

ELIMINATING POSITION BIAS OF LANGUAGE MODELS: A MECHANISTIC APPROACH

Anonymous authors

Paper under double-blind review

ABSTRACT

Position bias has proven to be a prevalent issue of modern language models (LMs), where the models prioritize content based on its position within the given context. This bias often leads to unexpected model failures and hurts performance, robustness, and reliability across various applications. A simple mechanistic analysis attributes the position bias to two components employed in nearly all state-of-the-art LMs: causal attention and position embedding. Based on the analyses, we propose to **eliminate** position bias (e.g., different retrieved documents' orders in QA affect performance) with a **training-free zero-shot** approach. Our method changes the causal attention to bidirectional attention between documents and utilizes model attention values to decide the relative orders of documents instead of using the order provided in input prompts, therefore enabling **Position-INvariant inferencE (PINE)** at the document level. By eliminating position bias, models achieve better performance and reliability in downstream tasks, including LM-as-a-judge, retrieval-augmented QA, molecule generation, and math reasoning. Notably, PINE is especially useful when adapting LMs for evaluating reasoning pairs: it consistently provides 8 to 10 percentage points performance gains, making Llama-3-70B-Instruct perform even better than GPT-4o-125-preview and GPT-4o-2024-08-06 on the RewardBench reasoning set.

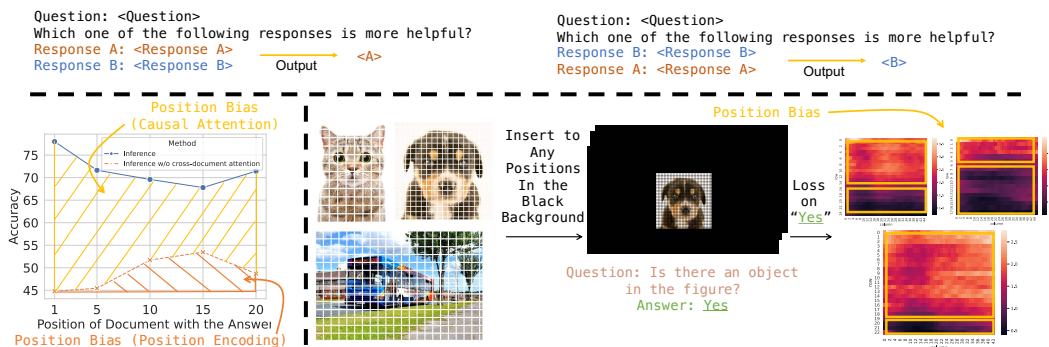


Figure 1: Motivating examples showing how position bias affects model outputs. **Upper:** LMs are asked to select a more helpful one from two given responses. The example shows that LMs are prone to prefer the response positioned at first. **Lower Left:** LMs (Llama-3-8B-Instruct) are presented with 20 documents to answer a question, with only one document (the gold-standard document) containing the correct answer. The blue curve represents normal inference, while the red curve represents inference without inter-document attention (RoPE position encodings are kept, a concrete implementation is shown in the middle of Figure 3). The height change of the yellow and orange area reflects the position bias brought by causal attention and RoPE: causal attention generally favors distant content, but RoPE prefers nearby content. **Lower Right:** We insert a real-world image to a large black background image at different positions and prompt VLMs (Fuyu-8B (Bavishi et al., 2023)) to compute the loss on the ground truth token. We observe a consistent pattern that models have lower losses (black color) when images are presented at the bottom.

1 INTRODUCTION

Language models (LMs) (Brown et al., 2020; Chowdhery et al., 2022; Touvron et al., 2023; Achiam et al., 2023) demonstrate impressive performance in general language tasks such as dialogue (Thoppilan et al., 2022), reasoning (Chowdhery et al., 2022), and schema induction Li et al. (2023). However, they tend to favor content at certain positions (Zheng et al., 2024b;a; Wang et al., 2023; Dominguez-Olmedo et al., 2023; Zhu et al., 2023; Chen et al., 2024b; Liu et al., 2024), which harms complex reasoning (Chen et al., 2024b), long-context understanding (Liu et al., 2024) and model-based evaluation (Zheng et al., 2024b). For example, LMs tend to favor the first when it is required to compare the quality of two candidate responses (Zheng et al., 2024b), which hurts their reliability when being used as evaluators (Figure 1 upper); vision-language models perform better in the recognition when the target content is presented at the bottom of the image (Figure 1 lower right, see more examples in Appendix A). Different from *ad hoc* solutions from previous works (Ratner et al., 2023; Cai et al., 2023; Hao et al., 2022; Junqing et al., 2023; Zhu et al., 2023), we seek to understand the causes of position bias and propose to eliminate (not just mitigate) the position bias without training and searching.

We start by analyzing the key components of state-of-the-art LMs – Casual Attention and Position Embedding. They are the key to the success of Transformers (Vaswani et al., 2017), and are also the only two operations in Transformers (Vaswani et al., 2017) that will bring undesirable position bias. This is because other operations do not change representations when position changes (Section 3.2). Moreover, we find an interesting phenomenon through simple experiments and give our hypothesis: the most popular Rotary Position Embedding (Su et al., 2024) is shown to have recency bias (Su et al., 2024; Peysakhovich & Lerer, 2023) due to its long-form attention weight decay w.r.t. the increase of relative positions, and the causal attention forces unidirectional information propagation, enabling models to pay more attention to distant content. Figure 1 lower left shows this retrieval-augmented QA (Liu et al., 2024) experiment. The height change of the yellow area and orange area reflects the position bias of causal attention and RoPE. Since the yellow area is mostly wider at the beginning and the orange area generally becomes wider at the end (except for the last data point), this shows that the causal attention generally tends to favor distant content, while RoPE generally tends to favor nearby content.¹

Since attention and position embedding are the causes of position bias, we propose PINE that can eliminate position bias by manipulating causal attention and RoPE to attend to different content equally. Take the retrieval augmented QA (Liu et al., 2024), a task requiring LMs to answer questions based on retrieved documents, for example. The orders of retrieved documents should not affect the final results. To achieve this, we make the inter-document attention bidirectional so that the attention mask will equally attend to all documents. Next, we compute importance scores between documents and use them to re-assign document positions so that positions in the original inputs are discarded. We prove resulting approach enables **Position-invariant inference** (PINE) w.r.t. documents in a **training-free/zero-shot** manner.

To justify the effectiveness of PINE, we select four tasks: LM-as-a-judge (Zheng et al., 2024b), which prompts LMs to choose the more helpful one from two given responses to a question; retrieval-augmented question-answering (Liu et al., 2024); molecule generation based on provided properties; and math reasoning. In different tasks, “documents” have different meanings: responses in LM-as-a-judge, properties in molecule generation, and conditions in math reasoning. Notably, we find our method especially useful when LMs are used to assess reasoning pairs: PINE with Llama-3-70B-Instruct perform even better than GPT-4-0125-preview and GPT-4o-2024-08-06 on the RewardBench (Lambert et al., 2024a) reasoning set.

To summarize, we:

- We first **revisit** the causes of position bias in transformers: causal attention and position encoding (Section 3.2), and then propose a training-free approach dubbed PINE that can eliminate (with proof) the position bias given documents presumed to be position-invariant (Section 3.3).
- Four popular tasks across the general domain to expert domains (chemistry and math) show PINE can bring performance gains consistently across different models and sizes.

¹More supporting experiments to this hypothesis in Section 4.3.

2 RELATED WORK

Position Bias in LMs. Position bias widely exists in LMs (Zheng et al., 2024b;a; Wang et al., 2023; Zhu et al., 2023; Chen et al., 2024b; Liu et al., 2024; Shi et al., 2024). The LM-as-a-judge task offers models two candidate responses to a question and asks models to select the more helpful one. It turns out that LM has a primacy bias that tends to favor the first response (Zheng et al., 2024b). Retrieval-augmented QA asks LM to answer a question based on retrieved documents. Liu et al. (2024); Peysakhovich & Lerer (2023) find that LMs are prone to answer correctly when the document that contains the correct answer is presented at the beginning and the end of retrieved documents. Zheng et al. (2024a) points out that models favor options at certain positions (e.g., prefer “A”) in multiple-choice QA. In the in-context learning task, Zhang et al. (2024a); Xu et al. (2024) find that the order of in-context examples affects the final performance. Recently, several papers have proposed to understand the nature of position bias through prompting (Zhang et al., 2024b) and calibration (Hsieh et al., 2024). Our paper analyzes the phenomenon from the mechanical perspective: the computation must be positional-invariant to eliminate position bias.

Position Bias Solutions in LMs. *Mitigating* position bias have been studied by many literature from many aspects, such as data augmentation with training (Junqing et al., 2023; Zhu et al., 2023), content sorting by attention value during inference (Peysakhovich & Lerer, 2023), searching (Yu et al., 2024; Adila et al., 2024), calibration (Hsieh et al., 2024), **or removing position encoding (Kazemnejad et al., 2024)**. Moving one step forward, some other solutions are designed to *eliminate* position bias. Zheng et al. (2024a;b) use permutation then average on classification tasks, which will have unacceptable $\mathcal{O}(k!)$ (k is the number of segments) computational overhead when k is large. Hsieh et al. (2024) assumes that the position bias and real relevance are linear combinations and propose solutions accordingly. Different from them, we aim to *eliminate* the position bias from the mechanical perspective without any assumption at a reasonable cost. Although several existing approaches are from the mechanical perspective (Ratner et al., 2023; Cai et al., 2023; Hao et al., 2022), they only perform well in classification tasks and fail in a more general setting: language generation.

3 METHODOLOGY

We start by running an example to illustrate position bias, followed by analyzing the cause of position bias, and end with our own approach PINE.

3.1 FORMULATION

We take retrieval-augmented QA as an example, where current LMs’ performance may greatly suffer from position bias (Liu et al., 2024). The task requires the model to answer a question based on a set of given retrieved documents, where only one of them contains the correct answer. The system prompt SYS is: “Write a high-quality one-sentence answer for the given question using only the provided search results (some of which might be irrelevant).”. Given a question Q, and three retrieved documents: \mathcal{D}_1 , \mathcal{D}_2 , and \mathcal{D}_3 , we can formulate several different inputs. For example, $[\text{SYS}|\text{Q}|\mathcal{D}_1|\mathcal{D}_2|\mathcal{D}_3]$, and $[\text{SYS}|\text{Q}|\mathcal{D}_2|\mathcal{D}_3|\mathcal{D}_1]$. We expect models to have the same output for these inputs because $\mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_1$ are **position-agnostic documents**: their relative order is not supposed to affect the final result. However, the current LMs answer differently when presented with these different inputs and tend to answer correctly when the document contains the answer at the beginning or at the end of all documents (Liu et al., 2024). The systematic differences of model outputs caused by relative positions of documents reflect the **position bias** of the model. Therefore, current LMs cannot conduct **inter-document position-invariant** inference, and our goal is to make the inference invariant w.r.t. relative document orders. In the rest of this section, we will use this running example. However, we emphasize that “documents” have different meanings in different tasks: responses in LM-as-a-judge, properties in molecule generation, and conditions in math reasoning. Therefore, readers should be aware that the method is not just designed for a single task. In the rest of the paper, we merge SYS and Q into SYS for simplicity.

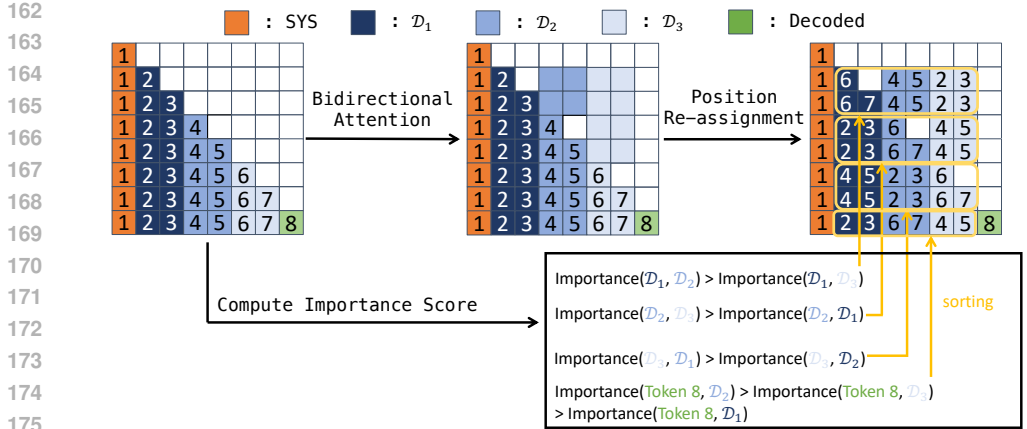


Figure 2: PINE: inter-document position-invariant inference via bidirectional attention. The attention matrix of the running example in Section 3.1 is at the left of the figure, the orange, different blue, and green colors denote system prompts (1 token), three different documents (2 tokens each) and decoded tokens (1 token), respectively. The number at (i, j) in the figure, p_{ij} , denotes the position of a token j when computing the attention from query \mathbf{q}_i . Therefore, $p_{.j}$ is equal for all i in vanilla inference. PINE enables inter-document bidirectional attention and then uses attention scores between documents to compute their importance. Then, documents are sorted by importances: more important documents are placed in closer positions. The computation of “importance score” is introduced in Section 3.3.

3.2 CAUSAL ATTENTION AND POSITION EMBEDDING ARE THE CAUSE OF POSITION BIAS

Feed-forward networks (FFNs), Query, Key, and Value (QKV) projections, and layer normalization in the Transformer architecture do not cause position bias, as they yield the same representations regardless of document positions. Rather, the attention computation that leads to the position bias:

$$\mathbf{Q}_{PE} = PE(\mathbf{Q}, \text{pos}_{\mathbf{Q}}), \mathbf{K}_{PE} = PE(\mathbf{K}, \text{pos}_{\mathbf{K}})$$

$$\mathbf{H} = \text{Softmax} \left(\mathbf{Q}_{PE} \mathbf{K}_{PE}^T / \sqrt{d} \right) \odot \mathbb{1}_{\text{causal}} \mathbf{V} \tag{1}$$

where $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{n \times d}$ are queries, keys, and values, PE denotes the position encoding, $\text{pos}_{\mathbf{Q}}$ and $\text{pos}_{\mathbf{K}}$ denote the position of queries and keys, and $\mathbb{1}_{\text{causal}}$ denotes the causal attention mask. Eq. 1 reveals that (1) the PE function yields different representations for documents if their orders changes, therefore affecting the importance score $Q_{PE} K_{PE}^T$ and hidden states; (2) the $\mathbb{1}_{\text{causal}}$ generates different attention masks for the input documents if we change their positions, resulting in different hidden states. To achieve inter-document position-invariant inference, **H must remain the same regardless of documents’ orders.**

3.3 PINE: INTER-DOCUMENT POSITION-INVARIANT INFERENCE VIA BIDIRECTIONAL ATTENTION

Our goal is to obtain an inter-document position-invariant hidden state \mathbf{H}_{PINE} , which does not change regardless of document orders. We can mechanistically eliminate the position bias by equally attending to all documents. Therefore, we propose PINE, an approach that uses bidirectional inter-segment attention and re-assigning positions by importance scores (computed from attention score) to eliminate position bias (Figure 2). We address that the “elimination” and “invariance” in our method are talked about from the input-output perspective, i.e., outputs remain unchanged regardless of the input-position orders. PINE still uses position encoding and does not eliminate position encoding itself.

Bidirectional Attention. We first change the attention mask so that documents can attend to each other. Specifically, we make the inter-document attention **bidirectional** but keep the intra-document attention **causal** (Figure 2, middle). Our goal is to eliminate “inter” position bias among different documents rather than “intra” position bias within each document. The latter will lose the order

information of tokens, and models can degenerate into bag-of-words models, which is not what we expect.

Re-assign Positions: Sorting By Importance Scores. Re-assigning positions must consider two folds: the position of queries and keys. Each token in conventional LMs has the same position when serving as both query and key. In the bidirectional attention we use, this assignment has to be reconsidered. First, LMs are trained causally, meaning the position of the query must be larger than the keys in the attention computation. Therefore, it is necessary to manipulate positions so that each document is the **last** document when serving as queries (the diagonal of the rightmost figure in Figure 2). For tokens before and after documents, their positions are not affected when serving as queries.

Re-assigning positions for keys must be redesigned to eliminate position bias. We determine the positions of documents based on importance scores when they serve as keys (numbers in the rightmost part of Figure 2). Specifically, we first compute the attentions without position embedding involved: $\text{Importance}_{\text{token}}(i, j) = \text{Softmax}(\mathbf{q}_i \mathbf{k}_j^T / \sqrt{d})$, where d is the hidden state dimension. Then, we obtain the importance score between documents by aggregation. For example, $\text{Importance}(\mathcal{D}_1, \mathcal{D}_2) = \sum_{i \in \mathcal{D}_1, j \in \mathcal{D}_2} \text{Importance}_{\text{token}}(i, j) / |\mathcal{D}_2|$. The length normalization is to prevent assigning higher importance scores to longer documents.² The importance score could also be computed between individual tokens (e.g., Token 8) and documents. Lastly, we re-assign positions by importance scores as shown in the rightmost part of Figure 2: more important documents will have closer positions to the query. The rightmost part of Figure 2 shows the concrete position re-assignment for keys (its diagonal also represents the position re-assignment for queries). To avoid confusion, we address the fact that we do not actually sort tokens and only re-assign them to different positions. In our position re-assignment, the position of keys may vary depending on the queries (numbers in column are different), which is the key difference between PINE and vanilla inference. Besides, our method is not limited to specific position embedding types.

Inter-Document Position Invariant Inference. Once we have new attention mask and position re-assignment, we can place them into Equation 1, and obtain \mathbf{H}_{PINE} . By applying \mathbf{H}_{PINE} to every layer, attention heads, and tokens, we reach our method PINE. We prove that:

Lemma 1. *If the input $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ are inter-document position-invariant representations, then \mathbf{H}_{PINE} are also inter-document position-invariant representations.*

Proof: To simplify the notation and without loss of generality (w.l.o.g), we still use examples in Section 3.1.

First, the SYS tokens already satisfy this lemma under the vanilla inference since they appear before documents, and PINE does not change their computation process. We only need to show PINE can make \mathcal{D}_i and Token 8 (i.e., tokens after documents) satisfy the lemma. W.l.o.g, we use \mathcal{D}_1 as a running example:

- PINE first obtains importance score between documents: $\text{Sim}(\mathcal{D}_1, \mathcal{D}_i) = \sum \text{Softmax}(\mathbf{Q}_1 \mathbf{K}_i^T / \sqrt{(d)}) / |\mathcal{D}_i|$, where $\mathbf{Q}_1 \in \mathbb{R}^{2 \times d}$, $\mathbf{K}_i \in \mathbb{R}^{2 \times d}$, 2 denotes the number of tokens in documents, and d denotes hidden states dimensions. Note that here the \mathbf{Q}, \mathbf{K} have not been applied to position embedding yet. Therefore, the importance score is not a function of input document positions.
- W.l.o.g, let's assume $\text{Sim}(\mathcal{D}_1, \mathcal{D}_2) > \text{Sim}(\mathcal{D}_1, \mathcal{D}_3)$, then we sort the document as follows $[\mathcal{D}_3 | \mathcal{D}_2 | \mathcal{D}_1]$ when they serve as keys and \mathcal{D}_1 as query. Concretely, $\mathbf{Q}_{\text{PE},1} = \text{PE}(\mathbf{Q}_1, 3)$ (3 denotes it is treated as the last, i.e., third, document), $\mathbf{K}_{\text{PE},1} = \text{PE}(\mathbf{K}_1, 3)$, $\mathbf{K}_{\text{PE},2} = \text{PE}(\mathbf{K}_2, 2)$, $\mathbf{K}_{\text{PE},3} = \text{PE}(\mathbf{K}_3, 1)$. Then we compute hidden states of \mathcal{D}_1 : $\mathbf{H}_1 = \text{Softmax}(\mathbf{Q}_{\text{PE},1} \mathbf{K}_{\text{PE}} / \sqrt{(d)})$, where \mathbf{K}_{PE} is the key values for the whole sequence $[\text{SYS} | \mathcal{D}_3 | \mathcal{D}_2 | \mathcal{D}_1]$. It is noted that this process does not use any variables that are dependent on the input document positions, nor directly use the input document positions. Therefore, \mathbf{H}_1 obtained by PINE is not a function of input document positions.
- Similarly, $\mathbf{H}_2, \mathbf{H}_3$, and Token 8's hidden states are not functions of input document positions. Their concatenation yields \mathbf{H}_{PINE} , which is not a function of input document positions.

²In our pilot experiments, we find summation makes models convert from position bias to length bias. We also try maximum instead of averaging and find this methods usually have noticeably worse performance than averaging possibly due to noises brought by unimportant tokens.

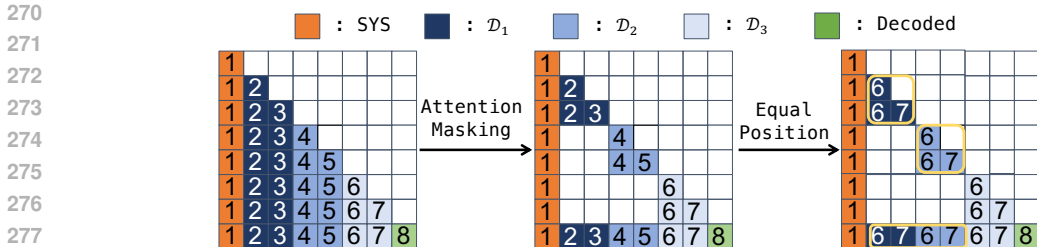


Figure 3: Previous work PCW (Ratner et al., 2023) eliminates position bias by first masking all inter-document attention and then re-assigning all documents the same positions. The notions are kept the same as Figure 2. Our experiment in Section 4 shows that PCW brings severe performance drop for tasks requiring language generation.

Proof ends.

Theorem 1. *Given an input, if H_{PINE} is applied to every layer, attention head, and token to replace the conventional attention computation, then the model outputs are inter-document position-invariant representations.*

The theorem can be proved by mathematical induction by (1) lemma, (2) FFN, QKV projection, and layer norm yield representations that are not a function of document positions, and (3) the embedding representation is not a function of document positions. We put the complete proof in Appendix B.

Some takeaways that are worth noting: (1) Both bidirectional attention mask and position re-assignment are needed to complete the proof. (2) PINE needs to be applied to every layer, attention heads, and tokens to complete the proof. (3) PINE is not limited to specific position embedding types.

3.4 DISCUSSION

Different Position Re-Assignment Methods. PINE puts documents with higher importance scores to a closer position to queries. Another option is to put documents with higher importance scores in a more distant position to the queries. Considering the recency bias brought by the most popular rotary position embedding (RoPE) (Su et al., 2024), this alternative approach makes RoPE “disrespect” the attention of models. Therefore, we believe this alternative choice is not optimal, which is justified by our experiments in Section 4.3.

Different Attention Masks. Previous work PCW (Ratner et al., 2023) adopts a different way: it masks the inter-document attention instead of making it bidirectional (Figure 3, middle and right). Accordingly, it adopts a simplified position re-assignment method of ours: putting all documents in the same positions. However, masking all inter-document attention loses contextual information (the white part surrounded by colored blocks in Figure 3). Moreover, some different tokens now share the same positions (Figure 3, right), which could confuse models. As a result, PCW performs poorly in language generation tasks (Section 4).

Inference Cost. PINE incurs additional computation overhead due to extra operations. Practically, the extra big \mathcal{O} computation complexity to obtain hidden states is $\mathcal{O}(nk \log k)$, where n and k denote text length and the number of input documents, respectively. The bidirectional attention does not bring extra cost, the position re-assignment brings $\mathcal{O}(k \log k)$ for each token since the sorting algorithms are involved. The real computation cost is acceptable since k is usually small (e.g., $k = 2$ in the LLM-as-a-judge task and $k = 20$ in the retrieval-augmented QA). Section 4.5 shows results of real-world wall time and memory cost.

4 EXPERIMENT

Our experiments aim to show PINE can improve model performance across diverse tasks and have superior performance than other approaches.

Table 1: The portion of data (%) that models have position bias in RewardBench, i.e., models change answers after swapping candidate responses orders. We color the subsets that have more than 25% data causing position bias with cyan.

Model	Size	Chat	Chat-Hard	Safety	Reasoning	Avg.	
LLaMa-3	8B	10.3	21.5	11.4	27.6	17.7	
	-Instruct	70B	3.6	16.0	5.8	15.2	10.2
Qwen-1.5	1.8B	33.5	37.9	24.7	13.3	27.4	
	4B	48.0	38.6	57.4	12.7	39.2	
	7B	17.0	20.6	10.9	26.5	18.8	
	-Chat	32B	7.8	20.0	9.6	26.4	16.0
	72B	10.9	22.6	9.6	24.7	17.0	
	110B	8.7	16.0	11.5	23.5	14.9	

4.1 SETTINGS

We select four tasks that pose position bias: LM-as-a-judge (Zheng et al., 2024b) that prompts LMs to select a better response out of two given a question, retrieval-augmented question-answering (Liu et al., 2024) that asks LMs to answer questions based on retrieved documents, molecule generation based on provided properties (Ramakrishnan et al., 2014), and math reasoning based on several given conditions Chen et al. (2024b). We follow previous work (Liu et al., 2024; Lambert et al., 2024a) and use temperature 0 in avoid variance.

LM-as-a-judge. We benchmark our method on 23 datasets in the RewardBench³ (Lambert et al., 2024b) that can be categorized into four types: Chat, Chat-Hard, Safety, and Reasoning. We use the official data split, prompts, and evaluation scripts to ensure reproducibility. We use LLaMa-3-Instruct models (AI, 2024) and Qwen-1.5-Chat models (Bai et al., 2023) for experiments. To show how positions affect results, we present four results: the ground-truth response is positioned at first, second, or shuffled, and PINE results (which yield the same results for all three scenarios above).

Retrieval-augmented QA. We follow the settings and use the prompts, data, and evaluation scripts of (Liu et al., 2024)⁴: Only one of the retrieved documents (10 or 20 in total) contains the ground-truth answer for the given question. We list prompts in Appendix D. We use LLaMa-3-70B-Instruct model (AI, 2024) for experiment. To show how positions affect results, we present several results: the ground-truth document is positioned at the beginning, middle, last, or shuffled, and PINE results (which yield the same results for all scenarios above).

Molecule Generation and Math reasoning. We also conduct two bonus experiments. Molecule generation based on given properties (property positions can be swapped), and math reasoning where conditions can be swapped.

More details of the four tasks can be found in Appendix D. Qualitative examples of the four tasks can be found in Appendix E.

Baselines. The goal of PINE is to eliminate position bias during inference mechanically. Therefore, we choose methods that have the same design principle as our baselines: (1) Vanilla inference (2) Vanilla inference with no inter-document attention (NIA for short, i.e., the middle figure in Figure. 3): The latter documents will have no attention to formers. (3) Parallel Context Window (PCW, rightmost in Figure. 3) (Ratner et al., 2023): PCW extends the baseline (2) by manipulating positions of documents. PCW allows all documents to share the same positions. (4) Structured Prompting (SP, a variant version of PCW) Hao et al. (2022): SP extends (3) by lowering attentions between decoded tokens and input documents to $\frac{1}{N}$ to solve the perplexity exploding problem in PCW. Similar to the proof in Section 3.3, we can know that (1) and (2) are not inter-document position invariant, whereas (3) and (4) are. Beyond these methods, we also introduce two other debiasing baselines: permutation (Zheng et al., 2024a) and calibration (Zhao et al., 2021).

³Apache-2.0 license. <https://github.com/allenai/reward-bench>

⁴MIT license. <https://github.com/nelson-liu/lost-in-the-middle>

Table 2: Main results of RewardBench. Vanilla denotes the normal inference, (GT at A) means the ground truth chosen response is presented at the first, and (GT at B) indicates the second. For the 72B model, we additionally benchmark the Qwen 2.5 model. PINE consistently improves LM’s performance across different models and sizes and is particularly useful when assessing reasoning pairs.

Method	Llama-3-Instruct		1.8B	4B	7B	Qwen-1.5-Chat		
	8B	70B				32B	72B / 72B (Qwen 2.5)	110B
RewardBench (Full set)								
Vanilla (GT at A)	67.5	78.0	36.3	29.5	61.4	74.2	79.6 / 87.2	87.2
Vanilla (GT at B)	66.3	76.5	66.2	76.6	59.6	74.8	69.5 / 80.5	75.7
Vanilla (Shuffle)	64.8	76.0	50.3	53.1	60.9	72.8	72.8 / 83.4	81.1
PINE	66.7_{+1.9}	77.4_{+1.4}	52.9_{+2.6}	58.2_{+5.1}	61.5_{+0.6}	74.8_{+2.0}	71.8_{-1.1} / 84.5_{+1.1}	82.9_{+1.7}
RewardBench (Reasoning set)								
Vanilla (GT at A)	80.3	87.8	43.3	42.8	62.1	78.3	83.0 / 93.7	90.0
Vanilla (GT at B)	66.0	80.3	57.2	62.3	54.3	73.6	68.7 / 76.0	73.0
Vanilla (Shuffle)	65.3	78.9	48.4	54.1	59.3	66.8	68.2 / 85.5	78.0
PINE	73.4_{+8.1}	87.6_{+8.7}	60.1_{+11.7}	61.0_{+6.9}	63.0_{+3.7}	76.7_{+9.9}	69.0_{+0.8} / 91.3_{+5.8}	86.2_{+8.2}

Table 3: Baseline performance on RewardBench. PINE achieves superior performance to baseline models, performing 4.8% and 4.7% better than the best performed baseline on two models.

Method	LLaMa-3-8B-Instruct		Qwen1.5-7B-Chat	
	Reasoning	Full Set	Reasoning	Full Set
NIA (GT at A)	43.7	56.3	60.7	61.3
NIA (GT at B)	66.7	65.8	44.1	52.2
NIA	55.9	61.9	51.4	56.8
PCW	56.5	61.7	53.4	55.2
SP	55.4	60.8	52.4	55.4
PINE	73.4_{+16.9}	66.7_{+4.8}	63.0_{+9.6}	61.5_{+4.7}

4.2 RESULTS ON LM-AS-A-JUDGE

Position bias exists across different models and sizes. Table 1 shows the statistics of position bias in RewardBench with different models. Position bias is quite common in RewardBench, and can be up to 48.0%. Larger models have less position bias, however, the position bias could still on average affect up to 10% data.

PINE consistently improve model performance across models and sizes. Table 2 shows the main results on RewardBench. We experiment with Llama-3 and Qwen-1.5 across different model sizes. The position of the ground truth chosen option is randomly shuffled. Therefore, the accuracy of the random guess method is expected to be 50%. First, the first two rows reveal that larger models tend to have a primacy bias, whereas smaller models tend to have a recency bias. By comparing the last two rows of each model size, we conclude that models across different sizes perform better with the help of PINE by eliminating position bias. The only exception is the Qwen-1.5-72B-Chat model. We suspect this model is not well-trained since Qwen-1.5-32B-Chat performs the same as the 72B model in vanilla inference, despite half of the model size. Qwen 2 report (Yang et al., 2024) also shows that the Qwen 1.5B 72B model performs even worse than 32B in reasoning. Moreover, Table 2 shows that Qwen 2.5 72B can obtain consistent performance gains. Overall, PINE improves performance from a statistical perspective and makes models more reliable when as evaluators. Full results are shown in Appendix C.

PINE is extremely useful when assessing reasoning problems in RewardBench. PINE consistently improves model performance on the “reasoning” subset by a large margin: from 8 to 10 percentage points in most cases. Specifically, LLaMa-3 Instruct 70B was originally ranked 22nd generative model in the reasoning subset of RewardBench. With PINE, it achieves the 7th rank (87.6%), **outperforming GPT-4-0125-preview (the previous 8th rank, 86.9%), GPT-4o-2024-08-06 (the previous 9th rank, 86.6%), and Llama-3.1-405B-Instruct-Turbo (the previous 7th rank, 87.1%).**⁵.

⁵Results are provided by the official leaderboard (as of Sep 17, 2024): <https://huggingface.co/spaces/allenai/reward-bench>

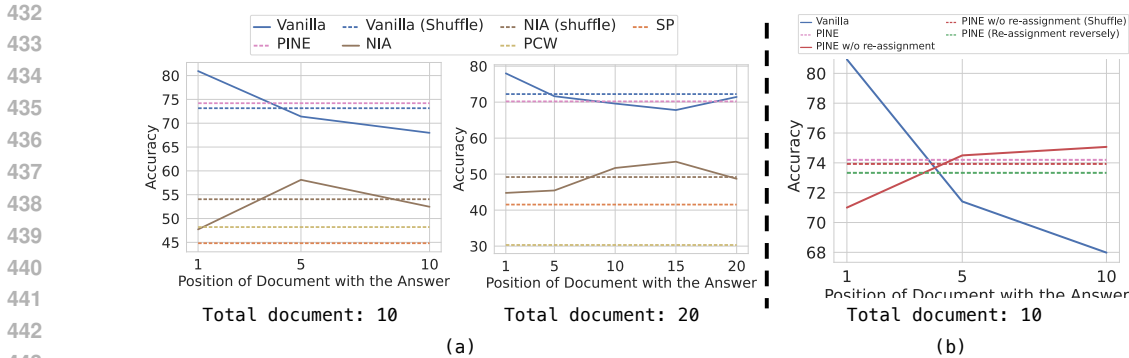


Figure 4: The results of retrieval-augmented QA on Llama-3-70B-Instruct. Dashed lines indicate that the method is either inter-document position-invariant or the result is obtained on the order-shuffled data (denoted in the legend). (a) shows results of PINE against baselines. (b) shows results of different designs of PINE.

PINE performs better than baseline models that adopt different attention masks. We then compare PINE with baseline models mentioned in Section 4.1 on Llama-3-8B-Instruct and Qwen1.5-7B-Chat model. They adopt a different attention mask: masking inter-document attention instead of making them bi-directional. Since NIA is not inter-document position-invariant, we also apply NIA with two extreme cases: the ground truth chosen response is always in the first or second place. Results on Table 3 show that PINE achieves the best performance and largely outperforms the best baselines by $\sim 5\%$, and outperforms NIA even if NIA is placed in the extreme case. On the reasoning subset, this performance gap becomes much even greater. The results reveal that masking inter-document attention mask is much less effective than bidirectional inter-document attention mask applied in PINE.

PINE performs better than permutation and calibration methods. Another two widely used debiasing methods are permutation (Zheng et al., 2024a) and calibration (Zhao et al., 2021). They are usually used in the logit-based evaluation or single-token generation. Their effectiveness in the open-ended generation is less explored. In our experiments, we find calibration methods generates rubbish responses, which we believe is because of the strong assumption in (Zhao et al., 2021): uniform distribution of all tokens in the generation task. For the permutation methods, we find LLaMa-3-8B-Instruct have 69.0% and 65.9% accuracy, Qwen1.5-7B-Chat has 58.2% and 61.3% accuracy on the reasoning set and fullest respectively, all underperforming PINE (numbers reported in Table 3).

4.3 RESULTS ON RETRIEVAL-AUGMENTED QUESTION-ANSWERING

PINE performs better than baselines, on-par with vanilla inference on average while not being affected by the worst case. Models tend to perform better when the gold-standard document is at the beginning and the end of all documents in retrieval-augmented question-answers. Figure 4 (a) shows the results on LLaMa-3-70B-Instruct when 10 or 20 documents were presented. First, it is easy to conclude that all baselines are much worse than PINE (the pink line), which is consistent to the previous experiment. Second, PINE achieves on-par performance on average compared with vanilla inference while being inter-document position invariant. Specifically, PINE is slightly better/worse than vanilla inference with the gap $+1.2/-2.0$ when there are 10 and 20 documents in total. We hypothesize that the slight performance drop of PINE for the 20 document setting is due to the performance drop of document importance score computation in PINE when presented with many documents. However, PINE is position-invariant, therefore does not be affected by the worst case (the bottom of blue solid curves). Third, the height generally becomes smaller between blue and brown solid lines in Figure 4 (a), and between the blue and red solid lines in Figure 4 (b) when the gold-standard document position increases, reflecting the causal attention generally prefers distant content, which is consistent to the hypothesis in Section 1. The brown line in Figure 4 (a) and red line (b) generally reflect recency bias brought by RoPE, which is consistent to previous works (Su et al., 2024; Peysakhovich & Lerer, 2023).

Table 4: The result of molecule generation on QM9 dataset. PINE improves model performance in 5 out of 6 criteria.

Model	α	ϵ_{HOMO}	ϵ_{LUMO}	$\Delta\epsilon$	μ	C_v
LLama	6.3997	103.93	53.4	99.13	3.4112	4.3785
Llama + PINE	6.3702	102.15	53.09	98.27	3.4917	4.2886

PINE performs better than other position assignment methods. So far, our experiments show that bidirectional inter-document attention is the better design choice than the masked one. However, there are still several design options for the position assignment, as discussed in Section 3.4. The first option is to re-assign position reversely, and the other is to use PINE without position re-assignment (i.e., use input document positions when they serve as keys). To gain a deeper understanding, we extend the retrieval-augmented QA experiments with the two mentioned alternative position assignment methods, and the results are presented in Figure 4 (b). The figure tells us that PINE is slightly better than PINE without position re-assignment on average (+0.3. The gap becomes larger when 20 documents are presented: +1.5). Position re-assignment reversely has relatively worse results, showing that PINE is a better design choice, which is consistent with the intuitive analysis mentioned in Section 3.4. Although position re-assignment seems only to bring less gains than bidirectional attention mask, it is required to complete the proof that PINE can *eliminate* the position bias. Therefore, PINE without position re-assignment may suffice if one does not aim to eliminate the position bias and cares more about efficiency (no extra $\mathcal{O}(nk \log k)$ sorting cost).

4.4 RESULTS ON MOLECULE GENERATION AND MATH REASONING

PINE improves model performance on 5 out of 6 criteria in molecule generation . Table 4 shows the results of molecule generation. The consistent gain in 5 out 6 criteria shows the effectiveness of PINE.

PINE improves math reasoning capabilities. Figure 5 shows the results of Qwen1.5 models on R-GSM dataset. It can be shown that PINE outperforms vanilla inference for both small 7B models and large 110B models.

4.5 COMPUTATIONAL OVERHEAD

Section 3.4 briefly discusses the computational overhead, with a conclusion that PINE’s efficiency is still doable. In our experiments, we find the wall time of PINE is $\sim 2x$ and $\sim 8x$ of the vanilla inference on the LM-as-a-judge task and retrieval-augmented QA task with 20 documents, which is acceptable at least during experiments. However, we did not specially optimize codes to accelerate PINE, and our implementation still contains a “for” loop. Therefore, we believe there is room to accelerate PINE. Compared with the time overhead, the memory overhead is small and PINE can be run with 70B models on 3x A100 80G on the retrieval-augmented QA task, which requires the same number of GPUs as the vanilla inference. Since efficiency is not the main focus of this paper, we leave this as our future work.

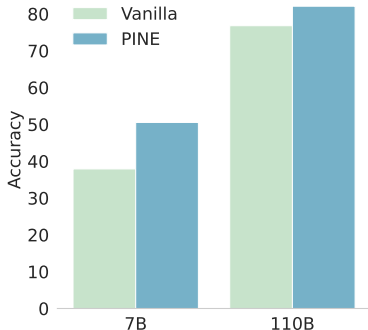


Figure 5: Math reasoning results of Qwen1.5 series on R-GSM subset. PINE improves the reasoning accuracy by 12.6% and 5.3% with 7B and 110B models respectively compared with vanilla inference.

5 CONCLUSION, LIMITATIONS AND FUTURE WORK

We propose a novel train-free zero-shot approach to eliminate the position bias mechanically. The core idea is to make every input documents equally affected by the attention mask and position embedding. However, PINE requires extra computation. We believe there is room to improve the efficiency with more efficient implementation, and we leave this as our future work.

540 REPRODUCIBILITY STATEMENT

541

542 Experiment details are described in Section 4.1 and Appendix D. Codes are uploaded in the Supple-
 543 mentary Material. A complete proof of the lemma and theorem occurred in Section 3.3 are presented
 544 in Appendix B.

545

546 REFERENCES

547

548 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
 549 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
 550 *arXiv preprint arXiv:2303.08774*, 2023.

551

552 Dyah Adila, Shuai Zhang, Boran Han, and Yuyang Wang. Discovering bias in latent space: An
 553 unsupervised debiasing approach. *arXiv preprint arXiv:2406.03631*, 2024.

554 Meta AI. Build the future of ai with meta llama 3, 2024. URL [https://llama.meta.com/
 555 llama3](https://llama.meta.com/llama3).

556

557 Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky,
 558 Bin Bao, Peter Bell, David Berard, Evgeni Burovski, et al. Pytorch 2: Faster machine learning
 559 through dynamic python bytecode transformation and graph compilation. In *Proceedings of the
 560 29th ACM International Conference on Architectural Support for Programming Languages and
 561 Operating Systems, Volume 2*, pp. 929–947, 2024.

562 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,
 563 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu,
 564 Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan,
 565 Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin
 566 Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng
 567 Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou,
 568 Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*,
 569 2023.

570 Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and
 571 Saḡnak Taşırlar. Introducing our multimodal models, 2023. URL [https://www.adept.ai/
 572 blog/fuyu-8b](https://www.adept.ai/blog/fuyu-8b).

573 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 574 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 575 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

577 Tianle Cai, Kaixuan Huang, Jason D. Lee, and Mengdi Wang. Scaling in-context demonstrations
 578 with structured attention. In *Workshop on Efficient Systems for Foundation Models @ ICML2023*,
 579 2023. URL <https://openreview.net/forum?id=jH580PKkPw>.

580 Xinyun Chen, Ryan A. Chi, Xuezhi Wang, and Denny Zhou. Premise order matters in reason-
 581 ing with large language models. *ArXiv*, abs/2402.08939, 2024a. URL [https://api.
 582 semanticscholar.org/CorpusID:267657940](https://api.semanticscholar.org/CorpusID:267657940).

583

584 Xinyun Chen, Ryan A Chi, Xuezhi Wang, and Denny Zhou. Premise order matters in reasoning with
 585 large language models. *arXiv preprint arXiv:2402.08939*, 2024b.

586 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
 587 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh,
 588 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam
 589 Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James
 590 Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Lev-
 591 skaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin
 592 Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph,
 593 Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M.
 Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon

- 594 Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark
595 Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean,
596 Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways, 2022.
597
- 598 Ricardo Dominguez-Olmedo, Moritz Hardt, and Celestine Mendler-Düner. Questioning the survey
599 responses of large language models. *arXiv preprint arXiv:2306.07951*, 2023.
- 600 Yaru Hao, Yutao Sun, Li Dong, Zhixiong Han, Yuxian Gu, and Furu Wei. Structured prompting:
601 Scaling in-context learning to 1,000 examples. *arXiv preprint arXiv:2212.06713*, 2022.
602
- 603 Cheng-Yu Hsieh, Yung-Sung Chuang, Chun-Liang Li, Zifeng Wang, Long T Le, Abhishek Kumar,
604 James Glass, Alexander Ratner, Chen-Yu Lee, Ranjay Krishna, et al. Found in the middle: Calibrat-
605 ing positional attention bias improves long context utilization. *arXiv preprint arXiv:2406.16008*,
606 2024.
- 607 He Junqing, Pan Kunhao, Dong Xiaoqun, Song Zhuoyang, Liu Yibo, Liang Yuxin, Wang Hao, Sun
608 Qianguo, Zhang Songxin, Xie Zejian, et al. Never lost in the middle: Improving large language
609 models via attention strengthening question answering. *arXiv preprint arXiv:2311.09198*, 2023.
610
- 611 Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Payel Das, and Siva Reddy.
612 The impact of positional encoding on length generalization in transformers. *Advances in Neural
613 Information Processing Systems*, 36, 2024.
- 614 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
615 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
616 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating
617 Systems Principles*, 2023.
- 618 Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu,
619 Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi.
620 Rewardbench: Evaluating reward models for language modeling. [https://huggingface.
621 co/spaces/allenai/reward-bench](https://huggingface.co/spaces/allenai/reward-bench), 2024a.
622
- 623 Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu,
624 Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. Rewardbench: Evaluating reward models
625 for language modeling. *arXiv preprint arXiv:2403.13787*, 2024b.
- 626 Sha Li, Ruining Zhao, Manling Li, Heng Ji, Chris Callison-Burch, and Jiawei Han. Open-domain
627 hierarchical event schema induction by incremental prompting and verification. In *Proc. The 61st
628 Annual Meeting of the Association for Computational Linguistics (ACL2023)*, 2023.
629
- 630 Xiner Li, Limei Wang, Youzhi Luo, Carl Edwards, Shurui Gui, Yuchao Lin, Heng Ji, and Shuiwang
631 Ji. Geometry informed tokenization of molecules for language model generation, 2024. URL
632 <https://arxiv.org/abs/2408.10120>.
- 633 Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and
634 Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the
635 Association for Computational Linguistics*, 12:157–173, 2024.
- 636 OpenAI. Gpt-4v(ision) system card. 2023. URL [https://api.semanticscholar.org/
637 CorpusID:263218031](https://api.semanticscholar.org/CorpusID:263218031).
- 638
- 639 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
640 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style,
641 high-performance deep learning library. *Advances in neural information processing systems*, 32,
642 2019.
- 643 Alexander Peysakhovich and Adam Lerer. Attention sorting combats recency bias in long context
644 language models. *arXiv preprint arXiv:2310.01427*, 2023.
645
- 646 Raghunathan Ramakrishnan, Pavlo O. Dral, Pavlo O. Dral, Matthias Rupp, and O. Anatole von
647 Lilienfeld. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific Data*, 1,
2014. URL <https://api.semanticscholar.org/CorpusID:15367821>.

- 648 Nir Ratner, Yoav Levine, Yonatan Belinkov, Ori Ram, Inbal Magar, Omri Abend, Ehud Karpas,
649 Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. Parallel context windows for large
650 language models. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of
651 the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,
652 pp. 6383–6402, Toronto, Canada, July 2023. Association for Computational Linguistics. doi:
653 10.18653/v1/2023.acl-long.352. URL [https://aclanthology.org/2023.acl-long.
654 352](https://aclanthology.org/2023.acl-long.352).
- 655 Víctor Garcia Satorras, Emiel Hooeboom, and Max Welling. E(n) equivariant graph neural net-
656 works. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Con-
657 ference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp.
658 9323–9332. PMLR, 18–24 Jul 2021. URL [https://proceedings.mlr.press/v139/
659 satorras21a.html](https://proceedings.mlr.press/v139/satorras21a.html).
- 660 Lin Shi, Weicheng Ma, and Soroush Vosoughi. Judging the judges: A systematic investigation of
661 position bias in pairwise comparative assessments by llms. *arXiv preprint arXiv:2406.07791*, 2024.
662
- 663 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced
664 transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
665
- 666 Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze
667 Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. Lamda: Language models for dialog
668 applications. *arXiv preprint arXiv:2201.08239*, 2022.
- 669 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
670 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
671 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 672 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
673 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing
674 systems*, 30, 2017.
675
- 676 Yiwei Wang, Yujun Cai, Muhao Chen, Yuxuan Liang, and Bryan Hooi. Primacy effect of ChatGPT.
677 In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference
678 on Empirical Methods in Natural Language Processing*, pp. 108–115, Singapore, December
679 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.8. URL
680 <https://aclanthology.org/2023.emnlp-main.8>.
- 681 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
682 Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick
683 von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gug-
684 ger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art
685 natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods
686 in Natural Language Processing: System Demonstrations*, pp. 38–45, Online, October 2020.
687 Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-demos.6. URL
688 <https://aclanthology.org/2020.emnlp-demos.6>.
- 689 Penghao Wu and Saining Xie. v^* : Guided visual search as a core mechanism in multimodal llms.
690 *arXiv preprint arXiv:2312.14135*, 2023.
691
- 692 Zhichao Xu, Daniel Cohen, Bei Wang, and Vivek Srikumar. In-context example ordering guided by
693 label distributions. *arXiv preprint arXiv:2402.11447*, 2024.
- 694 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
695 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint
696 arXiv:2407.10671*, 2024.
- 697 Yijiong Yu, Huiqiang Jiang, Xufang Luo, Qianhui Wu, Chin-Yew Lin, Dongsheng Li, Yuqing Yang,
698 Yongfeng Huang, and Lili Qiu. Mitigate position bias in large language models via scaling a single
699 dimension. *arXiv preprint arXiv:2406.02536*, 2024.
700
- 701 Kaiyi Zhang, Ang Lv, Yuhan Chen, Hansen Ha, Tao Xu, and Rui Yan. Batch-icl: Effective, efficient,
and order-agnostic in-context learning. *arXiv preprint arXiv:2401.06469*, 2024a.

702 Meiru Zhang, Zaiqiao Meng, and Nigel Collier. Attention instruction: Amplifying attention in the
703 middle via prompting. *arXiv preprint arXiv:2406.17095*, 2024b.
704

705 Tony Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving
706 few-shot performance of language models. In *International Conference on Machine Learning*,
707 2021. URL <https://api.semanticscholar.org/CorpusID:231979430>.

708 Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models
709 are not robust multiple choice selectors. In *The Twelfth International Conference on Learning*
710 *Representations*, 2024a. URL <https://openreview.net/forum?id=shr9PXz7T0>.
711

712 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
713 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and
714 chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024b.

715 Lianghui Zhu, Xinggang Wang, and Xinlong Wang. Judgelm: Fine-tuned large language models are
716 scalable judges. *arXiv preprint arXiv:2310.17631*, 2023.
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A ANOTHER EXAMPLE OF POSITION BIAS IN VLMS

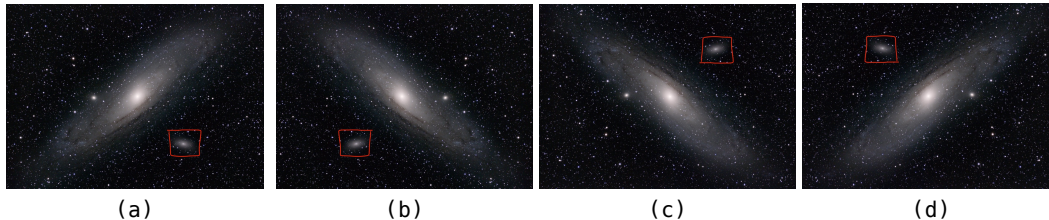


Figure 6: Image of Andromeda Galaxy. GPT-4V is asked to answer the galaxy of the redbox (M110). (b), (c), (d) images are generated by flipping (a). Results show that GPT-4V can answer correctly on (c) and (d) and incorrectly on (a) and (b), revealing position bias in VLMS.

To show VLMS suffer from position bias, we ask models to identify a small object appeared in the full image. We ask GPT-4V (OpenAI, 2023) to identify one satellite galaxy M110 of Andromeda Galaxy (M31). The central objective in the Figure 6 (a) is the Andromeda Galaxy and the galaxy in the red box is M110 satellite galaxy. Another satellite galaxy is M32 which is the brightest part near the Andromeda Galaxy.

We flip the image (a) and obtains images (b), (c), and (d). We prompt GPT-4V “What galaxy is in the red box?”, and here are the responses we collected:

(a): The galaxy within the red box is **M32**, also known as Messier 32. It’s a dwarf elliptical galaxy and is a satellite of the Andromeda Galaxy (M31), which is the prominent galaxy in the image. M32 is one of the more well-studied dwarf elliptical galaxies, largely due to its proximity to our own galaxy and its interactions with the Andromeda Galaxy.

(b): The object inside the red box in the image appears to be a smaller galaxy located near the larger, more prominent galaxy (Andromeda Galaxy, M31). This smaller galaxy is likely **M32**, a satellite galaxy of the Andromeda Galaxy. M32 is a compact elliptical galaxy and is often visible in images as a smaller, dense, bright spot near Andromeda, as seen in the boxed area of your image.

(c): The galaxy within the red box is **M110**, also known as NGC 205. M110 is a dwarf elliptical galaxy, which, like M32, is a satellite of the Andromeda Galaxy (M31). It is one of the members of the Local Group of galaxies, which includes the Milky Way, the Andromeda Galaxy, and about 54 other known galaxies.

(d): The galaxy in the red box appears to be a satellite galaxy of the Andromeda Galaxy, which is the prominent galaxy visible in the rest of the image. This smaller galaxy is likely **M110** (NGC 205), which is a dwarf elliptical galaxy and a companion to the Andromeda Galaxy, M31. It’s one of the several satellite galaxies gravitationally bound to Andromeda, visible here as a faint, elongated object in the outlined area.

We can find that models answer corrected when M110 is at the top of the image, revealing that VLMS also suffer from the position bias. The position bias may lead unreliable VLMS when fine-grained image analysis are needed (e.g., small object detection (Wu & Xie, 2023)).

B COMPLETE PROOF

This section provided a complete proof to show PINE can eliminate position bias.

To simplify the notation and without loss of generality (w.l.o.g), we still use examples in Section 3.1.

Theorem 1. *Given an input, if \mathbf{H}_{PINE} is applied to every layer, attention head, and token to replace the conventional attention computation, then the model outputs are inter-document position-invariant representations.*

First, the embedding layer is not a function of input documents positions. Suppose that the i th layer’s input hidden states are not a function of input documents positions, then within each layer:

- The attention hidden states are not a function of input documents positions (Lemma).
- The Layernorm, FFN outputs are not a function of input documents positions.
- Therefore, the output hidden states of i th transformer layer, i.e., the input hidden states of $i + 1$ th transformer layer, are not a function of input documents positions.

Using mathematical induction, we know the final outputs are not a function of input documents positions.

Proof ends.

Notes on the proof:

- PINE needs to be applied on each layer, attention heads, and tokens to satisfy the above proof.
- The extra big \mathcal{O} computation cost is purely come from the position re-assignment step: $\mathcal{O}(k \log k)$ for sorting k documents. Since we need to repeat this step for every token, the extra computation cost is $\mathcal{O}(nk \log k)$, where n is the number of tokens.
- Although position re-assignment brings an extra computational cost, it is a must to complete the proof. Removing this step will make PINE unable to “eliminate” position bias. Similarly, a bidirectional attention mask is also a must to complete the proof.
- PINE is not limited to specific position encoding algorithms.

C FULL RESULTS OF REWARD BENCH

Table 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15 present the full results of the reward bench. After inspecting the error cases, we categorize the performance drop in the Chat-Hard and Safety subsets into two main aspects:

- The instruction-following capabilities become a bit worse. For example, LLMs tend to solve the “user question” instead of comparing two responses, or LLMs do not output answers in requested formats, causing parsing failures when computing performance scores.
- LMs overly focus on helpfulness in safety prompts, therefore causing performance degradation in the Safety dataset.

However, the positive effect of PINE (i.e., eliminating position bias) is more significant than these negative effects; therefore, the overall PINE is still beneficial to models.

D IMPLEMENTATION DETAILS

D.1 EXPERIMENT SETTING

For reproducibility, the generation temperature is set to 0. We use PyTorch (Ansel et al., 2024; Paszke et al., 2019),⁶ Transformers (Wolf et al., 2020),⁷ and vLLM (Kwon et al., 2023) for our experiments.⁸ All experiments are launched with a single node of 8x A100 80G with SXM connection. 70B and 110B models are launched with 3x and 4x A100, and other model sizes can be launched with 1x A100.

⁶Customized license. <https://github.com/pytorch/pytorch>

⁷Apache-2.0 license. <https://huggingface.co/docs/transformers/en/index>

⁸Apache-2.0 license. <https://github.com/vllm-project/vllm>.

Table 5: Full results of Table 2. Vanilla denotes the normal inference, (GT at A) means the ground truth chosen response is presented at the first, and (GT at B) indicates at the second. PINE consistently improves LM’s performance across different model sizes. Consistent to Table 1, we color the subsets with severe position bias cyan. It can be observed that PINE generally improves performance on cyan subsets by a large margin, which is consistent to our motivation and goal.

Model	Size	Method	Chat	Chat-Hard	Safety	Reasoning	Avg.
LLaMa-3 -Instruct	8B	Vanilla (GT at A)	90.1	35.2	64.6	80.3	67.5
		Vanilla (GT at B)	85.3	48.7	65.3	66.0	66.3
		Vanilla	85.3	41.6	67.0	65.3	64.8
		PINE	85.6	41.5	66.5	73.4	66.7 _{+1.9}
	70B	Vanilla (GT at A)	98.6	52.0	73.6	87.8	78.0
		Vanilla (GT at B)	93.9	62.1	69.8	80.3	76.5
		Vanilla	97.4	58.3	69.6	78.9	76.0
		PINE	96.9	57.4	67.7	87.6	77.4 _{+1.4}
	1.8B	Vanilla (GT at A)	31.7	30.0	40.3	43.3	36.3
		Vanilla (GT at B)	69.4	72.6	65.7	57.2	66.2
		Vanilla	49.7	50.9	52.0	48.4	50.3
		PINE	30.0	59.9	61.4	60.1	52.9 _{+2.6}
Qwen-1.5 -Chat	4B	Vanilla (GT at A)	32.8	24.8	17.4	42.8	29.5
		Vanilla (GT at B)	86.6	74.5	82.9	62.3	76.6
		Vanilla	58.9	48.7	50.9	54.1	53.1
		PINE	73.0	45.2	53.7	61.0	58.2 _{+5.1}
	7B	Vanilla (GT at A)	85.5	35.9	62.4	62.1	61.4
		Vanilla (GT at B)	77.1	47.4	59.5	54.3	59.6
		Vanilla	77.5	44.2	62.6	59.3	60.9
		PINE	85.8	38.7	58.6	63.0	61.5 _{+0.6}
	32B	Vanilla (GT at A)	93.6	47.7	77.1	78.3	74.2
		Vanilla (GT at B)	91.9	52.2	81.6	73.6	74.8
		Vanilla	92.7	51.2	80.5	66.8	72.8
		PINE	93.0	49.8	79.7	76.7	74.8 _{+2.0}
72B	Vanilla (GT at A)	95.7	59.0	80.8	83.0	79.6	
	Vanilla (GT at B)	89.0	46.5	73.7	68.7	69.5	
	Vanilla	94.0	51.4	77.8	68.2	72.8	
	PINE	93.9	46.1	78.2	69.0	71.8 _{-1.1}	
72B (Qwen 2.5)	Vanilla (GT at A)	97.5	71.9	85.7	93.7	87.2	
	Vanilla (GT at B)	95.0	67.5	83.4	76.0	80.5	
	Vanilla	96.6	68.0	83.3	85.5	83.4	
	PINE	96.6	67.1	83.0	91.3	74.5 _{+1.1}	
110B	Vanilla (GT at A)	98.6	70.5	89.6	90.0	87.2	
	Vanilla (GT at B)	91.1	59.2	79.5	73.0	75.7	
	Vanilla	96.2	66.7	83.7	78.0	81.1	
	PINE	95.5	64.8	85.0	86.2	82.9 _{+1.7}	

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Table 6: Full version of Table 3. PINE achieves superior performance to baseline models, performing 4.8% and 4.7% better than the best performed baseline on two models.

Model	Method	Chat	Chat-Hard	Safety	Reasoning	Avg.
LLaMa-3 8B-Instruct	NIA (GT at A)	81.0	40.7	59.7	43.7	56.3
	NIA (GT at B)	81.0	49.7	65.8	66.7	65.8
	NIA	80.9	46.7	64.0	55.9	61.9
	PCW	78.6	46.8	64.8	56.5	61.7
	SP	79.6	43.3	65.0	55.4	60.8
	PINE	85.6	41.5	66.5	73.4	66.7 _{+4.8}
	NIA (GT at A)	67.7	57.2	59.6	60.7	61.3
Qwen-1.5 7B-Chat	NIA (GT at B)	67.9	35.9	61.0	44.1	52.2
	NIA	74.9	43.5	57.4	51.4	56.8
	PCW	67.2	42.0	58.3	53.4	55.2
	SP	69.4	41.8	58.0	52.4	55.4
	PINE	85.8	38.7	58.6	63.0	61.5 _{+4.7}

Table 7: Meta-Llama-3-8B-Instruct results on RewardBench

Dataset	Vanilla	PINE
alpacaeval-easy	91.0	90.0
alpacaeval-hard	91.6	90.0
alpacaeval-length	71.6	77.4
donotanswer	45.2	38.2
hep-cpp	78.7	82.6
hep-go	77.1	86.6
hep-java	73.5	82.9
hep-js	74.4	84.1
hep-python	79.0	85.7
hep-rust	74.4	81.4
llmbar-adver-GPTInst	23.4	24.5
llmbar-adver-GPTOut	63.8	67.0
llmbar-adver-manual	40.2	34.8
llmbar-adver-neighbor	20.9	16.8
llmbar-natural	66.0	74.5
math-prm	54.5	62.8
mt-bench-easy	92.9	87.5
mt-bench-hard	68.9	64.9
mt-bench-med	83.8	80.0
refusals-dangerous	71.5	74.0
refusals-offensive	76.0	73.5
xstest-should-refuse	70.5	71.8
xstest-should-respond	72.0	76.4

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Table 8: Meta-Llama-3-70B-Instruct results on RewardBench

Dataset	Vanilla	PINE
alpacaeval-easy	100.0	100.0
alpacaeval-hard	100.0	100.0
alpacaeval-length	91.1	89.5
donotanswer	47.1	48.2
hep-cpp	92.7	92.1
hep-go	89.9	97.0
hep-java	92.1	97.0
hep-js	93.3	95.1
hep-python	90.9	95.4
hep-rust	89.0	91.5
llmbar-adver-GPTInst	55.4	57.6
llmbar-adver-GPTOut	73.4	76.6
llmbar-adver-manual	53.3	47.8
llmbar-adver-neighbor	32.8	28.7
llmbar-natural	83.0	84.0
math-prm	66.4	80.5
mt-bench-easy	100.0	100.0
mt-bench-hard	78.4	75.7
mt-bench-med	97.5	97.5
refusals-dangerous	63.5	62.5
refusals-offensive	66.5	66.5
xstest-should-refuse	68.8	63.3
xstest-should-respond	96.8	96.4

Table 9: Qwen1.5-1.8B-Chat results on RewardBench

Dataset	Vanilla	PINE
alpacaeval-easy	47.5	17.0
alpacaeval-hard	56.3	13.7
alpacaeval-length	50.0	45.8
donotanswer	54.0	52.6
hep-cpp	49.4	51.5
hep-go	54.3	48.8
hep-java	52.7	49.4
hep-js	48.2	47.6
hep-python	49.4	52.1
hep-rust	54.6	50.3
llmbar-adver-GPTInst	44.0	76.6
llmbar-adver-GPTOut	55.3	40.4
llmbar-adver-manual	44.6	64.1
llmbar-adver-neighbor	56.7	66.8
llmbar-natural	49.0	51.5
math-prm	45.5	70.2
mt-bench-easy	39.3	57.1
mt-bench-hard	54.1	35.1
mt-bench-med	46.2	45.0
refusals-dangerous	48.5	88.0
refusals-offensive	49.5	54.5
xstest-should-refuse	53.2	53.9
xstest-should-respond	52.2	68.8

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Table 10: Qwen1.5-4B-Chat results on RewardBench

Dataset	Vanilla	PINE
alpacaeval-easy	59.5	77.5
alpacaeval-hard	62.1	80.5
alpacaeval-length	60.5	70.0
donotanswer	54.4	18.4
hep-cpp	50.0	50.0
hep-go	50.3	51.8
hep-java	49.1	51.2
hep-js	49.7	49.4
hep-python	50.0	53.0
hep-rust	50.6	50.6
llmbar-adver-GPTInst	36.4	38.6
llmbar-adver-GPTOut	54.3	51.1
llmbar-adver-manual	51.1	39.1
llmbar-adver-neighbor	52.2	42.5
llmbar-natural	48.0	53.5
math-prm	58.2	70.9
mt-bench-easy	44.6	75.0
mt-bench-hard	58.1	48.6
mt-bench-med	56.2	50.0
refusals-dangerous	43.0	31.0
refusals-offensive	47.0	71.0
xstest-should-refuse	53.6	64.3
xstest-should-respond	51.0	71.0

Table 11: Qwen1.5-7B-Chat results on RewardBench

Dataset	Vanilla	PINE
alpacaeval-easy	74.0	91.0
alpacaeval-hard	90.0	96.8
alpacaeval-length	65.8	74.7
donotanswer	19.9	11.0
hep-cpp	60.4	76.8
hep-go	61.9	69.2
hep-java	56.1	74.1
hep-js	59.5	68.6
hep-python	63.4	66.8
hep-rust	62.8	65.9
llmbar-adver-GPTInst	40.2	23.4
llmbar-adver-GPTOut	53.2	45.7
llmbar-adver-manual	35.9	35.9
llmbar-adver-neighbor	21.6	19.8
llmbar-natural	71.0	70.0
math-prm	57.9	55.8
mt-bench-easy	87.5	85.7
mt-bench-hard	62.2	55.4
mt-bench-med	77.5	72.5
refusals-dangerous	49.0	40.5
refusals-offensive	86.0	86.0
xstest-should-refuse	74.0	63.6
xstest-should-respond	75.6	86.6

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

Table 12: Qwen1.5-32B-Chat results on RewardBench

Dataset	Vanilla	PINE
alpacaeval-easy	97.0	97.0
alpacaeval-hard	98.9	98.9
alpacaeval-length	81.1	82.1
donotanswer	44.5	41.2
hep-cpp	87.8	91.5
hep-go	80.5	93.9
hep-java	88.7	96.3
hep-js	84.5	95.7
hep-python	86.3	93.6
hep-rust	82.0	88.1
llmbar-adver-GPTInst	43.5	34.8
llmbar-adver-GPTOut	68.1	57.4
llmbar-adver-manual	32.6	37.0
llmbar-adver-neighbor	25.0	28.4
llmbar-natural	83.0	85.0
math-prm	48.7	60.3
mt-bench-easy	96.4	92.9
mt-bench-hard	81.1	75.7
mt-bench-med	92.5	95.0
refusals-dangerous	80.0	80.0
refusals-offensive	99.0	99.0
xstest-should-refuse	90.6	89.3
xstest-should-respond	84.4	85.6

Table 13: Qwen1.5-72B-Chat results on RewardBench

Dataset	Vanilla	PINE
alpacaeval-easy	98.0	98.0
alpacaeval-hard	97.4	97.9
alpacaeval-length	85.3	84.2
donotanswer	39.0	38.2
hep-cpp	88.1	89.6
hep-go	85.4	92.1
hep-java	87.2	90.9
hep-js	90.9	90.9
hep-python	87.2	89.6
hep-rust	88.1	87.2
llmbar-adver-GPTInst	44.6	33.7
llmbar-adver-GPTOut	57.4	61.7
llmbar-adver-manual	41.3	39.1
llmbar-adver-neighbor	28.0	20.9
llmbar-natural	84.0	81.0
math-prm	48.5	47.9
mt-bench-easy	96.4	100.0
mt-bench-hard	70.3	62.2
mt-bench-med	95.0	92.5
refusals-dangerous	75.5	73.0
refusals-offensive	94.0	95.0
xstest-should-refuse	87.7	91.2
xstest-should-respond	86.8	85.0

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

Table 14: Qwen2.5-72B-Instruct results on RewardBench

Dataset	Vanilla	PINE
alpacaeval-easy	99.0	99.0
alpacaeval-hard	97.9	98.9
alpacaeval-length	91.6	89.5
donotanswer	48.5	52.9
hep-cpp	95.7	95.7
hep-go	97.0	98.8
hep-java	98.8	97.6
hep-js	94.5	98.2
hep-python	98.8	98.8
hep-rust	94.5	97.6
llmbar-adver-GPTInst	66.3	68.5
llmbar-adver-GPTOut	76.6	72.3
llmbar-adver-manual	65.2	63.0
llmbar-adver-neighbor	41.8	44.0
llmbar-natural	92.0	87.0
math-prm	74.5	84.8
mt-bench-easy	100.0	100.0
mt-bench-hard	94.6	91.9
mt-bench-med	97.5	100.0
refusals-dangerous	78.0	82.0
refusals-offensive	95.0	92.0
xstest-should-refuse	92.2	89.6
xstest-should-respond	95.2	93.6

Table 15: Qwen1.5-110B-Chat results on RewardBench

Dataset	Vanilla	PINE
alpacaeval-easy	95.0	97.0
alpacaeval-hard	98.9	98.9
alpacaeval-length	93.7	88.4
donotanswer	51.5	55.9
hep-cpp	87.8	92.1
hep-go	83.8	94.8
hep-java	86.6	94.8
hep-js	90.5	92.4
hep-python	83.8	93.9
hep-rust	85.7	90.9
llmbar-adver-GPTInst	70.1	65.2
llmbar-adver-GPTOut	72.3	61.7
llmbar-adver-manual	60.9	65.2
llmbar-adver-neighbor	44.8	41.4
llmbar-natural	86.5	90.0
math-prm	69.6	79.2
mt-bench-easy	98.2	100.0
mt-bench-hard	83.8	83.8
mt-bench-med	97.5	97.5
refusals-dangerous	76.0	84.0
refusals-offensive	97.0	97.0
xstest-should-refuse	91.6	91.9
xstest-should-respond	95.6	92.4

D.2 MOLECULE GENERATION AND MATH REASONING TASK DETAILS.

Molecule Generation. In this task, the input contains several properties that are interchangeable, and LMs are asked to generate molecules that satisfy these properties. We train such an LM with QM9 (Ramakrishnan et al., 2014) dataset. The QM9 dataset collects over 130k 3D molecules with 3D structures (Li et al., 2024) calculated by density functional theory (DFT). Each molecule in QM9 has less than 9 heavy atoms, and its chemical elements all belong to H, C, N, O, F. We take six quantum property values as the conditional input to LMs and train LMs to generate molecules with the conditioned quantum property values. We split the training dataset of QM9 to two subsets where each subset has 50k samples, and train LMs and an EGNN-based quantum property prediction models (Satorras et al., 2021) on these two subsets, respectively. The six quantum properties are polarizability (α), HOMO energy (ϵ_{HOMO}), LUMO energy (ϵ_{LUMO}), HOMO-LUMO gap ($\Delta\epsilon$), dipole moment (μ) and heat capacity at 298.15K (C_v). The LM is a 8-layer Llama model with 8 attention heads and 768 hidden dimensions. To evaluate the performance, we sample 10000 sets of 6-property conditions, randomize the property order in each condition, and generate molecules conditioned on these property values by the trained LM, and compute the mean absolute difference (MAE) between the given property values and the property values of the generated molecules. Note that we use the trained EGNN-based property prediction models to calculate the property values of the generated molecules.

Math Reasoning. We use R-GSM (Chen et al., 2024a), a subset of GSM8K. This small dataset (which contains 220 problems) is designed to test LMs’ performance with interchangeable premise orders. Problems in the dataset contain several conditions that do not have a progressive relationship. Therefore, their positions are interchangeable. We further clean this dataset to remove problems where conditions do not read smoothly after changing positions (e.g., use pronouns in the first condition but introduce names in the second condition), yielding a small set containing 95 problems. We test Qwen-1.5 models on this dataset.

D.3 PROMPTS

We use the prompts provided by RewardBench (Lambert et al., 2024a) official repo for the LM-as-a-judge task:

System prompt:

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user question displayed below. You should choose the assistant that follows the user’s instructions and answers the user’s question better. Your evaluation should consider factors such as the helpfulness, relevance, accuracy, depth, creativity, and level of detail of their responses. Begin your evaluation by comparing the two responses and provide a short explanation. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better.

User prompt:

[User Question]

 [The Start of Assistant A’s Answer]

 [The End of Assistant A’s Answer]

1242 [The Start of Assistant B’s Answer]
 1243
 1244
 1245 [The End of Assistant B’s Answer]
 1246 We use the prompts of official repo (Liu et al., 2024) for the retrieval-augmented QA experiments:
 1247
 1248 **User Prompt:**
 1249 Write a high-quality one-sentence answer for the given question
 1250 using only the provided search results (some of which might be
 1251 irrelevant).
 1252 Document (Title:):
 1253 Document (Title:):
 1254
 1255 Document (Title:):
 1256 Document (Title:):
 1257 Document (Title:):
 1258 Question:
 1259
 1260 Molecule generation does not need prompts. We use prompts in OpenAI/simple-evals⁹ to evaluate
 1261 R-GSM dataset.
 1262
 1263

1264 E QUALITATIVE EXAMPLES

1265 E.1 LM-AS-A-JUDGE

1266 [User Question]
 1267 Who created the Superman cartoon character?
 1268
 1269 [The Start of Assistant A’s Answer]
 1270 The Superman cartoon character was created by Jerry Siegel and Joe
 1271 Shuster in 1938.
 1272 [The End of Assistant A’s Answer]
 1273 [The Start of Assistant B’s Answer]
 1274 Superman, the iconic comic book superhero, was created by writer
 1275 Jerry Siegel and artist Joe Shuster. Superman first appeared
 1276 in Action Comics #1, which was published by Detective Comics,
 1277 Inc. (later DC Comics) in June 1938. The character’s immense
 1278 popularity established him as one of the most enduring and
 1279 recognizable figures in the superhero genre.
 1280 [The End of Assistant B’s Answer]
 1281 Here are outputs of Qwen1.5-7B-Chat:
 1282
 1283 **Vanilla:** Both AI assistants provide accurate and informative
 1284 responses to the user’s question. Assistant A’s answer is brief
 1285 and directly states the creators’ names and the year of creation,
 1286 while Assistant B offers additional context by mentioning the
 1287 first appearance, the publisher, and the character’s significance
 1288 in the superhero genre. However, since the user asked for the
 1289 creators, Assistant A’s answer is more directly helpful in this
 1290 specific context. Final verdict: [[A]]

1295 ⁹<https://github.com/openai/simple-evals>

1296 **Vanilla (Swapped):** [[B]]. Assistant B's response is more detailed
 1297 and informative, providing additional context about the origin
 1298 of the character in Action Comics #1, the publisher, and its
 1299 significance in the superhero genre. It also includes the names
 1300 of the creators, which is crucial information for the question.
 1301 Assistant A's response is shorter and lacks the historical
 1302 context and specific details. Therefore, Assistant B is better
 1303 in answering the user's question.

1304 **PINE:** [[B]]. Assistant B provides a more comprehensive and
 1305 detailed response to the user's question. It not only mentions
 1306 the creators but also provides additional information about
 1307 the character's debut and the company that published the comic.
 1308 This answer is not only accurate but also offers context that
 1309 highlights the significance of the creation. Assistant A's
 1310 response is brief and lacks the depth and historical context
 1311 provided by Assistant B.

1312

1313 E.2 RETRIEVAL-AUGMENTED QUESTION-ANSWERING

1314

1315 **User Prompt:**

1316 Document (Title: Pearl) south Pacific and Indian Ocean. The
 1317 largest pearl oyster is the "Pinctada maxima", which is roughly
 1318 the size of a dinner plate. South Sea pearls are characterized by
 1319 their large size and warm luster. Sizes up to 14 mm in diameter
 1320 are not uncommon. In 2013, Indonesia Pearl supplied 43 percent
 1321 of South Sea Pearls international market. The other significant
 1322 producers are Australia, Philippines, Myanmar and Malaysia. In
 1323 1914, pearl farmers began growing cultured freshwater pearls using
 1324 the pearl mussels native to Lake Biwa. This lake, the largest and
 1325 most ancient in Japan, lies near the city of Kyoto. The

1326 Document (Title: Laccadive Sea) the gulf as most productive in
 1327 the world. Although extraction of natural pearls is considered
 1328 too expensive in most parts of the world, it is still conducted
 1329 in the gulf. Also collected in large numbers are Shankha mollusks
 1330 ("Xancus pyrum") whose shells are used as a ritual and religious
 1331 object. Other mollusks of the sea are either too scarce or not
 1332 popular in the Indian society and therefore have no commercial
 1333 value. Another traditional occupation in the Laccadive Sea is
 1334 fishing. The annual fish catch is 2,000 to 5,000 tonnes from the
 1335 Lakshadweep islands, which is mostly constituted by tuna

1336 Document (Title: Pearl) including the Cook Islands and Fiji are
 1337 being extensively used for producing cultured pearls. The rarity
 1338 of the black cultured pearl is now a "comparative" issue. The
 1339 black cultured pearl is rare when compared to Chinese freshwater
 1340 cultured pearls, and Japanese and Chinese akoya cultured pearls,
 1341 and is more valuable than these pearls. However, it is more
 1342 abundant than the South Sea pearl, which is more valuable than
 1343 the black cultured pearl. This is simply because the black
 1344 pearl oyster "Pinctada margaritifera" is far more abundant than
 1345 the elusive, rare, and larger south sea pearl oyster "Pinctada
 maxima", which cannot

1346 Document (Title: Pearl powder) Pearl powder Pearl powder () is a
 1347 preparation of crushed pearls used in China and elsewhere for skin
 1348 care and in traditional Chinese medicine. Pearl powder is made
 1349 from freshwater pearls or saltwater pearls below jewellery grade.
 These are sterilised in boiling water and then milled into a fine

1350 powder using stainless steel grinding discs or by milling with
1351 small porcelain balls in moist conditions. The powder is sold
1352 as such or mixed into creams. Pearl powder is widely believed to
1353 help improve the appearance of the skin, and is used as a cosmetic
1354 by royal families in Asia. It

1355 Document (Title: Hyderabad pearl) with white pearls. Recently,
1356 several pearl makers are exporting processed pearls to markets
1357 in Europe and the US. With the capital that they gain from this
1358 marketing, they are able to purchase machinery for advanced
1359 refinement. In particular, equipment that uses enzymes present
1360 in thermophiles is able to substantially improve the process of
1361 refining pearls. Hyderabad pearl Hyderabad is considered the main
1362 pearl trading center in India. The most notable area devoted to
1363 the trade is the village called Chandanpet just outside Hyderabad,
1364 wherein almost the entire population is engaged in the delicate
1365 art of drilling pearls, a skill they

1366 Document (Title: Pearl) pearls". The correct definition of
1367 a South Sea pearl { as described by CIBJO and GIA { is a pearl
1368 produced by the "Pinctada maxima" pearl oyster. South Sea pearls
1369 are the color of their host "Pinctada maxima" oyster { and can
1370 be white, silver, pink, gold, cream, and any combination of these
1371 basic colors, including overtones of the various colors of the
1372 rainbow displayed in the pearl nacre of the oyster shell itself.
1373 South Sea pearls are the largest and rarest of the cultured pearls
1374 { making them the most valuable. Prized for their exquisitely
1375 beautiful órientór lustre,

1376 Document (Title: Chandrani Pearls) year 2007{08 Chandrani Pearls
1377 imported their pearls from Japan, China or Korea. Chandrani
1378 Pearls Chandrani Pearls is a prominent pearl jewelery brand of
1379 India. It pioneered the concept of pearls in India. Chandrani
1380 Pearls's headquarters is at Kolkata in West Bengal. Chandrani
1381 Pearls was started on 24 January 1985 by Mr. Kuldip Nayar, his
1382 wife Mrs. Lakshmi Nayar and his father late Mr. N.C. Nayar in
1383 Kolkata's up market Minto Park area. Chandrani Pearls management
1384 is now assisted by Nisheeth Nayar, sons of Mr. Kuldip Nayar.
1385 Chandrani Pearls have 63 showrooms across 9 states. From a modest
1386 turnover of Rs.

1386 Document (Title: Pearl) For thousands of years, seawater pearls
1387 were retrieved by divers in the Indian Ocean in areas such as the
1388 Persian Gulf, the Red Sea and the Gulf of Mannar. Evidence also
1389 suggest a prehistoric origin to pearl diving in these regions.
1390 Starting in the Han Dynasty (206 BC{220 AD), the Chinese hunted
1391 extensively for seawater pearls in the South China Sea. In the
1392 14th-century Arabian Sea, the traveller Ibn Battuta provided the
1393 earliest known description of pearl diving by means of attaching a
1394 cord to the diver's waist.

1395 Document (Title: Chandrani Pearls) Chandrani Pearls Chandrani
1396 Pearls is a prominent pearl jewelery brand of India. It pioneered
1397 the concept of pearls in India. Chandrani Pearls's headquarters
1398 is at Kolkata in West Bengal. Chandrani Pearls was started on
1399 24 January 1985 by Mr. Kuldip Nayar, his wife Mrs. Lakshmi
1400 Nayar and his father late Mr. N.C. Nayar in Kolkata's up market
1401 Minto Park area. Chandrani Pearls management is now assisted by
1402 Nisheeth Nayar, sons of Mr. Kuldip Nayar. Chandrani Pearls have
1403 63 showrooms across 9 states. From a modest turnover of Rs. 2.75
lakhs in 1990 to a turnover of Rs. 15.00 crores in

1404 Document (Title: Hyderabad pearl) Hyderabad pearl Hyderabad
 1405 is considered the main pearl trading center in India. The most
 1406 notable area devoted to the trade is the village called Chandanpet
 1407 just outside Hyderabad, wherein almost the entire population is
 1408 engaged in the delicate art of drilling pearls, a skill they have
 1409 practiced for generations. This practice also makes the city
 1410 one of the largest pearl drilling locations in India. The pearl
 1411 industry in Hyderabad flourished due to the patronage of the Qutub
 1412 Shahi kings and the Asaf Jahis, who were said to have an affinity
 1413 for sparkling jewels. The pearls were not only part of

1414 Question: in which sea pearl is found in india?

1415 Here are Llama-3-70B-Instruct results:

1416 **Vanilla:** According to the search results, pearls are found in the
 1417 **Laccadive** Sea in India.

1418 **PINE:** Pearls are found in the **Indian Ocean**, specifically in the
 1419 Gulf of Mannar and the Laccadive Sea, off the coast of India.

1420

1421 E.3 MOLECULE GENERATION

1422

1423 The task has properties as input, and molecule structure as output. Below is one example:

1424 **Input:** cv*-30.887* alpha*81.68* mu*2.0332* homo*-5.8749* gap*7.8668*
 1425 lumo*1.9919*

1426 **Output:** H 0.00 0.00° 0.00° O 0.96 1.57° 0.00° C 1.95 2.33° 0.00° H
 1427 2.35 2.50° 0.70° H 2.32 2.53° -0.71° C 3.22 2.00° -0.02° H 3.46
 1428 1.80° 0.24° C 4.50 2.23° 0.00° H 4.99 2.22° 0.26° H 4.55 2.46°
 1429 -0.12° C 5.58 2.08° -0.18° H 6.55 2.07° -0.09° N 5.89 2.05° -0.46°
 1430 C 5.57 2.18° -0.71° H 4.57 2.25° -0.78° H 6.29 2.32° -0.72° H 5.98
 1431 2.06° -0.86° C 5.26 1.85° -0.34° H 6.06 1.72° -0.33° C 3.81 1.74°
 1432 -0.36° H 3.94 1.47° -0.29° H 3.58 1.75° -0.65°

1433

1434 E.4 R-GSM

1435

1436 R-GSM is just a subset of GSM8K, with the premise order changes. Here is an example input:

1437 **Input:** Carmen goes to an auction to win an antique desk. The bids
 1438 on the desk rise by 50 each time and 3 other people each bid once.
 1439 She accepts the opening bid of 200 and continues bidding until she
 1440 wins. Carmen bids after each of the 3 other people and eventually
 1441 wins. How much money, in dollars, does the desk cost her?

1442 **Here** The bids on the desk rise by 50 each time and 3 other people
 1443 each bid once. **and** She accepts the opening bid of 200 and continues
 1444 bidding until she wins. **are interchangeable.**

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457