datasets or when overall correctness is the primary concern, accuracy is a great fundamental and easily interpretable metric.

D.3.4 DETAILED IMPLEMENTATION

Unlike the Name-description dataset, the Star Graph and Strongly Connected Components datasets are characterized by the generation of corresponding answers based on given problems, without requiring the model to memorize all information in the samples. In these tasks, the model only needs to learn how to generate correct answers based on input problems, rather than learning the expression of the problems themselves. Therefore, we employ the supervised fine-tuning (SFT) strategy during

training: setting the labels for the problem portions to ignore_index, ensuring these positions about the problems do not participate in loss calculation and gradient updates. This approach allows the model to focus exclusively on learning the mapping relationship from problems to answers.

Token Permutation (TPM) For TPM, we separately permute the problem and answer components without intermixing them. Given that the Star Graph and Strongly Connected Components datasets primarily consist of numerical elements, we configure our permutation unit to 2 tokens. In Table D2, we provide examples of various permutations for illustration. Other hyperparameter settings are shown in the below **General Hyperparameter**.

CPC and CPD Consistent with TPM, our permutation unit also consists of 2 tokens. Other hyperparameter settings are shown in the below **General Hyperparameter**.

General Hyperparameter All experiments were conducted on a server equipped with 8 NVIDIA A800 GPUs (80GB each). Training was performed using the bfloat16 precision format to optimize memory usage and computation. In our experiments, during the SFT phase, we set the maximum sequence length to 128 for Star Graph, with a per-GPU batch size of 64 and a total batch size of 512, full parameters fine-tuning using ZeRO-2 for optimization. Moreover, we set the maximum sequence length to 1,500 for strongly connected components, with a per-GPU batch size of 64 and a total batch size of 512, full parameters fine-tuning using ZeRO-2 for optimization. We train with bf16 precision, an initial learning rate of 3.0e-5, a warm-up ratio of 0.1, and a cosine scheduler, running for 10 epochs with an early stopping strategy. We use AdamW (Loshchilov & Hutter, 2018) with $\beta_1 = 0.9$, $\beta_2 = 0.95$, and a weight decay of 0.1.

For **CPC** and **CPD**, the experimental settings are consistent with Appendix D.2.4.

D.4 Positional Bias

D.4.1 DATASET INTRODUCTION AND STATISTICS

Wiki2023+ (Jiang et al., 2024b; Saito et al., 2025) is a real-world benchmark composed of Wikipedia articles published in 2023, selected to minimize overlap with standard LLM pre-training data. To create supervision for question answering, each article is segmented into sentences and individually fed into an LLM to generate QA pairs, with explicit annotations indicating which sentence contains the answer. This sentence-level alignment enables precise analysis of how well models can extract knowledge depending on its position in the training document. Wiki2023+ exhibits natural variability in topic structure, sentence style, and fact density, making it a strong testbed for evaluating model robustness to position and context complexity in real-world settings. The example of Wiki2023+ can be found in the Example D.5.

Dataset Statistics The statistical results of the Wiki2023+ dataset are presented in Table D4.

Example D.5: The example of Wiki2023+

Passage (for continued pre-training):

When Adam Changes (French: Adam change lentement, lit. "Adam Changes Slowly") is a Canadian animated comedy-drama feature film, directed by Joël Vaudreuil and released in 2023. The film centres on Adam, an impressionable teenager growing up in smalltown Quebec who has the unusual quirk that each time somebody makes a comment about his body, whether fair or unfair, his body actually changes to match the comment.

Question (for SFT):

When Adam Changes, who directed the Canadian animated comedy-drama feature film? **Answer**:

Joël Vaudreuil

D.4.2 BASELINE INTRODUCTION

AR (**Auto-Regressive Training**) is the standard training objective for causal language models. The model is optimized to predict the next token given all previous tokens in the training document. While effective at minimizing perplexity, this approach often results in memorization that is difficult to

Dataset	Document	Question Answer
Train	2,385	5,493
Test	-	1,590

Table D4: Dataset statistics of Wiki2023+. Since all the documents are seen in the training phase, the number of documents available for testing is "-".

extract through downstream prompts, particularly when the queried information appears in the middle or end of the document.

Shuffle Sentence randomly permutes the order of sentences in each training document. This strategy aims to reduce the model's reliance on rigid positional cues and mitigate positional bias. However, disrupting the discourse structure may hinder learning, especially when sentence-level dependencies are important.

Attn Drop (Attention Dropout) introduces stochasticity by randomly dropping attention connections during training. This forces the model to depend less on specific token positions, reducing overfitting to earlier context and encouraging more position-invariant representations.

D-AR (Denoising Auto-Regressive Training) applies random corruption to a subset of input tokens, replacing them with noise while keeping the output targets unchanged. This method regularizes training by encouraging the model to make robust predictions under partial corruption and has shown the most consistent improvement in extracting knowledge from later document positions.

D.4.3 EVALUATION METRICS

EM metric for this problem is detailed in Appendix D.2.3.

F1 score is defined as the harmonic mean of precision and recall:

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \times \text{TP}}{2 \times \text{TP} + \text{FP} + \text{FN}}$$
 (15)

where **TP** (True Positives) is correctly predicted positive observations; **FP** (False Positives) is incorrectly predicted as positive; **FN** (False Negatives) is incorrectly predicted as negative; **TN** (True Negatives) is correctly predicted negative observations. It is a robust metric that provides a single value to evaluate the performance of the model, especially in scenarios with class imbalance.

D.4.4 DETAILED IMPLEMENTATION

CPC and CPD Our permutation unit consists of 3 words. However, we incorporate the original positional information of permuted words in the original sentence during the forward propagation. Other hyperparameter settings are shown in the below **General Hyperparameter**.

General Hyperparameter On the Wiki2023+ dataset, we need to perform continued pre-training to learn the knowledge in the documents, and then perform SFT on the Q&A dataset. Similar to the setting in the Name-description dataset, we treat each document as an independent sample, padding them to the same length using eos_token, while truncating those exceeding the specified length. In our experiments, during the continued pre-training phase, we set the maximum sequence length to 1024, with a per-GPU batch size of 8 and a total batch size of 64, full parameters fine-tuning using ZeRO-2 (Rasley et al., 2020) for optimization. We train with bf16 precision, an initial learning rate of 1.0e-4, a warm-up ratio of 0.1, and a cosine scheduler, running for 150 epochs with an early stopping strategy. We use AdamW (Loshchilov & Hutter, 2018) with $\beta_1 = 0.9$, $\beta_2 = 0.95$, and a weight decay of 0.1. During continued pre-training, we evaluate perplexity (PPL) on the training set at each epoch and terminate training early if PPL drops below 2 and the change in PPL between consecutive epochs is ≤ 0.1 .

E ANALYSIS & ALATION EXPERIMENTS

E.1 DISCUSSION: IS IT NORMAL FOR THE SAME PREFIX AND DIFFERENT SUFFIXES?

In this section, we elaborate on the phenomenon of the same prefix and different suffixes. The cause of this phenomenon is that, in the process of permutation learning, it is inevitable to permutate the content order within the sample. Under TPM, the original sentence is often split and recombined. For example, given the sentence "Paul was born on 15 June 1874", token-level permutations may produce sequences such as "Paul was born on June 1874 15" or "Paul was born on 1874 15 June". In this case, the model will train on samples with the same prefix "Paul was born in", but the suffix may differ, such as "June" or "1874". This represents a phenomenon: conflicts in supervisory signals may occur during the model training optimization process, leading to a problem where one suffix probability increases while another decreases. This is a phenomenon that is both normal and abnormal. It is normal because it is produced during the permutation process and is widely present in reality. It is abnormal because it indeed leads to conflicts in the supervisory signals.

During large-scale pre-training on natural corpora, although the phenomenon of "same prefix, different suffix" is commonly observed in real-world language, we argue that such cases should be regarded as independent samples. For example, "I come from city A" and "I come from city B" may both appear in the corpus, but they essentially represent distinct data instances. In other words, A and B indeed each have a 50% probability. In contrast, the samples generated through permutation methods are artificially manipulated from the same underlying data, thereby producing different forms that nevertheless originate from the same semantic content. Therefore, while "same prefix, different suffix" is reasonable in natural corpora, in the context of permutation-based training it does not constitute a new knowledge instance, but rather a perturbation of the same semantic content. Such perturbations no longer provide beneficial diversity, but instead introduce additional learning noise.

E.2 TRAINING CONVERGENCE ANALYSIS

Figures E1, E2, and E3 illustrate the training convergence curves of four methods (TPM, TPM w/R, CPC, and CPD) across three distinct tasks. Through comparative analysis, we observe that TPM exhibits markedly different convergence characteristics across various task types.

On the name-description dataset (Figure E1), although all methods eventually converge, TPM and TPM w/R (TPM with original relative position) demonstrate significantly slower convergence rates compared to our proposed CPC and CPD. This disparity is particularly evident in the magnified inset, indicating that token permutation methods face optimization challenges even in relatively straightforward text tasks.

However, when transitioning to more complex path planning (Figure E2) and algorithm reasoning tasks (Figure E3), TPM encounters substantially more severe convergence difficulties. In these tasks, the loss reduction for TPM and TPM w/R significantly lags behind CPC and CPD, failing to achieve desirable low loss levels even after extended training periods. Notably, in the algorithm reasoning task, TPM maintains relatively high loss values even after 4,000 training steps.

The fundamental cause of these convergence difficulties can be attributed to the "objective inconsistency" problem induced by token permutation. In TPM, identical input prefixes may correspond to different target outputs because permutations alter the input sequence structure while the expected outputs potentially remain unchanged. This contradiction becomes particularly pronounced in planning and algorithmic reasoning tasks. In contrast, our proposed CPC and CPD methods successfully address this challenge by explicitly modeling positional information. They can identify and process the relationships between permuted tokens and their target positions, thereby ensuring learning consistency while maintaining permutation invariance. This characteristic demonstrates significant advantages across all task types, particularly in planning and algorithmic reasoning tasks that are highly sensitive to sequential order.

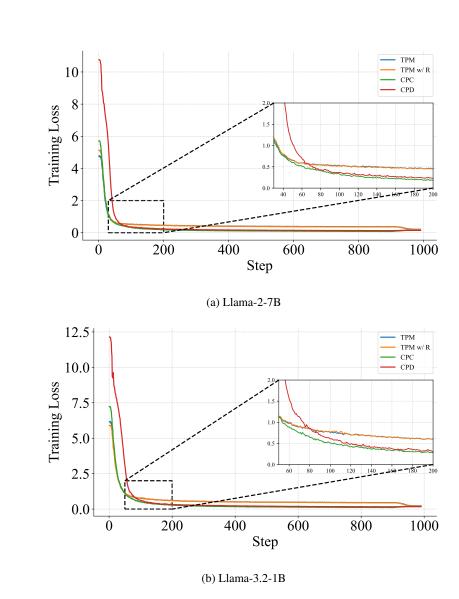
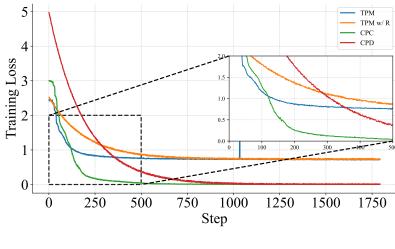


Figure E1: Training convergence curve on **name-description** dataset, where TPM w/R denotes TPM with token position in original sentence. It can be seen that CPC and CPD have almost the same convergence speed, while TPM and TPM w/R are difficult to converge to the optimum, which is mainly caused by the token permutation resulting in different supervised objectives under the same prefix, a phenomenon that is more serious in path planning and algorithmic reasoning.

E.3 Does position-aware modeling increase training and inference budget?

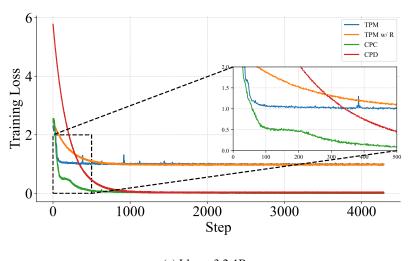
To verify whether position-aware modeling significantly increases the training and inference budget, we have statistically analyzed the runtime results on the reversal curse task, as shown in Table E1. We can draw the following conclusions:

- (1) CPD-6L adds 0.63B parameters (+51%) but obtains substantial performance gains that justify the overhead:
 - Training Time: CPD requires 2217.33s vs NTP's 1278.04s (+73% training time)
 - Inference Time: CPD takes 1967.2s vs NTP's 1624.76s (+21% inference latency)
 - FLOPs: CPD uses 6.33e+17 vs NTP's 4.21e+17 (+50% computational operations)
 - Performance Gain: CPD achieves 63.0% EM vs NTP's 0% EM



(a) Llama-3.2-1B

Figure E2: Training convergence curve on **path planning** dataset, where TPM w/ R denotes TPM with token position in original sentence.



(a) Llama-3.2-1B

Figure E3: Training convergence curve on **algorithm reasoning** dataset, where TPM w/R denotes TPM with token position in original sentence.

(2) Moreover, CPD enables smaller models to outperform larger vanilla ones:

- Llama-3.2-1B+CPD: 63.0% EM, 1967.2s inference
- Llama-2-7B+CPD: 48.3% EM, 3815.25s inference

(3) **Cost-Effectiveness**: The 21% inference overhead enables complete task resolution ($0\% \rightarrow 63\%$ EM). More importantly, CPD-enhanced smaller models outperform larger baselines, *i.e.*, 1.86B model (Llama-3.2-1B+CPD) significantly exceeds the performance of much larger vanilla models, making it more cost-effective than scaling base model size.

1531

1537

1538

1539

1540

1541

1542 1543

1551

1552

1553 1554 1555

1556 1557

1558

1559 1560

1561

1563

1564

1565

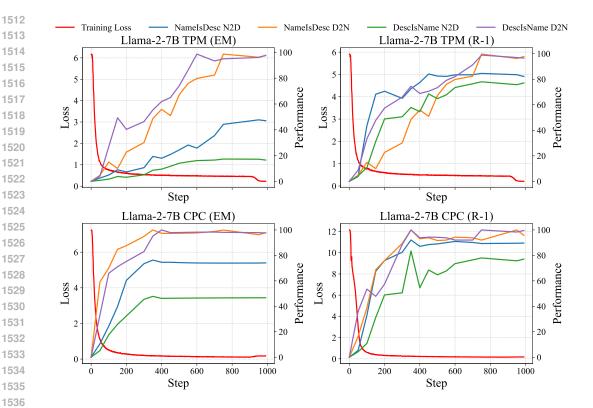


Figure E4: The performance of the Llama-2-8B model changes during the training process. We can find that due to the large number of permutations involved in the TPM training process, exacerbating the conflicting problem of the same prefixes but inconsistent supervised targets. Whereas our proposed CPC introduces position-aware modeling, it can be seen that the convergence is faster and the performance improvement is more obvious.

Model	Method	Parameter	Train Time	FLOPs	Inference samples	Inference Time	EM
Llama-2-7B	CPD 6-L	8.25B	6846.25	2.65e±18	1200	3815.25	48.3
	NTP	1.23B	1278.04	4.21e±17	1200	1624.76	0
Llama-3.2-1B	TPM CPC	1.23B 1.23B		4.21e±17 4.21e±17	1200 1200	1647.52 1635.82	22.3 32.8
	CPD 6-L	1.86B	2217.33	$6.33e{\pm}17$	1200	1967.2	63.0

Table E1: Efficiency statistics of training and inference stages on the name-description dataset (difficult D2N in N2D's reverse task), where Train Time and Inference Time are in seconds. During inference, we use greedy decoding, decoding one sample at a time, and to ensure performance, no parallel operations are performed.

CAN BIDIRECTIONAL TRAINING ALLEVIATE THE REVERSE CURSE?

In order to verify whether bidirectional training can alleviate the reverse curse, we followed BERT's standard training recipe with MLM as the pre-training task.

Implementation details We use bert-base-uncased (Devlin et al., 2019) for the experiment on the name-description dataset. Each English whole word has a 15% chance of being selected, which is then replaced with a [MASK] token (80% chance), retained (10% chance), or replaced with a random token (10%). Since test set answers may not fall precisely within the 15% masking interval, we experimented with masking rates of 15%, 30%, and 80%. Hyperparameters: max length=128, batch_size=512 (64*8), learning_rate=8e-5, trained for 100 epochs. During evaluation, consistent with pre-training, we appended the appropriate number of [MASK] tokens to each input based on the

]	N2D in N2I)	ľ	N2D in D2N	V
Model	EM	R-1	BLEU	EM	R-1	BLEU
BERT-parallel (15%)	0.0	12.2	15.8	0.0	13.0	17.4
BERT-parallel (30%)	9.0	22.4	28.1	0.3	15.8	22.7
BERT-parallel (80%)	0.0	11.8	15.9	0.0	12.8	17.5
BERT-AR (15%)	0.0	12.0	15.3	1.0	13.7	17.9
BERT-AR (30%)	2.3	14.3	18.0	2.0	14.8	19.6
BERT-AR (80%)	0.0	12.2	15.7	0.0	12.9	17.4
Llama-2-7B-CPC	76.2 ± 0.2	91.8 ± 0.8	93.2 ± 0.4	47.5 ± 0.3	83.2 ± 0.6	92.0 ± 0.4
Llama-2-7B-CPD-6L	78.1 ± 0.4	92.2 ± 0.5	94.2 ± 0.6	47.9 ± 0.7	85.4 ± 0.5	93.7 ± 0.4
Llama-3.2-1B-CPC	78.6 ± 0.2	91.5 ± 0.4	92.3 ± 0.2	32.6 ± 0.3	82.5 ± 0.7	89.5 ± 0.3
Llama-3.2-1B-CPD-6L	81.5±0.4	94.0 ± 1.2	95.6 ± 0.5	62.7 ± 0.5	84.9 ± 0.7	87.2 ± 0.9

Table E2: Comparison with the bidirectional training model BERT. To eliminate the problem of random error, we conducted five seed experiments on CPC and CPD, and the experimental results are expressed as $mean \pm standard$ deviation.

expected answer length. We evaluated BERT in two modes: (1) **BERT-parallel**: BERT predicts these masked positions simultaneously; (2) **BERT-AR**: simulating autoregressive generation by predicting tokens sequentially, where each step uses previously generated tokens as context.

Experimental Results The experimental results are displayed in Table E2. The results reveal several important insights: (1) BERT's bidirectional training struggles with the reversal curse: Despite its bidirectional nature, BERT achieves near-zero exact match scores across all masking rates, with the best performance at 30% masking (9.0% EM) still substantially lower than our methods. (2) Masking rate sensitivity: BERT shows optimal performance at 30% masking, suggesting that neither too sparse (15%) nor too dense (80%) masking effectively captures the required associations for this task. (3) Our methods' superiority: Both CPC and CPD significantly outperform BERT across all metrics, demonstrating that position-aware modeling in autoregressive frameworks is more effective than bidirectional attention for addressing permutation sensitivity.

E.5 DOES CPC&CPD TRAINING HURT PERFORMANCE ON STANDARD TASKS?

In our main experiments, we demonstrated that CPC and CPD achieve promising performance on three common failure modes of NTP. A natural concern, however, is that since the pre-training phase does not involve any position-aware training objectives, extensive permutation-based training might risk overfitting to these benchmark datasets of failure modes, potentially leading to catastrophic forgetting. To address this concern, we further investigate whether CPC and CPD disrupt zero-shot performance on eight standard evaluation tasks, including BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020), SIQA (Sap et al., 2019), HellaSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 2020), ARC (easy and challenge) (Clark et al., 2018), OpenBookQA (Mihaylov et al., 2018), and 5-shot aggregated MMLU (Hendrycks et al., 2020) dataset.

E.5.1 Dataset introduction and statistics

In this section, we introduce the datasets used to evaluate the LLMs' zero-shot and 5-shot performance, along with the prompt examples employed in the evaluation. We also present the corresponding dataset statistics in Table E3.

• **BoolQ**⁵ dataset is specifically designed for yes/no question answering tasks. Unlike artificially constructed queries, the questions in BoolQ originate from naturally occurring real-world scenarios, characterized by spontaneity and openness. Each instance in the dataset consists of three components: a question, a corresponding passage, and an answer. In terms of task formulation, the model is presented with a passage and required to answer the

⁵https://huggingface.co/datasets/google/boolq

given question based on that passage, with the answer constrained to either **True** or **False**. Since the test set does not have public answers, we use the validation set for evaluation.

Example E.1: The prompt of BoolQ

instruction:

Please answer the given 'Question' based on the following 'Passage', and only respond with 'True' or 'False'.

input:

1620

1621

1622

1623 1624

1625

1626

1627

1628

1629

1633

1635

1637

1639

1640

1641

1642

1643

1644 1645

1646

1647

1648

1652

1656

1657

1658 1659

1663

1664

1669

1671 1672

1673

Passage:

In mathematics, parity is the property of an integer's inclusion in one of two categories: even or odd. An integer is even if it is evenly divisible by two and odd if it is not even. For example, 6 is even because there is no remainder when dividing it by 2. By contrast, 3, 5, 7, 21 leave a remainder of 1 when divided by 2. Examples of even numbers include -4, 0, 82 and 178. In particular, zero is an even number. Some examples of odd numbers are -5, 3, 29, and 73.

Ouestion:

can an odd number be divided by an even number?

Answer:

• PIQA⁶ dataset is for physical commonsense reasoning. It contains questions about everyday scenarios that require practical knowledge of physical interactions, with answers often favoring unconventional but plausible solutions. In terms of task formulation, PIQA provides a context about a physical situation, and the model is required to choose the correct answer between two candidate solutions (A or B), where only one reflects valid physical commonsense.

Example E.2: The prompt of PIQA

instruction:

Please determine which of the two answers is more accurate and helpful for the following question. You must answer with either 'A' or 'B' only.

input:

Question:

dresser

A. replace drawer with bobby pin

B. finish, woodgrain with bobby pin

Answer:

• **SIQA**⁷(Social IQa) is a benchmark for social commonsense reasoning. Unlike datasets focused on physical or taxonomic knowledge, it centers on understanding people's actions and their social implications. Each instance presents an action and a question with multiple candidate answers (**A**, **B** or **C**), only one of which reflects plausible social reasoning.

Example E.3: The prompt of SIQA

instruction:

You are given a situation, a question, and three possible answers. Choose the best answer that most reasonably and socially fits the situation.

input:

Context:

Sasha protected the patients' rights by making new laws regarding cancer drug trials. Question:

What will patients want to do next?

A. write new laws

B. get petitions signed

C. live longer

Please respond with only the letter of the best answer (A, B, or C).

Answer:

⁶https://huggingface.co/datasets/ybisk/piqa

⁷https://huggingface.co/datasets/allenai/social_i_qa

• HellaS⁸ dataset is for commonsense natural language inference, specifically targeting the ability of models to select the most plausible continuation of a given context. Each instance presents a short context and four candidate endings (A, B, C, or D), only one of which is correct.

Example E.4: The prompt of HellaS

instruction:

You are given a context and four possible endings. Choose the best ending that most reasonably and logically completes the context.

input:

Context:

A boy is running down a track. the boy

A. runs into a car.

B. gets in a mat.

C. lifts his body above the height of a pole.

D. stands on his hands and springs.

Please respond with only the letter of the best answer (A, B, C, or D).

Answer:

• WinoG⁹ dataset is a commonsense reasoning benchmark inspired by the Winograd Schema Challenge, designed to address its limitations in scale and dataset-specific bias. Each instance presents a sentence with a blank and two candidate options (**A** or **B**), only one of which is correct.

Example E.5: The prompt of WinoG

instruction

You are given a sentence with a blank (_) and two possible options. Choose the option that best and most logically fills in the blank.

input:

Sentence:

The doctor diagnosed Justin with bipolar and Robert with anxiety. _ had terrible nerves recently.

A. Justin

B. Robert

Please respond with only the letter of the best answer (A or B).

Answer:

• **ARCe** and **ARCc**¹⁰ are two subsets of the AI2 Reasoning Challenge, a benchmark of grade-school science questions. The **Easy Set** (ARCe) contains questions solvable by simple retrieval or co-occurrence methods, whereas the **Challenge Set** (ARCc) consists of questions that these methods fail to answer, thus requiring deeper reasoning. Each instance is a multiple-choice question with four options (**A**, **B**, **C**, or **D**), only one of which is correct.

⁸https://huggingface.co/datasets/Rowan/hellaswag

⁹https://huggingface.co/datasets/allenai/winogrande

¹⁰https://huggingface.co/datasets/allenai/ai2_arc

Example E.6: The prompt of ARCe and ARCc instruction:

You are given a multiple-choice science question. Choose the best answer based on reasoning and knowledge.

input:

1728

1729 1730

1731

1732

1733

1734

1735

1736

1737

1738

1740

1741 1742

1743

1744

1746 1747

1748

1749

1750

1751

1752

1753

1755

1756

1757

1759

1760

1762

1763

1764

1765

1766 1767

1768

1769

1770

1771

1772

1773

1774

1775

1776

1781

Ouestion:

An astronomer observes that a planet rotates faster after a meteorite impact. Which is the most likely effect of this increase in rotation?

- A. Planetary density will decrease.
- B. Planetary years will become longer.
- C. Planetary days will become shorter.
- D. Planetary gravity will become stronger.

Please respond with only the letter of the best answer (A, B, C, or D).

Answer:

• **OBQA**¹¹ dataset is specifically designed to evaluate advanced question-answering abilities. Unlike simple fact-recall tasks, the questions in OpenBookQA require multi-step reasoning and the integration of both scientific knowledge and common sense. Each instance consists of a science question, several answer choices (**A**, **B**, **C**, or **D**), and access to a set of core science facts (the "open book") provided with the dataset.

Example E.7: The prompt of OBQA

instruction:

You are given a multiple-choice science question. Choose the best answer based on reasoning and knowledge.

input:

Question:

Predators eat

- A. lions
- B. humans
- C. bunnies
- D. grass

Please respond with only the letter of the best answer (A, B, C, or D).

Answer:

• MMLU¹² is a benchmark for evaluating multitask language understanding across a wide range of academic subjects. Each instance is a multiple-choice question with four candidate answers (A, B, C, or D), where the model must identify the correct option by combining world knowledge with reasoning ability. Given the difficulty and diversity of tasks, we randomly sample five validation examples of the same type as few-shot demonstrations when evaluating on the test set.

Example E.8: The prompt of MMLU

instruction:

The following are multiple choice questions (with answers) about {task type}. input:

Question:

Same type of task question 1, answer choice, and the corresponding answer. Same type of task question 2, answer choice, and the corresponding answer. Same type of task question 3, answer choice, and the corresponding answer. Same type of task question 4, answer choice, and the corresponding answer. Same type of task question 5, answer choice, and the corresponding answer. current question and answer choice.

Answer:

¹¹https://huggingface.co/datasets/allenai/openbookqa

¹²https://huggingface.co/datasets/cais/mmlu

Dateset	BoolQ	PIQA	SIQA	HellaS	WinoG	ARCe	ARCc	OBQA	MMLU
Eval Number	3,270	1,838	1,954	10,042	9,248	2,376	1,172	500	14,042

Table E3: Statistics of nine traditional natural language processing evaluation benchmarks.

	Methods	BoolQ	PIQA	SIQA	HellaS	WinoG	ARCe	ARCc	OBQA	MMLU	Avg
Original		63.9	45.5	32.9	25.0	50.3	24.0	22.4	27.6	24.2	35.5
NTP	Standard	47.3	49.5	33.5	25.0	52.3	24.8	22.5	28.0	24.3	33.9
1111	w PAE	9.5	49.5	32.9	25.0	50.3	23.9	22.4	27.6	23.5	29.4
TPM	Standard	0.0	49.4	32.7	25.0	49.7	24.3	22.1	25.6	24.3	28.1
	∦ All	4.2	48.5	33.6	24.6	49.6	22.7	23.5	22.0	25.4	28.2
CPC	M Embedding	0.0	48.4	32.9	24.5	49.8	21.6	23.8	22.0	24.4	27.5
	Transformers	8.8	49.2	32.5	25.3	50.6	24.1	24.6	21.6	23.6	28.9
	♦ All	0.0	48.6	33.1	24.5	50.3	24.0	23.0	25.6	24.2	28.1
	* Transformers	26.0	48.0	32.1	24.7	50.0	23.0	23.4	25.2	25.0	30.8
	Transformers (14-15)	48.6	49.5	33.6	25.0	50.5	24.1	22.6	27.6	25.6	34.1
CPD-6L	Transformers (13-15)	54.0	49.5	33.6	25.0	50.1	25.4	22.4	28.0	25.9	34.9
	Transformers (12-15)	52.7	49.5	33.6	25.0	50.3	24.2	22.4	27.8	25.2	34.5
	Transformers (11-15)	55.2	49.5	33.6	25.0	48.7	26.6	24.2	27.6	25.3	<u>35.1</u>
	Transformers (10-15)	54.2	49.7	33.8	25.1	50.2	24.2	23.5	27.2	25.1	34.8
	Transformers (9-15)	30.2	49.5	33.5	25.2	52.1	27.2	24.2	26.8	24.9	32.6

Table E4: Performance results of various fine-tuned versions of Llama-3.2-1B on standard benchmarks. Here, we investigate which part of the fine-tuned parameters has an impact on the original LLMs' ability. Original denotes the base model. All other models are fine-tuned on the name-to-description dataset. w/ PAE indicates the position-aware embedding introduced during fine-tuning. The $\mbox{\sc XX}$ signifies that only the parameters of component XX in the base model are trained. Transformers (i-j) refers to fine-tuning all Transformer blocks from layer i to layer j. If no specific range is indicated, the fine-tuning is applied to all Transformer layers.

E.5.2 IMPLEMENTATION DETAILS & EXPERIMENTAL RESULTS

Implementation details The proposed position-aware modeling is primarily designed to mitigate common failure modes of standard NTP, rather than to pre-train a LLM from scratch (which we leave for future work). Therefore, when evaluating whether the general performance is affected, we remove the position-aware modules at the testing stage, namely the position embeddings in CPC and the position-aware block layers in CPD. Specifically, for fine-tuned models, NTP and TPM introduce no additional components and can thus be directly evaluated with the fine-tuned model. For NTP (w/ PAE), the position-aware embeddings are incorporated during training but removed during evaluation. Similarly, for CPC and CPD variants, we retain only the original fine-tuned base model structure during evaluation, while the additional position-aware components are excluded.

Experimental Results The experimental results are summarized in Table E4 and Table E5, from which we draw the following conclusions:

- (1) Universality and controllability of catastrophic forgetting. Compared with the performance of the original model (35.5% on average), even standard NTP substantially degrades the general capabilities of the model (33.9% on average), indicating that catastrophic forgetting is a widespread issue. However, our CPD method can effectively mitigate this phenomenon by precisely controlling the degree of base model freezing. Specifically, for the Llama-3.2-1B model with 16 Transformer layers, when fine-tuning only the top few layers (*e.g.*, CPD-Transformers 11–15), the average performance drops by merely 0.4% (from 35.5% to 35.1%), demonstrating the effectiveness of our approach in preserving the model's original capabilities.
- (2) **Impact of coupling vs. decoupling content and position.** CPC introduces position-awareness by directly adding positional embeddings to the original input embeddings. This tight coupling of content and positional information leads to semantic drift in the learned representations. As a result,

			Nan	neIsDesc	cription		DescriptionIsName					
	Method	N2D		D2	D2N		N2D			2N		
			R-1	BLEU	EM	R-1	EM	R-1	BLEU	EM	R-1	
			Lla	ama-3.2-	1B-base							
	Transformers (14-15)	62.7	74.9	77.4	93.3	93.3	49.7	66.1	69.3	100.0	100.0	
	Transformers (13-15)	63.7	76.0	78.4	96.0	96.0	53.0	69.0	72.0	100.0	100.0	
	Transformers (12-15)	64.0	76.3	78.9	96.3	96.3	53.3	79.8	72.9	99.0	99.0	
CPD-6L	Transformers (11-15)	65.3	77.5	79.9	99.7	99.7	54.7	71.7	74.7	99.3	99.3	
	Transformers (10-15)	66.0	78.2	80.7	98.3	98.3	58.9	75.0	77.7	<u>99.7</u>	<u>99.7</u>	
	Transformers (9-15)	<u>70.3</u>	82.4	<u>84.5</u>	100.0	100.0	<u>59.0</u>	<u>75.2</u>	<u>77.9</u>	100.0	100.0	
	🔥 All	81.3	94.7	95.8	100.0	100.0	63.0	85.3	87.7	100.0	100.0	

Table E5: Performance of the CPD variant on the name-description dataset. Complementary to Table E4, the performance of downstream tasks needs to be guaranteed while retaining the performance of the original model.

different CPC configurations (All: 28.2%, Embedding: 27.5%, Transformers: 28.9%) all perform significantly worse than the original model, **underscoring the negative impact of inconsistent paradigms between pre-training and fine-tuning**. In contrast, CPD achieves a modular decoupling of content and positional information through dedicated position-aware blocks, while preserving the structural integrity of the base model. When fine-tuning only a subset of Transformer layers (*e.g.*, CPD-Transformers 11–15: 35.1%), the performance remains nearly identical to that of the original model, validating the advantage of the decoupled design.

- (3) Layer sensitivity and trade-offs in fine-tuning strategies. The results reveal a trade-off between adapting to new tasks and retaining pre-trained knowledge. When all base model parameters are fine-tuned (CPD-All: 28.1%), the model achieves the best performance on position-aware tasks but suffers from a sharp decline in general capabilities due to extensive parameter changes. Interestingly, as more layers are fine-tuned, we observe an improvement rather than a degradation: performance rises from 34.1% with CPD-Transformers (14–15) to 35.1% with CPD-Transformers (11–15). This suggests that moderate parameter fine-tuning, coupled with permutation-invariant training, allows the model to retain pre-trained knowledge while gaining additional position-aware abilities.
- (4) **Task-specific performance preservation.** Table E5 provides deeper insights into how our method maintains performance on the target position-aware tasks while preserving general capabilities. Notably, most CPD configurations show strong performance on the challenging name-description tasks, demonstrating robust position-invariant learning. The CPD-Transformers (11–15) configuration achieve an optimal balance, maintaining strong performance on both forward (N2D: 65.3% EM) and reverse (D2N: 99.7% EM) name-description tasks while achieving the best preservation of general capabilities (35.1% average). This verifies that our framework can both endow the model with permutation invariance and maintain the model's generalization ability, preventing excessive catastrophic forgetting from occurring.

E.6 ABLATION EXPERIMENT

E.6.1 THE NUMBER OF POSITION-AWARE BLOCKS

We conduct comprehensive ablation experiments to investigate the impact of the number of position-aware blocks on model performance in the reversal curse setting. As shown in Table E6, we evaluate CPD architectures with varying numbers of position-aware layers on NameIsDescription (N2D) and DescriptionIsName (D2N) tasks using two base models: Llama-2-7B and Llama-3.2-1B.

Our results reveal several key findings: (1) CPD consistently achieves perfect or near-perfect performance (EM scores of 100.0) on the reversed D2N task across most layer configurations, demonstrating their effectiveness in handling permutation-invariant tasks. (2) We observe a general trend of performance improvement as the number of position-aware layers increases, with the 6-L configuration

Method		Parameter	NameIsDescription					DescriptionIsName					
			N2D			D2N		N2D			D2N		
			EM	R-1	BLEU	EM	R-1	EM	R-1	BLEU	EM	R-1	
				Llama-2-7B-base									
CPC		6.74B	76.3	92.1	93.1	100.0	100.0	47.8	83.5	92.3	100.0	100.0	
CPD	1-L	7.07B	76.5	89.3	93.5	100.0	100.0	46.5	84.9	92.3	100.0	100.0	
	3-L	7.41B	77.2	91.3	93.2	98.3	98.3	47.9	84.2	92.8	<u>99.7</u>	<u>99.7</u>	
	6-L	7.92B	78.3	91.9	94.4	100.0	100.0	<u>48.3</u>	<u>85.7</u>	<u>93.6</u>	100.0	100.0	
	8-L	8.25B	<u>79.2</u>	<u>92.5</u>	<u>95.0</u>	100.0	100.0	47.6	84.3	92.8	100.0	100.0	
	12-L	8.93B	79.9	93.7	96.2	100.0	100.0	52.6	87.3	95.2	100.0	100.0	
			Llama-3.2-1B-base										
CPC		1.23B	78.7	91.8	92.8	82.7	83.6	32.8	82.9	<u>89.7</u>	100.0	100.0	
CPD	1-L	1.57B	79.6	91.9	93.0	86.7	86.7	31.5	83.0	88.6	100.0	100.0	
	3-L	1.68B	80.5	92.2	93.7	99.6	99.6	43.7	83.8	87.9	100.0	100.0	
	6-L	1.86B	81.3	94.7	95.8	100.0	100.0	63.0	85.3	87.7	100.0	100.0	
	8-L	1.98B	<u>81.9</u>	<u>95.3</u>	<u>96.2</u>	100.0	100.0	<u>63.4</u>	<u>85.8</u>	88.1	100.0	100.0	
	12-L	2.21B	82.8	95.9	96.3	100.0	100.0	65.8	87.2	90.1	100.0	100.0	

Table E6: Experimental results on the reversal curse setting. *i*-L denotes the number of position-aware layers, with CPD (6-L) serving as the default configuration throughout all experiments.

emerging as an optimal balance between performance and parameter efficiency. For instance, in the Llama-2-7B CPD model, BLEU scores on N2D improve from 91.3 (3-L) to 91.9 (6-L), while maintaining perfect scores on D2N tasks.

Notably, the performance gains begin to plateau beyond 6 layers, with diminishing returns observed in the 8-L and 12-L configurations. This suggests that 6 position-aware layers provide sufficient capacity to capture the necessary positional relationships for effective permutation-invariant learning. The consistent superiority of the 6-L configuration across both model sizes and task directions validates our choice of CPD (6-L) as the default setting throughout our experiments.

E.6.2 Whether to train the base AR model in CPD

In CPD, we append multiple layers of our proposed position-aware blocks after the output layer of the existing base AR model, effectively decoupling the target position and content representations, with target positions serving as query vectors. A natural question arises: can we train only the position-aware blocks while keeping the parameters of the base AR model fixed? To investigate this, we conducted comparative experiments on the name-description dataset using Llama-2-7B, with results presented in Table E7. The following conclusions can be drawn: (1) **Frozen ALL** (training only position-aware blocks while completely freezing base AR model parameters) demonstrates significantly degraded performance. On the NameIsDescription N2D task, performance drops precipitously from 79.9 (EM) and 93.7 (R-1) for CPD-12L to 48.7 (EM) and 72.4 (R-1). More severely, on the DescriptionIsName N2D task, performance almost completely collapses, declining from 52.6 (EM) and 87.3 (R-1) to merely 3.3 (EM) and 27.8 (R-1). This substantial performance deterioration primarily occurs because knowledge-related content representations are predominantly stored within the base AR model. When these parameters are frozen, the model cannot adjust its internal knowledge representations to accommodate the position-aware mechanism. Although position-aware blocks can theoretically store some knowledge information, their design primarily

		Parameter	NameIsDescription					DescriptionIsName				
	Method		N2D			D2N		N2D		D2N		
			EM	R-1	BLEU	EM	R-1	EM R-1	BLEU	EM	R-1	
	CPC	6.74B	<u>76.3</u>	92.1	<u>93.1</u>	100.0	100.0	<u>47.8</u> <u>83.</u>	<u>5 92.3</u>	100.0	100.0	
CPD	12-L	8.93B	79.9	93.7	96.2	100.0	100.0	52.687.	3 95.2	100.0	100.0	
	Frozen ALL	2.32B	48.7	72.4	76.3	28.3	29.6	3.3 27.	8 33.0	99.7	99.7	
	Frozen Embedding	8.93B	73.0	87.9	90.9	98.3	98.3	47.3 75.	7 79.6	99.0	99.0	

Table E7: Experimental results on the reversal curse setting with Llama-2-7B. *i*-L denotes the number of position-aware layers, **Frozen ALL** means freeze all parameters of the base AR model, and **Frozen Embedding** represents only freezing the parameters of the embedding layer in the base AR model.

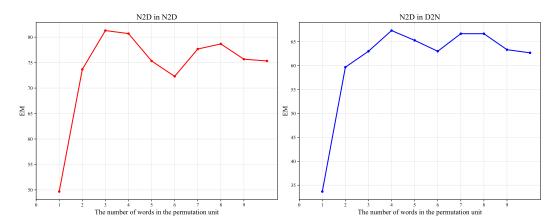


Figure E5: Effect of permutation unit size on reversal curse performance using Llama-3.2-1B-CPD (6-L). EM scores are shown for different word-level permutation unit sizes on the name-description dataset. Left: N2D task performance in NameIsDescription setting. Right: N2D task performance in DescriptionIsName setting.

focuses on processing positional information rather than content representation, resulting in limited knowledge storage capacity.

(2) In contrast, **Frozen Embedding** (freezing only the embedding layer while allowing updates to other parameters) exhibits performance more closely approximating the fully fine-tuned model. On the NameIsDescription task, this strategy achieves 73.0 (EM) and 87.9 (R-1), which, while slightly lower than the fully fine-tuned CPD-12L, significantly outperforms the **Frozen ALL**. On the DescriptionIsName task, **Frozen Embedding** approaches the performance of the fully fine-tuned model, with nearly identical results on the D2N task (98.3 vs. 100.0).

These results indicate that updating base AR model parameters (particularly parameters beyond the embedding layer) during training is crucial for effectively integrating positional information and content representations.

E.6.3 THE UNIT OF PERMUTATION

To confirm the impact of permutation unit granularity on model performance, we conducted experiments on permutation unit granularity under the reversal curse setting, and the results are shown in Figure E5. We can draw the following conclusions: (1) Both small and large permutation units are detrimental to model performance. When permutation units are too small (*e.g.*, 1-2 words), the model is forced to learn fragmented representations of common linguistic phrases and fixed collocations, which imposes an additional learning burden and disrupts the natural semantic coherence of language constructs. Conversely, when permutation units are too large (*e.g.*, 7+ words), the model cannot effectively perceive and adapt to different degrees of contextual variations, as the permutation granu-

larity becomes too coarse to provide meaningful positional diversity during training. (2) The results reveal that different task exhibit distinct optimal permutation unit sizes. For the N2D task within the NameIsDescription setting, peak performance is achieved around 3-4 words per permutation unit, while the N2D task within the DescriptionIsName setting shows optimal performance around 4-5 words per unit. This suggests that the complexity and structure of the underlying task influence the most effective permutation granularity. (3) The consistent decline in performance at both extremes suggests that maintaining an appropriate balance between providing positional diversity and preserving semantic coherence is essential for effective permutation-based training.